



## Management Science

Publication details, including instructions for authors and subscription information:  
<http://pubsonline.informs.org>

### Modeling Consumer Footprints on Search Engines: An Interplay with Social Media

Anindya Ghose, Panagiotis G. Ipeirotis, Beibei Li

To cite this article:

Anindya Ghose, Panagiotis G. Ipeirotis, Beibei Li (2019) Modeling Consumer Footprints on Search Engines: An Interplay with Social Media. *Management Science* 65(3):1363-1385. <https://doi.org/10.1287/mnsc.2017.2991>

Full terms and conditions of use: <https://pubsonline.informs.org/Publications/Librarians-Portal/PubsOnLine-Terms-and-Conditions>

This article may be used only for the purposes of research, teaching, and/or private study. Commercial use or systematic downloading (by robots or other automatic processes) is prohibited without explicit Publisher approval, unless otherwise noted. For more information, contact [permissions@informs.org](mailto:permissions@informs.org).

The Publisher does not warrant or guarantee the article's accuracy, completeness, merchantability, fitness for a particular purpose, or non-infringement. Descriptions of, or references to, products or publications, or inclusion of an advertisement in this article, neither constitutes nor implies a guarantee, endorsement, or support of claims made of that product, publication, or service.

Copyright © 2018, INFORMS

Please scroll down for article—it is on subsequent pages



With 12,500 members from nearly 90 countries, INFORMS is the largest international association of operations research (O.R.) and analytics professionals and students. INFORMS provides unique networking and learning opportunities for individual professionals, and organizations of all types and sizes, to better understand and use O.R. and analytics tools and methods to transform strategic visions and achieve better outcomes.

For more information on INFORMS, its publications, membership, or meetings visit <http://www.informs.org>

# Modeling Consumer Footprints on Search Engines: An Interplay with Social Media

Anindya Ghose,<sup>a</sup> Panagiotis G. Ipeirotis,<sup>a</sup> Beibei Li<sup>b</sup>

<sup>a</sup>Stern School of Business, New York University, New York, New York 10012; <sup>b</sup>Heinz College, Carnegie Mellon University, Pittsburgh, Pennsylvania 15213

Contact: aghose@stern.nyu.edu,  <http://orcid.org/0000-0002-6499-8944> (AG); panos@stern.nyu.edu,

 <http://orcid.org/0000-0002-2966-7402> (PGI); beibeili@andrew.cmu.edu,  <http://orcid.org/0000-0001-5466-7925> (BL)

Received: August 24, 2014

Revised: September 14, 2014; November 15, 2015; February 7, 2017; July 10, 2017

Accepted: September 22, 2017

Published Online in Articles in Advance:  
June 11, 2018

<https://doi.org/10.1287/mnsc.2017.2991>

Copyright: © 2018 INFORMS

**Abstract.** It is now well understood that social media plays an increasingly important role in consumers' decision making. However, an overload of social media content in product search engines can hinder consumers from efficiently seeking information. We propose a structural econometric model to understand consumers' preferences and costs on search engines to improve user experience under unstructured social media. Our model combines an optimal stopping framework with an individual-level random utility choice model and analyzes click behavior in conjunction with purchase choices. Our model accounts for three major constraints in a consumer's decision-making process: (1) interdependency in decision making for different alternatives, (2) sequential arrival of information revealed by click-throughs, and (3) nonnegligible search cost. Our approach allows us to jointly estimate consumers' heterogeneous preferences and search costs under the interplay of social media and search engines, and to predict search and purchase behavior for each consumer. We validate the model using an individual session-level data set of approximately seven million observations resulting in room bookings in 2,117 U.S. hotels. Interestingly, our analysis allows us to quantify the trade-off between consumers' benefits and cognitive costs from using large-scale unstructured social media information during decision making. Our policy experiments show that providing a carefully curated digest of social media content during the earlier stages of consumer search (i.e., on the search results summary page) can lead to a 12.01% increase in the overall search engine revenue.

**History:** Accepted by Anandhi Bharadwaj, information systems.

**Supplemental Material:** The online appendix is available at <https://doi.org/10.1287/mnsc.2017.2991>.

**Keywords:** consumer search • search cost • search engines • structural modeling • social media

## 1. Introduction

With the growing pervasiveness of social media, the volume and complexity of information product search engines need to access from their own platforms has been increasing rapidly. For example, websites such as Amazon.com, TripAdvisor.com, and Yelp.com can attract hundreds or even thousands of review postings that compete for users' attention. The onslaught of the exploding social media content can lead to significant information overload for consumers during product search. Such excess content can hinder consumers' efficiency in seeking information and making decisions (e.g., Iyengar and Lepper 2000). Even worse, it may discourage consumers from searching and cause unexpected termination of searches (e.g., session dropout).

During the past decade, product search engines have been trying to combine advanced techniques from information retrieval (e.g., Google Product Search) and recommender systems (e.g., Amazon.com) into their ranking design to improve the user search experience. Recent studies show that product search engines can

improve the ranking design and the user search experience based on the prediction of consumer preferences (e.g., Ghose et al. 2012, De los Santos and Koulayev 2017). Because consumers want the most desirable results early on, search engines can reorder the results by the predicted probabilities of consumer preferences (e.g., clicks or purchases).

Previous studies have examined how to estimate customer preferences based on online purchase information only (e.g., Ghose et al. 2012). However, consumer footprints on search engines provide us with a tremendous amount of information that reveals their preferences, even in the absence of purchases (e.g., Koulayev 2014, Kim et al. 2010, De los Santos and Koulayev 2017). When this search behavior is combined with purchases, the signals become even more comprehensive and useful. However, although many studies have worked on using either historical click-throughs or conversions separately to estimate consumer preferences, there is little work on jointly analyzing the search and purchase behavior to infer individual consumer

preferences and identify the products that satisfy most the user needs.

With the deluge of structured and unstructured social media content, consumers' cognitive costs in searching and evaluating product information become nonnegligible. As a result, search costs also play an important role in affecting consumers' choices in product search engines. Therefore, a major goal of our study is to better understand consumers' online footprints accounting for consumers' heterogeneous preferences and search costs, using both click and purchase information. However, this task can be challenging, because the cause of an observed search behavior by a consumer is hard to identify—e.g., the fact that a consumer prefers to click product A over product B may be because of a higher valuation for A, or because the consumer has incurred a lower search cost in searching for A than for B.

More generally, the challenge in predicting consumer choice with search cost is to simultaneously identify consumers' heterogeneous preferences and search costs (Hortacsu and Syverson 2004). A consumer may stop searching either because of a high valuation for the products already found or because of a high search cost. Either the preferences for product characteristics or the moments of the search cost distribution can explain the same observed search outcome (Koulayev 2014). Keeping the above in mind, another major goal of our study is to identify heterogeneous search costs under the social media context, examine their effect on consumer search behavior, and provide insights to product search engines on better design and management of social media content to improve user experience. The key identification strategy for consumer search cost in our study relies on the exclusion restriction that consumer preferences enter the decision-making processes of both search and purchase, whereas consumer search cost enters only the search decision-making process. Once the consideration set is generated after search, the conditional purchase decision should depend only on the consumer preferences. Our unique data set containing both consumer search data and purchase data allows us to identify these effects. In addition, we model search cost as a function of an exclusive set of variables. From an empirical identification perspective, we can simply view the search cost variables as additional product characteristics.

In summary, we propose a structural econometric model to understand consumers' preferences and search costs on product search engines to improve user experience under large-scale, unstructured social media. It combines an optimal stopping framework with an individual-level random utility choice model. It allows us to jointly estimate consumers' heterogeneous preferences and search costs. Based on the results, we are able to predict the probability that a

consumer clicks or purchases a certain product and provide a better understanding of what drives consumer engagement. Our analysis also allows us to quantify the trade-off between consumers' benefits and cognitive costs from using large-scale unstructured social media information during decision making. Our policy experiments offer insights to search engines on what product information they should display during different stages of consumer search (i.e., on the search result summary page versus product landing page), to improve user experience, click/purchase probabilities, and search engine revenues.

Our model is validated by a unique data set from the online hotel search industry. We have detailed individual consumer session-level search and transaction data from November 2008 through January 2009, containing approximately seven million observations resulting in room bookings in 2,117 hotels in the United States on Travelocity.com. Our model provides more precise measures of consumer price elasticity and heterogeneous preferences than does a static mixed logit model that does not account for consumer search cost or the sequence of the prior clicks. Our model also provides better predictive performance than does a click model that purely relies on the click information. More specifically, our model demonstrates the best performance in predicting the consumer click and purchase probabilities compared to other benchmark models. We see a 14.92% and an 18.77% improvement in the out-of-sample prediction using our model compared to the next-best-performing model, with respect to click-through and conversion probabilities, respectively.

Our policy experiments show that providing additional product information, especially the location-related information, on the travel search engine summary page will lead to a 22.16% increase in the overall search engine revenue. By contrast, although hiding all hotels' price information from the search summary page may lead to higher user "engagement" (when engagement is measured by number of clicks), it can hurt the travel search engine eventually by leading to a 7.08% drop in the overall search engine revenue. On the contrary, providing a carefully curated digest of social media textual content on the search results summary page can lead to a 12.01% increase in the overall search engine revenue. This finding suggests that it is important for product search engines to leverage the economic value of large-scale unstructured social media information, while in the meantime reducing the cognitive burden of consumers by automating the extraction of such information and presenting it to the consumers during the earlier stages of the decision-making process.

## 2. Prior Literature

Our paper draws from multiple streams of work. We summarize them in this section.

## 2.1. Search Cost and Consumer Information Search

First, our work builds on the literature on search cost and consumer information search. Recent studies have found that consumers have cognitive limitations, and search costs exist during the information search processes. Disregarding consumers' cognitive limitations and the limited nature of choice sets can lead to biased estimates of demand (e.g., Mehta et al. 2003, Kim et al. 2010, Brynjolfsson et al. 2010).

The existing literature in this field holds two different views of the nature of consumer search: non-sequential search and sequential search. The former strand of research follows Stigler's (1961) original model, assuming that consumers first sample a fixed number of alternatives and then choose the best from among them (e.g., Mehta et al. 2003, Honka 2014, Moraga-Gonzalez et al. 2017). By contrast, the other view, arising from the job-search literature (e.g., Mortensen 1970), argues the actual consumer search should follow a sequential model in which consumers keep searching until the marginal cost of an extra search exceeds the expected marginal benefit. Weitzman (1979), in single-agent scenarios, and Reinganum (1982), in multiagent scenarios, have laid theoretical foundations for sequential search models. In our paper, we assume consumers search sequentially on product search engines. This assumption is consistent with the mainstream research by the web search community (e.g., Chapelle and Zhang 2009). In addition, many recent studies in economics and marketing have also adopted the sequential search strategy for examining consumer search in an online environment (e.g., Kim et al. 2010, Koulayev 2014, Chen and Yao 2017).

With the growing interests and the recent development of information technologies that have made many intensive computation tasks more tractable today, empirical work to date has increased. Hong and Shum (2006) were the first to develop a structural methodology to recover search cost from price data only. Moraga-Gonzalez and Wildenbeest (2008) extend the approach of Hong and Shum to the oligopoly case and provide a maximum-likelihood estimate of the search cost distribution. Both papers focus on markets for homogeneous goods, using both sequential and nonsequential search models. Hortacsu and Syverson (2004) examine markets with differentiated goods and develop a sequential search model to recover search cost from the utility distribution. More recent empirical studies on nonsequential search tend to focus on the offline market with search frictions to study price dispersion (e.g., Wildenbeest 2011), endogenous choice sets and demand (e.g., Moraga-Gonzalez et al. 2017), or the identification of search cost from switching cost (Honka 2014). Recent empirical work on sequential

search examines consumers' limited search and the associated demand, with a focus on the online search market (Koulayev 2014, Kim et al. 2010). Meanwhile, De los Santos et al. (2012) use web browsing and purchasing behavior based on book-price distribution across 14 online bookstores to compare the extent to which consumers are searching under nonsequential and sequential search models.

One common practice in the existing empirical studies on both types of search models is that they typically model search cost as an inherent attribute of the consumer. Two exceptions are Kim et al. (2010), who model search cost as a function of the product's appearance frequency on Amazon.com, and Moraga-Gonzalez et al. (2017), who consider explanatory variables such as geographic distance from a consumer's home to different car dealerships. In our model, search cost is not only an inherent attribute of a consumer, but also a consequence of the social media context in which consumers of today are embedded. Note that, consistent with prior literature, the search cost in our study is modeled as exogenous to the consumer's search. By modeling consumer search cost as a random-coefficient function of the textual variables that are related to the unstructured social media content, we aim to examine the nature of search cost given the increasing interplay between product search engines and social media.

Finally, another related stream of consumer search literature has analyzed optimal search behavior when consumers are uncertain about the distribution of the product price or utility (e.g., Rothschild 1974, Rosenfield and Shapiro 1981, Bikhchandani and Sharma 1996, Koulayev 2013, De los Santos et al. 2017). For example, the recent work by De los Santos et al. (2017) has relaxed the assumption that consumers "know" the distribution of offerings while deciding on their search strategy, and allows for learning of the utility distribution. More specifically, consumers learn the utility distribution by Bayesian updating their Dirichlet process priors while sampling information about products and retailers. Our study is related to this stream of work in that we also consider the sequential arrival of information during different search stages, which allows for consumer update of the initial belief toward product utility via information search.

## 2.2. Search Engine Ranking and User-Generated Content

Our work is also related to the literature on search engine ranking. Examining the rank-position effect on the click-through rate and conversion rate on search engines has attracted a lot of attention. A number of recent studies focus on the context of search-engine-based keyword advertising and find significant empirical evidence on the rank-order effect (e.g., Ghose and Yang 2009, Goldfarb and Tucker 2011, Agarwal et al.

2011, Yao and Mela 2011). Other studies focus on search engine ranking for commercial products. For example, Baye et al. (2009) use a unique data set on clicks from one of Yahoo's price comparison sites to estimate the search engine ranking effect on clicks received by online retailers. Ellison and Ellison (2009) focus on the competition of retailers ranked on price search engines and find that the easy price search makes demand highly price sensitive for some products. Ghose et al. (2012) propose a utility-gain-based ranking (using data from past purchases, only, and not browsing behavior) that recommends products with the highest expected utility. The lab experiments indicate a strong preference for utility-based ranking compared to existing state-of-the-art alternatives. Ghose et al. (2014) combined a hierarchical Bayesian model and randomized user experiments to examine the search engine ranking and personalization effects from a causal perspective.

Finally, our work also relates to the stream of research on social media and user-generated content (e.g., Godes and Mayzlin 2004, Chevalier and Mayzlin 2006, Dellarocas et al. 2007, Duan et al. 2008, Forman et al. 2008). It especially builds on the recent research from a multidimensional view of the customer reviews (e.g., Archak et al. 2011, Ghose et al. 2012, Netzer et al. 2012, Chen et al. 2018). In this paper, we aim to examine the role of social media from multiple dimensions in affecting not only the product utility evaluation but also the search cost of consumers.

### 2.3. Comparisons with Recent Literature

Our model builds on Weitzman's (1979) optimal sequential search framework. To the best of our knowledge, five existing studies use similar methodologies to ours: Kim et al. (2010, 2014), Koulayev (2014), De los Santos and Koulayev (2017), and Chen and Yao (2017). However, our research differs from these studies in the following ways. (i) Our model incorporates not only consumers' search behavior, but also their purchases. Kim et al. (2010), De los Santos and Koulayev (2017), and Koulayev (2014) consider only consumers' search information as an approximation of their actual purchase decisions. (ii) Our observations include detailed click-throughs from each ranking position on a page, which allows us to precisely model the individual click probability for each product, rather than for a page with a bundle of products (i.e., a page of 15 hotels as in Koulayev 2014). More broadly speaking, Koulayev (2014) and our paper are complimentary: Koulayev models the costly process of discovering new hotels by flipping pages, but stops short of modeling what happens between click and booking. Our paper focuses on the second stage, starting from the costly click to the final booking. (iii) We conduct our analysis at the individual-consumer level as opposed

to at the aggregate market level (Kim et al. 2010, 2014). Such individual-level data allow us to leverage the detailed information of the *sequence of clicks* per session, rather than only the independent click-throughs. (iv) Chen and Yao (2017), De los Santos and Koulayev (2017), and Koulayev (2014) focus on constructing models that examine the joint use of search refinement tools (e.g., sorting) during consumer search. However, search refinement is not our focus in this paper. (v) Kim et al. (2010, 2014) and Chen and Yao (2017) assume a simpler information structure where consumers do not update their information set during search, whereas our paper allows for a more realistic information structure by allowing consumers to update their information set before and after click-through. (vi) Most importantly, our paper initiates a special focus on the interplay between consumer search and social media. Our goal is to use the structural econometric approach as a tool for analytics by product search engines to improve the user experience, especially under an overload of the unstructured social media content. We model the trade-off between the *value* and the *cognitive cost* associated with the large-scale unstructured social media information. We aim to examine how search engine policies regarding social media content, such as what information to show on the search summary page versus product landing page, may affect consumer search/purchase behaviors and search engine revenues.

