

Optimization for AI-Driven Sequential Decision-Making

Thesis by
Lin An

In Partial Fulfillment of the Requirements for the
Degree of
Doctor of Philosophy in Algorithms, Combinatorics, and Optimization



CARNEGIE MELLON UNIVERSITY
Tepper School of Business
Pittsburgh, Pennsylvania, USA

2026
Defended March 27

© 2026

Lin An

ORCID: 0009-0005-7840-2194

All rights reserved

ACKNOWLEDGEMENTS

I confess that I have rarely read the main body of other people's theses, but I have read quite a few acknowledgements. A person's academic contributions are recorded in the theorems, propositions, figures, and proofs that follow. But the acknowledgements are where I feel the presence of a real human being: their frustrations when a proof refuses to work, their hobbies and friendships outside of research, and their gratitude toward the people who helped them along the way.

When I felt frustrated or stuck during my PhD, reading these acknowledgements often gave me courage. They reminded me that I was not alone in this journey. And whenever I experienced moments that I knew I would cherish, I sometimes imagined the day I would write this acknowledgement myself. Still, now that the moment has arrived, it feels a little surreal.

To whoever is reading this: I hope these words can pass on some of the same encouragement, affection, and strength that others have given to me. And of course, I also hope that the hundreds of pages that follow might be helpful to you in some way.

Looking back, I can honestly say that the five years of my PhD have been the happiest five years of my life so far (though I certainly hope the rest of my life will be even more fulfilling). At some point I realized that feeling this way about one's PhD is itself a privilege.

First and foremost, I want to thank my advisors, Andrew Li and Benjamin Moseley. They are the two people who have helped me the most throughout my PhD, and I am deeply grateful for everything they have done for me.

Andrew is the person who shaped my research taste. Through working with him, I learned that research is not only about proving theorems. More importantly, it is about finding problems that truly matter: problems that are impactful for our field and, ideally, for the real world. Andrew has an extraordinary instinct for identifying such problems and modeling them in the right way. I have often felt that he is one of the very best researchers at discovering the right problems and framing them clearly. Watching him work has been one of the most valuable parts of my PhD training. My own research style has been profoundly influenced by Andrew's taste, and I expect it will continue to shape how I approach research throughout my career.

Outside of research itself, Andrew also taught me many of the softer skills needed to succeed in academia. These ranged from how to give a compelling presentation

to surprisingly practical advice about fly-outs, even down to what kind of shoes to wear. I still remember his recommendation: dark brown, pointed, real leather, with sole protection. At first I did not expect this kind of advice to matter very much, but it turned out to be unexpectedly valuable at many moments in my academic career. Andrew was always generous with his time and energy, and our meetings were frequent, lively, and extremely productive. Outside of research, he was also someone I could talk to openly about life and career. I am grateful for the many conversations we shared and for the support he gave me during both easy and difficult times.

Ben was the first professor in the group whom I interacted with, and if it were not for his encouragement, I might never have come to Tepper. From the very beginning, he welcomed me into the group and made me feel that this was a place where I could grow as a researcher. Throughout our collaborations, Ben's enthusiasm for mentoring students has had a strong influence on me. He genuinely cares about what his students want to pursue and always tries to find ways to support them.

One thing that always stood out to me about Ben is how approachable and encouraging he is when discussing research. Our meetings were often relaxed and conversational, yet somehow they were always productive. Even when I arrived with vague ideas or half-formed thoughts, Ben would patiently listen and then point out directions that were both insightful and practical to try. Many times, these suggestions turned out to be exactly what I needed to move the project forward. His ability to quickly see the core of a problem and suggest promising directions has been invaluable to me. I am very grateful for his generosity with his time, his thoughtful guidance, and the encouragement he has given me throughout my PhD.

I would also like to thank the rest of my committee members: Ravi, Joy, and Evelyn. Ravi played an important role early in my PhD by bringing together the group for my first summer project and by introducing me to Glance, which later became a fruitful collaboration. I am also grateful for his advice about both research and academic life. Joy introduced me to ideas and perspectives from marketing that were new to me and has been a wonderful collaborator. Evelyn provided thoughtful feedback during my proposal and defense, and I very much look forward to crossing paths again in our careers.

I am grateful to the faculty in the operations research community at Tepper, Gérard, John, Willem, Javier, Fatma, Karan, and Emily, for creating such a supportive environment. Their kindness and encouragement made this a wonderful place to

grow as a researcher.

I would also like to thank two mentors who influenced me before I arrived at CMU. Solomon Friedberg, my undergraduate advisor, inspired me to pursue research in the first place. The time and effort he devoted to mentoring an undergraduate student were extraordinary and had a lasting impact on me. Yuri Faenza, my research mentor during my master's studies, introduced me to the field of operations research and patiently guided me through my first research project, even though at that time I had very little experience in optimization.

Next, I want to thank my fellow PhD students at Tepper. Anthony, thank you for being such a great friend. You have listened patiently to my complaints about research and life and have always been there with encouragement. As you often joked, I could always "lean on" you, and that turned out to be very true. Thank you also for organizing so many gatherings that brought people together; without you, our PhD community would probably have hung out much less often.

I am grateful to Aidin and Weizhong for sharing the early struggles of classes and qualifying exams. I would also like to thank Kyra, Su, Thomas, Violet, Yuyan, Özgün, Daniel de Roux, Mik, Heather, Neha, Neda, and Savannah for being role models when I first joined the program. The time we spent together, especially the many game nights in the lounge, created memories that I will always treasure. To Nilsu, Siyue, Vrishabh, Seba, Maca, Helia, Stephen, Daniel Yamin, Emily, Aurora, Milkis, Mogu, Mian, and Juan, thank you for the friendship and the moments we shared together, and I wish you all the best in your future endeavors.

I would also like to thank my friends outside of Tepper. Chengyue He, I am incredibly lucky to have you as a friend. Since we met and worked together at Columbia, we have gone through the PhD application process, the PhD years, and the faculty job market alongside each other. Your support throughout all of this has meant a great deal to me. Our shared enthusiasm for poker has also been an important part of our friendship, from discussing poker strategy to sharing poker stories to going on poker trips together. Those trips brought me some of the most joyful and memorable moments of these years.

Ziyu Huang and Jieqi Di, thank you for your friendship since our undergraduate days. I feel lucky that we ended up working in the same field, which gave us many opportunities to reconnect through conferences and academic life. I am also grateful that in New York and Atlanta I always knew I had a place to stay because of you.

Hao Yang and Junjie Xiao, thank you for being my friends since high school and for being present at important milestones in my life. Shuai Guo, you have been my closest friend for many years. The laughter we shared, the challenges we faced together, and the important moments we experienced along the way mean more to me than I can fully express.

Finally, I want to thank my family.

To my grandparents, thank you for the love and care you gave me during my childhood. Before this year we had not seen each other for five years, and I know that distance was not easy. Seeing you again this year meant a great deal to me, and I hope we will have many more chances to spend time together.

To Vicky Han, thank you for your unwavering support, your care, your understanding, and your love throughout this journey. Since we have been together, you have been my closest companion, and I have always felt that we face life as a team. The academic job market was a difficult and uncertain period, and whenever I felt lost or discouraged, you were always there for me. Every obstacle I faced, you treated as if it were your own. The moments we shared, both joyful and difficult, are among the most meaningful in my life. Without your support, I could not have reached this point.

To my parents, Xuewei Dong and Chengbin An: words cannot fully express my gratitude for your unconditional love. Even though I know you would have preferred me to stay closer to home, you supported my decision to pursue opportunities abroad because you believed it would lead to a better future for me. You have always been my strongest foundation, giving me the freedom and confidence to pursue what I truly wanted. Everything I have achieved rests on the support you have given me throughout my life.

For the past five years, whenever I received kindness, support, or encouragement, from close friends or even from people I met only briefly, I sometimes imagined the day I would write this acknowledgement. Now that it is finally written, it feels like the closing of an important chapter of my life. And I could not have reached this moment without all of you.

ABSTRACT

Modern operations increasingly rely on artificial intelligence (AI) models, ranging from demand forecasts and learned value functions to transformer architectures, to guide sequential decision-making under uncertainty. While these models provide rich representations of demand, preferences, and resource value, incorporating them into optimization raises challenges. Predictive accuracy may be unknown or time-varying, model outputs may be misspecified, and expressive architectures can induce computationally intractable decision problems. This thesis develops optimization for AI-driven sequential decision-making, focusing on algorithms that translate predictive structure into decisions equipped with robustness, computational efficiency, and provable performance guarantees.

The first part studies real-time personalization under transformer architectures. Transformers capture economically meaningful interaction effects such as variety seeking, substitution, and complementarity, but optimizing recommendations under general attention models leads to NP-hard combinatorial problems. We identify a structured class of simple transformers whose induced objectives admit exploitable low-rank structure, and we develop near-optimal recommendation algorithms with sublinear-time complexity under mild assumptions. These results demonstrate that expressive AI models can be reconciled with tractable and theoretically grounded optimization.

The second part considers online resource allocation under uncertain arrivals, where AI systems provide predictions in the form of shadow prices for resources. We design algorithms that leverage predicted dual variables without requiring knowledge of their accuracy. The resulting policies achieve a principled balance between robustness and consistency: they maintain worst-case guarantees while converging to near-optimal performance as prediction quality improves.

The third part examines inventory control under nonstationary demand. We propose adaptive policies for the nonstationary newsvendor problem that achieve order-optimal regret without assuming a bound on demand drift. Incorporating generic AI-generated demand predictions, we show how to combine realized observations with predictive signals to obtain guarantees that remain robust to arbitrary prediction error while improving when predictions are informative.

Across personalization, resource allocation, and inventory control, this thesis provides

a unifying framework for integrating AI models into sequential decision-making. By identifying structural properties that enable tractable optimization and by designing algorithms that remain robust to imperfect model guidance, the results establish theoretical foundations for AI-driven operations and revenue management.

TABLE OF CONTENTS

Acknowledgements	iii
Abstract	vii
Table of Contents	ix
List of Illustrations	xi
List of Tables	xiii
Chapter I: Introduction	1
1.1 Motivation: AI models and sequential decisions	1
1.2 Thesis themes: predictive structure, adaptivity, and algorithmic guarantees	2
1.3 Chapter overview and contributions	3
1.4 Connections across chapters	6
1.5 Organization of the thesis	7
Chapter II: Near-Optimal Personalization with Simple Transformers	8
2.1 Introduction	8
2.2 Model	15
2.3 Main Result	28
2.4 Algorithm and Proof of Main Theorem	31
2.5 Experimental Results	35
2.6 Conclusion	41
Chapter III: Online Resource Allocation with Predictions under Unknown Arrival Model	42
3.1 Introduction	42
3.2 Model: The Online Resource Allocation with Predictions	50
3.3 Main Results	61
3.4 Algorithm and Proof of Main Result	64
3.5 Experiments	67
3.6 Conclusion	70
Chapter IV: The Nonstationary Newsvendor with (and without) Predictions	72
4.1 Introduction	73
4.2 Model: The Nonstationary Newsvendor (without Predictions)	80
4.3 Solution to the Nonstationary Newsvendor (without Predictions)	87
4.4 The Nonstationary Newsvendor with Predictions	93
4.5 Experiments	97
4.6 Conclusion	104
Chapter V: Conclusion	105
Bibliography	108
Appendix A: Appendix for Chapter 2	122
A.1 Proofs in Section 2.2	122
A.2 Proofs in Section 2.2	125

A.3 Relationship of Model 2 and Choice Models	134
A.4 Graph Interpretation of Self-attention Layers.	134
A.5 Proofs in Section 2.2.	135
A.6 Proofs in Section 2.3.	149
A.7 Pseudo-code of Main Algorithm	150
A.8 Proofs in Section 2.4.	150
A.9 Proof of Proposition 8.	156
Appendix B: Appendix for chapter 3	184
B.1 Preliminary Results and Proofs in Section 3.2	184
B.2 Details in Section 3.3.1	190
B.3 Proofs in Section 3.3.3	191
B.4 Description of Algorithm 18	192
B.5 Proofs in Section 3.4.1 and 3.4.2	193
B.6 Proofs in Section 3.4.3	202
B.7 Experiment Details	210
Appendix C: Appendix for chapter 4	215
C.1 Preliminary Observations and Proofs	215
C.2 Proof of Lemma 3	217
C.3 Proof of Theorem 2	220
C.4 Proof of Proposition 10, Proposition 11, and Observation 3b)	225
C.5 General Variation	228
C.6 Proof of Proposition 12	230
C.7 Proof of Theorem 3	231
C.8 Unknown E and Known 0	235
C.9 Experiment Details	237
C.10 Additional Experimental Results	238

LIST OF ILLUSTRATIONS

Number		Page
2.1	Architecture of the simple transformer used in the Spotify experiment.	37
2.2	Architecture of the transformer used in the Trivago experiment. The simple transformer only contained the decoder, that is, a single self-attention layer.	38
2.3	Performances of four algorithms. The G-axis is the number of candidate solutions generated by each algorithm, and the H-axis is the objective value of the current best candidate solution. Each figure is averaged across 100 instances.	40
2.4	Scatter plots of candidate solutions. Each point corresponds to a pair of matched candidate solutions produced by our algorithm and Beam Search. The G-axis represents our algorithm's objective value and the H-axis represents the Beam Search candidate solution's objective value. Each plot contains 2500 data points given by 25 candidate solutions in each of the 100 instances.	41
3.1	(Figure and caption from [5]) Daily number of customers (in blue), from September 2014 to January 2015, at two different stores in the Rossmann drug store chain. Predictions (in red), starting November 2014, are generated using Exponential Smoothing with the same fitting process. The store in the upper sub-figure has substantially more accurate predictions ($\hat{r}^2 = 0.88$) than that of the lower sub-figure ($\hat{r}^2 = 0.11$).	44
3.2	Two potential arrival sequences for an online resource allocation problem with a single resource (two lemons) and two time periods. The left (right) sequence falls under the stochastic (adversarial) arrival model.	46
3.3	Histograms of GAPs with different forecasting methods, each containing 100 instances.	70

- 4.1 Daily number of customers (in blue), from September 2014 to January 2015, at two different stores in the Rossmann drug store chain. Predictions (in red), starting November 2014, are generated using Exponential Smoothing with the same fitting process. The store in the upper sub-figure has substantially more accurate predictions ($r^2 = 0.88$) than that of the lower sub-figure ($r^2 = 0.11$). 75
- 4.2 The costs of NO-PRED, PURE-PRED, PERP and OPT when (a) the variation parameter α is fixed, and (b) the accuracy parameter β is fixed. Each dot represents the cost of the corresponding policy on a given instance. 90
- 4.3 An example of a single time series from each dataset. In (c), the red dashed line represents an additional feature: daily online reservations. 100
- 4.4 Histograms of GAPs across (a) 1,000 randomly-sampled instances on the Rossmann dataset, (b) 2,880 instances on the Wikipedia dataset, and (c) 740 instances on the Restaurant dataset. 102
- 4.5 Histograms of the GAPs for (a) 173 Rossmann instances (left), 522 Wikipedia instances (middle), 476 Restaurant instances (right) for which PURE-PRED has lower cost, and (b) 827 Rossmann instances (left), 2358 Wikipedia instances (middle), 264 Restaurant instances (right) for which NO-PRED has lower cost. 103
- C.1 Histograms of GAPs divided by forecasting methods across the Rossmann dataset, the Wikipedia dataset, and the Restaurant dataset. 239

LIST OF TABLES

Number	Page
2.1	Average accuracy of different machine-learning models on Spotify and Trivago. 37
2.2	Average accuracy of general (full encoder decoder) transformers with various numbers of self-attention layers on Trivago. 38
3.1	Summary of our main theoretical results (in bold). Each entry has a corresponding algorithm that, without knowing the underlying arrival model, achieves the stated performance simultaneously for both stochastic arrivals and adversarial arrivals. Each entry also has a matching lower bound. 49
3.2	Summary of synthetic experiments. For each of three levels of prediction quality (the rows), and each of three generative arrival models (the columns), two summary statistics are reported over a random ensemble of instances: (left) the median GAP, and (right) the proportion of instances for which the GAP was at least 0.5. (Top) Results over all instances. (Bottom) Results over instances for which the rewards of the Mirror Descent and Prediction algorithms differ by at least 25%. 69
4.1	Summary of main theoretical results. Each entry has a corresponding policy that achieves the stated regret, along with a matching lower bound. 78
4.2	Continuation of Figure 4.1: costs incurred by an optimal policy which makes no use of predictions, a policy which relies entirely on predictions, and our policy. 78
4.3	Description of experimental instances. 101
4.4	Summary of experimental results. 103

Chapter 1

INTRODUCTION

1.1 Motivation: AI models and sequential decisions

A broad set of operations problems can be viewed as sequential decision-making under uncertainty: a decision-maker repeatedly chooses actions (e.g., order quantities, accept/reject and allocation decisions, or recommended items) in the presence of evolving demand, arrivals, and user behavior. In recent years, these decisions have been increasingly guided by AI models, ranging from demand forecasts and learned value signals to modern neural architectures such as transformers. This creates an opportunity and a challenge.

The opportunity is that AI models can capture rich patterns from large-scale data and thereby provide powerful guidance. The challenge is that using AI guidance within decision-making is rarely straightforward. First, predictive accuracy is often unknown a priori and can vary over time, so blindly following a model can be harmful (e.g., when forecasts drift, arrivals become nonstationary, or user behavior shifts). Second, even when a model is accurate, its expressive power can turn the downstream optimization problem into a hard combinatorial object. A concrete example arises in personalization: moving from embedding-based objectives (which permit extremely fast retrieval at scale) to transformer-based objectives breaks the synergy between learning and optimization, and the induced recommendation problem can become NP-hard and difficult even to approximate efficiently. This thesis develops optimization foundations for AI-driven sequential decision-making: how to transform AI guidance into decisions that are (i) computationally efficient and (ii) supported by end-to-end performance guarantees.

A unifying viewpoint. Across the three chapters, AI enters decision-making in three increasingly expressive forms: (i) as a learned objective for set/sequence recommendation (transformers); (ii) as predicted dual variables (shadow prices) that summarize future opportunity costs in online allocation; and (iii) as generic predictions of demand statistics in inventory control. In each case, the central methodological question is the same:

How can operations fully leverage the predictive and representational

power of modern AI models, while ensuring that the resulting decisions remain robust, tractable, and equipped with performance guarantees under sequential uncertainty?

1.2 Thesis themes: predictive structure, adaptivity, and algorithmic guarantees

This thesis develops a unified perspective on AI-driven sequential decision-making through three recurring principles.

(1) Exploitable predictive structure. Modern AI models provide rich predictive and representational signals, but these signals are useful for decision-making only when their induced structure can be identified and exploited algorithmically. Across the three chapters, AI guidance enters in different forms: as a transformer-based objective for recommendation, as predicted shadow prices in online resource allocation, and as demand forecasts in nonstationary inventory control.

In Chapter 2, the transformer model induces a combinatorial objective over recommended sets. While general transformer objectives render the optimization problem NP-hard and difficult even to approximate, the chapter identifies a structured subclass, simple transformers, whose low-rank attention structure enables near-optimal optimization in sublinear time. Thus, expressive neural models can be reconciled with tractable decision rules by isolating and leveraging their latent algebraic structure.

In Chapter 3, structure appears in the dual formulation of the allocation problem. Predictions are interpreted as estimates of optimal Lagrangian multipliers (shadow prices). This dual viewpoint provides a principled interface between AI-generated predictions and sequential allocation decisions.

In Chapter 4, structure arises in the temporal evolution of demand. Nonstationarity is parameterized via a variation budget, capturing the intrinsic difficulty of learning in changing environments. This structural characterization enables tight regret bounds that separate the effects of demand drift and stochastic noise.

(2) Adaptivity to unknown model quality. A central challenge in AI-driven decision-making is that prediction accuracy is typically unknown and may vary over time. Blindly following model outputs can therefore lead to poor performance. Chapters 3 and 4 develop algorithms that adapt to prediction quality without requiring prior knowledge of its accuracy.

In online resource allocation, predictions are given as shadow prices, and accuracy is measured as distance to the true optimal dual variables. The proposed algorithm achieves a tight consistency robustness tradeo : it preserves worst-case guarantees when predictions are inaccurate, while converging to near-optimal performance as prediction error diminishes. In the nonstationary newsvendor, generic demand forecasts are incorporated into ordering decisions without assuming any bound on their error. The resulting policy matches the best of learning-only and prediction-only benchmarks up to logarithmic factors, achieving order-optimal regret under both nonstationarity and prediction error.

(3) End-to-end algorithmic guarantees. A recurring tension in operations is that behaviorally realistic models can induce computationally intractable decision problems. This thesis emphasizes algorithms that translate AI guidance into implementable policies with formal guarantees.

For personalization, Chapter 2 develops near-optimal algorithms with sublinear decision complexity, reconciling transformer expressiveness with real-time requirements. For online allocation and inventory control, Chapters 3 and 4 provide tight competitive and regret guarantees under unknown arrival models and unknown nonstationarity, respectively. Across all settings, the results establish end-to-end guarantees that integrate statistical learning, algorithmic design, and sequential decision-making.

Taken together, these principles demonstrate how operations can harness the full power of modern AI: by identifying exploitable predictive structure, designing algorithms that adapt to model quality, and maintaining tractability and provable performance under uncertainty.

1.3 Chapter overview and contributions

Chapter 2: Personalization with simple transformers

Chapter 2 studies the problem of selecting a set (or sequence) of ~~items~~ ^{items} for a user when preferences are modeled by a transformer. Embedding models enable fast optimization via approximate nearest neighbor methods, but they fail to capture interaction effects between items in the recommended set. Transformers can represent such interactions, yet they induce a downstream optimization problem that becomes NP-hard in general, undermining real-time decision-making. The chapter therefore asks two questions: (i) is there a nontrivial subclass of transformers that

admits fast (ideally sublinear) optimization with meaningful guarantees, and (ii) does restricting to such a subclass sacrifice modeling power?

The chapter resolves this tension through the notion of simple transformers, architectures consisting of a single self-attention layer that strike a deliberate balance between modeling expressiveness and optimization tractability. On the modeling side, we show that simple transformers can represent key behavioral patterns studied in economics, marketing, and operations management, including sequential variety-seeking behavior and pairwise complementarity and substitution effects. These interaction patterns cannot be captured by pure embedding-based models, highlighting the necessity of moving beyond independent-item representations.

On the algorithmic side, we design a two-phase retrieval-and-ranking framework that approximately solves the induced recommendation problem. Under mild rank assumptions on the attention structure, the algorithm achieves near-optimal recommendations with expected amortized runtime sublinear in the catalog size. Sublinear complexity is essential in large-scale personalization systems, where decisions must be computed instantly upon user interaction.

We demonstrate the effectiveness of our approach through an empirical study on datasets from Spotify and Trivago. Our experiment results show that (1) simple transformers can model/predict user preferences substantially more accurately than non-transformer models and nearly as accurately as more complex transformers, and (2) our algorithm completes simple-transformer-based recommendation tasks quickly and effectively.

Chapter 2 is based on [6] joint work with Andrew A. Li, Vaisnavi Nemala, and Gabriel Visotsky.

Chapter 3: Online resource allocation with predicted shadow prices

Chapter 3 considers the Online Resource Allocation problem, a general model for sequential decisions with limited resources. At each time period, a request arrives with an action set; choosing an action yields reward and consumes resources, and the objective is to maximize total reward subject to a global resource budget. This model captures applications across revenue management, online retail, and online advertising.

The contribution of Chapter 3 is to incorporate predictions in a way that is robust to unknown prediction accuracy and unknown arrival model. Predictions are introduced in the dual space as estimated resource shadow prices, and prediction accuracy

is defined as a distance between predicted and true shadow prices. The chapter develops a tight lower bound describing the best achievable use of predictions, informally, an optimal follow when accurate, ignore when inaccurate behavior without knowing accuracy in advance, and proposes an algorithm that matches this bound. The algorithm is then empirically validated using large-scale real data from a retail setting.

Chapter 3 is based on [4] joint work with Andrew A. Li, Benjamin Moseley, and Gabriel Visotsky. In addition to establishing tight theoretical guarantees, the chapter includes a large-scale empirical evaluation using retail data from H&M, highlighting the practical impact of prediction-augmented allocation policies.

Chapter 4: Nonstationary newsvendor with (and without) predictions

Chapter 4 studies inventory control under nonstationary demand. The problem consists of a sequence of ordering decisions, where the underlying demand distribution can evolve over time. Nonstationarity is captured through a variation-based measure that quantifies the magnitude of demand drift, and performance is evaluated via regret relative to a dynamic benchmark that knows the demand distribution in each period.

The chapter first provides a complete characterization of the nonstationary newsvendor without predictions. It establishes matching lower and upper bounds and develops a policy that achieves order-optimal regret without prior knowledge of the degree of nonstationarity.

The chapter then incorporates generic predictions of mean demand observed before each ordering decision. No structural assumptions are imposed on how predictions are generated, and prediction quality is parameterized through cumulative prediction error. The central challenge is to leverage predictions optimally without knowing their accuracy in advance. To address this, the chapter proposes a prediction-robust policy that achieves near-optimal regret across regimes: it remains robust when predictions are inaccurate, and improves automatically when predictions are sufficiently informative, matching the better of learning-only and prediction-only strategies up to logarithmic factors.

Chapter 4 is based on [5] joint work with Andrew A. Li, Benjamin Moseley, and R. Ravi. Extensive numerical experiments, including evaluations on real-world demand data from Wikipedia and Rossmann drug store, validate the theoretical guarantees and demonstrate the practical benefits of combining adaptivity to nonstationarity

with robustness to prediction error.

1.4 Connections across chapters

This section highlights how the three chapters form a coherent progression.

From expressive AI objectives to tractable optimization. Chapter 2 begins with the most expressive AI component: transformer-based objectives for recommendation. It formalizes the computational barrier for general transformers and shows that by identifying the right structural subclass (simple transformers), one can obtain both modeling richness (capturing interaction effects) and tractable near-optimal optimization in sublinear time.

From tractable structure to robust use of AI guidance. Chapters 3 and 4 focus on a different role of AI: as predictive guidance for online decisions. A key lesson is that predictions are valuable precisely when algorithms can adapt to their quality without knowing it. Chapter 3 implements this through predicted dual variables (shadow prices) and proves tight guarantees under unknown arrival models and prediction accuracy. Chapter 4 implements an analogous principle in inventory control under nonstationary demand, combining learning from realized outcomes with prediction signals.

A unified perspective. Taken together, the chapters develop a unified perspective on AI-driven sequential decision-making. Across the different settings studied in this thesis, a common methodological pattern emerges:

1. Identify structural properties of AI models (or their outputs) that preserve practical modeling power;
2. Design algorithms that exploit these structures to achieve computational tractability;
3. Ensure robustness by adapting to unknown and potentially time-varying model quality;
4. Provide end-to-end guarantees (approximation and sublinear-time decisions in personalization; tight regret or competitive guarantees in online allocation and inventory).

1.5 Organization of the thesis

The remainder of the thesis is organized as follows. Chapter 2 studies real-time personalization under transformer objectives and develops near-optimal sublinear-time algorithms for simple transformers. Chapter 3 studies online resource allocation with predicted shadow prices under unknown arrival models and derives tight robustness consistency guarantees. Chapter 4 studies the nonstationary newsvendor with and without predictions, providing order-optimal regret guarantees and prediction-error robustness. Chapter 5 concludes with directions for future research on optimization foundations for AI-driven sequential decision-making.

Chapter 2

NEAR-OPTIMAL PERSONALIZATION WITH SIMPLE TRANSFORMERS

Real-time personalization has advanced significantly in recent years, with platforms utilizing machine learning models to predict user preferences based on rich behavioral data on each individual user. Traditional approaches usually rely on embedding-based machine learning models to capture user preferences, and then reduce the personal optimization task to nearest-neighbors, which can be performed extremely fast. However, these models struggle to capture complex user behaviors, which are essential for making accurate recommendations. Transformer-based models, on the other hand, are known for their practical ability to model sequential behaviors, and hence have been intensively used in personalization recently to overcome these limitations. However, optimizing recommendations under transformer-based models is challenging due to their complicated architectures. In this paper, we address this challenge by considering a specific class of transformers, showing its ability to represent complex user preferences, and developing efficient algorithms for real-time personalization.

We focus on a particular set of transformers, called simple transformers, which contain a single self-attention layer. We show that simple transformers are capable of capturing complex user preferences. We then develop an algorithm that enables fast optimization of recommendation tasks based on simple transformers. Our algorithm achieves near-optimal performance in sub-linear time. Finally, we demonstrate the effectiveness of our approach through an empirical study on datasets from Spotify and Trivago. Our experiment results show that (1) simple transformers can model/predict user preferences substantially more accurately than non-transformer models and nearly as accurately as more complex transformers, and (2) our algorithm completes simple-transformer-based recommendation tasks quickly and effectively.

2.1 Introduction

Personalization today is already immensely sophisticated. Media platforms, online retailers, and subscription services (just to name a few) capture rich data on their users in the form of their behavior and interactions with individual items/products. There are then two key ingredients: (1) this data is used offline to build machine-learning

(ML) based models of users' preferences, and then (2) these models are used in real-time to make personalized recommendations.

Zooming out from personalization for just a moment, the most jarring improvements in ML models over the last few years have been in generative models for language, and specifically transformer-based models (e.g. the T in ChatGPT) that have proven to be extremely accurate in modeling sequential data. Perhaps unsurprisingly, these same models are well-equipped for personalization. To x a concrete example, suppose an Instacart user is in the process of shopping online for groceries. This user's behavior consists of interactions with grocery items: browsing through items, viewing a subset of these in more detail, and adding a subset of these to their shopping cart. The task of learning this user's preferences essentially amounts to predicting their future interactions. The important observation here is that the user's behavior is naturally sequential, and so this prediction task is similar to completing a sentence, where the words are the items themselves. This connection to language suggests that the same transformer-based models may succeed in learning preferences.

This is already being done in practice, often with substantial empirical success (e.g. by Alibaba [50], Amazon [105], Spotify [135], and Wayfair [131]). However, as examples of such successes become increasingly common, there is little principled guidance on how transformers should be used for real-time personalization. This is the problem we seek to address.

Real-Time Personalization, Before Transformers: To make the nature of this problem more precise, it is worth reviewing how real-time personalization is performed without transformers. Referring to the two key ingredients mentioned at the outset: first, the ML-based models of user preferences are, by and large, pure embedding-based models. Using past data, each item is mapped to some element \mathbb{R}^3 in such a way that (a) similar items are close together, and perhaps more formally, (b) each user can be represented as some \mathbb{R}^3 so that the inner products $\langle \mathbf{u}, \mathbf{d} \rangle$ fully represent the preference/a nity of the user for each item \mathbf{d} .

These pure embedding models are not necessarily the most accurate models that can be estimated from past data, but they enable fast execution of the second ingredient, which is to optimize a set (or sequence) of items for each user in real time:

$$(2.1) \quad \begin{aligned} \max \quad & \sum_{j \in S} \langle \mathbf{u}, \mathbf{d}_j \rangle \\ \text{s.t.} \quad & |S| = k \end{aligned}$$

The (extremely general) formulation above simply highlights that real-time personalization consists of solving cardinality-constrained set-optimization (or sequence-optimization, whose equivalence to set-optimization we will discuss later on) problems where (a) the objective function is user-dependent, (b) the number of items potentially quite large, often in the hundreds of millions, and (c) a solution must be found in real-time, often just milliseconds. This is of course hopeless in general, but feasible when the objective function results from a pure embedding models. In particular, $f(D)$ typically takes one of two forms:

1. An additive function:

$$f(D) = \sum_{i \in D} v_i$$

for some non-decreasing function v_i . In this case, the optimal set can be characterized as follows: if the number of items such that $v_i > 0$ is no greater than k , then S consists of all such items. Otherwise, S is the set of items with the largest values of v_i .

2. A monotone submodular function (in the argument for any D), such that each item is fully encoded by v_i , so that a constant-factor approximation can be found using greedy-style algorithms along with (multiple queries to) a black box which computes the item with largest inner product to a given point in \mathbb{R}^3 [67].

In both of the above cases, the pure embedding model essentially reduces the real-time optimization task to one of nearest neighbor algorithms, where the goal is to find the items whose embeddings are most similar (under inner product similarity) to the user embedding. This task can be performed extremely fast, both theoretically (using approximate nearest neighbor algorithms with runtime sub-linear) and practically (given the nonstop engineering and improvement of commercial vector databases).

An Attempt to Introduce Transformers: Returning to transformers now, the natural opportunity is to improve the accuracy of the pure embedding models in representing user preferences: the embedding models essentially fail to capture the effects that items may have on each other when present in the same recommended set. Such effects are well-known to exist in multiple fields of study, as we will discuss shortly, and often representable via transformers. Unfortunately, the just-described synergy between the upstream user preference model estimated from data, and the downstream optimization of user-dependent sets of items completely breaks down

here. As we will see in the next section, this is true in a formal sense: Problem (2.1) is NP-hard, and likely hard to even approximate in linear time, if the objective D^0 is transformer-based. As of now, any practical implementation using transformers either applies a pure nearest-neighbor or greedy-style algorithm (essentially ignoring this hardness), or a random search heuristic (such as Beam Search).

In the spirit of developing a principled approach to transformer-based real-time personalization, this raises two major questions:

1. **Fast Optimization:** The formal hardness results just alluded to imply that achieving fast (ideally sub-linear in n) optimization of Problem (2.1) with any meaningful optimality guarantee is impossible if the objective D^0 allows for all transformers. Naturally then, is there a non-trivial sub-class of transformers for which this is possible?
2. **Modeling Power:** Assuming a positive answer to the first question, i.e. assuming the existence of a subset of transformers which enable fast optimization, does restricting to this subset come at a substantial cost in terms of modeling user preferences? Put another way, can this smaller sub-class of transformers achieve the same predictive accuracy as the family of all transformers?

Our Contributions

In short, our contributions provide concrete, theoretically-backed answers to both of the above questions:

1. **Modeling User Preferences with Simple Transformers.** We focus our study on a sub-class of transformers that we refer to as simple transformers. These are transformers which contain a single self-attention layer (to be defined in the next section), whereas transformers as a whole may contain multiple attention layers. Addressing the pair of questions above in reverse, we first formally show that simple transformers are able to represent two known, popular parametric models of user preference (Proposition 3):

- ^ Sequential variety effects in the context of marketing;
- ^ Pairwise complementarity and substitution effects in the context of economics.

It should also be emphasized that none of these models are representable via pure embedding models.

2. Real-Time Personalization with Simple Transformers. Our main result (Theorem 1) is an algorithm (Algorithm 9) which approximately solves Problem (2.1) in sub-linear time, when the objective function is given by a simple transformer:

Theorem 1 (Informal). Let OPT denote the optimal objective value of Problem (2.1) when the objective function is given by a simple transformer. Under additional rank assumptions on the simple transformer, for any $\epsilon \in (0, 1]$ and $\delta \in (0, 1]$, Algorithm 9 returns a solution with objective value at least $(1 - \epsilon)\text{OPT}$, and runs in expected amortized time

$$\mathcal{O}(n^{1-\epsilon} \log n)$$

where $n^{1-\epsilon} \log n$ and n^{δ} depend only on ϵ and δ . Here, $\mathcal{O}(\cdot)$ hides factors of order $\leq n^{\delta}$.

We will present and discuss the rank assumptions in Section 2.2. We also show that these assumptions are necessary for approximately solving Problem (2.1) in sub-linear time (Proposition 4 and Proposition 5), and that our algorithm's expected amortized runtime dependence on n and δ is optimal (Proposition 5). In particular, we give the following lower bounds:

Proposition 1 (Informal). Let OPT denote the optimal objective value of Problem (2.1) when the objective function is given by a simple transformer.

- (a) Without the rank assumptions of Theorem 1, there is no exact algorithm for Problem (2.1) running in time $\leq n^{1-\epsilon}$, assuming the Exponential Time Hypothesis.
- (b) Even under the rank assumptions of Theorem 1, Problem (2.1) admits no $(1 - \epsilon)\text{OPT}$ -approximation algorithm with runtime $\leq n^{1-\delta}$ for any function $\delta \in (0, 1]$, assuming the Exponential Time Hypothesis.

3. Empirical Study. We empirically validated the theoretical results of the previous contributions on two large datasets from Spotify [17] (which includes 1–000–000 playlists with 2–262–292 unique songs) and the travel website Trivago [11] (which includes user sessions of searching for hotel bookings, with around 730 unique

¹The Exponential Time Hypothesis (ETH) asserts that the 3Sat problem cannot be solved in sub-exponential time. For a Boolean formula in conjunctive normal form with exactly three literals per clause, the 3Sat problem requires deciding if there exists a truth assignment to the variables that satisfies all clauses, and finding one if so.

users and around 40–100 unique hotels recorded in around 100–1000 different sessions).

In support of the first contribution, our first set of experiments demonstrates that, given data on past user behavior, simple transformers can effectively model and predict user preferences. They achieve (a) substantially higher accuracy than non-attention-based models (such as pure embedding models), and (b) performance that is nearly comparable to more complex transformer architectures. Specifically, simple transformers achieved, on average, 14.1% higher accuracy than non-attention models (e.g., logistic regression, random forest, support vector machine), and only 2.5% lower accuracy than general transformers with multiple self-attention layers.

In support of the second contribution, our second set of experiments demonstrates that our algorithm performs simple-transformer-based personalized recommendation tasks both efficiently and effectively. We solved instances of Problem (2.1) using the simple transformers trained in the first set of experiments and compared our algorithm to two widely used benchmark methods: Nearest Neighbor (for retrieval) and Beam Search (for ranking). Under a fixed candidate solution budget (to partially standardize runtime), our algorithm achieved objective values that were, on average, 20.86% higher than those obtained using Nearest Neighbor and 20.56% higher than those obtained using Beam Search. Therefore our algorithm achieved strong empirical performance in both the retrieval and ranking phases.

Literature Review

Transformers in Recommender Systems. Since their introduction by [65], transformer architectures have become central to recommender systems due to their ability to model long-range dependencies in user item interaction sequences. They have been widely adopted in practice by companies such as Alibaba, Amazon [105], Spotify [135], and Wayfair [131]. Prior work has primarily focused on developing specialized transformer architectures, including sequential self-attention models [90, 131, 174], single-layer attention models [1, 44, 50, 171], and deeper or more complex variants [49, 69, 105, 112, 159, 177, 182].

In contrast, our work focuses on the optimization problem that arises once the architecture is fixed. We study single-attention-layer transformers, which are both prominent and empirically successful in practice, and design algorithms for fast, near-optimal recommendation under these models.

Representational Power of Transformers. Transformer architectures are built on self-attention, which provides strong representational power in theory and practice. On the theoretical side, Yun et al [179] and Wei et al [173] established universal approximation results, analogous to classical results for feedforward networks [80]. Pérez et al [144] and Wei et al [173] further showed that transformers are (approximately) Turing-complete. More recently, Sanford et al [152] presented both positive and negative results, identifying a sparse averaging task where transformers scale logarithmically with input size, unlike recurrent or feedforward networks which scale polynomially, as well as a triple detection task requiring linear attention scaling. Related limitations for induction heads were shown by Bietti et al [38], Elhage et al. [66], Sanford et al. [151].

In the context of recommender systems and choice modeling, Ko et al [102] showed that classical models such as Halo-MNL can be represented by a single attention layer, and Wang et al [170] proposed transformer architectures for learning a broad class of choice models. We complement this line of work by showing that a single attention layer can also capture user preference models with sequential variety effects and pairwise complementarity or substitution.

Approximate Nearest Neighbor. Approximate nearest neighbor (ANN) methods are central to real-time personalization, enabling efficient retrieval of relevant items from embeddings when exact nearest neighbor search is computationally prohibitive at scale, especially in high dimensions. Common approaches include hashing-based methods such as Locality-Sensitive Hashing (LSH) [86], graph-based traversal of proximity graphs [27], and tree-based methods such as KD-trees and Ball Trees, which are effective in low dimensions but degrade as dimensionality increases [136]. In our work, ANN enables fast item retrieval, allowing real-time execution.

Binary Quadratic Optimization. Under transformer architectures, recommendation can be formulated as a binary quadratic optimization (quadratic knapsack) problem, which is NP-hard, subsumes maximum clique and densest subgraph, and is NP-hard to approximate within any finite factor [47]. As a result, prior work focuses on tractable cases, including FPTAS results for series-parallel and bounded-treewidth graphs [147], and a PTAS for planar graphs [161].

In our setting, we exploit the low non-negative rank of the softmax matrix to obtain a PTAS. Related low-rank results include FPTAS algorithms for quasi-concave objectives [73], low-rank objectives over polytopes [14], and PTAS results for

binary nonlinear programs [40]. Our approach instead relies on low non-negative rank [53], and is also related to multi-objective knapsack problems [28, 29, 66].

2.2 Model

Real-Time Personalization without Transformers

Pure Embedding Models. Before introducing the transformer-based models that will be the focus of this paper, we begin with a brief overview of pure embedding models, which are widely used in modern personalization systems. These models represent users and items as vectors in a shared low-dimensional space, enabling efficient modeling of interactions via simple operations such as the inner product.

Formally, let $\mathcal{I} = \{1, \dots, n\}$ denote a set of items, and let $V \in \mathbb{R}^{n \times d}$ be a matrix whose i -th row $v_i \in \mathbb{R}^d$ is the value vector of item i . These value vectors are designed so that items with similar embeddings are likely to be perceived similarly by users. Each user is likewise represented by a user vector $u \in \mathbb{R}^d$. Both the item and user vectors are typically learned from historical interaction data through matrix factorization, collaborative filtering, or more complex models trained on click or engagement feedback.

The utility (or reward) of recommending item i to user u is modeled as $\phi(v_i \cdot u)$, where $\phi : \mathbb{R} \rightarrow \mathbb{R}$ is a non-decreasing reward function specific to item i . In many applications, the same function ϕ is used for all items. However, in some cases it is important to allow heterogeneous reward functions that reflect item- or category-specific behavior. For example, different product categories often exhibit different click-through-rate (CTR) saturation patterns: a small increase in relevance (e.g., $v_i \cdot u$) might sharply increase CTR for breaking news articles, whereas product ads might exhibit more gradual, linear gains, and fashion items may plateau early due to browsing behavior. Modeling such differences requires different shapes of the function ϕ , even if all are monotonic. To maintain full generality, we therefore allow each item to have its own reward function ϕ_i . Moreover, in many applications the image of ϕ_i is often $(0, 1]$, in which case $\phi_i(v_i \cdot u)$ can be interpreted as the probability of a positive user action, such as a click or purchase. In these cases, ϕ_i is sometimes chosen to resemble functions like the logistic function.

One of the key advantages of pure embedding models is the efficiency of real-time personalization. Given a user vector u and a budget, the goal is to select a subset

(k = 1/4 of at most n items that maximizes the total reward:

$$\begin{aligned} \text{(Pure Embedding)} \quad & \max_{S \subseteq \mathcal{I}} \sum_{i \in S} \langle \mathbf{e}_i, \mathbf{D} \rangle \\ \text{s.t.} \quad & |S| \leq k \end{aligned}$$

Since the objective function is additive and the items are treated independently, this problem can be solved exactly via nearest neighbor. In our setting, this corresponds to identifying the k items whose value vectors are most aligned with the user vector \mathbf{D} under inner product similarity. That is, for a query $\mathbf{D} \in \mathbb{R}^d$, the goal is to find the top- k items maximizing $\langle \mathbf{e}_i, \mathbf{D} \rangle$.

Nearest neighbor arises in many applications across recommendation, vision, and language, where one often needs to retrieve items similar to a given input based on some feature representation. A naive solution evaluates inner products $\langle \mathbf{e}_i, \mathbf{D} \rangle$, which takes $\mathcal{O}(n)$ time. In our case, we compute $\langle \mathbf{e}_i, \mathbf{D} \rangle$ for every $i \in \mathcal{I}$ and then select the top k items with the largest values. The optimal solution can be described as follows: if there are at most k items with positive reward (i.e., $\langle \mathbf{e}_i, \mathbf{D} \rangle > 0$), then S^* includes all of them. Otherwise, S^* consists of the items with the largest values of $\langle \mathbf{e}_i, \mathbf{D} \rangle$. This greedy procedure yields an exact solution in linear time with respect to the total number of items.

Algorithms: Approximate Nearest Neighbor. The greedy implementation of nearest neighbor becomes computationally prohibitive as the number of points grows. In most applications such as personalization tasks, the number of items is on the scale of millions, and an algorithm with runtime linear in n cannot be performed in real-time. To address this problem, the notion of approximate nearest neighbor (ANN) has been widely adopted. Rather than finding the exact nearest point, ANN aims to return a point whose distance to the query is within a factor of the true minimum. Formally, for an additive approximation factor $\epsilon > 0$, the objective of ϵ -ANN is to find any point \mathbf{e}_i , where $i \in \mathcal{I}$ that approximately maximizes the inner product similarity to the query \mathbf{D} :

$$\mathbf{e}_i \in \arg \max_{\mathcal{I}} \langle \mathbf{e}_i, \mathbf{D} \rangle \cdot \epsilon$$

This relaxation enables much more efficient data structures and algorithms, often sub-linear in n , which can be implemented in real-time.² Various techniques have

²There are also other versions of ANN that concerns multiplicative approximation factors instead of additive ones. These two notions are equivalent under some boundedness assumptions on the points.

been applied to n-ANN algorithms. We give a detailed overview in Appendix A.1.

Because companies aim to recommend a set of at most k items to a user in their personalization task, their objective is not merely to identify a single item that is attractive to the user, but rather to efficiently retrieve a set of items that are collectively among the most attractive to the user. Therefore, they consider the notion of n-Approximate-Nearest Neighbor, which returns a set of items whose inner product similarities are within an additive error ϵ compared to the top-true nearest items. This is formally defined below.

Definition 1 (n-Approximate-Nearest Neighbor Algorithm). An n-Approximate-Nearest Neighbor algorithm builds a data structure on any given set of points $\{x_1, \dots, x_n\} \subseteq \mathbb{R}^d$, and takes any query $q \in \mathbb{R}^d$, any $\epsilon > 0$, and any $k \geq 0$ as inputs. Let $\pi = (\pi_1, \dots, \pi_n)$ be a permutation of the indices such that $\langle x_{\pi_1}, q \rangle \geq \langle x_{\pi_2}, q \rangle \geq \dots \geq \langle x_{\pi_n}, q \rangle$. The oracle outputs indices $\{i_1, \dots, i_k\}$ such that $\langle x_{i_j}, q \rangle \geq \langle x_{\pi_{k+j}}, q \rangle - \epsilon$ for each $j = 1, \dots, k$ with expected amortized runtime $O(n \log n)$.

Many n-ANN algorithms can be naturally modified to Approximate-Nearest Neighbor algorithms, with a similar expected amortized runtime, that is, sub-linear in n . As an example, we give an Approximate-Nearest Neighbor algorithm in Appendix A.1.

Lemma 1. There exists an Approximate-Nearest Neighbor algorithm, which we give in Appendix A.1, with expected amortized runtime

$$O(n \log n + k \log \frac{1}{\epsilon})$$

Finally, companies apply any given Approximate-Nearest Neighbor algorithm to approximately solve Problem (Pure Embedding). This can be done by first partitioning the items according to their reward functions. For each partition, they apply an Approximate-Nearest Neighbor algorithm to identify items that are collectively among the most attractive to the user within that partition. They then evaluate the rewards of all such candidate items across partitions and select the top- k items with the highest overall reward.

To analyze the approximation error incurred by this procedure, we introduce the following parametrization of the reward functions. For $\epsilon > 0$ and $k \geq 0$, the function f satisfies:

$$f(x) \geq \epsilon \cdot \mathbb{1}_{\|x\| \leq 1} - \epsilon \cdot \mathbb{1}_{\|x\| > 1} \quad \text{for all } x \in \mathbb{R}^d$$

where $\phi(x) = 1 - e^{-x}$ and $\psi(x) = \max\{0, x\}$ are non-negative functions. This parametrization captures the idea that a small additive error in the input leads to a controlled multiplicative and additive error in the output.

This parametrization captures the behavior of many reward functions commonly used in practice. In particular, we make the following observation:

Observation 1. If $\phi(x)$ is non-decreasing and Lipschitz, then it satisfies the above condition with $\alpha = 0$ and $\beta = 1$. In particular, many standard reward functions (also referred to as activation functions in machine learning), such as logistic, ReLU, leaky ReLU, PReLU, tanh, and softplus, satisfy this condition with $\alpha = 0$ and $\beta = 1$.

There are other activation functions that are not Lipschitz, yet still satisfy the condition with $\alpha = 1$ and $\beta = 1$. Examples include Fractional Power ReLU, Log-ReLU, and Soft Root Function.

With this parametrization, we can now state our main result on applying ANN algorithms to real-time personalization with pure embedding models. We begin by introducing notations that will be used throughout the paper. For an optimization problem P with objective function f , let G_P denote an optimal solution and define $\text{OPT}_P = f(G_P)$ as its optimal value. For an algorithm ALG applied to P , let G_P^{ALG} be the solution returned by ALG , and set $\text{ALG}_P = f(G_P^{\text{ALG}})$ to be the corresponding objective value.

Proposition 2. Suppose we have an Approximate-Nearest Neighbor algorithm with expected amortized runtime $\frac{1}{g} \log_2 g$ and suppose f is concave in \mathbb{R}^d . Let g be the number of distinct functions among $\{f_i\}_{i=1}^m$. Given any $\epsilon > 0$, we give an algorithm ALG for solving Problem (Pure Embedding) that satisfies

$$\text{ALG}_{\text{Pure Embedding}} \leq (1 + \epsilon) \text{OPT}_{\text{Pure Embedding}} + \epsilon$$

with expected amortized runtime

$$g \log_2 g \log_2 \frac{1}{\epsilon}$$

The proof of Proposition 2 appears in Appendix A.1. Notice that many Approximate-Nearest Neighbor algorithms have sub-linear and concave in d and we have given an example in Lemma 1. Therefore Problem (Pure Embedding)

can be approximately solved in sub-linear time, which is practical to implement in real-time.

Modeling Limitations. Proposition 2 shows that pure embedding models can be optimized efficiently. However, pure embedding models have a fundamental limitation: they treat each item independently and fail to capture how the value of an item may depend on the context in which it is presented. That is, the reward associated with item i is fixed once D and E_3 are known, regardless of what other items are presented alongside it. In many personalization applications, this is an unrealistic assumption. For example, the value of a song recommendation may drop if a similar song is already in the playlist; or the likelihood a user clicks on a hotel result may depend on what other hotels appear in the same search result page and how they compare. These effects, often referred to as set effects, are not captured by pure embedding models. Concretely, we present three common parametric models used in personalization that cannot be represented by pure embedding models.

First, we consider a famous parametric model of sequential variety effects in sequences. The concept of variety/diversity has been examined extensively in the marketing literature (see e.g., [129, 148]). This model proposes that the perceived utility of an item depends not only on its intrinsic quality but also on its novelty relative to previously seen items. In particular, repeated exposure to similar items leads to diminishing marginal utility, while introducing diverse or contrasting items can restore or amplify engagement. This intuition aligns closely with the notion of discounted utility over sequences, where the utility derived from an item is multiplicatively reduced based on its similarity to past items (see, e.g., [25, 26]).

Such models capture behavioral tendencies like satiation, boredom, and the desire for exploration, and have been influential in both economic theory and practical recommendation systems. Below we give a mathematical formulation of sequential variety effects.

Model 1 (Sequential Variety Effects). Let $\mathcal{I} := \{1, \dots, n\}$ denote the set of items. Each item i has a base utility $u_i \geq 0$ and a similarity embedding $\mathbf{g}_i \in \mathbb{S}^{3-1}$. Define pairwise similarity score between item i and item j by $B_{ij} := \langle \mathbf{g}_i, \mathbf{g}_j \rangle \in [0, 1]$. Fix a sequence length $T \geq 1$ and nonnegative lag weights $\lambda_1, \dots, \lambda_{T-1} \geq 0$.

For a sequence $\mathbf{c} = (c_1, \dots, c_T)$ of length T , the sequential variety adjusted utility

of the item at position C is

$$u_c = D_c \exp \left(V \sum_{i=1}^{C-1} B_{c-i} \right) \quad C = 1, \dots, n$$

where $V \in \mathbb{R}$ controls the strength and sign of the variety effect.

Model 1 adjusts the context-free utility D_c multiplicatively according to the similarity between the current item and recently shown items, with older instances discounted by $_$. When $V > 0$ (preference for variety), similarity to recent items $B_{c-i} > 0$ decreases the current utility, while dissimilarity $B_{c-i} < 0$ increases it; when $V < 0$ (preference for continuity), the effect reverses and thematic similarity boosts utility. The weights $_$ specify the memory profile: choosing $_ = d^{-i}$ with $d \in (0, 1]$ yields exponential decay in influence with lag, whereas setting $_ = 0$ and $_ = 1$ for $i = 1$ recovers a one-step effect based only on the immediate predecessor. The exponential factor ensures $B_{c-i} > 0$ and provides a smooth, multiplicative adjustment driven by recent sequence context.

Model 1 also connects closely to the notion of Maximal Marginal Relevance (MMR) introduced by Carbonell and Goldstein [42], a classical approach in information retrieval that balances relevance and diversity. MMR selects an item by trading off its intrinsic relevance score against its similarity to previously selected items. Our model can be viewed as a smooth, parametric generalization of this idea. When $_ = 0$ and the lag weights concentrate on the most recent item (e.g., $_ = 0$ for $i = 1$), the adjustment term penalizes items that are highly similar to those already chosen, mirroring the diversity-inducing term in MMR. More generally, Model 1 replaces the linear trade-off in MMR with a multiplicative, exponentially weighted discount that naturally extends to sequential settings with memory and graded similarity effects. Thus, the model provides a behavioral micro-foundation for diversity-aware ranking rules widely used in marketing and recommendation systems.

Second, we consider a model of pairwise complementarity and substitution effects. In economics, these effects describe how the presence of one item influences the desirability of another: complementary items enhance each other's attractiveness, while substitutes reduce it. A common modeling approach represents items as vectors in a feature space, with pairwise interactions captured via inner products (see, e.g., [32, 109, 150]).

³By convention, the empty sum equals 0, so $B_{c-0} = D_c$.

Model 2 (Pairwise Complementarity and Substitution Effects). Let \mathcal{I} denote the set of items. Each item i has a value vector $\mathbf{v}_i \in \mathbb{R}^3$. The pairwise complementarity and substitution effects is parametrized by a matrix $\mathbf{R} \in \mathbb{R}^{|\mathcal{I}| \times |\mathcal{I}|}$ where $r_{ii} = 0$ for every $i \in \mathcal{I}$. For a user $u \in \mathcal{U}$ and a subset of items $S \subseteq \mathcal{I}$, the interaction-adjusted utility of item $i \in S$ is defined as

$$u_i = \mathbf{v}_i^\top \mathbf{u} + \sum_{j \in S} r_{ij} \mathbf{v}_j^\top \mathbf{u}$$

In Model 2, the value vector \mathbf{v}_i captures the intrinsic preference of the user for item i independent of any other items shown. The matrix \mathbf{R} encodes all pairwise interaction effects between items: a positive entry r_{ij} means that the presence of item j increases the utility of item i (complementarity), while a negative entry r_{ij} means that the presence of j reduces the utility of i (substitution). The diagonal entries are zero by definition, so an item does not directly influence its own utility. When $r_{ij} = 0$, item j has no effect on the utility of item i .

Model 2 is closely related to choice models that generate conversion probabilities. We discuss this connection in Appendix A.3.

Because in both models the reward associated with an item can depend on the presence or absence of other recommended items, we make the (informal) observation that these models cannot be represented by pure embedding models. In contrast, we will later show that both models can be naturally expressed within the framework we study in this paper.

Simple Transformers

In this section, we formally state the model which will be our primary object of study: simple transformers, or neural networks with a single self-attention layer. First, some preliminary definitions: the row-wise softmax operator, which we denote as $\text{softmax} : \mathbb{R}^3 \rightarrow \mathbb{R}^3$, is given by

$$\text{softmax}(\mathbf{x})_i = \frac{\exp(x_i)}{\sum_{j=1}^3 \exp(x_j)}$$

In particular, each row of softmax sums up to 1, and can be interpreted as a vector of weights. The notion of attention is fundamental to transformer-based models (e.g. [152, 165]). For input dimension d , output dimension d' , embedding dimension d_e , and matrices $\mathbf{W}_k \in \mathbb{R}^{d \times d_e}$, and $\mathbf{W}_v \in \mathbb{R}^{d \times d_e}$, a self-attention layer is a function $\text{SA} : \mathbb{R}^d \rightarrow \mathbb{R}^{d'}$ given by

$$\text{SA}(\mathbf{x}) = \text{softmax}(\mathbf{x} \mathbf{W}_k) \mathbf{W}_v \mathbf{x}$$

Here, the matrices Q , K , and V are often called the query, key, and value matrix, respectively.

To better understand the self-attention layer, it is natural to view it as a function on subsets of items. Let S denote a set of items. Then the self-attention layer is fully parameterized by an individual query, key, and value vector for each item (forming the rows of the respective matrices Q , K , and V). For any subset $S \subseteq [n]$ we define the set-membership matrix $X_S \in \{0, 1\}^{n \times |S|}$ by

$$X_S = \begin{bmatrix} x_{11} & \dots & x_{1|S|} \\ \vdots & & \vdots \\ x_{n1} & \dots & x_{n|S|} \end{bmatrix}$$

The function $\text{SA}_{Q, K, V}^S$ is then applied to X_S . The output is an $|S| \times |S|$ matrix in which each row $S \ni i$ is interpreted as follows:

- If $x_{i1} = 1$, then the i th row is a weighted average of the value vectors v_j where the weight on each v_j is given by the softmax-normalized dot product between the query vector q_i and the key vector k_j for $j \in S$.
- If $x_{i1} = 0$, then the i th row is the uniform average of the value vectors v_j . This is because item i 's query vector is zeroed out, that is, q_i is the zero vector, and thus assigns equal softmax weight to all items in S .

Transformers are a broad family of functions (also known as neural networks) constructed by repeatedly applying self-attention layers along with simple transformations that operate independently on each item. These point-wise transformations typically consist of linear mappings followed by non-linear functions known as activation functions, which introduce flexibility and allow the model to capture complex behaviors. In our context, these activation functions can be naturally interpreted as the reward functions r representing how the utility of each item responds to its input.

Simple transformers are a subclass of transformers with a single self-attention layer, followed by point-wise activation functions r . Formally, we define them as follows:

Definition 2 (Simple Transformer). For matrices $Q \in \mathbb{R}^{n \times d_Q}$ and $V \in \mathbb{R}^{n \times d_V}$, a vector $D \in \mathbb{R}^{d_V}$, and non-decreasing activation function $r: \mathbb{R} \rightarrow \mathbb{R}$, a

simple transformer is a function $T: \mathbb{R}^D \rightarrow \mathbb{R}^D$ given by

$$T(x) = \text{softmax} \left(\frac{1}{D} \sum_{j=1}^D \text{SA}_{i,j} x_j \right) \cdot D$$

As a side note, a simple transformer can also include multiple attention heads, in which case several self-attention layers are computed in parallel, each called an attention head, using different learned W , V , and U matrices. The outputs of all attention heads are then concatenated and passed through point-wise transformations. Equivalently, a simple transformer with multiple attention heads can be written as one with a single head by arranging each W , V , and U in block form. All of our results extend naturally to this setting. For clarity of notation, however, throughout this paper we focus on simple transformers with a single attention head, as defined above.

Modeling Power

Simple transformers are already widely used in personalization [14, 50, 171], though more sophisticated multi-layer transformers are also common. This raises a key question: to what extent can simple transformers effectively model user preferences?

To address this, we build intuition for how simple transformers operate in personalization. A simple transformer can be viewed as modeling interactions between a user and a recommended set of items, while also capturing pairwise interactions among items, often called *set effects*. Concretely, let I denote a set of items, and let $S \subseteq I$ be a subset recommended to a user. Let $R \in \mathbb{R}^{|S| \times |I|}$ be the matrix of value vectors, where each row r_i represents item i in isolation. When $(i, j) \in S \times S$, the self-attention output is simply $r_i \cdot r_j^T$ and the reward is $r_i \cdot r_j$.

When $j \notin S$, the self-attention layer transforms each r_i into a convex combination of r_j with weights determined by the softmax of inner products between the query vector r_i and key vectors r_j , i.e.,

$$\text{softmax} \left(\frac{1}{D} \sum_{j=1}^D \text{SA}_{i,j} x_j \right) \cdot D$$

Intuitively, $\text{SA}_{i,j}$ captures how other items influence item i , while r_i captures how item i influences others. The resulting vector $r_i \cdot r_j^T$ is projected onto D and

passed through σ , yielding reward

$$\sum_{i \in S} \mathbf{A}_{i, t}^\top \mathbf{D}_t^{-1} \mathbf{1} - \mathbf{1}^\top \mathbf{D}_t^{-1} \mathbf{1}$$

which corresponds to the 8-th coordinate of $\mathbf{T}_{i, t}^{-1} \mathbf{1} - \mathbf{1}^\top \mathbf{D}_t^{-1} \mathbf{1}$.

This formulation is permutation invariant, reflecting the unordered nature of Sequence models can be incorporated by augmenting and with positional encodings, which tag items with their positions and allow the model to distinguish otherwise identical items appearing in different positions. Thus, sequential structure is captured within the same framework as set-based interactions.

Having discussed how simple transformers operate in the personalization setting, we now examine their ability to model user preferences. While our later experiments will empirically demonstrate that restricting the architecture to a single attention layer results in only a modest reduction in modeling/predictive power, we begin by supporting this claim through an analysis of the two common parametric models used in personalization presented in the previous section. We have already seen that these two models cannot be represented by pure embedding models. On the contrary, below we show that both of them can be represented by simple transformers.

Proposition 3. The sequential variety effects in Model 1 and the complementarity and substitution effects in Model 2 can both be represented by a simple transformer.

The proof of Proposition 3 appears in Appendix A.2.

Finally, we provide a graph interpretation of self-attention layers in Appendix A.4.

Optimization

The primary purpose of this paper is to study the problem of personalizing a set (or sequence, equivalently) of items in real-time, where the underlying model of user preferences is given by a simple transformer. Following the setup and terminology in the previous subsections, let \mathcal{I} index a set of items, for which the query, key, and value matrices $\mathbf{Q} \in \mathbb{R}^{3 \times |\mathcal{I}|}$ and $\mathbf{V} \in \mathbb{R}^{3 \times E}$ are fixed in advance (these should be thought of as having been learned from prior data). Each user is represented by a vector $\mathbf{D} \in \mathbb{R}^E$ in the same space as the value vectors. For a given set $S \subseteq \mathcal{I}$, the simple transformer's output that corresponds to S is given by $\mathbf{T}_{S, t}^{-1} \mathbf{1} - \mathbf{1}^\top \mathbf{D}_t^{-1} \mathbf{1} = \sum_{i \in S} \mathbf{A}_{i, t}^\top \mathbf{D}_t^{-1} \mathbf{1} - \mathbf{1}^\top \mathbf{D}_t^{-1} \mathbf{1}$. Intuitively, this can be thought of as the reward obtained from recommending item i .

As a user arrives, our goal is to choose a set of at most k items that maximizes the total reward. Formally, the optimization problem we study is:

Definition 3 (Simple-Transformer Based Recommendations).

$$\begin{aligned}
 \text{(Main)} \quad & \max_{S \subseteq [n], |S| \leq k} \sum_{i \in S} r_i \\
 \text{s.t.} \quad & \sum_{i \in S} c_i \leq C
 \end{aligned}$$

Let OPT denote the optimal objective value of Problem (Main). Notice that if $k = 0$, that is, if nothing is recommended to the user, then the objective value of Problem (Main) equals to 0. Therefore $\text{OPT} \geq 0$. To ensure that Problem (Main) is meaningful, from now on we assume $\text{OPT} > 0$, since otherwise the best decision would be recommending nothing.

Hardness

As a starting point toward solving Problem (Main), observe that it can be solved exactly in $O(n^k)$ time via brute-force evaluation of all feasible solutions. As mentioned earlier, we are motivated by settings in which the number of items is extremely large (possibly hundreds of millions), and Problem (Main) must be solved in real time (possibly milliseconds). Thus, our goal will be to find an algorithm, potentially approximate rather than exact, whose runtime is sub-linear, i.e., $O(n^W)$ for some $W < 1$. Moreover, while the budget on the number of items to recommend, k , is typically moderate in practice (often around ten), exponential dependence on k would still be impractical to be implemented in real time. Therefore, our algorithm's runtime should also be polynomial in k .

Before proceeding, it is useful to present some initial hardness results to temper our expectations. We provide two sets of results. The first addresses the requirement of sub-linear runtime in n , with hardness parametrized by d , the dimension of the key and query vectors. The second addresses the requirement of polynomial runtime in k , with hardness parametrized by the non-negative rank of the matrix K , $\text{rank}_+(K) := \text{softmax}(K) > 0$, denoted as $\text{rank}_+(K)$.

Hardness Parametrized by $\text{rank}_+(K)$. We first present a proposition that reduces Problem (Main) to graph problems involving cliques, and then discuss its implications.

Proposition 4.

- (a) If $\beta = 1$ and $\alpha \geq 4$, then Problem(Main) subsumes the β -Clique problem⁴ on graphs with n vertices.
- (b) For any constant $\beta \geq 3$, there exists a number $\alpha(\beta) \geq 0$ for which the following holds. For any β such that $\exp(2^{1-\alpha}) \beta = 1$ and any $\alpha \geq 1$, Problem(Main) subsumes the problem of finding a largest clique in a graph with n vertices, $\exp(2^{1-\alpha}) \beta$ disjoint cliques, and all cliques have size at least α and at most $\alpha + 1$.

The proof of Proposition 4 appears in Appendix A.5. Proposition 4 implies concrete limitations on the theoretical results we can expect for solving Problem (Main):

- ^ By Proposition 4 (a), Problem(Main) inherits the hardness of the Clique problem, which is known to be NP-hard (where n is allowed to grow with β) [93]. Thus, we should not expect to find an exact algorithm which runs in $\beta^{1-\alpha}$ for some $\alpha \geq 0$ independent of β .
- ^ Moreover, it is known [48] that no exact algorithm can run in time $\beta^{1-\alpha}$ assuming Exponential Time Hypothesis holds. Thus, absent additional assumptions, we can only expect to approximately solve Problem(Main) in sub-linear time with respect to n .
- ^ One natural assumption to make is that β is small (this is typically the case in practice), and indeed our main result will be parameterized by β and only non-trivial when $\beta = \Omega(\log n)$. By Proposition 4 (b), if $\beta = \Omega(\log n)$, Problem(Main) is at least as hard as finding the largest clique in a graph with n vertices and β disjoint cliques. In particular, each clique in such a graph is itself a candidate maximum clique, and any algorithm must effectively search over β disjoint candidates to determine the largest, since there is no structural overlap between the cliques that an algorithm could exploit to narrow the search space. This indicates that an exact algorithm in sub-linear time with respect to n cannot exist when $\beta = \Omega(\log n)$.

Hardness Parametrized by β, α . Following on the above discussion, we require some additional assumptions to ensure that Problem(Main) can be solved, even approximately, in sub-linear time with respect to n . One such assumption

⁴For a (undirected, unweighted) graph, the Clique problem requires deciding if a clique of size α exists, and finding one if so.

will be that β is small. We do not state this as a formal assumption, but rather our main result will be parameterized by β .

Similarly, we also require some additional assumptions to ensure that Problem (Main) can be solved, even approximately, in polynomial time with respect to n . Our main result will be parameterized by a rank-type quantity pertaining to the matrix $A = \text{softmax}(B)$. Now β is by definition of rank β , and while the softmax operator does not preserve rank exactly, it is known that it can be well-approximated by a matrix with rank polynomial in β (see [76]). Thus, for example, if β is constant, then A is approximately low-rank.

Due to the softmax operator, the matrix is entry-wise non-negative. This allows us to consider a structural parameter known as the non-negative rank, denoted $\text{rank}_{+,0}$. The non-negative rank of a matrix $A \in \mathbb{R}_+^{n \times n}$ is defined as the smallest integer $r_{+,0}$ such that A can be written as a product of two non-negative matrices:

$$A = UV^T \text{ where } U, V \in \mathbb{R}_+^{n \times r_{+,0}}.$$

Equivalently, A can be expressed as the sum of $r_{+,0}$ non-negative rank-one matrices. Such a representation is called a non-negative factorization. Clearly, this notion is stronger than standard matrix rank; in particular, we always have $\text{rank}_{+,0} \geq \text{rank}$. For more properties on non-negative rank, see [59]. Non-negative rank has many applications in various fields, including data mining, combinatorial optimization, quantum mechanics, and more. In our case, it turns out that the non-negative rank of A is a key parameter that quantifies the hardness of Problem (Main). Formally, we present the following proposition:

Proposition 5. If $\text{rank}_{+,0} = 2$, Problem (Main) admits no $(1-\epsilon)$ -approximation scheme

- (a) with runtime $5^{11 \cdot n^\epsilon}$ for any function ϵ , assuming the N-Clique problem is not Fixed-Parameter Tractable (Corollary of Theorem 6 in [103]), and
- (b) with runtime $5^{11 \cdot n^\epsilon}$ for any function ϵ , assuming Exponential Time Hypothesis holds (Corollary of Theorem 5.1 in [87]).

⁵A problem is said to be Fixed-Parameter Tractable (FPT) if it can be solved in time $f(k)n^O(1)$ for some function f , where k is a chosen problem parameter. A widely held conjecture is that the N-Clique problem is not FPT.

For general $\text{rank}_k^{1,0}$, Problem(Main) admits no $(1 + \frac{1}{n})$ -approximation scheme with runtime

$$: > \frac{\text{rank}_k^{1,0}}{n \log^2 \text{rank}_k^{1,0} \cdot n^0} \quad \text{or} \quad : > \frac{1}{\text{rank}_k^{1,0}}$$

assuming Gap Exponential Time Hypothesis holds (Corollary of [60]).

The proof of Proposition 5 appears in Appendix A.5. Our proof is based on a reduction from Problem(Main) to the well-known Multi-dimensional Knapsack Problem (MDKP), formally defined in Appendix A.5. Proposition 5 shows that, when viewing $\frac{1}{n}$ as the parameter in Problem(Main), any algorithm achieving a $(1 + \frac{1}{n})$ -approximation must incur runtime with exponential dependence that necessarily involves both $\text{rank}_k^{1,0}$ and $1/n$ in a non-trivial way.

2.3 Main Result

From the discussions in the previous section, recall that while our algorithm's runtime is parametrized by the non-negative rank of achieving sub-linear dependence on and polynomial dependence or requires that the non-negative rank of be small. Importantly, our main result does not require itself to have low non-negative rank, but only that A can be well approximated entry-wise by a low non-negative rank matrix.

Formally, suppose there exists a non-negative matrix $A \in \mathbb{R}^{n \times m}$ such that

$$1 - \frac{\epsilon}{89} \leq \frac{A_{ij}}{W} \leq 1 + \frac{\epsilon}{89} \quad \text{for all } i, j$$

with $\text{rank}_k^{1,0} = A$.⁸ Then our main result is parametrized by ϵ and W . Just as in the case of the parameter $\frac{1}{n}$, our guarantee is non-trivial when $A \in \mathbb{R}^{n \times m}$.

Before stating our main result, we introduce some notation. For a vector G , let G_{\max} and G_{\min} denote the largest and smallest entry of G , respectively. Likewise, for a matrix A , let A_{\max} and A_{\min} denote its largest and smallest entry, respectively. If

⁶The Gap Exponential Time Hypothesis (Gap-ETH) asserts that, for some constant ϵ , distinguishing between satisfiable 3-Sat formulas and those that are not even ϵ -satisfiable requires exponential time.

⁷Computing a non-negative matrix factorization is in general NP-hard [160]; there is a rich literature on this subject that has yielded multiple algorithms (see [172] for surveys). We will not be concerned with this runtime in analyzing Problem(Main), as A and W are given beforehand, and thus we view this as amortized across multiple instances of Problem(Main).

⁸Actually, we only require a particular sub-matrix to be of non-negative rank k . This sub-matrix is of significantly smaller size, corresponding to the entries which survives a certain pruning procedure.

has rows G , we define

$$\|G\|_{2,1} = \max_k \sum_i |G_{ik}|$$

that is, the matrix norm induced by the vector 2- and 1-norms.

We are now prepared to state our main result, which is that our algorithm (to be described in the next section) achieves the following:

Theorem 1. Let T_{ANN} be the expected amortized runtime of an Approximate-Nearest Neighbor algorithm, which we assume is concave. Let g be the number of distinct functions among $\{f_1, \dots, f_g\}$. Suppose there exists $R \in \mathbb{R}^{n \times n}$ such that $\|W\|_{2,1} \leq R$ for all W , and R has non-negative rank k with an explicit non-negative factorization. Let ALG denote the objective value returned by Algorithm 9 for Problem (Main).