In addition, in two recent papers, Ghose et al. (2012, 2014) also initialized their focus on the interplay of search engine and social media. This paper distinguishes from these two studies in the following ways. (i) Ghose et al. (2012) studied only the consumer purchase decisions, not search/click decisions, whereas this paper jointly studies the click and purchase decisions. (ii) Ghose et al. (2012, 2014) both focused on only the "benefit" of social media on consumer evaluation of product quality for the purchase decision, but did not consider the "cognitive cost" associated with processing social media information during consumer search. This is one major unique advantage of this paper. None of the previous work has studied the "cost" of social media content in affecting consumer search and purchase decisions on product search engines. (iii) Ghose et al. (2012, 2014) both used aggregated data on click/purchase share at product level, while this paper models consumer decision at individual level. (iv) From a methodology perspective, different from Ghose et al. (2012, 2014), this paper accounts for three unique constraints in the model: (1) interdependency in clicks/purchases among different products, (2) sequential arrival of information revealed by click-throughs, and (3) nonnegligible search costs.

A summary of the differences between this paper and the existing studies is shown in Table 1.

**Table 1.** Comparison with Recent Literature

	Kim et al. (2010)	Ghose et al. (2012)	Ghose et al. (2014)	Kim et al. (2014)	Koulayev (2014)	De los Santos and Koulayev (2017)	Chen and Yao (2017)	This paper
Data	Amazon, View-Rank, 18 months	Hotels, Purchase, ~8K observations, 3 months	Hotels, Click, Purchase, ~30K observations, 3 months	Amazon, View-Rank, Sale-Rank, 18 months	Hotels, Click, Search Refinement, 1 month, (Chicago)	Hotels, Click, Search Refinement, 1 month, (Chicago)	Hotels, Click, Search Refinement, Purchase, 215 sessions, 15 days	Hotels, Click, Purchase, ~7M observations, ~1M sessions, 2,117 hotels, 3 months, Individual clicks and purchases
Level of analysis	Market	Market	Market	Market	Individual clicks at page level	Individual clicks	Individual clicks and purchases	Individual clicks and purchases
Real transactions	No	Yes	Yes	No	No	No	Yes	Yes
Click sequence	No	No	No	No	Yes	Yes	Yes	Yes
Search refinements	No	No	Yes	No	Yes	Yes	Yes	No
Interplay with unstructured social media	No	Yes	No	No	No	No	No	Yes
Update of consumer information set	No	No	No	No	No	No	No	Yes
Major objectives	Consumer welfare, market structure	Design a novel consumer surplus-based ranking for search engine	Causal effect of search engine ranking and personalization	Market structure, innovation	Price sensitivity	Design search engine ranking by maximizing aggregate click-through rate	Search refinement → Welfare decrease	Interplay between search and social media, cognitive cost of unstructured social media

### 3. Data

**Clickstream and Transaction Data.** Our data set comes from Travelocity.com, a leading online travel search agency. The data set contains detailed information on session-level consumer search, click, and purchase events from November 2008 through January 2009, with a total of 7,059,122 observations from 969,033 individual sessions resulting in room bookings in 2,117 hotels in the United States.<sup>1</sup> A typical online session observed in our data set involves the following events: the initialization of the session, the search query, the hotel listings returned from that search query in a particular rank order, whether the consumer has used any special sorting criteria to rerank the hotels, clicks on any hotel listing, the login and actual transactions in a given hotel, and the termination of the session. We observe the hotel listings displayed to the consumer during the search session (regardless of whether any click occurs). If a click occurs, we observe hotel listings the consumer observed prior to that click. Moreover, we also observe the sequence of the clicks.

Notice that we also have detailed information associated with each event for every corresponding hotel, such as nightly room prices and the hotel's position in the set of listings returned by the search engine (i.e., "page" and "rank"). We have the detailed transaction-level information from Travelocity.com that is linked to the entire session-level consumer search data, including the final transaction price and the number of room units and nights purchased in each transaction. This information allows us to model consumer preferences for both the search and the purchase processes.

**Hotel General Information.** We collected hotel-related information from Travelocity.com, such as hotel class, hotel brand, number of amenities, number of rooms, reviewer rating, number of reviews, and the textual content of all of the reviews up to January 31, 2009 (the last date of transactions in our database).

**Hotel Location Information.** In addition, we have independently collected supplemental data on hotel location-related characteristics using automatic social geo-mapping techniques together with image data mining. We use geo-mapping search tools (in particular the Bing Maps API) and social geo-tags (from geonames.org) to identify the number of external amenities (e.g., shops, bars) in the area around the hotel. We use image classification methods together with human annotations (from Amazon Mechanical Turk, AMT) to extract whether a beach, lake, or downtown area is nearby, and whether the hotel is close to a highway or public transportation. We extract these characteristics from different zoom levels of the satellite images of a hotel location within a 0.25-, 0.5-, 1-, and 2-mile radius. We also collect local crime rates from FBI statistics.

**Hotel Service Quality Information Extracted from Social Media.** To fully exploit the information about hotel service quality, we combine text mining and sentiment analysis to examine the natural-language text of the customer reviews. For example, the helpfulness of the hotel staff is a service feature one can assess by reading the consumer opinions. Toward extracting such information, we build on the previous work of Archak et al. (2011) and Ghose et al. (2012). First, we extract the important hotel features. Following the automated approach introduced previously (Archak et al. 2011, Ghose et al. 2012), we use a part-of-speech tagger to identify the frequently mentioned nouns and noun phrases, which we consider candidate hotel features. We then use WordNet (Fellbaum 1998) and a context-sensitive hierarchical agglomerative clustering algorithm (Manning and Schütze 1999) to further cluster the identified nouns and noun phrases into clusters of similar nouns and noun phrases. The resulting set of clusters corresponds to the set of identified product features mentioned in the reviews. For our analysis, we kept the top six most frequently mentioned features, which were *hotel staff*, *food quality*, *bathroom*, *parking facilities*, *bedroom quality*, and *check-in/out front desk efficiency*.

For sentiment analysis, we extracted all of the evaluation phrases (adjectives and adverbs) that were used to evaluate the individual service features (for example, for the feature "hotel staff," we extracted phrases such as "helpful," "smiling," "rude," "responsive"). The process of extracting user evaluation phrases can be automated. To measure the meaning of these evaluation phrases, we used AMT to exogenously assign explicit polarity semantics to each word. To compute the scores, we used AMT to create our ontology, with the scores for each evaluation phrase. Our process for creating these "external" scores was done using the methodology of Archak et al. (2011). Finally, to handle the negation (e.g., "I didn't think the staff was helpful"), we built a dictionary database to store all of the negation words (e.g., "not," "hardly") using an approach similar to NegEx (<http://code.google.com/p/negex>; accessed September 10, 2015).

**Consumer Cognitive Cost Indicators Extracted from Social Media.** Although the textual content of customer reviews can reveal important information about hotel quality, there is a nonnegligible cognitive cost associated with processing such information. To capture consumers' cognitive costs in reading the user-generated reviews, we analyzed two sets of review text features that are likely to affect consumers' intellectual efforts in internalizing review content: "readability" (i.e., textual complexity, syllables, and spelling errors) and "subjectivity" (i.e., mean and standard deviation). Research has shown that both have had a significant impact on consumer online shopping behavior in the past (e.g., Ghose and Ipeirotis 2011). To derive the

probability of subjectivity in the review’s textual content, we apply text-mining techniques. In particular, we train a classifier using the hotel descriptions of each of the hotels in our data set as “objective” documents. We randomly retrieved 1,000 reviews to construct the “subjective” examples in the training set. We conduct the training process by using a 4-gram dynamic language model classifier provided by the LingPipe toolkit (<http://alias-i.com/lingpipe>; accessed September 10, 2015). Thus, we are able to acquire a subjectivity confidence score for each sentence in a review and then derive the mean and variance of this score, which represent the probability of the review being subjective.

In addition to review textual readability and subjectivity, we also extracted an additional cognitive cost indicator based on the topic complexity of the customer

reviews. In particular, built on prior literature (Gong et al. 2019), we analyzed the entropy value for the distribution of topics extracted from all customer reviews for each hotel (“Topic Entropy”). This entropy value measures the diversity of topics covered by the customer reviews for each hotel. Prior literature suggests the diversity in search results affects consumer search behavior (e.g., Weitzman 1979, Dellaert and Häubl 2012). In addition, consumer psychology theories suggest that as the information become noisier, users are more likely to abandon their search (e.g., Jacoby et al. 1974, Dhar and Simonson 2003), because users tend to get overwhelmed and discouraged by the complexity of information, and therefore lose their interest or trust in the search results. Therefore, we derived a topic entropy score using probabilistic topic models from machine learning and natural language processing to

**Table 2.** Definitions and Summary Statistics of Variables

Variable	Definition	Mean	Std. dev.	Min	Max
<i>PRICE_DISP</i>	Displayed price per room per night	230.98	179.76	16	2,849
<i>PRICE_TRANS</i>	Transaction price per room per night	148.08	108.18	52	2,252
<i>CLASS</i>	Hotel class	3.62	0.70	1	5
<i>AMENITYCNT</i>	Total number of hotel amenities	14.37	6.22	2	23
<i>ROOMS</i>	Total number of hotel rooms	210.12	258.27	12	2,900
<i>BRAND</i>	Dummies for 9 hotel brands: Accor, Best western, Cendant, Choice, Hilton, Hyatt, Intercontinental, Marriott, and Starwood	—	—	0	1
<i>PAGE</i>	Page number of the hotel	20.86	13.44	1	192
<i>RANK</i>	Screen position of the hotel	12.09	4.32	1	25
<i>SPECIALSORT</i>	Dummy for a special sorting method	0.10	0.30	0	1
<i>BEACH</i>	Beachfront within 0.6 miles	0.19	0.36	0	1
<i>LAKE</i>	Lake or river within 0.6 miles	0.23	0.44	0	1
<i>TRANS</i>	Public transportation within 0.6 miles	0.31	0.45	0	1
<i>HIGHWAY</i>	Highway exits within 0.6 miles	0.70	0.42	0	1
<i>DOWNTOWN</i>	Downtown area within 0.6 miles	0.66	0.45	0	1
<i>EXTAMENITY</i>	Number of external amenities within 1 mile, i.e., restaurants, shopping malls, or bars	4.63	7.99	0	27
<i>CRIME</i>	City annual crime rate	194.99	127.22	3	1,310
Social media variables (cognitive cost)					
<i>COMPLEXITY</i>	Average sentence length per review	17.50	3.77	4	44
<i>SYLLABLES</i>	Average number of syllables per review	246.81	50.53	76	700
<i>SPELLERR</i>	Average number of spelling errors per review	1.17	0.33	0	3.86
<i>SUB</i>	Review subjectivity—Mean	0.91	0.03	0.05	1
<i>SUBDEV</i>	Review subjectivity—Standard deviation	0.02	0.03	0	0.25
<i>TOPICENTROPY</i>	Entropy score to measure topic complexity	2.88	0.13	1.58	2.99
Social media variables (hotel quality)					
<i>REVIEWCNT</i>	Total number of reviews	13.56	25.60	0	202
<i>RATING</i>	Overall reviewer rating	3.94	0.39	1	5
<i>STAFF</i>	Sentiment score for helpfulness of staff	0.35	0.62	−3	3
<i>FOOD</i>	Sentiment score for food quality	0.69	0.66	−3	3
<i>BATHROOM</i>	Sentiment score for bathroom quality	0.42	0.74	−3	3
<i>PARKING</i>	Sentiment score for parking facilities	0.16	0.58	−3	3
<i>BEDROOM</i>	Sentiment score for bedroom quality	0.49	0.86	−3	3
<i>FRONTDESK</i>	Sentiment score for check-in/out front desk efficiency	0.54	0.55	−3	3
Model computed search cost (in U.S. dollars)					
$c_j$	Search cost for a hotel $j$ derived from the model estimation	6.18	0.38	3.43	7.75
Total number of sessions:	969,033	Total number of hotels:		2,117	
Total number of observations:	7,059,122				
Time period: November 1, 2008—January 31, 2009					



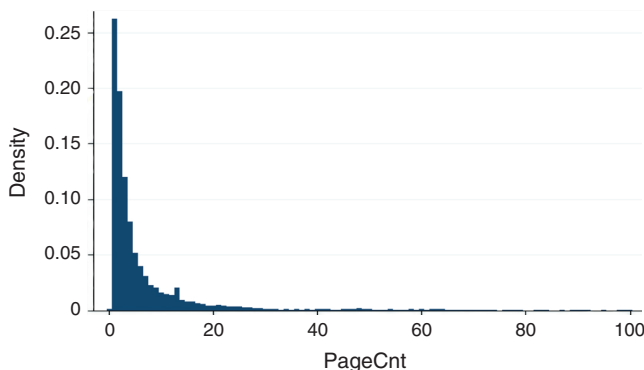
capture the “noisiness” of information provided by the customer reviews. Topic models are unsupervised algorithms that aim to extract hidden topics from unstructured text data. In particular, we measure the topic complexity of reviews for each product by estimating a topic model using a latent Dirichlet allocation model (LDA) (Blei et al. 2003) and subsequently computing the entropy (i.e., diversity) of the topic distribution of reviews for that product. We provide more technical details on the topic modeling in Online Appendix E.

For a better understanding of the variables, we present the definitions and summary statistics of all variables in Table 2. Note that the data set in this paper uses not only the transaction data (i.e., purchases), but also the complete session-level data (i.e., both clicks and purchases). The resulting data set contains approximately seven million observations from one million individual user sessions.

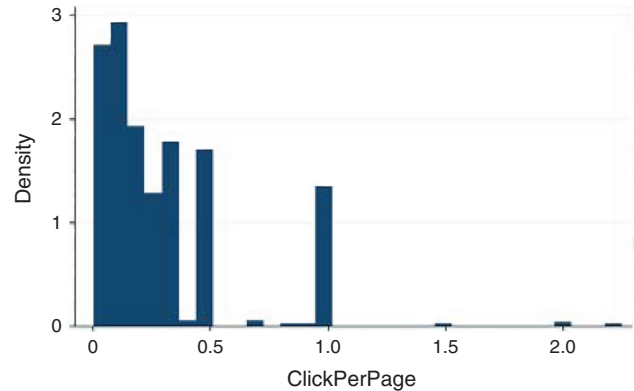
### 3.1. Model-Free Evidence of Limited Search by Consumers

Before we describe our model, we seek from the data suggestive evidence that could motivate our assumption of consumers’ limited search. First, we plot the distribution of the total number of pages a consumer browses in a search session. Figure 1(a) illustrates this distribution in detail, with the  $x$  axis representing the page counts and the  $y$  axis representing the density. We notice that over 25% of consumers browse only one page, over 50% of consumers browse less than three pages, and less than 10% of consumers browse more than 15 pages during their search for hotels. This finding is consistent with prior industry evidence that consumers seldom search more than three pages (e.g., Sterling 2008). Second, we further look into the distribution of the average number of click-throughs made per page during each search session. Figure 1(b) illustrates this distribution, with the  $x$  axis representing the click-throughs per page and the  $y$  axis representing the density. We find that on average, consumers

**Figure 1(a).** (Color online) Distribution of Number of Pages Browsed (Session Level)



**Figure 1(b).** (Color online) Distribution of Number of Click-Throughs per Page (Session Level)



click less than one hotel (out of a total of 25 hotels) per page during their search. In fact, a large majority of consumers click less than 0.5 hotels per page, on average. Besides, over 97% clicks occurred on the first page. These two figures provide us preliminary evidence that consumers incur nontrivial search costs and that consumer search is limited.<sup>2</sup>

## 4. A Structural Model of Consumer Sequential Search

Our data set contains the complete information on the browsing session (e.g., list of hotels displayed, sequence of clicks) and the purchasing decisions that consumers made. Consumers have three options for a hotel during a search session: (A) do not click on the hotel at all, (B) click on the hotel but do not purchase it, and (C) click on the hotel and also purchase it. To identify option A from options B and C, we need to model consumers’ click decision making. To identify option B from option C, we need to model consumers’ purchase decision making. As a key contribution of this analytical study, we build a holistic model of user behavior that models both the clicking and purchasing behavior. Our model, in summary, works as follows.