Given any $\epsilon > 0$, Algorithm 9 achieves

$$ALG \geq (1 - \epsilon) \text{OPT} - \epsilon n^2 \|W\|_{2,1}^2$$

Moreover, its expected amortized runtime is

$$T_{ANN} \leq \frac{1}{\epsilon} n^2 \|W\|_{2,1}^2 + \frac{1}{\epsilon} \log \frac{1}{\epsilon} n^2$$

where the Big-Oh depends only on k and $\max\{1, \frac{1}{\epsilon}\}$ and $\max\{1, \frac{1}{\epsilon}\} + D_{\max}^0$.

Some remarks will be useful for parsing and interpreting this result:

- ^ The concavity assumption on T_{ANN} is imposed for analytical convenience. It suffices that the runtime be sub-linear as any sub-linear function admits a concave sub-linear upper bound and hence can be handled within our analysis. Many Approximate-Nearest Neighbor algorithms satisfy these conditions; see Lemma 1 for one such example.
- ^ In practice, the number of distinct functions is typically very small, often just one.
- ^ In our algorithm's performance guarantee, the multiplicative dependence on n and the additive dependence on $\log \frac{1}{\epsilon}$ arise from the parametrization of ϵ .

$$\|G\|_{2,1} \leq R \text{ for all } G$$

The additional additive term comes from the use of the Approximate-Nearest Neighbor algorithm.

By Observation 1, if each \mathbf{B}_i is L -Lipschitz, our algorithm achieves

$$\text{ALG} \leq \text{OPT} + \epsilon n$$

for some universal constant $\epsilon > 0$.

If $\epsilon \geq \frac{1}{2} \log \frac{1}{\epsilon}$ and $A_i = \mathbf{B}_i^T \mathbf{B}_i$, the amortized runtime of our algorithm can be simplified to

$$\mathcal{O}(n \log \frac{1}{\epsilon})$$

Similar to how a small ϵ implies low (standard) rank approximation, it also implies that A_i can be approximated entry-wise by a matrix of low non-negative rank. In particular, we show that if the rows of \mathbf{B}_i and \mathbf{B}_j can be grouped into clusters, then such A_i can be constructed with $A_i = \mathbf{W}_i \mathbf{W}_i^T$ and a small ϵ .

Proposition 6. Let $\mathbf{A}_i = \text{softmax}(\mathbf{B}_i^T \mathbf{B}_i)$. Suppose $\mathcal{C}_1, \dots, \mathcal{C}_k$ form a partition of $[n]$ such that for every $i \in [k]$ and $j \in \mathcal{C}_i$, we have $\|\mathbf{B}_i(j) - \mathbf{B}_i(j')\| \leq \epsilon$ and $\|\mathbf{B}_i(j) - \mathbf{B}_i(j'')\| \leq \epsilon$. Then we can construct $A_i = \mathbf{W}_i \mathbf{W}_i^T$ such that

- $\|\mathbf{W}_i(j) - \mathbf{W}_i(j')\| \leq \epsilon$ for all $j, j' \in \mathcal{C}_i$, where $\mathbf{W}_i = \frac{1}{\sqrt{|\mathcal{C}_i|}} \sum_{j \in \mathcal{C}_i} \mathbf{B}_i(j)$
- and A_i has non-negative rank $\leq k$ with an explicit non-negative factorization.

To obtain the clusters, one approach is to directly partition the row spaces of \mathbf{B}_i and \mathbf{B}_j , and group the rows of \mathbf{B}_i and \mathbf{B}_j according to this partition. Since ϵ is assumed to be small, the number of clusters is also small. This yields the following corollary.

Corollary 1. Let $\mathbf{A}_i = \text{softmax}(\mathbf{B}_i^T \mathbf{B}_i)$. Given any $\epsilon > 0$, we can construct $A_i = \mathbf{W}_i \mathbf{W}_i^T$ such that $\|\mathbf{W}_i(j) - \mathbf{W}_i(j')\| \leq \epsilon$ for all $j, j' \in \mathcal{C}_i$ where $\mathbf{W}_i = \frac{1}{\sqrt{|\mathcal{C}_i|}} \sum_{j \in \mathcal{C}_i} \mathbf{B}_i(j)$, and A_i has non-negative rank

$$A_i = d \max_{j, j'} \|\mathbf{B}_i(j) - \mathbf{B}_i(j')\| \cdot \epsilon$$

with an explicit non-negative factorization.

The proofs of Proposition 6 and Corollary 1 appear in Appendix A.6. Notice that if $\epsilon \geq \frac{1}{2} \log \frac{1}{\epsilon}$, Corollary 1 shows that $A_i = \mathbf{B}_i^T \mathbf{B}_i$ for some small ϵ . In this case our main result is non-trivial.

- ^ In practice, the embedding dimension d of items is usually much smaller than the number of items n , so $d \ll n$ often holds. In addition under the condition that $A_i = \mathbb{R}^{d \times 1}$, for any fixed $n \geq 0$, our algorithm gives a $d \ll n$ approximation algorithm with expected amortized runtime that is sub-linear in the total number of items n , and polynomial in the number of items to recommend k .
- ^ Our algorithm's amortized runtime is practical. In reality companies and online platforms can easily have millions of products, so even algorithms with runtime polynomial in n would not be able to complete the personalized recommendation tasks in real-time. The number of items that companies can recommend to a single user is usually on the scale of tens (such as displaying products on a webpage), so algorithms with runtime exponential in k would not be able to complete the personalized recommendation tasks in real-time.

Real-time personalization algorithms typically operate in two phases: retrieval and ranking. In phase one (retrieval), the algorithm efficiently selects a small subset of promising candidate items from a large item pool. This is typically achieved using ANN algorithms over learned item embeddings. In phase two (ranking), an optimization problem, which is determined by the specific personalization model, is approximately solved over the retrieved subset to produce the final recommendations.

Our algorithm also follows the retrieval and ranking structure, but introduces novel methods in both the retrieval and ranking phases. The two pieces in our algorithm's expected amortized runtime directly correspond to the expected amortized runtime of our algorithm's two phases: the expected amortized runtime of phase one is

$$\mathcal{O}(n^{2.3} \log n) \text{ for ANN } \frac{d}{n} \ll 1$$

and the runtime of phase two is

$$\mathcal{O}(n^{2.3} \log n + W^A \log n)$$

2.4 Algorithm and Proof of Main Theorem

Before diving into the details of our algorithm, we rewrite Problem (Main) with slightly different notations. Let \mathbf{g}_g and \mathbf{z}_g denote the g th rows of \mathbf{G} and \mathbf{Z} , respectively. Let $\sigma = \text{softmax}(\mathbf{G} \mathbf{z}_g^T)$, and let F_g be the g th row of \mathbf{F} . For vectors $\mathbf{u}, \mathbf{v} \in \mathbb{R}^d$, let $\mathbf{u} \odot \mathbf{v} \in \mathbb{R}^d$ denote the element-wise product of \mathbf{u} and \mathbf{v} , i.e., $(\mathbf{u} \odot \mathbf{v})_i = u_i v_i$ for every $i \in [d]$. We make the following observation.

Observation 2. Problem (Main) is equivalent to the following Problem (P):

$$(P) \quad \max_{G \subseteq [n]} \sum_{i \in G} \frac{F_i + D_i \cdot G}{F_i} \\ \text{s.t. } |G| \leq k, \quad 1 \leq |G| \leq k :$$

where $\mathbf{1} \in \mathbb{R}^n$ is the all-ones vector, and $G_i = 1$ indicates that item i is selected into the set G .

The proof of Observation 2 appears in Appendix A.8, and the pseudo-code of our main algorithm appears in Appendix A.7. From this point onward, we will work with P, using its notation in place of Problem (Main).

Phase One (Retrieval)

In phase one, our algorithm aims to identify a small subset of items such that the optimal objective value of P does not decrease significantly when restricted to S . To achieve this, our algorithm first partitions items offline based on their query vectors q_i , key vectors k_i , and reward functions f_i . Since both q_i and k_i lie in a space of dimension d , the number of partitions is much smaller than items in the same partition are designed to behave similarly under the self-attention layer. That is, they produce similar outputs when interacting with other items. When a user arrives, the algorithm applies an Approximate-Nearest Neighbor algorithm (as defined in Section 2.2) within each partition to select at most k items whose value vectors have the highest approximate inner product similarity with u . Because items within the same partition respond similarly under attention, user preferences within each partition are primarily determined by similarity to the user vector. Thus, our algorithm retrieves a small subset containing high-reward items tailored to the user.

Proposition 7 (Phase One). Suppose we have an Approximate-Nearest Neighbor algorithm with expected amortized runtime $\text{ANN}^1 = \tilde{O}(n^2)$. Let g be the number of distinct functions among f_1, \dots, f_n . Given any $\epsilon > 0$, Algorithm 11 returns an index set $S \subseteq [n]$ such that

$$|S| \leq g \cdot \frac{140 \max_{i \in [n]} \|k_{i-1} - k_{i-2}\|_2^2}{1 + D_{\max} \cdot n} \cdot \epsilon^{-23}.$$

and the optimal value to the following Problem $P^{1-\epsilon}$

$$(P^{1-\epsilon}) \quad \max_{G \subseteq [n]} \sum_{i \in G} \frac{F_i + D_i \cdot G}{F_i} \\ \text{s.t. } |G| \leq k, \quad G_i = 0 \text{ for } i \in [n] \setminus S :$$

satisfies

$$OPT_{P^1} \leq (1 + \epsilon) OPT_P \leq (1 + \epsilon)^n$$

Moreover, suppose f is concave in x . Then the expected amortized runtime of Algorithm 11 is

$$\frac{140 \max_{k \leq k_{2-1}} \{k, k_{2-1} - k\} \cdot (k_{2-1} g^2 + D_{\max}^0)}{n} \cdot \frac{1}{g}$$

$$\leq \frac{1}{g} \cdot \frac{n}{35 \max_{k \leq k_{2-1}} \{k, k_{2-1} - k\} \cdot (k_{2-1} g^2 + D_{\max}^0)}$$

Suppose $\epsilon > 1/\log n$ and $k, k_{2-1} - k, k_{2-1} g^2 + D_{\max}^0$ are all viewed as constants, then for any given constant $\delta > 0$, we have $\frac{1}{g} = \delta^{1/\epsilon}$. Moreover, suppose $\frac{1}{g}$ is sub-linear in n when ϵ is fixed, then the expected amortized runtime of Algorithm 11 is also sub-linear in n .

Phase Two (Ranking)

In phase two, our algorithm approximately solves P , which is P over the retrieved subset of items S . Without loss of generality, assume $|S| \leq n/4$. By the first remark of Proposition 7, we may treat S as S^1 under mild assumptions. Then since $G_{ij} = 0$ for $i < j$, we may consider only the first entries of D and the top-left $\ell \times \ell$ principal sub-matrix of f . Therefore, with slight abuse of notation, we redefine $D \in \mathbb{R}^{\ell \times \ell}$ to include its first ℓ entries, and $f, R^{\ell \times \ell}$ to be the top-left $\ell \times \ell$ principal sub-matrices of the corresponding matrices, respectively. Moreover, note that the quantity

$$\frac{1^T f_{\ell} + D_{\ell}^0 \cdot G}{F_{\ell}^{\ell} \cdot G}$$

remains unchanged if the vector f_{ℓ} is multiplied by a non-zero constant. Thus, rescaling the rows of f does not change P^0 . Because $R^{\ell \times \ell}$ is now the $\ell \times \ell$ principal sub-matrix, the sum of its rows is not normalized to 1. So for simplicity of exposition, we assume each row of f is rescaled so that $\sum_{j=1}^{\ell} f_{ij} = 1$, and each row of $R^{\ell \times \ell}$ is rescaled accordingly so that $\sum_{j=1}^{\ell} R_{ij} = 1$, $W_{ij} = R_{ij} / f_{ij}$ for all i, j . Then we may rewrite P^1 as

$$(P^1)^0 \quad \max_{\{x_i\}_{i=1}^{\ell}} \sum_{i=1}^{\ell} x_i G^0 = \sum_{i=1}^{\ell} x_i \frac{1^T f_{\ell} + D_{\ell}^0 \cdot G}{F_{\ell}^{\ell} \cdot G}$$

$$\text{s.t. } \sum_{i=1}^{\ell} x_i = 1, \quad x_i \geq 0 \quad \forall i$$

From this point onward, we will work with this new form of P^1 .

Our algorithm begins by replacing with a low non-negative rank surrogate⁰ and showing that solving the problem under this approximation is sufficient. Rather than exhaustively enumerating all possible solutions, our algorithm then focuses on a restricted collection of partial solutions that retain the key structural information. The nonlinear terms in the objective are handled by introducing a family of auxiliary linearized problems, which can be further simplified through discretization. This reduction ensures that only a small number of auxiliary linearized problems need to be solved.

To address each auxiliary linearized problem, our algorithm employs a rounding procedure that converts fractional linear-programming solutions into valid discrete ones. At this stage, the central trade-off emerges: exploring too many partial solutions increases runtime beyond practical limits, while exploring too few places excessive burden on the rounding step, leading to higher approximation error. By carefully balancing this trade-off, the ranking phase achieves both computational efficiency through controlled exploration and strong accuracy by minimizing the loss introduced during rounding.

Proposition 8. Suppose there exists $R = R^0$ such that $W = W^0 + R$, $W^0 \succeq 0$, $R \succeq 0$, and R has non-negative rank k with an explicit non-negative factorization. Given any $\epsilon > 0$, Algorithm 16 achieves

$$\text{ALG}_{P^0} \leq (1 + \epsilon) \frac{W^0_{\max}}{W^0_{\min}} \text{OPT}_{P^0} + \epsilon \frac{W^0_{\max}}{W^0_{\min}} \frac{1}{\epsilon} \frac{W^0_{\max}}{W^0_{\min}}$$

where

$$\frac{W^0_{\max}}{W^0_{\min}} = \frac{\|W^0\|_F}{\lambda_{\min}(W^0)}$$

with runtime

$$\frac{4 \log_2 \frac{W^0_{\max}}{W^0_{\min}}}{\epsilon} \frac{1 + D^0_{\max}}{n} \frac{1 + D^0_{\min}}{n} \text{LP}$$

where $\frac{1 + D^0_{\max}}{n} = \frac{1 + \max_{i,j} |A_{ij}|}{n}$ and $\frac{1 + D^0_{\min}}{n} = \frac{1 + \max_{i,j} |A_{ij}|}{n}$. Here, LP = LP¹ < 3 < A, 2⁰, LP¹ < 2 < 2A, 2⁰ and LP¹ < = is the runtime of solving a linear program with < variables and = constraints.

Our proof of Proposition 8 appears in Appendix A.9.

2.5 Experimental Results

We conducted two sets of experiments with the following findings.

Representation. Simple transformers captured user preferences nearly as accurately as deeper transformer models. They achieved 14.1% higher accuracy than the best non-attention baselines (e.g., logistic regression, random forest, support vector machines), and only 2.5% lower accuracy than general multi-layer transformers, demonstrating strong representational power.

Optimization. Our two-phase algorithm (retrieval and ranking) efficiently optimizes simple-transformer-based recommendations. Compared against ~~Nearest Neighbors~~ Neighbors for retrieval and Beam Search for ranking under a fixed candidate budget, our algorithm achieved objective values 20.86% and 20.56% higher than the respective benchmark combinations, outperforming natural baselines in both phases.

We used two datasets. The first dataset was the Spotify Million Playlist Dataset [47]. Spotify is one of the largest music streaming platforms, with over 640 million monthly active users, including 252 million paying subscribers. The dataset comprises one million user-generated playlists created on the Spotify platform between January 2010 and October 2017. Each playlist includes features such as the playlist title and the titles of the tracks it contains.

The second dataset was the Trivago Session-based Hotel Recommendations Dataset [101]. Trivago is a global hotel search platform that operates 55 localized websites across more than 190 countries, providing access to over two million hotels. The dataset consists of user sessions related to hotel search and booking, encompassing approximately 730,000 unique users and 340,000 unique hotels across roughly 900,000 sessions. Each session includes information on user interactions with hotels, such as clicks and checkouts. In addition, the dataset contains various hotel attributes, including price, city, and other relevant features.

Representation

In this set of experiments, our goal was to show that simple transformers were able to learn from user behaviors and predict user preferences with high accuracy. More specifically:

Spotify. In the Spotify experiment, we designated playlists containing 20 songs as true playlists. To construct fake playlists, we took the first 15 songs from each true playlist and appended 5 randomly selected songs. The number of true and fake playlists was balanced to be equal. Given a playlist, the task was to classify it as either true or fake. The performance of an algorithm was evaluated based on the average classification accuracy.

Trivago. In the Trivago experiment, each session provided information on a user's interactions with the first 15 hotels. The task was to predict the user's interactions with the subsequent 5 hotels in the same session. The performance of an algorithm was evaluated based on the average prediction accuracy.

We compared three classes of machine-learning algorithms on these prediction tasks:

- ^ Non-Attention Models: This class included well-known machine learning algorithms that do not incorporate self-attention mechanisms, such as random guessing, logistic regression, support vector machines, and nearest neighbors. These models disregarded any potential sequential structure in the data.
- ^ Simple Transformers: This class consisted of transformer architectures with a single self-attention layer, followed by linear layers and activation functions.
- ^ General Transformers: This class contained more complex transformer architectures with multiple self-attention layers and potentially deeper network structures.

Below we give the architecture of the simple transformers used in both experiments.

Architecture used in Spotify. Let G_i denote the word2vec embedding of the i th song in a playlist, after being processed by a linear layer. Each user vector is modeled as the average of the embeddings of the first 15 songs in the playlist, that is, $\frac{1}{15} \sum_{i=1}^{15} G_i$. Let $SA_{\rightarrow}^{1 \times 1}$ be the self-attention layer with learned parameters θ , and let $\sigma^{1 \times 1} = \sigma^{1 \times 1}$ be an activation function composed of a linear transformation and a logistic function, both parameterized and learned during training. Let I be the set of indices corresponding to the 16th through 20th songs in the playlist. Then, the output of the simple transformer is given by

$$\tilde{G}_i = \sigma^{1 \times 1} \left(SA_{\rightarrow}^{1 \times 1} \left(\frac{1}{15} \sum_{j=1}^{15} G_j \right) \right)_{i \in I}$$

where \mathbf{e}_i denotes the input embeddings corresponding to songs in \mathcal{S} . General transformer models extend this architecture by incorporating additional self-attention layers.

Figure 2.1: Architecture of the simple transformer used in the Spotify experiment.

Architecture used in Trivago. Let \mathbf{e}_i denote the learned embedding of the i th hotel. Each user vector is modeled as the average of the embeddings of the first 15 hotels the user engaged within a given session, that is, $\mathbf{D} = \frac{1}{15} \sum_{i=1}^{15} \mathbf{e}_i$. Let $\text{SA}_{\rightarrow}^{(l)}$ denote a self-attention layer with learned parameters \mathbf{W}_s and \mathbf{W}_v , and let $\sigma^{(l)}$ be an activation function composed of a linear transformation and a logistic function, both parameterized and learned during training. Let \mathcal{I} be the set of indices corresponding to the 16th through 20th hotels in the session. Then, the output of the simple transformer is given by

$$\tilde{\mathbf{e}} = \sum_{i \in \mathcal{I}} \sigma^{(l)}(\mathbf{W}_s \mathbf{e}_i) \mathbf{W}_v \mathbf{e}_i + \mathbf{D}$$

where \mathbf{e}_i denotes the input embeddings corresponding to hotels in \mathcal{S} . The simple transformer architecture consisted solely of a decoder with one self-attention layer. General transformer models extended this architecture by incorporating an encoder, as well as additional self-attention layers in both the encoder and decoder.

The experimental results are presented in the tables below.

	Random Forest	Logistic Regression	Support Vector Machine	Simple Transformer	General Transformers
Spotify	0.518	0.520	0.334	0.702	0.726
Trivago	0.271	0.530	0.531	0.631	0.742

Table 2.1: Average accuracy of different machine-learning models on Spotify and Trivago.

In Table 2.1, the simple transformer outperformed non-attention models by average accuracy margins of 0.182 on Spotify and 0.20 on Trivago, while achieving performance comparable to general transformers. On Trivago, it outperformed several

Figure 2.2: Architecture of the transformer used in the Trivago experiment. The simple transformer only contained the decoder, that is, a single self-attention layer.

Dec. Layers \ Enc. Layers	1	2	4
1	0.590	0.602	0.596
2	0.654	0.692	0.700
4	0.724	0.742	0.675

Table 2.2: Average accuracy of general (full encoder decoder) transformers with various numbers of self-attention layers on Trivago.

more complex architectures and was only 2.4% below the average accuracy of all general transformers. Since the optimal architecture is not known a priori, simple transformers offer a strong and robust choice, delivering substantially higher accuracy than non-attention models while nearly matching general transformers.

Optimization

In the previous set of experiments, we demonstrated that simple transformers could empirically capture user preferences on both datasets. In this set of experiments, we turned to the task of personalized recommendation based on simple transformers. We treated the parameters θ , μ , and σ learned in the previous experiments as ground truth and solved Problem (Main). Each instance corresponded to an arriving user. More specifically:

Spotify. In the Spotify experiment, we were given 15 songs as input, and the task was to recommend an additional 5 songs to complete a 20-song playlist. The given 15 songs were treated as a representation of the user. The corresponding user vector D was computed by first averaging the word2vec embeddings of the 15 songs, and then applying a linear transformation to project the result into the value space. The

key, query, and value vectors of each song were obtained from the parameters learned in the previous experiments.

Trivago. In the Trivago experiment, we were given 15 user interactions as input, and the task was to recommend 5 additional hotels to maximize the booking rate. The user vector \mathbf{D} was computed by averaging the learned embeddings of the 15 hotels with which the user had interacted. The key, query, and value vectors of each hotel were similarly obtained from the parameters learned in the previous experiments.

Recall that our algorithm operates in two phases: retrieval and ranking. Our algorithm performs these two phases as described below:

Phase One (Retrieval): Our algorithm partitions the row space of the query matrix \mathbf{Q} and the key matrix \mathbf{K} into τ partitions, with the number of partitions treated as a tunable hyperparameter. Then, when a user vector \mathbf{D} arrives, for each group of items whose query and key vectors both belong to the same corresponding partition, we apply the τ -Nearest Neighbor algorithm and retain the τ items with the highest base reward $\mathbf{E}^T \mathbf{D}$. All retained items are then passed to the ranking phase.

Phase Two (Ranking): Our algorithm first enumerates all possible combinations of the top 2τ highest-reward items to include in a candidate solution, referred to as valid tuples in the proof of Proposition 8, where 2τ is a tunable hyperparameter. For the remaining items, the algorithm solves the residual subproblem by evaluating a fixed number of auxiliary problems, denoted as τ -PC in the same proof. The number of auxiliary problems solved was tuned to control the total number of candidate solutions generated.

We compared each phase of our algorithm against a natural benchmark:

τ -Nearest Neighbor (Retrieval): The τ -Nearest Neighbor algorithm served as the benchmark for the retrieval phase. It ignored any potential sequential effects and instead ranked items solely based on their base rewards, computed as $\mathbf{E}^T \mathbf{D}$. The algorithm then greedily selected the same number of items as our algorithm's retrieval phase to pass on to the ranking phase.

Beam Search (Ranking): Beam Search served as the benchmark for the ranking phase. Beam Search is a greedy-type heuristic commonly used in practice. Each candidate solution generated by Beam Search is specified by tuple (i_1, \dots, i_τ) .

To select the k -th item in the candidate solution, the algorithm evaluated each item not yet included by computing the incremental gain in the objective value if the item were added. It then selected the item corresponding to the highest increment and added it to the candidate solution. In particular, the tuple (k, \dots, k) corresponds to the fully greedy solution. The values k_1, \dots, k_n were tuned based on the desired number of candidate solutions.

Combining our retrieval and ranking phases with benchmark methods yields four algorithms, whose performance is shown in Figure 2.3. Our complete algorithm consistently outperformed all others across all candidate budgets.

Fixing the ranking phase (to either our Phase Two or Beam Search), our Phase One outperformed k -Nearest Neighbor by an average of 20.86%, demonstrating superior retrieval. Conversely, fixing the retrieval phase (to either our Phase One or k -Nearest Neighbor), our Phase Two outperformed Beam Search by 20.56%, highlighting its effectiveness in ranking. Overall, our algorithm achieved strong empirical performance in both phases, combining high efficiency with high accuracy.

(a) Spotify

(b) Trivago

Figure 2.3: Performances of four algorithms. The x -axis is the number of candidate solutions generated by each algorithm, and the y -axis is the objective value of the current best candidate solution. Each figure is averaged across 100 instances.

We further compare our Phase Two with Beam Search using scatter plots, fixing Phase One as the retrieval step. For ranking, both methods generated candidate solutions with the same computational budget: our Phase Two greedily solved a fixed number of auxiliary problems, while Beam Search branched the same number of times.

Each point in Figure 2.4 compares matched candidate solutions, with the x -axis showing our objective value and the y -axis that of Beam Search. Our Phase Two outperformed Beam Search in 90.92% of instances, with an average improvement of 29.01%, demonstrating substantial gains in the ranking phase.

(a) Spotify

(b) Trivago

Figure 2.4: Scatter plots of candidate solutions. Each point corresponds to a pair of matched candidate solutions produced by our algorithm and Beam Search. The X-axis represents our algorithm's objective value and the Y-axis represents the Beam Search candidate solution's objective value. Each plot contains 2500 data points given by 25 candidate solutions in each of the 100 instances.

2.6 Conclusion

In conclusion, we study real-time personalization using simple transformers, a single-self-attention-layer architecture that is both expressive and computationally efficient. We show that simple transformers capture key user preference structures, including sequential variety effects and pairwise complementarity and substitution, and develop a sub-linear-time algorithm achieving near-optimal performance. Empirically, simple transformers outperform non-transformer baselines, remain competitive with deeper transformers, and our algorithm improves objective values over standard methods such as ϵ -Nearest Neighbor and Beam Search.

Chapter 3

ONLINE RESOURCE ALLOCATION WITH PREDICTIONS UNDER UNKNOWN ARRIVAL MODEL

Online decision-makers often have access to predictions about future quantities such as arrivals, demand, or inventories, generated by methods ranging from simple time-series forecasts to modern machine-learning models. Because prediction accuracy is unknown a priori, blindly following such predictions can be harmful. We address this challenge by developing algorithms that leverage predictions in a manner robust to unknown prediction accuracy.

We study the Online Resource Allocation Problem, a general model for online decision-making in which limited resources must be allocated to a sequence of arrivals. Prior work has characterized optimal performance under both stochastic and adversarial arrivals, and designed algorithms that achieve these bounds without knowing the underlying arrival model. Building on this framework, we introduce predictions in the form of resource shadow prices and define prediction accuracy as the distance between predicted and true shadow prices.

We provide a tight lower bound characterizing how well any algorithm can exploit predictions by following them when accurate and ignoring them when inaccurate, without knowing either the prediction accuracy or the arrival model. We then propose an algorithm that matches this bound and empirically validate its performance using large-scale real data from the retailer H&M.

3.1 Introduction

Allocating a limited set of resources to satisfy different requests as they arrive is a key process in many operations problems. For example, airlines need to decide whether or not to accept a certain offer for a seat at a given price, while the total number of seats is limited [8, 160]; online retailers must choose which products to display to a browsing customer, taking into account inventory levels [121]; internet search engines auction off impressions to advertisers with limited budgets [30]. The Online Resource Allocation Problem is a generic model for all of these settings. In the problem, requests arrive sequentially, each request consisting of multiple actions to choose from, and each action generating some reward and consuming some subset

of resources. Actions are selected online, i.e. without knowing future requests. Resources are limited, and the objective is to maximize the total reward received across all time periods. While the Online Resource Allocation Problem is ubiquitous in practice, we briefly highlight several motivating examples:

- ^ Network Revenue Management: In airline revenue management, resources correspond to seats on future flights. Requests may involve multiple seats (e.g., group bookings or multi-leg itineraries) and have highly variable prices due to numerous fare classes.
- ^ Assortment Optimization: In online retail, a seller selects assortments to display during a customer's browsing session (e.g., search results or recommendations). Each display opportunity is a request, rewards are expected profits from customer choices, and resources correspond to product inventories.
- ^ Online Matching (AdWords, Online Auctions): Online matching models two-sided markets such as AdWords and auctions. Arrivals correspond to online nodes (e.g., impressions), while resources correspond to capacities of offline nodes (e.g., advertiser budgets).

Broadly, there are two main approaches to Online Resource Allocation. The traditional approach assumes a model for arrivals and designs algorithms with optimal worst-case guarantees, most commonly under either a stochastic model, where arrivals are drawn i.i.d. from an unknown distribution, or an adversarial model, where no assumptions are made. [22] shows that the best possible worst-case performance can be achieved simultaneously under both models without knowing the true arrival model. Interpreting stochastic and adversarial arrivals as stationary and nonstationary processes, respectively, their algorithm learns effectively in stationary settings while remaining robust to arbitrary nonstationarity. However, this notion of optimality is with respect to worst-case guarantees and may be overly pessimistic in practice.

The second, arguably more modern approach, is to utilize some sort of predictions on the future arrivals. Here we use the term prediction in the broadest possible sense, ranging from simple time-series forecasting models, to state-of-the-art machine learning algorithms based on large amounts of data, to human judgement, and even combinations of all of the above. The de facto approach in practice is to take these predictions as fact (in a way we will make formal momentarily). Naturally, the

performance of this approach relies heavily on the accuracy of the predictions, which is not guaranteed: Figure 3.1, taken from [5], shows this for the relatively simple task of forecasting daily visits to two stores.

Figure 3.1: (Figure and caption from [5]) Daily number of customers (in blue), from September 2014 to January 2015, at two different stores in the Rossmann drug store chain. Predictions (in red), starting November 2014, are generated using Exponential Smoothing with the same fitting process. The store in the upper sub-figure has substantially more accurate predictions ($\sigma^2 = 0.88$) than that of the lower sub-figure ($\sigma^2 = 0.11$).

To summarize so far, the Online Resource Allocation Problem admits algorithms with optimal worst-case guarantees (for both stochastic and adversarial arrival models, simultaneously), and these algorithms can be significantly better or worse than following predictions, depending on the prediction quality. This suggests the opportunity to design an algorithm that leverages predictions optimally, in the sense that the predictions are utilized when accurate, and ignored when inaccurate. Ideally, such an algorithm should operate without knowledge of (a) the accuracy of the predictions and (b) the method with which they are generated. This is precisely what we seek to accomplish in this paper.

Online Resource Allocation with Predictions

The primary purpose of this paper is to develop an algorithm that optimally incorporates predictions (defined in the most generic sense possible) into the Online

Resource Allocation Problem. Without predictions, the Online Resource Allocation Problem consists of a finite horizon of time periods and a limited number of types of resources. At each time period, a decision must be made which will consume a certain set of resources and yield a certain reward. The form of these individual decision problems changes over time and is unknown in advance.

Following [22], we consider both stochastic and adversarial arrival models. Under the stochastic model, performance is measured by regret, defined as the difference between the total reward of an optimal offline algorithm (which knows the entire arrival sequence) and that of the online algorithm. Our goal is to achieve sublinear regret ($\mathcal{O}(\sqrt{T})$), which ensures asymptotically optimal per-period performance. Under the adversarial model, sublinear regret is impossible in the worst case, and performance is instead measured by the competitive ratio, defined as the ratio between the optimal offline reward and the algorithm's reward; an α -competitive algorithm guarantees at least a $1/\alpha$ fraction of the offline optimum.

Without predictions, [10] showed that any algorithm incurs at least $\Omega(\sqrt{T})$ regret under stochastic arrivals, while [1] proved that under adversarial arrivals any algorithm has competitive ratio at least α , where α depends on simple problem parameters. These bounds are matched simultaneously by [2]. Our objective is to design algorithms that improve upon these worst-case guarantees when predictions are accurate, while retaining the same guarantees when predictions are inaccurate.

To that end, we introduce the notion of predictions. Our prediction is of the form of an κ -dimensional vector $\hat{\lambda}$ whose coordinates represent a predicted shadow price for each of the κ resources. In particular, $\hat{\lambda}$ may carry some information about the future (realized) arrival sequence. We will show that this form of prediction satisfies certain nice properties including that it (a) immediately translates to a decision policy, and (b) there always exists perfect predictions which achieve near-optimal reward.

We measure the quality of any prediction $\hat{\lambda}$ by its ℓ_1 distance to the closest perfect prediction λ^* .¹ Specifically, we use an accuracy parameter θ , defined as the largest θ such that $\|\hat{\lambda} - \lambda^*\|_1 \leq \theta \|\lambda^*\|_1$. Notice that when $\theta = 0$ the prediction is effectively useless, and as θ increases the prediction becomes more accurate. We call this problem Online Resource Allocation with Predictions. Our primary challenge will be to design algorithms with performances that are robust in the prediction quality without having access to θ .

¹The choice of the ℓ_1 norm follows naturally from our analysis, though any norm where $\|x\| \geq \|x\|_1 - 2\|x\|_2$ yields similar performance guarantees for our algorithms.

A Simple Example

Figure 3.2: Two potential arrival sequences for an online resource allocation problem with a single resource (two lemons) and two time periods. The left (right) sequence falls under the stochastic (adversarial) arrival model.

Before outlining our contributions, it is worth describing a simple example to illustrate the challenge we face in incorporating predictions of unknown accuracy. Consider the example in Figure 3.2, which depicts two potential arrival sequences for an online resource allocation problem with a single resource (two lemons) and two time periods. In both sequences, a single lemon may be sold for \$1 in the first time period. This same offer occurs in the second time period for the left sequence, but the right sequence offers \$2 for two lemons (this offer may not be split). Note that the left (right) sequence falls under the stochastic (adversarial) arrival model, and critically, an algorithm can not distinguish between the two sequences until the second time period. Still, a simple algorithm (accept all offers when feasible) achieves zero regret under the left sequence, and a competitive ratio of 2 under the right sequence (incidentally, $\frac{1}{2}$ for this problem instance).

However, suppose now we introduce a prediction, whose implication is that the first offer should be rejected. Under the right sequence, this constitutes a good prediction, and so an algorithm ideally would follow this prediction and collect the optimal \$2. Under the left sequence, this constitutes a bad prediction, and so an algorithm ideally would ignore this prediction but still achieve good regret as the arrival model is stochastic. It is of course impossible to do both of these.

More generally, there are essentially four worlds we must consider, depending on whether the arrival model is stochastic or adversarial, and whether the predictions are accurate or inaccurate. This example demonstrates that we can not hope to achieve the best of all four worlds simultaneously². Instead, we will find that just as the accuracy of the predictions is best characterized continuously between perfect and bad (via our accuracy parameter α), the arrival model is best characterized continuously along a carefully-defined interpolation between the stochastic and adversarial models.

Finally, we emphasize that under the stochastic arrival model, a prediction may legitimately carry information about the realized arrival sequence and thus be

²The language here is indeed in reference to the best of both worlds literature, e.g. [22].

statistically correlated with it. For example, suppose lemon demand depends on both predictable factors, such as weather or seasonality, and unpredictable idiosyncratic shocks, such as random household consumption variation. A prediction that accurately captures the predictable components, such as weather, will be correlated with the realized arrivals, even though it does not fully determine them. In this case, the remaining uncertainty arises only from the idiosyncratic shocks, which may have relatively low variance. This illustrates how informative predictions can naturally arise in stochastic environments and remain useful across all realizations, despite not being perfectly accurate.

Our Contributions

Our primary contributions can be summarized as follows.

1. A Nonstationary Arrival Model and a Lower Bound. We define a parameterized class of arrival models that interpolates between the stochastic and adversarial models. In particular, we define a precise measure of the stationarity of an arrival sequence (Definition 4), in terms of two values α and X , such that $\alpha = 0$ (loosely) corresponds to the stochastic model, and $\alpha = 1$ corresponds to the adversarial model (the two values are in general distinct, and have nice time-series interpretations in terms of trend and seasonality). Moreover, for any fixed arrival sequence, even X corresponds to some α for which the arrival sequence is α -stationary. Therefore X can be viewed as a parameter set by us. Later we will see that the choice of α and X affects our algorithm's performance and runtime. In particular, our algorithm needs to calculate the offline optimum $1/X$ times. Notably, this stationarity measure is defined for deterministic arrival sequences, and thus the corresponding (nonstationary) arrival models can be defined without positing a stochastic generative model.

The primary value of this new measure of stationarity is that it tightly characterizes the extent to which we can expect an algorithm to leverage predictions of unknown quality. Specifically, we prove the following lower bound:

Proposition 6 (Lower Bound, Informal). For any $0 \leq \alpha \leq 1$ and $0 < X \leq 1$, and any algorithm, at least one of the following holds:

- (1) Under the stochastic arrival model, the algorithm incurs α worst-case regret;

- (2) Under the adversarial arrival model, with $\lambda = X^0$ stationary arrivals, the algorithm's worst-case reward is at most

$$\frac{1}{U} \max_{\lambda} \text{OPT} - \text{PRD}(\lambda).$$

Here recall that U is the best competitive ratio (without predictions) which we will specify later. OPT denotes the optimal offline reward. Following the actions induced by the predictions also yields a certain amount of reward, which we denote PRD .

Now if $\lambda = X^0 = 1$, i.e. the adversarial model with no restrictions, then Proposition 7 implies that we can not simultaneously achieve sub-linear regret under the stochastic arrival model and a meaningful reward under the adversarial model (our lemon example was already evidence of this). However, for smaller values of X , we can hope for sub-linear regret and an adversarial reward close to the best value of $\max_{\lambda} \frac{1}{U} \text{OPT} - \text{PRD}$.

2. An Optimal Algorithm. We construct an algorithm that optimally leverages predictions, in the sense that it achieves the lower bound of Proposition 7, without knowing the underlying arrival model (stochastic or adversarial) and without knowing the prediction accuracy. In particular, our main theoretical result is the following:

Theorem 2 (Upper Bound, Informal). Given a prediction λ with (unknown) accuracy parameter ϵ and given $\lambda \in X$, there exists an algorithm such that, under mild (and tight) assumptions, both of the following hold:

- (1) Under the stochastic arrival model, the algorithm incurs $O(\frac{1}{\epsilon})$ worst-case regret;
- (2) Under the adversarial arrival model, with $\lambda = X^0$ stationary arrivals, the algorithm's worst-case reward is at least

$$\frac{1}{U} \max_{\lambda} \text{OPT} - \text{PRD}(\lambda).$$

Our theoretical results are summarized in bold in Table 3.1, with a comparison to the problem with no predictions and the problem with predictions of known accuracy.

³The $O(\cdot)$ notation hides logarithmic factors. Technically the regret should be $O(\frac{1}{\epsilon} \log(\frac{1}{\epsilon}))$ since if $\epsilon > \frac{1}{2}$ the regret bound should be a constant. For the simplicity of exposition we drop the obvious regret bound of $O(1)$ in the introduction section.

Arrival Model	Without Predictions	With Predictions of Known Accuracy	With Predictions of Unknown Accuracy
Stochastic (regret)	$\mathcal{O}(T^{1/2})$	$\mathcal{O}(T^{1/2})$	$\mathcal{O}(T^{1/2})$
Adversarial (reward)	$\frac{1}{T} \text{OPT}$	$\max_{\pi} \frac{1}{T} \text{OPT} - \text{PRDg}$	$\max_{\pi} \frac{1}{T} \text{OPT} - \text{PRDg}$

Table 3.1: Summary of our main theoretical results (in bold). Each entry has a corresponding algorithm that, without knowing the underlying arrival model, achieves the stated performance simultaneously for both stochastic arrivals and adversarial arrivals. Each entry also has a matching lower bound.

3. A Large-Scale Experiment. We evaluate the practical value of our model and algorithm using an H&M (Hennes & Mauritz AB) dataset containing two years of transactions for 105,542 products. The experiment corresponds to the assortment optimization problem and consists of three-month simulated horizons. We compare our algorithm against two natural baselines: the optimal algorithm without predictions and a policy that always follows the predictions. For each instance, these baselines define the worst- and best-achievable rewards, and we measure performance by the fraction of this gap captured by our algorithm; values close to 1 indicate near-best performance.

We generate predictions of varying quality using three popular forecasting methods. The average optimality gap achieved by our algorithm is 0.68 with Prophet forecasts, 0.58 with ARIMA forecasts, and 0.53 with Exponential Smoothing, demonstrating robust performance across prediction qualities.

Literature Review

Online Resource Allocation. Online resource allocation has been extensively studied under different arrival models. Under stochastic arrivals (i.i.d. or random order), the goal is to achieve sublinear regret, as in [57, 58, 68, 75, 97]. Under adversarial arrivals, sublinear regret is impossible, and performance is measured by competitive analysis; classic AdWords results include [130], which achieve a $\frac{1}{1+4^0}$ -competitive ratio.

Recent work studies algorithms that perform well under both models without knowing the arrival process. [33] achieved the optimal adversarial competitive ratio and improved stochastic performance for AdWords, while [22] proposed a mirror-descent algorithm achieving optimal stochastic regret and adversarial competitive ratio for Online Resource Allocation. Our algorithm also attains optimal performance under

both models, but differs fundamentally in design and analysis and is the first to robustly leverage predictions.

Our approach is based on gradient descent in the dual space, which has been widely applied in binary integer programming [14], online linear programming [13], bandits with knapsacks [15], online stochastic optimization [8], and online resource allocation [22]. We further draw on parameter-free optimization [43], using adaptive step sizes to exploit predictions without knowing their quality.

Algorithms with Predictions. Motivated by advances in data and machine learning, there is growing interest in augmenting online algorithms with predictions to go beyond worst-case guarantees. This paradigm has been applied to revenue optimization [20, 72, 137], caching [123, 149], online matching [9, 107], online scheduling [106, 146], the secretary problem [62, 63], and the nonstationary newsvendor [5]. Most of this literature uses competitive analysis and characterizes optimal consistency robustness trade-offs.

In contrast, under stochastic arrivals we perform regret analysis and show near-optimal worst-case regret without knowing prediction quality. Other regret-based analyses with predictions include [5, 81, 137].

The closest works are [20] and [72], which study strict special cases of Online Resource Allocation and treat prediction quality as binary. Under this assumption, they achieve optimal consistency robustness trade-offs. In contrast, we quantify prediction quality and provide tight guarantees for arbitrary accuracy.

3.2 Model: The Online Resource Allocation with Predictions

In this section, we first formally define the Online Resource Allocation with Predictions problem, and then describe two standard arrival models (stochastic and adversarial) as well as their respective performance metrics.

Problem Formulation

Online Resource Allocation. Consider a problem over time periods labeled $C = 1, \dots, T$. Assume there are d different types of resources. The total number of resources available is denoted by d , where $d \in \mathbb{R}^d_+$ is a non-negative d -dimensional vector. At each time period t , the decision-maker receives an arrival $W_t = (A_t, c_t) \in \mathcal{S}$. Here, $A_t: X_t \rightarrow \mathbb{R}_+$ is a non-negative reward function, $c_t: X_t \rightarrow \mathbb{R}^d_+$ is a non-negative resource consumption function, $X_t \subseteq \mathbb{R}^d$ is a

compact action space, and \mathcal{S} denotes the set of all possible arrivals. Note that we impose no convexity assumptions: A^c can be non-concave, G^c can be non-convex, and X_C can be non-convex or discrete. At each arrival W , without knowing any of the future arrivals, an action $G \in X_C$ must be selected, which yields $A^c(G)$ reward and consumes G^c resources. The objective is to maximize the total reward subject to the resource constraint. Finally, we assume that X_C always contains a 0 vector representing a void action that consumes no resources and yields no rewards: $A^c(0) = 0$ and $G^c(0) = 0$. This ensures that there is always a feasible action available. This is the problem we will refer to as Online Resource Allocation (without predictions).

We introduce some notations that will appear in our results later on (though our algorithms will not depend on these parameters). We denote by $\underline{d}_j = \min_{G \in X_C} G^c_j$ the lowest resource parameter and $\bar{d}_j = \max_{G \in X_C} G^c_j$ the highest resource parameter. Similarly, let A^c_0 be a constant which satisfies $A^c_0 \leq A^c(G)$ for every $G \in X_C$ (and let $\underline{d}_j \leq 0$ be constants satisfying $\underline{d}_j \leq G^c_j$ for every $G \in X_C$ and $G < 0$).

Primal and Dual. For any arrival sequence $W = \{W^1, \dots, W^T\}$ we use $\text{OPT}^1(W)$ to denote the online/hindsight optimum, which is the reward of the optimal solution when W is known in advance:

$$(3.1) \quad \text{OPT}^1(W) := \max_{G \in X_C} \sum_{C=1}^C A^c(G) \quad \text{s.t.} \quad \sum_{C=1}^C G^c \leq d$$

As we will describe momentarily, it will be natural to consider predictions in terms of the dual space, so the Lagrangian dual problem of Eq. 3.1 plays a key role. Let $\lambda \in \mathbb{R}^C_+$ be the vector of dual variables, where each λ_j can be thought of as the shadow price of resource j . We define

$$(3.2) \quad A^c_0(\lambda) := \sup_{G \in X_C} A^c(G) - \sum_{j=1}^C \lambda_j G^c_j$$

as the optimal opportunity-cost-adjusted reward of request c where the opportunity cost is calculated according to the shadow prices λ . Note that $A^c_0(\lambda)$ is a generalization

⁴We assume throughout the paper that the reward functions and the resource consumption functions are deterministic for any given action, but our algorithms also apply when the rewards and/or consumed resources are random.

⁵We will assume that the primal optimization problems in Eq. 3.2 admit an optimal solution. This is to simplify the exposition: our results still holds if we have an approximate of the optimal solution (see [22]).

of the convex conjugate of G^0 that takes the resource consumption function G^0 and the action space X_C into account. In particular, when $G^0 = G$ and X_C is the whole space $A^{1,0}$ becomes the standard convex conjugate. For fixed arrival rate define the Lagrangian dual function $\lambda \mapsto W^0 : \mathbb{R}^+ \rightarrow \mathbb{R}$ to be

$$(3.3) \quad \lambda \mapsto W^0 := \max_{C=1}^{\infty} A_C^{1,0}(\lambda, \gamma > d) \bullet$$

This allows us to move the constraints of Eq. 3.1 to the objective, which is easier to work with. By weak duality, $\lambda \mapsto W^0$ provides an upper bound on $\text{OPT}^1 W^0$. Moreover, even without any convexity assumptions, the duality gap of our problem is upper bounded by a constant that is independent from the time horizon. This can be shown via Shapley-Folkman Theorem (see Proposition 5.24 for a detailed explanation). We formally state and prove the weak duality result and the duality gap result in Appendix B.1. We equip the primal space of the resource constraints with the ℓ_1 norm $\|j\|_1$, and the Lagrangian dual space with the norm $\|j\|_1$. Such choices of norms come naturally from our analysis. Similar performance guarantees of our algorithms with the dependence on the number of resources can be obtained using the ℓ_2 norm for the primal space and the ℓ_1 norm for the primal space with $1 \leq \gamma, \lambda \leq 1$ and $\gamma \geq 2 - 1/\lambda$.

Predictions. So far, we have presented the problem of Online Resource Allocation without predictions. As described in the introduction, it is often the case that when this problem is faced in practice, some notion of a prediction can be made which might guide us in selecting actions. Such predictions can come from a diverse set of sources ranging from simple human judgment, to forecasting algorithms built on previous demand data, to more-sophisticated machine learning algorithms trained on feature information. The process of sourcing or constructing such predictions is orthogonal to our work. Instead, we treat these predictions as given to us endogenously (and in particular, we make no assumption on the accuracy of these predictions), and attempt to use these predictions optimally.

Notice that from Equation (3.2), at each time given a dual variable $\lambda \in \mathbb{R}^+$ there is a natural action to take, namely the action

$$C \in \arg \max_{C \in \mathcal{C}} \sum_{B=1}^C \lambda_B G_B^0(\lambda, \gamma > d) - \lambda C G^0(\lambda, \gamma > d)$$

⁶The number of resources C is viewed as a constant throughout the paper.

⁷We again assume this optimization problem and other similar-style optimization problems in this paper admit an optimal solution to simplify the exposition. When the right hand side contains more than one action, we naturally choose the action that has the highest reward.

In words, G_c is the greedy action that, subject to the resource constraint, maximizes the opportunity-cost-adjusted reward according to the shadow prices. Therefore each dual variable essentially corresponds to an algorithm, which is simply taking the greedy action G_c at each time period. Below we formally define this algorithm for any dual variable λ , which we call the Dual-Adjusted Greedy Algorithm (GRD):

Algorithm 1: Dual-Adjusted Greedy Algorithm GRD

Inputs: Dual variable λ , total time periods T , initial resources $r_1 = d$)

for $C = 1 \dots T$ do

 Receive request (c, X^c)

 Make the primal decision G_c and update the remaining resources

$$G_c = \arg \max_{c,1} \sum_{c \in C} \lambda^c X^c - \sum_{c \in C} \lambda^c G_c$$

end

Let $\text{GRD}(\lambda, j, W^0)$ denote the reward obtained by GRD with arrival sequence W and dual variable λ .⁸ For an arrival sequence W we say a dual variable λ is a perfect dual variable with respect to M , on arrival sequence W the reward obtained by GRD(λ, W) is at most an additive constant away from $\text{OPT}(W)$ (hence essentially optimal). To simplify notations, from now on we will use λ in place of λ^W . It can be shown that there always exists a perfect dual variable:

Proposition 1 (Perfect Dual Variable). For any arrival sequence W ,

$$\max_{\lambda \in \mathbb{R}^C} \text{GRD}(\lambda, j, W^0) \leq \text{OPT}(W^0) + \sum_{c \in C} \lambda^c$$

The proof of Proposition 1 appears in Appendix B.1, and utilizes the Shapley-Folkman Theorem. In words, Proposition 1 shows that there exists a dual variable λ that is essentially optimal to follow.

With the understanding of the key role that dual variables play in our problem, we formally introduce the notion of predictions. Because dual variables induce actions, they are natural quantities to predict. For an arrival sequence W we assume that before the first time period, the decision-maker receives a prediction $\lambda^W \in \mathbb{R}^C$ of the dual variable. Similar as λ , to simplify notations, from now on we will use λ in place of λ^W . We measure the prediction error of λ by $\|\lambda - \lambda^W\|_1$, its L_1 distance

⁸Later we will formally extend the notion of $\text{GRD}(\lambda, j, W^0)$ to any algorithm ALG.

from λ .⁹ We quantify the prediction error using the accuracy parameter, which is the smallest $\epsilon \in [0, 1/4]$ such that

$$|j_j^\lambda - j_{j_1}^\lambda| \leq \epsilon.$$

Here $\lambda \in [0, 1]$ is a scaling constant that we can choose, and that ensures $\epsilon \in [0, 1/4]$ can be chosen without affecting our performance bound asymptotically.¹⁰ Note that λ and ϵ both depend on n , so the definition makes sense even if λ is itself random. The two extreme cases of the prediction error are $\epsilon = 0$, in which case λ is almost a constant away from λ_1 , so the prediction is effectively useless; and $\epsilon = 1/4$, in which case $\lambda = \lambda_1$, so the prediction is perfect. We will always assume λ is unknown to the decision-maker.

In reality, a prediction λ is unlikely to be completely useless. We make the following technical assumption on the prediction quality:

Assumption 1 (Non-trivial Prediction). There exists a (known) function $\lambda_1(n) \in [0, 1/4]$ such that $|j_j^\lambda - j_{j_1}^\lambda| \leq \lambda_1(n)$.

Note that Assumption 1 does not eliminate the case $\lambda = 0$. In practice $\lambda_1(n)$ can be chosen to be a function close to 1 without hurting the algorithm's performance guarantee.

Aside: Predicting the Dual Variables. In many applications, historical data (e.g., arrivals from a previous month) is available and can be used to warm start dual variables, computed online without knowing the arrival order and interpreted as predictions (shadow prices) for resource values. While imperfect, such data often provides informative signals.

Predicting dual variables is a long-standing idea in both theory and practice. For example, in online matching under random order, [57, 167] compute dual variables from an initial prefix of arrivals and use them to obtain an approximate matching for the remainder, an approach once used by Yahoo! for ad matching. More broadly,

⁹In the case that multiple perfect dual variables exist, we take to be the perfect dual variable that is closest to λ .

¹⁰As a technical aside, there is a natural choice: Proposition 2 in [22] implies that it is enough to only consider dual variables that lie in the d -dimensional rectangle $[0, \frac{\max_j a_j}{1-\epsilon}] \times [0, \frac{\max_j a_j}{2\epsilon}] \times \mathbb{R}^{d-2}$ where $\epsilon = \frac{\max_j a_j}{A \cdot d}$, $\epsilon \in [0, 1]$. Therefore we may assume without loss of generality that the prediction we receive lies inside this rectangle (otherwise we could project to this rectangle). Thus setting $\lambda = \frac{\max_j a_j}{A \cdot d}$ ensures $\epsilon \in [0, 1/4]$.

dual prediction is common in dual-based algorithms for online scheduling [100], online matching [59, 107], online covering [23], and ad auctions [98].

In these problems, including ours, the online optimum depends only on the dual of the sample path, which lies in \mathbb{R}^n , rather than the full arrival sequence, whose dimension grows with the horizon. As a result, predicting the dual variable is often easier than predicting the entire sample path.

By contrast, under stochastic arrivals, prior work assumes predictions of per-period arrival distributions [19, 88], measured via total variation and Wasserstein distances. Such predictions are independent of realized arrivals and thus encode less information about the sample path. Indeed, even perfect distributional predictions in the sense of [19, 88] can still lead to $\Omega(\sqrt{P})$ regret due to inherent variance. In contrast, we allow predicted dual variables to depend on future arrivals, potentially encoding richer information through statistical association, for example via observable intermediary variables.

Arrival Models and Performance Metrics

An online algorithm ALG , at each time period t , takes an action G_t (potentially randomized, but deterministic here to save on notation) based on the prediction of the current request $A_t = G_t \cdot X_t^0$ and the previous history $H_{t-1} := \{A_{t-1}, G_{t-1}, X_{t-1}^0, \dots, G_1, X_1^0\}$, i.e., $G_t = \text{ALG}(A_t, G_t, X_t^0, H_{t-1})$. We denote the reward received by an algorithm on an arrival sequence W as

$$R(\text{ALG}; W) = \sum_{t=1}^T A_t \cdot G_t.$$

This notation is defined similarly to the notation $R(\text{GRD}; W)$ which we defined for the Dual-Adjusted Greedy Algorithm. For the prediction, we use the Prediction Algorithm to represent the special case of the Dual-Adjusted Greedy Algorithm where the dual variable is λ , and we let $R(\text{PRD}; W) = R(\text{GRD}; W)$. As stated in Proposition 1, if $\lambda = 1$, i.e., if $\lambda = \lambda^*$, then $R(\text{PRD}; W) \geq \frac{1}{2} \cdot R(\text{OPT}; W)$, which shows the Prediction Algorithm is essentially optimal if we have a perfect prediction.

Note that for any sequence of dual variables $\lambda_1, \dots, \lambda_T$, following λ_t at time period t gives a series of actions G_1, \dots, G_T . We define the depletion time of resource i by following $\lambda_1, \dots, \lambda_T$ to be the first time period such that the remaining amount of resource i is less than ϵ , that is, after this time period no actions that consumes

resource i is feasible (if this never happens we set the depletion time to be ∞). We will use the depletion time to quantify the behavior of λ . Intuitively, dual variables close to λ induce similar actions in most time periods as long as λ is not always on the boundary of decisions, and hence their depletion time should be similar. We make this idea formal using the following assumption on the depletion time.

Assumption 2 (Non-degenerate Prediction). Fix an arrival sequence W and a prediction λ for W there exists a constant $\bar{\epsilon} > 0$ satisfying the following: for any sequence of dual variables $\lambda^1, \dots, \lambda^T$ where $\lambda^t \in \mathbb{R}^C$ and $\|\lambda^t - \lambda\|_1 \leq \bar{\epsilon}$ for all t , the difference between the depletion time of resource i by following $\lambda^1, \dots, \lambda^T$ and by following λ is in $(-\bar{\epsilon}, \bar{\epsilon})$ for every resource i .

Assumption 2 roughly states that, for a sequence of dual variables that is close to the prediction, the action induced by the sequence of dual variables and the action induced by the prediction deplete resources at similar times. This assumption is reasonable and mild for the following reasons: in reality, most actions sets are discrete (such as accept, reject, N , etc.). Therefore for most λ , as long as it is not at the boundary (which is often a measure-zero set), dual variables close to λ induce the same action. Moreover, in practice it is also unlikely for the boundaries at each time period to be the same across a majority of time periods. Arrival rates λ^t vary over time, in which case Assumption 2 is satisfied with any prediction λ . Finally, perturbing each input with some small noise also turns a degenerate prediction into a non-degenerate one:

Proposition 2. Let the action set of each time period be $X = \{0, 1\}^C$. For any arrival sequence W and any prediction $\lambda < 0$, and any arbitrarily small $\epsilon > 0$, let W^ϵ be an arrival sequence obtained by randomly perturbing the reward functions to be A_C^0 in the following way: $A_C^0 = A_C^1 + \epsilon \lambda_C$ and $A_C^1 = A_C^0 - \epsilon \lambda_C$, where $\lambda_C = N^{-1} \mathbf{1} - \lambda$ are i.i.d. random variables. Then W^ϵ and λ satisfies Assumption 2 almost surely (with respect to the random perturbation). In particular, for sufficiently small $\epsilon > 0$, Assumption 2 is satisfied with probability $1 - \exp(-2\epsilon^2 Z)$ for some constant $Z > 0$.

The proof of Proposition 2 appears in Section B.1. Proposition 2 formalizes the intuition that perturbing each input with some small noise may turn a degenerate prediction into a non-degenerate one. Although Proposition 2 only considers $X = \{0, 1\}^C$ and perturbs the reward functions, similar results can be shown for more general X and perturbations on the resources consumption functions.

There are two primary arrival models when studying online problems: the stochastic (i.i.d) arrival model and the adversarial arrival model, both which we define formally below. Our goal is to design algorithms that have good performances under both arrival models, and for predictions of different qualities. Additionally, our algorithms should be oblivious to the arrival model and the prediction quality, i.e., they should have good performance without knowing the arrival model and the prediction quality.

Stochastic (i.i.d.) Arrival Model. The arrivals are drawn independently from an (unknown) underlying probability distribution $P \in \mathcal{P}(S)$, where S is the space of all probability distributions over \mathcal{S} . We measure the performance of an algorithm by its regret. Given an underlying arrival distribution P , the regret incurred by an algorithm ALG under P is defined as $E_{W \sim P} \sum_{j=1}^T W_j^0 - \sum_{j=1}^T W_j^{OPT}$. We will be concerned with the worst-case regret over all distributions $P \in \mathcal{P}(S)$: we define the regret of ALG to be

$$\text{Regret}^0(ALG) = \sup_{P \in \mathcal{P}(S)} E_{W \sim P} \sum_{j=1}^T W_j^0 - \sum_{j=1}^T W_j^{OPT}$$

Note that if the regret $\text{Regret}^0(ALG)$ is sub-linear in T , then algorithm ALG is essentially optimal on average as T goes to infinity.

Since $W \sim P$ is a random object, we assume that the prediction is defined on the same probability space as W , and that

$$k_\lambda \leq k_1 = \sum_{i=1}^n \lambda_i^2$$

holds for every realization of W .

To illustrate that this condition can hold for some $\lambda > 0$, suppose the optimal dual variable admits the decomposition

$$\lambda = \lambda_1 + \sum_{i=1}^n \lambda_i a_i + \sum_{i=1}^n \lambda_i^2$$

where $\lambda_1 \in \mathbb{R}^d$ is a fixed quantity in the dual space, $\lambda_i : \mathbb{R}^d \rightarrow \mathbb{R}^d$ are fixed functions, and a_i and λ_i^2 are random variables induced by the randomness of the arrival sequence W . We assume that the prediction \hat{w}_t has access to a_i so that

$$\lambda = \lambda_1 + \sum_{i=1}^n \lambda_i a_i$$

In this case, the prediction error $k_\lambda \leq k_1$ depends only on the randomness in λ_i^2 , and not on a_i . Consequently, if the variance of λ_i^2 is

small, we may have

$$k^{\lambda} \cdot k_1 = \epsilon^{1-2\epsilon}$$

for some $\epsilon > 0$, uniformly over all realizations of W .

In practice, the prediction \hat{w}_t is not assumed to depend causally on W , but is expected to be statistically related to W , for instance through shared observable co-variables.

Adversarial Arrival Model. The arrivals are arbitrary and adversarial, generated by an oblivious adversary, meaning that the arrival sequence is fixed prior to execution and does not depend on the algorithm's internal randomness or actions. Unlike the stochastic arrival model, regret here can be shown to grow linearly with any algorithm, so it is less meaningful to study the order of regret. Instead, we use competitive ratio as the performance metric. We say that an algorithm \mathcal{A} is asymptotically ϵ -competitive if ϵ^{-1} is the smallest number such that

$$\limsup_{\epsilon \rightarrow 0} \sup_{W \in \mathcal{W}} \left(\frac{1}{\epsilon} \text{OPT}(W) - \text{ALG}(W) \right) = o(\epsilon^{-1})$$

In words, if an algorithm is asymptotically ϵ -competitive, then it can obtain at least ϵ fraction of the optimal reward in hindsight as ϵ goes to infinity.

[21] proved that, without predictions, the lowest competitive ratio that any algorithm can achieve is $\frac{1}{2}$. [22] gave a mirror descent algorithm that achieves this competitive ratio. This is, loosely speaking, the best we might hope to achieve with bad predictions. On the other hand, we can always obtain $\text{PRD}(W)$ by following the prediction, which may exceed $\text{OPT}(W)$ with good predictions (indeed, as we have seen in Proposition 1, $\text{PRD}(W)$ can be as large as $\text{OPT}(W)$).

If we knew the prediction quality beforehand, we could obtain the maximum of the two by simply choosing the better approach (this is in fact the best we can hope for). Using this as the benchmark, we will compare an algorithm's reward to this maximum. That is, for an algorithm \mathcal{A} , we will analyze the following quantity:

$$\limsup_{\epsilon \rightarrow 0} \sup_{W \in \mathcal{W}} \left(\frac{1}{\epsilon} \max \left\{ \text{OPT}(W), \text{PRD}(W) \right\} - \text{ALG}(W) \right)$$

¹¹This definition serves as an implicit assumption that $\text{OPT}(W)$ grows linearly in ϵ . When $\text{OPT}(W)$ is sub-linear in ϵ , by our definition every algorithm is ϵ -competitive.

¹²We assume the arrival sequence $W = (w_1, \dots, w_T)$ is fixed in advance. Our results still hold if the arrival sequence is chosen by a non-oblivious or adaptive adversary who does not know the internal randomization of the algorithm.

A Measure of Stationarity

Ideally, one would hope to develop an algorithm that achieves the best performance under both stochastic and adversarial arrivals respectively without knowing the prediction quality and the underlying arrival model. However, we will show in the next section that this is provably not achievable by any algorithm.

For an arrival sequence W , its stationarity is closely related to the difficulty of solving the instance it created. As examples, an arrival sequence generated independently from the same underlying distribution can be considered as completely stationary, an arrival sequence that has certain seasonality/periodicity with small trend (e.g. generated from time series models) is less stationary, and an arrival sequence that is adversarially chosen (e.g. the lower bound instance) is completely nonstationary. Intuitively, an arrival sequence is more stationary if certain parts of the sequence with the same length are similar to each other. In this subsection we formalize this idea and develop a measure of arrival sequences' stationarity. We then use it to quantify algorithms' performances.

For a time interval from time period B to time period C let $W_{B:C} = W_B \dots W_C$ denote the arrival sequence from time period B to time period C . We define the $W_{B:C}$ subproblem to be the problem instance where the arrival sequence is $W_{B:C}$ and the total amount of resources is $d_{B:C} = d \cdot \frac{C-B+1}{T}$, i.e., scaled down proportionally. In particular, the (online) optimum of the $W_{B:C}$ subproblem is:

$$\text{OPT}_{W_{B:C}} := \max_{G \in \mathcal{X}_{B:C}} \sum_{t=B}^C A_t^T G_t \quad \text{s.t.} \quad \sum_{t=B}^C G_t \leq d_{B:C} \mathbf{1}.$$

Similarly, we use $\text{GRD}_{W_{B:C}}(\lambda)$ to denote the amount of reward obtained by the Dual-Adjusted Greedy Algorithm with dual variable λ on the $W_{B:C}$ subproblem.

Definition 4 (Measure of Stationarity). Given the total number of available resources d , an arrival sequence $W = W_1 \dots W_T$ is λ -stationary for some $\lambda \in [0, 1]$ and $0 \leq \lambda \leq 1$ if for every $\lambda \in [0, 1]$:

$$\min_{j=1, \dots, T} \text{GRD}_{W_{j:j+\lambda T}}(\lambda) \geq \lambda \text{GRD}_{W_{j:j+\lambda T}}(0) + (1-\lambda) \text{GRD}_{W_{j:j+\lambda T}}(1).$$

Intuitively, W being λ -stationary (roughly) means when we break W into two subproblems at any time period that is a multiple λT of the rewards obtained by these two subproblems sum up to be at least λ portion of the reward obtained by W .

A few remarks are in order:

- ϵ measures the number of possible partition time periods that makes the subproblems similar to the original problem. On the extreme, $\epsilon \rightarrow 0$ means that every subproblem is similar to the original problem, and imposes no restrictions on W . As examples, if W s generated i.i.d. from some underlying distribution, ϵ can be arbitrarily close to 0, and if W s periodic with small period, ϵ can be small.
- δ measures the loss in the partition, which can be viewed as the similarity of the subproblems to the original problem. On the extreme, $\delta = 0$ means W can be partitioned at time periods X without losing much rewards, and $\delta = 1$ imposes no restrictions on W . As an example, W having small trend (i.e., the infinity-norm of the vector of possible rewards is similar across all time periods) implies small δ .
- Smaller ϵ and smaller δ both represent more stationarity. Note that there is no single ϵ - δ pair for an arrival sequence, but rather each choice of X gives a corresponding δ , and smaller ϵ yields larger δ , i.e. the subproblems become less similar as the partition becomes more granular. The relation of ϵ and δ will become clear when we state our main theorem, and we will not need to know the value of δ in our algorithm.
- Unlike usual stochastic definitions of stationarity (such as [88] and [116]), here it is defined for deterministic arrival sequences. We show in the proposition below that if the arrivals are stochastic (i.i.d.), then the arrival sequence is ϵ -stationary for any $\delta > 0$ with high probability. This shows our definition of stationarity is compatible with stochastic definitions of stationarity.

Proposition 3. If an arrival sequence W is generated under the stochastic (i.i.d.) arrival model, then W is ϵ -stationary for any constant $\delta > 0$ with probability at least $1 - \delta^2$.

The proof of Proposition 3 appears in Appendix B.1.

- In most stochastic definitions of stationarity (such as [88] and [116]), the nonstationary measure is meant to characterize the instances that are possible to achieve sub-linear regret. All instances that incur linear regret are viewed as completely nonstationary. In our paper, the nonstationary measure is meant to characterize the instances that are impossible to achieve sub-linear regret

(or equivalently, that are impossible to achieve competitive ratio that is 1). All instances that incurs sub-linear regret are viewed as completely stationary.

- ^ Unlike usual definitions of stationarity, our definition is problem- (i.e. resource-) dependent. For example $d_i = 0$ and d sufficiently large both imply X can be arbitrarily small and $\underline{d} = 0$, since any partition of the arrival sequence gives the same reward.

3.3 Main Results

In this section, we first present previous results for the Online Resource Allocation problem without predictions, and then give our main results on the full problem (with predictions) along with matching lower bounds.

Prior Results: Online Resource Allocation without Predictions

[22] studied the no-prediction version of our problem and gave a mirror descent algorithm which achieves the best achievable performance under both arrival models without knowing the underlying arrival model. We discuss their algorithm in detail in Appendix B.2. They proved the following performance guarantee for their Mirror Descent Algorithm (MDA):

Proposition 4 (Theorem 1 and Theorem 2 in [22]). Consider the Mirror Descent Algorithm (MDA). It holds that:

- (1) If the arrivals are stochastic,

$$\text{Regret}^{\text{MDA}^0} = \mathcal{O}(U^{1/2});$$

- (2) If the arrivals are adversarial,

$$\limsup_{T \rightarrow \infty} \sup_{W \in \mathcal{W}} \left(\frac{1}{T} \text{OPT}^{\text{MDA}^0}(W) - \frac{1}{T} \text{MDA}^0(W) \right) = 0.$$

Proposition 4 shows that the Mirror Descent Algorithm achieves $\mathcal{O}(U^{1/2})$ regret and is U -competitive, which are both optimal [10, 21].

Prior Results: Lower Bounds

As a final step before describing our results, we present previous lower bounds for the full problem with predictions and known arrival model.

Stochastic Arrivals. Without predictions, the best achievable regret by any algorithm is $\frac{1}{2} \sum_{i=1}^T g_i^2$ [10]. With predictions, [41] gave the following lower bound on the best achievable regret with known accuracy parameter θ :

Proposition 5 (Corollary of Theorem 2 in [41]). Under stochastic arrival model, given a prediction λ with accuracy parameter θ , no algorithm can achieve regret better than $\frac{1}{2} \sum_{i=1}^T g_i^2 - 1/g^0$, even if θ is known.

Adversarial Arrivals. Without predictions, the best achievable reward by any algorithm (taken the worst-case across all problem instances) is $\frac{1}{U} \text{OPT}^1 W^0$ [21]. On the other hand, simply following the actions induced by the prediction at each time yields reward $\text{PRD}^1 W^0$. As we have seen in Proposition 1, for good predictions $\text{PRD}^1 W^0$ can be as high as $\frac{1}{U} \text{OPT}^1 W^0$. Hence we have the following lower bound under adversarial arrivals:

Proposition 6 (Corollary of Theorem 3.1 in [21]). Under adversarial arrival model, given a prediction with accuracy parameter θ , no algorithm can achieve (worst-case across all problem instances) reward higher than $\frac{1}{U} \text{OPT}^1 W^0 - \text{PRD}^1 W^0$, even if θ is known.

Our Results: Online Resource Allocation with Predictions

We are finally prepared to state our main result, which to develop a single algorithm that achieves the optimal performance using predictions, without knowing the underlying arrival model and the prediction accuracy.

Theorem 2 (Upper Bound). Assume that Assumptions 1 and 2 hold. Given a prediction λ with (unknown) accuracy parameter θ and given $\theta \leq \frac{1}{U} \sum_{i=1}^T g_i^2$, there exists an algorithm (MainALG) such that:

- (1) If the arrivals are stochastic,

$$\text{Regret}^1 \text{MainALG}^0 \leq \frac{1}{2} \sum_{i=1}^T g_i^2 - 1/g^0;$$

- (2) If the arrivals are adversarial and $\frac{1}{U} \sum_{i=1}^T g_i^2$ -stationary,

$$\limsup_{T \rightarrow \infty} \sup_{W \in \mathcal{W}} \left(\frac{1}{T} \sum_{t=1}^T \max_{j \in [K]} \left(\frac{1}{U} \text{OPT}^1 W^0 - \text{PRD}^1 W^0 \right) - \text{MainALG}^j W^0 \right) \leq \frac{1}{U} \sum_{i=1}^T g_i^2 - 1/g^0.$$

A few remarks are in order:

- ^ The performance guarantee under adversarial arrivals is better for smaller X , which matches the intuition that one can hope to achieve better performance with more stationary arrivals.
- ^ If the arrivals are known to be adversarial, which we will discuss in the next section (Algorithm 2 and Proposition 9), there exists an algorithm that achieves

$$\limsup_{T \rightarrow \infty} \sup_{W \in \mathcal{W}} \frac{1}{T} \max_{\text{ALG}} \frac{1}{U} \text{OPT}^1 W^0 - \text{PRD}^1 W^{01} \text{ALG} \leq \frac{1}{U} \text{OPT}^1 W^0 - \text{PRD}^1 W^{01} \text{ALG} \leq 0$$

That is, we are able to not suffer from nonstationarity. This is because the performance requirement is much higher for stochastic arrivals (sub-linear regret), which requires a conservative consumption of resources and hence obtains less rewards when the arrivals are highly nonstationary. This idea is elaborated in the lower bound construction below.

We provide a lower bound which shows Theorem 1 is tight in the sense that for any algorithm that achieves sub-linear regret under stochastic arrivals, one cannot replace $\frac{1}{U}$ with any number smaller and still get meaningful guarantees under adversarial arrivals. The proof of Proposition 7 appears in Appendix B.3. The lower bound construction consists of two instances, stochastic with bad prediction and adversarial with good prediction, that are provably indistinguishable for a certain period of time.

Proposition 7 (Lower Bound). For any $0 < \epsilon < 1$, $0 < \delta < 1$, and $\epsilon > \delta$, there exists a sequence of instances $\{I_n\}_{n=1}^{\infty}$ of increasing time horizon n that satisfies Assumptions 1 and 2, such that for any algorithm (ALG), at least one of the following holds:

- (1) The arrivals are stochastic, and

$$\text{Regret}^1 \text{ALG}^0 = o(n^\epsilon);$$

- (2) The arrivals are adversarial and $\frac{1}{U}$ -stationary, and

$$\limsup_{T \rightarrow \infty} \frac{1}{T} \max_{\text{ALG}} \frac{1}{U} \text{OPT}^1 W^0 - \text{PRD}^1 W^{01} \text{ALG} \leq \frac{1}{U} \text{OPT}^1 W^0 - \text{PRD}^1 W^{01} \text{ALG} \leq \epsilon \cdot XA.$$

3.4 Algorithm and Proof of Main Result

In this section, we first present two algorithms that utilize the prediction in an optimal way for the two arrival models respectively. Then we combine these two algorithms to a single algorithm that is oblivious to both the prediction quality and the arrival model, which completely solves the Online Resource Allocation with Prediction problem.

Our algorithms for each arrival model utilize mirror descent, which take an initial dual variable, a step-size, and a reference function as inputs. At each time period t , the algorithms take the action induced by the current dual variable and performs a first-order update on the dual variable. With prediction \hat{x}_t , a natural initialization of the dual variable is to set $\lambda_t = \lambda$, i.e., the algorithms start by assuming the prediction is accurate. Then, the algorithms use adaptive step sizes in mirror descent steps depending on the arrival model and the prediction's behavior.

Algorithm for the Stochastic Arrival Model

Let \hat{x} be a prediction with accuracy parameter θ , i.e., $\|\hat{x} - x^*\|_1 \leq \theta$. By Proposition 5, no algorithm can achieve regret better than $\frac{1}{2} \theta - 1g^*$ even if θ is known. As a comparison, we can show that the optimal fixed step size for the Mirror Descent Algorithm is $\frac{1}{2}$ using similar method as the proof of Proposition 4, which incurs $\frac{1}{2} \theta - 1g^*$ regret. Therefore, mirror descent with fixed step size is sub-optimal even if the prediction quality is known. This suggest us to use adaptive step sizes. The step size we use is drawn from their work in parameter-free optimization¹⁴.

Algorithm 18, which we call the Stochastic Arrival Algorithm (SA), initializes the dual variable at the prediction \hat{x} and updates the dual variable at each time period through mirror descent with fine-tuned step sizes. A detailed description of Algorithm 18 appears in Appendix B.4.

Proposition 8. Consider the Stochastic Arrival Algorithm (SA) under the stochastic arrival model. Given a prediction \hat{x} with (unknown) accuracy parameter θ , it holds that:

$$\text{Regret}^{\text{SA}} \leq \frac{1}{2} \theta - 1g^*.$$

¹³For completeness, in Appendix B.2 we state the standard assumptions on choosing the reference function ψ for mirror descent algorithms [30, 39, 119, 120].

¹⁴It follows the line of work in the more general online learning problem of parameter-free regret minimization [45, 54, 56, 132, 158].

The proof of Proposition 8 can be found in Appendix B.5. By Proposition 5, the Stochastic Arrival Algorithm achieves optimal worst-case regret up to logarithm factors.

Algorithm for the Adversarial Arrival Model

Different from the stochastic arrival model, under the adversarial arrival model it is impossible to achieve sub-linear worst-case regret. Instead, we directly compare the reward obtained by our algorithm to the maximum reward of two natural benchmark algorithms: the Mirror Descent Algorithm (which is optimal when the prediction quality is low) and the Prediction Algorithm (which is optimal when the prediction quality is high). That is, for an algorithm ALG , we will analyze the following quantity:

$$\limsup_{T \rightarrow \infty} \sup_{W \in \mathcal{W}} \frac{1}{T} \max \left\{ \frac{1}{U} \text{OPT}^1 W^0 - \text{PRD}^1 W^0, \text{ALG}^j W^0 \right\} \quad \bullet$$

We give Algorithm 2 for the adversarial arrival model, which we call the Adversarial Arrival Algorithm (AA). It performs mirror descent with fixed step size $\frac{1}{\sqrt{T}}$.

Algorithm 2: Adversarial Arrival Algorithm (AA)

Inputs: Prediction λ , total time period T , initial resources $\mathbf{1} = \mathbf{d}$, reference function $\mathbf{1}^0 : \mathcal{R}^d \rightarrow \mathbb{R}$, upper bound function $\mathbf{1}^0 = \mathbf{1}^0$, and step-size $\frac{1}{\sqrt{T}}$.

Initialize $\mathbf{1} = \lambda$

for C from 1 to T do

Receive request $\mathbf{1}_C = \mathbf{1}_C^0$

Make the primal decision \mathbf{c}_C and update the remaining resources

$$\mathbf{c}_C = \arg \max_{\mathbf{c} \in \mathcal{C}} \left\{ \mathbf{1}_C^0 \mathbf{c} - \frac{1}{\sqrt{T}} \langle \mathbf{q}_C, \mathbf{c} \rangle \right\}$$

Obtain a sub-gradient of the dual function:

$$\mathbf{q}_C \in \partial_{\mathcal{R}^d} \left(\mathbf{1}_C^0 \mathbf{c}_C \right)$$

Update the dual variable by mirror descent:

$$\mathbf{q}_C = \arg \min_{\mathbf{q} \in \mathcal{R}^d} \left\{ \langle \mathbf{q}, \mathbf{1}_C^0 \mathbf{c}_C \rangle + \frac{1}{\sqrt{T}} \text{Bregman}(\mathbf{q}, \mathbf{1}^0) \right\}$$

where $\text{Bregman}(\mathbf{q}, \mathbf{1}^0) := \mathbf{1}^0(\mathbf{q}) - \mathbf{1}^0(\mathbf{1}^0) - \langle \mathbf{1}^0, \mathbf{q} - \mathbf{1}^0 \rangle$ is the Bregman divergence.

end

Proposition 9. Assume that Assumptions 1 and 2 hold. Consider the Adversarial Arrival Algorithm (AA) under the adversarial arrival model. Given a prediction λ with (unknown) accuracy parameter ϵ , it holds that:

$$\limsup_{T \rightarrow \infty} \sup_{W \in \mathcal{W}} \frac{1}{T} \max \left\{ \frac{1}{U} \text{OPT}^1 W^0 - \text{PRD}^1 W^0, \text{AA}^j W^0 \right\} \leq \epsilon \quad \bullet$$

The proof of Proposition 9 can be found in Appendix B.5. Proposition 9 states that the Adversarial Arrival Algorithm achieves the maximum of the two benchmark algorithms without knowing the prediction quality. Proposition 9 is tight by Proposition 6.

Main Algorithm: Detection of Nonstationarity

With the Stochastic Arrival Algorithm and the Adversarial Arrival Algorithm, we are ready to present our main algorithm, which merges the two algorithms and works for both arrival models.

Algorithm 3: Main Algorithm

Inputs: Prediction λ , total time periods T , initial rewards $r_i^1 = 0$, initial resources $r_i = 0$, reference function $\phi^0 : \mathbb{R}^d \rightarrow \mathbb{R}$, upper bound function $\psi^0 : \mathbb{R}^d \rightarrow \mathbb{R}$, constant β which is specified in Equation (B.30) in Appendix B.6, constant $\epsilon \in (0, 1]$, and initial step-size η

for C from 1 to T do

Receive request (A_C, C, X_C)

if $C = \lfloor \beta X \rfloor c, 1$ for some $c = 0, \dots, d-1$ then

if $\sum_{t=1}^c r_{c,t} \geq \log^2(\beta) \cdot \text{OPT}_{c,1} - \beta \cdot C$ then

Release resources for the next βX time periods: $r_{c,t} = 0, t = c, \dots, c + \beta X - 1$

Take action a given by the Stochastic Arrival Algorithm with the following inputs: total time periods βX , initial dual variable λ , initial resources $r_{c,1}$, reference function ϕ^0 , and initial step-size η

Update resources: $r_{c,t} = r_{c,t} - \beta \cdot C$

Update rewards: $r_{c,1} = r_{c,1} + A_C$

else

break

end

end

Release all resources: $r_{c,t} = 0, t = c, \dots, c + \beta X - 1$

Use the Adversarial Arrival Algorithm with initial dual solution λ , remaining resources r_C and step size η .

end

The Main Algorithm initially assumes stochastic arrivals and runs the Stochastic Arrival Algorithm, cautiously releasing resources to avoid over-consumption. Simultaneously, it monitors observed arrivals to test this assumption: under stochastic arrivals, the accumulated reward should be comparable to the optimal offline reward under the underlying distribution. If the observed reward is significantly lower, the

algorithm infers with high probability that arrivals are not stochastic and switches to the Adversarial Arrival Algorithm for the remaining periods. If arrivals are adversarial but sufficiently stationary, this test may not trigger; however, in that case the Stochastic Arrival Algorithm remains effective, and no switch is needed.

Theorem 1 (Upper Bound). Consider the Main Algorithm (MainALG). Assume that Assumptions 1 and 2 hold. Given a prediction with (unknown) accuracy parameter θ and given $\delta \in (0, 1]$, it holds that:

(1) If the arrivals are stochastic,

$$\text{Regret}(\text{MainALG}, \theta) \leq \frac{1}{2} \delta - \delta g^0;$$

(2) If the arrivals are adversarial and δ -stationary,

$$\limsup_{T \rightarrow \infty} \sup_{W \in \mathcal{W}} \frac{1}{T} \delta \max_{j \in [K]} \frac{1}{U} \text{OPT}^j - \text{PRD}^j \leq \text{MainALG}^j - \delta \cdot \text{XA}^j.$$

The proof of Theorem 1 appears in Appendix B.6. Theorem 1 is tight by Proposition 7.

3.5 Experiments

Finally, we evaluate our algorithm empirically using both synthetic and real data. The main takeaway is that our algorithm is robust to prediction quality: its reward is consistently close to the better of the Mirror Descent Algorithm (which is worst-case optimal without predictions) and the Prediction Algorithm.

In both experiments, we consider sequential assortment optimization, modeling an online retailer making in-cart recommendations. Each arrival corresponds to a customer checkout, where the algorithm recommends a subset of products of fixed cardinality. Each recommended product is purchased with a customer-specific probability, generating revenue. In the Online Resource Allocation with Predictions framework, actions correspond to assortments, resources to product inventories, and rewards to expected profits.

Each experimental instance consists of fixed initial inventories, an arrival sequence, and a prediction of the shadow price for each product, with prediction quality varying across instances. For each instance, we compare three total rewards: those obtained by our Main Algorithm, the Mirror Descent Algorithm (MDA), and the Prediction Algorithm (PRD).

Our primary performance metric is the optimality gap. For an instance j , let R^j_{PRD} , R^j_{MDA} , and R^j_{MainALG} denote the rewards of PRD, MDA,

and the Main Algorithm, respectively. We define

$$\text{GAP}^{j, g} = \frac{\text{MainALG}^{j, g} - \min\{\text{PRD}^{j, g}, \text{MDA}^{j, g}\}}{\max\{\text{PRD}^{j, g}, \text{MDA}^{j, g}\} - \min\{\text{PRD}^{j, g}, \text{MDA}^{j, g}\}}.$$

This metric measures how close the Main Algorithm is to the better benchmark, normalized so that $\text{GAP} = 1$ corresponds to matching the better benchmark and $\text{GAP} = 0$ to matching the worse one (a random choice between PRD and MDA yields $\text{GAP} = 0.5$)¹⁵.

Synthetic Experiment

We began with a set of smaller, synthetic experiments with 25 products over 1000 time periods, and the task of recommending 2 products at a time. We assumed customers belong to one of 25 types. The process we used to randomly generate the product prices, the initial inventory levels, and the customer type-specific purchase probabilities, is described in Appendix B.7. Each time period corresponds to a single arriving customer (drawn uniformly from the 25 types), or no arrival. We used three types of arrival sequences, with the probability of a customer arrival changing over time: stationary with a fixed arrival probability of 0.7, nonstationary with an arrival probability linearly increasing from 0.4 to 1.0, and adversarial with an arrival probability of 0.0 during the first 300 periods and 1.0 afterward. We randomly generated predictions of varying qualities by computing the optimal shadow prices, and adding mean-zero Gaussian noise with standard deviations of 500 (bad), 5 (good), and 0 (perfect).

The results are summarized in the top half of Table 3.2, which for nine random ensembles of instances (depending on arrival model and prediction quality) reports both the median GAP, and the proportion of instances for which the GAP was at least 0.5 (for both values, higher is better). Recall that no algorithm can be expected to achieve high GAP values (say above 0.5) for all nine ensembles simultaneously. We see that our algorithm generally performs better with higher prediction quality and higher stationarity. Now one issue with GAP as a performance metric is that it can be quite erratic when the Mirror Descent and Prediction algorithms generate similar rewards (as their difference is the denominator in GAP), and these are arguably the instances in which GAP matters the least from a practical standpoint. Therefore, in the bottom half of Table 3.2 we removed the instances for which the two rewards are within 25% of each other. The remaining instances are exactly the instances

¹⁵Technically, GAP may lie outside $[0, 1]$.

Median	1	CDF ^{10•5⁰}	Stochastic		Nonstationary		Adversarial	
		Perfect Predictions	0.81	0.63	0.83	0.64	0.65	0.56
		Good Predictions	0.77	0.62	0.81	0.64	0.64	0.56
		Bad Predictions	0.54	0.52	0.45	0.49	0.22	0.40

Median	1	CDF ^{10•5⁰}	Stochastic		Nonstationary		Adversarial	
		Perfect Predictions	0.71	0.65	0.72	0.66	0.64	0.60
		Good Predictions	0.67	0.63	0.72	0.67	0.52	0.54
		Bad Predictions	0.58	0.55	0.49	0.49	0.36	0.40

Table 3.2: Summary of synthetic experiments. For each of three levels of prediction quality (the rows), and each of three generative arrival models (the columns), two summary statistics are reported over a random ensemble of instances: (left) the median GAP, and (right) the proportion of instances for which the GAP was at least 0.5. (Top) Results over all instances. (Bottom) Results over instances for which the rewards of the Mirror Descent and Prediction algorithms differ by at least 25%.

where robustness matters the most. We see that our algorithm is more robust to bad predictions on these instances.

H&M Experiment

We used a dataset from H&M (a fast-fashion clothing retailer), which contains the online transactions of 105,542 products from 2018 to 2020, along with product features. Because most products have zero or few transactions in two years, we selected the products with the top 5000 number of total transactions for our experiment, which includes 13,697,790 transactions. Our task was to recommend three products.

Each instance runs across three month's transactions from the data. The time horizon for each instance was the maximum number of transactions per day (103,473) multiplied by the total number of days (90), so that each day contained 103,473 time periods, each having zero or one arriving customer. To estimate customer-specific purchase probabilities, we used (customer, transaction time, product, price of product)-tuples and trained a random forest algorithm with the corresponding features of this tuple (a 209-dimensional vector after encoding) to estimate the probability that the customer, who brought product i at that certain time period with the certain price, would also buy product j if recommended.

To generate predictions, we used three popular forecasting methods ranging from classical algorithms to the state-of-the-art:

- ^ Prophet: A recent algorithm developed by Facebook [162] based on a (piecewise-linear) trend and seasonality decomposition, known to work well in practice with minimal tuning. Tuning parameters: software default.
- ^ Exponential Smoothing (Holt Winters): A classic algorithm based on a (linear) trend and seasonality decomposition, known for its simplicity and robust performance. It is frequently used as a benchmark in forecasting competitions [125]. Tuning parameters: seasonality of length.
- ^ ARIMA: Another classic algorithm that is rich enough to model a wide class of nonstationary time-series. Tuning parameters: ϕ, θ, A^0 .

These experiments were run on an N2D Series machine on Google Cloud's Compute Engine, with 224 vCPUs and 896GBs of memory. The total computation time was around 140 hours.

The results are summarized in Figure 3.3, which contains histograms of the GAPs across an ensemble of 100 instances (for varying three-month periods in the data), separately for each forecasting algorithm. The average GAP is 0.68 on instances with Prophet forecasts, 0.58 on instances with ARIMA forecasts, and 0.53 on instances with Exponential Smoothing forecasts. Because the average GAPs are large with all three forecasting methods, our algorithm performs close to the better one of the Prediction algorithm and the Mirror Descent Algorithm, showcasing its robustness to the unknown prediction accuracy.

(a) GAPs of Prophet

(b) GAPs of ARIMA

(c) GAPs of Exponential Smoothing

Figure 3.3: Histograms of GAPs with different forecasting methods, each containing 100 instances.

3.6 Conclusion

In this paper, we propose a new model for incorporating predictions into the Online Resource Allocation Problem. We first develop separate algorithms for stochastic and adversarial arrival models: under stochastic arrivals, the algorithm achieves nearly

optimal minimax worst-case regret, while under adversarial arrivals, it achieves a nearly optimal reward. Both algorithms operate without prior knowledge of prediction quality. Building on these, we propose a unified algorithm that attains the same guarantees under the appropriate arrival model, without knowing either the arrival model or the prediction quality in advance. The key idea is to initially treat arrivals as stochastic while continuously performing hypothesis tests to detect adversarial behavior.

Chapter 4

THE NONSTATIONARY NEWSVENDOR WITH (AND WITHOUT) PREDICTIONS

The classic newsvendor model yields an optimal decision for a newsvendor selecting a quantity of inventory, under the assumption that the demand is drawn from a known distribution. Motivated by applications such as cloud provisioning and staffing, we consider a setting in which newsvendor-type decisions must be made sequentially, in the face of demand drawn from a stochastic process that is both unknown and nonstationary. All prior work on this problem either (a) assumes that the level of nonstationarity is known, or (b) imposes additional statistical assumptions that enable accurate predictions of the unknown demand. Our research tackles the Nonstationary Newsvendor without these assumptions, both with and without predictions.

We first, in the setting without predictions, design a policy which we prove (via matching upper and lower bounds) achieves order-optimal regret. Ours is the first policy to accomplish this without being given the level of nonstationarity of the underlying demand. We then, for the first time, introduce a model for generic (i.e. with no statistical assumptions) predictions with arbitrary accuracy, and propose a policy that incorporates these predictions without being given their accuracy. We upper bound the regret of this policy, and show that it matches the best achievable regret had the accuracy of the predictions been known.

Our findings provide valuable insights on inventory management. Managers can make more informed and effective decisions in dynamic environments, reducing costs and enhancing service levels despite uncertain demand patterns. This study advances understanding of sequential decision-making under uncertainty, offering robust methodologies for practical applications with nonstationary demand. We empirically validate our new policy with experiments based on three real-world datasets containing thousands of time-series, showing that it succeeds in closing approximately 74% of the gap between the best approaches based on nonstationarity and predictions alone.

4.1 Introduction

The newsvendor problem is a century-old model that remains fundamental to the practice of operations management. In its original instantiation, a newsvendor is tasked with selecting a quantity of inventory before observing the demand for that inventory, with the demand itself randomly drawn from a known distribution. The newsvendor incurs a per-unit underage cost for unmet demand, and a per-unit overage cost for unsold inventory. The objective is to minimize the total expected cost, and the classic result is that the optimal inventory level is a certain problem-specific quantile (depending only on the underage and overage costs) of the demand distribution.

This paper is concerned with a modern instantiation of the same model, consisting of a sequence of newsvendor problems over time, each with unknown demand distributions that vary over time. While this version of the problem is arguably ubiquitous in practice today, it may be worth highlighting a few motivating examples:

- ^ Cloud Provisioning: Consider a website which provisions computational resources from a commercial cloud provider to serve its web requests. Such provisioning is typically done dynamically, say on an hourly basis, with the aim of satisfying incoming requests at a sufficiently high service level. Thus, the website faces a single newsvendor problem every hour, with an hourly demand that can (and does) vary drastically over time.
- ^ Sta ng: A more traditional example is sta ng, say for a brick-and-mortar retailer, a call center, or an emergency room. Each day (or even each shift) requires a separate newsvendor problem to be solved, with demand that is highly nonstationary.

Despite its ubiquity, this problem is far from resolved, precisely because the demand (or sequence of demand distributions) is both nonstationary and unknown indeed, the repeated newsvendor with stationary, but unknown demand was solved by [and the same setting with known, but nonstationary demand can be treated simply as a sequence of completely separate newsvendor problems. At present, there are by and large two existing approaches to this problem:

1. Limited Nonstationarity: One approach is to design policies which succeed under limited nonstationarity, i.e. the cost incurred by the policy should be parameterized by some carefully-chosen measure of nonstationarity (e.g. quadratic variation), and nothing else. This approach has proved fruitful

across a diverse set of problems ranging from dynamic pricing [51] to multi-armed bandit problems [6] to stochastic optimization [37]. Most relevant here, the recent work of [9] applies this lens to the newsvendor setting (we will discuss this work in detail momentarily). This approach yields policies with theoretical guarantees that are quite robust—no assumption on the demand (beyond the limited nonstationarity) is required. However, this is far removed from practice, where the next approach is more common.

2. Predictions: The second approach is to utilize some sort of predictions of the unknown demand. These predictions can be generated from simple forecasting algorithms for univariate time-series, all the way to state-of-the-art machine learning models that leverage multiple time-series and additional feature information. Therefore these predictions may contain much more information than past demand data points, such as various features/contexts, or even black-box type information that is non-identifiable. In addition to being the de facto approach in practice, the use of predictions in newsvendor-type problems is well-studied, and in fact provable guarantees exist for many specific prediction-based approaches [24, 83, 143, 181]. All such guarantees rely on (at the very least) the demand and potential features being generated from a known family of stochastic models, so that the framework and tools of statistical learning theory can be applied. Absent these statistical assumptions, it is unclear a priori whether the resulting predictions will be sufficiently accurate to outperform robust policies such as those generated in the previous approach. As a concrete example of this, see Figure 4.1, which demonstrates on a real set of retail data that prediction accuracy may vary drastically and unexpectedly, even when those predictions are generated according to the same procedure and applied during the same time period.

To summarize, the repeated newsvendor with unknown, nonstationary demand (which from here on we refer to as the Nonstationary Newsvendor) admits policies with nontrivial guarantees, which can be made significantly better or worse by following predictions. This suggests the opportunity to design a policy that uses predictions optimally, in the sense that the predictions are utilized when accurate, and ignored when inaccurate. Ideally, such a policy would run without knowledge of (a) the accuracy of the predictions and (b) the method with which they are generated. This is precisely what we accomplish in this paper.

Figure 4.1: Daily number of customers (in blue), from September 2014 to January 2015, at two different stores in the Rossmann drug store chain. Predictions (in red), starting November 2014, are generated using Exponential Smoothing with the same fitting process. The store in the upper sub-figure has substantially more accurate predictions ($r^2 = 0.88$) than that of the lower sub-figure ($r^2 = 0.11$).

The Nonstationary Newsvendor, with and without Predictions

The primary purpose of this paper is to develop a policy that optimally incorporates predictions (defined in the most generic sense possible) into the Nonstationary Newsvendor problem. Naturally, a prerequisite to this is a fully-solved model of the Nonstationary Newsvendor without predictions. At present this prerequisite is only partially satisfied (via the work of [6]), so a nontrivial portion of our contributions will be to fully solve this problem.

Without predictions, the Nonstationary Newsvendor consists of a sequence of newsvendor problems indexed by periods $1, \dots, T$, each with unknown demand distribution \mathcal{D}_t . The level of nonstationarity is characterized via a variation parameter $E \in [0, 1]$, where $E = 0$ essentially amounts to stationary demand, and $E = 1$ is effectively arbitrary (in a little more detail: a deterministic analogue of quadratic variation is applied to the sequence of means μ_1, \dots, μ_T and $E \in [0, 1]$ is the exponent such that this quantity equals E). Finally, we measure the performance of any policy using regret, which is the expected difference in the total cost incurred by the policy versus that of an optimal policy that knows the demand distributions. At minimum we aim to design a policy that achieves sub-linear ($\mathcal{O}(\sqrt{T})$) regret, as such a policy would incur a per-period cost that is on average no worse than the

optimal, as ϵ grows. We will in fact design policies which achieve order-optimal regret with respect to the variation parameter E .

To this base problem, we introduce the notion of predictions. In each period we receive a prediction \hat{d}_t of the mean demand $d_t = E\{D_t\}$ before selecting the order quantity. Our predictions are generic: no assumption is made on how they are generated. We measure the accuracy of the predictions through an accuracy parameter $\theta \in [0, 1]$ defined such that $\sum_{t=1}^T |d_t - \hat{d}_t| \leq \theta T$. Notice that when $\theta = 0$ the predictions are almost perfect, and when $\theta = 1$ the predictions are effectively useless. We will characterize a precise threshold θ^* (which depends on E) that determines when the predictions should be utilized. Our primary challenge will be to design a policy that makes use of the predictions only when they are sufficiently accurate, and without having access to E . As to the variation parameter E , we will separately consider policies which do and do not have access to E . This distinction will turn out to be the critical factor in classifying what is and is not achievable.

Our Contributions

Our primary contributions can be summarized as follows.

1. Nonstationary Newsvendor (without predictions): We completely solve the Nonstationary Newsvendor problem. This consists of first constructing a policy and proving an upper bound on its regret:

Theorem 2 (Informal). There exists a policy which achieves $O(\sqrt{E})$ regret without knowing E .

We then show that this regret is minimax optimal up to logarithmic factors:

Proposition 10 (Informal). No policy can achieve regret better than $\Omega(\sqrt{E})$, even if E is known.

As alluded to earlier, [6] previously initiated the study of the Nonstationary Newsvendor. Our results are distinct in terms of both modeling and theoretical contributions. We will expound these distinctions more carefully later on.

^ Modeling: The most crucial difference in our model is that we allow both the demand and the set of possible ordering quantities to be discrete. This is

¹The $\Omega(\cdot)$ notation hides logarithmic factors.

certainly of practical concern (e.g. physical inventory, employees, and virtual machines are all indivisible units of demand), but moreover we will show that the results of [96] require both the demand and set of feasible ordering quantities to be continuous. Thus, there is no overlap in our theoretical results.

- ^ Results: [96] succeed in designing a policy that achieves order-optimal regret, but crucially, their policy requires that the variation parameter E be known. In addition to being concerning from a practical standpoint, this leaves open the theoretical question of what exactly is achievable in settings for which E is unknown. Our results show that the same regret can be achieved without knowing E .

2. Nonstationary Newsvendor with Predictions: We construct a policy that optimally leverages predictions, i.e. it is robust to unknown prediction accuracy. To be precise, the previous contribution offers a policy that achieves $\mathcal{O}(E^{1/3}, E^{2/3} \cdot \theta)$ regret, and predictions yield a simple policy that achieves $\mathcal{O}(\theta)$ regret, so we would expect that the best possible regret is the minimum of these two quantities. We show this formally:

Proposition 11 (Informal). No policy can achieve regret better than $\mathcal{O}(\min\{E^{1/3}, E^{2/3} \cdot \theta\})$, even if E and θ are known.

Our main algorithmic contribution is a policy which achieves this lower bound (up to log factors) without knowing the prediction accuracy:

Theorem 3 (Informal). There exists a policy which achieves regret $\mathcal{O}(\min\{E^{1/3}, E^{2/3} \cdot \theta\})$, knowing E , and without knowing θ .

Finally, since our policy relies on knowledge of the variation parameter E , the remaining question is whether the same regret is achievable if E and θ are unknown. We show that in fact predictions cannot be incorporated in any meaningful way in this case:

Proposition 12 (Informal). If E and θ are unknown, then no policy can achieve regret better than $\mathcal{O}(\max\{E^{1/3}, E^{2/3} \cdot \theta\})$ for all $E, \theta \geq 0$.

Our theoretical results are summarized in the Table 4.1. Each entry has a corresponding policy that achieves the stated regret, along with a matching lower bound.

	Without predictions	With predictions of unknown accuracy
Known variation	$\$1 \cdot 10^4$	$\$1 \cdot \min_{f \in \mathcal{F}} \mathbb{E}[\text{cost}]$
Unknown variation	$\$1 \cdot 10^4$	$\$1 \cdot \max_{f \in \mathcal{F}} \mathbb{E}[\text{cost}]$

Table 4.1: Summary of main theoretical results. Each entry has a corresponding policy that achieves the stated regret, along with a matching lower bound.

	No Prediction	Prediction	Our Policy
Upper store	\$28,303	\$14,454	\$14,454
Lower store	\$23,460	\$35,600	\$23,899

Table 4.2: Continuation of Figure 4.1: costs incurred by an optimal policy which makes no use of predictions, a policy which relies entirely on predictions, and our policy.

3. Empirical Results: Finally, we demonstrate the practical value of our model (namely the Nonstationary Newsvendor with Predictions) and our policy via empirical results on three real-world datasets that span our motivating applications above: daily web traffic for Wikipedia.com (of various languages), daily foot traffic across the Rossmann store chain, and daily visitors at a certain Japanese restaurant. These datasets together contain over one thousand individual time-series on which we generate predictions of varying quality, using four different popular forecasting and machine learning algorithms. We apply our policy, and compare its performance against the two most-natural baseline policies: our optimal policy without predictions, and the simple policy which always utilizes the predictions (these correspond to the two existing approaches described previously). A snapshot of our results, for the Rossmann stores depicted in Figure 4.1, is given in Table 4.2.

More generally, on any given experimental instance (i.e. a time-series and a set of predictions), the minimum (maximum) of the costs incurred by these two baselines can be viewed as the best (worst) we can hope for. Thus we measure performance in terms of the proportion of the gap between these two costs incurred by our policy, so if this optimality gap is close to 0, our policy performs almost as good as the better one of the two baselines. Note that randomly selecting between the two baseline policies yields an (expected) optimality gap of 0.5. We find that in the Rossmann dataset, the average optimality gap is 0.26 when the predictions are accurate, and 0.28 when the predictions are inaccurate. In the Wikipedia dataset, the average optimality gap is 0.40 when the predictions are accurate, and 0.07 when the

predictions are inaccurate. In the Restaurant dataset, the average optimality gap is 0.10 when the predictions are accurate, and 0.39 when the predictions are inaccurate. This demonstrates that our policy performs well, irrespective of the quality of the predictions.

Literature Review

The earliest works on the newsvendor model assume that the demand distribution is known [11, 154]. This has since been relaxed, the resulting approaches being divided into parametric and nonparametric ones. Among parametric approaches much work is Bayesian, where a prior distribution is assumed over the parameters, applied the Bayesian approach to inventory models, and later this was studied in many works [15, 85, 91, 118]. [117] introduced another parametric approach called operational statistics which, unlike the Bayesian approach, does not assume any prior knowledge on the parameter values, instead using past demand observations to directly estimate the optimal ordering quantity.

Nonparametric approaches have been developed in recent years. The first example is the sample average approximation (SAA) method, first proposed by [157]. [110] applied SAA to the newsvendor problem, and [1] improved significantly upon their bounds. Other non-parametric approaches include stochastic gradient descent algorithms [1, 84, 104] and the concave adaptive value estimation (CAVE) method [71, 145]. With the development of machine learning, [24] and [143] propose machine learning/deep learning algorithms using demand features and historical data.

All the above studies treat the newsvendor in a static environment. There are two common approaches to nonstationarity. The first is to model (stochastically) the nonstationarity and utilize past demand observations according to the model. One common way is to model the nonstationarity as a Markov chain. For example, [162] applied this idea to inventory management, and [46] applied this idea to revenue management. Another approach is to bound the nonstationarity via a variation budget, which has been applied to stochastic optimization [37], dynamic pricing [95], multi-armed bandit [36], newsvendor problem [96], among others. Some of these works are applicable in the sense that our problem can be mapped to their settings (e.g. multi-armed bandit such as [52, 92, 122]), but these connections do not appear to be fruitful. In particular, the multi-armed bandit papers cited above typically consider a limited-feedback setting rather than the full-feedback setting

explored in this work. Related to feedback, while our study provides a complete characterization of the regret behavior for the nonstationary newsvendor problem with uncensored demand, practical applications often involve censored demand. The nonstationary newsvendor problem under censored demand is an interesting direction for future research.

Beyond the bandit literature, it is worth mentioning recent work on online convex optimization (OCO) with limited nonstationarity. When the level of nonstationarity is known, the standard first-order OCO algorithms can be modified with carefully chosen restarts and updating rules [178, [51]]. There are also recent works that concern unknown nonstationarity, such as [16], [17], and [82]. Finally, as mentioned before, [96] is particularly relevant, so we delay a careful comparison to Sections 4.2 and 4.3.

The second common practice is to use predictions. A recent line of work looks to help decision-making by incorporating predictions into online optimization problems such as revenue optimization [20, 137], caching [123, 149], online scheduling [106], and the secretary problem [62]. In this paper we combine the nonstationarity framework and the prediction framework on the newsvendor problem.

Finally, most previous works involving algorithms with predictions analyzed algorithms' performances using competitive analysis (e.g., [20, 89, 124]) and obtained optimal consistency-robustness trade-offs, where consistency is an algorithm's competitive ratio when the prediction is accurate, and robustness is the competitive ratio regardless of the prediction's accuracy. However, competitive ratio transfers to a regret bound that is linear in n . In contrast, we do regret analysis under this framework and design an algorithm that has near-optimal worst-case regret without knowing the prediction quality. Other papers with regret analyses under the prediction model include [37] (revenue optimization in auctions), [77] (Thompson sampling), [81] (constrained online two-stage stochastic optimization), and [41] (online resource allocation).

4.2 Model: The Nonstationary Newsvendor (without Predictions)

We begin this section with a formal description of the Nonstationary Newsvendor, along with a comparison to the problem of the same name [66]. Consider a sequence of newsvendor problems over time periods labeled $t = 1, \dots, T$. At the beginning of each time period t , the decision-maker selects a quantity q_t , where

\mathcal{C} is a fixed subset of \mathbb{R}^+ bounded above by a quantity we denote as c_{\max} .² Then the period's demand D_C is drawn from an (unknown) demand distribution D_C which depends on the time period t . These demand distributions are independent over time. Finally a cost is incurred. Specifically, there is a (known) per-unit underage cost $1 - c$ and a (known) per-unit overage cost c so that the total cost is equal to

$$G = \int_0^{\infty} (c - (c - 1)x) f(x) dx$$

where $G = \int_0^{\infty} x f(x) dx$. The decision-maker observes the realized demand D_C and thus the cost. Note that requiring \mathcal{C} does not impose any restriction on modeling, since \mathcal{C} could simply be selected to be \mathbb{R}^+ (as in much of the literature). In fact, introducing \mathcal{C} allows for modeling important practical concerns such as batched inventory or even simply the integrality of physical items. As we will discuss momentarily, this is a non-trivial concern insofar as theoretical guarantees are concerned.

To complete our description of the Nonstationary Newsvendor, we will need to (a) impose a few assumptions on the demand distributions, and then (b) describe how nonstationarity is quantified. These are, respectively, the subjects of the following two subsections.

Demand Distributions

We will assume that the demand distributions come from a known, parameterized family of distributions D :

Assumption 3. Every demand distribution D_C comes from a family of distributions D satisfying the following:

- (a) $D = f(D; \theta)$ where $\theta \in [\theta_{\min}, \theta_{\max}]$ that is D is parameterized by a scalar taking values in some bounded interval.
- (b) Each distribution $D \in D$ is sub-Gaussian.

Assumption 3 is fairly minimal. Parsing it in reverse: the sub-Gaussianity in part (b) allows for many commonly-used variables, such as the Gaussian distribution

²All of our results carry through if \mathcal{C} is allowed to depend on C .

³The demand is not censored here, as is the case in all of the motivating examples in the introduction. The censored version of our problem is an interesting, but separate subject.

⁴A random variable X is sub-Gaussian with sub-Gaussian norm $\|X\|_{\psi_2}$ if $\mathbb{P}(|X| > t) \leq 2 \exp(-G^2 t^2 / k^2)$ for all $G > 0$. For sub-Gaussian variables, we have $\mathbb{E}|X|^j \leq k^j$.

and any bounded random variable, while letting us eventually apply Hoeffding-type concentration bounds. Part (a) is particularly minimal at the moment, represents an arbitrary parameterization \mathcal{D} , but will become meaningful when combined with Assumption 4. The choice of the symbol μ might suggest that μ represents the mean of \mathcal{D} , and indeed this is what we will assume from here on. But it should be emphasized that our taking $\mu = \mathbb{E}[\mathcal{D}]$ is strictly for notational convenience (because we will frequently need to refer to the means of these distributions): if we were any other parameterization \mathcal{D} , we could simply define a mapping from μ to the mean values.

Now define $\mathcal{C}(\mu) = \mathbb{E}[\mathcal{C}(\mu)]$ be the expected newsvendor cost when selecting quantity Q , given underage/overage costs c_u and c_o , and demand distribution \mathcal{D}

$$\mathcal{C}(\mu) = \mathbb{E}_{\mathcal{D}}[\mathcal{C}(\mu)] = \int_{-\infty}^{\infty} \mathcal{C}(Q) f_{\mathcal{D}}(Q) dQ$$

The critical assumption, with respect to the parameterization in Assumption 3(a), is that the expected cost is well-behaved as a function of μ :

Assumption 4. For every $\mu_1, \mu_2 \in [\mu_{\min}, \mu_{\max}]$ and $Q \in \mathbb{R}$ the function $\mathcal{C}(\mu)$ is Lipschitz on its domain $[\mu_{\min}, \mu_{\max}]$, i.e. there exists $L \in \mathbb{R}$ such that for every $\mu_1, \mu_2 \in [\mu_{\min}, \mu_{\max}]$, we have

$$|\mathcal{C}(\mu_1) - \mathcal{C}(\mu_2)| \leq L |\mu_1 - \mu_2|$$

Note that in the above description, the Lipschitz constant may depend on c_u and c_o but by continuity, there exists a single L so that the above holds for all μ_1, μ_2 simultaneously.

Some useful examples of families \mathcal{D} satisfying Assumptions 3 and 4 are the following:⁵

1. $\mathcal{D} = \mathcal{N}(\mu, \sigma^2)$, the family of normal distributions with fixed variance σ^2 . In this case, $L = \frac{c_u + c_o}{\sigma}$. A relaxation is that the variances may vary (continuously) with μ .
2. $\mathcal{D} = \mathcal{D}_n$, where n is any mean-zero, sub-Gaussian variable.

⁵Unfortunately, Assumption 4 is not guaranteed to hold. For example, for the family of distributions

$$\mathcal{D} = \left\{ \text{Bernoulli}(\mu) \mid \mu \in [\mu_{\min}, \mu_{\max}] \right\}$$

the function $\mathcal{C}(\mu) = \mathbb{E}[\mathcal{C}(\mu)]$ is discontinuous (and thus not Lipschitz) at $\mu = 1$.

3. The Poisson distribution is frequently used to model demand (since arrivals are often modeled as a Poisson process). While the Poisson distribution is not sub-Gaussian, any reasonable truncation satisfies our assumptions. For example, $D \cdot \min\{\text{Poisson}(\lambda), \lambda\}$, for some constant λ . Here, λ can be taken to be large enough so that the truncation happens with small probability (in fact, this probability is $\leq 14^{-\lambda}$).

To understand the reasoning behind Assumption 4, consider the problem faced at some time C . The optimal choice for the decision-maker here is

$$(4.1) \quad \hat{c}_C = \arg\min_{c \in D} \mathbb{E}[\text{cost}(c, \mu_C)]$$

where μ_C is the mean of c_C (i.e. $c_C \sim D(\mu_C)$), and $\text{cost}(c, \mu) = \lambda c - \mu c$. To simplify the notation, since μ_C is unknown, it is likely that some $\hat{c}_C < \hat{c}_C^*$ will ultimately be selected, and we could measure the sub-optimality of this decision (i.e. regret, to be defined soon): $\mathbb{E}[\text{cost}(\hat{c}_C, \mu_C)] - \mathbb{E}[\text{cost}(\hat{c}_C^*, \mu_C)]$. It would be natural then to try to characterize this suboptimality as a function of μ_C , but in fact all of the algorithms we will consider work by making an estimate $\hat{\mu}_C$ of μ_C and then selecting $\hat{c}_C = \arg\min_{c \in D} \mathbb{E}[\text{cost}(c, \hat{\mu}_C)]$. So motivated, the purpose of Assumption 4 is to allow us to translate error in our estimate of μ_C to (excess) costs. The following structural lemma makes this precise, and will be used throughout the paper.

Lemma 2. Fix any ϵ and δ (we will suppress them from the notation). For any $\mu_1, \mu_2 \in D$, let $\hat{c}_1 = \arg\min_{c \in D} \mathbb{E}[\text{cost}(c, \mu_1)]$ and $\hat{c}_2 = \arg\min_{c \in D} \mathbb{E}[\text{cost}(c, \mu_2)]$. Then we have

$$\mathbb{E}[\text{cost}(\hat{c}_1, \mu_2)] - \mathbb{E}[\text{cost}(\hat{c}_2, \mu_2)] \leq \epsilon + \delta \|\mu_1 - \mu_2\|.$$

Lemma 2 states that estimation error of the mean translates linearly to excess cost. The proof of Lemma 2 appears in Appendix C.1.

Aside: Comparison to [96]: The main component in describing the Nonstationary Newsvendor is defining a proper quantification of nonstationarity. Before doing so, we delineate the modeling differences between our Nonstationary Newsvendor and that of [96]. There are two primary differences:

1. The demand distributions c_C in [96] are assumed to be of the form $c_C = \mu_C + \eta_C$ where μ_C is the mean of c_C that drifts across time and η_C is the noise distribution

⁶As a sanity check, the classical result for the newsvendor problem [54] states that if $\mu = R$, then \hat{c} is the $(1-\alpha)$ -th quantile of c .

that is i.i.d., continuous, and bounded. Effectively, the demand distributions fall into a non-parametric family of distributions with the same shape. In contrast, our demand distributions fall into a parametric family of distributions, though not necessarily of the same shape.

2. Our set of allowed order quantities is bounded, but otherwise arbitrary. In particular, it need not contain the optimal unconstrained order quantity $\arg\min_{x \in \mathbb{R}^+} f(x)$ for each μ (or any μ , for that matter). [96] assume $\mathcal{X} = \mathbb{R}^+$.⁷

Besides the practical reasons why discrete quantities arise in practice (non-divisible items, batched inventory, etc.), the primary consequence of either of the two differences above is that they preclude a critical lemma used in (and in fact by [110]) which states that $\mathbb{E}[\sum_{t=1}^T \ell_t(x_t^*) - \min_{x \in \mathcal{X}} \sum_{t=1}^T \ell_t(x)] \leq C \sqrt{\log T}$ (as defined in our Lemma 2) scales as $\sqrt{\log T}$. This scaling does not necessarily hold when either the demand distribution or \mathcal{X} is discrete. These relaxations in assumptions yield different lower bounds in the worst-case regret from [96], which we will discuss in detail later.

Demand Variation

Just as in [96] (and [95] before that), we measure the level of nonstationarity via a deterministic analogue of quadratic variation for the sequence of means $\{\mu_t\}_{t=1}^T$.

Specifically, define a partition of the time horizon $\mathcal{C} = \{C_0, \dots, C_P\}$ to be any subset of time periods $\{1, \dots, T\}$ where $1 = C_0 < \dots < C_P = T$. Here the subset can have any size between 1 and T , i.e. $0 < P \leq T$. Then for any sequence of means $\mu = \{\mu_1, \dots, \mu_T\}$, its demand variation is

$$(4.2) \quad \mathbb{V}(\mu) = \max_{\mathcal{C}} \max_{P \geq 1} \sum_{i=1}^P (\mu_{C_i} - \mu_{C_{i-1}})^2$$

where \mathcal{P} is the set of all partitions.

To motivate the use of partitions in the definition of $\mathbb{V}(\mu)$, it is worth contrasting with a measure that may feel more natural, namely the sum of squared differences (SSD) between consecutive terms $\sum_{t=2}^T (\mu_t - \mu_{t-1})^2$, which corresponds to taking the densest possible partition $\mathcal{C} = \{1, 2, \dots, T\}$. The maximum in the definition of $\mathbb{V}(\mu)$ is not necessarily achieved by selecting the densest possible partition, but rather by setting $\mathcal{C} = \{C_0, \dots, C_P\}$ to be the periods when the sequence $\mu = \{\mu_1, \dots, \mu_T\}$ changes direction. Thus, the demand variation penalizes trends, or consecutive increases/decreases, more so than the SSD. For example, the mean sequences $\mu = \{1, 2, 3, 4, 5\}$ and $\mu = \{1, 0, 1, 0, 1\}$

⁷While not stated explicitly, the results in [96] only require \mathcal{X} to contain points arbitrarily close to every optimal unconstrained order quantity.

respective variations $\sigma_{-1} = 15 \cdot 10^2 = 16$ and $\sigma_{-2} = 1^2, 1^2, 1^2, 1^2, 1^2 = 5$, despite having identical SSDs.

All of our theoretical guarantees (upper and lower bounds) will be parameterized by ϵ . This quantity of course depends on ϵ and so it is natural to allow ϵ to grow. It will turn out that the most natural parameterization of this growth is via what we will simply call the variation parameter $E \gg 0$, such that $\epsilon = \epsilon(E)$, where ϵ is some constant (which we take to be equal to one from here on). We denote the set of demand distribution sequences $\{f_{1-\dots}^g\}$ whose means $\mu = f_{1-\dots}^g$ satisfy $\epsilon(E)$ as

$$D^1 E^0 = \{f_{1-\dots}^g : \mu \in D \text{ for all } C \text{ and } \epsilon(E) \leq \epsilon\}$$

In the next section, we will show via a minimax lower bound that non-trivial guarantees are only achievable when $E \gg 1$, and provide an algorithm which achieves the same bound.

Aside: Time-Series Modeling: At this point, we have fully described our model for the Nonstationary Newsvendor. All that remains is to define our performance metric, which we will do in the next subsection. We conclude this subsection with an important practical consideration with respect to time-series models and our variation parameter.

Consider, as an example, the following class of time-series models:

$$(4.3) \quad \mathcal{Z}_C = \{^1C^0, (^1C^0, n^0\}$$

Here, $^1C^0$ represents a deterministic (and usually simple, e.g. linear) function representing some notion of trend, $^1C^0$ represents a deterministic, periodic function representing some notion of seasonality. Finally, all stochastic behavior is captured by the random variables which are assumed to be independent and mean-zero. This time-series model is classic, and yet drives forecasting algorithms (e.g. exponential smoothing) which are still competitive in modern forecasting competitions [125].

The above model raises an important practical issue: if there exists any (non-trivial) trend $^1C^0$ or seasonality $^1C^0$, then the demand variation of the sequence of means $\mu_C = \{^1C^0, (^1C^0\}$ would scale at least as $\epsilon(E) = 1$ and no meaningful guarantee will be achievable. Our main observation is that time-series effects like trend and seasonality are easily detected and estimated, so that in any practical

setting, estimates $\hat{\mu}^0$ and $\hat{\sigma}^0$ should be available, and used to de-trend and de-seasonalize the data. Concretely, the Nonstationary Newsvendor would take place on the sequence

$$\tilde{z}_C = \tilde{z}_1 C^0 + \hat{\mu}^0 C^0 \quad (\hat{\mu}^0 = \mu^1 C^0 + \hat{\mu}^{00}, \mu^1 C^0 + \hat{\mu}^{00}, \mu^0)$$

The resulting sequence of means $\mu^1 C^0 + \hat{\mu}^{00}, \mu^1 C^0 + \hat{\mu}^{00}$ does not stem from the trend and seasonality, but rather the error in estimating the trend and seasonality. It is this error that is assumed to be nonstationary, but with reasonable variation parameter.

Performance Metric: Regret

We conclude this section by formally defining our performance metric for any policy. A policy is simply a sequence of mappings $\pi = \{c_1, \dots, c_T\}$, where each c_C is a mapping from $\mathcal{D}_1, \dots, \mathcal{D}_T$ to an order quantity \mathbb{R}_+ & at time C (by convention, c_1 is a constant function). We measure the performance of a policy by its regret. Fix a sequence of demand distributions $J = \{J_1, \dots, J_T\}$. Following the earlier notation from Eq.(4.1), the regret incurred by a policy which selects order quantities

$$\pi = \{c_1, \dots, c_T\} \quad \mathbb{R}^+ \quad \#$$

$$E_J^c \sum_{C=1}^T c^1 c^C - c^1 c^C -$$

where the expectation is with respect to the randomness of the realized demands. Recall that the demand distributions are independent as defined in (4.1) depends only on c_C . In words, the regret measures the difference between the (expected) total cost incurred by the policy and that of a clairvoyant that knows the underlying demand distributions $J = \{J_1, \dots, J_T\}$.⁹

We will be concerned with the worst-case regret of a policy across families of instances (i.e. sequences of demand distributions) controlled by the variation parameter E :

$$R^{c_1}{}^0 = \sup_{J \in \mathcal{D}^1 E^0} E_J^c \sum_{C=1}^T c^1 c^C - c^1 c^C \cdot \#$$

Note that if the worst-case regret $R^{c_1}{}^0$ of some policy is sublinear in J , then that policy is essentially cost-optimal on average as J goes to infinity. In the next section,

⁸Note that we are not considering randomized policies here, but all of our theoretical results (the lower bounds, in particular) hold even when randomization is allowed.

⁹Note that this is different from a clairvoyant that knows the realized demands J_1, \dots, J_T . Such a clairvoyant would incur zero cost.

we will prove a lower bound on the achievable across all policies, and describe an algorithm which achieves this lower bound.

4.3 Solution to the Nonstationary Newsvendor (without Predictions)

This section contains a complete solution (i.e. matching lower and upper bounds on regret) to the Nonstationary Newsvendor. We begin with the lower bound:

Proposition 10 (Lower Bound: Nonstationary Newsvendor). For any variation parameter $E \geq 0$, and any policy π (which may depend on the knowledge of E), we have

$$R(\pi) \geq \frac{1}{2} E^2 - \frac{1}{4} E^4$$

where $\frac{1}{2} \geq 0$ is a universal constant.

Proposition 10 is a corollary of a more general lower bound (Proposition 11 in the next section) – it will turn out the Nonstationary Newsvendor is a special case of the Nonstationary Newsvendor with Predictions – so the proof is omitted. Proposition 10 states that the regret of any policy is at least $\frac{1}{2} E^2 - \frac{1}{4} E^4$. It is useful to contrast this with two existing results:

1. Stationary Newsvendor: In the special case of i.i.d. demand, it is known that the optimal achievable regret is $\frac{1}{2} E^2$. Example 1 of [35] demonstrates the lower bound, and the SAA method of [10, 11] achieves the upper bound. This point might appear to be incompatible with our result, which states a lower bound of $\frac{1}{2} E^2 - \frac{1}{4} E^4$ when $E = Q$ but in fact the case $E = 0$ is more general than i.i.d. demand since it allows $E = 0 = \frac{1}{2} E^2$ demand variation, while i.i.d. demand amounts to zero demand variation. Indeed, our proof of Proposition 10, for the special case $E = Q$ utilizes instances for which the demand distribution is allowed to change $\frac{1}{2}$ times (by an amount of $\frac{1}{4}$, resulting in $\frac{1}{2} E^2$ variation).

As an aside, this discussion raises a natural question: is the disconnect here between stationary (i.i.d.) demand and variation parameter E a consequence of our use quadratic variation, and would the same disconnect arise for other measures of demand variation? In Appendix C.5, we answer both questions in the affirmative by showing that if the exponent in the demand variation (Equation (4.2)) is instead some θ , then Proposition 10 generalizes to a lower bound of $\frac{1}{2} E^2 - \frac{1}{4} E^{2\theta}$. Thus, for any $\theta \geq 0$, the

case of variation parameter $E = 0$ is meaningfully more general than stationary demand.

2. Continuous Newsvendor: A similar story plays out in the setting of [96], which recall (among other key differences with our model, as described in Section 4.2) requires the additional assumption that both the demand distributions and the possible order quantities be continuous. [96] show an optimal achievable regret of $O(\sqrt{11, E^{0.2}})$, which can be contrasted with the stationary (i.i.d.) setting for which a $(\log)^0$ lower bound exists [5]. The following table summarizes these lower bounds:

	Continuous	General
Stationary (i.i.d.)	(\log)	$)^{1.2}$
Nonstationary	$)^{11, E^{0.2}}$	$)^{13, E^{0.4}}$

In the next two subsections, we will first analyze a simple algorithm which achieves the lower bound of Proposition 10 when the variation parameter E is known, and then use this as a building block for an algorithm which achieves the same bound when E is unknown.

Upper Bound with Known Variation Parameter E

If we assume that E is known, then designing a policy which achieves regret matching Proposition 10 is fairly straightforward. In fact, a simple policy based on averaging a fixed number of past demand observations does the job (use the same policy). That policy, which we call the Fixed-Time-Window Policy is defined in Algorithm 4.

The Fixed-Time-Window Policy uses a carefully-selected window size that is on the order of $\sqrt{11 E^{0.2}}$. At each time period t , it constructs an estimate of the mean by averaging the observed demands from the previous periods, and then selects the optimal order quantity corresponding to μ . Note that Algorithm 4 also includes a scaling constant $\hat{\lambda}$ this should be thought of as a practical tuning parameter, but for the coming theoretical result, it can be chosen arbitrarily (e.g. $\hat{\lambda} = 1$ suffices).

The following result bounds the worst-case regret of the Fixed-Time-Window Policy:

Lemma 3 (Upper Bound: Nonstationary Newsvendor with Known E). Fix any variation parameter $E \geq 0$. The Fixed-Time-Window Policy achieves

Algorithm 4: Fixed-Time-Window Policy

 Inputs: variation parameter $E \geq 0$ and scaling constant $\lambda > 0$

 Initialization: $\lambda_C = \frac{1}{C}$

 for $C = 1, 2, \dots$ do

 | select \mathcal{C} & arbitrarily

end

 for $C = 1, 2, \dots$ do

 | $\lambda_C = \frac{1}{C}$

 | if $\lambda_C \geq \frac{1}{\max\{g^1, X\}}$ then

 | | round λ_C to the nearest value in $\{\frac{1}{3}, \frac{1}{4}, \dots, \frac{1}{\max\{g^1, X\}}\}$

| end

 | $\mathcal{C} = \arg \min_{\mathcal{C}} \lambda_C$

 end

worst-case regret

$$R^{\text{Fixed-Time-Window}}(T) \leq \frac{1}{\lambda} \left(\frac{1}{\lambda} + \frac{1}{\lambda} \right) = \frac{2}{\lambda}$$

where

$$\frac{1}{\lambda} = \max\left\{ \frac{1}{\max\{g^1, X\}}, \frac{1}{\max\{g^0, X\}} \right\} = \frac{1}{\max\{g^1, X, g^0, X\}} = \frac{1}{\max\{g^1, g^0, X\}}$$

 and $X = \sup_{D \in \mathcal{D}} \max_{k \in \mathcal{K}} D_k$.

As promised, Lemma 3 shows that the Fixed-Time-Window Policy achieves regret that matches the lower bound in Proposition 10. Its proof can be found in Appendix C.2, and amounts to bounding the estimation error incurred by demand noise (which is worse for smaller time windows) and demand mean variation (which is worse for larger time windows). The exact time window used in the policy comes from balancing these two sources of error.

Upper Bound with Unknown Variation Parameter E

The lower bound in Proposition 10 holds for policies that know E . Naturally, it also holds for policies that do not know E , but an unanswered question at the moment is whether (a) the lower bound should be even larger when E is unknown, or (b) there exists a policy that matches Proposition 10 without knowing E . We show here that case (b) holds by constructing such a policy. Our policy, which we call the Shrinking-Time-Window Policy (Algorithm 5), at a high level uses the Fixed-Time-Window Policy with the smallest variation parameter that is consistent with the demand observed so far. In more detail:

¹⁰As a natural comparison to [96], they do not consider the unknown E setting.

1. It begins with a discrete set of candidate variation parameters E_1, \dots, E_9

$$(4.4) \quad E_9 = 1, \frac{1}{\log 2}, \dots, \frac{1}{\log 9} \quad 9 = 1 - \dots - :$$

where: is chosen so that $E_1 \leq \dots \leq E_9$. V is defined specifically so that the variation parameters are increasing $(E_1 \leq \dots \leq E_9)$, and so that it discretizes the interval $[0, 1]$ at a sufficiently fine granularity.

2. At any time period C , there is a current candidate parameter E_8 (initialized to be E_1 at $C = 1$) that is assumed to be the true variation parameter and so the corresponding Fixed-Time-Window Policy is applied: a time window of

$$(4.5) \quad \tau_8 = d \wedge \frac{1}{E_8^{0.2}} e$$

is used, and an estimate of μ is made:

$$(4.6) \quad \hat{\mu}_C = \frac{1}{\tau_8} \sum_{B=C-\tau_8}^{C-1} x_B \quad \text{rounded to the nearest value in } \{ \frac{1}{E_1}, \dots, \frac{1}{E_9} \}$$

3. The index of the current candidate parameter E_8 is incremented at any period in which the policy gathers sufficient evidence that $\mu \geq E_8$. This is possible due to the following observation: if $\mu \geq E_8$, then by Lemma 3 we have that for any $E_9 \geq E_8$, the regret incurred by the Fixed-Time-Window Policy corresponding to E_9 is $O(\frac{1}{E_9^{0.4}})$, and thus the cumulative difference between the estimated mean demand $(\hat{\mu}_C - \mu)$ cannot exceed $O(\frac{1}{E_8^{0.4}})$. Thus if this is observed for some $E_9 \geq E_8$, we can conclude that $\mu \geq E_8$, and 8 is incremented.

This policy's regret matches (up to log factors) the lower bound in Proposition 10:

Theorem 2 (Upper Bound: Nonstationary Newsvendor with Unknown μ). For any variation parameter $\mu \in [0, 1]$, the Shrinking-Time-Window Policy achieves worst-case regret

$$R^{C, \text{shrinking}}(\mu) = O(\frac{1}{\mu^{0.4}} \log^{5.2} \frac{1}{\mu})$$

where $C = \max\{1, \frac{1}{\mu}\}$, $\log^{5.2} \frac{1}{\mu} = 124 \cdot 1.4 \cdot W$, $W = \frac{1}{\mu^{0.4}}$, Lemma 3 and Lemma 3 is the constant in Lemma 3.

The proof of Theorem 2 can be found in Appendix C.3, but at a high level works as follows:

Proof Sketch of Theorem 2. First note that the total regret incurred during the first $\frac{1}{\mu^{0.4}}$ time periods is at most $O(\frac{1}{\mu^{0.4}})$. After that, the total regret incurred during

Algorithm 5: Shrinking-Time-Window Policy

Inputs: scaling constants $\alpha, \beta, \gamma, \delta, \epsilon, \zeta, \eta, \theta, \iota, \kappa, \lambda, \mu, \nu, \xi, \omicron, \pi, \rho, \sigma, \tau, \upsilon, \phi, \chi, \psi, \omega$, and W sufficiently large (Equation (C.2) in Appendix C.3)

Initialization: Set $V = \frac{1}{\epsilon} \sum_{j \in \mathcal{B}} \lambda_j^{\alpha} \lambda_j^{\beta}$ and $f = \frac{1}{\epsilon} \sum_{j \in \mathcal{B}} \lambda_j^{\alpha} \lambda_j^{\beta}$ according to Equations (4.4) and (4.5)

for $C = 1, 2, 3, \dots$ do

 | select \mathcal{C}_C & arbitrarily;

end

Initialize $\delta = 1$ and $G_C = \frac{1}{C}$

for $C = 1, 2, 3, \dots$ do

 | if $\frac{1}{C} \sum_{j \in \mathcal{B}} \lambda_j^{\alpha} \lambda_j^{\beta} \geq \frac{1}{W \log C} \sum_{j \in \mathcal{B}} \lambda_j^{\alpha} \lambda_j^{\beta}$ for some $j \in \mathcal{B}$ then

 | $\delta = \frac{1}{C}$

 | $G_C = C$

 end

 | $\mathcal{C}_C = \arg \min_{\mathcal{C} \subseteq \mathcal{B}} \sum_{j \in \mathcal{C}} \lambda_j^{\alpha} \lambda_j^{\beta}$

end

the time periods in which the if condition in Algorithm 5 is triggered is at most $\frac{1}{\delta}$, where: δ is the number of candidate variation parameters defined in Equation (4.4), and: $\log^2 C$. Thus, it suffices to bound the total regret incurred between successive triggerings of the if condition. As a final reduction before proceeding, consider the smallest candidate variation parameter that is at least defined to be the smallest index such that $\lambda_j \geq \frac{1}{\log C}$. Because $\lambda_j \geq \frac{1}{\log C}$, we have that λ_j^E is a constant multiple away from λ_j^E . Thus, it will suffice to bound the total regret by $\frac{1}{\delta} \sum_{j \in \mathcal{B}} \lambda_j^{\alpha} \lambda_j^{\beta} \log^{5 \cdot 2} C$.

To do this, we first show that (with high probability) throughout the algorithm, the running index δ never exceeds $\frac{1}{\delta}$. To see this, consider the following steps:

1. For every $j \in \mathcal{B}$, because $\lambda_j \geq \frac{1}{\log C}$, by Lemma 3 the Fixed-Time-Window Policy corresponding to the window size $\frac{1}{\log C}$ has worst-case regret $\frac{1}{\delta} \sum_{j \in \mathcal{B}} \lambda_j^{\alpha} \lambda_j^{\beta}$. In addition, we can show with high probability (via Hoeffding's inequality) that $\frac{1}{C} \sum_{j \in \mathcal{B}} \lambda_j^{\alpha} \lambda_j^{\beta} \leq \frac{1}{W \log C} \sum_{j \in \mathcal{B}} \lambda_j^{\alpha} \lambda_j^{\beta}$ for every C . We may assume this event occurs from now on.
2. For the sake of contradiction, suppose δ exceeds $\frac{1}{\delta}$ at some period C or equivalently, there is a period C in which $\delta > \frac{1}{\delta}$, and the if condition is triggered

with some $9 \leq 8$. Then

$$\begin{aligned}
 & \tilde{O}_{B=\emptyset} j \lambda_B^8 \wedge \lambda_{B_j}^9 \stackrel{10^0}{=} \tilde{O}_{B=\emptyset} j \lambda_B^8 \wedge \lambda_{B_j}, \quad \tilde{O}_{B=\emptyset} j \lambda_B^9 \wedge \lambda_{B_j} \\
 & \stackrel{11^0}{=} \tilde{O}_{B=\emptyset,1} j \lambda_B^8 \wedge \lambda_{B_j}, \quad \tilde{O}_{B=\emptyset,1} j \lambda_B^9 \wedge \lambda_{B_j} \\
 & \stackrel{12^0}{=} \frac{\rho}{W \log}, \rho_{\bar{\lambda}} \stackrel{13, E_8^{0.4}}{=} \frac{\rho}{W \log}, \rho_{\bar{\lambda}} \stackrel{13, E_9^{0.4}}{=} \\
 & \stackrel{13^0}{=} \frac{\rho}{2 W \log}, \rho_{\bar{\lambda}} \stackrel{13, E_9^{0.4}}{=}
 \end{aligned}$$

where 10^0 is the triangle inequality and 11^0 holds because $\epsilon_j \leq \epsilon$ and $\epsilon = \epsilon \wedge \epsilon^{1.2}$; since $9 \leq 8$, by our assumption in the previous step we get 12^0 ; 13^0 follows since $E_9 \leq E_8$. But this directly contradicts our assumption that the if condition is triggered at period C. Therefore δ never exceeds ϵ .

Now suppose two consecutive if conditions occur at times C^0 and C^1 . Between these periods, we may apply the negation of the if condition for any δ and since δ never exceeds ϵ , we specifically can take $9 = \epsilon$. This yields

$$\tilde{O}_{B=\emptyset,1} j \lambda_B^8 \wedge \lambda_{B_j} \leq \frac{\rho}{2 W \log}, \rho_{\bar{\lambda}} \stackrel{13, E^{0.4}}{=}$$

Then we have

$$\begin{aligned}
 & \tilde{O}_{B=\emptyset,1} j \lambda_B^8 \wedge \lambda_{B_j} \stackrel{10^0}{=} \tilde{O}_{B=\emptyset,1} j \lambda_B \wedge \lambda_{B_j}, \quad \tilde{O}_{B=\emptyset,1} j \lambda_B^8 \wedge \lambda_{B_j} \\
 & \stackrel{11^0}{=} \tilde{O}_{B=\emptyset,1} j \lambda_B \wedge \lambda_{B_j}, \quad \tilde{O}_{B=\emptyset,1} j \lambda_B^8 \wedge \lambda_{B_j} \\
 & \stackrel{12^0}{=} \frac{\rho}{W \log}, \rho_{\bar{\lambda}} \stackrel{13, E^{0.4}}{=} \frac{\rho}{2 W \log}, \rho_{\bar{\lambda}} \stackrel{13, E^{0.4}}{=} \\
 & = 3 \frac{\rho}{W \log}, \rho_{\bar{\lambda}} \stackrel{13, E^{0.4}}{=}
 \end{aligned}$$

where 10^0 is the triangle inequality and 11^0 holds because $\epsilon_j \leq \epsilon$ and $\epsilon = \epsilon \wedge \epsilon^{1.2}$; the first part of 12^0 follows by our high-probability assumption, and the second part of 12^0 follows since the if condition is not triggered between time C^0 and time C^1 ; Therefore between two consecutive if conditions, the worst-case regret incurred is $(\epsilon \frac{\rho}{W \log})^0$. Because the if condition can happen at most \log^2 times, the total worst-case regret of the Shrinking-Time-Window Policy is $(\epsilon \log^2)^0 = (\epsilon \log^2)^0$.

This concludes our discussion of the Nonstationary Newsvendor. In the next section, we turn to the second subject of this paper, which is the same problem with predictions.

4.4 The Nonstationary Newsvendor with Predictions

As described in the introduction, it is likely that when the Nonstationary Newsvendor is faced practice, some notion of a prediction of future demand will be made. Such predictions can come from a diverse set of sources ranging from simple human judgement, to forecasting algorithms built on previous demand data, to more-sophisticated machine learning algorithms trained on feature information. The process of sourcing or constructing such predictions is orthogonal to our work. Instead, we treat these predictions as given to us endogenously (and in particular, we make no assumption on the accuracy of these predictions), and attempt to use these predictions optimally.

Model

The Nonstationary Newsvendor with Predictions problem assumes all of the setup, assumptions, and notation of the previous Nonstationary Newsvendor problem. In addition, at each time period t we assume that the decision-maker receives a prediction \hat{d}_t before selecting an order quantity q_t .¹¹ This prediction is meant to be an estimate of d_t and so we measure the prediction error of a sequence $a = (a_1, \dots, a_T)$ with respect to a sequence of means simply as

$$\sum_{t=1}^T |a_t - \hat{d}_t|$$

Note that unlike demand variation, we have not used partitions here (and in fact, introducing partitions would not have any effect since we are measuring absolute rather than squared differences). Intuitively, we do not want to require the sequence of errors to be meaningful time-series: the predictions are generic, and their accuracy is allowed to change rapidly. Just as for the demand variation, the prediction error is expected to grow with the time horizon, and the proper parameterization of this growth is via an exponent: we call the accuracy parameter the small constant $\alpha > 0$ such that the prediction error is at most αT^α . We will always assume that α is unknown to the decision-maker.

¹¹We are taking the predictions to be entirely deterministic, so for example, not allowed to depend on the previously-observed demands d_1, \dots, d_{t-1} . Our results hold if we extend to the setting in which the predictions are stochastic (and adapted to the demand filtration).

Algorithm 6: Prediction Policy

```

for  $C = 1 \dots T$  do
   $\lambda_C = 0$ 
  if  $\lambda_C \geq \frac{1}{\sqrt{C}}$  round  $\lambda_C$  to the nearest value in  $\{\frac{1}{\sqrt{C}}, \dots, \frac{1}{\sqrt{2C}}\}$ 
   $c_C = \arg \min_{c \in \mathcal{C}} \sum_{t=1}^C \ell_t(c)$ 
end

```

Naturally, the notion of a policy c expands to include the predictions $\mathbf{s}_t = (s_{t1}, \dots, s_{tC})$, where each c_C is a mapping from $\mathcal{S}_1 \times \dots \times \mathcal{S}_C$ to an order quantity $c \in \mathcal{C}$. The simplest policy, which should be used if the prediction error is known to be sufficiently small, is to simply behave as if the predictions were perfect. We call this the Prediction Policy (Algorithm 6). The following observation collects a few (likely unsurprising) facts about the performance of this policy, with respect to worst-case regret (generalized in the obvious manner to incorporate prediction accuracy via the accuracy parameter θ):

Observation 3 (Upper and Lower Bounds: Prediction Policy). Fix any variation parameter $\epsilon \in (0, 1/4]$ and any accuracy parameter $\theta \in (0, 1/4]$.

- a) The Prediction Policy $c^{\text{prediction}}$ achieves worst-case regret

$$R^{c^{\text{prediction}}}(\theta) \leq \frac{1}{\theta} \epsilon$$

where $\epsilon \in (0, 1/4]$.

- b) For any policy c (which may depend on the knowledge of θ) that is solely a function of the predictions (i.e. does not depend on the observed demands), we have

$$R^c(\theta) \geq \frac{1}{2} \epsilon$$

where $\frac{1}{2} \epsilon > 0$ is a universal constant.

Observation 3a) states that the Prediction Policy translates prediction error directly to regret (incidentally, it does this without knowing θ). There are of course other ways in which the predictions could be used, but Observation 3b) essentially states that there is nothing to be gained by doing so (even if known). The proof of Observation 3a) appears in Appendix C.1. Observation 3b) is a direct corollary of Proposition 11, which is given in the next subsection.

Extreme Cases

What exactly is achievable for the Nonstationary Newsvendor with Predictions depends heavily on whether or not \mathbb{E} and $\mathbb{0}$ are known to the policy. To see this, it is worth first considering the two extremes.

Case 1: Known \mathbb{E} and $\mathbb{0}$. A simple policy is available when \mathbb{E} and $\mathbb{0}$ are both known. Compare the quantities \mathbb{E} , $\mathbb{E}^{\mathbb{0}\cdot\mathbb{4}}$ and $\mathbb{0}$. If the former quantity is smaller, apply the Fixed-Time-Window Policy. If the latter is smaller, apply the Prediction Policy. Lemma 3 and Observation 3 together imply that this achieves a worst-case regret of $\mathbb{1}) \min\{\mathbb{3}, \mathbb{E}^{\mathbb{0}\cdot\mathbb{4}} - \mathbb{0}\}$. This is optimal, as demonstrated by the following result:

Proposition 11 (Lower Bound: Known \mathbb{E} and $\mathbb{0}$). Fix any variation parameter $\mathbb{E} \geq \mathbb{0}$ and any accuracy parameter $\mathbb{0} \geq \mathbb{0}$. For any policy \mathbb{c} (which may depend on the knowledge of \mathbb{E} and $\mathbb{0}$), we have

$$R^{\mathbb{c}}(\mathbb{1}) \geq \mathbb{2}) \min\{\mathbb{3}, \mathbb{E}^{\mathbb{0}\cdot\mathbb{4}} - \mathbb{0}\}$$

where $\mathbb{2}$; $\mathbb{0}$ is a universal constant.

The proof of this result can be found in Appendix C.4, and relies on an explicit construction of a family of problem instances. Our construction breaks the total time horizon into cycles wherein the demand distribution is i.i.d.. We tune the length of each cycle to be small enough so that it is (provably) hard to detect the change in demand distributions and the predictions are essentially useless for most time periods in the cycle, and large enough so that the demand variation is within the prediction error is within $\mathbb{0}$.

Case 2: Unknown \mathbb{E} and $\mathbb{0}$. At the opposite extreme, \mathbb{E} and $\mathbb{0}$ are both unknown, is it still possible to achieve $\mathbb{1}) \min\{\mathbb{3}, \mathbb{E}^{\mathbb{0}\cdot\mathbb{4}} - \mathbb{0}\}$ worst-case regret? The answer is no:

Proposition 12 (Lower Bound: Unknown \mathbb{E} and $\mathbb{0}$). For any policy that does not depend on the knowledge of \mathbb{E} or $\mathbb{0}$, there exists a problem instance such that $\mathbb{0} < \mathbb{3}$, $\mathbb{E}^{\mathbb{0}\cdot\mathbb{4}}$, and the policy incurs regret at least $\mathbb{2}) \max\{\mathbb{3}, \mathbb{E}^{\mathbb{0}\cdot\mathbb{4}} - \mathbb{0}\}$ on the instance, where $\mathbb{2}$; $\mathbb{0}$ is a universal constant.