### Before Click:

1. A consumer session starts with consumer browsing hotels on the search results summary page. A consumer can obtain any hotel information provided on the search results summary page (with no clicks needed) at zero cost.

2. Before clicking on a hotel, the consumer does not observe the exact information shown on the “details” landing page for that hotel. Instead, she forms a belief about what information would appear on the landing page, conditional on the information observed in the search results summary page. Because no click is needed to form the belief, we assume the consumer incurs zero cost at this step.

3. Given the observed information on the search results summary page and the conditional belief of the unobserved information on the landing page, the consumer is able to infer the *expected utility* of each hotel before the click-through at zero cost.

4. Meanwhile, before clicking on a hotel, the consumer also forms a belief about what the *expected search cost* would be if she were to click on the hotel (e.g., due to the additional cognitive efforts needed for processing the unstructured information on the landing page), conditional on the information observed from the search results summary page. Again, no click is required to form the belief of search cost, and we assume the consumer incurs zero cost at this step.

#### After Click:

1. The consumer session continues with a series of clicks, where the consumer decides to click on the landing pages of some hotels and to find out the exact utilities from these hotels. The goal of search (i.e., via click-through) is to reveal any uncertainty in the utility (i.e., uncertainty in the landing-page characteristics as well as the unobserved error). The set of clicked hotels and the order of the clicks reveal information about the preferences and search costs of the user.

2. The consideration set is being generated during the search process. It contains all of the hotels the consumer has clicked. After the costly click-through, the consumer knows the *actual* utilities (rather than the expected utility) of the clicked hotels, which form the consideration set.

3. The consumer stops searching new hotels (and hence stops clicking) when the expected marginal benefit of doing so is less than the expected search cost. We adopt the concept of “reservation utility” from Weitzman (1979) to define when the consumer stops searching. The decision of whether to continue searching or to stop relies on the actual utilities of the hotels in the consideration set at that moment<sup>3</sup> and her expected utilities and expected search costs of the upcoming hotels.

4. Once the consumer stops searching, the consideration set is fixed. Based on the final consideration set, the consumer makes a purchase decision (or skips purchasing anything at all).

#### 4.1. Model Setting

**Product Utility.** Assume the utility of hotel  $j$  for consumer  $i$  to be a random-coefficient model as follows:

$$u_{ij} = V_{ij}^S + V_{ij}^L + e_{ij}, \quad (1)$$

where  $V_{ij} = V_{ij}^S + V_{ij}^L$  represents the hotel utility from the hotel characteristics displayed on the website. It consists of two conceptual components: (i) a deterministic component,  $V_{ij}^S$ , the exact utility from “summary-page” hotel characteristics consumers can directly observe on

the search summary page, and (ii) a stochastic component, the *additional utility*,  $V_{ij}^L$ , from “landing-page” hotel characteristics consumers cannot directly observe before the click-through but can observe after the click-through. To evaluate the overall expected utility before the click-through, a consumer  $i$  forms a belief of the distribution of the unobserved landing-page utility  $f(V_{ij}^L)$  based on  $V_{ij}^S$ . This belief comes from the consumer’s knowledge about the utility distribution for hotel  $j$  conditional on the observed summary-page characteristics for this hotel. The consumer makes the click decision based on the exact value of the summary-page utility  $V_{ij}^S$  and the *expected* value of the landing-page utility  $E(V_{ij}^L)$ . Once the consumer decides to click on the hotel, the click-through will reveal the actual value of the landing-page characteristics, and the consumer updates the expected value  $E(V_{ij}^L)$  with the actual value  $V_{ij}^L$ . Moreover, we let  $e_{ij}$  represent the unobserved uncertainty in the consumer’s evaluation. The consumer does not know the realization of  $e_{ij}$  unless she clicks on hotel  $j$  and visits its landing page. In particular, we assume  $e_{ij}$  to be i.i.d. across consumers and hotels, and to follow a type I extreme value distribution  $e_{ij} \sim \text{Type I EV}(0, 1)$ .<sup>4</sup>

In summary, our utility setting assumes that the consumer does not know the full realization of the utility of hotel  $j$  before the click-through. However, the consumer knows the distribution of the utility. This assumption is critical and is consistent with Weitzman (1979) and many recent studies that have examined consumers’ sequential search behavior in the online search contexts (e.g., Kim et al. 2010, Chen and Yao 2017, Koulayev 2014). Hence, the goal of search (i.e., click) is to solve the uncertainty in the consumer’s evaluation toward both the landing-page characteristics and the unobserved error to reveal the true utility of a hotel.

More specifically, let  $X_j$  be a vector of summary-page characteristics for hotel  $j$ . Let  $P_j$  represent the *Price* for hotel  $j$  that is also directly available to consumers on the search results summary page. Thus, we can model the summary-page utility as  $V_{ij}^S = X_j \beta_i - \alpha_i P_j$ , where  $\beta_i$  and  $\alpha_i$  are consumer-specific parameters capturing the heterogeneous preferences of consumers. We assume  $\beta_i \sim N(\bar{\beta}, \Sigma_\beta)$ , where  $\bar{\beta}$  is a vector containing the means of the random effects and  $\Sigma_\beta$  is a diagonal matrix containing the variances of the random effects. Similarly, we assume  $\alpha_i \sim N(\bar{\alpha}, \sigma_\alpha^2)$ .

Meanwhile, we model the expected value of the preclick stochastic part of the utility as  $E(V_{ij}^L) = \tilde{L}_j \lambda_i$ , where  $\tilde{L}_j$  represents the consumer expectation toward the landing-page characteristics for hotel  $j$  conditional on the observed summary-page characteristics  $(X_j, P_j)$ . Note that  $\tilde{L}_j$  may not equal the actual values of the landing-page characteristics. We use the tilde sign to

distinguish  $\tilde{L}_j$  from the realization of its deterministic value,  $L_j$ . Using a similar approach proposed by Koulayev (2014), we approximate  $\tilde{L}_j$  by taking the mean of the bootstrap samples from the actual information of the landing pages of the hotels that present the same summary-page characteristics. This approach allows consumers to infer knowledge about the utility distribution of a hotel based on the average knowledge from the population with similar experience (i.e., who are also exposed to  $(X_j, P_j)$ ). The consumer estimates the expected utility of the landing page based on  $\tilde{L}_j$ . She updates  $\tilde{L}_j$  with the deterministic value  $L_j$  only after she chooses to click on hotel  $j$  and reveals the actual deterministic values of the landing-page characteristics. Let  $\lambda_i$  represent consumer-specific parameter to capture the heterogeneity. Consistent with previous assumptions, we assume it follows a normal distribution  $\lambda_i \sim N(\bar{\lambda}, \Sigma_\lambda)$ .

Therefore, we have the overall utility function as follows. Before the click-through, the expected utility from hotel  $j$  for consumer  $i$  is

$$u_{ij} = X_j\beta_i - \alpha_i P_j + \tilde{L}_j\lambda_i + e_{ij}. \quad (2a)$$

After the click-through, the realization of the actual utility becomes

$$u_{ij} = X_j\beta_i - \alpha_i P_j + L_j\lambda_i + e_{ij}. \quad (2b)$$

**Search Cost.** We model a consumer's search cost as a result of the landing-page-evaluation behavior associated with a click (i.e., cognitive cost of processing additional unstructured information). More specifically, let  $Q_j$  denote the set of actual cognitive cost variables for evaluating the landing-page unstructured information of hotel  $j$ . We model the actual search cost of consumer  $i$  after clicking on hotel  $j$  to follow a log-normal distribution<sup>5</sup>:

$$c_{ij} = \exp(Q_j\gamma_i), \quad (3a)$$

where  $\gamma_i \sim N(\bar{\gamma}, \Sigma_\gamma)$ ,  $\bar{\gamma}$  is a vector containing the means of the random effects and  $\Sigma_\gamma$  is a diagonal matrix containing the variances of the random effects. To model the consumer's cognitive cost of evaluating the unstructured information on the landing page, we consider different dimensions in the cognitive-cost variables  $Q_j$ , including both the *readability* and the *subjectivity* of the textual content of online reviews.

However, because the landing-page information is not directly observable to the consumer before click, to decide whether to click on a hotel, the consumer needs to form a belief of her expected search cost conditional on the observed summary-page characteristics of that hotel. This means that in our model,  $Q_j$  is not directly observable to the consumer before the click-through. Similarly, the consumer forms an expectation based

on the observed summary-page characteristics. Let  $\tilde{Q}_j$  capture the consumer's expectation toward the unobserved cognitive-cost variables of hotel  $j$ . We approximate this expectation value by taking the mean of the bootstrap samples from the actual information of the hotels with the same summary-page characteristics.

Based on the discussion above, we can write the (expected) search cost of consumer  $i$  for hotel  $j$  before the click-through as the following:

$$c_{ij} = \exp(\tilde{Q}_j\gamma_i). \quad (3b)$$

Thus, before the click-through of hotel  $j$ , a consumer  $i$  makes the click decision based on the expected search cost for  $j$ .

Note that a consumer's search cost is a sunk cost. It enters only the consumer's click decision process but not the purchase decision process. Once the consumer forms an evaluation about the expected search cost, she will make a click decision based on this evaluation one time and will not need it again in the future. Therefore, the realized actual search cost after click-through in Equation (3a) does not enter either the click model or the purchase model in reality. Only the expected search cost before click-through in Equation (3b) will enter the model estimation process (i.e., click model). Hence, we can treat the consumer's expected search cost as a deterministic value in modeling her search (click) decision, which is consistent with Weitzman (1979). For simplicity of notation, we therefore keep the same notation  $c_{ij}$  to denote the expected search cost, although the expected search cost in Equation (3b) represents an expectation value (based on  $\tilde{Q}_j$ , not  $Q_j$ ).

## 4.2. Problem Description and the Optimal Search Framework

In general, our consumer search problem can be described as follows. Assume a consumer searches sequentially (i.e., examines alternatives one by one) to find a hotel. At each stage of the search, the consumer has two options: continue to search for the next alternative or stop and purchase the current best alternative (including purchasing nothing, i.e., an outside good). Consider that the consumer is forward looking. This situation implies that at any stage during her search, she always tries to choose an action that maximizes her *expected utility from the current stage going forward*—meaning she tries to maximize the marginal benefits from both the current stage and all potential future stages. Therefore, the key problem here is to determine the optimal point for the consumer to choose the “stop” option.

More formally, let  $S_i$  be the current search-generated consideration set (i.e., including all hotels consumer  $i$  has clicked). Let  $u_i^*$  denote the current highest value obtained by consumer  $i$  so far. We define

$$u_i^* = \max_{j \in S_i} \{u_{ij}, 0\}. \quad (4)$$

Note that we define  $u_i^*$  as the highest value  $u_{ij}$  consumer  $i$  obtains from the hotels in her consideration set. Given the current best value  $u_i^*$ , the expected marginal benefits for consumer  $i$  from searching  $j$  are

$$B_{ij}(u_i^*) = \int_{u_i^*}^{\infty} (u_{ij} - u_i^*) f_i(u_{ij}) du_{ij}, \quad (5)$$

where  $f_i(\cdot)$  is the probability density function of hotel utility  $u_{ij}$  and is individual specific. The expected marginal benefits  $B_{ij}(u_i^*)$  represent the expectation of the utility for hotel  $j$ , given that it is higher than  $u_i^*$ , multiplied by the probability that  $u_{ij}$  exceeds  $u_i^*$ . As we notice, the benefits of search depend only on the distribution of utility above  $u_i^*$ . Thus, for any hotel  $j$ , the reservation utility  $z_{ij}$  meets the following boundary condition, where the *expected search cost* equals the *expected marginal benefits* from searching the hotel:

$$c_{ij} = B_{ij}(z_{ij}) = \int_{z_{ij}}^{\infty} (u_{ij} - z_{ij}) f_i(u_{ij}) du_{ij}. \quad (6)$$

Note that in Equations (4)–(6), because the actual search cost and actual utility for an upcoming unsearched hotel  $j$  are not observable to consumer  $i$  before the click-through, her decision of whether to click on hotel  $j$  is based on her *expected search cost* and *expected utility*. By contrast,  $u_i^*$  is derived based on the *actual hotel utilities* because after the click-through, the consumer can observe the exact information about each hotel in her consideration set. Thus, when consumer  $i$ 's current best value is equal to the reservation utility of hotel  $j$ ,  $u_i^* = z_{ij}$ , she is indifferent between searching for  $j$  or stopping (and accepting  $u_i^*$ ). Consumer  $i$  will continue to search for hotel  $j$  if her current best value is lower than the reservation utility of hotel  $j$ ,  $u_i^* < z_{ij}$ , and she will stop otherwise. More details on the derivation of the optimal search strategy and the reservation utility are provided in Online Appendices B and C.

### 4.3. Click Probability

We define the click probability in a fashion similar to (Kim et al. 2010). Let  $r(j)$  denote the hotel with the  $j$ th highest-ranked reservation utility  $z_{i,r(j)}$ . Let  $\pi_{i,r(j)}$  be the probability that consumer  $i$  will click hotel  $r(j)$ . This probability equals the probability that the current highest value  $u_i^*$  within the consumer's current consideration set is lower than the reservation utility of hotel  $r(j)$ . Let  $S_{i,r(j)}$  be the current consideration set generated prior to hotel  $r(j)$ . It includes all hotels the consumer has clicked before hotel  $r(j)$ . For a consumer to click hotel  $r(j)$ ,  $z_{i,r(j)}$  has to exceed the maximum value from the clicked sets of hotels. Thus, we model the click probability of hotel  $r(j)$  for consumer  $i$  as

$$\begin{aligned} \pi_{i,r(j)} &= \Pr[r(j) \text{ is clicked by consumer } i] \\ &= \Pr[u_i^* < z_{i,r(j)}] \end{aligned}$$

$$\begin{aligned} &= \Pr \left[ \max_{m \in S_{i,r(j)}} (V_{i,r(m)}^S + V_{i,r(m)}^L + e_{i,r(m)}) < z_{i,r(j)} \right] \\ &= \prod_{m \in S_{i,r(j)}} F_{e_i}(z_{i,r(j)} - V_{i,r(m)}^S - V_{i,r(m)}^L), \quad j > 1, \quad (7) \end{aligned}$$

where  $F_{e_i}(\cdot)$  is the CDF of  $e_{ij}$ , which in our case  $e_{ij} \sim \text{Type I EV}(0, 1)$ .<sup>6</sup>

### 4.4. Conditional Purchase Probability

Conditional on the sequence of clicks consumer  $i$  has made in the search session, we can derive the conditional probability that she purchases hotel  $r(j)$  in her consideration set as the following:

$$\begin{aligned} \eta_{i,r(j)} &= \Pr(r(j) \text{ is booked by consumer } i | \text{all} \\ &\quad \text{clicks by consumer } i) \\ &= \Pr(u_{i,r(j)} \geq u_{i,r(j')} | \text{all clicks by consumer } i, \\ &\quad \forall r(j) \neq r(j'), r(j), r(j') \in S_i) \\ &= \Pr(V_{i,r(j)}^S + V_{i,r(j)}^L + e_{i,r(j)} \geq V_{i,r(j')}^S + V_{i,r(j')}^L \\ &\quad + e_{i,r(j')} | \text{all clicks by consumer } i, \\ &\quad \forall r(j) \neq r(j'), r(j), r(j') \in S_i), \quad (8) \end{aligned}$$

where  $S_i$  is the click-generated consideration set for consumer  $i$ . Note that because the consideration set  $S_i$  is selected by consumer  $i$  based on her search decisions,  $e_{ij}$  does not follow a full type I EV distribution. Instead, it follows a truncated type I EV distribution based on the optimality conditions used by the consumer. Unfortunately, under such circumstances the conditional choice probability does not have a close-form expression (e.g., Logit form). To address this issue, we applied a simulation approach. Similar methods have been adopted by the previous studies (Chen and Yao 2017, Honka 2014, McFadden 1989). McFadden (1989) proposed a kernel-smoothed frequency simulator to sample the random draws from a truncated type I EV distribution by smoothing the probabilities using a multivariate scaled logistic CDF (Gumbel 1961). Honka (2014) applied McFadden's approach to sample the error term from a truncated Type I EV distribution by accounting for the composition of the click-generated consideration set and the utility optimality of the final choice to model consumer simultaneous search. Chen and Yao (2017) applied a similar simulation approach to sample the error term from a truncated normal distribution by further accounting for not only the choice set composition and the utility optimality of the final choice, but also the sequence of the click-generated consideration set to model consumer sequential search. Our simulation approach builds on the methods from Chen and Yao (2017) and Honka (2014). It allows us to simulate the error term from a truncated Type I EV distribution by satisfying the follow three optimality conditions: (1) sequence of the click-generated consideration set, (2) composition

of the click-generated consideration set, and (3) utility optimality of the final choice. We provide the full details on how we use the simulated method to construct the conditional purchase probability in Online Appendix D.