Proposition 12 states that the best we can hope for, when \mathbb{E} and $\mathbb{0}$ are unknown, is a worst-case regret of at least $\mathbb{1}) \max\{\mathbb{3}, \mathbb{E}^{\mathbb{0}\cdot\mathbb{4}} - \mathbb{0}\}$.¹² Note that Proposition 12 shows

¹²Indeed, it implies that no algorithm can achieve regret $\mathbb{5} \cdot \mathbb{1} \cdot \mathbb{E}^{-\mathbb{0}\mathbb{0}}$ for a function $\mathbb{5} : \mathbb{0} - \mathbb{1} \cdot \mathbb{4} \rightarrow \mathbb{0} - \mathbb{1} \cdot \mathbb{4}$ satisfying $\mathbb{5} \cdot \mathbb{1} \cdot \mathbb{E} - \mathbb{0} \leq \max\{\mathbb{3}, \mathbb{E}^{\mathbb{0}\cdot\mathbb{4}} - \mathbb{0}\}$ for all $\mathbb{E} - \mathbb{0} \geq \mathbb{2} \cdot \mathbb{0}$ and $\mathbb{5} \cdot \mathbb{1} \cdot \mathbb{E} - \mathbb{0} \leq \max\{\mathbb{3}, \mathbb{E}^{\mathbb{0}\cdot\mathbb{4}} - \mathbb{0}\}$ for some $\mathbb{E} - \mathbb{0} \geq \mathbb{2} \cdot \mathbb{0}$.

there exists a pair ϵ and δ , and a corresponding problem instance, such that this lower bound holds. This is in contrast to a result showing that for any ϵ and δ , there exists a problem instance such that the lower bound holds, as is common in the literature (e.g. Theorem 1.1 in [1]). This lower bound is easily achieved, for example by applying the Shrinking-Time-Window Policy or the Prediction Policy (or any blind randomization of the two). The proof of Proposition 12 is in Appendix C.6. In contrast to Case 1, the lower bound construction here relies heavily on the fact that we do not know which one of ϵ , δ and 0 is smaller.

Final Solution

We have finally reached the problem which motivates this entire paper: designing an optimal policy for the Nonstationary Newsvendor with Predictions when the prediction error ϵ is unknown. We will assume that δ is known, since when δ is unknown, Proposition 12 rules out the possibility of using the predictions to improve on what is already achievable without predictions. On the other hand, by Proposition 11, the absolute best we could hope for is a policy which achieves a worst-case regret of $\min\{\epsilon, \delta\}$. In words, we would like a policy which, without knowing ϵ , achieves the same regret as if ϵ had been known. Our main result is the design of such a policy.

Our policy is called the Prediction-Error-Robust Policy (PERP), and is given in Algorithm 7. PERP utilizes the Fixed-Time-Window policy λ^{FTW} in Section 4.3 as an estimate of the true mean to track the quality of the predictions over time.

Theorem 3 (Upper Bound: Known δ and Unknown ϵ). For any variation parameter $\epsilon \geq 0$ and any accuracy parameter $\delta \geq 0$, the Prediction-Error-Robust Policy λ^{PERP} achieves worst-case regret

$$R^{\text{PERP}}(\epsilon, \delta) \leq \min\left\{ \epsilon, \delta, \frac{\delta}{\log\left(\frac{\delta}{\epsilon}\right)} \right\} + \epsilon$$

where Lemma 3 is the constant in Lemma 3 (and δ matches the constant in Observation 3a).

The intuition behind PERP is to follow the predictions until a time that is late enough to have evidence that the prediction quality is bad (compared to the Fixed-Time-Window Policy), while early enough to not incur much regret caused by the poor quality of the predictions. Because we do not observe the true past ϵ after time period C , we naturally use λ^{FTW} from the Fixed-Time-Window policy λ^{FTW} as an

Algorithm 7: Prediction-Error-Robust Policy (PERP)

Inputs: variation parameter $E \gg 0$ and scaling constants $\lambda_j, 0, W$ sufficiently large (Equation (C.3) in Appendix C.7)

Initialization: $c_C = d^{\lambda_j} E^{0.2}$

for $C = 1, \dots, T$ do

$c_C = c_C^{\text{prediction}}$

end

for $C = 1, \dots, T$ do

$\lambda_C^{\text{xed}} = \frac{1}{B=C} = \frac{1}{3B}$

$\lambda_C^0 = 0$

 if $\sum_{B=1}^C \lambda_B^0 \lambda_B^{\text{xed}} \leq W \log(C) + \frac{1}{E^{0.4}}$ then

$c_C = c_C^{\text{xed}}$

 break

 end

 else

$c_C = c_C^{\text{prediction}}$

 end

end

estimation of c_C and in turn keep tracking the cumulative difference the prediction quality $\sum_{j=1}^C d_j$. We carefully choose the parameters λ_B^{PERP} so that this estimation is not accurate only with a small probability, and we can identify the prediction quality is bad if this cumulative difference is too large. By Proposition 11, any policy can only achieve worst-case regret on the order of $\min\{13, E^{0.4}\}$, so PERP is order-optimal.

Aside: Unknown E and Known 0 . There are four possible scenarios depending on the knowledge of E and 0 : known/unknown E and known/unknown 0 . So far we have discussed three of them: known E and 0 (Proposition 11), unknown E and 0 (Proposition 12), and known E and unknown 0 (Theorem 3). For the sake of completeness, we discuss the remaining case of unknown E and known 0 in Appendix C.8, where we give a policy that achieves worst-case regret $\min\{13, E^{0.4}\}$. This is order-optimal by Theorem 3.

4.5 Experiments

Finally, we describe a set of experiments we performed to evaluate our policy (for the Nonstationary Newsvendor with Predictions). In all of our experiments, we compared PERP against the Shrinking-Time-Window Policy (STW-PERP) and the

Prediction Policy (PURE-PRED). The main takeaways are:

1. PERP performance is robust with respect to the quality of the predictions, without knowing the prediction quality beforehand. Specially, the (news vendor) cost it incurs is consistently close to the lower of the costs incurred by NO-PRED and PURE-PRED.
2. PERP performs especially well when the absolute difference between the costs of NO-PRED and PURE-PRED is large, i.e. when the stakes are highest.

Experiments on Synthetic Data

The objective of our first batch of experiments was to fix one of the two theoretical parameters (E or θ) and test PERP performance as the other parameters changes. To generate demand sequences, we used the parametric time-series model that corresponds to triple exponential smoothing (Holt Winters), a classic model for time-series data in the family of Equation (4.3). We give the exact formulas, and our choices of parameters, in Appendix C.9; for more on triple exponential smoothing, see [175]. In our experiment, each demand sequence consisted of the demands for the next 365 time periods, with the realized demands generated as Poisson variables with the corresponding means. We ran two sets of experiments:

- ^ Fixed E : We fixed a single set of parameters for the demand sequence and generated 1,000 different predictions of this demand sequence, each from a set of predicted parameters with different accuracy. Thus the variation parameter E was fixed, and the accuracy parameter θ varied across instances.
- ^ Fixed θ : We generated 1,000 demand sequences by selecting the parameters randomly. We then generated predictions by changing each parameter 10% and using the corresponding sequence. Thus the variation parameter E varied across instances, but the accuracy parameter θ was (roughly) fixed.

For each demand sequence and corresponding prediction, NO-PRED, PURE-PRED and PERP with equal overage and underage costs, and scaling constant $\alpha = 1$. Because the experiment was synthetic, the true underlying demand distribution was known at each time period. Therefore we also set a benchmark, which simply ordered the optimal quantile at each time period. The variation parameter E for PERP was calculated using the past demands of the pre-fixed 30 time periods by the definition in Section 2.2. We calculated the parameter θ by their definitions (given in Section 2.2 and Section 4.1, respectively), scaled appropriately to make

(a) Fixed E (b) Fixed θ

Figure 4.2: The costs of NO-PREP, PURE-PREP, PERP and OPT when (a) the variation parameter E is fixed, and (b) the accuracy parameter θ is fixed. Each dot represents the cost of the corresponding policy on a given instance.

them lie in $\approx 0 - 1/4$. The resulted scatter plots are shown in Figure 4.2. In (a), E is fixed, so the cost of NO-PREP (blue dots) is approximately the same for all instances. The cost of PURE-PREP (orange dots) is approximately exponential in θ , which follows by Observation 3a). In (b), θ is fixed, so the cost of PURE-PREP is approximately the same for all instances. The cost of NO-PREP is approximately exponential in E , which follows by Theorem 2. Note that in both (a) and (b), the cost of PERP (green dots) is close to the minimal cost of NO-PREP and PURE-PREP across all instances, showing that PERP's performance is robust in both E and θ .

Experiments on Real Data

We used real-world datasets to represent the demand sequences in our experiments. Figure 4.3 depicts example time series from each of these datasets. All datasets include multiple daily time series and are publicly available:

- ^ Rossmann:¹³ Daily number of customers that visited each of 1,115 stores in the Rossmann drug store chain during a 781-day period in 2013-2015.
- ^ Wikipedia:¹⁴ Daily web traffic across Wikipedia.com pages of 9 different languages for an 803-day period from 2015 to 2017.
- ^ Restaurant:¹⁵ Daily number of visitors and online reservations across 185 restaurants in Japan, during a 478-day period in 2016-2017. We treated

¹³Available at <https://www.kaggle.com/competitions/rossmann-store-sales/data>

¹⁴Available at <https://www.kaggle.com/competitions/web-traffic-time-series-forecasting/data>

¹⁵Available at <https://www.kaggle.com/competitions/recruit-restaurant-visitor-forecasting/>

the number of visitors as the demand, and the reservations as a predictive feature.

(a) Rossmann (b) Wikipedia (c) Restaurant

Figure 4.3: An example of a single time series from each dataset. In (c), the red dashed line represents an additional feature: daily online reservations.

Each instance of our experiment represented a single Nonstationary Newsvendor with Predictions problem, with the realized demands taken from a single time series in our data (a single Rossmann store, a single language on Wikipedia, or a single restaurant). The overage and underage costs were constant within each instance, and without loss of generality the two costs for an instance can be characterized by the corresponding critical quantile (specifically the ratio of the underage cost to the sum). The time horizon for each instance was a set number of days taken from the end of the time series, with the preceding days used to train one of four prediction methods. These predictions were also updated over the course of the instance at a set frequency. For the Wikipedia dataset, this yielded a total of 2,880 possible instances, all of which were tested. The Rossmann dataset has multiple orders of magnitude more instances, so we randomly sampled 1,000 from this set. For the Restaurant dataset, we used a single prediction method to generate two sets of predictions for each restaurant: one only utilized the number of past visitors and the other incorporated the number of reservations as a feature, which gave 740 instances. Table 4.3 describes all of the instances used.

For each instance, we applied \hat{D}_t -PURE-PRED and \hat{D}_t -PERP with scaling constants $\hat{\alpha} = W = 1$ and the variation parameter $\hat{\sigma}_t$ in PERP was calculated using the past demands of the training data by the definition in Section 2.2. To generate predictions, we used four popular forecasting methods ranging from classic to the state-of-the-art:

- Exponential Smoothing (Holt Winters): A classic algorithm based on a (linear) trend and seasonality decomposition as in Equation (4.3), known for

	Rossmann	Wikipedia	Restaurant
Number of time series	1,115	9	185
Critical quantiles (%)	30,40,50,60,70	95,98,99,99.9	50
Experimental period (days)	300,400,500,600	300,400,500,600,700	100
Prediction update frequency (days)	2,4,10,20	2,4,10,20	5, 10
Total number of instances	1,000 (sampled)	2,880 (exhaustive)	740 (exhaustive)

Table 4.3: Description of experimental instances.

its simplicity and robust performance. It is frequently used as a benchmark in forecasting competitions [125]. Tuning parameter: seasonality of length 50.

- ^ ARIMA: Another classic algorithm that is rich enough to model a wide class of nonstationary time-series. Tuning parameters: $\lambda = 13-2-5^0$.
- ^ Prophet: A recent algorithm developed by Facebook [62] based on a (piecewise-linear) trend and seasonality decomposition as in Equation (4.3), known to work well in practice with minimal tuning. Tuning parameters: software default.
- ^ LightGBM: A recent algorithm developed by Microsoft [94] based on tree algorithms. LightGBM formed the core of most of the top entries in the recent \$100,000 M5 Forecasting Challenge [25]. Tuning parameters: software default.

For the Restaurant dataset, we used Prophet as the forecasting method, with and without the reservations as an additive linear feature. We treated the outputs of these methods as predictions of the mean demand. To estimate the demand distribution around this mean, we used the empirical distribution of the residuals of the same predictions on the training period.¹⁶ In practice, even if the prediction quality is good,

¹⁶That is, if the training data consists of p_{train} periods, which without loss we index as $(t_{\text{train}, 1}, \dots, t_{\text{train}, p_{\text{train}}})$, then the demand distribution at any time C was estimated to be

$$\lambda_C, \text{Uniform } (f_{B, \lambda} : B = (t_{\text{train}, 1}, \dots, t_{\text{train}, p_{\text{train}}}))$$

the predictions of the first few days might incur large costs due to noise/instability of the predictions, which may cause PERP to misidentify the prediction quality. Therefore we restricted PERP to following the predictions for the first 20 days, only allowing switches afterward.

Results. Each instance yields three total costs: one incurred by and two incurred by the benchmark algorithms NO-PRED and PURE-PRED. The primary performance metric we report is a form of optimality gap. For an instance let $\text{cost}^{\text{PURE-PRED}}_i$ be the cost of PURE-PRED and similarly denote $\text{cost}^{\text{NO-PRED}}_i$, $\text{cost}^{\text{PERP}}_i$. Then the optimality gap (GAP) of PERP is defined as

$$\text{GAP}_i = \frac{\text{cost}^{\text{PERP}}_i - \min(\text{cost}^{\text{PURE-PRED}}_i, \text{cost}^{\text{NO-PRED}}_i)}{\max(\text{cost}^{\text{PURE-PRED}}_i, \text{cost}^{\text{NO-PRED}}_i)}$$

If we think of PERP as trying to achieve the minimum of the costs incurred by the two benchmark policies, then GAP measures the excess cost that occurs on top of this minimum, normalized so that $\text{GAP} = 0$ implies that the minimum has been achieved, and $\text{GAP} = 1$ implies that the maximum of the two costs was incurred.

(a) Rossmann (b) Wikipedia (c) Restaurant

Figure 4.4: Histograms of GAPs across (a) 1,000 randomly-sampled instances on the Rossmann dataset, (b) 2,880 instances on the Wikipedia dataset, and (c) 740 instances on the Restaurant dataset.

Experiments on the datasets yielded the histograms in Figure 4.4. For each instance i , the value on the horizontal axis is $\log_2 \frac{\text{cost}^{\text{PURE-PRED}}_i - \text{cost}^{\text{NO-PRED}}_i}{\max(\text{cost}^{\text{PURE-PRED}}_i, \text{cost}^{\text{NO-PRED}}_i)}$, which is greater than 0 if NO-PRED has a lower cost, and less than 0 if PURE-PRED has a lower cost. In the 1,000 Rossmann instances NO-PRED had a lower cost 82.7% of the time, in the 2,880 Wikipedia instances NO-PRED had a lower cost 81.9% of the time, and in the 740 Restaurant instances NO-PRED had a lower cost 64.3% of the time. The values on the vertical axis are the GAPs. Note that most GAPs are small when the absolute

¹⁷GAP may technically be outside of $[-1, 1]$.

values of the log difference are large. This shows PERP performs very well when the difference of costs between NO-PRED and PURE-PRED is large. On the other hand, there are instances where PERP has large GAPs, in particular there are instances with GAPs equal to 1 when the log difference of costs is close to 0. This happens because when the log difference of costs is close to 0, the cost of NO-PRED and the cost of PURE-PRED are close, so PERP may misidentify the prediction quality. Still, since the max cost and the min cost of the other two policies are close, even the GAPs are large in these instances, PERP does not perform badly.

	Rossmann	Wikipedia	Restaurant
Average GAP with good predictions	0.26	0.40	0.10
Average GAP with bad predictions	0.28	0.07	0.39

Table 4.4: Summary of experimental results.

We further divide the instances according to which of NO-PRED and PURE-PRED had lower cost in Table 4.4 and Figure 4.5. For comparison, if we did not know the

(a) GAP with high prediction accuracy. Left to right: Rossmann, Wikipedia, Restaurant.

(b) GAP with low prediction accuracy. Left to right: Rossmann, Wikipedia, Restaurant.

Figure 4.5: Histograms of the GAPs for (a) 173 Rossmann instances (left), 522 Wikipedia instances (middle), 476 Restaurant instances (right) for which PURE-PRED has lower cost, and (b) 827 Rossmann instances (left), 2358 Wikipedia instances (middle), 264 Restaurant instances (right) for which NO-PRED has lower cost.

prediction quality beforehand, uniformly random choosing between NO-PRED and PURE-PRED has an expected GAP of 0.5. Therefore PERP outperforms this natural benchmark in all cases of all datasets.

4.6 Conclusion

We proposed a new model incorporating predictions into the nonstationary newsvendor problem. We first gave a complete analysis of the Nonstationary Newsvendor (without predictions) by proving a lower regret bound and developing the Shrinking-Time-Window Policy, which was the first policy that achieves the lower bound up to log factors without knowing the variation parameter. We then considered the Nonstationary Newsvendor with Predictions and proposed the Prediction-Error-Robust Policy, which does not need to know the prediction quality beforehand, and achieves nearly optimal minimax worst-case regret.

Chapter 5

CONCLUSION

This thesis studies a central question in modern operations management: how can firms fully leverage the power of modern AI in sequential decision-making while retaining the rigor of optimization? Across a wide range of operational settings, AI models now provide increasingly rich forms of guidance, from demand forecasts and estimated opportunity costs to expressive neural architectures that capture complex behavioral interactions. Yet the availability of such models does not by itself resolve the underlying decision problem. Their outputs must still be translated into actions that are timely, scalable, and reliable. This thesis develops optimization foundations for AI-driven sequential decision-making, showing how modern AI can be incorporated into operational decisions in a way that is computationally tractable, robust to misspecification, and supported by formal performance guarantees.

A unifying theme across the thesis is that AI creates value for operations only when its outputs induce decision-relevant structure. In Chapter 2, this structure takes the form of a tractable subclass of transformer architectures. Transformers are powerful because they capture interaction effects across recommended items, but this same expressiveness can make downstream optimization computationally intractable. By identifying the class of simple transformers, the chapter shows that one can preserve key behavioral effects, including variety seeking, complementarity, and substitution, while still enabling near-optimal recommendation in sublinear time. The chapter further demonstrates, through both online experiments and large-scale field evidence, that these behavioral effects matter in practice and can generate measurable gains in user engagement. More broadly, the chapter illustrates that expressive AI models need not be black boxes for operations; when their structural properties are understood, they can be integrated into real-time optimization with theoretical rigor.

In Chapters 3 and 4, AI enters the decision process in a different way: not as the objective itself, but as a source of guidance for sequential decisions under uncertainty. Chapter 3 studies online resource allocation with predictions in the form of shadow prices. A central challenge in such settings is that prediction quality is unknown in advance, and algorithms that rely too heavily on inaccurate guidance may perform poorly. The chapter addresses this challenge by designing allocation policies that

adapt to prediction quality without requiring it as input. The resulting guarantees sharply characterize how one can balance robustness and consistency: the algorithm remains protected against poor predictions while improving automatically as the predictions become more accurate. Chapter 4 develops the same philosophy in inventory control under nonstationary demand. It first provides a complete regret characterization of the nonstationary newsvendor problem without predictions, and then shows how generic demand forecasts can be incorporated without assuming prior knowledge of their quality. The resulting policies adapt simultaneously to environmental change and to prediction error, remaining robust when forecasts are poor while benefiting when they are informative.

Taken together, these chapters suggest three broader principles for AI-driven operations. First, predictive or representational power alone is not enough; what matters is whether that power can be converted into actionable structure for optimization. Second, because model quality is typically unknown and can change over time, algorithms must be adaptive to the reliability of AI guidance rather than assuming it is correct. Third, progress in this area requires end-to-end guarantees that connect the AI model, the induced optimization problem, and the resulting sequential performance. These principles provide a common foundation across seemingly different settings, from personalization to revenue management to inventory control.

Beyond the specific technical results, the broader implication of this thesis is that AI and operations management should not be treated as separate layers of a system, with one producing predictions and the other consuming them passively. Instead, they should be designed jointly. The central opportunity is not merely to make better predictions, but to understand how AI models can reshape the structure of operational decisions themselves. This perspective becomes increasingly important as operational environments become more dynamic, as decision latency becomes more stringent, and as firms rely on AI models that are more expressive but also harder to optimize over directly.

Several directions for future research follow naturally from this work. One direction is to study richer model classes beyond the structured settings considered here, and to identify broader families of neural architectures that admit efficient optimization with guarantees. A second direction is to understand settings in which the quality of AI guidance evolves endogenously over time, for example because predictions are updated online or because decisions themselves affect the future data-generating process. A third direction is to extend these ideas to settings with strategic behavior,

partial observability, or multiple interacting decision-makers, where the interface between AI and optimization becomes even more subtle. From an applied standpoint, there is also substantial opportunity to bring these methods to a wider range of operational domains, including digital platforms, retail media, supply chains, and service systems, where firms increasingly require decisions that are both behaviorally informed and operationally scalable.

In summary, this thesis argues that the next stage of research in operations management is not simply to apply AI to existing decision problems, but to develop the optimization principles that make AI actionable. By studying transformer-based personalization, prediction-robust online allocation, and nonstationary inventory control, the thesis takes several steps toward that goal. I hope these results contribute to a broader research agenda at the interface of AI, optimization, and operations management, where modern predictive models are not treated as opaque tools, but as structured inputs to decision-making algorithms that are both practically effective and theoretically grounded.

BIBLIOGRAPHY

- [1] Shipra Agrawal, Zizhuo Wang, and Yinyu Ye. A dynamic near-optimal algorithm for online linear programming. *Operations Research*, 62(4):876–890, 2014.
- [2] Nir Ailon and Bernard Chazelle. The fast johnson lindenstrauss transform and approximate nearest neighbors. *SIAM Journal on computing*, 39(1):302–322, 2009.
- [3] Josh Alman and Zhao Song. Fast attention requires bounded entries. *Advances in Neural Information Processing Systems*, 36, 2024.
- [4] Lin An, Andrew A Li, Benjamin Moseley, and Gabriel Visotsky. Best of many in both worlds: Online resource allocation with predictions under unknown arrival model. *arXiv preprint arXiv:2402.13530*, 2024.
- [5] Lin An, Andrew A Li, Benjamin Moseley, and R Ravi. The nonstationary newsvendor with (and without) predictions. *Manufacturing & Service Operations Management*, 27(3):881–896, 2025.
- [6] Lin An, Andrew A Li, Vaisnavi Nemala, and Gabriel Visotsky. Real-time personalization with simple transformers. *ACM SIGMETRICS Performance Evaluation Review*, 53(2):116–118, 2025.
- [7] Alexandr Andoni and Piotr Indyk. Near-optimal hashing algorithms for approximate nearest neighbor in high dimensions. *Communications of the ACM*, 51(1):117–122, 2008.
- [8] Alexandr Andoni, Piotr Indyk, Thijs Laarhoven, Ilya Razenshteyn, and Ludwig Schmidt. Practical and optimal lsh for angular distance. *Advances in neural information processing systems*, 28, 2015.
- [9] Antonios Antoniadis, Themis Gouleakis, Pieter Kleer, and Pavel Kolev. Secretary and online matching problems with machine learned advice. *Advances in Neural Information Processing Systems*, 33:7933–7944, 2020.
- [10] Alessandro Arlotto and Itai Gurvich. Uniformly bounded regret in the multisecretary problem. *Stochastic Systems*, 9(3):231–260, 2019.
- [11] Kenneth Joseph Arrow et al. *Studies in the mathematical theory of inventory and production*. JSTOR, 1958.
- [12] Sunil Arya, David M Mount, Nathan S Netanyahu, Ruth Silverman, and Angela Y Wu. An optimal algorithm for approximate nearest neighbor searching in high dimensions. *Journal of the ACM (JACM)*, 45(6):891–923, 1998.

- [13] Peter Auer, Nicolo Cesa-Bianchi, Yoav Freund, and Robert E Schapire. The nonstochastic multiarmed bandit problem. *SIAM journal on computing*, 32(1):48–77, 2002.
- [14] Yossi Aviv and Amit Pazgal. A partially observed markov decision process for dynamic pricing. *Management science*, 51(9):1400–1416, 2005.
- [15] Katy S Azoury. Bayes solution to dynamic inventory models under unknown demand distribution. *Management science*, 31(9):1150–1160, 1985.
- [16] Dheeraj Baby and Yu-Xiang Wang. Online forecasting of total-variation-bounded sequences. *NIPS*, 32, 2019.
- [17] Yong Bai, Yu-Jie Zhang, Peng Zhao, Masashi Sugiyama, and Zhi-Hua Zhou. Adapting to online label shift with provable guarantees. *NIPS*, 35:29960–29974, 2022.
- [18] Michael O Ball and Maurice Queyranne. Toward robust revenue management: Competitive analysis of online booking. *Operations Research*, 57(4):950–963, 2009.
- [19] Santiago Balseiro, Haihao Lu, and Vahab Mirrokni. Dual mirror descent for online allocation problems. In *International Conference on Machine Learning*, pages 613–628. PMLR, 2020.
- [20] Santiago Balseiro, Christian Kroer, and Rachitesh Kumar. Single-leg revenue management with advice. *arXiv preprint arXiv:2202.10939*, 2022.
- [21] Santiago R Balseiro and Yonatan Gur. Learning in repeated auctions with budgets: Regret minimization and equilibrium. *Management Science*, 65(9):3952–3968, 2019.
- [22] Santiago R Balseiro, Haihao Lu, and Vahab Mirrokni. The best of many worlds: Dual mirror descent for online allocation problems. *Operations Research*, 71(1):101–119, 2023.
- [23] Etienne Bamas, Andreas Maggiori, and Ola Svensson. The primal-dual method for learning augmented algorithms. *Advances in Neural Information Processing Systems*, 33:20083–20094, 2020.
- [24] Gah-Yi Ban and Cynthia Rudin. The big data newsvendor: Practical insights from machine learning. *Operations Research*, 67(1):90–108, 2019.
- [25] Salvador Barberà, Peter Hammond, and Christian Seidl. *Handbook of utility theory: volume 2 extensions*. Springer Science & Business Media, 2004.
- [26] Manel Baucells and Rakesh K Sarin. Satiation in discounted utility. *Operations research*, 55(1):170–181, 2007.

- [27] Heinz H Bauschke, Jonathan M Borwein, and Patrick L Combettes. Essential smoothness, essential strict convexity, and legendre functions in banach spaces. *Communications in Contemporary Mathematics*, 3(04):615–647, 2001.
- [28] Cristina Bazgan, Hadrien Hugot, and Daniel Vanderpooten. Implementing an efficient fptas for the 0–1 multi-objective knapsack problem. *European Journal of Operational Research*, 198(1):47–56, 2009.
- [29] Cristina Bazgan, Hadrien Hugot, and Daniel Vanderpooten. Solving efficiently the 0–1 multi-objective knapsack problem. *Computers & Operations Research*, 36(1):260–279, 2009.
- [30] Amir Beck and Marc Teboulle. Mirror descent and nonlinear projected subgradient methods for convex optimization. *Operations Research Letters*, 31(3):167–175, 2003.
- [31] Walid Bendada, Théo Bontempelli, Mathieu Morlon, Benjamin Chapus, Thibault Cador, Thomas Bouabça, and Guillaume Salha-Galvan. Track mix generation on music streaming services using transformers. In *Proceedings of the 17th ACM Conference on Recommender Systems*, pages 112–115, 2023.
- [32] Steve Berry, Ahmed Khwaja, Vineet Kumar, Andres Musalem, Kenneth C Wilbur, Greg Allenby, Bharat Anand, Pradeep Chintagunta, W Michael Hanemann, Przemek Jeziorski, et al. Structural models of complementary choices. *Marketing Letters*, 25:245–256, 2014.
- [33] Dimitri P Bertsekas. *Nonlinear programming*. *Journal of the Operational Research Society*, 48(3):334–334, 1997.
- [34] Dimitri P Bertsekas. *Constrained optimization and Lagrange multiplier methods*. Academic press, 2014.
- [35] Omar Besbes and Alp Muharremoglu. On implications of demand censoring in the newsvendor problem. *Management Science*, 59(6):1407–1424, 2013.
- [36] Omar Besbes, Yonatan Gur, and Assaf Zeevi. Stochastic multi-armed-bandit problem with non-stationary rewards. *Advances in neural information processing systems*, 27, 2014.
- [37] Omar Besbes, Yonatan Gur, and Assaf Zeevi. Non-stationary stochastic optimization. *Operations research*, 63(5):1227–1244, 2015.
- [38] Alberto Bietti, Vivien Cabannes, Diane Bouchacourt, Herve Jegou, and Leon Bottou. Birth of a transformer: A memory viewpoint. *Advances in Neural Information Processing Systems*, 36, 2024.
- [39] Sébastien Bubeck et al. *Convex optimization: Algorithms and complexity*. *Foundations and Trends® in Machine Learning*, 8(3-4):231–357, 2015.

- [40] Niv Buchbinder, Kamal Jain, and Joseph Naor. Online primal-dual algorithms for maximizing ad-auctions revenue. In European Symposium on Algorithms, pages 253 264. Springer, 2007.
- [41] Apostolos N Burnetas and Craig E Smith. Adaptive ordering and pricing for perishable products. *Operations Research*, 48(3):436 443, 2000.
- [42] Jaime Carbonell and Jade Goldstein. The use of mmr, diversity-based reranking for reordering documents and producing summaries. In Proceedings of the 21st Annual International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR), pages 335 336, New York, NY, 1998. ACM.
- [43] Yair Carmon and Oliver Hinder. Making sgd parameter-free. In Conference on Learning Theory, pages 2360 2389. PMLR, 2022.
- [44] Marjan Celikik, Ana Peleteiro Ramallo, and Jacek Wasilewski. Reusable self-attention recommender systems in fashion industry applications. In Proceedings of the 16th ACM Conference on Recommender Systems, pages 448 451, 2022.
- [45] Kamalika Chaudhuri, Yoav Freund, and Daniel J Hsu. A parameter-free hedging algorithm. *Advances in neural information processing systems*, 22, 2009.
- [46] Boxiao Chen, Xiuli Chao, and Hyun-Soo Ahn. Coordinating pricing and inventory replenishment with nonparametric demand learning. *Operations Research*, 67(4):1035 1052, 2019.
- [47] Ching-Wei Chen, Paul Lamere, Markus Schedl, and Hamed Zamani. Recsys challenge 2018: Automatic music playlist continuation. In Proceedings of the 12th ACM Conference on Recommender Systems, pages 527 528, 2018.
- [48] Jianer Chen, Xiuzhen Huang, Iyad A Kanj, and Ge Xia. Strong computational lower bounds via parameterized complexity. *Journal of Computer and System Sciences*, 72(8):1346 1367, 2006.
- [49] Jingyuan Chen, Hanwang Zhang, Xiangnan He, Liqiang Nie, Wei Liu, and Tat-Seng Chua. Attentive collaborative filtering: Multimedia recommendation with item-and component-level attention. In Proceedings of the 40th International ACM SIGIR conference on Research and Development in Information Retrieval, pages 335 344, 2017.
- [50] Qiwei Chen, Huan Zhao, Wei Li, Pipei Huang, and Wenwu Ou. Behavior sequence transformer for e-commerce recommendation in alibaba. In Proceedings of the 1st international workshop on deep learning practice for high-dimensional sparse data, pages 1 4, 2019.

- [51] Xi Chen, Yining Wang, and Yu-Xiang Wang. Nonstationary stochastic optimization under l_p, q -variation measures. *Operations Research*, 67(6): 1752–1765, 2019.
- [52] Wang Chi Cheung, David Simchi-Levi, and Ruihao Zhu. Hedging the drift: Learning to optimize under nonstationarity. *Man. Sci.*, 68(3):1696–1713, 2022.
- [53] Joel E Cohen and Uriel G Rothblum. Nonnegative ranks, decompositions, and factorizations of nonnegative matrices. *Linear Algebra and its Applications*, 190:149–168, 1993.
- [54] Ashok Cutkosky. Artificial constraints and hints for unbounded online learning. In *Conference on Learning Theory*, pages 874–894. PMLR, 2019.
- [55] Ashok Cutkosky and Kwabena Boahen. Online learning without prior information. In *Conference on learning theory*, pages 643–677. PMLR, 2017.
- [56] Ashok Cutkosky and Francesco Orabona. Black-box reductions for parameter-free online learning in banach spaces. In *Conference On Learning Theory*, pages 1493–1529. PMLR, 2018.
- [57] Nikhil R Devanur and Thomas P Hayes. The adwords problem: online keyword matching with budgeted bidders under random permutations. In *Proceedings of the 10th ACM conference on Electronic commerce*, pages 71–78, 2009.
- [58] Nikhil R Devanur, Kamal Jain, Balasubramanian Sivan, and Christopher A Wilkens. Near optimal online algorithms and fast approximation algorithms for resource allocation problems. In *Proceedings of the 12th ACM conference on Electronic commerce*, pages 29–38, 2011.
- [59] Michael Dinitz, Sungjin Im, Thomas Lavastida, Benjamin Moseley, and Sergei Vassilvitskii. Faster matchings via learned duals. *Advances in neural information processing systems*, 34:10393–10406, 2021.
- [60] Ilan Doron-Arad, Ariel Kulik, and Pasin Manurangsi. Fine grained lower bounds for multidimensional knapsack. *arXiv preprint arXiv:2407.10146*, 2024.
- [61] Ilya Dumer. Covering spheres with spheres. *Discrete & Computational Geometry*, 38:665–679, 2007.
- [62] Paul Dütting, Silvio Lattanzi, Renato Paes Leme, and Sergei Vassilvitskii. Secretaries with advice. In *Proceedings of the 22nd ACM Conference on Economics and Computation*, pages 409–429, 2021.
- [63] Paul Dütting, Evangelia Gergatsouli, Rojin Rezvan, Yifeng Teng, and Alexandros Tsigonias-Dimitriadis. Prophet secretary against the online optimal. *arXiv preprint arXiv:2305.11144*, 2023.

- [64] Benjamin Edelman, Michael Ostrovsky, and Michael Schwarz. Internet advertising and the generalized second-price auction: Selling billions of dollars worth of keywords. *American economic review*, 97(1):242–259, 2007.
- [65] Francis Y Edgeworth. The mathematical theory of banking. *Journal of the Royal Statistical Society*, 51(1):113–127, 1888.
- [66] Nelson Elhage, Neel Nanda, Catherine Olsson, Tom Henighan, Nicholas Joseph, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, et al. A mathematical framework for transformer circuits. *Transformer Circuits Thread*, 1(1):12, 2021.
- [67] Vivek F Farias, Andrew A Li, and Deeksha Sinha. Optimizing order sets in sub-linear time. In *Proceedings of the 21st ACM Conference on Economics and Computation*, pages 639–640, 2020.
- [68] Jon Feldman, Monika Henzinger, Nitish Korula, Vahab S Mirrokni, and Eli Stein. Online stochastic packing applied to display ad allocation. In *European Symposium on Algorithms*, pages 182–194. Springer, 2010.
- [69] Mingsheng Fu, Hong Qu, Dagmawi Moges, and Li Lu. Attention based collaborative filtering. *Neurocomputing*, 311:88–98, 2018.
- [70] Guillermo Gallego, Garud Iyengar, Robert Phillips, and Abhay Dubey. Managing flexible products on a network. Available at SSRN 3567371, 2004.
- [71] Gregory A Godfrey and Warren B Powell. An adaptive, distribution-free algorithm for the newsvendor problem with censored demands, with applications to inventory and distribution. *Management Science*, 47(8):1101–1112, 2001.
- [72] Negin Golrezaei, Patrick Jaillet, and Zijie Zhou. Online resource allocation with convex-set machine-learned advice. *arXiv preprint arXiv:2306.12282*, 2023.
- [73] Vineet Goyal and Ramamoorthi Ravi. An FPTAS for minimizing a class of low-rank quasi-concave functions over a convex set. *Operations Research Letters*, 41(2):191–196, 2013.
- [74] Martin Grötschel, László Lovász, and Alexander Schrijver. The ellipsoid method and its consequences in combinatorial optimization. *Combinatorica*, 1:169–197, 1981.
- [75] Anupam Gupta and Marco Molinaro. How the experts algorithm can help solve Lips online. *Mathematics of Operations Research*, 41(4):1404–1431, 2016.
- [76] Insu Han, Rajesh Jayaram, Amin Karbasi, Vahab Mirrokni, David P Woodruff, and Amir Zandieh. Hyperattention: Long-context attention in near-linear time. *arXiv preprint arXiv:2310.05869*, 2023.

- [77] Botao Hao et al. Leveraging demonstrations to improve online learning: Quality matters. In ICML, 2023.
- [78] Elad Hazan et al. Introduction to online convex optimization. *Foundations and Trends® in Optimization*, 2(3-4):157–325, 2016.
- [79] Stephen J Hoch, Eric T Bradlow, and Brian Wansink. The variety of an assortment. *Marketing Science*, 18(4):527–546, 1999.
- [80] Kurt Hornik, Maxwell Stinchcombe, and Halbert White. Multilayer feedforward networks are universal approximators. *Neural networks*, 2(5):359–366, 1989.
- [81] Piao Hu, Jiashuo Jiang, Guodong Lyu, and Hao Su. Constrained online two-stage stochastic optimization: Algorithm with (and without) predictions. arXiv preprint arXiv:2401.01077, 2024.
- [82] Chengpiao Huang and Kaizheng Wang. A stability principle for learning under non-stationarity. arXiv preprint arXiv:2310.18304, 2023.
- [83] Jakob Huber, Sebastian Müller, Moritz Fleischmann, and Heiner Stuckenschmidt. A data-driven newsvendor problem: From data to decision. *European Journal of Operational Research*, 278(3):904–915, 2019.
- [84] Woonghee Tim Huh and Paat Rusmevichientong. A nonparametric asymptotic analysis of inventory planning with censored demand. *Mathematics of Operations Research*, 34(1):103–123, 2009.
- [85] Donald L Iglehart. The dynamic inventory problem with unknown demand distribution. *Management Science*, 10(3):429–440, 1964.
- [86] Piotr Indyk and Rajeev Motwani. Approximate nearest neighbors: towards removing the curse of dimensionality. In *Proceedings of the thirtieth annual ACM symposium on Theory of computing*, pages 604–613, 1998.
- [87] Klaus Jansen, Felix Land, and Kati Land. Bounding the running time of algorithms for scheduling and packing problems. *SIAM Journal on Discrete Mathematics*, 30(1):343–366, 2016.
- [88] Jiashuo Jiang, Xiaocheng Li, and Jiawei Zhang. Online stochastic optimization with wasserstein based non-stationarity. arXiv preprint arXiv:2012.06961, 2020.
- [89] Billy Jin and Will Ma. Online bipartite matching with advice: Tight robustness-consistency tradeoffs for the two-stage model. *Advances in Neural Information Processing Systems*, 35:14555–14567, 2022.
- [90] Wang-Cheng Kang and Julian McAuley. Self-attentive sequential recommendation. In *2018 IEEE international conference on data mining (ICDM)*, pages 197–206. IEEE, 2018.

- [91] Samuel Karlin. Dynamic inventory policy with varying stochastic demands. *Management Science*, 6(3):231–258, 1960.
- [92] Zohar S Karnin and Oren Anava. Multi-armed bandits: Competing with optimal sequences. *NIPS*, 29, 2016.
- [93] Richard M Karp. *Reducibility among combinatorial problems*. Springer, 2010.
- [94] Guolin Ke, Qi Meng, Thomas Finley, Taifeng Wang, Wei Chen, Weidong Ma, Qiwei Ye, and Tie-Yan Liu. Lightgbm: A highly efficient gradient boosting decision tree. *Advances in neural information processing systems*, 30, 2017.
- [95] N Bora Keskin and Assaf Zeevi. Chasing demand: Learning and earning in a changing environment. *Mathematics of Operations Research*, 42(2):277–307, 2017.
- [96] N Bora Keskin, Xu Min, and Jing-Sheng Jeannette Song. The nonstationary newsvendor: Data-driven nonparametric learning. Available at SSRN 3866171, 2023.
- [97] Thomas Kesselheim, Andreas Tönnis, Klaus Radke, and Berthold Vöcking. Primal beats dual on online packing lps in the random-order model. In *Proceedings of the forty-sixth annual ACM symposium on Theory of computing*, pages 303–312, 2014.
- [98] Enikő Kevi and Kim Tháng Nguyen. Primal-dual algorithms with predictions for online bounded allocation and ad-auctions problems. In *International Conference on Algorithmic Learning Theory*, pages 891–908. PMLR, 2023.
- [99] Leonid G Khachiyan. Polynomial algorithms in linear programming. *USSR Computational Mathematics and Mathematical Physics*, 20(1):53–72, 1980.
- [100] Anton J Kleywegt, Alexander Shapiro, and Tito Homem-de Mello. The sample average approximation method for stochastic discrete optimization. *SIAM Journal on optimization*, 12(2):479–502, 2002.
- [101] Peter Knees, Yashar Deldjoo, Farshad Bakhshandegan Moghaddam, Jens Adamczak, Gerard-Paul Leyson, and Philipp Monreal. Recsys challenge 2019: Session-based hotel recommendations. In *Proceedings of the Thirteenth ACM Conference on Recommender Systems, RecSys '19, New York, NY, USA, 2019*. ACM. ISBN 978-1-4503-6243-6/19/09. doi: 10.1145/3298689.3346974. URL <https://doi.org/10.1145/3298689.3346974>.
- [102] Joohwan Ko and Andrew A Li. Modeling choice via self-attention. *arXiv preprint arXiv:2311.07607*, 2023.
- [103] Ariel Kulik and Hadas Shachnai. There is no eptas for two-dimensional knapsack. *Information Processing Letters*, 110(16):707–710, 2010.

- [104] Sumit Kunnumkal and Huseyin Topaloglu. Using stochastic approximation methods to compute optimal base-stock levels in inventory control problems. *Operations Research*, 56(3):646–664, 2008.
- [105] Thom Lake, Sinead A Williamson, Alexander T Hawk, Christopher C Johnson, and Benjamin P Wing. Large-scale collaborative filtering with product embeddings. *arXiv preprint arXiv:1901.04321*, 2019.
- [106] Silvio Lattanzi, Thomas Lavastida, Benjamin Moseley, and Sergei Vassilvitskii. Online scheduling via learned weights. In *Proceedings of the Fourteenth Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 1859–1877. SIAM, 2020.
- [107] Thomas Lavastida, Benjamin Moseley, R Ravi, and Chenyang Xu. Using predicted weights for ad delivery. In *SIAM Conference on Applied and Computational Discrete Algorithms (ACDA21)*, pages 21–31. SIAM, 2021.
- [108] Daniel Lee and H Sebastian Seung. Algorithms for non-negative matrix factorization. *Advances in neural information processing systems*, 13, 2000.
- [109] Sanghak Lee, Jaehwan Kim, and Greg M Allenby. A direct utility model for asymmetric complements. *Marketing Science*, 32(3):454–470, 2013.
- [110] Retsef Levi, Robin O Roundy, and David B Shmoys. Provably near-optimal sampling-based policies for stochastic inventory control models. *Mathematics of Operations Research*, 32(4):821–839, 2007.
- [111] Retsef Levi, Georgia Perakis, and Joline Uichanco. The data-driven newsvendor problem: New bounds and insights. *Operations Research*, 63(6):1294–1306, 2015.
- [112] Chengxi Li, Yejing Wang, Qidong Liu, Xiangyu Zhao, Wanyu Wang, Yiqi Wang, Lixin Zou, Wenqi Fan, and Qing Li. Strec: Sparse transformer for sequential recommendations. In *Proceedings of the 17th ACM Conference on Recommender Systems*, pages 101–111, 2023.
- [113] Xiaocheng Li and Yinyu Ye. Online linear programming: Dual convergence, new algorithms, and regret bounds. *Operations Research*, 70(5):2948–2966, 2022.
- [114] Xiaocheng Li, Chunlin Sun, and Yinyu Ye. Simple and fast algorithm for binary integer and online linear programming. *Advances in Neural Information Processing Systems*, 33:9412–9421, 2020.
- [115] Xiaocheng Li, Chunlin Sun, and Yinyu Ye. The symmetry between arms and knapsacks: A primal-dual approach for bandits with knapsacks. In *International Conference on Machine Learning*, pages 6483–6492. PMLR, 2021.

- [116] Shang Liu, Jiashuo Jiang, and Xiaocheng Li. Non-stationary bandits with knapsacks. *Advances in Neural Information Processing Systems*, 35:16522–16532, 2022.
- [117] Liwan H Liyanage and J George Shanthikumar. A practical inventory control policy using operational statistics. *Operations Research Letters*, 33(4):341–348, 2005.
- [118] William S Lovejoy. Myopic policies for some inventory models with uncertain demand distributions. *Management Science*, 36(6):724–738, 1990.
- [119] Haihao Lu. relative continuity for non-lipschitz nonsmooth convex optimization using stochastic (or deterministic) mirror descent. *INFORMS Journal on Optimization*, 1(4):288–303, 2019.
- [120] Haihao Lu, Robert M Freund, and Yurii Nesterov. Relatively smooth convex optimization by first-order methods, and applications. *SIAM Journal on Optimization*, 28(1):333–354, 2018.
- [121] R Duncan Luce. *Individual choice behavior: A theoretical analysis*. Courier Corporation, 2012.
- [122] Haipeng Luo, Chen-Yu Wei, Alekh Agarwal, and John Langford. Efficient contextual bandits in non-stationary worlds. In *Conference On Learning Theory*, pages 1739–1776. PMLR, 2018.
- [123] Thodoris Lykouris and Sergei Vassilvitskii. Competitive caching with machine learned advice. *Journal of the ACM (JACM)*, 68(4):1–25, 2021.
- [124] Mohammad Mahdian, Hamid Nazerzadeh, and Amin Saberi. Online optimization with uncertain information. *ACM Transactions on Algorithms (TALG)*, 8(1):1–29, 2012.
- [125] Spyros Makridakis and Michele Hibon. The m3-competition: results, conclusions and implications. *International journal of forecasting*, 16(4):451–476, 2000.
- [126] Spyros Makridakis, Evangelos Spiliotis, and Vassilios Assimakopoulos. M5 accuracy competition: Results, findings, and conclusions. *International Journal of Forecasting*, 38(4):1346–1364, 2022.
- [127] Yu A Malkov and Dmitry A Yashunin. Efficient and robust approximate nearest neighbor search using hierarchical navigable small world graphs. *IEEE transactions on pattern analysis and machine intelligence*, 42(4):824–836, 2018.
- [128] Reza Yousef Maragheh, Alexandra Chronopoulou, and James Mario Davis. A customer choice model with halo effect. *arXiv preprint arXiv:1805.01603*, 2018.

- [129] Leigh McAlister. A dynamic attribute satiation model of variety-seeking behavior. *Journal of consumer research*, 9(2):141–150, 1982.
- [130] Aranyak Mehta, Amin Saberi, Umesh Vazirani, and Vijay Vazirani. Adwords and generalized online matching. *Journal of the ACM (JACM)*, 54(5):22–es, 2007.
- [131] M Jeffrey Mei, Cole Zuber, and Yasaman Khazaeni. A lightweight transformer for next-item product recommendation. In *Proceedings of the 16th ACM Conference on Recommender Systems*, pages 546–549, 2022.
- [132] Zakaria Mhammedi and Wouter M Koolen. Lipschitz and comparator-norm adaptivity in online learning. In *Conference on Learning Theory*, pages 2858–2887. PMLR, 2020.
- [133] Vahab S Mirrokni, Shayan Oveis Gharan, and Morteza Zadimoghaddam. Simultaneous approximations for adversarial and stochastic online budgeted allocation. In *Proceedings of the twenty-third annual ACM-SIAM symposium on Discrete Algorithms*, pages 1690–1701. SIAM, 2012.
- [134] Shashi Mittal and Andreas S Schulz. An fptas for optimizing a class of low-rank functions over a polytope. *Mathematical Programming*, 141:103–120, 2013.
- [135] Dmitrii Moor, Yi Yuan, Rishabh Mehrotra, Zhenwen Dai, and Mounia Lalmas. Exploiting sequential music preferences via optimisation-based sequencing. In *Proceedings of the 32nd ACM International Conference on Information and Knowledge Management*, pages 4759–4765, 2023.
- [136] Marius Muja and David G Lowe. Scalable nearest neighbor algorithms for high dimensional data. *IEEE transactions on pattern analysis and machine intelligence*, 36(11):2227–2240, 2014.
- [137] Andres Munoz and Sergei Vassilvitskii. Revenue optimization with approximate bid predictions. *Advances in Neural Information Processing Systems*, 30, 2017.
- [138] Arkadij Semenov, Nemirovskij and David Borisovich Yudin. Problem complexity and method efficiency in optimization. 1983.
- [139] Yurii Nesterov and Arkadij Nemirovskii. Interior-point polynomial algorithms in convex programming. SIAM, 1994.
- [140] Trung Thanh Nguyen and Khaled Elbassioni. A ptas for a class of binary non-linear programs with low-rank functions. *Operations Research Letters*, 49(5):633–638, 2021.
- [141] Francesco Orabona. Dimension-free exponentiated gradient. *Advances in Neural Information Processing Systems*, 26, 2013.

- [142] Francesco Orabona and Ashok Cutkosky. Icml tutorial on parameter-free stochastic optimization. ICML, 2020. URL: <https://parameterfree.com/icml-tutorial/>.
- [143] Afshin Oroojlooyjadid, Lawrence V Snyder, and Martin Takáč. Applying deep learning to the newsvendor problem. IIE Transactions, 52(4):444–463, 2020.
- [144] Jorge Pérez, Javier Marinković, and Pablo Barceló. On the turing completeness of modern neural network architectures. arXiv preprint arXiv:1901.03429, 2019.
- [145] Warren Powell, Andrzej Ruszczyński, and Huseyin Topaloglu. Learning algorithms for separable approximations of discrete stochastic optimization problems. Mathematics of Operations Research, 29(4):814–836, 2004.
- [146] Manish Purohit, Zoya Svitkina, and Ravi Kumar. Improving online algorithms via ml predictions. Advances in Neural Information Processing Systems, 31, 2018.
- [147] David J Rader Jr and Gerhard J Woeginger. The quadratic 0/1 knapsack problem with series parallel support. Operations Research Letters, 30(3):159–166, 2002.
- [148] Omid Rezaei. Optimizing user engagement through adaptive ad sequencing. Marketing Science, 42(5):910–933, 2023.
- [149] Dhruv Rohatgi. Near-optimal bounds for online caching with machine learned advice. In Proceedings of the Fourteenth Annual ACM-SIAM Symposium on Discrete Algorithms, pages 1834–1845. SIAM, 2020.
- [150] Francisco JR Ruiz, Susan Athey, and David M Blei. Shopper. The Annals of Applied Statistics, 14(1):1–27, 2020.
- [151] Clayton Sanford, Daniel Hsu, and Matus Telgarsky. One-layer transformers fail to solve the induction heads task. arXiv preprint arXiv:2408.14332, 2024.
- [152] Clayton Sanford, Daniel J Hsu, and Matus Telgarsky. Representational strengths and limitations of transformers. Advances in Neural Information Processing Systems, 36, 2024.
- [153] Herbert Scarf. Bayes solutions of the statistical inventory problem. The annals of mathematical statistics, 30(2):490–508, 1959.
- [154] Herbert Scarf, K Arrow, S Karlin, and P Suppes. The optimality of (s, s) policies in the dynamic inventory problem. Optimal pricing, in action, and the cost of price adjustment, pages 49–56, 1960.
- [155] A Schrijver. Theory of linear and integer programming. John Wiley & Sons, 1998.

- [156] Shai Shalev-Shwartz and Shai Ben-David. Understanding machine learning: From theory to algorithms. Cambridge university press, 2014.
- [157] Alexander Shapiro. Monte carlo sampling methods. Handbooks in operations research and management science, 10:353 425, 2003.
- [158] Matthew Streeter and H Brendan McMahan. No-regret algorithms for unconstrained online convex optimization. arXiv preprint arXiv:1211.2260, 2012.
- [159] Sainbayar Sukhbaatar, Jason Weston, Rob Fergus, et al. End-to-end memory networks. Advances in neural information processing systems, 28, 2015.
- [160] Kalyan T Talluri and Garrett J Van Ryzin. The theory and practice of revenue management, volume 68. Springer Science & Business Media, 2006.
- [161] Richard Taylor. Approximation of the quadratic knapsack problem. Operations Research Letters, 44(4):495 497, 2016.
- [162] Sean J Taylor and Benjamin Letham. Forecasting at scale. The American Statistician, 72(1):37 45, 2018.
- [163] James T Treharne and Charles R Sox. Adaptive inventory control for nonstationary demand and partial information. Management Science, 48(5):607 624, 2002.
- [164] Pravin M Vaidya. A new algorithm for minimizing convex functions over convex sets. Mathematical programming, 73(3):291 341, 1996.
- [165] A Vaswani. Attention is all you need. Advances in Neural Information Processing Systems, 2017.
- [166] Stephen A Vavasis. On the complexity of nonnegative matrix factorization. SIAM journal on optimization, 20(3):1364 1377, 2010.
- [167] Erik Vee, Sergei Vassilvitskii, and Jayavel Shanmugasundaram. Optimal online assignment with forecasts. In Proceedings of the 11th ACM conference on Electronic commerce, pages 109 118, 2010.
- [168] Jean-Louis Verger-Gaugry. Covering a ball with smaller equal balls in Discrete & Computational Geometry, 33:143 155, 2005.
- [169] Roman Vershynin. High-dimensional probability: An introduction with applications in data science, volume 47. Cambridge university press, 2018.
- [170] Hanzhao Wang, Xiaocheng Li, and Kalyan Talluri. Transformer choice net: A transformer neural network for choice prediction. arXiv preprint arXiv:2310.08716, 2023.

- [171] Shoujin Wang, Liang Hu, Longbing Cao, Xiaoshui Huang, Defu Lian, and Wei Liu. Attention-based transactional context embedding for next-item recommendation. In Proceedings of the AAAI conference on artificial intelligence, volume 32, 2018.
- [172] Yu-Xiong Wang and Yu-Jin Zhang. Nonnegative matrix factorization: A comprehensive review. *IEEE Transactions on knowledge and data engineering*, 25(6):1336–1353, 2012.
- [173] Colin Wei, Yining Chen, and Tengyu Ma. Statistically meaningful approximation: a case study on approximating turing machines with transformers. *Advances in Neural Information Processing Systems*, 35:12071–12083, 2022.
- [174] Timo Wilm, Philipp Normann, Sophie Baumeister, and Paul-Vincent Kobow. Scaling session-based transformer recommendations using optimized negative sampling and loss functions. In Proceedings of the 17th ACM Conference on Recommender Systems, pages 1023–1026, 2023.
- [175] Peter R Winters. Forecasting sales by exponentially weighted moving averages. *The Use of MMR, Diversity-Based Reranking for Reordering Documents and Producing Summaries*, 6(3):324–342, 1960.
- [176] Yihong Wu. Lecture 14: Covering and packing numbers. <http://www.stat.yale.edu/~yw562/teaching/598/lec14.pdf>, 2017. STAT 598: High Dimensional Statistics, Yale University.
- [177] Boming Yang, Dairui Liu, Toyotaro Suzumura, Ruihai Dong, and Irene Li. Going beyond local: Global graph-enhanced personalized news recommendations. In Proceedings of the 17th ACM Conference on Recommender Systems, pages 24–34, 2023.
- [178] Tianbao Yang, Lijun Zhang, Rong Jin, and Jinfeng Yi. Tracking slowly moving clairvoyant: Optimal dynamic regret of online learning with true and noisy gradient. In International Conference on Machine Learning, pages 449–457. PMLR, 2016.
- [179] Chulhee Yun, Srinadh Bhojanapalli, Ankit Singh Rawat, Sashank J Reddi, and Sanjiv Kumar. Are transformers universal approximators of sequence-to-sequence functions? *arXiv preprint arXiv:1912.10077*, 2019.
- [180] Lijun Zhang, Shiyin Lu, and Zhi-Hua Zhou. Adaptive online learning in dynamic environments. *NIPS*, 31, 2018.
- [181] Luhao Zhang, Jincheng Yang, and Rui Gao. Optimal robust policy for feature-based newsvendor. *Management Science*, Forthcoming, 2023.
- [182] Zhi Zheng, Ying Sun, Xin Song, Hengshu Zhu, and Hui Xiong. Generative learning plan recommendation for employees: A performance-aware reinforcement learning approach. In Proceedings of the 17th ACM Conference on Recommender Systems, pages 443–454, 2023.

Appendix A

APPENDIX FOR CHAPTER 2

A.1 Proofs in Section 2.2

First, we provide some previous results on n -ANN algorithms. [8] gives an n -ANN algorithm by using Locality-Sensitive Hashing:

Proposition 13 (Corollary 1 in [8]). For any given $\epsilon > 0$, n -ANN on the unit sphere $S^{3-1} \subset \mathbb{R}^3$ can be solved with expected amortized runtime $O(n^{1+\epsilon})$, where $W^{1+\epsilon} = \frac{1}{1-2\epsilon}$, where $2 - \epsilon$ is a universal constant.

[2] gives an n -ANN algorithm by applying fast Johnson-Lindenstrauss transform (FJLT):

Proposition 14 (Theorem 1 in [2]). For any given $\epsilon > 0$, n -ANN on the unit sphere $S^{3-1} \subset \mathbb{R}^3$ can be solved with expected amortized runtime $O(n^3 \log^{2+\epsilon} n)$.

[12] gives an n -ANN algorithm by using a tree-based data structure:

Proposition 15 (Theorem 3.1 in [12]). For any given $\epsilon > 0$, n -ANN on the unit sphere $S^{3-1} \subset \mathbb{R}^3$ can be solved with expected amortized runtime $O(n \log^3 n)$.

Proof of Lemma 1. We prove that this can be done by applying Approximate Nearest Neighbor (ANN) algorithm from [2]. Given a set of n data points $\{x_1, \dots, x_n\} \subset \mathbb{R}^3$, the goal of the ϵ -Approximate Nearest Neighbor algorithm is to build a data structure that, given a query $q \in \mathbb{R}^3$, returns a data point x_i such that $\|x_i - q\| \leq (1 + \epsilon) \|x_{i^*} - q\|$ for all $x_{i^*} \in \{x_1, \dots, x_n\}$.

Lemma 4 (Theorem 1 in [2]). Algorithm 1 in [2] solves the ϵ -ANN problem in query time $O(n^3 \log^{2+\epsilon} n)$.

Moreover, the Algorithm 1 in [2] can be easily modified to return data points $\{x_{i_1}, \dots, x_{i_k}\}$ such that $\|x_{i_j} - q\| \leq (1 + \epsilon)^j \|x_{i^*} - q\|$ for all $j = 1, \dots, k$ and all $x_{i^*} \in \{x_1, \dots, x_n\}$. This can be done by running their Algorithm 1 for k times, where each time we ignore the data points that have already been returned in previous runs. More specifically, in their Algorithm 1 we change the "else if

... else if ... or ?⁰ has been returned in previous rounds', which gives our desired performance.

Because the guarantee of Lemma 1 is a multiplicative guarantee on distance, below we show how to convert it into an additive guarantee on the inner products.

We first reduce E_1, \dots, E_n and D to the unit sphere. Let $E_{\max} = \max_k \|E_k\|_2$. Let

$$E_k^0 = \frac{1}{E_{\max}} E_k \quad \text{for } k = 1, \dots, n$$

for each $k = 1, \dots, n$, where we append a scalar 1 and rescale the vector. Let $D^0 = \frac{1}{D_{\max}} D \in \mathbb{R}^n$. Then $\|E_k^0\|_2 \leq 1$. Moreover,

$$\begin{aligned} \|D^0\|_2 &= \|D\|_2, \|E_k^0\|_2 \leq \|E_k\|_2 / E_{\max} \\ &= 2 \|D\|_2 / E_{\max} \end{aligned}$$

Therefore $\|D^0\|_2 \leq \frac{2}{E_{\max}} \|D\|_2$. Hence if $\|D^0\|_2 \leq \frac{1}{2} \|D\|_2$, then $\|D\|_2 \leq 2 \|D^0\|_2$. Hence, we can run the Algorithm 1 [2] with inputs E_1^0, \dots, E_n^0 and $\lambda = \frac{8^p 2n}{1 + D_{\max}^0}$, which outputs indices i_1, \dots, i_n such that $\|D^0\|_2 \leq \frac{1}{2} \|D\|_2$ for each $i = 1, \dots, n$. Then

$$\begin{aligned} \|D^0\|_2 &= \|D^0\|_2 \\ &= \|D^0\|_2 \\ &= \|D^0\|_2 \\ &= \frac{8^p 2n}{1 + D_{\max}^0} \end{aligned}$$

where the first equality follows from the definition of D^0 , the second equality and the second inequality follows since $\|D^0\|_2 \leq \|D\|_2$. Therefore we have $\|D^0\|_2 \leq \|D\|_2$ for each $i = 1, \dots, n$: as desired. ..

Proof of Proposition 2. The algorithm ALG operates as follows:

1. Partition the items according to their reward functions.
2. For each partition, apply the given Approximate-Nearest Neighbor algorithm to identify items that are collectively among the most attractive to the user within that partition.

Algorithm 8: n -Approximate ϵ -Nearest Neighbor in Lemma 1

PREPROCESS

Input: $E_1, \dots, E_n \subseteq \mathbb{R}^3$ and $n \geq 0$;

$E_{\max} = \max_{q \in \mathbb{R}^3} \max_{1 \leq i \leq n} \|E_i - q\|_2$;

$E_{\delta}^0 = \max_{q \in \mathbb{R}^3} \max_{1 \leq i \leq n} \|E_i - q\|_2 - \delta \cdot E_{\max}$ for each $\delta = 1, \dots, \frac{1}{\epsilon}$;

$\lambda = \frac{1}{\epsilon} \cdot \frac{1}{2n+2}$, where $\frac{1}{2}$ is an upper bound on $\max_{i,j} |D_{ij}|$ for all possible D ;

Run PREPROCESS of the λ -ANN algorithm in [2] with inputs $E_1, \dots, E_n \subseteq \mathbb{R}^3$ and $\lambda \geq 0$.

QUERY

Input: $D \subseteq \mathbb{R}^3$;

$D^0 = \max_{q \in \mathbb{R}^3} \|D - q\|_2$;

$\rho = 1$;

while $\rho > \epsilon$: do

Run PREPROCESS of the λ -ANN algorithm in [2] with input D with the modification of changing λ to $\lambda \cdot \frac{1}{\rho}$ if $\rho > \frac{1}{2}$ or $\lambda \cdot \frac{1}{\rho} > \frac{1}{2}$ or $\lambda \cdot \frac{1}{\rho} < \frac{1}{2}$;

ρ_j the index of the output of the λ -ANN algorithm in [2];

Return ρ_j .

- Evaluate the rewards of all such candidate items across partitions and select the top- ϵ items with the highest overall reward.

We analyze the performance of ALG. Let $\{g_1, \dots, g_\ell\}$ be a partition of the index set $[n]$ such that $|g_i| = \frac{n}{\ell}$ for every $g_i \subseteq [n]$ and $\ell \geq \frac{1}{\epsilon}$. We first construct an index set $g^0 = \{g_i^0\}$ for each $g_i \subseteq [n]$ and then combine them to obtain the set of all candidate items \mathcal{C} .

For each index set g_i^0 , we only choose indices out of it to include in g_i^0 , namely the ϵ indices that are approximately the highest indices in $E_{\delta}^0 \cap g_i^0$.¹ Specifically, for each $g_i^0 \subseteq [n]$, we run the given n -Approximate ϵ -Nearest Neighbor oracle with given set of points $E_{\delta}^0 \cap g_i^0 \subseteq \mathbb{R}^3$, and query D , numbers λ and n as inputs. We let g_i^0 be the collection of all output indices for each $g_i^0 \subseteq [n]$. Then $\mathcal{C} = \bigcup_{g_i^0} g_i^0$: for each g_i^0 . Let $\mathcal{C} = \bigcup_{g_i^0} g_i^0$, then $|\mathcal{C}| \leq \frac{n}{\epsilon}$.

Let $(x^* = f_{1-\epsilon, \dots, \epsilon})$ be the optimal solution to Problem (Pure Embedding) (for simplicity we assume $\epsilon = \frac{1}{2}$, and other cases can be handled similarly). We show that $\text{ALG}_{\text{Pure Embedding}} \leq \frac{1}{1-\epsilon} \cdot \frac{1}{\epsilon} \cdot \text{OPT}_{\text{Pure Embedding}} \leq \frac{1}{1-\epsilon}$. For each $\epsilon = 1-\dots, \frac{1}{\epsilon}$, let ρ_j be the index such that x_{ρ_j} and x_{ρ_j} are in the same g_i^0 and

¹For simplicity we assume $|g_i^0| \geq \frac{n}{\epsilon}$. Otherwise we simply choose all indices in g_i^0 .

$M = f^1 9 - < - 0^0 g$ These rows have non-zero keys and scalar values. They have zero queries. They encode the item q_i at position $<$ and its similarity embedding $G_{<-0}$

$D = f g$: This is a dummy row that dominates the softmax denominator, so that the softmax vector behaves like division by a constant.

$B = f^{18} - C^{base} g$: These rows encode the base utility of each item q_i .

^ The embedding dimension is set $d_{emb} = d_{pos} + d_{comp} + d_{base} + 1$. This is split into a position block of length d_{pos} , a component block of length d_{comp} , a single dummy column (denoted d_{base}), and a single base column (denoted d_{base}).

^ Let $\epsilon > 0$ be a sufficiently large constant and $\eta \in \mathbb{R}$ a fixed constant. Later we will take ϵ large enough so that the simple transformer approximates $6^1(-\epsilon)$ to arbitrary precision.

^ For each position $C \in \{1, \dots, 4\}$, define the row vector A_C by

$$A_C = \begin{bmatrix} \log_{<-C} 1 \\ \vdots \\ \log_{<-C} 2 \\ \vdots \\ \log_{<-C} 4 \end{bmatrix} \in \mathbb{R}^{<}$$

i.e.

$$A_C = \begin{bmatrix} \log_{<-C} 1 & \log_{<-C} 2 & \dots & \log_{<-C} 4 \end{bmatrix} \in \mathbb{R}^{<}$$

For each $C \in \{1, \dots, 4\}$, let $R^C \in \mathbb{R}^{< \times <}$ be the matrix with all rows equal to A_C

$$R^C = \begin{bmatrix} A_C \\ \vdots \\ A_C \end{bmatrix} \in \mathbb{R}^{< \times <} \quad (= \text{rows}).$$

Then the position block of the query matrix is the vertical stacking of these blocks:

$$\&_{pos} = \begin{bmatrix} R^1 \\ R^2 \\ R^3 \\ R^4 \end{bmatrix} \in \mathbb{R}^{< \times <}$$

Next, define the $n \times 3$ matrix collecting the (component-wise) logarithms of the similarity embeddings $G^1, G_{-1}, \dots, G_{-3}$ ($n \times 3$):

$$= \begin{pmatrix} \log G^0 & \log G_{-1} & \log G_{-2} & \log G_{-3} \\ \log G^0 & \log G_{-1} & \log G_{-2} & \log G_{-3} \\ \vdots & \vdots & \vdots & \vdots \\ \log G^0 & \log G_{-1} & \log G_{-2} & \log G_{-3} \end{pmatrix}$$

The component block of the query matrix is identical copies of stacked vertically:

$$\&_{\text{cmp}} = \begin{pmatrix} 1 \\ 1 \\ 1 \\ \vdots \\ 1 \\ 1 \\ 1 \\ 1 \end{pmatrix} 2 R^{1 \times 3}$$

Let $\mathbf{1} \in \mathbb{R}^{n \times 1}$ be the all-ones column, and $\mathbf{0} \in \mathbb{R}^{n \times 1}$ be the all-zeros column. Define the dummy and base query columns as

$$\& = \begin{pmatrix} \mathbf{1} \\ \mathbf{0} \\ 0 \\ \mathbf{0} \end{pmatrix} \quad \& = \begin{pmatrix} \mathbf{1} \\ \mathbf{0} \\ 0 \\ \mathbf{0} \end{pmatrix} 2 R^{\# 1}$$

Putting the query blocks together, we define $\&$ to be

$$\& = \begin{pmatrix} \&_{\text{pos}} & \&_{\text{cmp}} & \& & \& \end{pmatrix} 2 R^{\# 1, 3, 2^0}$$

^ The position block of the key matrix is defined as

$$\text{pos} = \begin{pmatrix} \mathbf{1}_{=3} & 0 & 0 & 0 \\ 0 & \mathbf{1}_{=3} & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \mathbf{1}_{=3} \end{pmatrix} 2 R^{1 \times 3^0}$$

i.e., a block-diagonal matrix with diagonal blocks, each block the column $\mathbf{1}_{=3}$.

Let $E \in \mathbb{R}^{3 \times 3}$ be the identity matrix. The component block of the key matrix is defined as

$$cmp = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} \in \mathbb{R}^{3 \times 3}$$

Thus, in the block associated with position and item, the 3 consecutive rows equal the identity I_3 .

The dummy and base key columns are

$$= \begin{pmatrix} 0 \\ 0 \\ 1 \\ 0 \end{pmatrix} \quad = \begin{pmatrix} 0 \\ 0 \\ 0 \\ 1 \end{pmatrix} \in \mathbb{R}^4$$

Putting the key blocks together, we define to be

$$= \begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix} \in \mathbb{R}^{4 \times 3}$$

Let $E \in \mathbb{R}^{3 \times 3}$:

$$E = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} \in \mathbb{R}^{3 \times 3}$$

We define α (one scalar per row) by placing zeros on and D , the item component entries on M , and the base utilities on B :

$$+ = \begin{pmatrix} \alpha_1 & \dots & \alpha_3 \\ \alpha_1 & \dots & \alpha_3 \\ \alpha_1 & \dots & \alpha_3 \end{pmatrix} \in \mathbb{R}^{3 \times 3}$$

where $4^{n-3} \frac{1}{V} \log D \dots \frac{1}{V} \log D$ is repeated times. Set $D = 1$ so that $\frac{1}{8} D = \frac{1}{8}$ for each row s .

- ^ We set $\frac{1}{8} G^0 = \exp^{1/V} G^0$ for all $s = 1, 2, \dots, 3, 1, \dots, 1/4$.
- ^ For a length- n sequence $x = x_1 \dots x_n$ let $(s = 1, 2, \dots, 3, 1, \dots, 1/4 =$
 $\dots; M; D; B/4$ such that $(s$ contains $\frac{1}{8} G^0 : C = 2 \dots 1/4 g$ from the first
 set of rows, $\frac{1}{8} G^0 : C = 2 \dots 1/4 - 0 = 2 \dots 3/4 g$ from the second set of
 rows, the dummy row D , and $\frac{1}{8} G^0$ from the set of base
 rows.

We show that the simple transformer above approximates σ to arbitrary precision. The intuition of our construction is given below:

- ^ $\&$ encodes lags ($_$) and similarity embeddings ($\log G$)
- ^ $_$ encodes positions (C) in the sequence.
- ^ $+$ encodes similarity embeddings ($\log G$) and base utilities ($\frac{1}{8}$)
- ^ The dummy row makes softmax act like a constant divider.

Fix a length- n sequence $x = x_1 \dots x_n$. We first compute $\frac{1}{8} G^0$ for each $\frac{1}{8} G^0$

(context) $\frac{1}{8} G^0 = \frac{1}{8} G^0$

(dummy) $\frac{1}{8} G^0 = \frac{1}{8}$

(base) $\frac{1}{8} G^0 = 0$

Hence the denominator of the softmax operation on row $s = 1, 2, \dots, 3, 1, \dots, 1/4$ is

$$Z_C = 4^{n-3} \left(\frac{1}{8} G^0 + \frac{1}{8} G^0 + \dots + \frac{1}{8} G^0 \right)$$

Define the following (nite) constants

$$c = \frac{1}{8} G^0, \quad c^0 = \frac{1}{8} G^0, \quad c^1 = \frac{1}{8} G^0$$

Set

$$X_C = 4^{n-3} \left(\frac{1}{8} G^0 + \frac{1}{8} G^0 \right)$$

Then

$$/C = 4^{n-3} 1, X_C^0.$$

Therefore the attention weights on row \mathbf{e}^0 are

$$\begin{aligned} \text{softmax}_{\mathbf{e}^0}(\mathbf{A}^1 - \mathbf{C}^0) &= \frac{\exp(\mathbf{A}^1 \mathbf{e}^0, \log G_{\mathbf{e}^0})}{/C} \\ &= \frac{-C < G_{\mathbf{e}^0}}{4^{n-3} 1, X_C^0} \end{aligned}$$

and

$$\text{softmax}_{\mathbf{e}^0}(\mathbf{A}^1 - \mathbf{C}^0) = \frac{\exp(10^0)}{/C} = \frac{4^{n-3}}{1, X_C}$$

Recall that

$$+_{\mathbf{e}^0} = 4^{n-3} G_{\mathbf{e}^0} + = 0 +_{\mathbf{e}^0} = 4^{n-3} \frac{1}{V} \log D_{\mathbf{e}^0}$$

So we have

$$\begin{aligned} SA_{\mathbf{e}^0} &= \text{softmax}_{\mathbf{e}^0}(\mathbf{A}^1 - \mathbf{C}^0) +_{\mathbf{e}^0} \\ &= \frac{1}{1, X_C} \text{softmax}_{\mathbf{e}^0}(\mathbf{A}^1 - \mathbf{C}^0) +_{\mathbf{e}^0} \\ &= \frac{1}{1, X_C} \frac{-C < G_{\mathbf{e}^0}}{4^{n-3}} +_{\mathbf{e}^0} \\ &= \frac{1}{1, X_C} \frac{-C < G_{\mathbf{e}^0}}{4^{n-3}} +_{\mathbf{e}^0} \\ &= \frac{1}{1, X_C} \frac{-C < G_{\mathbf{e}^0}}{4^{n-3}} +_{\mathbf{e}^0} \end{aligned}$$

Fix any $n \geq 0$. Since $\binom{10}{C} < 1$ and $\binom{10}{C} < 3 \cdot 4^{n-3}$, we can set n large enough so that

$$4^{n-3} > 1, \quad \binom{10}{C} < 3 \cdot 4^{n-3},$$

which gives $X_C > n$ for every C . Finally, we have

$$\begin{aligned} T_{\mathbf{e}^0} &= \exp \frac{V}{1, X_C} \frac{-C < G_{\mathbf{e}^0}}{4^{n-3}} +_{\mathbf{e}^0} \\ &= D_{\mathbf{e}^0}^{\frac{1}{X_C}} \exp \frac{V}{1, X_C} \frac{-C < G_{\mathbf{e}^0}}{4^{n-3}} +_{\mathbf{e}^0} \end{aligned}$$

We have

$$\text{softmax} - \left(\mathbf{1}^{10} \& \mathbf{1}^{10} \mathbf{1} - \left(\mathbf{1}^{10} \mathbf{1}^{10} \right) \right)_{89}$$

$$= \begin{matrix} \frac{\exp^{1 \ 8g}}{92(\exp^{1 \ 8g}, \exp^{1^{10}}, \mathbf{1} = j(j^0))} & 8 \ 2 \ (- \ 9 \ 2 \ (- \\ \frac{\exp^{1^{10}}}{92(\exp^{1 \ 8g}, \exp^{1^{10}}, \mathbf{1} = j(j^0))} & 8 \ 2 \ (- \ 9 \ = \ , \ 1- \\ \frac{1}{92(\exp^{1 \ 8g}, \exp^{1^{10}}, \mathbf{1} = j(j^0))} & 8 \ 2 \ (- \ 9 \ 2 \ \gg \frac{1}{4} \ n \ (- \\ \frac{1}{\exp^{1^{10}}} & 8 \ = \ , \ 1- \ 9 \ 2 \ \gg \frac{1}{4} - \\ \frac{\exp^{1^{10}}}{\exp^{1^{10}}} & 8 \ = \ , \ 1- \ 9 \ = \ , \ 1- \\ \frac{1}{\mathbf{1}} & 8 \ 2 \ \gg \frac{1}{4} \ n \ (- \ 9 \ 2 \ \gg \ , \ 1 \frac{1}{4} \bullet \end{matrix}$$

Therefore, for every $8 \ 2 \ ($, we have

$$SA_{\& \mathbf{1}^{10} - \mathbf{1}^{10} \rightarrow \mathbf{1}^{10} \mathbf{1} - \left(\mathbf{1}^{10} \right)_{8} D^{10}} = \frac{\int \frac{92(\exp^{1 \ 8g}}{92(\exp^{1 \ 8g}, \exp^{1^{10}}, \mathbf{1} = j(j^0))} \bullet$$

For any given $n \ j \ 0$, we can take n to be sufficiently large such that

$$\frac{\int \frac{92(\exp^{1 \ 8g} \ \frac{n}{2}}{\exp^{1^{10}}} \ SA_{\& \mathbf{1}^{10} - \mathbf{1}^{10} \rightarrow \mathbf{1}^{10} \mathbf{1} - \left(\mathbf{1}^{10} \right)_{8} D^{10}} \ \frac{\int \frac{92(\exp^{1 \ 8g}}{\exp^{1^{10}}} \bullet$$

Finally, we get

$$\tilde{O} \ \frac{\exp^{1 \ 8g}}{92(\ \ \frac{n}{3} \ T_{\& \mathbf{1}^{10} - \mathbf{1}^{10} \rightarrow \mathbf{1}^{10} \mathbf{1} - \left(\mathbf{1}^{10} \right)_{8} D^{10}} \ \tilde{O} \ \frac{\exp^{1 \ 8g}}{92(\ \ \bullet$$

Head 2. The second self-attention head has exactly the same parameters, except we replace $\mathbf{1}$ with $\mathbf{0}$ and we flip the sign of $\mathbf{5}$

$\hat{\&}^{120} - \mathbf{1}^{20} \ 2 \ R^{1,10 \ 1,10}$ are set such that

$$\&^{120 \ 1 \ \mathbf{1}^{20} \ 0} = \begin{matrix} \begin{matrix} 2 \\ \vdots \\ \vdots \\ \vdots \\ 4 \end{matrix} & \left| \begin{matrix} \mathbf{1} \ 3 \\ \vdots \\ \vdots \\ \vdots \\ \mathbf{1} \ 5 \end{matrix} \end{matrix}$$

$\hat{\&}^{120} \ 2 \ R^{1,10 \ 1,10}$ is set to $\mathbf{1}^{20} = 1$ for each $8 \ 2 \ \gg \ \mathbf{1}$ and $\mathbf{1}^{20} = 0$. Set $D^{120} = 1$ so that $\mathbf{1}^{10} \ 0 \ D^{120} = \mathbf{1}^{20}$.

^ Set $\frac{1}{8}^{2^0} G^0 = \exp^{1^{1^0}} G$ for every $8 \geq 2 \gg \frac{1}{4}$.

^ For a subset $(\gg \frac{1}{4}$, we set $\tilde{f}^{2^0} = \frac{1}{8}$ where $\tilde{f}^{2^0} = (\frac{1}{8} \exp^{1^{1^0}} G)$.

Then, similar to head 1, we get

$$\tilde{O} \exp^{1^{1^0}} \frac{n}{3} \tilde{T}_{\&^{1^{2^0} - 1^{2^0} + 1^{2^0} - \frac{1}{5}^{2^0} \dots \frac{1}{5}^{2^0} - D^{2^0} 1^{1^0} - (1^{2^0} \frac{1}{8}} \tilde{O} \exp^{1^{1^0}} \frac{n}{3}$$

Head 3. The third self-attention head has the following parameters:

^ $\&^{1^{3^0} - 1^{3^0}} \frac{2}{R^3}$ are set such that

$$\&^{1^{3^0} - 1^{3^0}} \frac{2}{R^3} = \begin{matrix} 2 & 0 & " & " & 3 \\ " & 0 & " & " & \frac{1}{3} \\ \vdots & \vdots & \cdot & \cdot & \frac{1}{3} \\ 4 & " & " & 0 & \frac{1}{3} \\ & & & & 5 \end{matrix}$$

^ $\frac{1}{8}^{1^{3^0}} \frac{2}{R^3}$ is set to $\frac{1}{8}^{1^{3^0}} = \frac{1}{8}$ for each $8 \geq 2 \gg \frac{1}{4}$.

^ $D^{1^{3^0}} \frac{2}{R^3}$ is set to $D^{1^{3^0}} = D$.

^ Set $\frac{1}{8}^{1^{3^0}} G^0 = G$ for every $8 \geq 2 \gg \frac{1}{4}$.

^ For a subset $(\gg \frac{1}{4}$, we set $\tilde{f}^{1^{3^0}} = (\frac{1}{8} \exp^{1^{1^{3^0}}} G)$.

For every $8 \geq 2$, we have

$$SA_{\&^{1^{3^0} - 1^{3^0} + 1^{3^0} - \frac{1}{5}^{3^0} \dots \frac{1}{5}^{3^0} - D^{3^0} 1^{1^0} - (1^{3^0} \frac{1}{8}} D^{1^{3^0}} = \frac{1^{1^{3^0}} \frac{1}{8} \frac{1}{3} \exp^{1^{1^0}}}{1^{1^0} \frac{1}{8} \frac{1}{3} \exp^{1^{1^0}}}$$

For any given $n \geq 0$, we can take n to be sufficiently large such that

$$\frac{1^{1^{3^0}} \frac{1}{8} \frac{1}{3} \exp^{1^{1^0}}}{1^{1^0} \frac{1}{8} \frac{1}{3} \exp^{1^{1^0}}} \tilde{T}_{\&^{1^{3^0} - 1^{3^0} + 1^{3^0} - \frac{1}{5}^{3^0} \dots \frac{1}{5}^{3^0} - D^{3^0} 1^{1^0} - (1^{3^0} \frac{1}{8}} \frac{n}{3} \tilde{O} \exp^{1^{1^0}} \frac{n}{3}$$

Complete the Proof. Because $\exp^{1^{1^0}} \frac{n}{3} \exp^{1^{1^0}} \frac{n}{3} = \frac{n}{3}$ Combining the above three self-attention heads, we have

$$\tilde{O} \exp^{1^{1^0}} \frac{n}{3} \tilde{T}_{\&^{1^{3^0} - 1^{3^0} + 1^{3^0} - \frac{1}{5}^{3^0} \dots \frac{1}{5}^{3^0} - D^{3^0} 1^{1^0} - (1^{3^0} \frac{1}{8}} \tilde{O} \exp^{1^{1^0}} \frac{n}{3}$$

Therefore the three attention heads approximate 6^{1^0} to arbitrary precision.

"

A.3 Relationship of Model 2 and Choice Models

Model 2 is closely related to choice models that generate conversion probabilities. Choice models are a fundamental input to many now-canonical optimization problems in the field of operations management, including assortment, inventory, and price optimization. At a high level, a choice model maps an ordered set \mathcal{S} to conversion probabilities $\{p_j\}_{j \in \mathcal{S}}$ that sum to one. Interpreting β_j from Model 2 as an unnormalized log-weight, a choice model can naturally set purchase probabilities proportional to $\exp(\beta_j)$. This precisely recovers a particular type of choice model, called the Halo Multinomial Logit (Halo MNL) choice model [128].

The Halo MNL choice model is a generalization of the simplest and most widely used multinomial logit (MNL) choice model, which assumes that customers choose among a set of items according to a fixed utility associated with each option. While the MNL choice model provides a tractable and interpretable framework, it has well-known limitations, such as the independence of irrelevant alternatives property, which can fail to capture context-dependent or irrational choice behaviors observed in practice. The Halo MNL choice model addresses this by allowing the utility of each item to depend on the presence or absence of other items in the ordered set. In particular, certain items can impose a positive or negative halo effect on the utility of other items, enabling the model to capture behaviors that violate classical rationality assumptions.

The halo effect is completely parametrized by a matrix $\mathbf{R} \in \mathbb{R}^{n \times n}$, where each off-diagonal entry R_{ij} quantifies the halo effect that the presence of item i imposes on the utility of item j : positive values capture complementarity (item i makes j more attractive), while negative values capture substitution (item i detracts from j). Crucially, these interactions are precisely represented by Model 2: the interaction-adjusted scores serve as the weights of the conversion probabilities. Normalizing these weights over the ordered set (e.g., via a softmax) yields the Halo MNL choice model with parameter β . Moreover, when \mathbf{R} is the zero matrix, there are no interaction terms and the normalization reduces to the standard MNL choice model. Therefore the MNL choice model is a special case nested within our framework.

A.4 Graph Interpretation of Self-attention Layers.

Let

$$\mathbf{p} = \text{softmax}(\mathbf{1} - \mathbf{R}^{-1} \mathbf{1}) \in \mathbb{R}^n$$

Because each row of A sums up to one, A can be interpreted as the weighted adjacency matrix of a directed graph with edge weights equal to A_{ij} . Then a single self-attention layer performs one round of information passing:

$$SA_{i-} = \sum_{j=1}^n \tilde{O}_{ij} A_{ij} v_j$$

i.e., vertex i aggregates value vectors from its (out-)neighbors according to the edge weights in A .

This perspective also clarifies what stacking self-attention layers does and does not do. For $k=1, 2$, define

$$v_i^{(k)} = \text{softmax}_i(A^{(k-1)} v^{(k-1)})$$

where we ignore point-wise activations functions for intuition. Expanding $v_i^{(2)}$ shows that information propagates along two-hop paths:

$$v_i^{(2)} = \sum_{j=1}^n \tilde{O}_{ij} v_j^{(1)} = \sum_{j=1}^n \tilde{O}_{ij} \sum_{k=1}^n \tilde{O}_{jk} v_k^{(0)} = \sum_{k=1}^n \left(\sum_{j=1}^n \tilde{O}_{ij} \tilde{O}_{jk} \right) v_k^{(0)}$$

so the contribution of vertex k to vertex i through vertex j factorizes as $\tilde{O}_{ij} \tilde{O}_{jk}$.

This is a composition of two pairwise interactions, not an arbitrary triplet interaction. Indeed, representing a general triplet interaction requires $\Theta(n^3)$ free parameters, whereas two self-attention layers only provide $\Theta(n^2)$ free parameters. Therefore, stacked self-attention layers naturally model multi-hop effects rather than unrestricted set interactions.

A.5 Proofs in Section 2.2.

Proof of Proposition 4 (a). Fix any instance of the k -Clique problem: let G be the (undirected, unweighted) graph with n vertices, and let $A \in \{0, 1\}^{n \times n}$ be its adjacency matrix (where we follow the convention that diagonal entries are set equal to 1). We will create an explicit instance of Problem (Main) in the following way (using the formulation in P):

$$\hat{A} = A + I, \text{ i.e. the } I \text{ is identity matrix, so that } \hat{A}_{ij} = \text{softmax}_i(A_{ij}) + \delta_{ij} = \text{softmax}_i(A_{ij}) + \delta_{ij}$$

$\hat{A} \in \mathbb{R}^{n \times n}$ is set according to

$$\hat{A}_{ij} = \begin{cases} 0 & \text{for } i = j \text{ or } j = i \\ \frac{1}{\#_{89}} & \text{for } 8-9 < =- \end{cases}$$

where ϵ_j is a constant large enough that $\epsilon_j \exp^{-\epsilon_j} > \frac{1}{2}$. In block notation, this is

$$A = \begin{array}{c|c} \begin{array}{c} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \\ 0 \dots \dots \dots 0 \end{array} & \begin{array}{c} 0 \\ \epsilon_1 \\ \vdots \\ \epsilon_n \\ 0 \\ 0 \end{array} \end{array}$$

Let $\epsilon = 1$, and we set $D = 1$, so that $A + D = A + I$.

Let $\epsilon = 1$ is set according to $\epsilon_j = 1$ for $j = 1, \dots, n$, and $\epsilon_n = \epsilon$, where we take any constant $\epsilon > 2$.

For $\epsilon = 1, \dots, n - 1$, we set ϵ to be:

$$\epsilon_j = \begin{cases} 0 & \text{if } 0 < G < \frac{\epsilon_j - 1}{2} \\ \frac{6G - 3\epsilon_j + 1}{\epsilon_j - 1} & \text{if } \frac{\epsilon_j - 1}{2} < G < \frac{\epsilon_j + 1}{3} \\ 1 & \text{if } G > \frac{\epsilon_j + 1}{3} \end{cases}$$

Because $\epsilon_j > 2$, we have $\frac{\epsilon_j - 1}{2} < \frac{\epsilon_j + 1}{3} < \epsilon_j - 1$, so ϵ_j is continuous piece-wise linear for every $\epsilon_j = 1, \dots, n - 1$. We set $\epsilon_n = 0$ for every G .

Note that the quantity

$$\frac{1F_8 + D^{0 > G}}{F_8^{> G}}$$

remains unchanged if the vector F_8 is multiplied by a non-zero constant. Thus, rescaling the rows of A does not change P . So for simplicity of exposition, we replace A with A^0 , defined to be

$$A^0 = \begin{cases} 1 & \text{for } \epsilon_j = \epsilon \text{ or } \epsilon_j = \epsilon \\ \exp^{-\epsilon_j} & \text{for } \epsilon_j - \epsilon < \epsilon \end{cases}$$

As a sanity check, A^0 is simply A with each row rescaled to sum to one.

Now notice that because $\epsilon_j > 1$ for every $\epsilon_j = 1, \dots, n - 1$, and $\epsilon_n = 0$ the optimal value of this instance of P is at most 1. It will suffice to show that G has a clique of size $n - 1$ if and only if the optimal value is 1. We prove both directions separately.

If G has a clique of size k , then the optimal value is $\frac{1}{k}$: Suppose G has a clique of size k , and let $S = \{1, \dots, k\}$ be the vertex set of one such clique. Consider the solution G , where $G_{ij} = 1$ if and only if $i, j \in S$. We will show that the objective value at G is $\frac{1}{k}$ (and thus the optimal value is $\frac{1}{k}$):

Because $\frac{1}{k} = \frac{1}{k} \exp^{1/k}$ whenever $k \geq 2$, we have

$$\begin{aligned} \tilde{O} &= \sum_{S \subseteq [n]} \frac{1F_S^0 + D^0 \cdot G}{F_S^0 \cdot G} = \sum_{S \subseteq [n]} \frac{1F_S^0 + D^0 \cdot G}{F_S^0 \cdot G} \\ &= \sum_{S \subseteq [n]} \frac{1, 1: \exp^{1/k}}{1, 1: \exp^{1/k}} \cdot \end{aligned}$$

Since k was chosen to be large enough that $1: \exp^{1/k} > 1 + \frac{1}{k}$, it follows that

$$\frac{1, 1: \exp^{1/k}}{1, 1: \exp^{1/k}} > \frac{1, \frac{1}{k}}{1, \frac{1}{k}} = \frac{2, 1}{3}.$$

Therefore,

$$\tilde{O} = \sum_{S \subseteq [n]} \frac{1F_S^0 + D^0 \cdot G}{F_S^0 \cdot G} = \sum_{S \subseteq [n]} \frac{1, 1: \exp^{1/k}}{1, 1: \exp^{1/k}} = \sum_{S \subseteq [n]} 1 = \frac{1}{k}.$$

If the optimal value is $\frac{1}{k}$, the G has a clique of size k : Suppose the optimal value is $\frac{1}{k}$. Let G be an optimal solution, and let $S = \{1, \dots, k\}$ be the index set such that $G_{ij} = 1$ for $i, j \in S$. Notice that if a solution G has $G_{ij} = 0$, then since $1F_S^0 + D^0 \cdot G = \frac{1}{k}$ for all $S \subseteq [n]$, we have

$$\sum_{S \subseteq [n]} \frac{1F_S^0 + D^0 \cdot G}{F_S^0 \cdot G} = \sum_{S \subseteq [n]} 1 = \frac{1}{k}.$$

Therefore, it must be the case that $G_{ij} = 1$. Also, because

$$\frac{1}{k} = \sum_{S \subseteq [n]} \frac{1F_S^0 + D^0 \cdot G}{F_S^0 \cdot G} = \sum_{S \subseteq [n]} \frac{1}{k} = \frac{1}{k}.$$

we must have $\sum_{S \subseteq [n]} 1 = \frac{1}{k}$ and

$$\sum_{S \subseteq [n]} \frac{1F_S^0 + D^0 \cdot G}{F_S^0 \cdot G} = 1$$

for every $S \subseteq [n]$.

Let d_j be the degree of vertex j in the induced subgraph $G[S]$ with vertex set S .

Then

$$\begin{aligned} \sum_{S \ni j} \frac{d_j}{|S|} &= \sum_{S \ni j} \frac{d_j^0 + d_j^1}{|S|} \\ &= \sum_{S \ni j} \frac{d_j^0}{|S|} + \sum_{S \ni j} \frac{d_j^1}{|S|} \\ &= \sum_{S \ni j} \frac{d_j^0}{|S|} + \sum_{S \ni j} \frac{d_j^1}{|S|} \end{aligned}$$

where the last inequality follows since $d_j \leq 2$. Because $\sum_j d_j = 2|E|$, we have $\sum_j d_j^0 \leq 2|E|$. Suppose for the sake of contradiction that $\sum_j d_j^0 \geq 3$ for some j .

Then

$$\sum_{S \ni j} \frac{d_j^0}{|S|} \geq \sum_{S \ni j} \frac{1}{2} = 0$$

which contradicts that

$$\sum_{S \ni j} \frac{d_j^0 + d_j^1}{|S|} = 1$$

for every j . Therefore we must have $d_j^0 \leq 2$ for every j . Hence S corresponds to the vertex set of a clique of size ≤ 2 .

Proof of Proposition 4 (b). Fix any constant $\epsilon > 0$. We construct an instance of Problem(Main) (using the formulation in P) by applying the Johnson-Lindenstrauss Lemma:

Lemma 5 (Johnson-Lindenstrauss Lemma). For any $0 < \epsilon < 1$ and any set of n points in \mathbb{R}^d , there exists a universal constant $\epsilon_0 > 0$ and a linear function $f : \mathbb{R}^d \rightarrow \mathbb{R}^k$ with $k = 2n \log \frac{1}{\epsilon_0}$ such that

$$(1 - \epsilon) \|x - y\|_2^2 \leq \|f(x) - f(y)\|_2^2 \leq (1 + \epsilon) \|x - y\|_2^2$$

for all $x, y \in \mathbb{R}^d$ (and

$$(1 - \epsilon) \|x - y\|_2^2 \leq \|f(x) - f(y)\|_2^2 \leq (1 + \epsilon) \|x - y\|_2^2$$

for all $x, y \in \mathbb{R}^d$).

Take n to be large enough such that $\epsilon_0 \geq \frac{1}{2}$. Take $\epsilon < \epsilon_0$ to be small enough such that $\frac{1}{1 + \epsilon} \geq \frac{1}{2}$. Then n and k can be chosen to be both only depend on ϵ . We obtain the following corollary:

Corollary 2 (Corollary of Lemma 5.). There exists a number $\epsilon > 0$ and a set of $\exp(2^{1/\epsilon})$ unit vectors in \mathbb{R}^3 such that $\langle D_i, D_j \rangle \leq \epsilon$ for every D_i, D_j ($i \neq j$) and $D_i < D_j$.

Proof of Corollary 2. Let $\epsilon > 0$ be the universal constant in Lemma 5 and let $\delta = \exp(2^{1/\epsilon})$. Let $2^{1/\epsilon} = 2^{1/\epsilon} \cdot 4^\delta$, then $\delta = \exp(2^{1/\epsilon})$. Because $\epsilon > 0$ only depends on δ , we have $2^{1/\epsilon} > 0$ also only depends on δ . Consider \mathbb{R}^3 where e_1 is the unit vector where the first entry equals to 1 and all other entries equal to 0. Then we have $\|e_1\|_2 = 1$ for every $\delta < 4$. By Lemma 5, there exists a linear function $f: \mathbb{R}^3 \rightarrow \mathbb{R}$ such that

$$\left| \langle f, e_1 \rangle - \frac{\epsilon}{4} \right| \leq \frac{\epsilon}{4}$$

for all δ and

$$\left| \langle f, e_2 \rangle - \frac{\epsilon}{4} \right| \leq \frac{\epsilon}{4}$$

for all $\delta < 4$.

For every $\delta < 4$, because

$$\|e_1 - e_2\|_2 = \sqrt{2}, \quad \|e_1 + e_2\|_2 = \sqrt{2}$$

we have

$$\begin{aligned} \langle f, e_1 - e_2 \rangle &= \langle f, e_1 \rangle - \langle f, e_2 \rangle \\ &= \frac{\epsilon}{4} - \frac{\epsilon}{4} \\ &= 0 \end{aligned}$$

and

$$\begin{aligned} \langle f, e_1 + e_2 \rangle &= \langle f, e_1 \rangle + \langle f, e_2 \rangle \\ &= \frac{\epsilon}{4} + \frac{\epsilon}{4} \\ &= \frac{\epsilon}{2} \end{aligned}$$

Let $D_i = \frac{1}{\sqrt{2}}(e_1 - e_2)$ for all $i \in \{1, \dots, \delta\}$. Then each D_i is a unit vector. Moreover, for every $D_i < D_j$,

$$\langle D_i, D_j \rangle = \frac{\langle f, D_i \rangle - \langle f, D_j \rangle}{\sqrt{2}} = \frac{\frac{\epsilon}{4} - \frac{\epsilon}{4}}{\sqrt{2}} = 0$$

^ 3E = 1, and we set D = 1, so that +D = +.

^ + 2 R= 1 is set according to g= 1 for 8 = 1-•••- = 1, and ≠ 2.

^ For 8 = 1-•••- = 1, we set 5 to be:

$$5^1 G^0 = \begin{cases} 0 & \text{if } 0 \leq G \leq \frac{2, \exp^1 \#^0, \frac{1}{4}}{1, \exp^1 \#^0, \frac{1}{4}} \\ \frac{G \frac{2, \exp^1 \#^0, \frac{1}{4}}{1, \exp^1 \#^0, \frac{1}{4}}}{\frac{2, \exp^1 \#^0, \frac{1}{4}}{1, \exp^1 \#^0, \frac{1}{4}} - \frac{2, \exp^1 \#^0, \frac{1}{4}}{1, \exp^1 \#^0, \frac{1}{4}}} & \text{if } \frac{2, \exp^1 \#^0, \frac{1}{4}}{1, \exp^1 \#^0, \frac{1}{4}} \leq G \leq \frac{2, \exp^1 \#^0, \frac{1}{4}}{1, \exp^1 \#^0, \frac{1}{4}} \\ 1 & \text{if } G \geq \frac{2, \exp^1 \#^0, \frac{1}{4}}{1, \exp^1 \#^0, \frac{1}{4}} \end{cases}$$

Because $\frac{2, \exp^1 \#^0, \frac{1}{4}}{1, \exp^1 \#^0, \frac{1}{4}} \leq \frac{2, \exp^1 \#^0, \frac{1}{4}}{1, \exp^1 \#^0, \frac{1}{4}}$, we have $5^1 G^0$ is continuous piece-wise linear for every $8 = 1-•••- = 1$. We set $G = 0$ for every G .

Now notice that because $5^1 G^0 = 1$ for every $8 = 1-•••- = 1$, and $5^1 G^0 = 0$ the optimal value of this instance of Problem P is at most 1. It will suffice to show that the largest clique has size 0 if and only if the optimal value is 0. We prove both directions separately.

If \mathcal{G} has a clique of size 0, then the optimal value is at least 0. Suppose \mathcal{G} has a clique of size 0. Without loss of generality, let $t_1, \dots, t_{|C|}$ correspond the vertex set of a clique of size 0 in \mathcal{G} . Consider the solution G , where $G_{ij} = 1$ for $8 \in \{1, \dots, |C|\} [f=g]$. Because $\mathcal{G} \in \mathcal{G}$, the solution G is feasible. We will show that the objective value of P at G is 0 (and thus the optimal value is at least 0).

Because $G_{ij} = \exp^1 \#^0$ whenever $8 \in \{1, \dots, |C|\}$, we have

$$\begin{aligned} \tilde{G}_{8=1} &= \frac{1 F_8^0 + D^0 > G}{F_8^0 G} = \frac{1 F_8^0 + D^0 > G}{F_8^0 G} \\ &= \frac{2, :^0 \exp^1 \#^0}{1, :^0 \exp^1 \#^0} \end{aligned}$$

Since $:^0 \leq 1$. It follows that

$$\frac{2, :^0 \exp^1 \#^0}{1, :^0 \exp^1 \#^0} \leq \frac{2, : \exp^1 \#^0}{1, : \exp^1 \#^0}$$

Therefore

$$\tilde{G}_{8=1} = \frac{1 F_8^0 + D^0 > G}{F_8^0 G} = \frac{2, :^0 \exp^1 \#^0}{1, :^0 \exp^1 \#^0} = j \quad |j| = :^0$$

If the largest clique in G has size $\omega(G)$, then the optimal value is at most $\omega(G)$. By our assumption we have $\omega(G) \geq \frac{1}{4}$. Suppose for the sake of contradiction that the optimal value is $\frac{1}{4}$. Let G be an optimal solution to P . Notice $\omega(G) = 1$, then since $\frac{1}{4} \omega(G) + D^0 \geq \frac{1}{4}$ for all $9 \leq 2 \leq n$, we have

$$\sum_{8=1}^n \frac{1}{4} \omega(G) + D^0 \geq \frac{1}{4} \sum_{8=1}^n \omega(G) = 0.$$

Therefore we must have $\omega(G) \geq 1$.

Let $S \subseteq [n]$ be the index set where $\omega(G) = 1$ for $8 \in S$. Because $\omega(G) \geq 1$ for every $8 = 1, \dots, n$, we have $|S| \geq \frac{1}{4}n$. Let 3^{18^0} be the degree of vertex 8 in the induced subgraph of G with vertex set S . Because $G[S]$ consists of disjoint cliques with size at most $\omega(G)$, we have $3^{18^0} \leq \omega(G)$ for every $8 \in S$.

Fix any $8 \in S$ such that

$$\frac{1}{4} \omega(G) + D^0 \geq \frac{1}{4} \omega(G)$$

or, equivalently,

$$\frac{1}{4} \omega(G) + D^0 \geq \frac{1}{4} \omega(G) \iff \frac{1}{4} \omega(G) \geq \frac{1}{4} \omega(G)$$

Because $\exp^1 \#j D_j^0 \leq \exp^1 X \#^0$ for every $8 - 9 \leq 2 \leq n$ and $\omega(G) \geq 1$, we have

$$\begin{aligned} \frac{1}{4} \omega(G) + D^0 \geq \frac{1}{4} \omega(G) &= \frac{2 \sum_{9 \leq 2 \leq n} \exp^1 \#j D_j^0 \omega(G) + 3^{18^0} \omega(G)}{1 \sum_{9 \leq 2 \leq n} \exp^1 \#j D_j^0 \omega(G) + 3^{18^0} \omega(G)} \\ &\geq \frac{2 \sum_{9 \leq 2 \leq n} \exp^1 X \#^0 \omega(G) + 3^{18^0} \omega(G)}{1 \sum_{9 \leq 2 \leq n} \exp^1 X \#^0 \omega(G) + 3^{18^0} \omega(G)} \\ &= \frac{2 \sum_{9 \leq 2 \leq n} \exp^1 X \#^0 \omega(G) + 3^{18^0} \omega(G)}{1 \sum_{9 \leq 2 \leq n} \exp^1 X \#^0 \omega(G) + 3^{18^0} \omega(G)} \end{aligned}$$

where the second inequality follows since $\omega(G) \geq 1$, and the last inequality follows since $\exp^1 X \#^0 \geq \exp^1 \#j D_j^0$ and $3^{18^0} \geq 1$. Therefore

$$\frac{2 \sum_{9 \leq 2 \leq n} \exp^1 X \#^0 \omega(G) + 3^{18^0} \omega(G)}{1 \sum_{9 \leq 2 \leq n} \exp^1 X \#^0 \omega(G) + 3^{18^0} \omega(G)} \geq \frac{2 \sum_{9 \leq 2 \leq n} \exp^1 \#j D_j^0 \omega(G) + 3^{18^0} \omega(G)}{1 \sum_{9 \leq 2 \leq n} \exp^1 \#j D_j^0 \omega(G) + 3^{18^0} \omega(G)}$$

So we have

$$\exp^1 X \#^0 \geq \exp^1 \#j D_j^0 \iff \exp^1 X \#^0 \geq \exp^1 \#j D_j^0$$

Rearrange the above inequality gives

$$\exp^1 X \#^{01}: : \infty \ddot{Y} \exp^1 \#^{01}: : \infty, \frac{1}{4} \bullet$$

Because : : 0 , 1 and : 0 : " , we get

$$\exp^1 X \#^0 \ddot{Y} \exp^1 \#^{0"} , \frac{1}{4} \bullet$$

However, since $\exp^1 \#^0 = 1 \bullet 2^"$ and $\exp^1 X \#^0 = 1 = 1 \bullet 2^"$, we have

$$\exp^1 X \#^0 = 1 = \frac{1}{2^"} = \frac{5}{6} \geq \exp^1 \#^{0"} , \frac{1}{4} \bullet$$

A contradiction. Therefore : : 0 .

Proof of Proposition 5. Our proof is based on a reduction from Problem (Main) to the well-known Multi-dimensional Knapsack Problem (MDKP), defined as follows:

Definition 5 (Multi-dimensional Knapsack Problem). The Multi-dimensional Knapsack Problem (MDKP) with 2 items and 3 dimensions is defined as:

$$\begin{aligned} \text{(MDKP)} \quad & \max \quad \sum_{MDKP}^1 G^0 = ? G \\ & \text{s.t.} \quad \sum_{i=1}^2 G_i \leq 882 \gg 3\frac{1}{4} - \\ & \quad G_2 \leq f_0 - 1^2 g \end{aligned}$$

Here, $? \in \mathbb{N}^2$, and for all $82 \gg 3\frac{1}{4}$, we have $2 \in \mathbb{Z}_0^2$ and $6 \in \mathbb{Z}_0^2$.

We present the reduction in the following proposition:

Proposition 16. Consider Problem (Main), written in the equivalent form P as given in Observation 2, reproduced below, which has parameters α and A , where A is the non-negative rank of α :

$$\begin{aligned} \text{(P)} \quad & \max \quad \sum_{i=1}^8 \alpha_i G^0 = \sum_{i=1}^8 \frac{1}{F_i} \frac{F_i + D^0 > G}{F_i > G} \\ & \text{s.t.} \quad G_2 \leq f_0 - f_g = 1 \quad 4^> G : \bullet \end{aligned}$$

Now consider an instance of MDKP with parameters 2 and 3. Suppose there exists an algorithm ALG for solving P with parameters $\alpha = : = 2, 3, 1$ and

²We assume $\alpha_{\min} \geq 0$ without loss of generality. If $\alpha_9 = 0$ for some $9 \in \{2, 3, 1\}$, there always exists an optimal solution with $\alpha_9 = 0$. Hence, we can safely ignore the 9-th entry of α , G , and each H

minf2-3g A 2, 3, 1, such that for any sufficiently small $\epsilon > 0$, the algorithm satisfies

$$ALG_P \leq (1 + \epsilon) OPT_P$$

with runtime $O(n^3)$. Then we can construct an algorithm ALG^0 for solving MDKP with parameters 2 and 3, such that for the same n , it satisfies

$$ALG^0_{MDKP} \leq (1 + \epsilon) OPT_{MDKP}$$

with runtime $O(n^3)$.

Proof of Proposition 16. Fix an MDKP instance. We create an instance of P as follows:

$\hat{m} = 3 + 2 + 1 = 6$. Here, out of the \hat{m} total variables, the first 3 variables will correspond to the 3 constraints of MDKP, the next 2 variables after that will correspond to the 2 variables of MDKP, and the last variable will be a dummy variable that must be selected in any optimal solution of P.

For each $8 = 1, \dots, 3$, set

$$F_{89} = \begin{cases} 0 & \text{for } 9 = 1, \dots, 3 \\ H_{8-9,3} & \text{for } 9 = 3, 1, \dots, 3, 2 \\ 1 & \text{for } 9 = 3, 2, 1 \end{cases}$$

That is, in block notation,

$$F_8^> = \begin{matrix} h & i \\ 0 & 0 \end{matrix} \Big| \begin{matrix} H_8 \\ 1 \end{matrix}$$

Then we have

$$F_8^> G = \begin{matrix} \tilde{Q} \\ H_{8,9} \\ G_{3,2,1} \end{matrix}$$

Moreover, since $H_8 \in Z_0^2$ for each $8 \in \{1, 2, 3\}$, the non-negative rank of the sub-matrix of G , consisting of its first 3 rows is at most $\min\{2, 3\}$.

For each $8 = 3, 1, \dots, 3, 2, 1$, set $F_{8,9} \in R_0^{3,2,1}$ to be any non-zero vectors such that G has the desired non-negative rank. This is possible since $\min\{2, 3\} = 2$.

$3_E = 1$, and we set $D = 1$, so that $+D = +$.

$\hat{+} + 2 R^{13,2,1^0 1}$ is set to be $+_8 = 0$ for $8 = 1, \dots, 7$, and $+_8 = 1$ for $8 = 3, 1, \dots, 3, 2$, and $+_{3,2,1} = 2$.

$\hat{+}$ For each $8 = 1, \dots, 3$, set

$$\hat{+}^1 G^0 = \begin{cases} \frac{G_{9,1}^2}{G_{9,1}^1} & \text{for } G_{9,1}^2 \geq G_{9,1}^1 \\ \frac{G_{9,1}^1}{G_{9,1}^2} & \text{for } G_{9,1}^1 \geq G_{9,1}^2 \\ 0 & \text{for } G_{9,1}^2 = G_{9,1}^1 \end{cases}$$

Then $\hat{+}^1 G^0$ is continuous piecewise-linear.

$\hat{+}$ For each $8 = 3, 1, \dots, 3, 2$, set $\hat{+}^1 G^0 = \frac{G_{3,2,1}^3}{G_{3,2,1}^2}$. Set $\hat{+}^1 G^0 = 0$.

Because $\hat{+}^1 G^0 \geq 0$ and $\hat{+}^1 G^0 \leq 1$ is a feasible solution to MDKP, we have $\text{OPT}_{\text{MDKP}} = 0$ if and only if $\hat{+}^1 G^0$ is the only feasible solution to MDKP. Moreover, if $\hat{+}^1 G^0$ is not the only feasible solution to MDKP, then $\text{OPT}_{\text{MDKP}} = \min$.

Our proof relies on the following lemma.

Lemma 6. Let $G \geq 0$ be a feasible solution to P such that $\hat{+}^1 G^0 \geq 0$. Then we can construct a feasible solution $I \geq 0$ to MDKP such that $\hat{+}^1 I^0 = \hat{+}^1 G^0$.

Conversely, let $I \geq 0$ be a feasible solution to MDKP such that $\hat{+}^1 I^0 \geq 0$. Then we can construct a feasible solution $G \geq 0$ to P such that $\hat{+}^1 G^0 = \hat{+}^1 I^0$.

Lemma 6 gives a correspondence between solutions to P and solutions to MDKP. In particular, as a corollary, Lemma 6 implies the relationship between the optimal objective values of P and MDKP.

Corollary 3 (Corollary of Lemma 6.). $\text{OPT}_P = 0$ if and only if $\text{OPT}_{\text{MDKP}} = 0$. Moreover, suppose $\text{OPT}_{\text{MDKP}} > 0$. Then $\text{OPT}_P = \text{OPT}_{\text{MDKP}}$.

Proof of Corollary 3. Suppose $\text{OPT}_P > 0$. Let G be an optimal solution to P. By Lemma 6 there exists a feasible solution I to MDKP such that $\hat{+}^1 I^0 \geq 0$ and

$$\text{OPT}_P = \hat{+}^1 G^0 = \hat{+}^1 I^0 = \text{OPT}_{\text{MDKP}}$$

Conversely, suppose $\text{OPT}_{\text{MDKP}} > 0$. Let I be an optimal solution to MDKP. By Lemma 6 there exists a feasible solution G to MDKP such that

$$\text{OPT}_{\text{MDKP}} = \hat{+}^1 I^0 = \hat{+}^1 G^0 = \text{OPT}_P$$

We also get

$$\epsilon^1 G^0 = \epsilon_{MDKP}^1 | \epsilon^0 | 0$$

”

Before proving Lemma 6, we first show that it is sufficient to prove Lemma 6. Assume we have an algorithm ALG for solving P such that, for any sufficiently small $\epsilon > 0$, we have ALG satisfies

$$ALG_P \leq (1 + \epsilon) \cdot OPT_P$$

Then our proposed algorithm ALG^0 for solving $MDKP$ works as follows:

- ^ If $ALG_P = 0$ – ALG^0 outputs $\epsilon^2 f_0 - 1g$.
- ^ If $ALG_P > 0$ – let G be the solution to P given by ALG_P . Then ALG^0 outputs the solution $\epsilon^2 f_0 - 1g$ to $MDKP$ given by Lemma 6. That is, $\epsilon^2 f_0 - 1g$ that satisfies

$$\epsilon^1 G^0 = \epsilon_{MDKP}^1 | \epsilon^0 |$$

Performance Guarantee of ALG^0 : Suppose $OPT_{MDKP} = 0$. Then by Corollary 3,

$$ALG_P \leq OPT_P = 0$$

Therefore ALG^0 correctly outputs $\epsilon^2 f_0 - 1g$.

On the other hand, suppose $OPT_{MDKP} > 0$. Then by Corollary 3,

$$ALG_{MDKP}^0 = \epsilon_{MDKP}^1 | \epsilon^0 | = \epsilon^1 G^0 \leq (1 + \epsilon) \cdot OPT_P = (1 + \epsilon) \cdot OPT_{MDKP}$$

This proves the desired performance guarantee of ALG^0 . To finish the proof, we prove Lemma 6 and analyze the runtime of ALG

Proof of Lemma 6. Recall for each $\delta = 1 - \dots - 3$, we set

$$F_{\delta, \epsilon} = \begin{cases} 0 & \text{for } \delta = 1 - \dots - 3 \\ H_{\delta - 3} & \text{for } \delta = 3, 1 - \dots - 3, 2 \\ 1 & \text{for } \delta = 3, 2, 1 \end{cases}$$

Also, $\epsilon^1 D_{\delta} = 0$ for $\delta = 1 - \dots - 3$, and $\epsilon^1 D_{\delta} = 1$ for $\delta = 3, 1 - \dots - 3, 2$, and $\epsilon^1 D_{3,2,1} = 2$. Therefore for $\delta = 1 - \dots - 3$ we have

$$\frac{\epsilon^1 F_{\delta} + D_{\delta} \epsilon^0 G}{F_{\delta} \epsilon^0 G} = \frac{\sum_{\delta=1}^2 H_{\delta, \epsilon} G_{\delta, \epsilon} + 2G_{3,2,1}}{\sum_{\delta=1}^2 H_{\delta, \epsilon} G_{\delta, \epsilon} + G_{3,2,1}}$$

Recall for each $8 = 3, 1, \dots, 3$, we set $5^1 G^0 = 3_3$, and we set $5_{3,2,1}^1 G^0 = 0$. Therefore we have

$$\begin{aligned} 5^1 G^0 &= \sum_{8=1}^2 \tilde{G}_{8,1} \sum_{9=1}^2 \frac{1 F_{8,9} + D^0 > G}{F_{8,9}^G} \\ &= \sum_{8=1}^2 \tilde{G}_{8,1} \left(\sum_{9=1}^2 \frac{H_{8,9} G_{8,9} + 2G_{3,2,1}}{H_{8,9} G_{8,9} + G_{3,2,1}} \right) \tilde{G}_{8,1} \end{aligned}$$

First, let $G \geq 0$ be a feasible solution to P^0 such that $5^1 G^0 \leq 0$. Let $I \geq 0$ where $I_9 = G_{3,9}$ for $9 \geq 2 \gg 3/4$. We claim that I is a feasible solution to MDKP and

$$5^1 G^0 = 5_{MDKP}^1 I^0.$$

Because $\sum_{8=1}^2 \tilde{G}_{8,1} \sum_{9=1}^2 \frac{H_{8,9} G_{8,9} + 2G_{3,2,1}}{H_{8,9} G_{8,9} + G_{3,2,1}} \leq 0$, we must have $G_{3,2,1} = 0$ for every $8 \geq 2 \gg 3/4$. Moreover, if $G_{3,2,1} = 0$, then

$$\sum_{8=1}^2 \left(\sum_{9=1}^2 \frac{H_{8,9} G_{8,9} + 2G_{3,2,1}}{H_{8,9} G_{8,9} + G_{3,2,1}} \right) = 5^1 I^0 \leq 0$$

for every $8 \geq 2 \gg 3/4$. Hence we must have $G_{3,2,1} = 0$. Because $I \geq 0$ and $G \geq 0$ for all $8 \geq 2 \gg 3/4$, if $\sum_{9=1}^2 H_{8,9} G_{8,9} \leq C_8$, we must have $\sum_{9=1}^2 H_{8,9} I_9 \leq C_8$. Then

$$\sum_{8=1}^2 \left(\sum_{9=1}^2 \frac{H_{8,9} G_{8,9} + 2G_{3,2,1}}{H_{8,9} G_{8,9} + G_{3,2,1}} \right) \sum_{9=1}^2 \frac{G_{3,9}}{G_{3,9}} = \sum_{9=1}^2 \frac{G_{3,9}}{G_{3,9}}$$

Therefore we must have $\sum_{9=1}^2 H_{8,9} I_9 \leq C_8$ for all $8 \geq 2 \gg 3/4$. Hence $\sum_{9=1}^2 H_{8,9} I_9 \leq C_8$ for all $8 \geq 2 \gg 3/4$, which shows I is a feasible solution to MDKP. Finally, because

$$\sum_{9=1}^2 \frac{H_{8,9} G_{8,9} + 2G_{3,2,1}}{H_{8,9} G_{8,9} + G_{3,2,1}} \leq \frac{G_{3,2,1}}{G_{3,1}}$$

for all $8 \geq 2 \gg 3/4$, we have

$$\sum_{8=1}^2 \left(\sum_{9=1}^2 \frac{H_{8,9} G_{8,9} + 2G_{3,2,1}}{H_{8,9} G_{8,9} + G_{3,2,1}} \right) = 0$$

for all $8 \geq 2 \gg 3/4$. Therefore

$$\begin{aligned} 5^1 G^0 &= \sum_{8=1}^2 \tilde{G}_{8,1} \left(\sum_{9=1}^2 \frac{H_{8,9} G_{8,9} + 2G_{3,2,1}}{H_{8,9} G_{8,9} + G_{3,2,1}} \right) \tilde{G}_{8,1} \\ &= \sum_{8=1}^2 \tilde{G}_{8,1} \tilde{G}_{8,1} \\ &= 5_{MDKP}^1 I^0. \end{aligned}$$

Conversely, let $(f_0 - 1)g$ be a feasible solution to MDKP such that $\sum_{i=1}^n f_i = 1$. Let G be a solution to P where $G_{i,1} = 1$ for $i = 1, \dots, n$ and $G_{i,2} = 1$ for $i = 1, \dots, n$. We claim that G is a feasible solution to MDKP and

$$\sum_{i=1}^n G_{i,1} = \sum_{i=1}^n f_i = 1.$$

Because $(f_0 - 1)g$ is a feasible solution to MDKP, we have $\sum_{i=1}^n f_i g_{i,2} \leq C_2$ for all $2 \leq i \leq n$. Therefore $\sum_{i=1}^n f_i G_{i,2} \leq C_2$ for all $2 \leq i \leq n$. Then

$$\sum_{i=1}^n f_i G_{i,2} - \sum_{i=1}^n f_i g_{i,2} \leq C_2 - C_2 = 0$$

for all $2 \leq i \leq n$. Therefore we have

$$\sum_{i=1}^n f_i (G_{i,2} - g_{i,2}) \leq 0$$

for all $2 \leq i \leq n$. Hence

$$\begin{aligned} \sum_{i=1}^n f_i G_{i,2} &= \sum_{i=1}^n f_i g_{i,2} + \sum_{i=1}^n f_i (G_{i,2} - g_{i,2}) \\ &\leq C_2 + 0 \\ &= C_2 \\ &= \sum_{i=1}^n f_i g_{i,2} \\ &= \sum_{i=1}^n f_i = 1. \end{aligned}$$

”

Given an MDKP instance, our construction of the corresponding P instance takes runtime. Also, the procedure of the construction of a feasible solution to MDKP given a feasible solution to P described in Lemma 6 takes runtime: we simply set $(f_0 - 1)g$ where $g_{i,1} = G_{i,1}$ for $1 \leq i \leq n$. Therefore, if ALG has runtime $\mathcal{O}(n^k)$, our ALG⁰ has runtime $\mathcal{O}(n^k)$.

”

By Proposition 16, any $(1 - \epsilon)$ -approximation scheme for P can be directly translated into a $(1 - \epsilon)$ -approximation scheme for MDKP with comparable runtime. Therefore, many hardness results for MDKP carry over directly to P . We cite several such results below:

Proposition 17 (Theorem 6 in [103]). If $\beta = 2$, MDKP admits no $(1 - \epsilon)n^0$ -approximation scheme with runtime $5^{11 \cdot n^0} 2^{11^0}$ for any function ϵ , assuming the Clique problem is not Fixed-Parameter Tractable.

Proposition 18 (Theorem 5.1 in [87]). If $\beta = 2$, MDKP admits no $(1 - \epsilon)n^0$ -approximation scheme with runtime $5^{11 \cdot n^0} 2^{11^0}$ for any function ϵ , with runtime $5^{11 \cdot n^0} 2^{11 \cdot n^0}$ for any function ϵ , assuming Exponential Time Hypothesis holds.

Proposition 19 (Corollary of [60]). For general β , MDKP admits no $(1 - \epsilon)n^0$ -approximation scheme with runtime

$$: > \frac{3}{n \log^{2+3 \cdot n^0}} \quad \text{or} \quad : > 1^{P \cdot 3^0}$$

assuming Gap Exponential Time Hypothesis holds.

These results, along with Proposition 16, gives our desired lower bound statement in Proposition 5.

”

A.6 Proofs in Section 2.3.

Proof of Proposition 6. Let $\delta \in (0, 1)$. Let $c : \{1, \dots, 2\} \rightarrow \{0, 1\}$ be any bijection. Let $\gamma \in (0, 1)$ where

$$\delta - c^1 < -\epsilon = \begin{cases} \max_{B=1} \frac{\exp^1 @ : c^0 \cdot \prod_{B=1} \exp^1 @ : B^0}{\dots} & \text{for } \delta = < \\ \dots & \text{otherwise,} \end{cases}$$

and

$$\delta - c^1 < -\epsilon = \begin{cases} \max_{B=1} \frac{\exp^1 @ : \gamma^0 \cdot \prod_{B=1} \exp^1 @ : B^0}{\dots} & \text{for } \delta = < \\ \dots & \text{otherwise.} \end{cases}$$

Let $\gamma = \delta$. Then

$$\begin{aligned} \delta - c^1 < -\epsilon &= \max_{B=1} \frac{\exp^1 @ : \gamma^0 \cdot \prod_{B=1} \exp^1 @ : B^0}{\dots} \\ &= \max_{c^1 < -\epsilon} \frac{\delta - c^1 < -\epsilon \cdot \delta - c^1 < -\epsilon}{\dots} \\ &= \delta - c^1 \delta - \gamma \delta - c^1 \delta - \gamma^0 \\ &= \exp^1 @ : \gamma^0 \cdot \prod_{B=1} \exp^1 @ : B^0 \end{aligned}$$

Finally, because for every $\delta > 0$ we have $\mathbb{Q} \subseteq \mathbb{Q}_{\delta}^{k_2}$ and $k_{\delta} : \mathbb{Q}_{\delta}^{k_2} \rightarrow X$, we have exactly the same setup as in the proof of Proposition 7. Therefore the exact same proof gives $\frac{1}{\delta} W$, $\frac{1}{\delta} \mathbb{Q}$, $\frac{1}{\delta} \mathbb{Q}_{\delta}^{k_2}$ for all $\delta > 0$, where $W = \frac{1}{\delta} \max_{k \in \mathbb{Q}_{\delta}^{k_2}} \|k - k_{k_2-1}\|_g$.

”

Proof of Corollary 1. In the proof of Proposition 7, we constructed a partition $I = \{I_1, \dots, I_g\}$ with $\frac{1}{\delta} \max_{k \in \mathbb{Q}_{\delta}^{k_2}} \|k - k_{k_2-1}\|_g \leq \epsilon$ that satisfies for every $\delta > 0$ we have $\mathbb{Q} \subseteq \mathbb{Q}_{\delta}^{k_2}$ and $k_{\delta} : \mathbb{Q}_{\delta}^{k_2} \rightarrow X$. The result then follows from Proposition 6.

”

A.7 Pseudo-code of Main Algorithm

Below we give the pseudo-code of our main algorithm.

Algorithm 9: Main Algorithm

Input: Number of items n , maximum number of recommended items k , key matrix $\mathbb{K} \in \mathbb{R}^{n \times k}$, query matrix $\mathbb{Q} \in \mathbb{R}^{k \times k}$, value matrix $\mathbb{V} \in \mathbb{R}^{n \times k}$, reward functions $f_{\delta=1}^g$, user vector $\mathbb{D} \in \mathbb{R}^k$, parameter $n \geq 0$, n -Approximate δ -Nearest Neighbor oracle (ANN), attention matrix $\mathbb{A} = \text{softmax}(\mathbb{K} \mathbb{Q}^{-1}) \in \mathbb{R}^{n \times k}$, low-rank approximation $\mathbb{W} \in \mathbb{R}^{n \times k}$ with factorization $\mathbb{W} = \mathbb{U} \mathbb{V}^T$ and element wise guarantee $\frac{1}{\delta} W$, $\frac{1}{\delta} \mathbb{Q}$, $\frac{1}{\delta} \mathbb{Q}_{\delta}^{k_2}$

Output: Solution G to Problem (Main).

// Phase One: Preprocess (Algorithm 10)

$\mathbb{X} \leftarrow I \leftarrow S \leftarrow$ preprocessed ANN⁰ Run Algorithm 10 with inputs

$\frac{1}{\delta} \mathbb{Q} \leftarrow \mathbb{Q} \leftarrow \mathbb{Q}_{\delta}^{k_2} \leftarrow n \leftarrow$ ANN⁰

// Phase One: Query (Algorithm 11)

Run Algorithm 11 with inputs $\mathbb{D} \leftarrow \mathbb{X} \leftarrow I \leftarrow S \leftarrow$ preprocessed ANN⁰

// Phase Two (Algorithm 16)

$G \leftarrow$ Run Algorithm 16 with inputs $\frac{1}{\delta} \mathbb{Q} \leftarrow \mathbb{Q} \leftarrow \mathbb{Q}_{\delta}^{k_2} \leftarrow W \leftarrow \mathbb{D} \leftarrow n$

return G

A.8 Proofs in Section 2.4.

Proof of Observation 2. Fix (\mathbb{K}, \mathbb{Q}) and let $G \in \mathbb{R}^{n \times k}$ where $G_{\delta} = \mathbb{K} \mathbb{Q}^{-1}$. We prove that the objective values of Problem (Main) and P are the same. By our construction of \mathbb{Q} , we get $\mathbb{Q} \in \mathbb{R}^{k \times k}$ where $\mathbb{Q}_{\delta} = \mathbb{Q}_{\delta}$, and similarly for $\mathbb{K} \in \mathbb{R}^{n \times k}$ and $\mathbb{V} \in \mathbb{R}^{n \times k}$. Then for $\delta > 0$ (we get

$\text{softmax}^{11-\&^{01}-0} >^0_{8-B} = \exp^1_{\&: B^0} \cdot \prod_{B^0} \exp^1_{\&: B^0}$. Therefore, for $\delta \geq 2$ (we have

$$\begin{aligned} \frac{1F_8 + D^0 > G}{F_8^> G} &= \frac{\prod_{B^0} (F_{8-B} + D_B)}{\prod_{B^0} F_{8-B}} \\ &= \frac{\prod_{B^0} (\exp^1_{\&: B^0} \cdot \prod_{B^0} \exp^1_{\&: B^0})}{\prod_{B^0} (\exp^1_{\&: B^0} \cdot \prod_{B^0} \exp^1_{\&: B^0})} \\ &= \frac{\prod_{B^0} (\exp^1_{\&: B^0})}{\prod_{B^0} (\exp^1_{\&: B^0})} \\ &= \frac{\prod_{B^0} (\text{softmax}^{11-\&^{01}-0} >^0_{8-B})}{\prod_{B^0} (\text{softmax}^{11-\&^{01}-0} >^0_{8-B})} \\ &= \text{softmax}^{11-\&^{01}-0} >^0_{8-B} \\ &= \text{SA}_{\&-} \rightarrow^1- (\&^0_{8} D-) \end{aligned}$$

where the fifth equality follows since $\prod_{B^0} \text{softmax}^{11-\&^{01}-0} >^0_{8-B} = 1$ for every matrix $\&$. Because this holds for every $\&$, the objective values of Problem (Main) and P are the same. ”

Proof of Proposition 7. Fix any $n \geq 0$. Let $X \geq 0$ be a parameter such that

$$(A.1) \quad X \frac{1}{140 \max\{k_{\&-1} - k_{\&-2}, 1\}} \quad \text{and} \quad n \leq 34X \max\{k_{\&-1} - k_{\&-2}, 1\} + D_{\max}^0 X$$

First we partition the rows $\&$ and $\&$. Because we can cover a ball with radius $k_{\&-1}$ in $\mathbb{R}^{3,\&}$ using $d_4 k_{\&-1} \cdot X e^{3,\&}$ number of balls with radius $X/2$ (see e.g. [61, 156, 168]), we can create a partition of the index set $\&$ such that the size of the partition is $d_4 k_{\&-1} \cdot X e^{3,\&}$, and $k_{\&-1} \leq k_{\&-2} + X$ for every $\&$ in the same partition.

³ Similarly, we can create a partition of the index set $\&$ such that the size of the partition is $d_4 k_{\&-1} \cdot X e^{3,\&}$, and $k_{\&-1} \leq k_{\&-2} + X$ for every $\&$ in the same partition.

Let $I = \&_1 \dots \&_g$ be the product partition of the above two partitions. Then

$$|I| = d_4 \max\{k_{\&-1} - k_{\&-2}, 1\} \cdot X e^{23,\&}, \quad \text{and for every } \& \in I \text{ we have } k_{\&-1} \leq k_{\&-2} + X \text{ and } k_{\&-1} \leq k_{\&-2} + X.$$

³We can create such a partition via the packing-covering duality in the following way (see e.g. Theorem 14.1, Theorem 14.2, and Example 14.11 [61]). First we create a maximal packing of the ball $\& \in \mathbb{R}^{3,\&}$ using balls with radius $X/4$ greedily, where we greedily choose points $\&_1, \dots, \&_g$ such that the balls $\&_i \in \mathbb{R}^{3,\&}$ are disjoint. We stop when no more such points can be added in $\& \in \mathbb{R}^{3,\&}$. Then, by the packing-covering duality, the balls $\&_i \in \mathbb{R}^{3,\&}$ cover $\& \in \mathbb{R}^{3,\&}$. This is because any uncovered point can be added to the packing we created before, which contradicts the maximality of the packing.

have

$$\begin{aligned}
 j_{g^0} &= j_{g^0} : g^0, j_{g^1} = j_{g^1} : g^1 \\
 j_{g^1} &= j_{g^1} : g^1, j_{g^2} = j_{g^2} : g^2 \\
 k_{g^2} &= k_{g^2} : g^2, k_{g^1} = k_{g^1} : g^1 \\
 2X \max_{k \in k_{g^2}} & k_{g^2} - k_{g^1}
 \end{aligned}$$

Then, because $2X \max_{k \in k_{g^2}} k_{g^2} - k_{g^1} \leq 1$, we have

$$\begin{aligned}
 \frac{\exp^{j_{g^0}}}{\exp^{j_{g^1}}} &\leq \exp^{j_{g^1} - j_{g^0}} \\
 j_{g^1} - j_{g^0} &\leq 2X \max_{k \in k_{g^2}} k_{g^2} - k_{g^1} \\
 4X \max_{k \in k_{g^2}} &k_{g^2} - k_{g^1}
 \end{aligned}$$

where the last inequality follows by $\exp^{G^0} \leq 2G$ for $0 \leq G \leq 1$. Therefore, because $4X \max_{k \in k_{g^2}} k_{g^2} - k_{g^1} \leq 1$, we have

$$\begin{aligned}
 \frac{F_{g^0}}{F_{g^1}} &\leq \frac{\exp^{j_{g^0}}}{\exp^{j_{g^1}}} \\
 &\leq \frac{1 + 4X \max_{k \in k_{g^2}} k_{g^2} - k_{g^1} g^0}{1 + 4X \max_{k \in k_{g^2}} k_{g^2} - k_{g^1} g^1} \\
 &\leq \frac{4 + 35^2}{34^2} 4X \max_{k \in k_{g^2}} k_{g^2} - k_{g^1} \\
 &\leq 17X \max_{k \in k_{g^2}} k_{g^2} - k_{g^1}
 \end{aligned}$$

where the first inequality follows since

$$\begin{aligned}
 1 + 4X \max_{k \in k_{g^2}} k_{g^2} - k_{g^1} &\leq \frac{\exp^{j_{g^0}}}{\exp^{j_{g^1}}} \\
 &\leq 1 + 4X \max_{k \in k_{g^2}} k_{g^2} - k_{g^1}
 \end{aligned}$$

, and the second inequality follows since $\frac{4 + 35^2}{34^2} \leq 1 + 140 \cdot 10^{-20} G$ for every $0 \leq G \leq 1$ and $0 \leq G \leq 1$.

Now we construct our desired index set. Let (g_1, \dots, g_ℓ) be a partition of the index set \mathcal{I} such that $g_i = g_j$ for every $g^0 \geq g^1$ and $g_i \geq g_j$. We first construct an index set $g^0 = \{g^0\}$ for each $g^0 \geq g^1$, and then combine them to obtain \mathcal{I} .

Fix any $g^0 \geq g^1$. For each index set g^0 , we only choose indices out of it to include in g^0 , namely the indices that are approximately the highest indices in $f^1 + D^0_{g^0}$.⁴ Specifically, for each $g^0 \geq g^1$, we run the given X -Approximate

⁴For simplicity we assume $(g^0 \setminus g^1) \cap g^1 = \emptyset$. Otherwise we simply choose all indices $(g^0 \setminus g^1)$.

Let \mathcal{D} be the set of points $\{x_i\}_{i=1}^n \in \mathbb{R}^3$, and query \mathcal{D} , numbers α and X as inputs. We let \mathcal{J}_g be the collection of all output indices for each $g \in \mathcal{G}$. Then because $\alpha = d \max_{k \in \mathcal{K}} \|k - k_{2-1}\|_g \cdot X e^{23 \cdot \alpha}$, we have

$$j \in \mathcal{J}_g = \{i : d \max_{k \in \mathcal{K}} \|k - k_{2-1}\|_g \cdot X e^{23 \cdot \alpha} \leq \alpha\}$$

and the expected amortized runtime of constructing each \mathcal{J}_g is

$$d \max_{k \in \mathcal{K}} \|k - k_{2-1}\|_g \cdot X e^{23 \cdot \alpha} \cdot |\mathcal{J}_g| \leq \alpha \cdot |\mathcal{J}_g| \leq \alpha \cdot n$$

Let $\mathcal{G} = \{g^0, g^1, \dots, g^L\}$. Then we have

$$|\mathcal{J}_g| = g : d \max_{k \in \mathcal{K}} \|k - k_{2-1}\|_g \cdot X e^{23 \cdot \alpha}$$

and the expected amortized runtime of constructing \mathcal{G} is

$$\sum_{g^0=1}^L d \max_{k \in \mathcal{K}} \|k - k_{2-1}\|_{g^0} \cdot X e^{23 \cdot \alpha} \cdot |\mathcal{J}_{g^0}| \leq \alpha \cdot \sum_{g^0=1}^L |\mathcal{J}_{g^0}| \leq \alpha \cdot n$$

The inequality follows since $\sum_{g^0=1}^L |\mathcal{J}_{g^0}| = n$ and $\alpha \cdot n$ is concave in α .

Finally, we show that $\text{OPT}_{P^0} \leq 11 \cdot 6^{1/n} \cdot \text{OPT}_P \leq 11 \cdot n^{1/n}$ with our choice of X .

Let $\mathcal{S} = \{s_1, \dots, s_m\}$ be the non-zero coordinates of an optimal solution to the original problem (for simplicity we assume \mathcal{S} has m non-zero entries, and other cases can be handled similarly). For each $s_i = 1 - \epsilon_i$, let \mathcal{K}_i be the index such that \mathcal{K}_i and \mathcal{K}_i are in the same $\mathcal{G}^0 \setminus \mathcal{G}^1$ and $1 + D_{\mathcal{K}_i}^0 \leq 1 + D_{\mathcal{K}_i}^1 \leq X$. Then $\mathcal{S}_{\mathcal{K}_i} = \mathcal{S}_i$. Let $G \in \mathcal{G}$ such that $G_{\mathcal{K}_i} = 1$, then G is feasible to P^0 . We prove that the objective value of P^0 at G is at least $11 \cdot 6^{1/n} \cdot \text{OPT}_P \leq 11 \cdot n^{1/n}$. Indeed, because $n \leq 34X \max_{k \in \mathcal{K}} \|k - k_{2-1}\|_g + D_{\max}^0 \cdot X$, then for each $s_i = 1 - \epsilon_i$: we have

$$\begin{aligned}
 & \frac{5_{\&}}{5_{\&}} \frac{1F_{\&} + D^0 > G}{F_{\&} > G} \\
 = & \frac{5_{\&}}{5_{\&}} \frac{\prod_{g=1}^{\&} 1F_{\&}^0 \&_g^1 + D_{\&}^0}{\prod_{g=1}^{\&} 1F_{\&}^0 \&_g} \\
 & \frac{5_{\&}}{5_{\&}} \frac{\prod_{g=1}^{\&} 1F_{\&}^0 \&_g^1 + D_{\&}^0}{\prod_{g=1}^{\&} 1F_{\&}^0 \&_g} X \\
 & \frac{5_{\&}}{5_{\&}} \frac{11 \ n^0 \prod_{g=1}^{\&} 1F_{\&}^0 \&_g^1 + D_{\&}^0}{11, n^0 \prod_{g=1}^{\&} 1F_{\&}^0 \&_g} X \\
 & \frac{5_{\&}}{5_{\&}} \frac{\prod_{g=1}^{\&} 1F_{\&}^0 \&_g^1 + D_{\&}^0}{\prod_{g=1}^{\&} 1F_{\&}^0 \&_g} 34X \max\{k_{2-1} - k_{2-1}g^1 + D_{\max}^0\} X \\
 & \frac{5_{\&}}{5_{\&}} \frac{\prod_{g=1}^{\&} 1F_{\&}^0 \&_g^1 + D_{\&}^0}{\prod_{g=1}^{\&} 1F_{\&}^0 \&_g} n \\
 & 11 \ 6^{1n^0} 5_{\&} \frac{\prod_{g=1}^{\&} 1F_{\&}^0 \&_g^1 + D_{\&}^0}{\prod_{g=1}^{\&} 1F_{\&}^0 \&_g} 1n^0
 \end{aligned}$$

where the first inequality follows since $\prod_{g=1}^{\&} 1F_{\&}^0 \&_g^1 + D_{\&}^0 \geq X$, the second inequality follows since $11 \ n^0 \prod_{g=1}^{\&} 1F_{\&}^0 \&_g^1 + D_{\&}^0 \geq 11, n^0 \prod_{g=1}^{\&} 1F_{\&}^0 \&_g$ for every $\& - \&$ that are in the same partition in I and $\&_g - \&$ that are in the same partition in I_n and the third inequality follows since $11 \ n^0 \geq 11, n^0 \geq 12n$ and $\prod_{g=1}^{\&} 1F_{\&}^0 \&_g^1 + D_{\&}^0 \geq \prod_{g=1}^{\&} 1F_{\&}^0 \&_g^1 + D_{\max}^0$. Then

$$\begin{aligned}
 \text{OPT}_{P1} & \leq \frac{5_{\&}}{5_{\&}} \frac{1F_{\&} + D^0 > G}{F_{\&} > G} \\
 & \leq \frac{5_{\&}}{5_{\&}} \frac{\prod_{g=1}^{\&} 1F_{\&}^0 \&_g^1 + D_{\&}^0}{\prod_{g=1}^{\&} 1F_{\&}^0 \&_g} 1n^0 \\
 & \leq 11 \ 6^{1n^0} \text{OPT}_P : 1n^0
 \end{aligned}$$

as desired. By our choice of X, we have

$$j \ j = : \frac{140 \max\{k_{2-1} - k_{2-1}g^{23}\}}{1 + D_{\max}^0 \ n}$$

and the expected amortized runtime of finding i is

$$\frac{140 \max\{k_{2-1} - k_{2-1}g^{23}\}}{n} \cdot g$$

$$\text{-ANN} = \frac{n}{g} \cdot \frac{n}{35 \max\{k_{2-1} - k_{2-1}g^{23}\} + D_{\max}^0}$$

”

Below we give the pseudo-code of our phase one algorithm in Proposition 7.

Algorithm 10: Phase One: Preprocess

// Preprocess

Input: Number of items n , maximum number of recommended items k , key matrix $R \in \mathbb{R}^{n \times k}$, query matrix $Q \in \mathbb{R}^{n \times k}$, value matrix $V \in \mathbb{R}^{n \times k}$, reward functions $f_{i,j}$, parameter $\epsilon > 0$, and n -Approximate: -Nearest Neighbor oracle

Output: Parameter $X \in \mathbb{R}^{n \times k}$, partitions $I = \{I_1, \dots, I_g\}$ and $S = \{S_1, \dots, S_g\}$, and preprocessed n -Approximate: -Nearest Neighbor oracle with set of points $f_{i,j} : \mathbb{R}^2 \rightarrow \mathbb{R}$ for each $i \in [n]$ and $j \in [k]$

Set $X \in \mathbb{R}^{n \times k}$ according to Eq. (A.1)

Let $\gamma = \max_{k \in [k]} \|k_{2-1} - k_{k-1}\|_2$

Form a maximal $\gamma X \cdot 2^0$ -separated set $I_0 = \{i_0\}$

Set balls $B_j = \{i \in I_0 : \|i - i_j\|_2 \leq \gamma\}$ with $i_j \in I_0$

Create partition $I = \{I_1, \dots, I_g\}$ of $[n]$: each I_j contains indices whose i_0 and i_j fall in the same pair of balls

Let g be the number of distinct functions among $f_{i,j} : i \in I_j, j \in [k]$

Create partition $S = \{S_1, \dots, S_g\}$ of $[n]$ where $S_j = \{i \in I_j : f_{i,j} = f_{i_0,j}\}$ for all $j \in [g]$

for $g^0 = 1$ to g do

for $g^0 = 1$ to g do

if $\|i_0 - i_j\|_2 < \gamma$; then

Preprocess the n -Approximate: -Nearest Neighbor oracle with set of points $f_{i,j} : \mathbb{R}^2 \rightarrow \mathbb{R}$

end

end

end

return Parameter $X \in \mathbb{R}^{n \times k}$, partitions $I = \{I_1, \dots, I_g\}$ and $S = \{S_1, \dots, S_g\}$, and preprocessed n -Approximate: -Nearest Neighbor oracle with set of points $f_{i,j} : \mathbb{R}^2 \rightarrow \mathbb{R}$ for each $i \in [n]$ and $j \in [k]$

Algorithm 11: Phase One: Query

// Query

Input: User vector $D \in \mathbb{R}^E$, maximum number of recommended items τ ,
 outputs of Algorithm 10: parameter $X \geq 0$, partitions $I = \{1, \dots, g\}$
 and $S = \{1, \dots, g\}$, and preprocessed n -Approximate τ -Nearest
 Neighbor oracle with set of points $\{x_8 \in \mathbb{R}^E \mid x_8 \in I\}$ for each $g \in S$
 and $\tau \geq \tau_0 \gg \frac{1}{4}$

Output: Index set of candidate items

Initialize \mathcal{C} ;

for $g^0 = 1$ to g do

```

  for  $g^0 = 1$  to  $\tau$  do
    if  $g^0 \in S$ ; then
      Query then-Approximate:  $\tau$ -Nearest Neighbor oracle with set of
      points  $\{x_8 \in \mathbb{R}^E \mid x_8 \in I\}$ , query  $D$ , and numbers  $\tau$  and  $X$  as inputs
      [
    end
  end
end

```

end

return

A.9 Proof of Proposition 8.

The proof is completed according to the following steps:

1. Low Non-negative Rank Approximation: In Lemma 7, we prove that in order to approximately solve P^0 , it is sufficient to approximately solve P^1 , where P^1 is obtained by replacing τ with τ_0 .
2. Enumeration of Partial Solutions: We took guesses on some index sets $I = \{1, \dots, g\}$, which corresponds to the non-negative factorization of D . Let $I = \{1, \dots, g\}$ and let P^1 denote the problem where P^0 has additional constraints that $c_8 = 1$ for all $8 \in I$. We showed that the total number of guesses is bounded above, so it is sufficient to solve P^1 for each guess.
3. Linearization of Fractional Objective Terms: In order to solve P^1 , we linearize the fractional terms in the objective function by defining a set of auxiliary problems $P^1 - C$ for each $C \in \mathbb{R}^E$. These problems are parameterized by the denominator terms in the objective function of P^1 . In Lemma 9, we prove it suffices to find a C for which $P^1 - C$ has the highest optimal value.

among all $P \rightarrow C$, as the corresponding optimal is an optimal solution to $P^1 \rightarrow C^0$.

4. Dimensionality Reduction and Discretization of Auxiliary Problems: In order to approximately solve $P \rightarrow C$ for all $C \in \mathbb{R}^n$, we discretize C space and show in Lemma 10 that it suffices to solve $P^1 \rightarrow C^0$ for a small number of C 's.
5. Complete Linearization of Auxiliary Problems: Fix a given C , the objective functions of $P \rightarrow C$ inside \mathcal{S} has rank k . We discretized the value space of those objective functions. In Lemma 11, we showed that in order to solve $P^1 \rightarrow C^0$, it is sufficient to give an oracle that, for each discretization of the value space, identify whether there exists a feasible solution to $P \rightarrow C$ with objective values that are approximately inside the discretization.
6. Approximation of Linearized Auxiliary Problems via LP Rounding: Finally, we gave such an oracle by a rounding procedure. Lemma 12 and Lemma 13 proved that the oracle is correct by using the properties of our guess \tilde{G} .