#### 4.5. Likelihood Function

We model the overall likelihood as the product of the probabilities of all of the observed consumer clicks and purchases:

$$\begin{aligned} \text{Likelihood} &= \prod_i \Pr(\text{CLICK}_i, \text{PURCHASE}_i) \\ &= \prod_i \Pr(\text{CLICK}_i) \Pr(\text{PURCHASE}_i | \text{CLICK}_i), \quad (9) \end{aligned}$$

where  $\text{PURCHASE}_i$  represents the observed purchase event by consumer  $i$ , and  $\text{CLICK}_i$  represents the observed sequence of all click events by consumer  $i$ .

We can then model  $\Pr(\text{CLICK}_i)$  and  $\Pr(\text{PURCHASE}_i | \text{CLICK}_i)$  as follows. First, let  $N$  be the total number of hotels that consumer  $i$  has clicked (i.e., size of the consideration set) and  $J$  be the total number of hotels available in the market. We can model the joint probability of the sequence of click events for consumer  $i$  as the following:

$$\begin{aligned} \Pr(\text{CLICK}_i) &= \Pr[\text{click}_{i,r(1)}, \text{then click}_{i,r(2)}, \dots, \text{then click}_{i,r(N)}, \\ &\quad \text{then all\_unclicks}_i] \\ &= \prod_{r(j) \in S_i^{\text{clicked}}} \Pr(u_{i,n} \leq z_{i,r(j)}, \forall n \in S_i^{\text{clicked\_before\_}r(j)}) \\ &\quad \cdot \prod_{r(m) \in S_i^{\text{unclicked}}} \Pr(u_{i,N}^* > z_{i,r(m)}) \\ &= \prod_{r(n) \in S_i^{\text{clicked}}} \Pr(u_{i,n-1}^* < z_{i,r(n)}) \\ &\quad \cdot \prod_{r(m) \in S_i^{\text{unclicked}}} [1 - \Pr(u_{i,N}^* < z_{i,r(m)})] \\ &= \prod_{r(n) \in S_i^{\text{clicked}}} \pi_{i,r(n)} \prod_{r(m) \in S_i^{\text{unclicked}}} [1 - \pi_{i,r(m)}], \quad (10) \end{aligned}$$

where  $S_i^{\text{clicked}}$  represents the set of all hotels that have been clicked by consumer  $i$ ,  $S_i^{\text{unclicked}}$  represents the set of all hotels that have not been clicked by consumer  $i$ , and  $S_i^{\text{clicked\_before\_}r(j)}$  represents the set of hotels that have been clicked by consumer  $i$  before  $r(j)$ .

Second, conditional on the sequence of click events, we can derive the conditional probability of the purchase event from Equation (8). Again,  $S_i$  is the click-generated consideration set for consumer  $i$ .

$$\begin{aligned} \Pr(\text{PURCHASE}_i | \text{CLICK}_i) &= \Pr(u_{i,r(j)} \geq u_{i,r(j')} | \text{CLICK}_i, \forall r(j) \neq r(j'), r(j), r(j') \in S_i) \\ &= \eta_{i,r(j)}. \quad (11) \end{aligned}$$

Finally, based on Equations (10) and (11), we can rewrite the likelihood function as follows:

*Likelihood*

$$= \prod_i \left\{ \eta_{i,r(j)} \prod_{r(n) \in S_i^{\text{clicked}}} \pi_{i,r(n)} \prod_{r(m) \in S_i^{\text{unclicked}}} [1 - \pi_{i,r(m)}] \right\}. \quad (12)$$

With this model setting, we are able to account for the fact that the decision-making processes for the hotels in the same session are not completely independent from each other. Instead, the click and purchase decisions for a hotel depend not only on its own utility, but also on the prior sequence of clicks associated with the consideration set.<sup>7</sup>

#### 4.6. Estimation

To model the utility of a hotel, we consider  $X$  to contain all hotel characteristics that are directly available on the search summary page, including *Hotel Class*, *Hotel Brand*, *Customer Rating*, *Total Review Count*, *Page*, and *Rank*. We consider  $L$  to contain all additional characteristics that can only be revealed from the hotel landing page, including *Amenity Count*, *Number of Rooms*, *Number of External Amenities*, the top-6 service characteristics extracted from the social media textual content including *hotel staff*, *food quality*, *bathroom*, *parking facilities*, *bedroom quality* and *check-in/out front desk efficiency*, as well as locational factors such as *Beach*, *Lake*, *Downtown*, *Highway*, *Public Transportation*, and *Crime Rate*.

To analyze consumers' search costs, we consider  $Q$  to contain different factors that capture the cognitive cost of the unstructured hotel information on the landing page. In particular, we consider both the *Readability* (i.e., complexity, syllables, and spelling errors) and *Subjectivity* (i.e., mean and standard deviation of the linguistic subjectivity) of the textual content of online reviews.<sup>8</sup>

Note that the website also provides a sorting mechanism for consumers to refine their search by sorting the results under criteria other than the default sorting algorithm. Technically, if a consumer chooses to customize the sorting algorithm, her search cost for each hotel in the ranking list may also change, hence becoming endogenous to her own search behavior. However, in reality, we find that with approximately one million online search sessions in our data set, more than 90% of these sessions do not involve any customized sorting behavior at all. This finding is consistent with previous randomized experimental results (Ghose et al. 2014) that the majority of users tend to stick with the default sorting method during online product search. Therefore, for simplicity in this study, we focus on only those search sessions conducted under the default sorting algorithm.<sup>9</sup>

To estimate our model, we derive the overall log-likelihood function as the following:

$$LL(\theta) = \sum_i \left\{ \ln(\eta_{i,r(j)}) + \sum_{r(n) \in S_i^{\text{clicked}}}^{1..N} \ln(\pi_{i,r(n)}) + \sum_{r(m) \in S_i^{\text{unclicked}}}^{N+1..J} \ln(1 - \pi_{i,r(m)}) \right\}, \quad (13)$$

where  $\theta$  represents the set of parameters of the random coefficients we aim to estimate:

$$\{\theta\} = \left\{ (\bar{\alpha}, \sigma_\alpha), \left( \bar{\beta}, \sum_\beta \right), \left( \bar{\lambda}, \sum_\lambda \right), \left( \bar{\gamma}, \sum_\gamma \right) \right\}.$$

We iteratively estimate the model using a maximum simulated likelihood (MSL) method. In particular, we apply the Monte Carlo method for numerical simulation, where for each individual observation, we simulate 250 random draws from the joint distribution of the individual heterogeneous parameters  $\{\theta\}$  and compute the corresponding individual-level click probability  $\pi_{i,j}$  and conditional purchase probability  $\eta_{i,j}$ . To maximize the log-likelihood function  $LL(\theta)$ , we use a non-derivative-based optimization algorithm (i.e., the Nelder–Mead simplex method) for heuristic search.<sup>10</sup> This procedure iteratively searches for the optimal set of parameters  $\{\theta^*\}$  until the log-likelihood function is maximized:

$$\{\theta^*\} = \arg \min_{\{\theta^*\}} LL(\theta). \quad (14)$$

The main computational complexity of the estimation comes from the calculation of the reservation values. During each iteration of the optimization algorithm, for each observation and each value of the search cost, we need to solve  $z_{ij} = B_{ij}^{-1}(c_{ij})$  numerically. To improve the estimation efficiency, we apply an interpolation-based method to compute the reservation values (Kim et al. 2010, Koulayev 2014). We provide more details of this computation procedure in Online Appendix C.

#### 4.7. Identification

One of the major challenges is to simultaneously identify consumers' heterogeneous preferences and search costs. A person may stop searching either because she has a high valuation for the products already found or because she has a high search cost. Therefore, either the preferences for product characteristics or the moments of the search cost distribution can explain an observed search outcome. In our study, we need to identify four major effects: consumer preferences (mean and heterogeneity) and consumer search cost (mean and heterogeneity). The key identification strategy of our estimation relies on the exclusion restriction

that consumer preferences enter the decision-making processes of both search and purchase, whereas consumer search cost enters only the search decision-making process. Once the consideration set is generated after search, the conditional purchase decision should depend only on the consumer preferences. Our unique data set containing both consumer search data and purchase data allows us to identify these effects.

Moreover, the identification also relies on search-cost “shifters.” Recall that we model search cost as a function of a completely different set of variables compared to the consumer-preferences variables. When we choose different sets of covariates for search cost and consumer preferences, the covariates enter search cost function but not utility function serve as the exclusion restrictions for identification. Conditioned on the exclusion restrictions, the utility and the search cost can be separated. From an empirical identification perspective, we can simply view the search-cost variables as additional product characteristics (i.e., similar to Kim et al. 2010). Thus, we can identify the search-cost and consumer-preferences variables simultaneously.<sup>11</sup> We provide more detailed discussions below.

**Mean Effects.** We identify the mean effects of consumer preferences variables based on the correlation between the observed click/purchase frequencies and the frequencies of underlying preferences variables. In other words, we measure the mean effect of a consumer preference variable by how often the same (or similar) variable appears in the hotels consumers click or purchase. For example, if on average people tend to click (or purchase) low price hotels, we may conclude that people have a high price sensitivity. This identification is similar to the one in most traditional choice models, except that it takes into consideration not only the observed purchases, but also the clicks, to infer the mean effect of consumer preferences.

We identify the mean effect of search cost partially based on the observed average size of the consumer's search-generated consideration set. Importantly, note that we model the search cost as a function of completely different variables compared to the consumer-preferences variables, which can be viewed simply as additional hotel characteristics. Thus, similar to the identification of consumer mean preferences, we can identify the mean search cost coefficients based on the correlation between the observed click frequencies and the frequencies of underlying search cost characteristics.

**Heterogeneous Effects.** Note that across both purchase data and search data, we have multiple observations per consumer. For a given consumer and her search cost, we observe the deviation of observed purchase and searches from those predicted decisions

based on the mean preferences and search cost parameters. The distribution of these deviations across individual consumers allows us to identify the heterogeneity distribution parameters.

More specifically, we identify consumer heterogeneous preferences from two perspectives. First, we partially identify them from the search data based on the distribution of the deviations across individual consumers between our model's predicted click probabilities (based solely on the mean effects) and individual consumers' observed click probabilities. Second, since we also observe individual consumers' final purchases, the purchase data allow us to identify the heterogeneous preferences based on the distribution of the deviations across individual consumers between the model's predicted purchase probabilities (based solely on the mean effects) and individual consumers' observed purchase probabilities.

We identify the heterogeneous search cost through two sources. First, our identification relies on the exclusion restriction that search cost variables do not enter purchase decision processes. After identifying the consumer heterogeneous preferences through the conditional purchase probabilities, we can then identify the heterogeneous search cost by the joint variation of the consideration set size and the click probabilities. In particular, at each point during a consumer's search, based on mean parameters, her reservation utility, and the products already searched in the consideration set, we can predict the mean probability of her stopping the search. The deviation of her search activities from the predicted values give us the information of one's heterogeneity in search cost. The distribution of these deviations across individual consumers identifies the search cost heterogeneity distribution parameters. Second, the nonlinear functional form in the reservation utility (i.e., Equation (6)) can also help identify consumer preference and search cost parameters (Kim et al. 2010). Since the consumer preferences enter the equation in a nonlinear manner (i.e., need to integrate over the utility), whereas the search cost enters the equation in a linear manner, this mathematical nonlinearity also helps us separately identify consumer heterogeneous preferences and search cost.

## 5. Empirical Results

### 5.1. Main Results

Our main results are shown in Table 3. First, we find that the majority of the coefficients are statistically significant at the  $p \leq 5\%$  level, including both the mean effects  $(\bar{\alpha}, \bar{\beta}, \bar{\lambda}, \bar{\gamma})$  and the heterogeneity parameters  $(\sigma_{\alpha}, \Sigma_{\beta}, \Sigma_{\lambda}, \Sigma_{\gamma})$ , for price, summary hotel characteristics, landing page hotel characteristics, and cost of absorbing social media content, respectively). Consistent with theory, *PRICE* has a negative effect on

hotel demand. *CLASS*, *AMENITYCNT*, *ROOMS*, *RATING*, and *REVIEWCNT* have a positive effect on hotel demand. For hotel-location characteristics, we find that *BEACH*, *TRANS*, *HIGHWAY*, and *DOWNTOWN* have a positive effect on hotel demand, whereas *LAKE* and *CRIME* show a negative effect. Consistent with prior literature, online position has a significant effect on consumer click and demand (e.g., Yao and Mela 2011, Ghose and Yang 2009). In particular, *PAGE* and *RANK* each leads to a decrease in the hotel demand. Moreover, we find that three service variables that are extracted from social media textual content demonstrate significant effect on hotel demand. In particular, food quality presents the highest positive impact, followed by hotel staff and parking.

On the other hand, we find that the additional unstructured information from the landing page indeed leads to an increase in consumer search cost. In particular, the readability-related review features such as *COMPLEXITY*, *SYLLABLES*, and *SPELLERR* have a positive sign, suggesting that long and complex sentences, words with many syllables, or spelling errors in user reviews discourage consumers from continuing to search on product search engines. Moreover, *SUB* has a positive sign, implying that highly subjective and opinionated content that lacks objective information creates a cognitive burden for consumers during hotel search and may lead to early termination of their search. Finally, *SUBDEV* also has a positive sign, which suggests that a mixture of both objective and subjective messages is likely to lead to higher cognitive costs. In other words, *SUBDEV* represents the standard deviation of the subjectivity value and it captures the level of heterogeneity in the type of information provided in the reviews. The higher the heterogeneity, the higher the cognitive cost associated with processing such information (i.e., when a review is a mix of both subjective and objective messages, it adds to the cognitive costs because readers might have to incur additional effort when switching between different types of information).

To get a handle on the actual magnitude of the search cost, we quantitatively derive the dollar value of different search cost variables. This value represents how much a certain variable effect can be translated into price. We find that, on average, the effort of continuing to search an additional hotel costs \$6.18. The search costs differ across hotels from \$3.43 to \$7.75. Our findings are consistent with previous findings suggesting a nontrivial search cost in online markets. For example, Koulayev (2014) found that the page-level median search costs rise from \$4 per first search to \$16 per fifth on a travel search engine. Brynjolfsson et al. (2010) found that the benefits from searching lower screens equal \$6.55 for the median consumer.