Step 1: Low Non-negative Rank Approximation

First we bound the loss incurred by replacing G , with \tilde{G} :

Lemma 7. Let Problem $P^0 \rightarrow G$ be defined as

$$(P^0 \rightarrow G) \quad \max_{G \succeq 0} \sum_{g=1}^k \frac{F_g^0 + D^0 \cdot G}{F_g^0 \cdot G}$$

s.t. $G \succeq 0, \text{rank}(G) \leq k$

Let G be a feasible solution to $P^0 \rightarrow G$ (and hence also a feasible solution to P), and suppose G satisfies

$$\sum_{g=1}^k F_g^0 \cdot G \succeq \frac{1}{\epsilon} U^0 \text{OPT}_{P^0 \rightarrow G} \cdot V$$

Then we have

$$\begin{aligned} \sum_{g=1}^k F_g^0 \cdot G &\succeq \frac{1}{\epsilon} U^0 \left(\sum_{g=1}^k W^0 F_g^0 + D_{\max}^0 \right) \text{OPT}_{P^0 \rightarrow G} \\ &\succeq \frac{1}{\epsilon} \sum_{g=1}^k W^0 F_g^0 + D_{\max}^0 \cdot \frac{1}{\epsilon} U^0 \left(\sum_{g=1}^k W^0 F_g^0 + D_{\max}^0 \right) \text{OPT}_{P^0 \rightarrow G} \\ &\succeq \frac{1}{\epsilon} \sum_{g=1}^k W^0 F_g^0 + D_{\max}^0 \cdot \frac{1}{\epsilon} V \end{aligned}$$

Proof of Lemma 7. Let G be a feasible solution to $P^0 \rightarrow G$. Because $\sum_{g=1}^k W^0 F_g^0 + D_{\max}^0 \succeq \frac{1}{\epsilon} U^0 \left(\sum_{g=1}^k W^0 F_g^0 + D_{\max}^0 \right) \text{OPT}_{P^0 \rightarrow G}$, we have

$$\frac{\sum_{g=1}^k F_g^0 + D^0 \cdot G}{\sum_{g=1}^k F_g^0 \cdot G} \succeq \frac{\sum_{g=1}^k W^0 F_g^0 + D_{\max}^0}{\sum_{g=1}^k W^0 F_g^0} \succeq \frac{1}{\epsilon} \sum_{g=1}^k W^0 \frac{F_g^0 + D^0 \cdot G}{F_g^0 \cdot G}$$

Therefore, for any feasible solution G , we have

$$\begin{aligned}
 \sum_{i=1}^n G_i &= \sum_{i=1}^n \frac{F_i^0 + D^0 > G}{F_i^0 > G} \\
 &\leq \sum_{i=1}^n \frac{F_i^0 + D^0 > G}{F_i^0 > G} \\
 &\leq \sum_{i=1}^n \frac{F_i^0 + D^0 > G}{F_i^0 > G} \cdot 2W^1 + D_{\max}^0 \\
 &\leq (2W^1 + D_{\max}^0) \sum_{i=1}^n \frac{F_i^0 + D^0 > G}{F_i^0 > G}
 \end{aligned}
 \tag{A.2}$$

where the third inequality follows since $\frac{F_i^0 + D^0 > G}{F_i^0 > G} \leq 2W^1 + D_{\max}^0$.

Similarly, we also have

$$\frac{F_i^0 + D^0 > G}{F_i^0 > G} \leq 2W^0 \frac{F_i^0 + D^0 > G}{F_i^0 > G}$$

which gives

$$\sum_{i=1}^n G_i \leq (2W^1 + D_{\max}^0) \sum_{i=1}^n G_i \leq (2W^1 + D_{\max}^0) \sum_{i=1}^n G_i
 \tag{A.3}$$

Now Let G_{p^1} be an optimal solution to P^1 . Then by Eq. (A.2), we have

$$\begin{aligned}
 \text{OPT}_{P^1} &= \sum_{i=1}^n G_{p^1,i} \\
 &\leq (2W^1 + D_{\max}^0) \sum_{i=1}^n G_{p^1,i} \\
 &= (2W^1 + D_{\max}^0) \text{OPT}_{P^1}
 \end{aligned}$$

Finally, applying Eq. (A.3), we conclude that

$$\begin{aligned}
 \sum_{i=1}^n G_i &\leq (2W^1 + D_{\max}^0) \sum_{i=1}^n G_i \\
 &\leq (2W^1 + D_{\max}^0) \text{OPT}_{P^1} \\
 &\leq (2W^1 + D_{\max}^0) \sum_{i=1}^n G_i \\
 &\leq (2W^1 + D_{\max}^0) \sum_{i=1}^n G_i \\
 &\leq (2W^1 + D_{\max}^0) \sum_{i=1}^n G_i
 \end{aligned}$$

”

Lemma 7 shows that, in order to approximately solve P^1 it is enough to approximately solve P^1 .

Let $\mathbf{D} = \mathbf{D}^0 \geq \mathbf{0}$ be the known non-negative factorization, where $\mathbf{D} \in \mathbb{R}^{n \times n}$. Let $\mathbf{R}_8^A \in \mathbb{R}^{1 \times n}$ be the 8th row of \mathbf{A} and $\mathbf{C}_9 \in \mathbb{R}^{1 \times n}$ be the 9th column of \mathbf{C} . Then $\mathbf{F}_8^0 = \mathbf{R}_8^A \mathbf{D} \mathbf{C}_9^T$. Let

$$(A.4) \quad \mathbf{z}_9 = \mathbf{1}_9 + \mathbf{D}^0 \mathbf{e}_9$$

Then \mathbf{P}^{01} can be rewritten as:

$$(P^{01}) \quad \max_{\mathbf{z}_9} \quad \mathbf{c}_9^T \mathbf{z}_9 = \mathbf{c}_9^T (\mathbf{1}_9 + \mathbf{D}^0 \mathbf{e}_9) = \mathbf{c}_9^T \mathbf{1}_9 + \mathbf{c}_9^T \mathbf{D}^0 \mathbf{e}_9 = \mathbf{c}_9^T \mathbf{1}_9 + \mathbf{F}_8^0 \mathbf{z}_9$$

$$\text{s.t. } \mathbf{z}_9 \geq \mathbf{0}, \mathbf{z}_9 \leq \mathbf{1}_9$$

Step 2: Enumeration of Partial Solutions

Our algorithm enumerates a set of partial solutions (where a partial solution fixes the values of a subset of variables), and then for each partial solution, solves the remaining problem near-optimally. In this step we bound the total number of partial solutions, and in the next steps we show that for each partial solution, the remaining problem can be solved sufficiently fast.

Let

$$(A.5) \quad \mathbf{z}_9 = \mathbf{D}^0 \mathbf{e}_9 + \mathbf{1}_9 + \mathbf{D}^0 \mathbf{e}_9$$

Each partial solution that our algorithm considers is specified by a tuple $(\mathbf{z}_9, \mathbf{A}_9^0)$, where each $\mathbf{z}_9 \in \{0, 1\}^n$ is an index set such that $\mathbf{z}_9(j) = 1 \iff j \in \mathbf{A}_9^0$. For each $\mathbf{z}_9 \in \{0, 1\}^n$, let

$$(A.6) \quad \mathbf{A}_9^0 = \{j \in \{1, \dots, n\} \mid \mathbf{z}_9(j) = 1 \text{ and } \mathbf{c}_9(j) > \min_{j \in \mathbf{A}_9^0} \mathbf{c}_9(j)\}$$

In words, \mathbf{A}_9^0 consists of the indices outside of \mathbf{z}_9 whose coefficients in \mathbf{c}_9 are strictly greater than the minimum coefficient in \mathbf{z}_9 across indices in \mathbf{z}_9 .

We say a tuple $(\mathbf{z}_9, \mathbf{A}_9^0)$, with corresponding index sets \mathbf{z}_9 and \mathbf{A}_9^0 defined according to (A.6), is valid if $\mathbf{z}_9(j) = 1 \iff j \in \mathbf{A}_9^0$. Then every feasible solution to \mathbf{P}^{01} corresponds to a valid tuple $(\mathbf{z}_9, \mathbf{A}_9^0)$ in the following way: Let $\mathbf{z}_9 = \mathbf{z}_9^0 \in \{0, 1\}^n$. For each $\mathbf{z}_9 \in \{0, 1\}^n$, we define \mathbf{z}_9^0 to be the set of indices $\mathbf{z}_9(j) = 1$ such that $\mathbf{z}_9(j)$ is among the \mathbf{z}_9 highest values in \mathbf{z}_9 . That is, let $\mathbf{z}_9^0 = \{j \in \{1, \dots, n\} \mid \mathbf{z}_9(j) = 1 \text{ and } \mathbf{z}_9(j) \geq \mathbf{z}_9(j') \text{ for } j' \in \mathbf{z}_9\}$. Then $\mathbf{z}_9^0 = \mathbf{z}_9 \setminus \mathbf{A}_9^0$. Notice that $\mathbf{z}_9^0(j) = 1 \iff j \in \mathbf{z}_9^0$. Also, we claim that $\mathbf{z}_9^0 \setminus \mathbf{A}_9^0 = \mathbf{z}_9^0$; for each $\mathbf{z}_9 \in \{0, 1\}^n$. Supposing otherwise that $\mathbf{z}_9^0 \setminus \mathbf{A}_9^0 \neq \mathbf{z}_9^0$, then

by construction $\exists j \in \{1, \dots, n\}$ such that $c_j = 1$. Because $\delta_j \in \{0, 1\}$, we have that δ_j is in the image of c_j , so there exists δ_j such that $c_j \delta_j = \delta_j$. Therefore $c_j \delta_j \in \{0, 1\}$, which shows $\delta_j \in \{0, 1\}$. This contradicts $\delta_j \in \{0, 1\}$. Therefore $\delta_j \in \{0, 1\}$ for each $j \in \{1, \dots, n\}$. Hence we have $(\delta_1, \dots, \delta_n) \in \{0, 1\}^n$. Therefore $(\delta_1, \dots, \delta_n)$ is indeed a valid tuple.

The notion of correspondence to valid tuples forms a partition of the set of feasible solutions to P^0 , so it suffices to solve P^0 separately for each subset of this partition. This is formally stated in the following result:

Lemma 8. Suppose we are given an oracle ALG^0 that takes P^0 , any valid tuple $(\delta_1, \dots, \delta_n)$, and any $X_j \geq 0$ as inputs, and outputs a solution $G_{(\delta_1, \dots, \delta_n)}$ of P^0 that satisfies

1. $G_{(\delta_1, \dots, \delta_n)}$ corresponds to $(\delta_1, \dots, \delta_n)$, and
2. for any $G_{(\delta_1, \dots, \delta_n)}$ that corresponds to $(\delta_1, \dots, \delta_n)$, we have

$$\sum_{j=1}^n c_j G_{(\delta_1, \dots, \delta_n)} \leq \sum_{j=1}^n c_j X_j + \sum_{j=1}^n \delta_j X_j$$

where $0 \leq X_j \leq 1$ and $X_j = 0$,

with runtime $O(n)$. Then there exists an algorithm ALG that takes P^0 and any $X_j \geq 0$ as inputs, and outputs a solution of P^0 that satisfies

$$ALG_{P^0} \leq \sum_{j=1}^n c_j X_j + \sum_{j=1}^n \delta_j X_j$$

with runtime

$$O(n \log_2 \sum_{j=1}^n X_j)$$

Proof of Lemma 8. We will construct ALG explicitly. Now because every feasible solution I to P^0 corresponds to exactly one valid tuple, we can solve P^0 by enumerating all valid tuples, and applying ALG^0 to each valid tuple. It turns out that pre-sorting the vector c_j allows for more-efficient enumeration. Let ALG take the following steps:

1. Sort c_j for each $j \in \{1, \dots, n\}$. This takes runtime $O(n \log_2 n)$.
2. Enumerate all valid tuples with $j = 1, \dots, n$, and record the unique corresponding feasible solutions. We will show momentarily that when $j = 1, \dots, n$, there is a unique corresponding feasible solution, and as a result this step takes runtime $O(n)$.

3. Enumerate all valid tuples with $j - 1j = _$, and record the solution output by $ALG_{P^{01} \ 0}^0$ for each such valid tuple. We will show that it takes runtime $_A^2 < -^A$ to enumerate all such valid tuples, and then because there are at most $< -^A$ such valid tuples, the total runtime of this step is $_A^2 < -^A < -^A) 1X^0$.
4. Output a solution that is recorded with the highest objective value in $P^{01} \ 0$.

Therefore the total runtime of ALG is

$$A < \log_2 < _ < _ _ A^2 < -^A _ < -^A) 1X^0$$

Finally, let $G_{_1 \ \dots \ _A}^0$ be an optimal solution of P^0 where $_1 \ \dots \ _A^0$ is the valid tuple that it corresponds to. Let $G_{_1 \ \dots \ _A}^0$ be the solution that $ALG_{P^{01} \ 0}^0$ outputs with input $P^{01} \ 0$, valid tuple $_1 \ \dots \ _A^0$, and $X_j = 0$. Then

$$\begin{aligned} ALG_{P^{01} \ 0} &= 5_{P^{01} \ 0} 1G_{_1 \ \dots \ _A}^0 \\ &= 11 \ 6 \ 01X^{00} 5_{01 \ 0} 1G_{_1 \ \dots \ _A}^0 \ 01X^0 \\ &= 11 \ 6 \ 01X^{00} OPT_{P^{01} \ 0} \ 01X^0 \end{aligned}$$

It remains to analyze Steps 2 and 3 of ALG.

Step 2: Fix any valid tuple $_1 \ \dots \ _A^0$ such that $j - 1j = _$. Assume there exists a feasible solution to $P^{01} \ 0$ that corresponds to the valid tuple, and let $_ = f8 \ 2 \ \gg < _4 \ j \ _ = 1g$. Then because $_1j = \minf _ - j / jg = j / j$ and $_1 \ / \$, we have $_1 = _ /$. Similarly, $_9 = _ /$ for all $9 \ 2 \ \gg _4$. Therefore, there exists a feasible solution to $P^{01} \ 0$ that corresponds to $_1 \ \dots \ _A^0$ only if $_1 = _ = _ A$. There are at most $_1 \ _1 < _ < _$ such tuples. Moreover, there is a unique feasible solution $_1$ that corresponds to $_1 \ \dots \ _A^0$, namely $_8 = 1$ for every $8 \ 2 \ _1$ and $_8 = 0$ for every $8 \ 8 \ _1$. Thus, the runtime of enumerating all corresponding feasible solutions is bounded by $_ <$.

Step 3: Fix any valid tuple $_1 \ \dots \ _A^0$ such that $j - 1j = _$. Then by construction we must have $_8 = 1$ for every $8 \ 2 \ [9 \ _9$, and $_8 = 0$ for every $8 \ 2 \ [9 \ _9$. Therefore every feasible solution $_1$ that corresponds to $_1 \ \dots \ _A^0$ must lie in the following set:

$$\begin{aligned} f \ 1 \ 2 \ f0 \ - \ 1g \ _8 = 1 \ 8 \ 2 \ [9 \ _9 \\ _8 = 0 \ 8 \ 2 \ [9 \ _9 \\ 1 \ 4 \ _1 \ :g \end{aligned}$$

There are $\binom{A}{j} \binom{A}{A-j}$ tuples such that $j_1 = j$. We enumerate all such tuples, and check each for validity according to the following procedure:

0. From Step 1 of ALG, let $c_g: \mathbb{R}^A \rightarrow \mathbb{R}$ be a sorting of \mathbb{R}^A according to c_g such that $c_{g-1} > c_g > c_{g+1}$.
1. Fix any given tuple (c_1, \dots, c_A) . For each $g \in \{1, \dots, A\}$, let g_1, \dots, g_{2-g} be an index such that $c_{g-1} > c_{g_1} > \dots > c_{g_{2-g}} > c_g$ for all g_1, \dots, g_{2-g} . Without loss of generality, if $c_{g-1} > c_{g_1}$ and $c_{g_1} > c_{g_2}$ for some index g_2 , we let $c_{g-1} > c_{g_1} > c_{g_2}$ for tie-breaking in the sorting. Also, if $c_{g-1} > c_{g_1}$ and $c_{g_1} > c_{g_2}$ for some index g_2 , we let $c_{g-1} > c_{g_1} > c_{g_2}$ for tie-breaking in the sorting. Then by our construction

$$\begin{aligned} \hat{c}_g &= \min_{j \in \{g_1, \dots, g_{2-g}\}} c_j \\ &= \min_{j \in \{g_1, \dots, g_{2-g}\}} c_{g-1} \\ &= c_{g-1} \end{aligned}$$

Therefore $\hat{c}_g = c_{g-1}$.

2. Note that $\sum_{g=1}^A j_g = A$. Therefore, in order to check whether $j_g = j$, we just need to count the number of overlaps among \mathcal{S}_g . Set a counter $2 = 0$ to count the overlaps. For each $g \in \{1, \dots, A\}$ and for each $g_1 \in \{g_1, \dots, g_{2-g}\}$, we check if $c_{g-1} > c_{g_1}$, and do the following:

- If $c_{g-1} > c_{g_1}$, then we have $g_1 \in \mathcal{S}_g$. We do nothing in this case.
- If $c_{g-1} > c_{g_1}$ and $c_{g-1} > c_{g_2}$, then g appears in both \mathcal{S}_g and \mathcal{S}_{g_2} . Therefore we increase 2 by 1 .
- If $c_{g-1} > c_{g_1}$ and $c_{g-1} > c_{g_2}$, then we must have $c_{g-1} > c_{g_2}$. Therefore $\sum_{g=1}^A j_g > A$, so we can terminate the process and declare that (c_1, \dots, c_A) is not a valid tuple.

We iterate all $g \in \{1, \dots, A\}$ and $g_1 \in \{g_1, \dots, g_{2-g}\}$. Notice that 2 counts the number of overlaps (with multiplicity) of elements in \mathcal{S}_g , we have $\sum_{g=1}^A j_g = \sum_{g=1}^A j_g = A$. Therefore we can check if $j_g = j$. Also, if the above procedure never encounters $c_{g-1} > c_{g_2}$, then we have $\sum_{g=1}^A j_g = A$. Therefore this procedure allows us to check the validity of (c_1, \dots, c_A) .

Because each \mathcal{S}_g is sorted, the above procedure takes a constant runtime for each $g \in \{1, \dots, A\}$ so the runtime for each (c_1, \dots, c_A) is A . Because

$\sum_{j=1}^A j = \frac{A(A+1)}{2}$, there are $\frac{A(A+1)}{2}$ combinations of g_1, \dots, g_A and g_1, \dots, g_A . Therefore the runtime to check the validity is $O(A^2)$ for any given tuple (g_1, \dots, g_A) . Because there are at most $\frac{A(A+1)}{2}$ such tuples, the runtime of enumerating all such tuples is $O(A^3)$.

”

Below we give the pseudo-code of the algorithm ALG in Lemma 8.

Algorithm 12: Phase Two: Enumeration of Partial Solutions

Input: Instance of Problem P^0 from Lemma 7, oracle ALG from Lemma 8,
parameter $X \geq 0$

Output: Solution G to P^0

Set $\epsilon \geq 0$ according to Eq. (A.5)

// Enumerate valid tuples with $j = \lfloor \epsilon \rfloor$

for $\epsilon = 1$ to $\epsilon - 1$ do

 for each $\epsilon \in \mathcal{E}$ with $1 \leq \epsilon \leq \epsilon$ do

 Set $G = f(\epsilon)$ where $f(\epsilon) = \text{ALG}(P^0, \epsilon)$; record $\text{value}(G)$

 end

end

// Enumerate valid tuples with $j = \epsilon$

for $\epsilon \in \mathcal{E}$ do

$\epsilon = \epsilon + D^0$

ϵ permutation sorting \mathcal{E} by ϵ descending

end

for each tuple $\epsilon_1, \dots, \epsilon_{\epsilon}$ where $\epsilon_j = \epsilon$ for all j do

 for $\epsilon \in \mathcal{E}$ do

$\epsilon = \arg \min_{\epsilon \in \mathcal{E}} \text{value}(\epsilon)$

$\epsilon = \text{value}(\epsilon) - \epsilon$; $\epsilon = \text{value}(\epsilon) - \epsilon$

 end

$\epsilon = 0$; valid = true

 for $\epsilon \in \mathcal{E}$ and $\epsilon \in \mathcal{E}$ do

 for $\epsilon \in \mathcal{E}$ do

 if $\text{value}(\epsilon) < \text{value}(\epsilon)$ then $\epsilon = \epsilon, 1$

 if $\text{value}(\epsilon) < \text{value}(\epsilon)$ and $\epsilon \in \mathcal{E}$ then valid = false; break

 end

 if not valid then break

 end

 if valid and $\epsilon = \epsilon$ then

 Let G be the output of ALG with inputs P^0 where valid tuple

$\epsilon = \epsilon_1, \dots, \epsilon_{\epsilon}$, and parameter X

 record $\text{value}(G)$

 end

end

return recorded G with maximum $\text{value}(G)$

The consequence of this result is that we have reduced to the task of, for each valid tuple $\alpha = (\alpha_1, \dots, \alpha_n)$ with $\sum_{j=1}^n \alpha_j = 1$, solving P^{α} with the additional constraint that the solution must correspond to the valid tuple, as stated in Problem 9 below:

$$\begin{aligned}
 (P^{\alpha}) \quad & \max_{G \in \mathbb{R}^n} \sum_{g=1}^n \alpha_g G_g \\
 & \text{s.t. } G \geq f_0 - \tau g \\
 & \quad 1 - 4^> G \geq - \\
 & \quad G_g = 1 \quad \text{if } g \in \mathcal{G} \\
 & \quad G_g = 0 \quad \text{if } g \in \mathcal{G}^c
 \end{aligned}$$

Step 3: Linearization

In order to solve P^{α} , we define the following auxiliary problem $P^{\alpha} - C$ for each $C \in \mathbb{R}$:

$$\begin{aligned}
 (P^{\alpha} - C) \quad & \max_{G \in \mathbb{R}^n} \sum_{g=1}^n \alpha_g G_g - C \\
 & \text{s.t. } G \geq f_0 - \tau g \\
 & \quad 4^> G \geq - \\
 & \quad F_g^{\alpha} G_g \leq C \quad \forall g \in \mathcal{G} \\
 & \quad G_g = 1 \quad \text{if } g \in \mathcal{G} \\
 & \quad G_g = 0 \quad \text{if } g \in \mathcal{G}^c
 \end{aligned}$$

Note that we have dropped the constraint $4^> G \geq -$ in $P^{\alpha} - C$. This is inconsequential: because we assume OPT_P is positive, the solution G with all entries equal to zero is not an optimal solution to P . Indeed, the only reason we have maintained the $1 - 4^> G \geq -$ constraint until now has been to rule out notational edge cases (such as dividing by zero).

We prove an important property regarding the relationship between optimal solutions of P^{α} and those of $P^{\alpha} - C$.

Lemma 9. Fix any valid tuple α . Let $C \in \mathbb{R}$ and let G^* be an optimal solution to $P^{\alpha} - C$. Then, $G^* \geq C$, and G^* is an optimal solution to P^{α} .

Proof of Lemma 9. First, we show that $C^0 = C$. Suppose otherwise, and let $C^0 = \tilde{C}$. Then $\tilde{C}_8 < C_8$ for every $8 \in \{1, \dots, n\}$ and $\tilde{C}_8 > C_8$ for some 8 . Therefore

$$\text{OPT}_{P^1-C^0} = \sum_{8=1}^n \frac{\tilde{C}_8}{C_8} \sum_{9=1}^A \frac{0_{89} \tilde{C}_9}{C_9} > \sum_{8=1}^n \frac{C_8}{C_8} \sum_{9=1}^A \frac{0_{89} C_9}{C_9} = \text{OPT}_{P^1-C}$$

contradicting the definition of C

Now we show that C is an optimal solution to the P^0 . For the sake of contradiction, suppose that \tilde{C} gives a higher objective value than C to P^1-C^0 , that is,

$$\sum_{8=1}^n \frac{\tilde{C}_8}{C_8} \sum_{9=1}^A \frac{0_{89} \tilde{C}_9}{C_9} > \sum_{8=1}^n \frac{C_8}{C_8} \sum_{9=1}^A \frac{0_{89} C_9}{C_9} = \text{OPT}_{P^1-C^0}$$

Let $C^0 = \tilde{C}$. Then \tilde{C} is a feasible solution to P^1-C^0 . Therefore, we have

$$\text{OPT}_{P^1-C^0} = \sum_{8=1}^n \frac{\tilde{C}_8}{C_8} \sum_{9=1}^A \frac{0_{89} \tilde{C}_9}{C_9} > \text{OPT}_{P^1-C^0}$$

contradicting the definition of C

Because $C^0 = \tilde{C}$ and $\tilde{C}_8 < C_8$ for every $8 \in \{1, \dots, n\}$, we have $\tilde{C}_8 < C_8$ for every $8 \in \{1, \dots, n\}$. Thus, from here on we will only consider the problems P^1-C with $C_8 > \frac{0_{89} C_9}{C_9}$.

Step 4: Dimensionality Reduction and Discretization of Auxiliary Problems:

Lemma 9 shows that, in order to solve P^1-C , it is enough to solve for

$$\arg \max_{C \in \mathbb{R}^n} \text{OPT}_{P^1-C}$$

Below we show that it suffices to solve the auxiliary problem P^1-C for a smaller, discretized set of C 's.

Lemma 10. Suppose we are given an oracle ALG^0 that takes P^1-C with any $C \in \mathbb{R}^n$ and any $X \geq 0$ as inputs, and outputs a solution of P^1-C that satisfies

$$\text{ALG}_{P^1-C}^0 \geq (1 - \epsilon) \text{OPT}_{P^1-C} + \epsilon X$$

with runtime $O(n)$, where $0 < \epsilon < 1$ and $X \geq 0$. Then there exists an algorithm ALG that takes P^1-C and any $X \geq 0$ as inputs, and outputs a solution of P^1-C that satisfies

$$\text{ALG}_{P^1-C} \geq (1 - \epsilon) \text{OPT}_{P^1-C} + \epsilon \frac{1}{n} \sum_{9=1}^A \frac{0_{89} W^0 X}{C_9}$$

with runtime

$$O\left(\frac{1}{\epsilon} \sum_{9=1}^A \frac{0_{89} W^0}{C_9}\right)$$

Proof of Lemma 10. Recall that we have known non-negative factorization $P = UV^T$, where $U \in \mathbb{R}^{n \times k}$. First, we make the following observation on the scale of U and V :

Observation 4. There exists $U_0 \in \mathbb{R}^{n \times k}$, where, $U_0 = \frac{1}{\|U_0\|} U$, such that $\|U_0\|_1 = 1$, $\|U_0\|_2 \leq \frac{1}{\sqrt{k}}$ and $\|U_0\|_F = 1$ for every $n \geq 2$ and $k \geq 1$.

Proof of Observation 4. We construct U_0 and V_0 explicitly. Let the rows of U_0 be the rows of U that are rescaled so that $\|U_0\|_1 = 1$. That is, let $U_{0i} = \frac{1}{\|U_{i\cdot}\|_1} U_{i\cdot}$ for every $i \in [n]$ and $\|U_{i\cdot}\|_1 = \sum_{j=1}^k U_{ij}$. Let the columns of U_0 be the columns of U_0 that are rescaled accordingly. That is, $U_{0j} = \frac{1}{\|U_{\cdot j}\|_1} U_{\cdot j}$ for every $j \in [k]$ and $\|U_{\cdot j}\|_1 = \sum_{i=1}^n U_{ij}$. Then we have $\|U_0\|_1 = 1$. Therefore, $U_0 = \frac{1}{\|U_0\|_2} U$. Finally, since $\|U\|_F = 1$, $\|U_0\|_F = 1$ and $\|U_0\|_2 = \frac{1}{\sqrt{k}}$ for every $n \geq 2$ and $k \geq 1$, we have $\|U_0\|_2 \leq \frac{1}{\sqrt{k}}$. Therefore for every $n \geq 2$ and $k \geq 1$, we have

$$\|U_0\|_2 = \frac{\|U_0\|_F}{\|U_0\|_1} = \frac{1}{\|U_0\|_1} = \frac{1}{\sum_{j=1}^k \|U_{\cdot j}\|_1} = \frac{1}{\sum_{j=1}^k \sum_{i=1}^n U_{ij}} = \frac{1}{\|U\|_1} = \frac{1}{\sum_{i=1}^n \sum_{j=1}^k U_{ij}} = \frac{1}{\sum_{i=1}^n \|U_{i\cdot}\|_1} = \frac{1}{\sum_{i=1}^n \frac{1}{\|U_{i\cdot}\|_1} \|U_{i\cdot}\|_1} = \frac{1}{\sum_{i=1}^n 1} = \frac{1}{n} \leq \frac{1}{\sqrt{k}}$$

By Observation 4, we may assume $\|U\|_1 = 1$ and $\|U\|_2 = \frac{1}{\sqrt{k}}$ for every $n \geq 2$ and $k \geq 1$ from now on. Let

$$(A.7) \quad X = \frac{1}{\|U\|_2} U V^T \text{ is a feasible solution to } \mathcal{P}_g \text{ in } \mathbb{R}^A.$$

We will partition the space \mathbb{R}^A by partitioning \mathcal{P}_g . Let X^0 be the quantity

$$(A.8) \quad X^0 = \frac{1}{\|U\|_2} \frac{1}{\sum_{j=1}^k \|U_{\cdot j}\|_1} U V^T$$

where the reason for this choice will be specified momentarily. Notice that $\|X^0\|_2 = \frac{1}{\sqrt{k}}$ for every $k \geq 1$. As seen in the proof of Proposition 7, we can create a cover of a ball with radius $\frac{1}{\sqrt{k}}$ in \mathbb{R}^A using $d^A \frac{1}{\sqrt{k}}$ balls with radius $\frac{1}{\sqrt{k}}$. Therefore we can create a partition $\mathcal{Y} = \{Y_1, \dots, Y_H\}$ of \mathcal{P}_g such that $\|Y_i\|_2 \leq \frac{1}{\sqrt{k}}$ and $\|Y_i\|_1 \leq \frac{1}{\sqrt{k}}$ for every $i \in [H]$.

all G . As a consequence of Lemma 9, we have $OPT_{P^1,0} = OPT_{P^1,0} - C^0$. Then

$$\begin{aligned}
 OPT_{P^1,0} &= OPT_{P^1,0} - C^0 \\
 &= \sum_{g=1}^8 \frac{G^5 \cdot \prod_{g=1}^A \frac{0_{8g} 3^g}{G}}{1 \cdot 6^{11} X^{11} W^0, \min^0} \\
 &\leq \frac{\sum_{g=1}^8 \frac{G^5 \cdot \prod_{g=1}^A \frac{0_{8g} 3^g}{G}}{1 \cdot 6^{11} X^{11} W^0, \min^0}}{OPT_{P^1,0} - C^0} \cdot \frac{1 \cdot 6^{11} X^{11} W^0, \min^0}{1 \cdot 6^{11} X^{11} W^0, \min^0}
 \end{aligned}$$

Rearranging the above, we have

$$OPT_{P^1,0} \cdot \frac{1 \cdot 6^{11} X^{11} W^0, \min^0}{OPT_{P^1,0} - C^0} \leq \frac{1 \cdot 6^{11} X^{11} W^0, \min^0}{OPT_{P^1,0} - C^0}$$

Therefore,

$$\begin{aligned}
 &ALG_{P^1,0} \\
 &ALG_{P^1,0}^0 \\
 &1 \cdot 6^{0_1} X^{0_0} OPT_{P^1,0} \cdot 0_1 X^0 \\
 &1 \cdot 6^{0_1} X^{0_0} \cdot 1 \cdot 6^{\frac{11, W^0 X}{0, \min}} OPT_{P^1,0} : \frac{1 \cdot 6^{11, W^0 X}}{0, \min} \cdot 0_1 X^0
 \end{aligned}$$

”

Below we give the pseudo-code of the algorithm ALG in Lemma 10.

Algorithm 13: Phase Two: Discretization of Auxiliary Problems

Input: Instance of Problem P^{1-0} , oracle ALG from Lemma 10, parameter

X_{j0}

Output: Solution G to P^{1-0}

// Discretize the space of auxiliary variable

Set X_{j0} according to Eq. (A.8)

Construct maximal X_{j0} -separated set $\{G^j\}_{j=1}^d$ with $\|G^j - G^k\| \geq d^{-1/4} X_{j0}$

Let $Y = f_1 - \dots - g_1 - \dots - g_d$

// Solve each discretized auxiliary problem

Set $X_1 = X$ and $X_2 = X$

for $j = 1$ to d do

$G^j \in H : H \subseteq X_j$

Choose arbitrary $G^j \in H$

Let G be the output of ALG with inputs P^{1-0} and parameters X_1, X_2

Record $G^j = f_1 - \dots - g_1 - \dots - g_d$

end

return $\arg \max_{G^j} f_1 - \dots - g_1 - \dots - g_d$ among recorded solutions

Step 5: Complete Linearization of Auxiliary Problems

Lemma 10 shows that, to approximately solve P , it suffices to construct an oracle that approximately solves $P - C^0$ for any given $C \geq \frac{0}{\min} - 1$, $W \leq \frac{1}{4}$. Below we give such an oracle. Let $\epsilon = 0$, $C \geq R_0^A$ for each $\epsilon \geq \frac{1}{4}$. Then $P - C^0$ can be equivalently formulated as

$$\begin{aligned}
 (P^{1-0} - C^0) \quad & \max_{G^1, \dots, G^d} f_1 - \dots - g_1 - \dots - g_d \\
 & \text{s.t. } G^j \in H^j \text{ for } j=1, \dots, d \\
 & G^j \geq 0 \text{ for } j=1, \dots, d \\
 & F_j(G^j) \leq C_j \text{ for } j=1, \dots, d \\
 & G^j = 1 \text{ for } j=1, \dots, d \\
 & G^j = 0 \text{ for } j=1, \dots, d
 \end{aligned}$$

To solve $P - C^0$, we partition the space of possible values $\{G^1, \dots, G^d\} \subseteq \mathbb{R}^A$, as well as the space of the objective value $f_1 - \dots - g_1 - \dots - g_d$.

Lemma 11. Fix any $C \in \mathbb{R}^n$, $W \in \mathbb{R}^m$. Suppose we are given an oracle with runtime $O(n^2)$ that takes $P \subseteq \mathbb{R}^n$, any $\lambda \in \mathbb{R}^n$, any $Z \in \mathbb{R}^m$, and any $X_1, X_2 \geq 0$ as inputs, and either

1. Scenario one: correctly declares that there is no feasible $P \subseteq \mathbb{R}^n$ such that $\sum_{g=1}^n g \lambda_g \geq C$ and $\sum_{g=1}^n g X_g \leq Z$, or
2. Scenario two: outputs a feasible $P \subseteq \mathbb{R}^n$ such that $\sum_{g=1}^n g \lambda_g \geq C$ and $\sum_{g=1}^n g X_g \leq Z$.

Then there exists an algorithm ALG that satisfies

$$\text{ALG}_{P \subseteq \mathbb{R}^n, C} \leq \frac{\sum_{g=1}^n g \lambda_g \geq C}{\sum_{g=1}^n g X_g \leq Z} \cdot \text{OPT}_{P \subseteq \mathbb{R}^n, C} \leq \frac{\sum_{g=1}^n g \lambda_g \geq C}{\sum_{g=1}^n g X_g \leq Z}$$

with runtime

$$O(n^2) \cdot \max_{g \in \{1, \dots, n\}} \lambda_g \leq \frac{1}{\min_{g \in \{1, \dots, n\}} \lambda_g} \cdot \max_{g \in \{1, \dots, n\}} \lambda_g \leq \frac{1}{\min_{g \in \{1, \dots, n\}} \lambda_g} \cdot \max_{g \in \{1, \dots, n\}} \lambda_g$$

Proof of Lemma 11. Let $d = \min_{g \in \{1, \dots, n\}} \lambda_g$ be the minimum entry of D (possibly negative). Because $\lambda_g \geq d$ and $\sum_{g=1}^n g \lambda_g \geq C$ for every $g \in \{1, \dots, n\}$, we have $\sum_{g=1}^n g \lambda_g \geq C \geq \sum_{g=1}^n g d \geq n d$ for every $g \in \{1, \dots, n\}$. Also, because each λ_g is non-decreasing, $\text{OPT}_{P \subseteq \mathbb{R}^n, C} \geq \sum_{g=1}^n g d \geq n d$. Similar to the proof of Lemma 10, we will create a partition of the space of possible values $\lambda \in \mathbb{R}^n$, as well as the space of the objective value $\sum_{g=1}^n g \lambda_g$. We then show that it is sufficient to solve $P \subseteq \mathbb{R}^n$ in each subset of the partition. Let $\lambda = \min_{g \in \{1, \dots, n\}} \lambda_g$, $X_1 \geq 0$ for $\lambda = 1 - \dots - d$, $\sum_{g=1}^n g \lambda_g \geq C$. Let $X_2 \geq 0$ for $\lambda = 1 - \dots - d$, $\sum_{g=1}^n g \lambda_g \geq C$. Consider all tuples $\lambda = 1 - \dots - d$. There are in total

$$O(n^2) \cdot \max_{g \in \{1, \dots, n\}} \lambda_g \leq \frac{1}{\min_{g \in \{1, \dots, n\}} \lambda_g} \cdot \max_{g \in \{1, \dots, n\}} \lambda_g$$

such tuples. Moreover, there exists a tuple $\lambda = 1 - \dots - d$ such that $\sum_{g=1}^n g \lambda_g \geq C$ for each $g \in \{1, \dots, n\}$ and $\sum_{g=1}^n g X_g \leq Z$.

For each tuple $\lambda = 1 - \dots - d$, our desired algorithm uses the given oracle to determine whether there exists a feasible $P \subseteq \mathbb{R}^n$ that satisfies the conditions in scenario two with $\lambda_g = \lambda$ for each $g \in \{1, \dots, n\}$ and $Z = Z$. Then our desired algorithm returns the $P \subseteq \mathbb{R}^n$ with the highest objective value of $P \subseteq \mathbb{R}^n$ among all tuples. Note that $P \subseteq \mathbb{R}^n$ satisfies the conditions in scenario two on the tuple $\lambda = 1 - \dots - d$. Therefore the given oracle would return some feasible $P \subseteq \mathbb{R}^n$ that satisfies the conditions in scenario two

Algorithm 14: Phase Two: Complete Linearization of Auxiliary Problems

Input: Instance of Problem $P^1 - C^0$, oracle from Lemma 11, parameters

$$X_1 - X_2 \geq 0$$

Output: Solution G to $P^1 - C^0$

// Discretize value and objective spaces

Let $a_{\min} = \min_{f \in D^0} g$ and $a_{\max} = \max_{f \in D^0} g$

Let $d = \frac{a_{\max} - a_{\min}}{8}$ and $(\delta = \frac{1}{8}, \epsilon = \frac{1}{8})$

for $i = 1$ to n do

 | $a_{\min} + X_i \delta$

end

for $B = 1$ to (δ) do

 | $\frac{0}{B} + X_2 \delta B$

end

// Enumerate and solve linearized problems

for each tuple $(i_1, \dots, i_n) \in \{1, \dots, B\}^n$ where $B \geq \frac{1}{\delta}$ and $B \geq \frac{1}{\epsilon}$ do

 | Set λ_{i_j} for each $j \in \{1, \dots, n\}$ and $Z = \frac{0}{B}$

 | Run oracle from Lemma 11 with inputs $(P - C^0)$, value tuple $(\lambda_1, \dots, \lambda_n) - Z$,

 | and parameters (X_1, X_2)

 | if oracle returns scenario two with feasible solution G then

 | Record $(G, P - C^0, G^0)$

 | end

end

return $\arg \max_{G \in \{P - C^0\}} G^0$ among recorded solutions

Step 6: Approximation of Linearized Auxiliary Problems via LP Rounding

Lemma 11 shows that, to solve $P^1 - C^0$ for any given $C \geq \frac{0}{\min} - 1$, $W \geq \frac{1}{4}$ it is enough to give an oracle described in Lemma 11. Similar to the idea of enumerating partial solutions by constructing valid tuples based on the values of λ_{i_j} , we further construct index sets based on the values $\lambda_{i_j} \delta$ and enumerate all possible index sets. Recall that via the valid tuple (i_1, \dots, i_n) , we have already fixed at least indices of any feasible solution to $P - C^0$ to be equal to i_j , namely the indices in $[i_j - \epsilon, i_j + \epsilon]$.

$$(A.9) \quad \lambda_{i_j} = d \frac{1}{2} A_{i_j}, \quad \delta = \frac{1}{8}, \quad \epsilon = \frac{1}{8}, \quad \max_{f \in D^0} g \leq \frac{1}{8} + D_{\max}^0 g \cdot \epsilon$$

Let $\mathcal{I} \subseteq \{1, \dots, n\}$ be an index set such that $0 \in \mathcal{I}$. Let

$$(A.10) \quad \mathcal{I}^0 = \{j \in \mathcal{I} \mid c_j > \min_{i \in \mathcal{I}} c_i\}$$

In words, \mathcal{I}^0 consists of the indices outside \mathcal{I} whose corresponding values c_j are strictly greater than the minimum value of c_j across indices in \mathcal{I} . Then every feasible solution to $P^1 - C^0$ corresponds to an index set \mathcal{I} in the following way: let $\mathcal{I} = \{j \in \mathcal{I} \mid c_j = 1\}$. We define \mathcal{I}^0 to be the set of indices $j \in \mathcal{I}$ such that c_j is among the $\min_{j \in \mathcal{I}} c_j$ highest values in \mathcal{I} . That is, let $c^0 : \mathcal{I} \rightarrow \mathbb{R}$ be a sorting of c_j according to c_j such that $c_{j^0} \geq c_{j^0}$. Then $\mathcal{I}^0 = \{j \in \mathcal{I} \mid c_j = c_{j^0}\}$. Similar to before, if \mathcal{I} corresponds to \mathcal{I}^0 , we must have $l_j = 1$ for $j \in \mathcal{I}^0$ and $l_j = 0$ for $j \in \mathcal{I} \setminus \mathcal{I}^0$.

As in Lemma 8, the notion of correspondence to index sets forms a partition of the set of feasible solutions to $P - C^0$. Therefore, in order to give an oracle described in Lemma 11, it suffices to give an oracle described in Lemma 11 separately for each subset of this partition. This is formally stated in the following result:

Observation 5. Suppose we are given an oracle with runtime $\mathcal{O}(n^2)$ that takes $P^1 - C^0$, any $\mathcal{I} \subseteq \{1, \dots, n\}$, any $Z \subseteq \mathbb{R}^A$, any $X_1 - X_2 \subseteq \mathbb{R}^B$, and any index set $\mathcal{I} \subseteq \{1, \dots, n\}$ as inputs, and either

1. Scenario one: correctly declares that there is no feasible $P^1 - C^0$ that corresponds to \mathcal{I} such that $c_j > 1$ for every $j \in \mathcal{I}$ and $\sum_{j \in \mathcal{I}} c_j \leq G$, or
2. Scenario two: outputs a feasible $P^1 - C^0$ that corresponds to \mathcal{I} such that $c_j > 1$, $X_1 \subseteq Z$ for every $j \in \mathcal{I}$ and $\sum_{j \in \mathcal{I}} c_j \leq G$, $X_2 \subseteq Z$.

Then there exists an oracle described in Lemma 11 with runtime $\mathcal{O}(n^2)$.

Proof of Observation 5. Because every feasible solution to $P^1 - C^0$ corresponds to exactly one index set $\mathcal{I} \subseteq \{1, \dots, n\}$ such that $0 \in \mathcal{I}$, we can give an oracle described in Lemma 11 by enumerating all such index sets and applying the oracle in Observation 5.

More specifically, for the oracle described in Lemma 11 with inputs \mathcal{I}, Z, X

- If the oracle in Observation 5 outputs \mathcal{I} in scenario two with inputs $\mathcal{I} - Z - X$ for some \mathcal{I} , then the oracle described in Lemma 11 also outputs this G in scenario two.

^ If the oracle in Observation 5 declares scenario one with inputs $\mathbf{z}, \mathbf{x}, \mathbf{x}^0$ for all \mathbf{z}^0 , then the oracle described in Lemma 11 also declares scenario one.

We can enumerate all index sets $S \subseteq \{1, \dots, n\}$ such that $|S| \leq \frac{1}{4}n$ by simply enumerating all combinations of indices from $\{1, \dots, n\}$. Because there are at most $\sum_{i=0}^{\lfloor \frac{1}{4}n \rfloor} \binom{n}{i} < 2^{\frac{1}{4}n}$ such index sets, the runtime of the oracle described in Lemma 11 is $O(2^{\frac{1}{4}n})$.

By Observation 5, it suffices to give an oracle as described. Below we give such an oracle.

First, suppose $\mathbf{z}^0 \in \mathcal{C}^0$. Assume there exists a feasible solution to $P^1 - C^0$ that corresponds to \mathbf{z}^0 . Let $\mathbf{g} = \mathbf{f}(\mathbf{z}^0) \in \mathbb{R}^m$. Then because $\mathbf{z}^0 = \arg \min_{\mathbf{z} \in \mathcal{C}^0} \|\mathbf{z}\| = \mathbf{z}^0$ and $\mathbf{g} \in \mathbb{R}^m$, we have $\mathbf{g} = \mathbf{f}(\mathbf{z}^0)$. Therefore, there is a unique feasible solution of $P^1 - C^0$ that corresponds to \mathbf{z}^0 , namely $\mathbf{g} = 1$ for every $\mathbf{z} \in \mathcal{C}^0$ and $\mathbf{g} = 0$ for every $\mathbf{z} \notin \mathcal{C}^0$. Therefore, if $\mathbf{z}^0 \in \mathcal{C}^0$, the oracle in Observation 5 can directly check this unique feasible solution of $P^1 - C^0$ that corresponds to \mathbf{z}^0 , and outputs the correct scenario accordingly.

From now on, we assume $\mathbf{z}^0 \notin \mathcal{C}^0$. Fix any $\mathbf{z} \in \mathbb{R}^n$, any $\mathbf{z}^0 \in \mathbb{R}^n$, any $\mathbf{x}_1, \mathbf{x}_2 \in \mathbb{R}^n$, and any \mathbf{z}^0 . The oracle essentially needs to determine the existence of a feasible binary solution to a system of linear inequalities, which is NP-hard in general. However, the linear constraints of $P - C^0$ lie in a lower-dimensional subspace, a structure we can exploit by solving a relaxation of the system, obtained by replacing binary variables with continuous ones, and rounding its solution back to a binary solution. Because of the rounding, it is possible that the values $\mathbf{g} = \mathbf{f}(\mathbf{z})$ and $\mathbf{g} = \mathbf{f}(\mathbf{z}^0)$ of the rounded solution are out of the desired ranges. However, the valid tuple and the index set S we fixed before ensures that the gaps between the values and the desired ranges are within small constants.

Note that $\|l\|_1 \leq 1$; since the projection of l on \mathbb{R}^n is in $\|l\|_1 \leq 1$. Because $\|l\|_1 \leq 1$ has $2n$ linear inequalities other than the inequality $\|l\|_1 \leq 1$ for $n \geq 2$, we can compute a vertex of $\|l\|_1 \leq 1$ with at most $2n$ fractional components with runtime $O(n^2)$ (see a standard textbook on linear programming, e.g., [155]).

Let $l \in \mathbb{R}^n$ where

$$l_j = \begin{cases} 1 & \text{if } j \in S_1 \\ 0 & \text{if } j \in S_2 \\ -1 & \text{if } j \in S_3 \end{cases}$$

Then l has at most $2n$ fractional components. We show that l is feasible. Because $l_j = 1$ for $j \in S_1$ and $l_j = 0$ for $j \in S_2$, the last three sets of constraints of P is satisfied. By the first set of constraints of P we have $\sum_{j \in S_1} l_j = 1$. Because $l_j = 0$ for every $j \in S_2$, the first set of constraints of P is satisfied. Similarly the second constraint of P is also satisfied. By the third set of constraints of P we have

$$\sum_{j \in S_1} l_j = \sum_{j \in S_2} l_j = \sum_{j \in S_3} l_j = 0$$

so the third set of constraints of P is satisfied. Similarly the fourth constraint of P is also satisfied. Therefore $l \in P$ is the desired point. \square

Let l be the point obtained in Lemma 12. This satisfies all the constraints of P except the integrality constraints. We round down to obtain a feasible solution: let $\tilde{l} \in \mathbb{R}^n$ where $\tilde{l}_j = \lfloor l_j \rfloor$ for each j . Notice that since $l_j \in [0, 1]$ for every j , we have $\tilde{l}_j \in \{0, 1\}$ and $\tilde{l}_j \leq l_j \leq C_j$ for every $j \in [n]$. Therefore \tilde{l} is feasible to P . Moreover, since $l_j = 1$ for $j \in S_1$ and $l_j = 0$ for $j \in S_2 \cup S_3$, we have that \tilde{l} corresponds to l^0 .

In the final step, we show that by setting \tilde{l} appropriately \tilde{l} satisfies the conditions in scenario two of Observation 5, hence completing the oracle in Observation 5.

Lemma 13. Set

$$\tilde{l} = d^1 A_{\cdot, 2}^{0,1} + D_{\max}^0 \cdot X_1 e$$

and

$$\tilde{l}^0 = d^1 A_{\cdot, 2}^{0,0} : \max_{\|l\|_1 \leq 1} f^0 + D_{\max}^0 g \cdot X_2 e$$

Then $\tilde{l}_j \in \{0, 1\}$ for every $j \in [n]$, and $\tilde{l}_j \leq C_j$ for every $j \in [n]$.

Proof of Lemma 13. Because $2 \leq \dots$, we have $3 \geq | \dots |$ for every $9 \geq 2 \dots$ and $\sum_{g=1}^{\dots} \dots \leq \dots$. Fix $9 \geq 2 \dots$ and let $2 - 9$ be an index where $3 \geq | \dots | = \min_{0 \leq j < 9} \dots$. Then since $j \geq \dots$, we have

$$3 \geq | \dots | \geq \sum_{0 \leq j < 9} \dots$$

On the other hand, \dots is obtained by rounding down. Notice that $\dots \leq 2 \dots - 1$ for all $0 \leq j < 9$, that is, for all $0 \leq j < 9$ such that $3 \geq | \dots | \geq 3 - j$. Therefore for all $0 \leq j < 9$ such that $3 \geq | \dots | \geq 3 - j$, we have $\dots = \lfloor \dots \rfloor$. By Lemma 12, \dots has at most $2 \dots$ fractional components. Therefore, because $3 \geq | \dots | \leq 1 + D_{\max}^0$, we have

$$\begin{aligned} 3 \geq | \dots | &\leq \sum_{0 \leq j < 9} \dots \\ &\leq \sum_{0 \leq j < 9} \dots \\ &\leq \sum_{0 \leq j < 9} \dots \\ &\leq \dots \end{aligned}$$

Similarly, let $7 \geq 2 - 0$ be an index where $5 \geq | \dots | = \min_{0 \leq j < 7} \dots$. Then since $j - 0 \geq \dots$, we have

$$\sum_{8=1}^{\dots} \dots \geq \sum_{0 \leq j < 7} \dots$$

On the other hand, \dots is obtained by rounding down. Notice that $\dots \leq 2 \dots - 1$ for all $0 \leq j < 7$, that is, for all $0 \leq j < 7$ such that $5 \geq | \dots | \geq 5 - j$. Therefore for all $0 \leq j < 7$ such that $5 \geq | \dots | \geq 5 - j$, we have $\dots = \lfloor \dots \rfloor$. By Lemma 12, \dots has at most $2 \dots$ fractional components. Therefore, because $\sum_{8=1}^{\dots} \dots \leq 1 + D_{\max}^0$, we have

$$\begin{aligned} \sum_{8=1}^{\dots} \dots &\leq \sum_{8=1}^{\dots} \dots \\ &\leq \sum_{8=1}^{\dots} \dots \\ &\leq \sum_{8=1}^{\dots} \dots \\ &\leq \dots \end{aligned}$$

Below we give the pseudo-code of the oracle described in Lemma 11.

Algorithm 15: Phase Two: Approximation of Linearized Auxiliary Problems

Input: Instance of Problem $P^1 \text{---} C^0$, thresholds $\delta, \epsilon, \gamma, \eta, \theta, \tau, \rho, \alpha, \beta, \lambda, \mu, \nu, \xi, \zeta, \eta, \theta, \tau, \rho, \alpha, \beta, \lambda, \mu, \nu, \xi, \zeta$, parameters $\gamma, \eta, \theta, \tau, \rho, \alpha, \beta, \lambda, \mu, \nu, \xi, \zeta$

Output: One of two scenarios:

1. Scenario 1: Certify no feasible G to $P^1 \text{---} C^0$ satisfies $\|G\|_1 \leq \delta$ and $\sum_{i=1}^n G_i \geq \epsilon$.
2. Scenario 2: Return feasible G to $P^1 \text{---} C^0$ with $\|G\|_1 \leq \delta$ and $\sum_{i=1}^n G_i \geq \epsilon$.

Set $\epsilon_j = 0$ according to Eq. (A.9)

Let $M = \lfloor n \rfloor \cdot \lfloor \epsilon \rfloor \cdot \lfloor \delta \rfloor$

// Check small solutions

for each M with $0 \leq j \leq M - 1$ do

Set ϵ_j according to Eq. (A.10)

Define $l_j = 1$ if $\lfloor \epsilon_j \rfloor \cdot \lfloor \delta \rfloor \geq \epsilon$, else $l_j = 0$

if l_j feasible and satisfies thresholds with slack $\leq \delta$ then
 return Scenario 2 with solution l_j

end

end

// Solve via LP rounding for larger solutions

for each M with $j = 0$ do

Construct polyhedron P_j according to Eq. (A.11)

if $P_j \neq \emptyset$; then

Choose arbitrary $l_j \in P_j$

Construct l_j^0 according to Eq. (A.12)

Compute vertex H_j of P_j with 2 fractional components

Set $H_j = \lfloor n \rfloor \cdot \lfloor \epsilon \rfloor \cdot \lfloor \delta \rfloor$

Define $l_j^0 = 0$ if $\lfloor \epsilon \rfloor \cdot \lfloor \delta \rfloor \geq \epsilon$
 $l_j^0 = H_j$ if $\lfloor \epsilon \rfloor \cdot \lfloor \delta \rfloor < \epsilon$

return Scenario 2 with solution l_j^0

end

end

return Scenario 1

Completing the Proof

Finally, we analyze our algorithm's overall performance and runtime. Let $n_1 = n$ and $n_2 = \lfloor n/2 \rfloor$, and in order to apply Lemma 13, we set $\alpha = 1/2$, $\beta = 2^{0.1} + D_{\max}^0/n$ and $\gamma = 1/2$. We will treat the performance guarantee and runtime analysis separately.

Performance Guarantee: Let

$$Z_{n-W_{\min}^0} = \frac{11}{6} \frac{W^0 n}{n - W_{\min}^0}$$

The algorithm ALG (for solving $P^1 - C^0$) in Lemma 11 satisfies

$$\begin{aligned} \text{ALG}_{P^1 - C^0} &\leq \frac{11}{6} \frac{X_1 11}{n - W_{\min}^0} \frac{W^0}{n} \text{OPT}_{P^1 - C^0} + 2X_2^0 : \frac{X_1 11}{n - W_{\min}^0} \frac{W^0}{n} \\ &= \frac{11}{6} \frac{12}{n - W_{\min}^0} \text{OPT}_{P^1 - C^0} : \frac{11}{6} \frac{12}{n - W_{\min}^0} \frac{W^0}{n}, \frac{12}{n - W_{\min}^0} \frac{W^0}{n} \end{aligned}$$

This gives the ALG (for solving $P^1 - C^0$) in Lemma 10 with

$$6^0 n^0 = \frac{11}{6} \frac{12}{n - W_{\min}^0} \frac{W^0}{n}$$

and

$$0_1 n^0 = \frac{11}{6} \frac{12}{n - W_{\min}^0} \frac{W^0}{n}, \frac{12}{n - W_{\min}^0} \frac{W^0}{n}$$

Therefore, the algorithm ALG (for solving $P^1 - C^0$) in Lemma 10 satisfies

$$\begin{aligned} \text{ALG}_{P^1 - C^0} &\leq \frac{11}{6} \frac{0_1 n^0}{n - W_{\min}^0} \frac{11}{6} \frac{12}{n - W_{\min}^0} \frac{W^0}{n} \text{OPT}_{P^1 - C^0} : \frac{11}{6} \frac{12}{n - W_{\min}^0} \frac{W^0}{n} \frac{0_1 n^0}{n - W_{\min}^0} \\ &= \frac{11}{6} \frac{12}{n - W_{\min}^0} \frac{0_1 n^0}{n - W_{\min}^0} \frac{11}{6} \frac{12}{n - W_{\min}^0} \frac{W^0}{n} \text{OPT}_{P^1 - C^0} : \frac{11}{6} \frac{12}{n - W_{\min}^0} \frac{W^0}{n} \frac{0_1 n^0}{n - W_{\min}^0} \\ &\quad : \frac{11}{6} \frac{12}{n - W_{\min}^0} \frac{W^0}{n}, \frac{12}{n - W_{\min}^0} \frac{W^0}{n} \\ &= \frac{11}{6} \frac{12}{n - W_{\min}^0} \frac{0_1 n^0}{n - W_{\min}^0} \text{OPT}_{P^1 - C^0} : \frac{11}{6} \frac{12}{n - W_{\min}^0} \frac{W^0}{n} \frac{0_1 n^0}{n - W_{\min}^0} \\ &\quad , \frac{12}{n - W_{\min}^0} \frac{W^0}{n} \end{aligned}$$

Therefore, the algorithm ALG (for solving P^0) in Lemma 8 satisfies

$$\begin{aligned} \text{ALG}_{P^0} &\leq \frac{11}{6} \frac{12}{n - W_{\min}^0} \frac{0_1 n^0}{n - W_{\min}^0} \text{OPT}_{P^0} \\ &\quad : \frac{11}{6} \frac{12}{n - W_{\min}^0} \frac{W^0}{n}, \frac{12}{n - W_{\min}^0} \frac{W^0}{n} \end{aligned}$$

Finally, we apply Lemma 7 by plugging in $U = \frac{11}{6} \frac{12}{n - W_{\min}^0} \frac{0_1 n^0}{n - W_{\min}^0}$ and $V = \frac{11}{6} \frac{12}{n - W_{\min}^0} \frac{W^0}{n}, \frac{12}{n - W_{\min}^0} \frac{W^0}{n}$. This gives

$$\begin{aligned} \text{ALG}_{P^0} &\leq \frac{11}{6} \frac{12}{n - W_{\min}^0} \frac{0_1 n^0}{n - W_{\min}^0} \text{OPT}_{P^0} \\ &\quad : \frac{11}{6} \frac{12}{n - W_{\min}^0} \frac{W^0}{n} \frac{0_1 n^0}{n - W_{\min}^0}, \frac{12}{n - W_{\min}^0} \frac{W^0}{n} \frac{0_1 n^0}{n - W_{\min}^0} \\ &\quad , \frac{11}{6} \frac{12}{n - W_{\min}^0} \frac{0_1 n^0}{n - W_{\min}^0} \frac{12}{n - W_{\min}^0} \frac{W^0}{n}, \frac{12}{n - W_{\min}^0} \frac{W^0}{n} \end{aligned}$$

Runtime Analysis: By Lemma 12, we give an oracle described in Observation 5 with runtime

$$)_{LP} := LP^{1 < -3 < , A , 2^0 , LP^{1 < -2 < , 2A , , 2^0 \bullet}$$

Therefore, by Observation 5, we give an oracle described in Lemma 11 with runtime

$$_{< -^0})_{LP} \bullet$$

Therefore, the algorithm ALG⁰ (for solving $P^{1 - 0}$) in Lemma 10 has runtime

$$\begin{aligned} & d^{11+D_{\max}^0} \min_{f_0-1+D_{\min}^0 g^0 \bullet \chi e^A} : \max_{82 \gg < 1/4} f_{5^{11+D_{\max}^0 g^0 \bullet \chi} }_{< -^0})_{LP} \\ & = d^{11+D_{\max}^0} \min_{f_0-1+D_{\min}^0 g^0 \bullet n e^A} \max_{82 \gg < 1/4} f_{5^{11+D_{\max}^0 g^0 \bullet n} }_{< -^0})_{LP} \bullet \end{aligned}$$

Therefore, the algorithm ALG⁰ (for solving $P^{1 - 0}$) in Lemma 8 has runtime

$$\begin{aligned} & \frac{4^{11} , W^0 : A}{n , \min} d^{11+D_{\max}^0} \min_{f_0-1+D_{\min}^0 g^0 \bullet n e^A} \\ & \max_{82 \gg < 1/4} f_{5^{11+D_{\max}^0 g^0 \bullet n} }_{< -^0})_{LP} \bullet \end{aligned}$$

Finally, by Lemma 8, our algorithm's runtime is

$$\begin{aligned} & A , < \log_2 < , - < - , - A^2 < - A , \frac{4^{11} , W^0 : A}{n , \min} \frac{1+D_{\max}^0 \min_{f_0-1+D_{\min}^0 g^0 \bullet A}}{n} \\ & \frac{\max_{82 \gg < 1/4} f_{5^{11+D_{\max}^0 g^0 \bullet A}}}{n} \quad \quad \quad)_{LP} \bullet \end{aligned}$$

Below we give the pseudo-code of our phase two algorithm in Proposition 8.

Algorithm 16: Phase Two

Input: Attention matrix $A \in \mathbb{R}^{m \times n}$, $\sigma = \text{softmax}(A^T A)$, $U \in \mathbb{R}^{m \times k}$, low-rank approximation $\hat{A} = U \Sigma U^T$ with factorization $\Sigma = \text{diag}(\sigma_1, \dots, \sigma_k)$ and element wise guarantee $\sigma_i \geq \epsilon$, $W \in \mathbb{R}^{m \times n}$, value matrix $V \in \mathbb{R}^{n \times k}$, user vector $D \in \mathbb{R}^k$, maximum number of recommended items n_r , candidate index set S , parameter $n \geq 0$.

Output: Solution G to Problem (Main).

Construct the instance of Problem (P^{1-0}) from the inputs

// Low Non-negative Rank Approximation

Form the instance of Problem (P^0) by replacing A with \hat{A} in (P^{1-0})

// Enumeration of Partial Solutions

G Run Algorithm 12 on instance (P^{1-0}) with parameter n , which internally invokes:

 // Discretization of Auxiliary Problems

- Algorithm 13 with Input: Instance (P^{1-0}) , parameter n ; Calls: Algorithm 14

 // Complete Linearization of Auxiliary Problems

- Algorithm 14 with Input: Instance (P^{1-0}) , parameters $X, n, X_2 : n$; Calls: Algorithm 15

 // Approximation of Linearized Auxiliary Problems via LP Rounding

- Algorithm 15 with Input: Instance (P^{1-0}) , thresholds $\epsilon_1, \dots, \epsilon_k$, Z^0 , parameters $X, n, X_2 : n$

return G

Appendix B

APPENDIX FOR CHAPTER 3

B.1 Preliminary Results and Proofs in Section 3.2

We first make a remark on Assumption 1 and Assumption 2. Both of the assumptions can be partly motivated by [113]. Their paper considered the stochastic arrival model (without predictions) where the reward functions and the resource consumption functions are all linear. Suppose that, in our problem under the stochastic arrival model and with the linearity assumptions, the prediction is obtained by observing historical arrivals, solving the problem offline, and taking the optimal offline dual variable to be λ . Then by Theorem 1 in [113], we have

$$\lambda_j \leq \lambda_j^* = \frac{r_j}{\log \log \frac{r_j}{\lambda_j}}$$

which provides a guideline to set λ_j in Assumption 1. For example, if $r_j = \frac{p_j}{n}$, that is, we have observed $\frac{p_j}{n}$ historic arrivals, then λ_j can be set as $\lambda_j = \frac{p_j}{n} \cdot \frac{1}{\log \log \frac{p_j}{n}}$. More generally, if $r_j = \frac{p_j}{n}$ for some function p_j such that $p_j \geq 1$ as $n \rightarrow \infty$, then λ_j can be set as $\lambda_j = \frac{p_j}{n} \cdot \frac{1}{\log \log \frac{p_j}{n}}$.

[113] solved their problem by performing dual gradient descent (in their Algorithm 3). More specifically, at each time period they used gradient information to update a hypothetical dual variable λ_j^t . In the proof of Theorem 5 in [113], they showed that the difference of the depletion time of each resource by following λ_j^t and following the optimal dual variable is upper bounded by $\frac{1}{\lambda_j} \log \log \frac{r_j}{\lambda_j}$. This shows if the prediction is perfect in our problem, the λ_j^t obtained by their Algorithm 3 is a particular sequence of dual variables for which Assumption 2 is satisfied.

We state two structural results regarding the duality of the offline problem.

Lemma 14 (Weak Duality). $\text{OPT}^1 W^0 \leq \sum_{j \in J} W_j^0$ for every $\lambda \in \mathbb{R}_+^J$.

Lemma 15 (Duality Gap). $\min_{\lambda \in \mathbb{R}_+^J} \sum_{j \in J} W_j^0 - \text{OPT}^1 W^0 \leq \frac{1}{\lambda} \cdot \frac{1}{\log \log \frac{r_j}{\lambda}}$.

Lemma 14 is the standard weak duality result. Lemma 15 states that, even without any convexity assumptions, the duality gap of our problem is upper bounded by

a constant that is independent from the time horizon. This can be shown via Shapley-Folkman Theorem (see Proposition 5.2.3 for a detailed explanation).

Proof of Lemma 14. This proof appears in [2]. We include it for the sake of completeness. It holds for any $\epsilon > 0$ that

$$\begin{aligned}
 \text{OPT}^1 W^0 &= \max_{G \in \mathcal{X}_C} \sum_{C=1}^n A_C^1 G^0 \\
 &\text{s.t. } \sum_{C=1}^n 6_C^1 G^0 \leq d \\
 &\quad \tilde{0} \\
 &= \max_{G \in \mathcal{X}_C} \sum_{C=1}^n A_C^1 G^0, \sum_{C=1}^n 6_C^1 G^0 \leq d \\
 &= \sum_{C=1}^n A_C^1 \tilde{0} \\
 &= \sum_{C=1}^n W^0
 \end{aligned}$$

where the first inequality is because we relax the constraint $\sum_{C=1}^n 6_C^1 G^0 \leq d$ and $\tilde{0} \geq 0$, and the last equality utilizes the definition of A^0 .

Proof of Proposition 1. Let $\text{conv}^1 \mathcal{X}_C \subseteq \mathbb{R}^3$ denote the convex hull of \mathcal{X}_C . For each C , define the function $A_C: \text{conv}^1 \mathcal{X}_C \rightarrow \mathbb{R}$ by

$$\begin{aligned}
 A_C^1(G^0) &= \sup_{\substack{U^i \in \mathcal{U}^i \\ U^i = 1 - \sum_{j=1}^n U^j \geq 0}} \sum_{i=1}^3 U^i A_C^i(G^0) \\
 &= \sum_{i=1}^3 U^i G^i - G^2 \in \mathcal{X}_C
 \end{aligned}$$

A_C is concave regardless of whether \mathcal{X}_C is concave or not, and it can be viewed as a concavification of A_C on $\text{conv}^1 \mathcal{X}_C$. Similarly, for each C define the function $6_C: \text{conv}^1 \mathcal{X}_C \rightarrow \mathbb{R}$ by

$$\begin{aligned}
 6_C^1(G^0) &= \inf_{\substack{U^i \in \mathcal{U}^i \\ U^i = 1 - \sum_{j=1}^n U^j \geq 0}} \sum_{i=1}^3 U^i 6_C^i(G^0) \\
 &= \sum_{i=1}^3 U^i G^i - G^2 \in \mathcal{X}_C
 \end{aligned}$$

6_C is convex regardless of whether \mathcal{X}_C is convex or not.

Let P^0 denote the optimization problem in Equation (3.1). Consider the following convex relaxation \tilde{P}^0 of the optimization problem in Equation (3.1):

$$\max_{\substack{\tilde{O} \\ G \in \text{conv}^1 X^0}} \sum_{C=1}^n A_C^1 \tilde{O}_C \quad \text{s.t.} \quad \sum_{C=1}^n \tilde{O}_C \leq d;$$

and its Lagrangian dual problem D^0 :

$$\min_{\lambda \in \mathbb{R}^+} \sum_{C=1}^n A_C^1 \tilde{O}_C, \lambda \geq d \quad \text{where} \quad A_C^1 \lambda = \sup_{G \in \text{conv}^1 X^0} f(A_C^1 G^0 - \lambda \sum_{C=1}^n G_C^0)$$

Because $G \in \text{conv}^1 X^0$, $G_C^0 = 0$ for all C and $\lambda \geq 0$, P^0 satisfies Slater's condition. Therefore by strong duality $\sup P^0 = \inf D^0$. By an application of the Sharpley-Folkman Theorem ([4], Proposition 5.26), there exists an optimal solution $(\tilde{O}_C)_{C=1}^n$ of \tilde{P}^0 with the following property: let $J \subseteq \{1, \dots, n\}$ be the set of time periods where $\tilde{O}_C > 0$ for $C \in J$, then $|J| \leq 1$ and $\sum_{C \in J} \tilde{O}_C = d$. Let $\tilde{\lambda}$ be the optimal dual variable of D^0 that induces $(\tilde{O}_C)_{C=1}^n$ and let $(\tilde{g}_C)_{C=1}^n$ be the actions induced by $\tilde{\lambda}$ in the original primal P^0 . We prove that

$$\sum_{C=1}^n A_C^1 \tilde{g}_C^0 = \sum_{C=1}^n A_C^1 \tilde{O}_C \leq d \leq \sum_{C=1}^n A_C^1 \tilde{O}_C$$

Let $J \subseteq \{1, \dots, n\}$ be the set of time periods such that $\tilde{O}_C > 0$ for $C \in J$, and let $I = \{1, \dots, n\} \setminus J$. Then I is the set of time periods where the resource constraint becomes active when choosing the action induced by $\tilde{\lambda}$. Because $|J| \leq 1$ and $\tilde{O}_C = 0$ for $C \in I$, $J \cap I = \emptyset$

$$\sum_{C \in I} \tilde{O}_C = \sum_{C \in I} \tilde{O}_C \leq d \leq \sum_{C \in I} \tilde{O}_C$$

Therefore

$$|J \cap I| = 0 \text{ for } |J| \leq 1, \sum_{C \in I} \tilde{O}_C = d$$

so

$$\sum_{C \in I} \tilde{O}_C = \sum_{C \in I} \tilde{O}_C \leq d \leq \sum_{C \in I} \tilde{O}_C$$

Further, we also have

$$\sum_{C \in I} \tilde{O}_C = \sum_{C \in I} \tilde{O}_C \leq d \leq \sum_{C \in I} \tilde{O}_C$$

and $\tilde{O}_C = 0$ for $C \in I \setminus J$. These together gives

$$\sum_{C=1}^n A_C^1 \tilde{g}_C^0 = \sum_{C=1}^n A_C^1 \tilde{O}_C \leq d \leq \sum_{C=1}^n A_C^1 \tilde{O}_C$$

Finally, since \mathbb{P}^0 is a relaxation of $\mathbb{1}P^0$,

$$\sup_{C=1}^{\tilde{O}} A_C^1 G^0 \text{OPT}^1 W^0 -$$

so we have

$$\max_{2R^c} \mathbb{1}GRD \cdot j W^0 \stackrel{\tilde{O}}{=} A_C^1 G^0 \text{OPT}^1 W^0 \cdot 16 \cdot \underline{\epsilon}, 1^{01} < \cdot, 1^0 A \bullet$$

”

Proof of Proposition 2. Let $G_C^{01} = \arg \max_{G \geq 2f_0 - 1} \{A_C^0 G^0 \cdot \lambda > 6_C^1 G^0\}$. Note that $G_C^{01} = 0$ if and only if $A_C^1 0^0 \cdot \lambda > 6_C^1 0^0$; $A_C^1 1^0, n_C \cdot \lambda > 6_C^1 1^0$, which after rearranging gives $\lambda > 6_C^1 1^0$; $A_C^1 0^0, n_C$. Then for any C such that $\|j\| \cdot \lambda \cdot \|dj\| \cdot Z$, we have

$$\begin{aligned} P^1(G_C^{01} \cdot C = 0) \cdot j G_C^{01} \cdot \lambda = 0^0 &= P \cdot \lambda > 6_C^1 1^0 ; A_C^1 1^0 \cdot A_C^1 0^0, n_C \\ &= P \cdot \lambda > 6_C^1 1^0 ; A_C^1 1^0 \cdot A_C^1 0^0, n_C \\ &= \frac{1 - \Phi\left(\frac{\lambda - 6_C^1 1^0 \cdot f^0}{\lambda > 6_C^1 1^0 \cdot f^0}\right)}{1 - \frac{Z < k \cdot 6_C^1 1^0 \cdot k_1}{\lambda > 6_C^1 1^0}} \end{aligned}$$

where Φ is the cumulative distribution function of the standard normal distribution, and the second equality follows since $\epsilon \sim N(0, \sigma^2)$. The case where $G_C^{01} = 1$ can be handled similarly. Therefore the probability of following λ and following λ consume different amount of resources is Z^0 , so by independency Assumption 2 is satisfied with probability $1 - \exp(-2Z^0) = 1 - \exp(-2Z)$ for some $2 > 0$. ”

We make the following observation, which follows since Proposition 1 shows there exists a perfect dual variable for every arrival sequence.

Observation 6. If an arrival sequence W is $\mathbb{1} - X^0$ -stationary, then

$$\min_{C=1}^{\tilde{O}} \mathbb{1}OPT^1 W_{(X)}^0, \text{OPT}^1 W_{(X), (1)}^0 \leq \mathbb{1} \text{OPT}^1 W^0 \bullet$$

Proof of Proposition 3. Fix any $X > 0$. If $\mathbb{1}GRD \cdot j W^0 = > 1$, then since the total amount of available resources scales linearly in λ and every single action consumes constants amount of resources, the Dual-Adjusted Greedy Algorithm with dual variable λ never depletes resources. Therefore

$$\mathbb{1}GRD \cdot j W_{(X)}^0, \mathbb{1}GRD \cdot j W_{(X), (1)}^0 = \mathbb{1}GRD \cdot j W^0$$

for every $1 \leq j \leq n$, which shows W_j is stationary for every $j \geq 0$. From now on we assume that $\mathbb{P}(\text{GRD}_j W^0 = 1) \leq \epsilon$.

For any time period $t \in [0, T]$ and any amount of resources $d_t \in \mathbb{R}_+^n$, we use GRD_j to denote the amount of reward obtained by the Dual-Adjusted Greedy Algorithm with dual variable λ on the subproblem with available amount of resources d_t .

Fix an integer $1 \leq j \leq n$. If $\mathbb{P}(W_j \geq 1) \leq \epsilon$, then

$$\begin{aligned} \mathbb{P}(\text{GRD}_j W_{(X),1} \geq 1) &\leq \mathbb{P}(\text{GRD}_j W_{(X),1} - d) \geq 1) \leq 6A \cdot \epsilon \\ \mathbb{P}(\text{GRD}_j W^0 \geq 1) &\leq \mathbb{P}(\text{GRD}_j W^0 \geq 1) \leq 6A \cdot \epsilon \end{aligned}$$

where the first inequality follows since any algorithm can consume at most amount of resources in ℓ_1 -norm in a single time period, which can be translated to at most $6A \cdot \epsilon$ amount of reward; the second inequality follows since any algorithm can obtain at most $(X)A$ amount of rewards in the first $(X)A$ time periods. Since $(X) \geq 2$, this shows

$$\mathbb{P}(\text{GRD}_j W_{(X),1} \geq 1) \leq \mathbb{P}(\text{GRD}_j W_{(X),1} \geq 1) \leq \mathbb{P}(\text{GRD}_j W^0 \geq 1) \leq \epsilon.$$

Similarly, we can also show this if $\mathbb{P}(W_j \geq 1) \leq \epsilon$.

Suppose $\mathbb{P}(W_j \geq 1) \leq \epsilon$. Let d_t be the amount of resources that is consumed by the Dual-Adjusted Greedy Algorithm with dual variable λ and arrivals W_t , then

$$\mathbb{P}(\text{GRD}_j W^0 = 1) = \mathbb{P}(\text{GRD}_j W^0 = 1)$$

and

$$\mathbb{P}(\text{GRD}_j W_{(X),1} \geq 1) = \mathbb{P}(\text{GRD}_j W_{(X),1} \geq 1) \leq \epsilon.$$

Condition on W^0 , for each time period $t \in [0, T]$ let $c_t \in \mathbb{R}_+^n$ be the amount of resources consumed at time period t . Then c_t 's are independent with $\mathbb{P}(c_t = 0) = \epsilon$ and $0 \leq c_{t,j} \leq 1$. By Hoeffding's inequality we have

$$\mathbb{P}_{W^0} \left(\sum_{C=1}^{\tilde{C}} 1 - c_{C,j} \geq (X) d_{C,j} \right) \leq 6 \cdot 2^{\tilde{C}} \frac{1}{(X) \log(X)} \leq \frac{1}{2}$$

where $1 - c_{C,j}$ denotes the j th coordinate of c_C . So by union bound we have

$$(B.1) \quad \mathbb{P}_{W^0} \left(\sum_{C=1}^{\tilde{C}} c_{C,j} \geq (X) d_{C,j} \right) \leq 6 \cdot 2^{\tilde{C}} \frac{1}{(X) \log(X)} \leq \frac{1}{2}.$$

Therefore

$$\begin{aligned}
 & \mathbb{E}_{W \sim P} \left[\sum_{j=1}^h \text{GRD} \cdot j W_{j, :X} - \sum_{c=1}^C \mathbb{E}_{W \sim P} \left[\sum_{j=1}^h \text{GRD} \cdot j W_{j, :X} - \sum_{c=1}^C \text{GRD} \cdot j W_{j, :X} \right] \right] \\
 &= \mathbb{E}_{W \sim P} \left[\sum_{j=1}^h \text{GRD} \cdot j W_{j, :X} - \sum_{c=1}^C \mathbb{E}_{W \sim P} \left[\sum_{j=1}^h \text{GRD} \cdot j W_{j, :X} - \sum_{c=1}^C \text{GRD} \cdot j W_{j, :X} \right] \right] \\
 &\leq \frac{1}{2} \mathbb{E}_{W \sim P} \left[\sum_{j=1}^h \text{GRD} \cdot j W_{j, :X} - \sum_{c=1}^C \text{GRD} \cdot j W_{j, :X} \right] + 6 \sqrt{2^p \sum_{j=1}^h \text{GRD} \cdot j W_{j, :X} \log \sum_{j=1}^h \text{GRD} \cdot j W_{j, :X}} \\
 &\leq \frac{1}{2} \mathbb{E}_{W \sim P} \left[\sum_{j=1}^h \text{GRD} \cdot j W_{j, :X} - \sum_{c=1}^C \text{GRD} \cdot j W_{j, :X} \right] + \frac{1}{2} \mathbb{E}_{W \sim P} \left[\sum_{j=1}^h \text{GRD} \cdot j W_{j, :X} - \sum_{c=1}^C \text{GRD} \cdot j W_{j, :X} \right] \\
 &\leq \frac{1}{2} \mathbb{E}_{W \sim P} \left[\sum_{j=1}^h \text{GRD} \cdot j W_{j, :X} - \sum_{c=1}^C \text{GRD} \cdot j W_{j, :X} \right] + \frac{1}{2} \mathbb{E}_{W \sim P} \left[\sum_{j=1}^h \text{GRD} \cdot j W_{j, :X} - \sum_{c=1}^C \text{GRD} \cdot j W_{j, :X} \right] \\
 &= \mathbb{E}_{W \sim P} \left[\sum_{j=1}^h \text{GRD} \cdot j W_{j, :X} - \sum_{c=1}^C \text{GRD} \cdot j W_{j, :X} \right]
 \end{aligned}$$

where the first inequality follows by conditioning on two cases given by Eq. (B.1) and noticing that any algorithm can obtain at most $6A$ amount of rewards in the first $\sum_{j=1}^h \text{GRD} \cdot j W_{j, :X}$ time periods, and the second inequality follows since any algorithm can consume at most 6 amount of resources in ℓ_1 -norm in a single time period, which can be translated to at most $6A$ amount of reward. Similarly, we also have

$$\mathbb{E}_{W \sim P} \left[\sum_{j=1}^h \text{GRD} \cdot j W_{j, :X} - \sum_{c=1}^C \mathbb{E}_{W \sim P} \left[\sum_{j=1}^h \text{GRD} \cdot j W_{j, :X} - \sum_{c=1}^C \text{GRD} \cdot j W_{j, :X} \right] \right] = \mathbb{E}_{W \sim P} \left[\sum_{j=1}^h \text{GRD} \cdot j W_{j, :X} - \sum_{c=1}^C \text{GRD} \cdot j W_{j, :X} \right]$$

Hence

$$\mathbb{E}_{W \sim P} \left[\sum_{j=1}^h \text{GRD} \cdot j W_{j, :X} - \sum_{c=1}^C \mathbb{E}_{W \sim P} \left[\sum_{j=1}^h \text{GRD} \cdot j W_{j, :X} - \sum_{c=1}^C \text{GRD} \cdot j W_{j, :X} \right] \right] = \mathbb{E}_{W \sim P} \left[\sum_{j=1}^h \text{GRD} \cdot j W_{j, :X} - \sum_{c=1}^C \text{GRD} \cdot j W_{j, :X} \right]$$

Note that $\sum_{j=1}^h \text{GRD} \cdot j W_{j, :X} - \sum_{c=1}^C \mathbb{E}_{W \sim P} \left[\sum_{j=1}^h \text{GRD} \cdot j W_{j, :X} - \sum_{c=1}^C \text{GRD} \cdot j W_{j, :X} \right]$ is a function from \mathcal{P}^h to \mathbb{R} such that each W_j is drawn independently. Moreover, it satisfies the bounded differences property with bound $6A$ since any algorithm can consume at most 6 amount of resources in ℓ_1 -norm in a single time period, which can be translated to at most $6A$ amount of reward. Therefore by McDiarmid's inequality we have

$$\begin{aligned}
 & \mathbb{P}_{W \sim P} \left[\sum_{j=1}^h \text{GRD} \cdot j W_{j, :X} - \sum_{c=1}^C \mathbb{E}_{W \sim P} \left[\sum_{j=1}^h \text{GRD} \cdot j W_{j, :X} - \sum_{c=1}^C \text{GRD} \cdot j W_{j, :X} \right] \geq \frac{1}{3} \right] \\
 & \leq \frac{1}{3} \sum_{j=1}^h \text{GRD} \cdot j W_{j, :X} + \frac{1}{3} \sum_{c=1}^C \mathbb{E}_{W \sim P} \left[\sum_{j=1}^h \text{GRD} \cdot j W_{j, :X} - \sum_{c=1}^C \text{GRD} \cdot j W_{j, :X} \right]
 \end{aligned}$$

This implies $\sum_{j=1}^h \text{GRD} \cdot j W_{j, :X} - \sum_{c=1}^C \mathbb{E}_{W \sim P} \left[\sum_{j=1}^h \text{GRD} \cdot j W_{j, :X} - \sum_{c=1}^C \text{GRD} \cdot j W_{j, :X} \right] \geq \frac{1}{3}$ with probability at least $\frac{1}{3}$ for any $\sum_{j=1}^h \text{GRD} \cdot j W_{j, :X}$. Since $\sum_{j=1}^h \text{GRD} \cdot j W_{j, :X}$ can take at most $\sum_{j=1}^h \text{GRD} \cdot j W_{j, :X}$ values, by union bound

$$\min_{\sum_{j=1}^h \text{GRD} \cdot j W_{j, :X}} \left[\sum_{j=1}^h \text{GRD} \cdot j W_{j, :X} - \sum_{c=1}^C \mathbb{E}_{W \sim P} \left[\sum_{j=1}^h \text{GRD} \cdot j W_{j, :X} - \sum_{c=1}^C \text{GRD} \cdot j W_{j, :X} \right] \geq \frac{1}{3} \right]$$

with probability at least $\frac{1}{3}$. Because $\sum_{j=1}^h \text{GRD} \cdot j W_{j, :X} = \sum_{j=1}^h \text{GRD} \cdot j W_{j, :X}$, this implies $\sum_{j=1}^h \text{GRD} \cdot j W_{j, :X}$ is $\frac{1}{3}$ -stationary for every $\sum_{j=1}^h \text{GRD} \cdot j W_{j, :X} \geq \frac{1}{3}$.

B.2 Details in Section 3.3.1

For completeness, we discuss the Mirror Descent Algorithm (MDA) given in [22].

Algorithm 17: Mirror Descent Algorithm (MDA)

Inputs: Initial dual solution λ^0 , total time period T , initial resources $\mathbf{1} = \mathbf{d}$, reference function $\phi^0 : \mathbb{R}^n \rightarrow \mathbb{R}$, and step-size γ

for C from 1 to T do

 Receive requests $\mathbf{A}_C = \mathbf{G}_C \mathbf{X}_C^0$

 Make the primal decision \mathbf{c}_C and update the remaining resources

$$\mathbf{c}_C = \arg \max_{\mathbf{c} \in \mathcal{C}} \sum_{j=1}^n \mathbf{c}_j \mathbf{G}_j^0 \mathbf{X}_j^0 - \sum_{j=1}^n \mathbf{c}_j \mathbf{G}_j^0 \mathbf{1}$$

 Obtain a sub-gradient of the dual function:

$$\mathbf{q}_C = \sum_{j=1}^n \mathbf{c}_j \mathbf{G}_j^0, \mathbf{d}$$

 Update the dual variable by mirror descent:

$$(B.2) \quad \lambda_{C,1} = \arg \min_{\lambda \in \mathbb{R}^n} \langle \mathbf{q}_C, \lambda \rangle + \frac{1}{\gamma} \phi^0(\lambda)$$

 where $\phi^0(\lambda) = \phi^0(\lambda) - \langle \mathbf{1}, \lambda \rangle$ is the Bregman divergence.

end

The Mirror Descent Algorithm takes an initial dual variable, a step-size, and a reference function as inputs. At each time period C , the algorithm takes the action induced by the current dual variable $\lambda_{C,1}$ and performs a first-order update on the dual variable. For the updating step, note we can write the dual function in Equation (3.3) as $\phi^0(\lambda) := \sum_{j=1}^n \lambda_j W_j^0$ where the j th term of the dual function is given by $\lambda_j W_j^0 = \sum_{C=1}^T \lambda_j \mathbf{c}_C^T \mathbf{G}_j^0 \mathbf{X}_j^0 - \lambda_j \mathbf{G}_j^0 \mathbf{1}$. Then it follows that $\mathbf{q}_C := \sum_{j=1}^n \mathbf{c}_j \mathbf{G}_j^0, \mathbf{d}$ is a sub-gradient of ϕ^0 at λ_C under our assumptions by Danskin's Theorem (see, e.g., Proposition B.25 in [3]), and the algorithm uses \mathbf{q}_C to update the dual variable by performing a mirror descent step in Equation (B.2) with step-size γ and reference function ϕ^0 . Intuitively, the Mirror Descent Algorithm tries to find dual variables via gradient information such that these dual variables induce actions with good primal performances. For more on mirror descent algorithms in general, see [30, 78, 120, 138].

We state the standard assumptions on choosing the reference function for mirror descent algorithms [30, 39, 119, 120]. These assumptions are applicable to all algorithms in our paper.

- (a) ϕ^0 is either differentiable or essentially smooth [27] and Lipschitz in \mathbb{R}^n

(b) f^0 is f -strongly convex with respect to the l_1 -norm in \mathbb{R}^d , i.e.,

$$f^0(\mathbf{x}_1) - f^0(\mathbf{x}_2) \geq \langle \mathbf{g}, \mathbf{x}_1 - \mathbf{x}_2 \rangle - \frac{f}{2} \|\mathbf{x}_1 - \mathbf{x}_2\|_1^2$$

for all $\mathbf{x}_1, \mathbf{x}_2 \in \mathbb{R}^d$.

(c) f^0 is coordinately-wise separable, i.e., $f^0(\mathbf{x}) = \sum_{j=1}^d g_j(x_j)$ where $g_j: \mathbb{R} \rightarrow \mathbb{R}$ is an univariate function. Moreover, for every resource j , the function g_j is f_j -strongly convex with respect to the l_1 -norm over $[-\frac{1}{2}, \frac{1}{2}]$ where $d_j := A \cdot d_j \leq 1$.

B.3 Proofs in Section 3.3.3

Proof of Proposition 7. Let 2 be a positive integer such that $2 \geq \max_{j \in [d]} \frac{1}{d_j} g_j$. Set $d = 1$, $U = 1 \cdot 2X$, and $A = 1 \cdot U \cdot 1 \cdot 1 - 1 \cdot U \cdot 0$, then by our choice of 2 we have $\frac{1}{2} \leq 1 \cdot U$ and $A \leq U$. Consider two different types of arrivals $W^1 = 1 \cdot A \cdot 1 - 0 \cdot X^0$ and $W^2 = 1 \cdot A \cdot 0 - X^0$, where $X = f \cdot 0 - 1 \cdot g$ (one can think of this as {reject, accept}). Set $A^{11} = 1 - 0 \cdot 1 \cdot 1 = 1 - A$, and $A^{21} = U$. Let $\lambda = 1 \cdot 1 \cdot \log(1)$ be the prediction, then following λ means taking action 0 for W^1 and taking action 1 for W^2 . Because $A \cdot U \leq \lambda$, one can verify that Assumptions 1 and 2 are satisfied.

Consider the following two instances.

- ^ Instance one: the arrivals are stochastic where the state space is \mathbb{W} , i.e., $W_C = W^1$ for every $C = 1 \cdot \dots \cdot 1$. In this instance the optimum is to take action 1 for all arrivals. Note that following λ would take action 0 for all arrivals, which means the prediction has bad quality.
- ^ Instance two: the arrivals are adversarial where $W_C = W^1$ for $C = 1 \cdot \dots \cdot \frac{U-1}{U}$ and $W_C = W^2$ for $C = \frac{U-1}{U}, 1 \cdot \dots \cdot 1$. In this instance the optimum is to take action 0 for W^1 and take action 1 for W^2 . We have $PRD^1 W^0 = OPT^1 W^0 \cdot \frac{A}{U}$, which means the prediction is perfect. Moreover, since $1 \cdot 2U$, one can verify that

$$\begin{aligned} & \min_{\lambda \in [0,1]} OPT^1 W_{(1 \cdot \dots \cdot 1)}^0, OPT^1 W_{(\frac{U-1}{U}, 1)}^0 \\ &= OPT^1 W_{(1 \cdot \dots \cdot \frac{U-1}{U})}^0, OPT^1 W_{(\frac{U-1}{U}, 1)}^0 \\ &= \left(\frac{U-1}{U}, \frac{A}{U^2} \right) \cdot \end{aligned}$$

Then we get

$$\begin{aligned}
 & 1 - \frac{OPT^1(W_{(1, \frac{U-1}{U})})}{OPT^1(W_{(\frac{U-1}{U}, 1)})} \leq OPT^1(W^0) \\
 & = 1 - \frac{U-1}{U} \cdot \frac{A}{U^2} \leq \frac{A}{U} \\
 & = 1 - \frac{1}{U} \cdot \frac{U-1}{A} \\
 & = \frac{1}{U}
 \end{aligned}$$

where the last equality follows by our choice of $A = U - 1$. Therefore W is $\frac{1}{U}$ -nonstationary with respect to X .

Note that no algorithm can distinguish instance one and instance two before time period $C = \frac{U-1}{U}$. For any algorithm, assume in instance one it satisfies $\text{Regret}^1(\text{ALG}) \leq \epsilon$, then since the optimum is to take action 1 for all arrivals, at time period $C = \frac{U-1}{U}$ the amount of resources left is at most ϵ . Therefore in instance two the algorithm can take action 1 for at most $\frac{\epsilon}{A}$ time periods, so the total rewards gained in instance two satisfies $\text{ALG}_j(W^0) \leq \frac{U-1}{U} + \frac{\epsilon}{A}$. Because $\text{PRD}^1(W^0) = \frac{A}{U}$, in instance two we have

$$\begin{aligned}
 & \limsup_{j \rightarrow \infty} \frac{1}{j} \sum_{i=1}^j \max \left\{ \frac{1}{U} \text{OPT}^1(W^0) - \text{PRD}^1(W^0), \text{ALG}_j(W^0) \right\} \\
 & = \limsup_{j \rightarrow \infty} \frac{1}{j} \sum_{i=1}^j \frac{A}{U} \\
 & = \frac{A}{U} \\
 & = \frac{1}{U}
 \end{aligned}$$

where the last inequality follows since $\frac{1}{j} \sum_{i=1}^j \frac{A}{U} \leq \frac{A}{U}$ and $X = U - 1$.

”

B.4 Description of Algorithm 18

We list some notations used in Algorithm 18, which follow notations in [48]. Given initial dual solution λ_1 and step-size η :

- (a) Let $c_t^1 \in [0, 1]$ and $\lambda_t^1 \in [0, 1]$ be the action we take and the dual variable we get after $C = 1$ iterations of the Mirror Descent Algorithm with initial dual solution λ_1 , step-size η , and the same requests as the requests that Algorithm 18 received so far;

- (b) Define $\lambda_{C^1} := \max_{B^1} \|j_{C^1} - j_{B^1}\|_1$ to be the maximum l_1 -distance from any updated dual variable used in the Mirror Descent Algorithm before time C to the initial dual variable;
- (c) Define $\lambda_{C^1}^2 := \sum_{B=1}^C \|j_{B^1} - j_{C^1}\|_1^2$ to be the running sum of squared l_1 -norms of the dual functions' sub-gradients.

A high-level intuition behind the choices of step sizes is that, it is well-known [42] that the hindsight (asymptotically) optimal step size is λ that satisfies

$$\lambda = \frac{\sum_{B=1}^C \|j_{B^1} - j_{C^1}\|_1}{\sum_{B=1}^C \|j_{B^1} - j_{C^1}\|_1^2}.$$

Because $\sum_{B=1}^C \|j_{B^1} - j_{C^1}\|_1$ and $\sum_{B=1}^C \|j_{B^1} - j_{C^1}\|_1^2$ are unknown a priori, at each time period C we use λ_{C^1} as an approximation of $\sum_{B=1}^C \|j_{B^1} - j_{C^1}\|_1$ and use $\lambda_{C^1}^2$ as an approximation of $\sum_{B=1}^C \|j_{B^1} - j_{C^1}\|_1^2$, and these approximations can be proven to be accurate. Then we use bisection to find an approximate solution of the implicit function

$$\lambda_C = \frac{\lambda_{C^1}}{U \lambda_{C^1}^2 + V}$$

where U, V are damping parameters. For a more detailed explanation, see Note 4.3. Note that the Stochastic Arrival Algorithm does not need to know the accuracy parameter ϵ .

B.5 Proofs in Section 3.4.1 and 3.4.2

Proof of Proposition 8. The proof technique is similar to the proof of Theorem 1 in [22], which we largely borrow. We break down the proof in three steps.

Step 1 (Primal performance.) First, we define the stopping time of Algorithm 18 as the first time less than T that there exists resources such that $\sum_{C=1}^g \mathbb{1}_{\{C \leq T\}} \mathbb{1}_{\{d_C \geq \epsilon\}}$. Notice that g is a random variable, and moreover, we will not violate the resource constraints before the stopping time. We here study the primal-dual gap until the stopping time. Notice that before the stopping time, Algorithm 18 performs the mirror descent steps on the dual function with ϵ -tuned step sizes.

Consider a time $C \leq g$ so that actions are not constrained by resources. Then the algorithm takes the action $C^1 \in \arg \max_{C^1 \in \mathcal{C}^1} \langle C^1, G^0 \rangle - \frac{1}{\epsilon} \langle C^1, G^0 \rangle$, so we have that

$$\langle C^1, G^0 \rangle = \langle C^1, G^0 \rangle - \frac{1}{\epsilon} \langle C^1, G^0 \rangle.$$

Algorithm 18: Stochastic Arrival Algorithm (SA)

Inputs: Prediction λ , total time period \mathcal{T} , initial resources $\mathbf{x}_1 = \mathbf{d}$, reference function $\mathbf{H}^0 : \mathbb{R}^C \rightarrow \mathbb{R}$, and initial step-size μ_1

Initialize $\lambda_1 = \lambda$

for C from 1 to C do

 Receive request $\mathbf{x}_C = \mathbf{c}$

 Make the primal decision \mathbf{c} and update the remaining resources

$$\mathbf{c}_C = \arg \max_{\mathbf{c} \in \mathcal{C}} \mathbf{c}^T \mathbf{G}^0 - \mu_C \sum_{c=1}^C \mathbf{c}^T \mathbf{G}^0$$

 Obtain a sub-gradient of the dual function:

$$\mathbf{q}_C = \mathbf{c}^T \mathbf{G}^0, \mathbf{d}$$

 Update the dual variable by mirror descent:

$$\lambda_{C,1} = \arg \min_{\lambda \in \mathbb{R}^C} \langle \mathbf{q}_C, \lambda \rangle + \frac{1}{\mu_C} \mathcal{D}(\lambda | \lambda_{C-1})$$

 where $\mathcal{D}(\lambda | \lambda^0) := \mathbf{H}^0(\lambda) - \mathbf{H}^0(\lambda^0) - \langle \mathbf{H}^0(\lambda^0), \lambda - \lambda^0 \rangle$ is the Bregman divergence.

 Tune the step size:

 for $\tau = 2, 4, 8, 16, \dots$ do

$C = bC + \tau$

$\mathbf{U}^{\tau,0} = \frac{1}{\tau} \sum_{c=1}^{\tau} \mathbf{c}^T \mathbf{G}^0, \mathbf{V}^{\tau,0} = \frac{1}{\tau} \sum_{c=1}^{\tau} \mathbf{c}^T \mathbf{G}^0, \mathbf{d}^{\tau,0} = \mathbf{d}$ where

$$\mathbf{c}^{\tau,0} := \frac{1}{\tau} \sum_{c=1}^{\tau} \mathbf{c}^T \mathbf{G}^0$$

 if Root Finding Bisection $\mathcal{R}(\mathbf{c}^{\tau,0} | \mathbf{c}^{\tau,0} - \mathbf{U}^{\tau,0} - \mathbf{V}^{\tau,0}) \neq 1$ then

$\mathcal{R}(\mathbf{c}^{\tau,0} | \mathbf{c}^{\tau,0} - \mathbf{U}^{\tau,0} - \mathbf{V}^{\tau,0})$

 end

 end

end

Function Root Finding Bisection($\mathbf{c}^{\tau,0} | \mathbf{c}^{\tau,0} - \mathbf{U}^{\tau,0} - \mathbf{V}^{\tau,0}$):

$k^{\tau,0} := \lfloor \frac{\mathbf{c}^{\tau,0} - \mathbf{U}^{\tau,0} - \mathbf{V}^{\tau,0}}{\mathbf{c}^{\tau,0}} \rfloor$

 if $\mathbf{c}^{\tau,0} \leq k^{\tau,0} \mathbf{c}^{\tau,0}$ then return 1

 if $\mathbf{c}^{\tau,0} \geq k^{\tau,0} \mathbf{c}^{\tau,0}$ then return $\lfloor k^{\tau,0} \rfloor$

 while $\mathbf{c}^{\tau,0} \geq 2 \lfloor k^{\tau,0} \rfloor$ do

$\mathbf{c}^{\tau,0} = \lfloor k^{\tau,0} \rfloor$

 if $\mathbf{c}^{\tau,0} \leq k^{\tau,0} \mathbf{c}^{\tau,0}$ then $\mathbf{c}^{\tau,0} = \mathbf{c}^{\tau,0}$ else $\mathbf{c}^{\tau,0} = \mathbf{c}^{\tau,0} - \mathbf{c}^{\tau,0}$

 end

 if $\mathbf{c}^{\tau,0} \leq \mathbf{c}^{\tau,0} \lfloor k^{\tau,0} \rfloor$ then return $\mathbf{c}^{\tau,0}$ else return $\lfloor k^{\tau,0} \rfloor$.

End Function

Let

$$j P^0 = \frac{1}{g} E_{W|P} \gg j W^0 = E_{A-g-X^0|P} A^1 c^0 \gg d$$

be the expected dual objective when requests are drawn i.i.d. from S^0 .

Let $b_C = f(W_C, \dots, W_C)$ and \mathcal{F}_C be the sigma-algebra generated by adding the last two equations and taking expectations conditional on \mathcal{F}_C we obtain, because $c^2 f^1 b_C$ and $A_C - X_C^0 \in P$, that

$$\begin{aligned} E_{A_C^1} j f^1 b_C &= E_{A-g-X^0|P} E_{A_C^1} j f^1 b_C \\ &= E_{A_C^1} j f^1 b_C \\ (B.3) \quad &= j P^0 E_{A_C^1} j f^1 b_C - \end{aligned}$$

where the second equality follows the definition of the dual function.

Consider the process

$$/C = \sum_{B=1}^C \sum_{B=1}^C E_{A_C^1} j f^1 b_B -$$

which is martingale with respect to \mathcal{F}_C (i.e., $E_{A_C^1} j f^1 b_C = /C$). Since g is a stopping time with respect to \mathcal{F}_C and g is bounded, the Optional Stopping Theorem implies that $E_{g} = 0$. Therefore,

$$E_{g} \sum_{C=1}^g \sum_{C=1}^C E_{A_C^1} j f^1 b_C = E_{g} \sum_{C=1}^g \sum_{C=1}^C E_{A_C^1} j f^1 b_C \cdot$$

Using a similar martingale argument for A_C^1 and summing Eq. (B.3) from $C = 1$ to g we obtain that

$$\begin{aligned} E_{g} \sum_{C=1}^g A_C^1 &= E_{g} j P^0 E_{g} \sum_{C=1}^g \sum_{C=1}^C E_{A_C^1} j f^1 b_C \\ (B.4) \quad E_{g} \sum_{C=1}^g A_C^1 &\leq E_{g} \sum_{C=1}^g \sum_{C=1}^C E_{A_C^1} j f^1 b_C \cdot \end{aligned}$$

where the inequality follows from denoting $g = \frac{1}{g} \sum_{C=1}^g g_C$ to be the average dual variable and using that the dual function is convex.

Step 2 (Complementary slackness). Consider the sequence of functions $F_C^1 = \sum_{C=1}^C E_{A_C^1} j f^1 b_C$, which capture the complementary slackness at C . The subgradients are given by $\nabla F_C^1 = \sum_{C=1}^C E_{A_C^1} j f^1 b_C$, which are bounded as follows $\|\nabla F_C^1\|_1 \leq k \sum_{C=1}^C E_{A_C^1} j f^1 b_C\|_1 \leq k g, d$. Therefore, Algorithm 18 applies online mirror descent to the sequence of functions F_C^1 with the ϵ -tuned step sizes. To analyze the performance, we use the following lemma from [43].

G_C $2 \arg \max_{G_C} A_C^1 G_C^0 \geq 6_C^1 G_C^0$, and thus $A_C^1 G_C^0 - A_C G_C \geq 6_C G_C - 6_C^1 G_C^0$ and $0 = A^1 G^0 - A_C^1 G_C^0 \geq 6_C^1 G_C^0$. Therefore

$$\begin{aligned} U A_C^1 G_C^0 &= A_C^1 G_C^0, \quad U^1 A_C^1 G_C^0 \\ &A_C G_C, \quad \geq 6_C^1 G_C^0 \geq 6_C G_C, \quad U^1 \geq 6_C^1 G_C^0 \\ &= A_C G_C - U \geq 1 d \quad 6_C^1 G_C^0, \quad U \geq d \geq 6_C G_C \\ &A_C G_C - U \geq 1 d \quad 6_C^1 G_C^0 - \end{aligned}$$

where the second inequality is because $U \geq d \geq 6_C G_C = 0$ by our definition of U and the fact that $\tilde{c} = 0$. Summing up over $C = 1, \dots, C$ yields

$$(B.7) \quad U \sum_{C=1}^{\tilde{C}} A_C^1 G_C^0 - \sum_{C=1}^{\tilde{C}} A_C G_C - U \sum_{C=1}^{\tilde{C}} \geq 1 d \quad 6_C^1 G_C^0.$$

Step 2 (Complementary slackness). Denoting, as before, $F_C^1 = \geq 1 d \quad 1 \quad C^1 G_C^0$. As we have seen in the step 2 in the proof of Proposition 8 (the analysis is deterministic in nature), Algorithm 2 performs online mirror descent to the sequence of functions F_C^1 with step size $\epsilon = 2n^1)^0$ where $2 \quad 0$ is an arbitrary scaling constant. By our assumption that the reference function is Lipschitz, there exists a constant $\lambda \geq 0$ such that $\|j^1 - j^0\| \leq \lambda \|j_1 - j_0\|$ for all $j^0, j^1 \in \mathbb{R}^C$. By a standard result on online mirror descent (see, e.g., Appendix G of [22]), we have

$$(B.8) \quad \sum_{C=1}^{\tilde{C}} F_C^1 \leq \sum_{C=1}^{\tilde{C}} F_C^1 + \frac{16 \lambda^2 d^0^2}{2f} g + \frac{1}{\epsilon} \sum_{C=1}^{\tilde{C}} \|j^1 - j^0\|^2 \\ \sum_{C=1}^{\tilde{C}} F_C^1 \leq \frac{2 \cdot 16 \lambda^2 d^0^2}{2f} n^1)^0 + \frac{\lambda^2 (k^1)^0}{2n^1)^0}$$

where the first inequality is the standard online mirror descent result, and the second inequality follows by the step size $\epsilon = 2n^1)^0$ and the fact that $\|j^1 - j^0\| = \|j^1 - j_1 - (j_1 - j^0)\| \leq \|j^1 - j_1\| + \|j_1 - j^0\| = \|j^1 - j_1\| + \|j_1 - j^0\|$.

Step 3 (Putting it all together). We have

$$\begin{aligned} \text{OPT}^1 W^0 - U^1 \text{AA} j W^0 &\leq \sum_{C=1}^g \sum_{c=1}^d A_c^1 G_c^0 - U \sum_{C=1}^g \sum_{c=1}^d A_c^1 G_c^0 \\ &\leq \sum_{C=g,1}^g \sum_{c=1}^d A_c^1 G_c^0 + U \sum_{C=1}^g F_C^1 \cdot c^0 \\ &\leq \sum_{C=1}^g \sum_{c=1}^d A_c^1 G_c^0 + U \sum_{C=1}^g F_C^1 \cdot c^0 \\ &\leq \sum_{C=1}^g \sum_{c=1}^d A_c^1 G_c^0 + U \left(\frac{2^6 d^2}{2f} n^1 \right)^0 + \frac{\wedge!k^1)^0}{2n^1)^0} \cdot \end{aligned}$$

where the first inequality follows because $\sum_{C=1}^g \sum_{c=1}^d A_c^1 G_c^0 \geq 0$, the second inequality is from Eq. (B.7), and the third inequality utilizes $\sum_{C=1}^g \sum_{c=1}^d A_c^1 G_c^0 \leq \sum_{C=1}^g \sum_{c=1}^d A_c^1 G_c^0$ and Eq. (B.8). Similar to the proof of Proposition 8, we note that $\sum_{C=1}^g F_C^1 \cdot c^0 \leq \sum_{C=1}^g F_C^1 \cdot c^0$ for every $c \in \{1, \dots, d\}$ and discuss the choice of c in order to upper bound $\sum_{C=1}^g F_C^1 \cdot c^0$. If $c = 9$, then set $\tilde{c} = 0$, and the result follows. If $c \neq 9$, then there exists a resource $9 \leq k \leq d$ such that $\sum_{C=1}^g \sum_{c=1}^d A_c^1 G_c^0 \geq \sum_{C=1}^g \sum_{c=1}^d A_c^1 G_c^0$. Set $\tilde{c} = A \cdot U d_g - 4_g$ where 4_g is the g -th unit vector and repeat the steps of the stochastic arrivals case to obtain:

$$\text{OPT}^1 W^0 - U^1 \text{AA} j W^0 \leq \frac{A6}{d} \sum_{C=1}^g \sum_{c=1}^d A_c^1 G_c^0 + U \left(\frac{2^6 d^2}{2f} n^1 \right)^0 + \frac{\wedge!k^1)^0}{2n^1)^0} -$$

which finishes the proof by noticing that $\sum_{C=1}^g \sum_{c=1}^d A_c^1 G_c^0$ and $\frac{\wedge!k^1)^0}{2n^1)^0} \cdot n^1)^0$ are both sub-linear in n .

Lemma 18. Consider the Adversarial Arrival Algorithm (AA) under the adversarial arrival model. Given a prediction λ with accuracy parameter θ , it holds that:

$$\limsup_{n \rightarrow \infty} \sup_{W \in \mathcal{S}} \frac{1}{n} \text{PRD}^1 W^0 - U^1 \text{AA} j W^0 \leq \theta$$

Proof of Lemma 18. Recall the updating rule

$$q_{C,1} \geq \arg \min_{2R^c} q_C^{\geq}, \frac{1}{\Gamma} + \lambda^1 - \lambda^c$$

where $q_C = \sum_{c=1}^d A_c^1 G_c^0$, d . Note that $q_C^{\geq}, \frac{1}{\Gamma} + \lambda^1 - \lambda^c$ is convex in λ^c , and set its gradient of λ^c to zero yields

$$q_C, \frac{1}{\Gamma} \lambda^1 - \lambda^c - r^1 - c^0 = 0$$

where ϕ^0 is the reference function. Because

$$\|j_j - j_{j-1}\| \leq \frac{1}{f} \|j_j - j_{j-1}\| \leq \frac{1}{f} \|j_j - j_{j-1}\|$$

and by our assumption ϕ^0 is f -strongly convex with respect to the norm in \mathbb{R}^d , we have

$$\|j_j - j_{j-1}\| \leq \frac{1}{f} \|j_j - j_{j-1}\| \leq \frac{2^j \delta^2 d^0 n^0}{f}$$

Therefore

$$\|j_j - j_{j-1}\| \leq \sum_{B=2}^j \|j_B - j_{B-1}\| \leq \frac{2^j \delta^2 d^0 n^0}{f} C$$

and hence

$$(B.9) \quad \sum_{C=1}^{\tilde{O}} \|j_C - j_{C-1}\| \leq \frac{2^{\tilde{O}} \delta^2 d^0 n^0}{2f} C$$

Let \hat{G}_C be the actions taken by the Prediction Algorithm at time C . Then $\text{PRD}^1 W^0 = \sum_{C=1}^{\tilde{O}} A_C^1 \hat{G}_C^0$. Because $Z_j \geq 0$ is a constant and $n^0 \geq 2^j$,

$$\|j_C - j_{C-1}\| \leq \frac{2^C \delta^2 d^0 n^0}{f} C \leq Z$$

for all $C \leq \tilde{O}$. Therefore $\{j_C\}_{C=1}^{\tilde{O}}$ is a sequence of dual variables that satisfies Assumption 2. Let g_j^0 be the depletion time of resource j of Algorithm 2 and g_j^1 be the depletion time of resource j of Algorithm 1. Then by Assumption 2 we have $g_j^0 \geq g_j^1 \geq \frac{1}{2} g_j^0$ for all resource j . Moreover, since there are resources, outside of all times between each g_j^0 and g_j^1 is partitioned into at most $\frac{1}{\epsilon}$ consecutive time blocks, say $[t_{j,k}^0, t_{j,k}^1]$ for some $\epsilon < \frac{1}{2}$. Note that the set of feasible actions $\{G_j \in \mathcal{G}^0 \mid \sum_{C=1}^k G_C^0 \leq Z\}$ amount of remaining resource at time period C is the same for Algorithm 2 and Algorithm 1 for all $C \in [t_{j,k}^0, t_{j,k}^1]$. Therefore both algorithms perform online mirror descent during time periods $[t_{j,k}^0, t_{j,k}^1]$. Therefore similar to Eq. (B.8) we have

$$(B.10) \quad \sum_{C \in [t_{j,k}^0, t_{j,k}^1]} \|F_C^1 - F_C^0\| \leq \sum_{C \in [t_{j,k}^0, t_{j,k}^1]} \|F_C^1 - F_C^0\| \leq \frac{2^C \delta^2 d^0 n^0}{2f} C + \frac{\epsilon^2 k^2 n^0}{2n^0}$$

for each $j \in [1, d]$. Also, because $G_C^2 \in \arg \max_{G \in \mathcal{G}^0} \langle G, \phi^0 \rangle - \frac{\epsilon}{2} \|G\|^2$, for $C \in [t_{j,k}^0, t_{j,k}^1]$ we have

$$(B.11) \quad \langle G_C^2, \phi^0 \rangle - \frac{\epsilon}{2} \|G_C^2\|^2 \geq \langle G_C^1, \phi^0 \rangle - \frac{\epsilon}{2} \|G_C^1\|^2$$

Because $F_{C^1} \leq C^1 d^{-1} G^{00}$ and $F_{C^1} = \lambda^1 d^{-1} G^{\lambda 00}$, for each ϵ we get

$$\begin{aligned} & \tilde{O}_{C^{2,0}} \left(A_{C^1} G^{\lambda 0} + A_{C^1} G^{00} \right) \\ &= \tilde{O}_{C^{2,0}} \left(C^1 d^{-1} G^{00} + \lambda^1 d^{-1} G^{\lambda 00} \right) \\ & \leq \tilde{O}_{C^{2,0}} \left(C^1 \lambda^{-1} d^{-1} G^{\lambda 00} + C^1 d^{-1} G^{00} \right) \\ & \leq \tilde{O}_{C^{2,0}} \left(F_{C^1} \leq C^1 d^{-1} G^{00} + F_{C^1} = \lambda^1 d^{-1} G^{\lambda 00} \right) \\ & \leq \frac{2^1 6^1 d^{0^2} n^1)^0}{2f} + \frac{\lambda^1 k^1)^0}{2n^1)^0} + \frac{2d^1 6^1 d^0}{2f} + d^{-1} n^1)^0 \bullet \end{aligned}$$

where the first inequality follows from Eq. (B.11), the second inequality is because by Hölder's inequality $\lambda^{-1} d^{-1} G^{\lambda 00} \leq \lambda^{-1} d^{-1} G^{\lambda 00} \leq \lambda^{-1} d^{-1} G^{\lambda 00}$ and $d^{-1} G^{00} \leq d^{-1} G^{00}$, and the third inequality follows from Eq. (B.10) and Eq. (B.9). Therefore

$$\begin{aligned} \text{PRD}^1 W^0 \leq \sum_{j=1}^n \tilde{O}_{C^{2,0}} \left(A_{C^1} G^{\lambda 0} + A_{C^1} G^{00} \right) \\ \leq \sum_{j=1}^n \tilde{O}_{C^{2,0}} \left(C^1 d^{-1} G^{\lambda 00} + C^1 d^{-1} G^{00} \right) \\ \leq \sum_{j=1}^n \left(\frac{2^1 6^1 d^{0^2} n^1)^0}{2f} + \frac{\lambda^1 k^1)^0}{2n^1)^0} + \frac{2d^1 6^1 d^0}{2f} + d^{-1} n^1)^0 \right) \\ \leq \sum_{j=1}^n \left(A_{C^1} G^{\lambda 0} + A_{C^1} G^{00} \right) \\ \Rightarrow \dots \bullet \end{aligned}$$

where the first inequality is because $A_{C^1} G^{\lambda 0} + A_{C^1} G^{00} \leq A$ for each C and the second inequality is by noting that $\frac{2^1 6^1 d^{0^2} n^1)^0}{2f} + \frac{\lambda^1 k^1)^0}{2n^1)^0} + \frac{2d^1 6^1 d^0}{2f} + d^{-1} n^1)^0 \leq A$, and $\sum_{j=1}^n \left(A_{C^1} G^{\lambda 0} + A_{C^1} G^{00} \right) \leq \sum_{j=1}^n A$. This shows

$$\limsup_{\epsilon \rightarrow 0} \sup_{W \in \mathcal{W}} \frac{1}{\epsilon} \text{PRD}^1 W^0 \leq \sum_{j=1}^n A = 0 \bullet$$

”

Combine Lemma 17 and Lemma 18 gives Proposition 9.

”

B.6 Proofs in Section 3.4.3

Proof of Theorem 1. We divide the proof into three cases. The first case is that the underlying arrival model is stochastic and the algorithm never switches to the Adversarial Arrival Algorithm (i.e., the for loop in the algorithm is never broken), and in this case we show that $\text{Regret}^1(\text{MainALG}^0) \leq \frac{1}{2} \max_{j \in \mathcal{W}^0} \sum_{t=1}^n w_t^j - 1$. The second case is that the underlying arrival model is stochastic and yet the algorithm switches to the Adversarial Arrival Algorithm at some point, and we prove that this case happens with low probability. The third case is that the underlying arrival model is adversarial, and in this case we show that

$$\limsup_{n \rightarrow \infty} \sup_{W \in \mathcal{W}^0} \frac{1}{n} \text{Regret}^1(\text{MainALG}^0) \leq \max_{j \in \mathcal{W}^0} \frac{1}{U} \sum_{t=1}^n w_t^j - 1$$

regardless of whether the algorithm switches to the Adversarial Arrival Algorithm or not. To simplify the notation, throughout the proof we will assume n and $1/X$ are integers. The roundings $\lfloor n/X \rfloor$ and $\lfloor 1/X \rfloor$ in our algorithm will not affect the result of our analysis.

Case 1: Suppose the underlying arrival model is stochastic where each arrival W_t is drawn i.i.d. from an underlying probability distribution \mathcal{P}^0 on \mathcal{W}^0 , and the algorithm never switches to the Adversarial Arrival Algorithm. Then the algorithm decomposes n time periods into $\lfloor n/X \rfloor$ time blocks, where each time block contains X time periods and has at least d amount of resources available. During each time block the algorithm performs the Stochastic Arrival Algorithm. Therefore, by our definition of $\text{OPT}_{B^0}^1(W_{(X), 1:1, 1^0 X})$ and the performance guarantee of the Stochastic Arrival Algorithm given by Proposition 8, we have

$$\begin{aligned} \mathbb{E}_{W \in \mathcal{P}^0} \left[\frac{1}{n} \text{Regret}^1(\text{MainALG}^0) \right] &= \mathbb{E}_{W \in \mathcal{P}^0} \left[\frac{1}{n} \sum_{t=1}^n w_t^j - 1 \right] \\ &= \mathbb{E}_{W \in \mathcal{P}^0} \left[\frac{1}{X} \sum_{i=1}^{\lfloor n/X \rfloor} \sum_{t=(i-1)X+1}^{iX} w_t^j - 1 \right] \\ &= \mathbb{E}_{W \in \mathcal{P}^0} \left[\frac{1}{X} \sum_{i=1}^{\lfloor n/X \rfloor} \text{OPT}_{B^0}^1(W_{(X), 1:1, 1^0 X}) - 1 \right] \end{aligned} \tag{B.12}$$

For each time period $t \in [1, n]$, let $c^t \in \mathcal{W}^0 := \mathcal{A}_c^0, c^t \geq d$ be the term of the Lagrangian dual function Equation (3.3), then every $c^t \in \mathcal{W}^0$ is also i.i.d. By Lemma 14 and Lemma 15, for every arrival sequence W we have

$$\text{OPT}_{B^0}^1(W) \leq \min_{C=1}^{\tilde{\mathcal{C}}^B} \sum_{t=1}^n c^t \quad c^t \in \mathcal{W}^0, c^t \geq d \tag{B.13}$$

and

$$\min_{C=1}^{\tilde{\mathcal{C}}^B} \sum_{t=1}^n c^t \leq \text{OPT}_{B^0}^1(W) + 1, c^t \in \mathcal{A}_c^0 \tag{B.14}$$

for every time period t . Setting $B = \frac{1}{2} R^2$, taking \tilde{c}_j to be the minimizer, and taking the expected value of Eq. (B.13) gives

$$(B.15) \quad E_{W,P} \left[\sum_{j \in C} c_j^t \right] \leq E_{W,P} \left[\sum_{j \in C} \tilde{c}_j^t \right] + \frac{1}{2} R^2$$

Taking expected value on both sides of Eq. (B.14) yields

$$(B.16) \quad E_{W,P} \left[\sum_{j \in C} c_j^t \right] \leq E_{W,P} \left[\sum_{j \in C} \tilde{c}_j^t \right] + \frac{1}{2} R^2$$

Therefore

$$(B.17) \quad E_{W,P} \left[\sum_{j \in C} c_j^t \right] \leq E_{W,P} \left[\sum_{j \in C} \tilde{c}_j^t \right] + \frac{1}{2} R^2$$

Combine Eq. (B.15) and Eq. (B.17) we have

$$(B.18) \quad E_{W,P} \left[\sum_{j \in C} c_j^t \right] \leq E_{W,P} \left[\sum_{j \in C} \tilde{c}_j^t \right] + \frac{1}{2} R^2$$

We conclude the proof of case 1 by combining Eq. (B.12) and Eq. (B.18) and noting that $\frac{1}{2} R^2$ is a constant:

$$\text{Regret}(\text{MainALG}) = E_{W,P} \left[\sum_{j \in C} c_j^t \right] - \min_{j \in C} \sum_{t=1}^T c_j^t \leq \frac{1}{2} R^2$$

Case 2: Suppose the underlying arrival model is stochastic where each c_j^t is drawn i.i.d. from an underlying probability distribution \mathcal{D} . We show that the probability that the algorithm switches to the Adversarial Arrival Algorithm is low. More specifically, we show that this probability is no more than $\frac{3X}{2\epsilon}$.

First we prove a Chernoff-like bound for sums with stopping times.

Lemma 19 (Stopping Time Chernoff). Consider a discrete-time random sequence with states (s_1, s_2, \dots) where each state s_t determines two values G_t and H_t with

$G_C - H_C \geq 0 - 2\epsilon$ for some constant $\epsilon > 0$. Suppose $E \gg G_j (C \geq 1) \quad E \gg H_j (C \geq 1)$. Then for every $0 \leq n \leq 1$ and every $\epsilon > 0$ we have

$$P \exists g \text{ such that } \prod_{C=1}^g \frac{G_C}{1+n} \prod_{C=1}^g \frac{H_C}{1+n} \leq \exp(-n^{2-\epsilon})$$

Proof of Lemma 19. Let $q_0 = 1$, and for $g = 1, 2, \dots$, let $q_g = \prod_{C=1}^g \frac{G_C}{1+n} \prod_{C=1}^g \frac{H_C}{1+n}$. Then $q_0 - q_1, q_1 - q_2, \dots$ is a non-negative super-martingale. Indeed, for g we have

$$\frac{q_g}{q_{g-1}} = \frac{G_g}{1+n} \frac{H_g}{1+n} \leq \frac{G_g - H_g}{1+n} \leq -\epsilon$$

where the first inequality is because $G_C - H_C \geq 0 - 2\epsilon$ for every C . Because $E \gg G_j (C \geq 1) \quad E \gg H_j (C \geq 1)$ we get $E \gg q_j (j \geq 0) \leq 1$, which shows $q_0 - q_1, q_1 - q_2, \dots$ is a non-negative super-martingale.

If the event in the statement happens at some g , then

$$\exp \left(\prod_{C=1}^g \frac{G_C}{1+n} \prod_{C=1}^g \frac{H_C}{1+n} \right) \leq \exp(-n^{2-\epsilon})$$

Using $4^{g-1} n^g \leq 1, n$ we get

$$q_g = \prod_{C=1}^g \frac{G_C}{1+n} \prod_{C=1}^g \frac{H_C}{1+n} \leq \exp(-n^{2-\epsilon})$$

Therefore

$$P \exists g \text{ such that } \prod_{C=1}^g \frac{G_C}{1+n} \prod_{C=1}^g \frac{H_C}{1+n} \leq \exp(-n^{2-\epsilon}) \leq \exp(-n^{2-\epsilon})$$

where the second inequality follows by Doob's martingale inequality. „

To analyze the reward obtained so far by the algorithm at a certain time period, we revisit the proof of Proposition 8 and inherit all notations are ed from the proof of Proposition 8. Recall in Eq. (B.3) we have

$$E_{W,P} \gg \sum_{C=1}^g \frac{G_C}{1+n} \sum_{C=1}^g \frac{H_C}{1+n} \leq E_{W,P} \quad \sum_{C=1}^g \frac{G_C}{1+n} \sum_{C=1}^g \frac{H_C}{1+n} \leq \exp(-n^{2-\epsilon})$$

For any $\mathcal{C} \subset \mathcal{X}_C$ and $\mathcal{C}' \subset \mathcal{R}_C^<$ we have $\mathbb{P} \left(\sum_{C=1}^C \sum_{j \in \mathcal{C}} \sum_{\mathcal{C}' \in \mathcal{C}'} P_j^0 F_{\mathcal{C}'}^1 \leq \frac{2n^0 A}{1+n^0} \right) \leq \frac{2n^0 A}{1+n^0}$. Therefore, for the stopping time τ defined in the proof of Proposition 8, we can apply Lemma 19 on $\sum_{C=1}^C \sum_{j \in \mathcal{C}} P_j^0 F_{\mathcal{C}'}^1$ and $\sum_{C=1}^C A_C^1 G_C^0$, which gives

$$\begin{aligned} P_{W(P)} &\leq \sum_{C=1}^C \sum_{j \in \mathcal{C}} P_j^0 F_{\mathcal{C}'}^1 \leq \sum_{C=1}^C A_C^1 G_C^0 \\ &\leq \frac{2n^0 A}{1+n^0} + \sum_{C=1}^C \sum_{j \in \mathcal{C}} P_j^0 F_{\mathcal{C}'}^1 \\ &\leq \sum_{C=1}^C A_C^1 G_C^0 \frac{1+n^0}{1+n^0} + \sum_{C=1}^C \sum_{j \in \mathcal{C}} P_j^0 F_{\mathcal{C}'}^1 \\ P_{W(P)} &\leq \frac{\sum_{C=1}^C \sum_{j \in \mathcal{C}} P_j^0 F_{\mathcal{C}'}^1}{1+n^0} \leq \sum_{C=1}^C \frac{A_C^1 G_C^0}{1+n^0} n^0 A \\ &\leq \exp(n^0) \sum_{C=1}^C A_C^1 G_C^0. \end{aligned}$$

where the first inequality follows because $\sum_{C=1}^C A_C^1 G_C^0 \leq n^0 A$ and $\sum_{C=1}^C \sum_{j \in \mathcal{C}} P_j^0 F_{\mathcal{C}'}^1 \leq \frac{2n^0 A}{1+n^0}$, so

$$\sum_{C=1}^C A_C^1 G_C^0 \leq \frac{2n^0 A}{1+n^0} + \sum_{C=1}^C \sum_{j \in \mathcal{C}} P_j^0 F_{\mathcal{C}'}^1 \leq \frac{2n^0 A}{1+n^0} + \frac{2n^0 A}{1+n^0} = \frac{4n^0 A}{1+n^0};$$

the second inequality is obtained by dividing n^0 on both sides of the inequality; the third inequality utilizes Lemma 19. Plug in $\beta = \frac{1}{2}$ and $\gamma = \frac{1}{2} \log(1+n^0)$ yields

$$(B.19) \quad P_{W(P)} \leq \sum_{C=1}^C \sum_{j \in \mathcal{C}} P_j^0 F_{\mathcal{C}'}^1 \leq \sum_{C=1}^C A_C^1 G_C^0 \left(4A + 2A \log(1+n^0) \right) \leq \frac{1}{2} \sum_{C=1}^C A_C^1 G_C^0.$$

We will use Eq. (B.19) later in bounding the concentration of $\sum_{C=1}^C \sum_{j \in \mathcal{C}} W_j^0$.

Then we look to bound the concentration of $\text{OPT}^1 W^0$. Because $\sum_{C=1}^C A_C^1 \leq A$ for every $\mathcal{C} \subset \mathcal{R}_C^<$, by Hoeffding's inequality we have

$$(B.20) \quad P_{W(P)} \left(\sum_{C=1}^C \sum_{j \in \mathcal{C}} W_j^0 \geq E_{W(P)} + \sqrt{H} \right) \leq \exp \left(-\frac{2H^2}{A^2} \right)$$

and

$$(B.21) \quad P_{W(P)} \left(\sum_{C=1}^C \sum_{j \in \mathcal{C}} W_j^0 \leq E_{W(P)} - \sqrt{H} \right) \leq \exp \left(-\frac{2H^2}{A^2 B} \right) \leq 8B$$

for every $\mathcal{C} \subset \mathcal{R}_C^<$ and $H \geq 0$. Apply Eq. (B.13) and Eq. (B.16) to Eq. (B.20) and take \mathcal{C} to be the minimizer on the left hand side of Eq. (B.16) gives

$$(B.22) \quad P_{W(P)} \left(\text{OPT}^1 W^0 \geq E_{W(P)} + \sqrt{H} \right) \leq \exp \left(-\frac{2H^2}{A^2} \right).$$

Take $H = \frac{P}{A^2} \log^1)^{0 \cdot 2}$ yields

$$(B.23) \quad R_{W(P)} \text{OPT}^1 W^0 \ll E_{W(P)} \gg \text{OPT}^1 W^0 \frac{P}{A^2} \log^1)^{0 \cdot 2}, 1 <, 1^0 A \quad \frac{1}{\cdot}$$

Recall in the steps of Eq. (B.6), we have $\sum_{C=1}^g A^1 C^0$ and $\sum_{C=1}^g c_j P^0$. Combine Eq. (B.19) and Eq. (B.23) gives that, for every $l_j 0$,

$$\begin{aligned} & P_{W(P)} \text{OPT}^1 W^0 \ll \sum_{C=1}^g A^1 C^0 \\ & \quad \left(\sum_{C=1}^g c_j P^0 \right), \frac{P}{A^2} \log^1)^{0 \cdot 2}, 1 <, 1^0 A \\ &^{10} P_{W(P)} \ll E_{W(P)} \gg \text{OPT}^1 W^0 \frac{P}{A^2} \log^1)^{0 \cdot 2}, 1 <, 1^0 A \quad \frac{1}{\cdot} \\ &^{11} P_{W(P)} \ll E_{W(P)} \gg \text{OPT}^1 W^0 \frac{P}{A^2} \log^1)^{0 \cdot 2}, 1 <, 1^0 A \quad \frac{1}{\cdot} \\ &^{12} P_{W(P)} \ll E_{W(P)} \gg \text{OPT}^1 W^0 \frac{P}{A^2} \log^1)^{0 \cdot 2}, 1 <, 1^0 A \quad \frac{2}{\cdot} \\ &^{13} P_{W(P)} \ll E_{W(P)} \gg \text{OPT}^1 W^0 \frac{P}{A^2} \log^1)^{0 \cdot 2}, 1 <, 1^0 A \quad \frac{2}{\cdot} \\ &^{14} P_{W(P)} \ll E_{W(P)} \gg \text{OPT}^1 W^0 \frac{P}{A^2} \log^1)^{0 \cdot 2}, 1 <, 1^0 A \quad \frac{2}{\cdot} \\ &^{15} P_{W(P)} \ll E_{W(P)} \gg \text{OPT}^1 W^0 \frac{P}{A^2} \log^1)^{0 \cdot 2}, 1 <, 1^0 A \quad \frac{3}{\cdot} \end{aligned}$$

Here¹⁰ follows by Eq. (B.23);¹¹ is because $\sum_{C=1}^g A^1 C^0$; ¹² follows by Eq. (B.19);¹³ is because $\sum_{C=1}^g c_j P^0$; ¹⁴ holds since the last three steps of Eq. (B.6) is deterministic in nature; ¹⁵ follows from the last paragraph of the proof of Proposition 8. Take $l = \frac{A^6}{d}$ and note that

$$\left(\sum_{C=1}^g c_j P^0 \right), \frac{P}{A^2} \log^1)^{0 \cdot 2}, 1 <, 1^0 A \ll \frac{A^6}{d} \log^1)^{0 \cdot 2}, 1 <, 1^0 A$$

i.e., there exists a constant $\epsilon_j 0$ such that

$$\left(\sum_{C=1}^g c_j P^0 \right), \frac{P}{A^2} \log^1)^{0 \cdot 2}, 1 <, 1^0 A \ll \epsilon_j \log^1)^{0 \cdot 2}, 1 <, 1^0 A$$

This gives

$$(B.24) \quad R_{W(P)} \text{OPT}^1 W^0 \ll \sum_{C=1}^g A^1 C^0 \epsilon_j \log^1)^{0 \cdot 2}, 1 <, 1^0 A \quad \frac{3}{\cdot}$$

Suppose the algorithm does not switch to the Adversarial Arrival Algorithm before time period: 0^X for some: $0 \leq t_0 \dots - 1 \cdot X \leq 1$. For $t = 0, 1, \dots, 0: 1$, the

algorithm performs the Stochastic Arrival Algorithm during each time block between time periods X_{j-1} and X_j . Apply Eq. (B.14) to Eq. (B.21) over each time block gives that, for every $H_j > 0$,

$$(B.25) \quad P_{W(P)} \mathbb{E}_{W(P)} \left[\sum_{C=X_{j-1}}^{X_j} c^j W_C^0 \right] \leq \text{OPT}^1 W_{(X),1:1:,1^0 X} \leq H_j \exp \left(\frac{2H_j^2}{A^2 X} \right) \cdot$$

Let Z_j be the random variable such that

$$Z_j = \mathbb{E}_{W(P)} \left[\sum_{C=X_{j-1}}^{X_j} c^j W_C^0 \right] - \text{OPT}^1 W_{(X),1:1:,1^0 X} \leq H_j$$

where $W_{(X),1:1:,1^0 X} \sim \mathcal{N}(0, P(X))$. Then by Eq. (B.25) each Z_j is an independent sub-Gaussian random variable with parameter $\sqrt{2A^2 X}$. Therefore Z_j is also a sub-Gaussian random variable with parameter at most $\sqrt{2A^2 X}$. Hence we get

$$(B.26) \quad P_{W(P)} \mathbb{E}_{W(P)} \left[\sum_{C=1}^{X} c^j W_C^0 \right] \leq \text{OPT}^1 W_{(X),1:1:,1^0 X} \leq H_j \exp \left(\frac{2H_j^2}{A^2 X} \right) \cdot$$

$$(B.27) \quad \exp \left(\frac{2H_j^2}{A^2 X} \right) \cdot$$

Note that

$$\mathbb{E}_{W(P)} \left[\sum_{C=1}^{X} c^j W_C^0 \right] = \mathbb{E}_{W(P)} \left[\sum_{C=1}^{X} c^j W_C^0 \right]$$

so combining Eq. (B.22) from time period $C = 1$ to time period $C = X$ and Eq. (B.26) and using union bound we get

$$P_{W(P)} \left[\sum_{C=1}^{X} c^j W_C^0 \right] \leq \text{OPT}^1 W_{(X),1:1:,1^0 X} \leq H_j \exp \left(\frac{2H_j^2}{A^2 X} \right) \cdot$$

Take $H = \frac{P}{2A^2X^3 \log^2(1)} \cdot 2$ yields

$$(B.28) \quad R_{WP} - \text{OPT}^1 W_{:,0X} \leq \sum_{i=0}^{\infty} \frac{P}{2A^2X^3 \log^2(1)} \cdot \frac{1}{2^i} < \frac{1}{2} \cdot \frac{1}{2} = \frac{1}{4}$$

For $i = 0, 1, \dots, \infty$, let $r_{i,j}^{\text{SA}}(W_{:,1:1:,1^0X})$ denote the reward obtained by the Stochastic Arrival Algorithm during each time block between time periods $i, 1$ and $i+1, 1^0X$, then

$$r_{i,j}^{\text{SA}} = \sum_{t=i}^{i+1} r_{j,t}^{\text{SA}}(W_{:,1:1:,1^0X})$$

where $r_{j,t}$ is the total amount of reward obtained between time periods $t, 1$ as defined in the algorithm. Apply Eq. (B.24) on each time block shows that for each we have

$$R_{WP} - \text{OPT}^1 W_{:,1:1:,1^0X} \leq \sum_{i=0}^{\infty} \frac{P}{2A^2X^3 \log^2(1)} \cdot \frac{3}{2^i}$$

and therefore

$$(B.29) \quad R_{WP} - \text{OPT}^1 W_{:,1:1:,1^0X} \leq \sum_{i=0}^{\infty} \frac{P}{2A^2X^3 \log^2(1)} \cdot \frac{3}{2^i} = \frac{3}{2} \cdot \frac{P}{2A^2X^3 \log^2(1)}$$

Let ϵ be a constant such that

$$(B.30) \quad \epsilon \log^2(1) \leq \frac{P}{2A^2X^3 \log^2(1)} \cdot \frac{3}{2} < \frac{1}{2} \cdot \frac{1}{2} = \frac{1}{4}$$

Combine Eq. (B.28) and Eq. (B.29) gives

$$R_{WP} - \text{OPT}^1 W_{:,0X} \leq \sum_{i=0}^{\infty} \frac{P}{2A^2X^3 \log^2(1)} \cdot \frac{1}{2^i} + \frac{3}{2} \cdot \frac{P}{2A^2X^3 \log^2(1)} = \frac{3}{2} \cdot \frac{P}{2A^2X^3 \log^2(1)}$$

Therefore

the algorithm switches to the Adversarial Arrival Algorithm incorrectly

$$= \sum_{i=0}^{\infty} \frac{P}{2A^2X^3 \log^2(1)} \cdot \frac{1}{2^i} + \frac{3}{2} \cdot \frac{P}{2A^2X^3 \log^2(1)} = \frac{3}{2} \cdot \frac{P}{2A^2X^3 \log^2(1)}$$

Case 3: Suppose the underlying arrival model is adversarial and the algorithm switches to the Adversarial Arrival Algorithm at time period $t_0 + 1$ for some $t_0 \in \{0, 1, \dots, T-1\}$. Here, to simplify the notation, we set $t_0 = 1$ if the algorithm never switches to the Adversarial Arrival Algorithm. For time periods $t \in \{1, \dots, T\}$, let W_t be the amount of rewards that the algorithm obtained between time periods $t-1$ and t .

Because the algorithm does not switch at time period $t_0 + 1$, we have

$$\begin{aligned}
 & \mathbb{E}[\text{MainALG } j \text{ } W_t] \geq (1 - \frac{1}{2^{\log t}}) \frac{P}{2} \\
 & \text{OPT}(W_{1:t_0}) \leq \\
 & \max \frac{1}{U} \text{OPT}(W_{1:t_0}) - \text{PRD}(W_{1:t_0}) \\
 \text{(B.31)} \quad & \max \frac{1}{U} \text{OPT}(W_{1:T}) - \text{PRD}(W_{1:T}) \leq XA
 \end{aligned}$$

where the last inequality follows since the total rewards obtained in time periods $t \in \{1, \dots, T\}$ is upper bounded by XA .

Because the algorithm releases the remaining $(T - t_0)$ amount of resources for the remaining $(T - t_0)$ time periods and performs the Adversarial Arrival Algorithm, by Theorem 9

$$\begin{aligned}
 \text{(B.32)} \quad & \max \frac{1}{U} \text{OPT}(W_{(t_0+1):T}) - \text{PRD}(W_{(t_0+1):T}) \geq \mathbb{E}[\text{MainALG } j \text{ } W_t : t \geq t_0 + 1] \\
 & \geq \frac{P}{2}
 \end{aligned}$$

Combining Eq. (B.31) and Eq. (B.32) gives

$$\begin{aligned}
 & \mathbb{E}[\text{MainALG } j \text{ } W_t] \geq XA \\
 = & \mathbb{E}[\text{MainALG } j \text{ } W_t] \geq (1 - \frac{1}{2^{\log t}}) \frac{P}{2} + \mathbb{E}[\text{MainALG } j \text{ } W_t : t \geq t_0 + 1] \\
 & \max \frac{1}{U} \text{OPT}(W_{1:T}) - \text{PRD}(W_{1:T}) \geq \frac{P}{2} \\
 & \max \frac{1}{U} \text{OPT}(W_{(t_0+1):T}) - \text{PRD}(W_{(t_0+1):T}) \geq \frac{P}{2} \\
 & \mathbb{E}[\text{MainALG } j \text{ } W_t] \geq \frac{P}{2}
 \end{aligned}$$

where the last inequality follows since W_t is stationary (Definition 4 and

Observation 6). Hence

$$\limsup_{T \rightarrow \infty} \sup_{W \in \mathcal{W}} \left(\frac{1}{T} \max_{j \in [1, \dots, n]} \text{OPT}^j(W) - \text{PRD}^j(W) \right) \leq \text{MainALG}^j(W) \leq \frac{1}{U} \max_{j \in [1, \dots, n]} \text{OPT}^j(W) - \text{PRD}^j(W)$$

Putting it all together. If the underlying arrival model is stochastic, combining case 1 and case 2 gives

$$\begin{aligned} \text{Regret}^j(\text{MainALG}^j) &= \mathbb{E}_{W \in \mathcal{W}} \left[\text{OPT}^j(W) - \text{MainALG}^j(W) \mid j \text{ never switches} \right] + \mathbb{E}_{W \in \mathcal{W}} \left[\text{OPT}^j(W) - \text{MainALG}^j(W) \mid j \text{ switches} \right] \\ &\leq \mathbb{E}_{W \in \mathcal{W}} \left[\text{OPT}^j(W) - \text{MainALG}^j(W) \mid j \text{ never switches} \right] + \mathbb{E}_{W \in \mathcal{W}} \left[\text{OPT}^j(W) - \text{MainALG}^j(W) \mid j \text{ switches} \right] \end{aligned}$$

By case 1, $\mathbb{E}_{W \in \mathcal{W}} \left[\text{OPT}^j(W) - \text{MainALG}^j(W) \mid j \text{ never switches} \right] \leq \frac{1}{2} \max_{j \in [1, \dots, n]} \left(\frac{1}{U} \text{OPT}^j(W) - \text{PRD}^j(W) \right)$

By case 2, $\mathbb{E}_{W \in \mathcal{W}} \left[\text{OPT}^j(W) - \text{MainALG}^j(W) \mid j \text{ switches} \right] \leq \frac{3}{X} \max_{j \in [1, \dots, n]} \left(\frac{1}{U} \text{OPT}^j(W) - \text{PRD}^j(W) \right)$. Since $\text{OPT}^j(W) \geq 2 \max_{j \in [1, \dots, n]} \left(\frac{1}{U} \text{OPT}^j(W) - \text{PRD}^j(W) \right)$, we have

$$\text{Regret}^j(\text{MainALG}^j) \leq \frac{1}{2} \max_{j \in [1, \dots, n]} \left(\frac{1}{U} \text{OPT}^j(W) - \text{PRD}^j(W) \right)$$

If the underlying arrival model is adversarial, case 3 shows

$$\limsup_{T \rightarrow \infty} \sup_{W \in \mathcal{W}} \left(\frac{1}{T} \max_{j \in [1, \dots, n]} \text{OPT}^j(W) - \text{PRD}^j(W) \right) \leq \text{MainALG}^j(W) \leq \frac{1}{U} \max_{j \in [1, \dots, n]} \text{OPT}^j(W) - \text{PRD}^j(W)$$

This completes the proof. ”

B.7 Experiment Details

Synthetic Experiment

The detailed setups were the following. There were 25 products, where each product was randomly assigned a unique integer price in the range $[1, 25]$ and an embedding that lied randomly in \mathbb{R}^3 . There were 26 types of customers, consisting of 25 customer types that each corresponded to exactly one unique product, and one no-customer type, corresponding to no product being selected in that time interval. For each customer type (apart from the no-customer type), the probability that it would buy product i if recommended was $\text{sigmoid}(4_8 \cdot 10 - 4_9)$, where 4_8 and 4_9 were the 3-dimensional embeddings for product i and 9 , respectively. For the no-customer type, the probabilities were zero - we could not recommend anything.

One instance contained 1000 time periods. To build the arrival sequence, we had a function $\#$ that maps each time period t to a probability of observing a no-customer

type at that time. If a customer did arrive, we chose its type uniformly at random. The initial inventory level was controlled by d . Modeling inventory shortages or excess inventory can be done by changing d . The price of each product was fixed at the start of the experiment and is held constant. To generate predictions, we first calculated the true item counts in the demand sequence for each product. Depending on the arrival model, we applied various amounts of zero-mean Gaussian noise with variance σ^2 to each of the true item counts. We then took these counts, compared them to our inventory, and determined the predicted shadow prices for each product. By changing f , we were able to simulate predictions of different qualities.

The synthetic experiment was run on a MacBook Pro equipped with Apple's M2 Chip. The total compute time was under 20 hours. All offline optimization problems in the algorithms were solved by Gurobi.

We list all the (hyper)parameters used:

- ^ Low inventory level: $d = 0.015$ medium inventory level $d = 0.03$ and high inventory level: $d = 0.06$;
- ^ Root finding bisection parameters: $U = 10^6, V = 0, \epsilon = 10^{-4}, \delta = 1$;
- ^ Perfect predictions $f = 0$, good predictions $f = 5$, and bad predictions $f = 500$;
- ^ Stochastic arrivals $\#^1 C^0 = 0.7$ nonstationary arrivals $\#^1 C^0 = 0.4, \frac{3C}{5000}$, and adversarial arrivals: $\#^1 C^0 = 1^1 C_i 300^0$;
- ^ Parameters for the Main Algorithm (Algorithm 3): $X_{\frac{1}{20}}$ and $! = 7$.

H&M Experiment

The H&M dataset contains two years of online purchase data from H&M customers, consisting of dates, purchase prices, customer IDs, and product ID. For each product, there are basic categorical information about its type, appearance, and department. For computational reasons, we only considered the 5000 most purchased products during this experiment. Our goal was to simulate 90 days of the online marketplace where when a customer selects a product, we recommend three other products in return. Encoding the days using a start day B . We started by building a sequence of customer/no-customer arrivals for the 90 day window: $\{x_t\}_{t=0}^{89}$. Let $\max_{0 \leq t \leq 89} f$ Amount of customers in day B_t . We initialized an empty array of size 90 . For a day B_t , for every product that was purchased in that day, we

randomly placed this product in the array between indices i and $i + 1$. We call this sequence of customer/no-customer interactions our demand sequence. Note that each entry in the tuple contained the product and the price for which it was purchased.