**Table 3.** Estimation Results—Main Model

Variable (Preferences)	Mean effect (Std. err.) <sup>M</sup> $\bar{\alpha}, \bar{\beta}, \bar{\lambda}$	Heterogeneity (Std. err.) <sup>M</sup> $\sigma_{\alpha}, \Sigma_{\beta}, \Sigma_{\lambda}$	Variable (Preferences)	Mean effect (Std. err.) <sup>M</sup> $\bar{\alpha}, \bar{\beta}, \bar{\lambda}$	Heterogeneity (Std. err.) <sup>M</sup> $\sigma_{\alpha}, \Sigma_{\beta}, \Sigma_{\lambda}$
<i>PRICE</i> <sup>(L)</sup>	-1.252* (0.022)	0.417* (0.074)	<i>DOWNTOWN</i>	1.198* (0.061)	0.471* (0.093)
<i>PAGE</i>	-0.239* (0.003)	0.080 (0.133)	<i>CRIME</i>	-0.173* (0.043)	0.015 (0.034)
<i>RANK</i>	-0.314* (0.008)	0.132* (0.067)	<i>RATING</i>	2.661* (0.015)	1.308* (0.091)
<i>CLASS</i>	1.516* (0.023)	0.935* (0.181)	<i>REVIEWCNT</i> <sup>(L)</sup>	1.230* (0.107)	0.369* (0.069)
<i>AMENITYCNT</i> <sup>(L)</sup>	0.146* (0.034)	0.066 (0.070)	<i>STAFF</i>	0.139* (0.027)	0.034 (0.088)
<i>ROOMS</i> <sup>(L)</sup>	0.394* (0.024)	0.195 (0.287)	<i>FOOD</i>	0.225* (0.038)	0.136* (0.002)
<i>EXTAMENITY</i> <sup>(L)</sup>	0.165* (0.036)	0.041 (0.046)	<i>BATHROOM</i>	0.290 (0.271)	0.060 (0.103)
<i>BEACH</i>	1.539* (0.028)	0.561* (0.099)	<i>PARKING</i>	0.097* (0.008)	0.075* (0.011)
<i>LAKE</i>	-0.663* (0.116)	1.560* (0.389)	<i>BEDROOM</i>	-0.175 (0.232)	0.253 (0.269)
<i>TRANS</i>	1.336* (0.140)	0.192* (0.064)	<i>FRONTDESK</i>	0.065 (0.103)	0.021 (0.076)
<i>HIGHWAY</i>	0.447* (0.093)	0.068 (0.061)			
<i>BRAND</i>		Yes			
(Search cost)	$\bar{\gamma}$	$\Sigma_{\gamma}$	(Search cost)	$\bar{\gamma}$	$\Sigma_{\gamma}$
<i>Search Base Cost (Constant)</i>	-7.511* (0.089)	0.971* (0.176)	<i>SPELLERR</i> <sup>(L)</sup>	0.329* (0.082)	0.033 (0.101)
<i>COMPLEXITY</i>	0.541* (0.094)	0.398* (0.115)	<i>SUB</i>	0.196* (0.045)	0.057 (0.229)
<i>SYLLABLES</i> <sup>(L)</sup>	0.678* (0.115)	0.721* (0.106)	<i>SUBDEV</i>	0.342* (0.056)	0.119 (0.273)
<i>Maximum LL</i>		-405,418			
<i>Price Elasticity</i>		-1.619			

Note. L, logarithm of the variable; M, main model.  
 \*Statistically significant at 5%.

Hann and Terwiesch (2003) quantified rebidding costs to be \$4.00–\$7.50 in a reverse-auction channel. Hong and Shum (2006) found consumers’ median search costs to be \$1.31–\$2.90 for a sample of text books. In addition, De los Santos (2008) found search costs ranging from \$0.90 to \$1.80 per search in the online book industry. Meanwhile, a one-word increase in the average sentence length increases consumer search cost by \$0.44. One more syllable or one more spelling error per review can cost consumers \$0.56 or \$0.28, respectively, during the product search.

Importantly, our empirical analysis allows us to quantify the trade-off between consumers’ benefits and costs toward leveraging social media information for decision making. Our results indicate that more social media information (especially textual content) may not always improve consumer decision making. Certain service- and quality-related information extracted from review textual content can indeed facilitate consumer decision making and impact product demand. However, because of the size and the unstructured nature of such information, it also brings in nonnegligible cognitive costs to the consumers. Our study aims to explore a more effective and scalable way of managing social media information that can help search engines extract and provide useful information to consumers without introducing high cognitive costs. Moreover, our model and policy experiments (in Section 6) allow us to evaluate the associate economic outcome on consumers as well as on product search engine revenues.

To further analyze the robustness of our model performance, and how social media and consumer heterogeneity (e.g., travel purposes) may affect the search cost and decisions of a consumer, we conduct three robustness tests by (1) excluding the social media variables from the main model, (2) including additional topic entropy variable into the main model, and (3) adding interaction effects between consumer travel purposes and summary-page variables. We find that the estimated coefficients are qualitatively consistent with the main results. Interestingly, we notice that the model that does not account for social media textual variables presents significantly higher price elasticity. This result indicates that the unstructured social media information plays an important role in consumer decision making, and that consumers’ cognitive costs to digest such information are nonnegligible. Without accounting for such unstructured information during consumer search can lead to an overestimation of price elasticity. We provide more details on the robustness tests in Online Appendix F.

### 5.2. Model Comparisons

Furthermore, to understand how the type and scale of data or modeling mechanisms may affect the performance of our analysis, we conducted model comparison analyses with a set of alternative benchmark models using different data sets or modeling mechanisms.

**5.2.1. Alternative Models.** In particular, we considered four alternative benchmark models. (1) *Alternative Model I:* Use the purchase data only (Mixed Logit



Model). (2) *Alternative Model II*: Use the purchase data only (Mixed Logit Model + Additional Search Cost Variables). (3) *Alternative Model III*: Use the click data only (Click Model).<sup>12</sup> (4) *Alternative Model IV*: Use both the click and the purchase data (Joint Probabilistic Model of Click and Purchase + Additional Search Cost Variables, But No Click Sequence Information).<sup>13</sup> Due to space limitation, we provide the details on the alternative model mechanisms in Online Appendix G.

Overall, we find that the estimation results are qualitatively consistent with our main findings. Interestingly, we find that using a static model without accounting for consumers' search behavior can lead to an overestimation of the price elasticity. The interpretation of this finding can be attributed to the nature of the hotel search market. A model that captures consumers' actual search behaviors finds lower price elasticity, implying that consumers in the hotel search market tend to highly evaluate the quality of hotels and put weight on nonprice factors during search (e.g., class, amenities, or reviews). Our finding on price elasticity is consistent with prior findings by Koulayev (2014) and Brynjolfsson et al. (2010). Both studies show that when consumers face a highly differentiated market (e.g., product differentiation or retailer differentiation), they are more likely to focus on nonprice factors during search. Hence, the estimated price elasticity is lower when incorporating consumers' search behaviors into the model. On the contrary, when a market is less differentiated, consumers become more price sensitive and focus more on price search. Thus, a search model that incorporates consumers' search behaviors may find a higher price elasticity of demand than a static model (e.g., De los Santos et al. 2012).

Furthermore, from Alternative Models (I)–(III), we find that use of only the click data or only the purchase data is likely to overestimate the price elasticity, and therefore it is important to consider both click and purchase decisions when modeling consumer preferences. However, interestingly, from Alternative Model (IV), we find that although incorporating both click and purchase decisions information can improve the model estimation, the joint probabilistic model without considering the click sequence information can still lead to an overestimation of price elasticity. This result indicates that the final click or purchase decisions are not all that matter, but that the sequential click path is critical in revealing consumer preferences. Failing to capture consumers' search paths can lead to an overestimation of price elasticity in the online search market. For more details on the alternative model results, we illustrate them in Tables G1 and G2 in Online Appendix G.

**5.2.2. Model Prediction Experiments.** Based on the model-estimated coefficients, our final goal is to predict the click and purchase probabilities for a hotel by

an individual consumer. The prediction of the two individual probabilities can be achieved by substituting the model-estimated coefficients into the Equations (7) and (8). To obtain individual-level consumer heterogeneity, we apply the Monte Carlo simulation method. In particular, we use the same random draws we simulated previously (i.e., in Section 4.6) from the joint distribution of the individual heterogeneous parameters. Based on the steps above, we are able to compute the corresponding individual click and purchase probabilities for each hotel for an individual consumer.

To examine the predictive performance of our model, we conduct a set of model-prediction experiments. We first compute the predicted individual click and purchase probabilities for each hotel as described above. Then, we compare the predicted individual click and purchase probabilities with the observed click and purchase probabilities (i.e., observed search and choice shares for the hotels). We calculate the prediction error for each hotel at individual-session level for both click and purchase probabilities. Then, we compute the root mean square error (RMSE) and mean absolute deviation (MAD). We consider all of the four alternative models discussed above as our baseline models. Furthermore, we are interested in examining how the use of unstructured data (social media textual variables) may affect the model's predictive power. Therefore, we consider a fifth baseline model: main model without the social media textual variables (*Robustness Test 1*) for both click- and purchase-probability predictions.

We randomly partition our data set into two subsets: one with 70% of the total observations as the estimation sample and the other with 30% of the total observations as the holdout sample. To minimize any potential bias from the partition process, we perform a 10-fold cross-validation. We conduct both in-sample and out-of-sample estimation using our model and the two baseline models. We then compare the predictive performance of both the click and the purchase probabilities of a hotel. The prediction results are illustrated in Tables A.1(a) and A.1(b) (click probability) and Tables A.2(a) and A.2(b) (purchase probability) in Appendix A. Our model-prediction results demonstrate that our model has the overall highest predictive power. Our model outperforms the baseline models in both in-sample and out-of-sample predictive power for both click and purchase predictions. Similar trends in improvement in the predictive power occur with respect to RMSE and MAD.

For example, with regard to the click-probability prediction, the out-of-sample results in Table A.1(b) show that with respect to the RMSE, our proposed model can improve the prediction performance by 14.92% compared to the search model without the social media textual variables. It can improve the prediction performance by 33.20% compared to the click model with

only click data. We find a similar trend with regard to the purchase-probability prediction. For example, the out-of-sample results in Table A.2(b) demonstrate that with respect to the RMSE, our main model can improve the prediction performance by 18.77% compared to the next-best model—search model without the social media textual variables. It can improve the prediction performance by 37.36% compared to the mixed logit model with a limited consideration set, and by 32.81% compared to the mixed logit model with a limited consideration set plus the additional search cost variables.

Notice that the model-prediction experiments indicate that our model is better able to predict the individual click and purchase probabilities for each hotel than the click model and the static mixed logit model. Even after considering various extensions of the mixed logit models accounting for the limited consideration set and the additional search cost variables, or considering other alternative behavioral models, our search model still provides the best predictive performance. The potential reasons are the following. Our proposed search model is a holistic model that captures both the click and the purchase decision-making processes for a consumer. Therefore, our model is able to account for the following three unique features of consumer search. (1) *Interdependency in decision making*: the search model predicts that a click decision depends on the ordered list of previously clicked products, and consequently, a purchase decision depends on the previous click-generated consideration set; however, static models assume independent decision making during consumer evaluation. (2) *Information arrives sequentially*: our model assumes that detailed product landing-page attributes can only become available to a consumer after she clicks on the product; however, static models tend to ignore the fact that information arrives sequentially and assume both landing- and summary-page attributes are available to the consumer at the beginning. (3) *Nonnegligible search cost*: the search model predicts that a click decision depends on the (expected) search cost associated with this click, and the formation of the final consideration set depends on the search cost toward each product; however, the static model ignores such opportunity cost.

In summary, three major indications from our model comparison experiments are as follows. (i) Both the click and the purchase data reveal significant information about consumer preferences and search cost, and both are critical to improving the model predictive power. (ii) The *sequence* of the clicks reveals significant information about consumer preferences and search cost. Our main search model incorporates not only the click decisions, but also the sequential order of these clicks. However, the static mixed logit models ignore the sequence of the clicks and simply

take the final consideration set as exogenously given. (iii) Unstructured social media data play an important role in consumer decision making. Incorporating such information into the model can lead to a significant improvement in the model's predictive power.

## 6. Policy Experiment

Based on our model estimation results, we conduct counterfactual analyses under various policy experiments to explore the what-if type of questions. More specifically, considering the amount and type of different information, we are interested in what information product search engines should present during different stages of consumer search (i.e., on the search results summary page versus product landing page).

### 6.1. Information Shown on the Search Summary Page vs. Landing Page

As we notice in our data set, most online products contain a large number of characteristics. However, due to the limitation in screen space, search engines are unable to show all product information on the search summary page. Instead, search engines choose to highlight a snapshot of some product information on the summary page, while leaving the majority of information to the landing page. The information selected for the search summary page for a product becomes critical because it can influence both consumers' perceptions of the utility of the product and their expectations regarding the search costs associated with further evaluation of the product (i.e., via click-through).<sup>14</sup>

To explore what information should be shown on the search summary page, we conduct a policy experiment using our model. In particular, we assume that search engines show different sets of hotel characteristics on the summary page—meaning that these chosen characteristics are directly observable to consumers before the click-through. We reestimate consumers' conditional belief regarding the unobserved characteristics using bootstrap samples and then compute the individual session-hotel-level predicted click and purchase probabilities based on the parameter estimates from the original model estimation. We compute the overall click and purchase probabilities for a hotel by taking an average of the click and purchase probabilities across all sessions for that hotel. Finally, we sum over all hotels based on the prices and the predicted purchase probabilities to compute the predicted search engine revenue.

By doing so, we aim to examine the following question: Holding consumers' preferences for product characteristics and search cost variables consistent, how would consumers' click and purchase behavior change if the search engine websites were to provide different sets of information on the search results summary page? Moreover, we are interested in exploring a better strategy for search engines to design the search

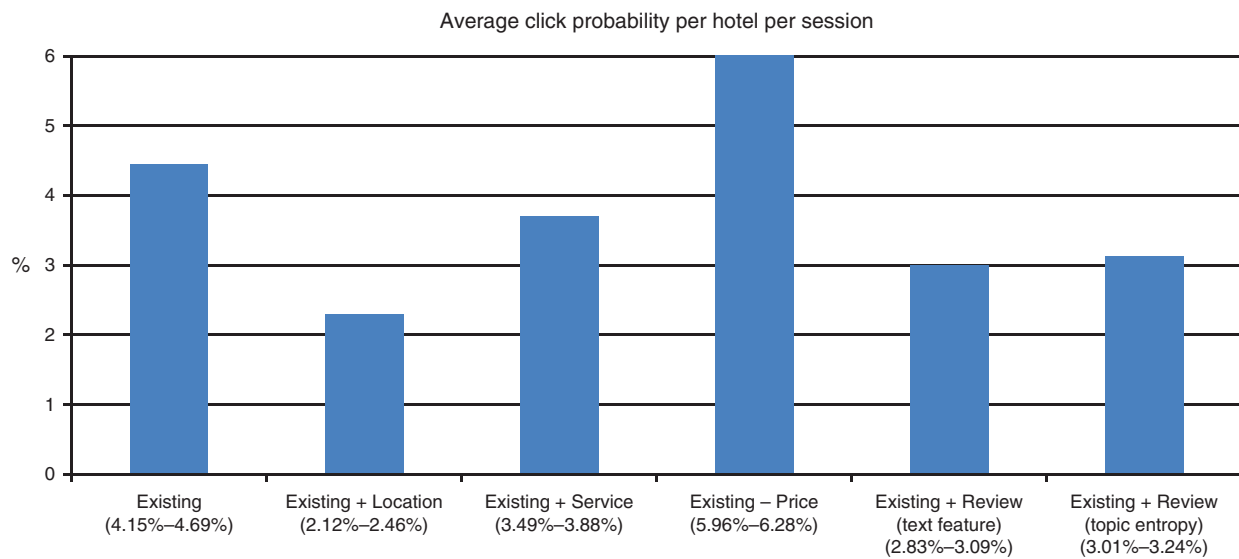
results summary page such that it improves the overall click/purchase probabilities and the search engine revenue.

More specifically, we focus on six alternative sets of product information that may be potentially useful to show on the search results summary page: (1) existing summary-page characteristics (i.e., price, hotel class, hotel brand, customer rating, review count, page, rank), (2) existing summary-page characteristics plus additional location-related characteristics (i.e., number of external amenities, beach, lake, downtown, highway, public transportation, crime

rate), (3) existing summary-page characteristics plus additional service-related information (i.e., amenity count), (4) existing summary-page characteristics plus additional review-text-related information (i.e., textual review features), (5) existing summary-page characteristics plus additional review-topic-related information (i.e., topic entropy score derived from the entropy measurement), and (6) existing summary-page characteristics minus the product price information.

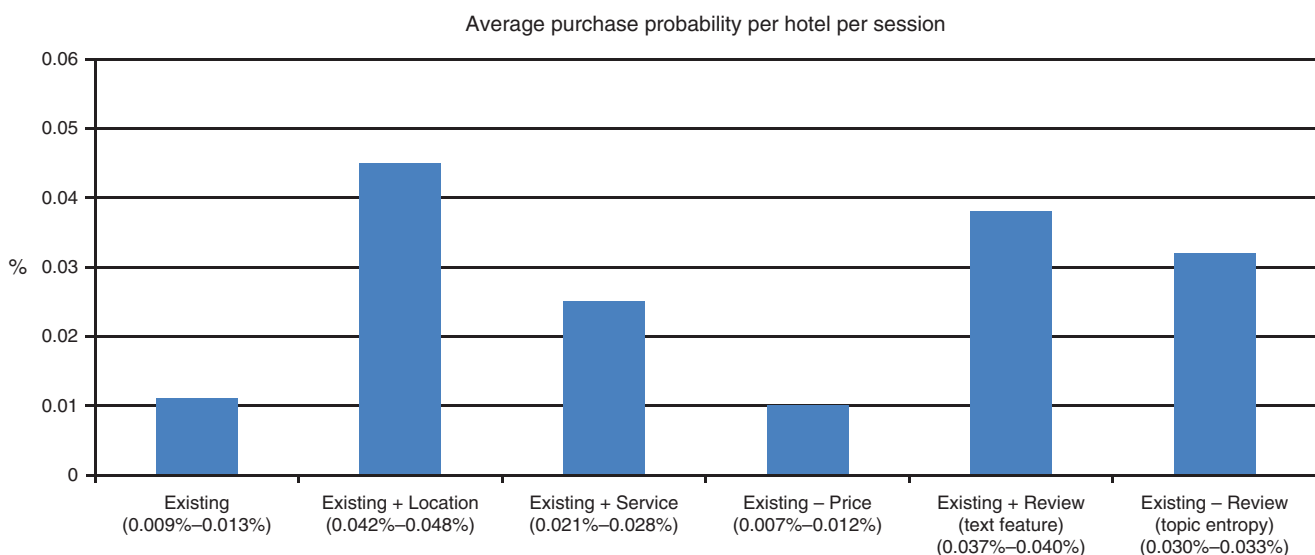
We compute the average predicted click and purchase probabilities per hotel per session under each of the above six assumptions. We provide our results

**Figure 2(a).** (Color online) Predicted Click Probabilities with Different Information Provided on Search Summary Page



Note. Confidence intervals in parentheses.

**Figure 2(b).** (Color online) Predicted Purchase Probabilities with Different Information Provided on Search Summary Page



Note. Confidence intervals in parentheses.

in Figures 2(a) and 2(b). Our findings demonstrate that the type of information search engines choose to show on the summary page has a statistically significant effect on consumers' click and purchase probabilities. In particular, we find that providing additional location-, service-, or review-related information for products on the search results summary page will lead to a significant decrease in click probability. This finding is intuitive. Because providing more information on the summary page will reduce the variance of the product utility (i.e., reducing uncertainty in consumer expectation) before click, it lowers the reservation utility of the product (hence making the product less attractive for a consumer to click).

However, interestingly, we find that providing additional product information, especially the location-related information, on the travel search engine summary page will lead to a significant increase in the purchase probability. A potential reason for this finding is that providing additional product information on the search summary page can reduce the potential error in consumers' expectation toward product utility and search costs before click. As a consequence, consumers are more likely to click on the best set of products that will provide them the highest utility. Hence, the maximum utility discovered from this click-generated consideration set is more likely to exceed the utility of the outside good. As a result, consumers are less likely to miss a good-value deal (i.e., leave without purchase).

Meanwhile, we find that although excluding price information from the search summary page can lead to a significantly higher click probability, it does not seem to increase the purchase probability at the end. This finding indicates that strategically hiding price information (i.e., price obfuscation) from the search summary page can make further searching (i.e., clicking) for products on a search engine more attractive. However, this strategy may not increase the overall purchase probability.

Finally, we compute the overall search engine revenue based on the hotel prices and the predicted purchase probabilities. Our results show that the location-related information is the most influential, compared to the service- and review-related information, when the travel search engine presents this information on the search summary page. It can lead to a 22.16% increase in the overall search engine revenue. Providing service-related information, such as the total number of hotel amenities, on the search summary page can lead to a 3.22% increase in the overall search engine revenue. By contrast, strategically hiding price information from the search summary page can hurt the search engine revenue, leading to a 7.08% drop in the overall revenue. We provide more details on the corresponding results in Table 4.<sup>15</sup>

**Table 4.** Predicted Overall Search Engine Revenue with Different Information on Search Summary Page

	Overall search engine revenue (\$)	95% confidence interval <sup>a</sup> (\$)
<i>Existing</i>	452,781	445,263–458,260
<i>Existing + Location Information</i>	553,136	538,026–561,989
<i>Existing + Service Information</i>	467,369	460,031–474,112
<i>Existing – Price Information</i>	420,132	411,585–429,203
<i>Existing + Review Information (Text Features)</i>	507,160	500,327–514,278
<i>Existing + Review Information (Topic Entropy)</i>	490,063	481,314–498,157

<sup>a</sup>Confidence interval is calculated based on bootstrapping the policy simulation experiments 200 times.

Interestingly, providing a carefully curated digest of social media textual content on product summary page (e.g., the six most frequently mentioned product features extracted from the customer reviews, customers' attitudes toward these popular features, the readability of the review's textual content) can lead to a 12.01% increase in the overall search engine revenue. Meanwhile, providing an overall "topic entropy" score of the review content (i.e., derived from topic models to measure the complexity of the review topic content) can lead to an 8.23% increase in the overall search engine revenue. These findings suggest that it is important for product search engines to leverage the economic value of large-scale unstructured social media information, while at the same time reducing the cognitive burden of consumers by automating the extraction of such information and providing it to consumers during the earlier stages of decision making.

Furthermore, to examine where the revenue increase came from, we conducted an additional analysis on the breakdown of the revenue in the simulation. Interestingly, we found that the revenue increase came from both existing consumers and expansion of market coverage. In addition, we also found that the revenue increase occurred for both existing hotels and new hotels. This finding provides further supports that with carefully designed information on search summary page, search engine can improve the market coverage of consumers as well as the diversity of products consumed, which can lead to a potential increase in consumer surplus. For more details, we provide the complete revenue breakdown analysis in Online Appendix H.

In sum, our policy experiment offers critical insights on the potential of analyzing large historical user behavioral data for search engines to improve the landing-page design strategy for better user experience and higher overall business revenues.

## 7. Managerial Implications and Conclusion

In this paper, we propose a structural econometric model for product search engines to understand consumers' search and purchase behavior as well as to quantify the search costs incurred by consumers. Our model combines an optimal stopping framework with an individual-level random utility choice model. It allows us to jointly estimate consumers' heterogeneous preferences and search costs in a product search engine context where unstructured social media information is quite pervasive, and to identify the key driver of a consumer's decision at each stage of the search and purchase process. Our final results suggest that both the historical clicking decisions and the purchase decisions reveal significant information of consumer preferences and search costs. Moreover, the paths of searches (i.e., sequence of clicks) also reveal significant information of consumer preferences and search costs. Our analyses can help search engines predict consumer online footprints and design the search results summary page to improve user experience and search engine revenues.

On a broader note, our research makes two key contributions. *First*, we show the advantage of incorporating multiple and large-scale data sources in analyzing how humans search, evaluate information, and make decisions under cognitive constraints in response to the emerging interplay between social media and search engines. Moreover, we are able to quantify the effects of unstructured social media content on user search cost. Our empirical analysis aims to provide an approach on which future studies can build, with the goal of exploring the potential of "Big Data" and sophisticated customer analytics tools for managerial decision making. *Second*, we demonstrate the value of using digital analytics by search engines based on structural econometric methods in finding solutions for important business problems. Our structural model of consumer search combines the optimal stopping framework with an individual-level random utility choice model. It allows us to harness the advantage of *multistage* consumer behavioral data on search engines to identify the drivers of consumer decisions in electronic markets. It enables the prediction of consumers' future search behavior on search engines. Moreover, it offers insights to search engines on the design of the search results summary page (i.e., what information to show on the summary page versus the landing page) to improve the user experience and the search engine revenues. Importantly, this approach can be generalized to any electronic market with an in-house search engine (e.g., Amazon.com), especially in a mobile search environment (e.g., Apple's iTunes or App store), given the commonality in the goal of improving user experience.

Our work has several limitations, some of which can serve as fruitful areas for future research. First, our

model assumes that the consumer knows the general distribution of utilities of alternatives, and each alternative follows the same distribution. However, when the alternatives are sorted on search engines under certain criteria including the default method, they are presented in the order of their predicted attractiveness to a consumer. Such recommendations can alter the distribution of the expected utilities of alternatives and may induce a shift in consumers' decision making (Dellaert and Häubl 2012). Examining this fact from an empirical perspective would be interesting. Second, testing other alternative consumer behavioral models would be interesting. For example, instead of searching sequentially, consumers may search in a nonsequential fashion by first choosing a fixed size of consideration set (e.g., Honka 2014). Comparing the differences in the corresponding model prediction of consumer search strategy would be interesting. Third, in this study, we assume each online consumer session to be an independent search process. Due to the data limitation, we cannot identify the possibility that a consumer may leave a session without booking but come back at a later time to resume the search. In this case, we treat these searches as two separate results in our estimation. Distinguishing such repeated searchers and more precisely estimating the search costs would be an interesting avenue for future research. Meanwhile, due to the data limitation, we do not have the consumer-level demographic information. Because the search cost is likely to relate to the opportunity cost of time, including such information (e.g., age, income) in the future would be useful. Finally, it would be very interesting for future research to consider the supply side (e.g., how the hotels/advertisers may respond to the search engine's policy change) in addition to the demand side to examine the effects of policy change on search engines.

## Acknowledgments

Author names are in alphabetical order. The authors thank Travelocity for providing data and support for this research. The authors thank the department editor, the anonymous associate editor, and three anonymous reviewers for their constructive comments. For valuable feedback and suggestions, the authors thank Ravi Bapna, Pradeep Chintagunta, Daria Dzyabura, Alok Gutpa, Elisabeth Honka, Sam Hui, Sergei Koulayev, Ramayya Krishnan, Michael D. Smith, Rahul Telang, Russell S. Winer, and audiences at the 2012 International Conference on Information Systems, the 2012 Workshop for Information Systems and Economics, the 2012 INFORMS Annual Conference, the 2014 Marketing Science Conference, Carnegie Mellon University, the Stern School of Business at New York University, and the Carlson School of Management at University of Minnesota. An earlier conference version of this paper won the Best Theme Paper Award at the International Conference on Information Systems in 2012.

## Appendix A. Model Comparisons

**Table A.1(a).** In-Sample Model Prediction Results (Click Probability)

	Main model	Main model without social media textual variables	Click model (with only click data)	Joint model of click and purchase (no click sequence)
RMSE	0.0514	0.0588	0.0627	0.0613
MAD	0.0197	0.0221	0.0278	0.0262

**Table A.1(b).** Out-of-Sample Model Prediction Results (Click Probability)

	Main model	Main model without social media textual variables	Click model (with only click data)	Joint model of click and purchase (no click sequence)
RMSE	0.1163	0.1367	0.1741	0.1541
MAD	0.0526	0.0614	0.0712	0.0658

**Table A.2(a).** In-Sample Model Prediction Results (Purchase Probability)

	Main model	Main model without social media textual variables	Mixed logit model		Joint model of click and purchase (no click sequence)
			(Limited consideration set)	(Limited consideration set + Additional search cost variables)	
RMSE	0.0833	0.0912	0.1107	0.1074	0.0942
MAD	0.0274	0.0292	0.0392	0.0359	0.0312

**Table A.2(b).** Out-of-Sample Model Prediction Results (Purchase Probability)

	Main model	Main model without social media textual variables	Mixed logit model		Joint model of click and purchase (no click sequence)
			(Limited consideration set)	(Limited consideration set + Additional search cost variables)	
RMSE	0.1251	0.1540	0.1997	0.1862	0.1662
MAD	0.0670	0.0729	0.0951	0.0845	0.0819

### Endnotes

<sup>1</sup> In our data set, 2,117 hotels had at least one booking during the data collection period. A total of 13,546 hotels had at least one display in consumer search sessions.

<sup>2</sup> For some cities, the number of hotels might be small and therefore no additional page is available for searching. We find that 56 out of 4,845 cities (approximately 1.15%) in our data have less than 25 hotels (which means only one page is available for searching). After excluding these small cities, the model-free evidence shows a similar trend that consumer search is highly limited.

<sup>3</sup> In particular, the decision of whether to continue searching or to stop depends on the actual utility of the hotel with the maximum utility in the consideration set.

<sup>4</sup> Note that different from Kim et al. (2010), who assume standard normal distribution of the error, we allow for logit distribution of the error term in our model, as we assume that the consumers may optimize their utility over unobserved (to the econometrician) variables. In our estimation, we transform the logit error into standard normal disturbances using an inverse standard normal cumulative distribution function (CDF) function. This transformation approach was proposed and widely used by previous studies to compute the

inverse Mill's ratio for logit distribution (e.g., Lee 1983, Greene 2002). We provide more details in Online Appendix C. In addition, we have also tried the normal distribution assumption for the error term. We find that our final results stay very consistent.

<sup>5</sup> The log-normal assumption of search cost is consistent with the prior literature (e.g., Kim et al. 2010, Wildenbeest 2011).

<sup>6</sup> Note that  $u_i^*$  is the maximum utility value from the current consideration set  $S_{i,r(j)}$ . Hence, the value of  $u_i^*$  depends on what products are included in the current stage of the consideration set.

<sup>7</sup> Note that based on the model framework, we do not explicitly model the selection rule for the search order, but take it as precalculated (i.e., based on the Weitzman 1979 optimal search model, this search order is precalculated based on the descending order of the reservation utility of each product).

<sup>8</sup> Note that search cost on product search engine might be partly associated with the search engine design. To better account for this factor, in the main analysis, we have controlled for the online positions of a hotel (i.e., page and rank on the search engine website), by which we aim to control for the search engine design efficiency to a large extent. Under such circumstance, our model estimated search costs indicate, conditional on the same online position, the cognitive cost

of searching for a certain product. In addition, we have conducted additional robustness tests by controlling for the sorting methods in a consumer's search session. This is to control for additional factors introduced by search engine design. We find in all of these cases that our model estimated results remain highly consistent.

<sup>9</sup>In principle, consumers can choose from various search strategies to customize their search results. To study this direction, one needs to separately look into the data under each different strategy. In this paper, our main focus is not on the search refinement strategies. We refer the readers to Chen and Yao (2017) and Koulayev (2014) for a more in-depth analysis on that front.

<sup>10</sup>As a robustness check, we also tried the derivative-based optimization algorithms (e.g., the Broyden–Fletcher–Goldfarb–Shanno (BFGS) algorithm and the nested fixed point algorithm (NFXP)). We found that different optimization algorithms are able to recover consistent structural parameters.

<sup>11</sup>One important fact to note is that we also observe rich variation in the characteristics of hotels that enter the consumer's consideration set. In particular, we find that among all of the sessions in which consumers incur click-throughs, 8,731 sessions are associated with a size of (click-generated) consideration set that is larger than 5, and 3,506 sessions are associated with a size that is larger than 10. These observations are critical for our model identification.

<sup>12</sup>In particular, we include only the click-sequence-related information in the likelihood function using the click data only. We estimate this click model using a similar simulated maximum-likelihood approach based on only the click probability.

<sup>13</sup>The major difference between this joint probabilistic model and our main search model is that instead of capturing the sequence of clicks and allowing clicks to be interdependent, the joint model assumes each click decision to be independent. Correspondingly, it models the click decisions independently as following a discrete choice process (e.g., logit model).

<sup>14</sup>Note that we do not focus on the supply-side model in the paper, and we make the implicit assumption that their information-providing practices are exogenous and not necessarily optimal. We believe that this is a reasonable assumption for our model, but more research in that direction could potentially shed more light on the information-provision decisions of the hotels and examine whether there is any strategic rationale for their actions on that front.

<sup>15</sup>We conducted additional analysis to examine the statistical significance in the difference across the simulated revenues in the policy experiments. In particular, given each different set of information on the search summary page, we replicated our simulation experiments 200 times (i.e., via bootstrapping) to acquire the confidence interval of the corresponding simulated platform revenue. We found that the predicted revenues under different scenarios are statistically different from the existing case (i.e., confidence intervals do not overlap).

## References

- Agarwal A, Hosanagar K, Smith M (2011) Location, location, location: An analysis of profitability of position in online advertising markets. *J. Marketing Res.* 48(6):1057–1073.
- Archak N, Ghose A, Ipeirotis PG (2011) Deriving the pricing power of product features by mining consumer reviews. *Management Sci.* 57(8):1485–1509.
- Baye MR, Gatti JR, Kattuman P, Morgan J (2009) Clicks, discontinuities, and firm demand online. *J. Econom. Management Strategy* 18(4):935–975.
- Bikhchandani S, Sharma S (1996) Optimal search with learning. *J. Econom. Dynam. Control* 20(1–3):333–359.
- Blei DM, Ng AY, Jordan MI (2003) Latent Dirichlet allocation. *J. Machine Learn. Res.* (3):993–1022.
- Brynjolfsson E, Dick A, Smith M (2010) A nearly perfect market? Differentiation vs. price in consumer choice. *Quant. Marketing Econom.* 8(1):1–33.
- Chapelle O, Zhang Y (2009) A dynamic Bayesian network click model for web search ranking. *Proc. 18th Internat. Conf. World Wide Web (ACM, New York)*, 1–10.
- Chen P-Y, Hong Y, Liu Y (2018) The value of multidimensional rating systems: Evidence from a natural experiment and randomized experiments. *Management Sci.* 64(10):4629–4647.
- Chen Y, Yao S (2017) Sequential search with refinement: Model and application with click-stream data. *Management Sci.* 63(12):4345–4365.
- Chevalier JA, Mayzlin D (2006) The effect of word of mouth on sales: Online book reviews. *J. Marketing Res.* 43(3):345–354.
- De los Santos B (2008) Consumer search on the Internet. Ph.D. dissertation, Chicago University.
- De los Santos B, Koulayev S (2017) Optimizing click-through in online rankings with endogenous search refinement. *Marketing Sci.* 36(4):542–564.
- De los Santos B, Hortacsu A, Wildenbeest M (2012) Testing models of consumer search using data on web browsing and purchasing behavior. *Amer. Econom. Rev.* 102(6):2955–2980.
- De los Santos B, Hortacsu A, Wildenbeest MR (2017) Search with learning for differentiated products: Evidence from e-commerce. *J. Bus. Econom. Statist.* 35(4):626–641.
- Dellaert BGC, Häubl G (2012) Searching in choice mode: Consumer decision processes in product search with recommendations. *J. Marketing Res.* 49(2):277–288.
- Dellarocas C, Awad N, Zhang M (2007) Exploring the value of online product reviews in forecasting sales: The case of motion pictures. *J. Interactive Marketing* 21(4):23–45.
- Dhar R, Simonson I (2003) The effect of forced choice on choice. *J. Marketing Res.* 40(2):146–160.
- Duan W, Gu B, Whinston AB (2008) Do online reviews matter? An empirical investigation of panel data. *Decision Support Systems* 45(4):1007–1016.
- Ellison G, Ellison SF (2009) Search, obfuscation, and price elasticities on the Internet. *Econometrica* 77(2):427–452.
- Fellbaum C (1998) *Wordnet: An Electronic Lexical Database* (MIT Press, Cambridge, MA).
- Forman C, Ghose A, Wiesenfeld B (2008) Examining the relationship between reviews and sales: The role of reviewer identity disclosure in electronic markets. *Inform. Systems Res.* 19(3):291–313.
- Ghose A, Ipeirotis PG (2011) Estimating the helpfulness and economic impact of product reviews: Mining text and reviewer characteristics. *IEEE Trans. Knowledge Data Engrg.* 23(10):1498–1512.
- Ghose A, Yang S (2009) An empirical analysis of search engine advertising: Sponsored search in electronic markets. *Management Sci.* 55(10):1605–1622.
- Ghose A, Ipeirotis P, Li B (2012) Designing ranking systems for hotels on travel search engines by mining user-generated and crowd-sourced content. *Marketing Sci.* 31(3):493–520.
- Ghose A, Ipeirotis PG, Li B (2014) Examining the impact of ranking on consumer behavior and search engine revenue. *Management Sci.* 60(7):1632–1654.
- Godes D, Mayzlin D (2004) Using online conversations to study word-of-mouth communication. *Marketing Sci.* 23(4):545–560.
- Goldfarb A, Tucker C (2011) Search engine advertising: Channel substitution when pricing ads to context. *Management Sci.* 57(3):458–470.
- Gong J, Abhishek V, Li B (2019) Examining the impact of contextual ambiguity on search advertising keyword performance: A topic model approach. *MIS Quart.* Forthcoming.
- Greene WH (2002) *Limdep Manual*, Version 8.0 (Econometric Software, New York).
- Gumbel EJ (1961) Bivariate logistic distributions. *J. Amer. Statist. Assoc.* 56(294):335–349.
- Hann I, Terwiesch C (2003) Measuring the frictional cost of online transactions: The case of a name-your-own-price channel. *Management Sci.* 49(11):1563–1579.

- Hong H, Shum M (2006) Can search cost rationalize equilibrium price dispersion in online markets? *RAND J. Econom.* 37(2): 258–276.
- Honka E (2014) Quantifying search and switching costs in the US auto insurance industry. *RAND J. Econom.* 45(4):847–884.
- Hortacsu A, Syverson C (2004) Product differentiation, search costs, and competition in the mutual fund industry: A case study of S&P 500 index funds. *Quart. J. Econom.* 119(May):403–456.
- Iyengar SS, Lepper MR (2000) When choice is demotivating: Can one desire too much of a good thing? *J. Personality Soc. Psych.* 79(6):995–1006.
- Jacoby J, Speller DE, Berning CK (1974) Brand choice behavior as a function of information load: Replication and extension. *J. Consumer Res.* 1(1):33–42.
- Kim JB, Albuquerque P, Bronnenberg BJ (2010) Online demand under limited consumer search. *Marketing Sci.* 29(6):1001–1023.
- Kim JB, Albuquerque P, Bronnenberg B (2014) The effects of product innovation on online search and choice. Working paper, Hong Kong University of Science and Technology, Clear Water Bay.
- Koulayev S (2013) Search with Dirichlet priors: Estimation and implications for consumer demand. *J. Bus. Econom. Statist.* 31(2): 226–239.
- Koulayev S (2014) Search for differentiated products: Identification and estimation. *RAND J. Econom.* 45(3):553–575.
- Lee L-F (1983) Generalized econometric models with selectivity. *Econometrica* 51(2):507–512.
- Manning C, Schutze H (1999) *Foundations of Statistical Natural Language Processing* (MIT Press, Cambridge, MA).
- McFadden D (1989) A method of simulated moment for estimation of discrete response models without numerical integration. *Econometrica* 57(5):905–1026.
- Mehta N, Rajiv S, Srinivasan K (2003) Price uncertainty and consumer search: A structural model of consideration set formation. *Marketing Sci.* 22(1):58–84.
- Moraga-Gonzalez JL, Wildenbeest MR (2008) Maximum likelihood estimation of search costs. *Eur. Econom. Rev.* 52(5):820–848.
- Moraga-Gonzalez JL, Sandor Z, Wildenbeest MR (2017) Consumer search and prices in the automobile market. Working paper, Vrije Universiteit Amsterdam, Amsterdam.
- Mortensen DT (1970) Job search, the duration of unemployment, and the Phillips curve. *Amer. Econom. Rev.* 60(5):847–862.
- Netzer O, Feldman R, Goldenberg J, Fresko M (2012) Mine your own business: Market-structure surveillance through text mining. *Marketing Sci.* 31(3):521–543.
- Reinganum JF (1982) Strategic search theory. *Internat. Econom. Rev.* 23(1):1–17.
- Rosenfield DB, Shapiro RD (1981) Optimal adaptive price search. *J. Econom. Theory* 25(1):1–20.
- Rothschild M (1974) Searching for the lowest price when the distribution of prices is unknown. *J. Political Econom.* 82(4):689–711.
- Sterling G (2008) iProspect: Blended search resulting in more clicks on news, images, and video. *Search Engine Land* (April 7), <https://searchengineland.com/iprospect-blended-search-resulting-in-more-clicks-on-news-images-and-video-13708>.
- Stigler GJ (1961) The economics of information. *J. Political Econom.* 69(3):213–225.
- Weitzman ML (1979) Optimal search for the best alternative. *Econometrica* 47(3):641–654.
- Wildenbeest MR (2011) An empirical model of search with vertically differentiated products. *RAND J. Econom.* 42(4):729–757.
- Yao S, Mela CF (2011) A dynamic model of sponsored search advertising. *Marketing Sci.* 30(3):447–468.



**Modeling Consumer Footprints on Search Engines:  
An Interplay with Social Media  
(Online Appendices)**

Anindya Ghose

Stern School of Business, New York University  
[aghose@stern.nyu.edu](mailto:aghose@stern.nyu.edu)

Panagiotis G. Ipeirotis

Stern School of Business, New York University  
[panos@stern.nyu.edu](mailto:panos@stern.nyu.edu)

Beibei Li

Heinz College, Carnegie Mellon University  
[beibeli@andrew.cmu.edu](mailto:beibeli@andrew.cmu.edu)

## Online Appendix B.

### Optimal Search Framework

Our model builds on the optimal sequential search framework. Consider that consumers are forward-looking and trying to maximize the *expected present value of utility* over a planning horizon (e.g., Erdem and Keane 1996). The expected present value in our setting can be computed as follows. First, we partition the set of available alternatives into  $S_i \cup \bar{S}_i$ , with  $S_i$  containing all the ones that have been searched and  $\bar{S}_i$  containing all the non-searched ones. Let  $u_i^*$  be the highest net value searched so far, thus, we have

$$u_i^* = \max_{j \in S_i} \{u_{ij}, 0\}. \quad (\text{B1})$$

Note that Equation (B1) is the same as Equation (4) in the paper.

The state of the system at any time during the search is given by  $(u_i^*, \bar{S}_i)$ . Define  $\Psi(u_i^*, \bar{S}_i)$  as the expected present discounted value of following an optimal search policy, from the current state  $(u_i^*, \bar{S}_i)$  going forward. Therefore, for each  $u_i^*$  and  $\bar{S}_i$ , the state valuation function  $\Psi(u_i^*, \bar{S}_i)$  must satisfy the Bellman equation:

$$\Psi(u_i^*, \bar{S}_i) = \max \left( u_i^*, \max_{j \in \bar{S}_i} \left( -c_{ij} + d_i \cdot \left[ \underbrace{\Psi(u_i^*, \bar{S}_i - \{j\}) \cdot \int_{-\infty}^{u_i^*} f(u_{ij}) du_{ij}}_{u_{ij} \leq u_i^*} + \underbrace{\int_{u_i^*}^{\infty} \Psi(u_{ij}, \bar{S}_i - \{j\}) f(u_{ij}) du_{ij}}_{u_{ij} > u_i^*} \right] \right) \right), \quad (\text{B2})$$

where  $F(\bullet)$  is the CDF of  $u_{ij}$  and  $f(\bullet)$  is the probability density function of  $u_{ij}$ . Therefore, at current state  $(u_i^*, \bar{S}_i)$ , the consumer can either terminate search and collect reward  $u_i^*$ , or search any  $j \in \bar{S}_i$  to maximize  $\Psi(u_i^*, \bar{S}_i)$ . Given the short time span in online search, we set the discount rate  $d_i$  to 1. Equation (B2) is the principle of optimality for dynamic programming.

As pointed out by Weitzman (1979) and Lippman & McCall (1976) in the classical economic literature of search, the optimal solution to this dynamic programming has a *myopic* solution: Namely, the consumer needs only compare her return from stopping and accepting reward  $u_i^*$  with the expected return from exactly one more search.

More formally, let the expected marginal utility for consumer  $i$  from the search of product  $j$  be

$$B_{ij}(u_i^*) = \int_{u_i^*}^{\infty} (u_{ij} - u_i^*) f(u_{ij}) du_{ij}. \quad (\text{B3})$$

Thus, consumer  $i$  will continue to search if there exists at least one  $j$  such that the expected marginal benefit from searching product  $j$  exceeds its corresponding search cost

$$c_{ij} < B_{ij}(u_i^*). \quad (\text{B4})$$

Therefore, the optimal search strategy for a consumer is to continue searching until a value  $u_i^*$  is found that violates Equation (B4).

## Online Appendix C.

### Computation of Reservation Utility

Our model builds on the optimal sequential search framework by Weitzman (1979). Define the *reservation utility*  $z_{ij}$  as the utility value that satisfies the following boundary condition, where the search cost equates the expected marginal utility from searching product  $j$  (same as Equation (6)).

$$c_{ij} = B_{ij}(z_{ij}) = \int_{z_{ij}}^{\infty} (u_{ij} - z_{ij})f(u_{ij})du_{ij}. \quad (C1)$$

The optimal search strategy for a consumer is to continue searching until she finds a value  $u_i^*$  larger than the boundary solution  $z_{ij}$ .

The reservation utility  $z_{ij}$  can be solved from Equation (C1) given the search cost,  $z_{ij} = B_{ij}^{-1}(c_{ij})$ . Weitzman (1979) has proved the function  $B_{ij}(z_{ij})$  is continuous and monotonic. Therefore, there exists a unique solution  $z_{ij}$  to the equation  $c_{ij} = B_{ij}(z_{ij})$ .

Let  $\bar{u}_{ij}$  be the mean and  $\sigma_{ij}^2$  be the variance of the utility distribution  $f(u_{ij})$ . Based on our model setting, before the click-through we can write down the mean and the variance of the expected utility as the following:

$$\begin{aligned} \bar{u}_{ij} &= E[u_{ij}] = E(X_j\beta_i - \alpha_i P_j + \tilde{L}_j\lambda_i + e_{ij}) \\ &= X_j\beta_i - \alpha_i P_j + \tilde{L}_j\lambda_i + E(e_{ij}), \end{aligned} \quad (C2)$$

and

$$\begin{aligned} \sigma_{ij}^2 &= \text{VAR}[u_{ij}] = \text{VAR}(X_j\beta_i - \alpha_i P_j + \tilde{L}_j\lambda_i + e_{ij}) \\ &= \text{VAR}(e_{ij}). \end{aligned} \quad (C3)$$

$\tilde{L}_j$  is the expected value of the unobserved landing-page characteristics before click. It can be estimated based on the mean of the bootstrapping samples drawn from the empirical distribution of landing-page characteristics for hotel  $j$  conditional on the observed summary-page characteristics  $(X_j, P_j)$ . Meanwhile, based on the assumption  $e_{ij} \sim \text{Type I EV}(0,1)$ , we can derive  $E(e_{ij}) = 0.5772$  (“Euler-Mascheroni constant”) and  $\text{VAR}(e_{ij}) = \pi^2 / 6$ . Therefore,  $\bar{u}_{ij}$  and  $\sigma_{ij}^2$  can be derived accordingly.

We rewrite Equation (C1) as follows:

$$\begin{aligned} c_{ij} &= \int_{z_{ij}}^{\infty} (u_{ij} - z_{ij})f(u_{ij})du_{ij} \\ &= (1 - F(z_{ij})) \left[ \int_{z_{ij}}^{\infty} (u_{ij} - z_{ij}) \frac{f(u_{ij})}{1 - F(z_{ij})} du_{ij} \right], \end{aligned} \quad (C4)$$

where  $F(z_{ij})$  is the CDF of  $u_{ij}$  evaluated at  $z_{ij}$ .

We compute the reservation utility using a similar approach as Kim et al. (2010).<sup>1</sup> Let  $\eta_{ij} = \frac{z_{ij} - \overline{u_{ij}}}{\sigma_{ij}}$

and  $\tau_{ij} = \frac{c_{ij}}{\sigma_{ij}}$ , we can rewrite Equation (C4) as follows:

$$\tau_{ij} = \frac{c_{ij}}{\sigma_{ij}} = \frac{(1 - F(\eta_{ij})) \left[ \overline{u_{ij}} - z_{ij} + \sigma_{ij} \frac{f(\eta_{ij})}{1 - F(\eta_{ij})} \right]}{\sigma_{ij}}$$

$$\Leftrightarrow \tau_{ij} = g(\eta_{ij}) = (1 - F(\eta_{ij})) \left[ \frac{f(\eta_{ij})}{1 - F(\eta_{ij})} - \eta_{ij} \right], \quad (C5)$$

If we can solve  $\eta_{ij} = g^{-1}(\tau_{ij})$  from the above Equation (C5), then we can solve  $z_{ij} = \eta_{ij}\sigma_{ij} + \overline{u_{ij}}$ .

Note that Equation (B5) does not involve any model parameters. Therefore, in practice we only need to solve it once and use the results in the model estimation (Kim et al. 2010, Koulayev 2014). For computational tractability, we apply an interpolation approach to solve Equation (C5).

More specifically, we compute the reservation utility  $z_{ij}$  in the following four steps:

- 1) Pre-construct a lookup table for each pair  $(\eta_{ij}, \tau_{ij})$  based on Equation (C5).
- 2) At any stage in the estimation, use the current values of  $c_{ij}$  and  $\sigma_{ij}$  to compute  $\tau_{ij} = \frac{c_{ij}}{\sigma_{ij}}$ .
- 3) Based on the value of  $\tau_{ij}$ , look up the corresponding value  $\eta_{ij}$  in the pre-constructed table.
- 4) Based on the current values of  $\eta_{ij}$ ,  $\sigma_{ij}$  and  $\overline{u_{ij}}$ , solve the reservation utility  $z_{ij} = \eta_{ij}\sigma_{ij} + \overline{u_{ij}}$ .

---

<sup>1</sup> Note that different from Kim et al. (2010), who assume standard normal distribution of the error, we allow for logit distribution of the error term in our model. To calculate the function with regard to the inverse Mill's ratio ( $f(u_{ij}) / [1 - F(z_{ij})]$ ) from Equation (C4) to Equation (C5), we first need to transform the logit error into standard normal disturbances using an inverse standard normal CDF function. This transformation approach was proposed and widely used by previous studies to compute the inverse Mill's ratio for logit distribution (e.g., Lee (1983), Greene (2002)).

## Online Appendix D.

### More Details on Using the Simulated Approach to Construct the Conditional Purchase Probability

As we discussed in section 4.4, conditional on the sequence of clicks consumer  $i$  has made in the search session, we can derive the conditional probability that she purchases hotel  $r(j)$  in her consideration set as the following:

$$\begin{aligned}
 \eta_{i,r(j)} &= P(r(j) \text{ is booked by consumer } i) \\
 &= \Pr\left(u_{i,r(j)} \geq u_{i,r(j')}, \quad \forall r(j) \neq r(j'), \quad r(j), r(j') \in S_i\right) \\
 &= \Pr\left(\begin{array}{c} V_{i,r(j)}^S + V_{i,r(j)}^L + e_{i,r(j)} \geq V_{i,r(j')}^S + V_{i,r(j')}^L + e_{i,r(j')}, \\ \forall r(j) \neq r(j'), \quad r(j), r(j') \in S_i \end{array}\right),
 \end{aligned} \tag{D1}$$

where  $S_i$  is the click-generated choice set for consumer  $i$ . Equation (D1) is identical to Equation (8) in the paper.

Note that because the consideration set  $S_i$  is selected by consumer  $i$  based on her search decisions,  $e_{ij}$  does not follow a full Type I EV distribution. Instead, it follows a truncated Type I EV distribution based on the optimality conditions used by the consumer. Unfortunately, under such circumstance the conditional choice probability does not have a close-form expression (e.g., Logit form). To address this selection issue, we applied a simulation approach. Similar methods have been adopted by the previous studies (Chen and Yao 2016, Honka 2014, McFadden 1989).

Our simulation approach builds on the methods from Chen and Yao (2016) and Honka (2014). It allows us to simulate the error term from a truncated Type I EV distribution by satisfying the follow three optimality conditions: 1) Sequence of the click-generated choice set; 2) Composition of the click-generated choice set; 3) Utility optimality of the final choice. More specifically, the simulated purchase probability has to satisfy the following conditions:

- 1) At any moment during the consumer search process, the utility of the currently being clicked product  $j$ ,  $u_{i,r(j)}$ , is smaller than the reservation utilities of those products clicked after  $j$ . This is because the consumer continues to search afterwards. Here,  $S_i^{clicked}$  denotes the set of all products that have been clicked by consumer  $i$ .

$$u_{i,r(j)} < \min(z_{i,r(j')}, \forall r(j') \in S_i^{clicked} \text{ and } r(j') > r(j)) \tag{D2}$$

- 2) The utility of the final purchased product,  $u_{i,r(j^*)}$ , is greater than the reservation utilities of all the remaining unsearched products,  $z_{i,r(j')}$ . This is because the consumer stops searching afterwards. Here,  $S_i^{unclicked}$  represents the set of all products that have not been clicked by consumer  $i$ .

$$u_{i,r(j^*)} \geq z_{i,r(j')}, \forall r(j') \in S_i^{unclicked} \quad (D3)$$

- 3) The utility of the final purchased product,  $u_{i,r(j^*)}$ , is greater than the utility of any other product in the click-generated choice set,  $u_{i,r(j')}$ . This is the final choice utility optimality condition.

$$u_{i,r(j^*)} \geq u_{i,r(j')}, \forall r(j') \in S_i^{clicked} \text{ and } r(j^*) \neq r(j') \quad (D4)$$

Hence, when simulate the error term in the utility function to construct the conditional purchase probability, we need to draw from the truncated Type I EV distribution by taking into consideration all the three optimality conditions (D2)-(D4) above. As discussed in Chen and Yao (2016), for different products in the click-generated choice set, the error terms are truncated differently. For clicked products that are not purchased, the error term is right truncated. For the purchased product, the error term is left truncated if it is the final click during the search process; if it is not the final search, then it is truncated on both sides.

To construct the conditional purchase probability, an intuitive approach is to draw the error term from the Type I EV distribution based on the three truncation optimality conditions in (D2)-(D4), by counting the frequency that the three optimality conditions are satisfied. More specifically, we adopted a similar approach as used by Chen and Yao (2016). The step-by-step implementation of our simulation method can be summarized as follows:

- i) Conditional on a given set of other parameters in our model, for each consumer  $i$  and each product  $j$  in the consideration set, draw 200  $e_{i,j}$  from Type I EV distribution, depending on the three truncation optimality conditions;
- ii) Count the frequency of the three optimality conditions (D2)-(D4) being satisfied across the 200 random draws of  $e_{i,j}$ 's;
- iii) Iterate 100 times the above two steps, repeatedly making new draws of other parameters during each round of iteration;
- iv) Average the simulated frequencies from step ii) across the 100 iterations to calculate the final simulated purchase probability of consumer  $i$ .

However, this approach can be computationally expensive. As a robustness check, we also tried an alternative method with a kernel-smoothed frequency simulator which was proposed by McFadden (1989) and was suggested by Honka (2014). In this approach, we smoothed the probabilities using a multivariate scaled logistic CDF (Gumbel 1961) with all the scaling factors equal to 15.<sup>2</sup> Notice that due to data limitation, Honka (2014) considers only the composition of the click-generated choice set but not the sequence of clicks as the optimality condition, whereas Chen and Yao (2016) observe the sequence of search process which allows them to consider both the composition and the sequence of the click-generated choice set. In our study, similarly as Chen and Yao (2016) we observe both types of information. Therefore, we are able to account for the additional optimality condition regarding the sequence of consumer clicks compared to Honka (2014).

---

<sup>2</sup> For more details on the kernel-smoothed frequency simulator, we refer interested readers to Online Appendix B in Honka (2014). The main idea of this simulator is to calculate the smoothed conditional purchase probability using a multivariate scaled logistic CDF (Gumbel 1961). In our estimation, we allow all the scaling factors to be equal to 15 as suggested by Honka (2014). However, we have also tried other values for the scaling factors ranging from 10-30. We found our estimation results stay qualitatively consistent.

## Online Appendix E.

### Using Topic Modeling to Generate the Topic Entropy Score for Each Hotel

In addition to review textual readability and subjectivity, we also extracted an additional cognitive cost indicator based on the topic complexity of the customer reviews. In particular, built on prior literature (Gong et al. 2016) we analyzed the entropy value for the distribution of topics extracted from all customer reviews for each hotel. This topic entropy measures the diversity of topics covered by the customer reviews for each hotel. Prior literature suggests the diversity in search results affects consumer search behavior (e.g., Weitzman 1979, Dellaert and Haubl 2012). In addition, consumer psychology theories suggest that as the information become noisier, users are more likely to abandon their search (e.g., Jacoby et al. 1974; Dhar and Simonson 2003), because users tend to get overwhelmed and discouraged by the complexity of information, and therefore lose their interest or trust in the search results. Therefore, we derived a “Topic Entropy” score using probabilistic topic models from machine learning and natural language processing to capture the “noisiness” of information provided by the customer reviews.

Topic models are unsupervised algorithms that aim to extract hidden topics from unstructured text data. The intuition behind topic models is that a topic is a cluster of words that frequently occur together, and that documents, consisting of words, may belong to multiple topics with different probabilities. A probabilistic topic model tries to discover the underlying topic structure in a statistical framework. In particular, we measure the topic complexity of reviews for each product by estimating a topic model using Latent Dirichlet Allocation model (LDA; Blei et al. 2003), and subsequently computing the entropy (i.e. diversity) of the topic distribution of reviews for that product. We discuss the details how we use topic modeling to generate the Topic Entropy score for each hotel below.

#### **1. Corpus Construction and Document Pre-processing.**

We first construct a corpus of documents that describe the information content conveyed by the hotel reviews. In particular, we collect all the customer reviews for each hotel. Hence, each hotel is associated with a review document, and all documents together construct the overall corpus. After constructing the corpus of review documents, we pre-process the documents following a standard procedure (e.g., Aral et al. 2011, Gong et al. 2016). We first remove annotations and tokenize the sentence into distinct terms. Then we remove stop words using a standard dictionary.

#### **2. Latent Dirichlet Allocation (LDA).**

We use topic models to automatically infer semantic interpretations of keyword meanings. The most widely used topic model is the Latent Dirichlet Allocation model (LDA; Blei et al. 2003), which is a hierarchical Bayesian model that describes a generative process of document creation. Previous research shows that humans tend to agree with the coherence of the topics generated by LDA, which provides strong support for the use of topic models for information retrieval applications.



The goal of LDA is to infer topics as latent variables from the observed distribution of words in each document. In particular, a topic is defined as a multinomial distribution over a vocabulary of words, a document is a collection of words drawn from one or more topics, and a corpus is the set of all documents. Based on the discussion above, we construct a document for each hotel that best reflects the information of the hotel review. We now discuss how we use LDA to infer the topics from the corpus of documents.

Formally, let  $T$  be the number of topics related to the corpus, let  $D$  be the number of documents in the corpus, and let  $W$  be the total number of words in the corpus. We assume that each document in the corpus is generated according to the following process:

Step 1. For each topic  $t$ , choose  $\phi_t = (\phi_{t1}, \dots, \phi_{tW}) \sim \text{Dirichlet}(\varphi)$ , where  $\phi_t$  describes the word distribution of topic  $t$  over the vocabulary of words.

Step 2. For each document  $d$ , choose  $\theta_d = (\theta_{d1}, \dots, \theta_{dT}) \sim \text{Dirichlet}(\omega)$ , where  $\theta_{dt}$  is the probability of topic  $t$  to which document  $d$  belongs.

Step 3. For each word  $n$  in document  $d$ , (1) choose a topic  $t_{dn} \sim \text{Multinomial}(\theta_d)$ , and (2) choose a word  $w_{dn} \sim \text{Multinomial}(\phi_{t_{dn}})$ .

$\varphi$  and  $\omega$  are hyper-parameters for the two prior distributions -  $\text{Dirichlet}(\varphi)$  as the prior distribution of  $\phi$  (word distribution in a topic) and  $\text{Dirichlet}(\omega)$  as the prior distribution of  $\theta$  (topic distribution in a document). We use the values suggested by Steyvers and Griffiths (2007) ( $\varphi = 0.01$  and  $\omega = 50/T$ ).

Based on the generative process described above, we use a Markov chain Monte Carlo (MCMC) algorithm to estimate  $\phi$  and  $\theta$ . Specifically, we use a collapsed Gibbs sampler to sequentially sample the topic of each word token in the corpus conditional on the current topic assignments of all other word tokens. We run a collapsed Gibbs sampler using MALLETT (McCallum 2002) with 2,000 iterations. For each hotel, we obtain the posterior topic probabilities inferred from its corresponding document of customer reviews. In our study, we estimate the LDA model with a different number of topics,  $T = 20, 50, \text{ and } 100$ .

### 3. Topic Entropy as a Measure for Keyword Ambiguity.

We propose using Topic Entropy to measure the complexity of hotel reviews. It captures the uncertainty of a document's topic distribution. In information theory, entropy measures the unpredictability of a random variable. In our context, each hotel is associated with its own review topic distribution inferred from the hotel-specific document. Therefore, we treat the topic assignment as a multinomial random variable, and use topic entropy to quantify how “noisy” the customer reviews for a hotel are in terms of underlying topics. The higher the entropy is, the more complex or noisier the reviews for that hotel. In other words, hotels with higher Topic Entropy tend to relate to a broader range of topics (more complex), whereas hotel with lower Topic Entropy tend to relate to fewer dominant topics (less complex).

More formally, let  $\tilde{\theta}_{kt}$  denote the posterior probability that hotel  $k$  belongs to topic  $t$ . We therefore define the topic entropy of hotel  $k$  as follows:

$$TopicEntropy_k = -\sum_{t=1}^T \tilde{\theta}_{kt} \log(\tilde{\theta}_{kt}), \quad (E1)$$

where  $T$  is the total number of topics.

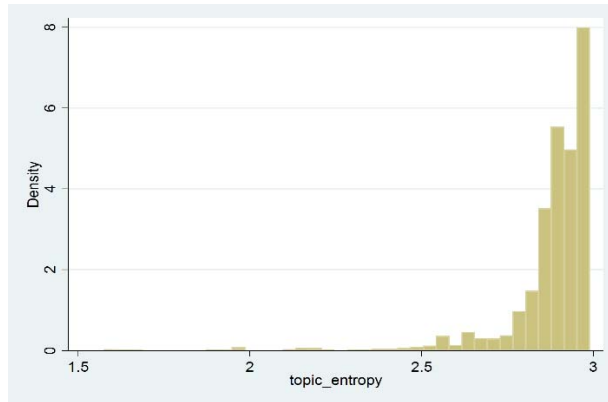
We present the summary statistics for the estimated Topic Entropy in Table E1. As we can see, the maximum entropy value depends on the number of topics chosen. Simple calculation also shows that with  $T$  topics, entropy ranges from 0 to  $\ln(T)$ .<sup>3</sup> The high correlations among entropy values derived based on a different number of topics also suggest entropy seems to be fairly robust to the number of topics specified in the LDA model.

**Table E1. Summary Statistics of Topic Entropy**

	Mean	Std. Dev.	Min.	Max.	Correlation	
					20 Topics	50 Topics
20 Topics	2.78	0.13	1.58	2.99		
50 Topics	3.03	0.17	1.62	3.91	0.89	
100 Topics	3.16	0.20	1.68	4.60	0.88	0.93

In our model estimation (i.e., Robustness Test 2) and policy experiment, we used the Topic Entropy values derived based on 20 topics. We illustrate the distribution of the Topic