A product's price on a given day was set to be the price of that product purchased by some customer on a given day. To ensure that this process was deterministic, as there could be multiple customers purchasing the same product for different prices, we defined the product's price on that day to be the first time that product was purchased by a customer on that day. If no customer purchased that product, we made the assumption that the product was unavailable and took this under consideration when recommending products during the experiment, as we could not recommend a product that is not available. In order to facilitate this experiment, we built an accurate model which took in two products, along with their prices, and determined the probability that those two products were bought together. We did this using sklearn's Random Forest model. First, we created a 50-dimensional embedding for each product. This was done by creating a matrix where each entry represented that the i th product was bought by the j th customer. Using a matrix factorization collaborative filtering algorithm, we were able to obtain a 50-dimensional embedding for each product. Next, for each product, we created a one-hot vector for "product_group_name", "graphical_appearance_no", "perceived_colour_value_id", "perceived_colour_master_id", "index_code", "index_group_no", and "garment_group_no", and concatenated these one-hot vectors to form a vector of length 102 that contains exactly 7 ones. Given two products, p_1 and p_2 , we created the final 207-dimensional vector we fed into the Random Forest model by concatenating p_1 and p_2 's one-hot vectors, adding in the dot product similarity metric between p_1 and p_2 's embeddings, and finally adding the prices for both items on that specific day. To train this model, we generated 100,000 positive instances, meaning a customer bought products p_1 and p_2 together on the same day, and 1,000,000 negative instances, where we randomly selected a product p_1 purchased by customer i and a product p_2 that was available on that day but not bought by i . The trained model had an AUC of 0.78. For any two products, p_1 and p_2 , that also contained correct price information for that day, we referred to this probability function as $P_{A>1}(p_1, p_2)$, giving us the probability that items p_1 and p_2 were bought together on that specific day.

In executing the Main Algorithms, we modeled the random nature of recommending products to customers, that is, we did not know whether or not a customer would

select the products we recommended. To remedy this, we performed the following procedure to closely model real-world customer decision making. At each time step, we either saw a no-customer, which we would recommend no products, or we observed a product that the customer selected, say A_{42} . Then the value of recommending some product A_{42} was given by $\lambda_{A_{42}} - p_{A_{42}}$, where $\lambda_{A_{42}}$ represented the current price of the recommended product, and $p_{A_{42}}$ was treated as being since the customer would only consume one unit of the recommended product, and $\lambda_{A_{42}}$ was the current shadow price of the recommended item as predicted by the dual variable at time period t . We then recommended the top three products according to this above metric that also satisfied the inventory constraint. Note that if no three products existed to recommend, then we recommended no products. Once we recommended three products to the customer, the customer would pick each product with probability $\frac{1}{3}$ and we in turn received the product's value along with the decrease in inventory only if the customer ended up buying the product. The customer could select anywhere from none to all of the recommended products, and the selections were assumed to be independent of each other.

To generate the prediction for each instance, we used 365 days of data before the starting day of our testing window. For every 5-day span, we added up all the products that were purchased within that interval. We took this and converted it into counts for the embedding vectors. This gave us 50 streams of 73 data points each (one stream per embedding dimension and one data point for each of the 50 combined points). From here, we ran our prediction algorithm (FB Prophet, ARIMA, and Exponential Smoothing) to generate 18 more data points (a total of 91 days) for each of the 50 streams, converted these data points back into counts for the products themselves, and determined the shadow price for each product using these predicted demands for each product. This was done by solving the online problem where the inventory levels are given by the predicted demand. The vector of shadow prices became our prediction.

Here, the prediction \hat{W} is constructed from historical data that reflect such persistent demand patterns. Even though individual arrivals are modeled as stochastic, the realized arrival sequence is influenced by common underlying factors that induce correlation between past observations and future arrivals. As a result, for an arrival sequence W^P , the prediction \hat{W} and the optimal dual variable W^0 are not independent, and the prediction error may scale sub-linearly consistent with

our theoretical assumptions. More broadly, our results highlight that the stochastic model relevant for prediction-based decision making is often best viewed as one with structured randomness: while arrivals may be independent across periods conditional on latent parameters, predictions can exploit partial information about these parameters, leading to meaningful correlation between the predicted and the realized arrival sequence.

When running the Main Algorithm (Algorithm 3), we performed sequential hypothesis testing to determine whether or not the arrival sequence was stochastic or not. We began by assuming the arrival sequence was stochastic. Then we performed the following online hypothesis test: after allowing for a burn-in period of 20 days, for every $t \in \{25, 30, \dots, 85\}$, we performed a one-sided one sample t-test on the number of arrivals in \mathcal{C}_t compared to the average number of arrivals $\bar{C} = 51$. For sufficiently low p -value, chosen to be 0.05, the algorithm switched to be adversarial. Additionally, due to the large amounts of data used, using the bisection algorithm as written in the Stochastic Arrival Algorithm (Algorithm 18) was too computationally inefficient, so instead we used an approximation of this by selecting a step size ϵ , and computing

$$c := \arg \min_{c \in \mathbb{Z}} \left[\frac{\|c - \bar{C}\|^2}{U \|c - \bar{C}\|^2 + V} \right]$$

This method allowed quicker computation, as we were only running a constant number of versions of the Mirror Descent Algorithm for each instance. The larger our set \mathcal{C} was, the closer we would get to the solution outputted by the root finding bisection algorithm.

We list all the (hyper)parameters used:

- ^ Prophet and Exponential Smoothing: default;
- ^ ARIMA parameters: $\alpha = 5$, $\beta = 2$, $\gamma = 1$;
- ^ Random Forest classifier: $\text{num_estimators} = 100$, $\text{max_depth} = 18$;
- ^ $\epsilon = 10^{-11}$, $\delta = 30^{-9}$, $\eta = 8 \times 10^{-3}$;
- ^ Stochastic Arrival Algorithm parameters: $U = 1$, $V = 0$.

Appendix C

APPENDIX FOR CHAPTER 4

C.1 Preliminary Observations and Proofs

Proof of Lemma 2. By Assumption 4, we have

$$\begin{aligned}
 \mathbb{E} \sum_{t=1}^T \ell_t(\hat{c}_t) - \mathbb{E} \sum_{t=1}^T \ell_t(c_t^*) &= \mathbb{E} \sum_{t=1}^T \ell_t(\hat{c}_t) - \mathbb{E} \sum_{t=1}^T \ell_t(c_t^*), \\
 &\stackrel{10^\circ}{=} \mathbb{E} \sum_{t=1}^T \ell_t(\hat{c}_t) - \mathbb{E} \sum_{t=1}^T \ell_t(c_t^*), \\
 &\stackrel{11^\circ}{=} \mathbb{E} \sum_{t=1}^T \ell_t(\hat{c}_t) - \mathbb{E} \sum_{t=1}^T \ell_t(c_t^*), \\
 &\stackrel{12^\circ}{=} 2 \sum_{j=1}^d \lambda_j.
 \end{aligned}$$

Here 10° and 12° follow from $\ell_t(\cdot)$ being λ -Lipchitz in \cdot , and 11° uses the definition of \hat{c}_t . ”

Proof of Observation 3a). By Lemma 2,

$$\mathbb{E} \sum_{t=1}^T \ell_t(\hat{c}_t) - \mathbb{E} \sum_{t=1}^T \ell_t(c_t^*) \leq 2 \sum_{j=1}^d \lambda_j$$

where λ_j is the Lipschitz constant of $\ell_t(\cdot)$. Therefore,

$$\begin{aligned}
 R^{C^{\text{prediction}}}(T) &= \sup_{J \in \mathcal{D}^1} \mathbb{E} \sum_{t=1}^T \ell_t(\hat{c}_t) - \mathbb{E} \sum_{t=1}^T \ell_t(c_t^*) \\
 &\leq 2 \sum_{j=1}^d \lambda_j.
 \end{aligned}$$

Finally, by the construction of \hat{c}_t

$$\mathbb{E} \sum_{t=1}^T \ell_t(\hat{c}_t) - \mathbb{E} \sum_{t=1}^T \ell_t(c_t^*) \leq 2 \sum_{j=1}^d \lambda_j.$$

Taking $\lambda_j = 2^{-j}$ gives the desired result. ”

For any time periods $0 \leq t \leq T-1$, let

$$R^{C^1}(T) \stackrel{0 \leq t \leq T-1}{=} \sup_{J \in \mathcal{D}^1} \mathbb{E} \sum_{t=0}^{T-1} \ell_t(\hat{c}_t) - \mathbb{E} \sum_{t=0}^{T-1} \ell_t(c_t^*)$$

be the regret incurred by the policy from time 0 to time $T-1$. We make the following two useful observations:

Observation 7. For any policy c , $R^{(c)} \geq 0 - 1/4 - \epsilon$ for some universal constant $\epsilon > 0$.

Proof of Observation 7. Because c and μ are both bounded, $c - \mu$ is also bounded in $[-1, 1]$ where

$$\min = \inf_{c \in [-1, 1]} (c - \mu)^2$$

and

$$\max = \sup_{c \in [-1, 1]} (c - \mu)^2$$

Thus,

$$R^{(c)} \geq 0 - 1/4 = \sup_{J \subseteq [D]} E_J^c \left(\sum_{C=0}^T (c - \mu_C)^2 - \min_{C=0}^T (c - \mu_C)^2 \right)$$

Take $\epsilon = \max_{C=0}^T \min_{j \in J} \mu_j$ gives the desired result. As a remark, note that by footnote 3 we have $\max_{C=0}^T \min_{j \in J} \mu_j \leq \max_{C=0}^T \max_{j \in J} \mu_j$, so ϵ can be taken as $\max_{C=0}^T \max_{j \in J} \mu_j$.

Observation 8. For any policy c such that, from time t to time $t+1$, c estimates the mean at time t to be $\lambda_{2t} \in [-1, 1]$ for every $0 \leq t < T$ and then order $T/2$ $\arg \min_{C=0}^T (c - \mu_C)^2$ we have $R^{(c)} \geq 0 - 1/4 - \epsilon$ for some universal constant $\epsilon > 0$.

Proof of Observation 8. Let L be the Lipschitz constant of μ , then by Lemma 2

$$|c - \mu_C| \leq L |C - t|$$

Therefore,

$$R^{(c)} \geq 0 - 1/4 = \sup_{J \subseteq [D]} E_J^c \left(\sum_{C=0}^T (c - \mu_C)^2 - \min_{C=0}^T (c - \mu_C)^2 \right) \leq 2 \sup_{J \subseteq [D]} E_J^c \sum_{C=0}^T L^2 |C - t|^2$$

Taking $\epsilon = 2L^2$ gives the desired result.

Observation 8 implies that any policy that estimates the mean accurately at each time achieves low regret.

Finally, we will make use of standard sub-Gaussian concentration:

Lemma 20 (Hoeffding's Inequality, e.g. [169]). Let $\{n_i\}_{i=1}^B$ be independent, mean-zero, sub-Gaussian variables with sub-Gaussian norm at most C :

$$P\left\{\sum_{i=1}^B n_i \geq C\right\} \leq 2 \exp\left(-\frac{C^2}{2}\right) \quad \text{for all } C \geq 0.$$

Then for some universal constant C ,

$$P\left\{\sum_{i=1}^B \tilde{O}_i \geq C\right\} \leq 2 \exp\left(-\frac{C^2}{2}\right) \quad \text{for all } C \geq 0.$$

Moreover, \tilde{O}_i is upper bounded by 1444 .

C.2 Proof of Lemma 3

Proof of Lemma 3. Because our bounds are all asymptotic, we ignore the rounding and write $\tilde{O}_i = \tilde{O}_i^{11E^{0.2}}$ to simplify the notation.

By Observation 7 $P\left\{\sum_{i=1}^B \tilde{O}_i \geq C\right\} \leq 2 \exp\left(-\frac{C^2}{2}\right)$ for some universal constant $C = 3 \max\{1, \max_{i \in [B]} \tilde{O}_i\}$.

From now on we consider the time period from $C = 1$ to $C = \frac{1}{\epsilon}$. We first upper bound the total estimation error of the mean $\sum_{i=1}^B \tilde{O}_i$. Note that

$$\begin{aligned} \sum_{i=1}^B \tilde{O}_i &\stackrel{10^0}{=} \sum_{i=1}^B \tilde{O}_i^{11E^{0.2}} \\ &= \sum_{i=1}^B \tilde{O}_i^{11E^{0.2}} \mathbb{1}_{\tilde{O}_i \leq C} + \sum_{i=1}^B \tilde{O}_i^{11E^{0.2}} \mathbb{1}_{\tilde{O}_i > C} \\ &= \sum_{i=1}^B \tilde{O}_i^{11E^{0.2}} \mathbb{1}_{\tilde{O}_i \leq C} + \sum_{i=1}^B \tilde{O}_i^{11E^{0.2}} \mathbb{1}_{\tilde{O}_i > C} \\ &\stackrel{10^0}{=} \sum_{i=1}^B \tilde{O}_i^{11E^{0.2}} \mathbb{1}_{\tilde{O}_i \leq C} + \sum_{i=1}^B \tilde{O}_i^{11E^{0.2}} \mathbb{1}_{\tilde{O}_i > C} \end{aligned}$$

where 10^0 follows because $\mathbb{1}_{\tilde{O}_i \leq C}$ is the projection of $\mathbb{1}_{\tilde{O}_i \leq C}$ on $\mathbb{1}_{\tilde{O}_i \leq C}$.

We bound these two parts separately through the following two lemmas.

Lemma 21. There exists a universal constant C such that

$$\sum_{i=1}^B \tilde{O}_i^{11E^{0.2}} \mathbb{1}_{\tilde{O}_i \leq C} \leq 2 \sum_{i=1}^B \tilde{O}_i^{11E^{0.4}}.$$

Proof of Lemma 21. For any $\epsilon, 1 \leq C$ we have

$$\frac{1}{n} \sum_{B=C}^{\tilde{O}} \max_{C=B, C+1} j_B \cdot d \leq \frac{1}{n} \sum_{B=C}^{\tilde{O}} \max_{C=B, C+1} j_B \cdot d$$

Also, by bounded demand variation,

$$\begin{aligned} \frac{1}{n} \sum_{B=C}^{\tilde{O}} \max_{C=B, C+1} j_B \cdot d^2 &\stackrel{10^0}{=} \sum_{9=1}^{\tilde{O}} \sum_{8=1}^{\tilde{O}-\epsilon} \max_{9=8, 9} j_B \cdot d^2 \\ &\stackrel{11^0}{=} \sum_{8=1}^{\tilde{O}} \sum_{9=1}^{\tilde{O}-\epsilon} \max_{9=8, 9} j_B \cdot d^2 \\ &\stackrel{12^0}{=} \dots \\ &\stackrel{\wedge)}{=} \dots \end{aligned}$$

where 10^0 is obtained by partitioning the sum into time windows, 11^0 is obtained by exchanging the summations, and 12^0 follows by the definition of demand variation. Since $C_9 = 9, 8, 1 : 9 = 0 - 1 - \dots - d$ is a partition off $1 - \dots - 1$ for all $8 = 1 - \dots - \dots$.

By Cauchy Schwarz inequality,

$$\frac{1}{n} \sum_{B=C}^{\tilde{O}} \max_{C=B, C+1} j_B \cdot d \leq \sqrt{\frac{1}{n} \sum_{B=C}^{\tilde{O}} \max_{C=B, C+1} j_B \cdot d^2} \leq \sqrt{\frac{1}{n} \sum_{B=C}^{\tilde{O}} \max_{C=B, C+1} j_B \cdot d^2} \stackrel{\wedge)}{=} \dots$$

Take $\epsilon = \sqrt{\dots}$ gives the desired result. "

Lemma 22. There exists a universal constant $\epsilon > 0$ such that

$$E_J^{C, \text{red}} \left(\frac{1}{n} \sum_{B=C}^{\tilde{O}} \max_{C=B, C+1} j_B \cdot d \right) \leq \dots$$

Proof of Lemma 22. Because each X_B is sub-Gaussian, each X_B is sub-Gaussian. Therefore there exists a constant $\epsilon > 0$ where $X_B = \inf_{0 \leq t \leq 1} E \exp(t X_B - \frac{1}{2} t^2 G^2) \leq \exp(-\frac{1}{2} t^2 G^2 + \epsilon t^2 G^2)$. Let $X = \max_{B=1, \dots, n} X_B$, then by Hoeffding's inequality, for any $\epsilon, 1 \leq C$ we have

$$P \left(\frac{1}{n} \sum_{B=C}^{\tilde{O}} \max_{C=B, C+1} j_B \cdot d \geq G \right) \leq 2 \exp \left(-\frac{d^2 = G^2}{C \sum_{B=C}^{\tilde{O}} j_B} \right) \leq 2 \exp \left(-\frac{d = G}{X^2} \right)$$

where $d_j \geq 0$ is a universal constant. Therefore we have

$$\begin{aligned}
 E_J^{C_{xed}} \left(\frac{1}{B=C} \tilde{O}^1 \right) &= P \left(\frac{1}{B=C} \tilde{O}^1 \right) G \cdot 3G \\
 &= 2 \exp \frac{d}{X^2} G^2 \cdot 3G \\
 &= \frac{cX^2}{d} \\
 &= \left(\frac{cX^2}{d^\lambda} \right)^{1E^{10 \cdot 4}}.
 \end{aligned}$$

Hence

$$E_J^{C_{xed}} \left(\frac{\tilde{O}}{C_{=,1}} \frac{1}{B=C} \tilde{O}^1 \right) = \left(\frac{cX^2}{d^\lambda} \right)^{1E^{10 \cdot 4}} \left(\frac{cX^2}{d^\lambda} \right)^{13, E^{0 \cdot 4}}.$$

Take $\frac{3}{12} = \frac{q}{12} \frac{cX^2}{d^\lambda}$ gives the desired result. Note by Lemma 20 1444, so
 $\frac{3}{12} \frac{X^p}{c} \frac{1}{4^\lambda}$ ”

Now we finish the proof of Lemma 3. With Lemma 21 and Lemma 22, we conclude that

$$\begin{aligned}
 E_J^{C_{xed}} \left(\frac{\tilde{O}}{C_{=,1}} \right) &= E_J^{C_{xed}} \left(\frac{\tilde{O}}{C_{=,1}} \frac{1}{B=C} \tilde{O}^1 \right) \\
 &= E_J^{C_{xed}} \left(\frac{\tilde{O}}{C_{=,1}} \frac{1}{B=C} \tilde{O}^1 \right) \\
 &= \left(\frac{cX^2}{d^\lambda} \right)^{13, E^{0 \cdot 4}}.
 \end{aligned}$$

Therefore by Observation 8 there exists a universal constant 2^{10-10} such that

$$\begin{aligned}
 R^{C_{xed}}(1)^0 &= R^{C_{xed}}(1)^0 \gg 1 - \frac{1}{4}, R^{C_{xed}}(1)^0 \gg (1 - \frac{1}{4}) \\
 &= \left(\frac{cX^2}{d^\lambda} \right)^{13, E^{0 \cdot 4}}, \left(\frac{cX^2}{d^\lambda} \right)^{13, E^{0 \cdot 4}} \\
 &= \left(\frac{cX^2}{d^\lambda} \right)^{13, E^{0 \cdot 4}}.
 \end{aligned}$$

Take $\frac{1}{12} = \frac{1}{12} \frac{cX^2}{d^\lambda}$ we get $R^{C_{xed}}(1)^0 = \left(\frac{cX^2}{d^\lambda} \right)^{13, E^{0 \cdot 4}}$ ”

C.3 Proof of Theorem 2

Proof of Theorem 2. Because our bounds are all asymptotic, we ignore the roundings and write $\epsilon = \wedge^{11 E_8^{0.2}}$ to simplify the notation. By Observation 7 $R^{C^{\text{shrinking}}_1) \gg 1 - \epsilon^{3.4}$ for some universal constant

$$1 = 3 \max_{1 \leq i \leq n} \max_{g \in X_i} \epsilon^{10 - 10 \cdot i}$$

From now on we consider the time periods after $t^{3.4}$. Let i be the smallest index such that $E_i \leq \epsilon$, then $E_i \leq \frac{1}{\log} \epsilon$, so

$$(C.1) \quad \epsilon^{E_i} \leq \epsilon^{11 \cdot \log \epsilon} = 4 \epsilon^{E_i - 4} \epsilon$$

First, we show that \tilde{c} is close to c when ϵ is small via the following lemma:

Lemma 23. For every ϵ with the corresponding window size $\epsilon = \wedge^{11 E_8^{0.2}}$, there exists a universal constant $W > 0$ such that

$$P_{\substack{W \\ C=\epsilon, 1}}^{\tilde{c}} \tilde{c} \leq \frac{P_{\substack{W \\ \log}}{c}}{P_{\substack{W \\ \log}}{\tilde{c}}} \leq \frac{2}{3.2}$$

Proof of Lemma 23. Same as in the proof of Lemma 22, by Hoeffding's inequality for any $\epsilon, 1 > C > \epsilon$ we have

$$P_{\substack{W \\ =9 \\ B=C \geq \epsilon}}^{\tilde{c}} \frac{1}{n_B} \tilde{c} \leq \frac{P_{\substack{W \\ =9 \\ B=C \geq \epsilon}}{c}}{P_{\substack{W \\ =9 \\ B=C \geq \epsilon}}{\tilde{c}}} \leq 2 \exp \left(- \frac{d^2}{C^2 X_B^2} \right) \leq 2 \exp \left(- \frac{d^2}{X^2} \right)$$

where d and X are the same as in the previous proof. Let ϵ be large enough so that

$$(C.2) \quad \frac{d^2}{X^2} \geq \frac{5}{2}$$

Take $G = \frac{P_{\substack{W \\ \log}}{c}}{P_{\substack{W \\ \log}}{\tilde{c}}}$ and plug in $\epsilon = \wedge^{11 E_8^{0.2}}$ yields

$$P_{\substack{W \\ =9 \\ B=C \geq \epsilon}}^{\tilde{c}} \frac{1}{n_B} \tilde{c} \leq \frac{P_{\substack{W \\ \log}}{c}}{P_{\substack{W \\ \log}}{\tilde{c}}} \leq 2 \exp \left(- \frac{d^2}{X^2} \log \right) \leq \frac{2}{5.2}$$

Then we get

$$\begin{aligned}
 & \left(\tilde{O}_{C=\mathfrak{g},1} \frac{1}{\mathfrak{g}} \tilde{O}_{B=C=\mathfrak{g}}^1 n_B W^{\overline{p \log}} \right)^{13, E_{\mathfrak{g}^0 \cdot 4}} \\
 & \left(\max_{C=\mathfrak{g},1} \frac{1}{\mathfrak{g}} \tilde{O}_{B=C=\mathfrak{g}}^1 n_B W^{\overline{p \log}} \right)^{1 E_{\mathfrak{g}^0 \cdot 4}} \\
 & \left(\tilde{O}_{C=\mathfrak{g},1} \frac{1}{\mathfrak{g}} \tilde{O}_{B=C=\mathfrak{g}}^1 n_B W^{\overline{p \log}} \right)^{1 E_{\mathfrak{g}^0 \cdot 4}} \\
 & \frac{2}{3 \cdot 2}
 \end{aligned}$$

where 10^0 follows by union bound. Note that since $\tilde{O}_{C=\mathfrak{g},1} \frac{1}{\mathfrak{g}} \tilde{O}_{B=C=\mathfrak{g}}^1 n_B W^{\overline{p \log}} \leq \tilde{O}_{C=\mathfrak{g},1} \frac{1}{\mathfrak{g}} \tilde{O}_{B=C=\mathfrak{g}}^1 n_B W^{\overline{p \log}}$, by Lemma 21 we know $\tilde{O}_{C=\mathfrak{g},1} \frac{1}{\mathfrak{g}} \tilde{O}_{B=C=\mathfrak{g}}^1 n_B W^{\overline{p \log}} \leq \tilde{O}_{C=\mathfrak{g},1} \frac{1}{\mathfrak{g}} \tilde{O}_{B=C=\mathfrak{g}}^1 n_B W^{\overline{p \log}}$, and in the proof of Lemma 3 we have

$$\tilde{O}_{C=\mathfrak{g},1} \lambda_C^{\mathfrak{g}} \cdot \tilde{O}_{C=\mathfrak{g},1} \frac{1}{\mathfrak{g}} \tilde{O}_{B=C=\mathfrak{g}}^1 n_B W^{\overline{p \log}} \leq \tilde{O}_{C=\mathfrak{g},1} \frac{1}{\mathfrak{g}} \tilde{O}_{B=C=\mathfrak{g}}^1 n_B$$

Therefore

$$\begin{aligned}
 & \left(\tilde{O}_{C=\mathfrak{g},1} \lambda_C^{\mathfrak{g}} \cdot \tilde{O}_{C=\mathfrak{g},1} \frac{1}{\mathfrak{g}} \tilde{O}_{B=C=\mathfrak{g}}^1 n_B W^{\overline{p \log}} \right)^{13, E_{\mathfrak{g}^0 \cdot 4}} \\
 & \left(\tilde{O}_{C=\mathfrak{g},1} \frac{1}{\mathfrak{g}} \tilde{O}_{B=C=\mathfrak{g}}^1 n_B W^{\overline{p \log}} \right)^{13, E_{\mathfrak{g}^0 \cdot 4}} \\
 & \left(\tilde{O}_{C=\mathfrak{g},1} \frac{1}{\mathfrak{g}} \tilde{O}_{B=C=\mathfrak{g}}^1 n_B W^{\overline{p \log}} \right)^{13, E_{\mathfrak{g}^0 \cdot 4}} \\
 & \frac{2}{3 \cdot 2}
 \end{aligned}$$

For each \mathfrak{g} , let \mathfrak{g}^0 be the event $\tilde{O}_{C=\mathfrak{g},1} \lambda_C^{\mathfrak{g}} \cdot \tilde{O}_{C=\mathfrak{g},1} \frac{1}{\mathfrak{g}} \tilde{O}_{B=C=\mathfrak{g}}^1 n_B W^{\overline{p \log}} \geq \tilde{O}_{C=\mathfrak{g},1} \frac{1}{\mathfrak{g}} \tilde{O}_{B=C=\mathfrak{g}}^1 n_B W^{\overline{p \log}}$, then $P(\mathfrak{g}^0) \leq \frac{2}{3 \cdot 2}$ for each \mathfrak{g} . First we assume that \mathfrak{g}^0 does not happen for any \mathfrak{g} . We break the proof into three lemmas.

Lemma 24. Each time an if condition happens, for the current we have $\delta \leq \delta_0$. Therefore $\delta \leq \delta_0$ throughout the algorithm.

Proof of Lemma 24. Suppose an if condition happens at time t triggered by some index $i \in S$, then $\delta(t) \leq \delta_0$. Suppose for the sake of contradiction that $\delta(t) > \delta_0$, then

$$\begin{aligned}
 & \delta(t) \leq \delta_0 \quad \text{by assumption} \\
 & \delta(t) \leq \delta_0 \quad \text{by assumption} \\
 & \delta(t) \leq \delta_0 \quad \text{by assumption} \\
 & \delta(t) \leq \delta_0 \quad \text{by assumption} \\
 & \delta(t) \leq \delta_0 \quad \text{by assumption}
 \end{aligned}$$

where $\delta(t) \leq \delta_0$ comes from the triangle inequality and $\delta(t) \leq \delta_0$ is because $\delta(t) \leq \delta_0$ and $\delta(t) \leq \delta_0$; since $i \in S$, by our assumption neither i nor j occurs, so we get $\delta(t) \leq \delta_0$; $\delta(t) \leq \delta_0$ follows since $i \in S$. This contradicts with $\delta(t) > \delta_0$. Therefore, we must have $\delta(t) \leq \delta_0$. Because the index i never decreases and can only increase by 1 each time an if condition happens, throughout the algorithm. ”

Lemma 25. Suppose two consecutive if conditions occur at times t_1 and t_2 , then $R^{C^{shrinking}}(t_1) \geq R^{C^{shrinking}}(t_2) - \epsilon$ for some universal constant $\epsilon > 0$.

Proof of Lemma 25. At time t_1 where $\delta(t_1) \leq \delta_0$, by Lemma 24 $\delta(t_1) \leq \delta_0$. We have

$$\begin{aligned}
 & \delta(t_1) \leq \delta_0 \quad \text{by Lemma 24} \\
 & \delta(t_1) \leq \delta_0 \quad \text{by Lemma 24} \\
 & \delta(t_1) \leq \delta_0 \quad \text{by Lemma 24} \\
 & \delta(t_1) \leq \delta_0 \quad \text{by Lemma 24} \\
 & \delta(t_1) \leq \delta_0 \quad \text{by Lemma 24}
 \end{aligned}$$

where $\delta(t_1) \leq \delta_0$ comes from the triangle inequality and $\delta(t_1) \leq \delta_0$ is because $\delta(t_1) \leq \delta_0$ and $\delta(t_1) \leq \delta_0$; the first part of $\delta(t_1) \leq \delta_0$ follows by our assumption that i does not occur, and the second part $\delta(t_1) \leq \delta_0$ follows since the if condition is not triggered between

time C and time C^0 ; 13^0 follows by Equation (C.1). Then by Observation 8 there exists a universal constant $\epsilon = 2^{-2} 10^{-1}$ such that

$$R^{C^{shrinking_1})^0} \gg C^{-\epsilon} 1^{1/4} \left(\sum_{B=C}^{C^1} j \lambda_B^8 \cdot B \right)^{1/4} \left(W \log \right)^{1/4} \left(\frac{P}{\bar{\lambda}} \right)^{13, E^{0.4}}.$$

Take $\epsilon = 34^{1.4} \cdot W \cdot \frac{P}{\bar{\lambda}}$ gives the desired result. „

The above proof also works for $R^{C^{shrinking_1})^0} \gg 3^{.4} \cdot (1 - C_{st})^{1/4}$ if the if condition happens at time C_{st} , or $R^{C^{shrinking_1})^0} \gg 3^{.4} \cdot (1 -)^{1/4}$ if the if condition never happens. Note that the index s never decreases and increases by 1 if and only if the if condition happens. Suppose that the last if condition happens at time C_{last} . On we have

$$\begin{aligned} R^{C^{shrinking_1})^0} &\gg 3^{.4} \cdot (1 - C_{last})^{1/4} \cdot \left(\sum_{B=C_{last}}^{C^0} j \lambda_B^8 \cdot B \right)^{1/4} \left(W \log \right)^{1/4} \left(\frac{P}{\bar{\lambda}} \right)^{13, E^{0.4}} \\ &= 2 \frac{\log}{\log^{11} \left(\frac{1}{\log} \right)^0} \left(\sum_{B=C_{last}}^{C^0} j \lambda_B^8 \cdot B \right)^{1/4} \left(W \log \right)^{1/4} \left(\frac{P}{\bar{\lambda}} \right)^{13, E^{0.4}} \\ &= 2 \frac{\log}{\log^{11} \left(\frac{1}{\log} \right)^0} \left(\sum_{B=C_{last}}^{C^0} j \lambda_B^8 \cdot B \right)^{1/4} \left(W \log \right)^{1/4} \left(\frac{P}{\bar{\lambda}} \right)^{13, E^{0.4}} \log^{5.2}. \end{aligned}$$

for some universal constant $\epsilon = 2^{-2} 10^{-1}$. Here 10^0 follows by Lemma 25, 11^0 is because the maximum index B is C^0 , 12^0 is because $\sum_{B=C_{last}}^{C^0} j \lambda_B^8 \cdot B \leq \sum_{B=C_{last}}^{C^0} j \lambda_B^8 \cdot C^0$, and 13^0 follows by $\log^{11} \left(\frac{1}{\log} \right)^0 \geq G^0$ when G is small.

Finally, we analyze the time periods after C_{last} .

Lemma 26. $R^{C^{shrinking_1})^0} \gg C_{last}^{-\epsilon} \left(\sum_{B=C_{last}}^{C^0} j \lambda_B^8 \cdot B \right)^{1/4} \left(W \log \right)^{1/4} \left(\frac{P}{\bar{\lambda}} \right)^{13, E^{0.4}}$ for some universal constant $\epsilon = 2^{-2} 10^{-1}$.

Proof of Lemma 26. Because the if condition happens at C_{last} by Lemma 24 either

$\delta \dot{Y}$ or $\delta = \dots$. Suppose $\delta \dot{Y}$, then similar to Lemma 25 we have

$$\begin{aligned}
 & \tilde{O}_{B=C_{fast}} j \lambda_B^8 \cdot \tilde{O}_{B=C_{fast}}^{10^0} j \lambda_B \cdot \tilde{O}_{B=C_{fast}} j \lambda_B^8 \wedge \tilde{O}_{B=C_{fast}} \\
 & \tilde{O}_{B=C_{fast}}^{11^0} j \lambda_B \cdot \tilde{O}_{B=C_{fast}} j \lambda_B^8 \wedge \tilde{O}_{B=C_{fast}} \\
 & \tilde{Y} \left(\frac{P}{W \log} \right), P_{\bar{\lambda}} \Big)^{13, E^{0.4}}, 2 \left(\frac{P}{W \log} \right), P_{\bar{\lambda}} \Big)^{13, E^{0.4}} \\
 & 34^{1.4} \left(\frac{P}{W \log} \right), P_{\bar{\lambda}} \Big)^{13, E^{0.4}}
 \end{aligned}$$

where 10^0 comes from the triangle inequality and 11^0 is because $C_{fast} i) \cdot 3.4$ and $\approx 8 \wedge) \cdot 1.2$; the first part of 12^0 follows by our assumption that δ does not occur, and the second part of 12^0 follows since the if condition is never triggered after C_{fast} ; 13^0 follows by Equation (C.1). By Observation 8, we get

$$R^{C_{shrinking} 1)^0} \gg C_{fast}^{-1/4} 34^{1.4} \left(\frac{P}{W \log} \right), P_{\bar{\lambda}} \Big)^{13, E^{0.4}}$$

for some universal constant $\epsilon = 2 \cdot 2^{10} - 1^0$.

On the other hand, suppose $\delta = \dots$. Note $\dots \Big)^ E \Big)^ E$ and after time C_{fast} the Shrinking-Time-Window Policy just performs the Fixed-Time-Window Policy with variation parameter ϵ so by Lemma 3

$$R^{C_{shrinking} 1)^0} \gg C_{fast}^{-1/4} \left(\frac{P}{W \log} \right), P_{\bar{\lambda}} \Big)^{13, E^{0.4}} \left(\frac{P}{W \log} \right), P_{\bar{\lambda}} \Big)^{13, E^{0.4}}$$

for some universal constant $\epsilon = 2 \cdot 10 - 1^0$, where the last inequality again follows by Equation (C.1).

Taking $\epsilon = 34^{1.4} \max \left(\frac{P}{W}, P_{\bar{\lambda}} \right) - \epsilon$ we get

$$R^{C_{shrinking} 1)^0} \gg C_{fast}^{-1/4} \left(\frac{P}{W \log} \right), P_{\bar{\lambda}} \Big)^{13, E^{0.4}} \left(\frac{P}{W \log} \right).$$

”

Combining everything above, in the case where δ doesn't happen for any ϵ , we get

$$\begin{aligned}
 R^{C_{shrinking} 1)^0} &= R^{C_{shrinking} 1)^0} \gg 1 - \epsilon)^{3.4} \cdot \frac{1}{4}, R^{C_{shrinking} 1)^0} \gg \left(\frac{P}{W \log} \right), P_{\bar{\lambda}} \Big)^{13, E^{0.4}} \cdot \frac{1}{4} \\
 &\cdot R^{C_{shrinking} 1)^0} \gg C_{fast}^{-1/4} \\
 &1) \cdot 3.4, 2) \left(\frac{P}{W \log} \right), P_{\bar{\lambda}} \Big)^{13, E^{0.4}} \log^{5.2}, 3) \left(\frac{P}{W \log} \right).
 \end{aligned}$$

Now let us consider the case where θ happens for some θ . Because $\theta \leq \frac{2}{3 \cdot 2}$ for each θ , by union bound

$$P(\theta \text{ happens for some } \theta) \leq \sum_{\theta} \frac{2}{3 \cdot 2} \cdot \frac{2}{\theta}$$

where the last inequality follows $\frac{1}{\theta} = \log \frac{2}{\theta}$. By Observation 7 we always have $R^{c^{shrinking}_1} \leq \frac{1}{4}$ for some universal constant $\frac{1}{4} \leq 10^{-10}$. Therefore in summary we have

$$R^{c^{shrinking}_1} \leq P(\theta \text{ doesn't happen for any } \theta) + P(\theta \text{ happens for some } \theta) \leq \frac{1}{4} + \sum_{\theta} \frac{2}{3 \cdot 2} \cdot \frac{2}{\theta} \leq \frac{1}{4} + \sum_{\theta} \frac{1}{3 \cdot \theta} \leq \frac{1}{4} + \frac{1}{3} \sum_{\theta} \frac{1}{\theta} \leq \frac{1}{4} + \frac{1}{3} \cdot \frac{1}{2} = \frac{1}{4} + \frac{1}{6} = \frac{5}{12}$$

Therefore, there exists a constant $\frac{1}{4} \leq 10^{-10}$ such that

$$R^{c^{shrinking}_1} \leq \frac{1}{3} \sum_{\theta} \frac{1}{\theta} \leq \frac{1}{3} \cdot \frac{1}{2} = \frac{1}{6}$$

”

C.4 Proof of Proposition 10, Proposition 11, and Observation 3b)

Note that Proposition 10 and Observation 3b) can be easily deduced from Proposition 11 by setting $\theta = 1$ and $E = 1$ respectively, so it suffices to prove Proposition 11.

Proof of Proposition 11. We construct the following worst-case problem instance: divide the time horizon into cycles of length $\frac{1}{E}$ (since the analysis below is compatible with scaling, for simplicity we assume $\frac{1}{E}$ is an integer), so there are $\frac{1}{E}$ cycles. Assume that the demand distribution within each cycle is a Bernoulli distribution that equals $\frac{1}{2}$ with probability $\frac{1}{2}$ and equals $\frac{1}{20}$ with probability $\frac{1}{2}$. At the beginning of each cycle, we set the order of the upcoming cycle to be either $\frac{1}{2}$ or $\frac{1}{20}$, each with probability $\frac{1}{2}$. Set $c = R$, and $c_C = c = 1$ for all C so the optimal ordering amount is the median of the demand distribution at each time.

First we show that the demand variation is at most $\frac{1}{E}$. Note that since $c_C = 1$ is fixed within each cycle, only the times between cycles contribute to the demand variation. Suppose the cycle changes between time C and $C + 1$, then

$$|c_{C+1} - c_C| \leq \left| \frac{1}{20} - \frac{1}{2} \right| \leq \frac{1}{20} \leq \frac{1}{5} \cdot \frac{1}{E}$$

Because there are $\frac{1}{5}E^{0.2}$ cycles,

$$\sum_{c=1}^{\frac{1}{5}E^{0.2}} \left(\frac{1}{5}\right)^{1E^{0.2}} = \frac{1}{5}E.$$

Then we add predictions into the instance. We divide the analysis into two cases.

Case 1: $\frac{3E}{4}$. For each c we set θ_c to be either $\left(\frac{1}{2}, \frac{1}{20}\right)^{1E^{0.4}}$ or $\left(\frac{1}{2}, \frac{1}{20}\right)^{1E^{0.4}}$, each with probability $\frac{1}{2}$. Then because each θ_c and θ'_c are i.i.d., θ_c provides no information about θ'_c . Hence the predictions are useless in this instance. Note that

$$\sum_{c=1}^{\frac{3E}{4}} \left(\frac{1}{2}, \frac{1}{20}\right)^{1E^{0.4}} = \left(\frac{1}{5}\right)^{13E^{0.4}} \left(\frac{1}{5}\right)^{0.2}$$

so the prediction accuracy is within $\frac{1}{5}$.

Then we analyze the amount of regret incurred. For $1 \leq t \leq \frac{1}{5}E^{0.2}$, let θ_t and θ'_t be i.i.d. distributions respectively, where θ_t is Bernoulli $\left(\frac{1}{2}, \frac{1}{20}\right)^{1E^{0.4}}$ and θ'_t is Bernoulli $\left(\frac{1}{2}, \frac{1}{20}\right)^{1E^{0.4}}$. Then the Kullback-Leibler divergence D_{KL} from θ_t to θ'_t is

$$D_{KL}(\theta_t \parallel \theta'_t) = \frac{1}{2} \left(\frac{1}{2}, \frac{1}{20}\right)^{1E^{0.4}} \log \frac{\left(\frac{1}{2}, \frac{1}{20}\right)^{1E^{0.4}}}{\left(\frac{1}{2}, \frac{1}{20}\right)^{1E^{0.4}}} + \frac{1}{2} \left(\frac{1}{2}, \frac{1}{20}\right)^{1E^{0.4}} \log \frac{\left(\frac{1}{2}, \frac{1}{20}\right)^{1E^{0.4}}}{\left(\frac{1}{2}, \frac{1}{20}\right)^{1E^{0.4}}}$$

We show that $D_{KL}(\theta_t \parallel \theta'_t) \geq \frac{13}{20} \frac{1}{G^2}$. Let $G = \left(\frac{1}{20}\right)^{1E^{0.4}}$, then because $E \geq 2 \gg \frac{1}{20}$, $\frac{1}{G^2} \gg \frac{1}{20}$. Note that the $D_{KL}(\theta_t \parallel \theta'_t) = \frac{13}{20} \frac{1}{G^2}$ for $G = 0$ and we have

$$\frac{3}{3G} D_{KL}(\theta_t \parallel \theta'_t) \Big|_{G=0} = 0 = \frac{3}{3G} \frac{13}{20} \frac{1}{G^2} \Big|_{G=0}$$

and

$$\frac{3^2}{3G} D_{KL}(\theta_t \parallel \theta'_t) = \frac{1}{G^2} \frac{1}{0.25^{0.2}} \frac{1}{26} = \frac{3^2}{3G^2} \frac{13}{20}$$

for $G \geq \frac{1}{20}$. This shows $D_{KL}(\theta_t \parallel \theta'_t) = 0$ at $G = 0$ the first derivative is 0 at $G = 0$ and the second derivative is non-negative for $G \geq \frac{1}{20}$. This implies

$$D_{KL}(\theta_t \parallel \theta'_t) \geq \frac{13}{20} \frac{1}{G^2}.$$

Let $\theta = \prod_{s=1}^{\frac{1}{5}E^{0.2}} \theta_s$ and $\theta' = \prod_{s=1}^{\frac{1}{5}E^{0.2}} \theta'_s$ be the two possible demand distributions within a cycle, then because θ_s are i.i.d. and θ'_s are i.i.d.,

$$D_{KL}(\theta \parallel \theta') = \sum_{s=1}^{\frac{1}{5}E^{0.2}} D_{KL}(\theta_s \parallel \theta'_s) \geq \frac{13}{20} \frac{1}{G^2} = \frac{13}{20}.$$

We claim that within each cycle, any attempt to distinguish between θ and $\bar{\theta}$ has at least a constant probability of making a mistake, i.e., one cannot effectively estimate the θ value within each cycle. Let $C : \{0, 1\}^{11E^{0.2}}$ be any classifier that takes the demand observations within a cycle as inputs and determine the true demand distribution of this cycle. Let $\tau = C^{-1}(0)$ be the event where C classifies the demand observations as from demand distribution θ . Then by Pinsker's inequality,

$$j\theta^{1-\theta} + \bar{\theta}^{1-\theta} \geq \frac{1}{2} \text{KL}(\theta \| \bar{\theta}) \geq \frac{13}{40}$$

where $\theta^{1-\theta}$ is the probability of τ happening under the condition that the true demand distribution is θ , and the same for $\bar{\theta}^{1-\theta}$. Therefore we have

$$P\{C \text{ makes a mistake}\} = \theta^{1-\theta} + \bar{\theta}^{1-\theta} \geq \frac{13}{40} = 1 - \frac{13}{40}.$$

Hence for any classifier C the probability of making the wrong guess τ in the current cycle is at least $1 - \frac{13}{40}$.

Note for demand distribution $\theta = 1$ and

$$\begin{aligned} 1 - \theta^{1-\theta} - \bar{\theta}^{1-\theta} &= \frac{1}{2} \left(\frac{1}{20} \right)^{1E^{10.4}} - \left(\frac{1}{2} \right)^{1E^{10.4}} \\ &= \left(\frac{1}{5} \right)^{1E^{10.4}}. \end{aligned}$$

and similarly for demand distribution $\bar{\theta} = 0$ and

$$\begin{aligned} 1 - \bar{\theta}^{1-\theta} - \theta^{1-\theta} &= 1 - \frac{1}{2} \left(\frac{1}{20} \right)^{1E^{10.4}} - \left(\frac{1}{2} \right)^{1E^{10.4}} \\ &= \left(\frac{1}{5} \right)^{1E^{10.4}}. \end{aligned}$$

Therefore each wrong guess τ incurs a difference between $1 - C(\theta)$ and $1 - C(\bar{\theta})$ by $\left(\frac{1}{5} \right)^{1E^{10.4}}$, which incurs a regret of $1 - C(\theta) - 1 - C(\bar{\theta}) = \left(\frac{1}{5} \right)^{1E^{10.4}}$. Therefore over T time periods the expected total regret of any General Policy c satisfies

$$\begin{aligned} R^c(T) &\geq 1 - \frac{13}{40} \left(\frac{1}{5} \right)^{1E^{10.4}} \\ &= 1 - \frac{13}{40} \left(\frac{1}{5} \right)^{13E^{0.4}} \\ &= 1 - \frac{13}{40} \left(\frac{1}{5} \right)^{\min\{13, E^{0.4}\} - 0.9} \end{aligned}$$

where the last equality follows from $\frac{3E}{4}$. Take $T = 1 - \frac{13}{40} \left(\frac{1}{5} \right)^{13E^{0.4}}$ gives the desired result.

Case 2: $0 \leq \frac{3E}{4}$. For the first $0 \leq \frac{3E}{4}$ time periods of each cycle, we set θ to be either $\frac{1}{2}$ or $\frac{1}{20}$ with probability $\frac{1}{2}$. Note that $0 \leq \frac{3E}{4}$ implies $0 \leq \frac{13E}{4} \leq \frac{1E}{2}$, so time periods of length $0 \leq \frac{3E}{4}$ are indeed contained in each cycle, which has length $11E^0 \cdot 2$. For the other time periods we set $\theta = \frac{1}{5}$. Then, similar as in the case above, the predictions are useless for the first $0 \leq \frac{3E}{4}$ time periods of each cycle in this instance. Since there are $\frac{11E^0 \cdot 2}{0 \leq \frac{3E}{4}}$ number of cycles,

$$\sum_{C=1}^{\frac{11E^0 \cdot 2}{0 \leq \frac{3E}{4}}} \left(\frac{1}{2} \left(\frac{1}{2} \right)^{11E^0 \cdot 2} + \frac{1}{2} \left(\frac{1}{20} \right)^{11E^0 \cdot 2} \right) = \frac{1}{5}$$

so the prediction accuracy is within $\frac{1}{5}$

Again, the same analysis as the case above shows that for the first $0 \leq \frac{3E}{4}$ time periods of each cycle, the probability of making the wrong guess in the current cycle is at least $\frac{13}{40}$. Also, each wrong guess incurs a regret of $\frac{1}{5}$. Because there are $\frac{11E^0 \cdot 2}{0 \leq \frac{3E}{4}}$ number of cycles, over time periods the expected total regret of any General Policy satisfies

$$\begin{aligned} R^C &\geq \frac{13}{40} \frac{1}{5} \frac{11E^0 \cdot 2}{0 \leq \frac{3E}{4}} \\ &= \frac{13}{40} \frac{1}{5} \\ &= \frac{13}{40} \frac{1}{5} \min\{13, E^0 \cdot 4 - 0\} \end{aligned}$$

where the last equality follows from $0 \leq \frac{3E}{4}$. Take $\frac{13}{40} \frac{1}{5}$ gives the desired result. ”

C.5 General Variation

Recall that for any sequence of means μ_1, \dots, μ_T , we define the demand variation to be its quadratic variation

$$+ = \max_{P \in \mathcal{P}} \sum_{i=1}^T (\mu_i - \mu_{i-1})^2$$

where \mathcal{P} is the set of all partitions. The choice of quadratic variation follows from previous literature [95, 96]. We can also define the demand variation using what we will call λ -variation for some $0 < \lambda \leq 1$:

$$+ = \max_{P \in \mathcal{P}} \sum_{i=1}^T (\mu_i - \mu_{i-1})^\lambda$$

In the special case $\delta = 1$, this is the total variation, which has also been used extensively in previous literature [37, 52, 92, 122]. We prove the following lower bound, which generalizes Proposition 10.

Proposition 20 (Lower Bound: Nonstationary Newsvendor with δ -variation). Suppose the demand variation is defined using δ -variation and $\delta \in (0, 1]$. For any variation parameter $\delta \in (0, 1]$ and any policy c (which may depend on the knowledge of δ), we have

$$R^c(1) \geq \frac{1}{2} \delta^{1/2} E^{0,1/2,1/2} -$$

where $\frac{1}{2} \delta^{1/2}$ is a universal constant.

Proof of Proposition 20. Similar to the proof of Proposition 10, we construct the following worst-case problem instance: divide the time horizon into cycles of length $\frac{1}{2} \delta^{1/2} E^{0,1/2,1/2}$ (since the analysis below is compatible with scaling, for simplicity we assume $\frac{1}{2} \delta^{1/2} E^{0,1/2,1/2}$ is an integer), so there are $\frac{1}{2} \delta^{1/2} E^{0,1/2,1/2}$ cycles. Assume that the demand distribution within each cycle is a Bernoulli distribution that equals to 1 with probability $\frac{1}{2}$ and equals to 0 with probability $\frac{1}{2}$. At the beginning of each cycle, we set the order of the upcoming cycle to be either $\frac{1}{2} \delta^{1/2} E^{0,1/2,1/2}$ or $\frac{1}{2} \delta^{1/2} E^{0,1/2,1/2} + 1$, each with probability $\frac{1}{2}$. Set $c = R_c$ and $c = 1$ for all C so the optimal ordering amount is the median of the demand distribution at each time.

First we show that the demand variation is at most $\delta^{1/2} E^{0,1/2,1/2}$. Note that since $c = 1$ is fixed within each cycle, only the times between cycles contribute to the demand variation. Suppose the cycle changes between time C and $C + 1$, then

$$|c_{C+1} - c_C| \leq \frac{1}{2} \delta^{1/2} E^{0,1/2,1/2} + 1 \leq \frac{1}{2} \delta^{1/2} E^{0,1/2,1/2} + \frac{1}{2} \delta^{1/2} E^{0,1/2,1/2} = \delta^{1/2} E^{0,1/2,1/2}.$$

Because there are $\frac{1}{2} \delta^{1/2} E^{0,1/2,1/2}$ cycles, $\delta^{1/2} E^{0,1/2,1/2} \cdot \frac{1}{2} \delta^{1/2} E^{0,1/2,1/2} = \frac{1}{2} \delta E^{0,1/2,1/2}$.

Then, following the calculations in the proof of Proposition 11, the probability of making the wrong guess δ in the current cycle is at least $\frac{13}{40}$. Also, each wrong guess of δ incurs a regret of $\frac{1}{5} \delta^{1/2} E^{0,1/2,1/2}$ at each time period. Therefore over $\frac{1}{2} \delta^{1/2} E^{0,1/2,1/2}$ time periods the expected total regret of any General Policy c satisfies

$$R^c(1) \geq \frac{1}{2} \delta^{1/2} E^{0,1/2,1/2} \left(\frac{13}{40} \cdot \frac{1}{5} \delta^{1/2} E^{0,1/2,1/2} \right) = \frac{1}{2} \delta^{1/2} E^{0,1/2,1/2} \cdot \frac{13}{200}.$$

Take $\frac{1}{2} \delta^{1/2} E^{0,1/2,1/2} = \frac{13}{40} \delta^{1/2} E^{0,1/2,1/2}$ gives the desired result. ”

C.6 Proof of Proposition 12

Proof of Proposition 12. Set $\mathcal{C} = \mathbb{R}_+$ and $1_C = c = 1$ for all C so the optimal ordering amount is the median of the demand distribution at each time. We construct the following two problem instances:

Instance 1: for all C set $\mathcal{C}^{11^0} = 8.5 \cdot 10^{-2}$. Then $\mathcal{C}^{11^0} = 1$. Set $0_C^{11^0} = 3_C^{11^0}$, where $3_C^{11^0}$ is the realization of \mathcal{C}^{11^0} . Since $\mathcal{C}^{11^0} = 1$ for all C the variation parameter $E^{11^0} = 0$. Because $0_C^{11^0}$ is a constant away from \mathcal{C}^{11^0} with probability 1, $0^{11^0} = 1$. Because the optimal order amount is $\mathcal{C}^{11^0} = 1$,

$$1 \cdot \mathcal{C}^{11^0} - \mathcal{C}^{11^0} = E \left[1_C^{11^0} \cdot \mathcal{C}^{11^0} \right] = E \left[1_C^{11^0} \cdot 1 \right] = 1 \cdot 1 = 1.$$

Instance 2: for all C set $\mathcal{C}^{12^0} = \mathcal{C}^{12^0} = 8.5 \cdot 10^{-2}$ i.e., \mathcal{C}^{12^0} is a one-point distribution. Set $0_C^{12^0} = 3_C^{12^0}$, where $3_C^{12^0}$ is the realization of \mathcal{C}^{12^0} . Since \mathcal{C}^{12^0} changes by a constant amount from \mathcal{C}^{12^0} with probability 1 throughout $\mathcal{C} = 2 \dots$, the variation parameter $E^{12^0} = 1$. Because $0_C^{12^0} = \mathcal{C}^{12^0}$ for every C , $0^{12^0} = 0$. Because the optimal order amount is $\mathcal{C}^{12^0} = \mathcal{C}^{12^0}$,

$$1 \cdot \mathcal{C}^{12^0} - \mathcal{C}^{12^0} = E \left[1_C^{12^0} \cdot \mathcal{C}^{12^0} \right] = E \left[1_C^{12^0} \cdot \mathcal{C}^{12^0} \right] = 0.$$

Note that a General Policy can only observe information 0_C and 3_C 's. Because $0_C^{11^0} = 3_C^{11^0}$ and $0_C^{12^0} = 3_C^{12^0}$ have the same distribution for every C , no General Policies can distinguish between the two instances. For any General Policy \mathcal{C} be its output at time C . Without loss of generality, we may assume $\mathcal{C} \geq 0$ since any other ordering amount is clearly sub-optimal. Then

$$1 \cdot \mathcal{C}^{11^0} - \mathcal{C} = E \left[1_C^{11^0} \cdot \mathcal{C} \right] = \frac{\mathcal{C}^{11^0} \cdot \mathcal{C}^{11^0}}{4} = \frac{1}{2}$$

and

$$1 \cdot \mathcal{C}^{12^0} - \mathcal{C} = E \left[1_C^{12^0} \cdot \mathcal{C} \right] = \frac{2 \cdot \mathcal{C}^{12^0}}{2} = 1.$$

Let $R^C \gg 1/4$ denote the regret of on instance 1 and $R^C \gg 2/4$ denote the regret

of C on instance 2, then

$$\begin{aligned}
 R^C \gg (1/4)^0 \gg R^C \gg (2/4)^0 &= \sum_{\ell=1}^{\tilde{O}} \left(\frac{1}{\ell} \cdot \theta^{\ell^0} - \theta^{\ell^0} \right) + \sum_{\ell=1}^{\tilde{O}} \left(\frac{1}{\ell} \cdot \theta^{\ell^{1^0}} - \theta^{\ell^{1^0}} \right) \\
 &\quad + \sum_{\ell=1}^{\tilde{O}} \left(\frac{1}{\ell} \cdot \theta^{\ell^{2^0}} - \theta^{\ell^{2^0}} \right) + \sum_{\ell=1}^{\tilde{O}} \left(\frac{1}{\ell} \cdot \theta^{\ell^{2^0}} - \theta^{\ell^{2^0}} \right) \\
 &= \sum_{\ell=1}^{\tilde{O}} \left(\frac{1}{2} - \frac{1}{2} \right) + \sum_{\ell=1}^{\tilde{O}} \left(1 - 0 \right) \\
 &= \tilde{O}.
 \end{aligned}$$

Because $R^C \gg (1/4)^0 \gg R^C \gg (2/4)^0$, any General Policy C incurs regret at least $\tilde{O}(\sqrt{L})$ on at least one of the instances. Since no General Policy can distinguish between the two instances, we can always choose the worse one of the two instances to feed to the policy. Also, since $E^{1^0} = 0, \theta^{1^0} = 1, E^{2^0} = 1, \theta^{2^0} = 0$, in both instances we have $\theta < \frac{3-E}{4}$. Therefore for any General Policy C there always exists a problem instance with $\theta < \frac{3-E}{4}$ such that $R^C \gg (1/4)^0 \gg R^C \gg (2/4)^0 = \tilde{O}(\sqrt{L})$ on the instance. "

C.7 Proof of Theorem 3

Proof of Theorem 3. Because our bounds are all asymptotic, we ignore the rounding and write $\tilde{O}(\sqrt{L})$ to simplify the notation. We slightly abuse the notation and reuse $\tilde{O}(\sqrt{L})$ as constants that differ from the statement of Theorem 3.

Because $c_t^{\text{PERP}} = c_t^{\text{prediction}}$ for t from 1 to $\tilde{O}(\sqrt{L})$, by Observation 3

$$R^C \gg (1/4)^0 \gg R^C \gg (1/4)^0 = \tilde{O}(\sqrt{L})$$

for some universal constant $\tilde{O}(\sqrt{L})$. Also, by Observation 7, since

$$\tilde{O}(\sqrt{L}) \ll \tilde{O}(\sqrt{L})$$

there exists some universal constant $\tilde{O}(\sqrt{L})$ such that

$$R^C \gg (1/4)^0 \gg R^C \gg (1/4)^0 = \tilde{O}(\sqrt{L})$$

Take $\tilde{O}(\sqrt{L}) = \max\{\tilde{O}(\sqrt{L}), \tilde{O}(\sqrt{L})\}$ we get

$$R^C \gg (1/4)^0 \gg R^C \gg (1/4)^0 = \tilde{O}(\sqrt{L})$$

Now we consider the time periods after time $\tilde{O}(\sqrt{L})$. First, same as in the proof of Lemma 22, by Hoeffding's inequality for any $\tilde{O}(\sqrt{L}) \leq t \leq L$ we have

$$\mathbb{P} \left(\sum_{B=\ell}^{\tilde{O}} n_B G \geq 2 \exp \left(-\frac{d^1 = G^2}{\sum_{B=\ell} \chi_B^2} \right) \right) \leq 2 \exp \left(-\frac{d^1 = G^2}{\chi^2} \right)$$

where d and X are the same as in the previous proof. Set $W \geq 0$ to be large enough so that

$$(C.3) \quad \frac{d^2 W^2}{\chi^2} \leq 2.$$

Take $G = \frac{W^2}{\log W}$ and plug in $\lambda = \frac{1}{W}$ yields

$$\mathbb{P} \left(\frac{1}{n_B} \sum_{B=\ell}^{\infty} \tilde{\mathcal{O}}_B \left(\frac{W^2}{\log W} \right)^{E-1} \leq 2 \exp \left(\frac{d^2 W^2}{\chi^2} \log W \right) \right) \geq \frac{2}{3}.$$

Then we get

$$\begin{aligned} & \mathbb{P} \left(\frac{1}{n_B} \sum_{B=\ell}^{\infty} \tilde{\mathcal{O}}_B \left(\frac{W^2}{\log W} \right)^{E-1} \leq 2 \exp \left(\frac{d^2 W^2}{\chi^2} \log W \right) \right) \\ & \geq \mathbb{P} \left(\max_{\ell=1}^{\infty} \frac{1}{n_B} \sum_{B=\ell}^{\infty} \tilde{\mathcal{O}}_B \left(\frac{W^2}{\log W} \right)^{E-1} \leq 2 \exp \left(\frac{d^2 W^2}{\chi^2} \log W \right) \right) \\ & \geq \mathbb{P} \left(\frac{1}{n_B} \sum_{B=\ell}^{\infty} \tilde{\mathcal{O}}_B \left(\frac{W^2}{\log W} \right)^{E-1} \leq 2 \exp \left(\frac{d^2 W^2}{\chi^2} \log W \right) \right) \\ & \geq \frac{2}{3}. \end{aligned}$$

where (a) follows by union bound. Note Lemma 21 says

$$\frac{1}{n_B} \sum_{B=\ell}^{\infty} \tilde{\mathcal{O}}_B \left(\frac{W^2}{\log W} \right)^{E-1} \leq \frac{1}{n_B} \sum_{B=\ell}^{\infty} \tilde{\mathcal{O}}_B \left(\frac{W^2}{\log W} \right)^{E-1}.$$

and in the proof of Lemma 3 we have

$$\frac{1}{n_B} \sum_{B=\ell}^{\infty} \tilde{\mathcal{O}}_B \left(\frac{W^2}{\log W} \right)^{E-1} \leq \frac{1}{n_B} \sum_{B=\ell}^{\infty} \tilde{\mathcal{O}}_B \left(\frac{W^2}{\log W} \right)^{E-1}.$$

Therefore

$$\begin{aligned} & \mathbb{P} \left(\frac{1}{n_B} \sum_{B=\ell}^{\infty} \tilde{\mathcal{O}}_B \left(\frac{W^2}{\log W} \right)^{E-1} \leq 2 \exp \left(\frac{d^2 W^2}{\chi^2} \log W \right) \right) \\ & \geq \mathbb{P} \left(\frac{1}{n_B} \sum_{B=\ell}^{\infty} \tilde{\mathcal{O}}_B \left(\frac{W^2}{\log W} \right)^{E-1} \leq 2 \exp \left(\frac{d^2 W^2}{\chi^2} \log W \right) \right) \\ & \geq \mathbb{P} \left(\frac{1}{n_B} \sum_{B=\ell}^{\infty} \tilde{\mathcal{O}}_B \left(\frac{W^2}{\log W} \right)^{E-1} \leq 2 \exp \left(\frac{d^2 W^2}{\chi^2} \log W \right) \right) \\ & \geq \frac{2}{3}. \end{aligned}$$

where \mathcal{O}^1 is because $c^{\text{PERP}} = c^{\text{prediction}}$ before time B . Hence $R^{c^{\text{PERP}}(1)^0} =_{\mathcal{O}^1} 1 - B^{-1/4}$ is on the order of $\mathcal{O}^{13, E^0 \cdot 4} \left(\frac{P}{\log} \right)$, so there exists some universal constant $\mathcal{O}^2 \geq 10^{-10}$ such that

$$R^{c^{\text{PERP}}(1)^0} =_{\mathcal{O}^1} 1 - B^{-1/4} \geq \mathcal{O}^2 \left(\frac{P}{\log} \right)^{13, E^0 \cdot 4}.$$

Also, since after time B we have $c^{\text{PERP}} = c^{\text{fixed}}$, by Lemma 3 there exists some universal constant $\mathcal{O}^3 \geq 10^{-10}$ such that

$$R^{c^{\text{PERP}}(1)^0} \geq_{\mathcal{O}^3} R^{c^{\text{fixed}}(1)^0} \geq \mathcal{O}^3 \left(\frac{P}{\log} \right)^{13, E^0 \cdot 4}.$$

In summary, if

$$\mathcal{O}^{\ell=1} \left(\frac{P}{\log} \right)^{13, E^0 \cdot 4} \geq \mathcal{O}^{\ell=1} \left(\frac{P}{\log} \right)^{13, E^0 \cdot 4} \cdot \mathcal{O}^{\ell=1} \left(\frac{P}{\log} \right)^{13, E^0 \cdot 4}.$$

we have

$$\begin{aligned} R^{c^{\text{PERP}}(1)^0} &= R^{c^{\text{PERP}}(1)^0} \geq_{\mathcal{O}^1} R^{c^{\text{PERP}}(1)^0} \geq_{\mathcal{O}^2} 1 - B^{-1/4} \geq_{\mathcal{O}^3} R^{c^{\text{PERP}}(1)^0} \geq_{\mathcal{O}^3} \left(\frac{P}{\log} \right)^{13, E^0 \cdot 4} \\ &= \mathcal{O}^1 \left(\frac{P}{\log} \right)^{\min\{13, E^0 \cdot 4\}}. \end{aligned}$$

where the last equality is because $\min\left\{\frac{3-E}{4}, 0\right\} = \frac{3-E}{4}$.

Second, suppose $\mathcal{O}^{\ell=1} \left(\frac{P}{\log} \right)^{13, E^0 \cdot 4} < \mathcal{O}^{\ell=1} \left(\frac{P}{\log} \right)^{13, E^0 \cdot 4}$. By Observation 7 there exists some universal constant $\mathcal{O}^4 \geq 10^{-10}$ such that $R^{c^{\text{PERP}}(1)^0} \geq \mathcal{O}^4 \left(\frac{P}{\log} \right)^{13, E^0 \cdot 4}$.

Therefore combining the above two scenarios we get

$$\begin{aligned} R^{c^{\text{PERP}}(1)^0} &\geq \mathcal{O}^{\ell=1} \left(\frac{P}{\log} \right)^{\min\{13, E^0 \cdot 4\}} \\ &\geq \mathcal{O}^{\ell=1} \left(\frac{P}{\log} \right)^{\min\{13, E^0 \cdot 4\}} \geq \mathcal{O}^{\ell=1} \left(\frac{P}{\log} \right)^{\min\{13, E^0 \cdot 4\}} \cdot \mathcal{O}^{\ell=1} \left(\frac{P}{\log} \right)^{\min\{13, E^0 \cdot 4\}} \\ &= \mathcal{O}^1 \left(\frac{P}{\log} \right)^{\min\{13, E^0 \cdot 4\}} \geq \mathcal{O}^2 \left(\frac{P}{\log} \right)^{\min\{13, E^0 \cdot 4\}}. \end{aligned}$$

Hence $R^{c^{\text{PERP}}(1)^0} \geq \mathcal{O}^2 \left(\frac{P}{\log} \right)^{\min\{13, E^0 \cdot 4\}}$ for some universal constant $\mathcal{O}^2 \geq 10^{-10}$.

Case 2: the if condition does not happen. In this case $c^{\text{PERP}} = c^{\text{prediction}}$, so by Observation 3a) we immediately have $R^{c^{\text{PERP}}(1)^0} \geq \mathcal{O}^5 \left(\frac{P}{\log} \right)^{\min\{13, E^0 \cdot 4\}}$ for some universal

constant $c_5 \geq 2^{10} - 1$. Also, following the same analysis as the part of Case 1 where an if condition has not happened, i.e., between time $t = s - 1$ and time $B - 1$, we get

$$R^{C^{\text{PERP}}_1}(s) \leq c_5 \cdot 2^{10} \cdot \frac{P}{\log} \cdot 2^{P \cdot \frac{1}{s}} \cdot 1^{13 \cdot E \cdot 4} \cdot \frac{P}{\log}$$

for some universal constant $c_6 \geq 2^{10} - 1$. Then we have

$$R^{C^{\text{PERP}}_1}(s) \leq R^{C^{\text{PERP}}_1}(s) \cdot \frac{1}{s} \leq c_5 \cdot R^{C^{\text{PERP}}_1}(s) \cdot \frac{1}{s} \leq c_6 \cdot \frac{P}{\log}$$

Hence take $\epsilon = \max\{c_5 - 1, c_6\}$ we have

$$R^{C^{\text{PERP}}_1}(s) \leq \epsilon \cdot \frac{P}{\log}$$

"

C.8 Unknown E and Known θ

We design a policy for the case of unknown E and known θ . Our policy utilizes the famous Exp3 algorithm (Exponential-weight Algorithm for Exploration and Exploitation) as a subroutine. For the sake of completeness we state Exp3 in our setting and its regret bound below. One can refer to [13] for a more detailed discussion on Exp3.

Proposition 21 (Corollary 3.2 in [13]). *Exp3 achieves worst-case regret*

$$R^{\text{Exp3}}(s) \leq \min\{R^{C^{\text{shrinking}}_1}(s), R^{C^{\text{prediction}}_1}(s)\} + \frac{P}{2^{1 \ln 2^{14}} \cdot 1^{\max}} \cdot \frac{P}{\log}$$

We refer the readers to [13] for the proof. Note that Exp3 incurs an additive $\frac{P}{\log}$ regret on top of $\min\{R^{C^{\text{shrinking}}_1}(s), R^{C^{\text{prediction}}_1}(s)\}$, which is a lower order term if $\theta \geq \frac{1}{2}$. On the other hand, if $\theta < \frac{1}{2}$, we have $\min\{R^{C^{\text{shrinking}}_1}(s), R^{C^{\text{prediction}}_1}(s)\} = R^{C^{\text{prediction}}_1}(s)$, so we can simply apply the Prediction Policy. This idea gives the policy of unknown E and known θ .

Observation 9 (Upper Bound: Unknown E and Known θ). *For any variation parameter $E \geq 2^{10} - 1$ and any accuracy parameter $\theta \geq 2^{10} - 1$, the Divide-Into-Cases Policy C^{Divide} achieves worst-case regret*

$$R^{C^{\text{Divide}}_1}(s) \leq \min\{R^{C^{\text{shrinking}}_1}(s), R^{C^{\text{prediction}}_1}(s)\} \cdot \log^{5 \cdot 2}$$

where c is a universal constant.

Proof of Theorem 9. If $\theta \geq \frac{1}{2}$, then $\min\{R^{C^{\text{shrinking}}_1}(s), R^{C^{\text{prediction}}_1}(s)\} = R^{C^{\text{prediction}}_1}(s)$. The result follows from Observation 3. If $\theta < \frac{1}{2}$, then $\min\{R^{C^{\text{shrinking}}_1}(s), R^{C^{\text{prediction}}_1}(s)\} = R^{C^{\text{shrinking}}_1}(s)$. The result follows from Proposition 21 and Theorem 2.

"

C.9 Experiment Details

In the synthetic experiment we used triple exponential smoothing (Holt Winters) to generate the demand sequences. Triple exponential smoothing takes in the following parameters: data smoothing factor $U \in (0, 1]$, trend smoothing factor $V \in (0, 1]$, seasonal change smoothing factor $W \in (0, 1]$, and season length $l \in \mathbb{N}$. Given historical observations G_1, \dots, G_{l-1} , triple exponential smoothing outputs $G_{l,t}$ for $t \geq 1$, which is an estimate of $G_{l,t}$, according to the following formula:

$$\begin{aligned} B_0 &= G_0 \\ B_\ell &= U \frac{G_\ell}{G_{\ell-l}} + (1-U) B_{\ell-1} + 1_\ell - 1_{\ell-1} \\ 1_\ell &= V B_\ell - B_{\ell-1} + (1-V) 1_{\ell-1} \\ 2_\ell &= W \frac{G_\ell}{B_\ell} + (1-W) 2_{\ell-l} \\ G_{l,t} &= B_{\ell-1} + 1_{\ell-1} + 2_{\ell-1} + 1_{\ell-1} \cdot 1_{\ell-1} \cdot 1_{\ell-1} \end{aligned}$$

where B_ℓ is the smoothed value of the constant part for time ℓ , 1_ℓ is the sequence of best estimates of the linear trend that are superimposed on the seasonal changes, and 2_ℓ is the sequence of seasonal correction factors.

In our experiment, we first generated a demand sequence for 30 time periods where the demand at each time period is drawn uniformly between 80 and 120. This was treated as the historical observations and was fixed throughout the experiment. Then for each set of parameters (U, V, W, l) , which we will specify later in each case, we used triple exponential smoothing to generate the mean of demands for 365 time periods, where at each time period we also added a random Gaussian noise with mean equals to 0 and variance equals to 5. The true demand at each time period was generated as a Poisson variable with the corresponding mean.

We ran two sets of experiments:

- **Fixed E :** We fixed a single set of parameters $(U, V, W, l) = (0.5, 0.5, 0.5, 30)$ for the demand sequence, and generated 1,000 different predictions of this demand sequence, each from a set of “predicted” parameters $(\hat{U}, \hat{V}, \hat{W}, \hat{l})$ where each $\hat{U}, \hat{V}, \hat{W}$ was sampled uniformly at random from $(0.2, 0.8]$ and \hat{l} from $\{10, 20, 30\}$. Thus the variation parameter E was fixed, and the accuracy parameter O varied across instances.
- **Fixed O :** We generated 1,000 demand sequences by selecting the parameters (U, V, W, l) uniformly at random where each U, V, W was sampled uniformly at

random from $\mathcal{U}(0.2, 0.8)$ and β from $\mathcal{U}(10, 20, 30)$. We then generated predictions by changing each parameter 10% (e.g., U becomes either $1.1U$ or $0.9U$) and using the corresponding sequence. Thus the variation parameter E varied across instances, but the accuracy parameter O was (roughly) fixed.

C.10 Additional Experimental Results

In our experiments on real data, we used four popular forecasting methods to generate predictions for each instance: Exponential Smoothing (Holt Winters), ARIMA, Prophet, and LightGBM. Experiments on the datasets yielded the histograms in Figure 4.4. To show the performance of PERP is robust to the forecasting method, in Figure C.1 we further divide the three histograms in Figure 4.4 into twelve histograms separated by the four forecasting methods.

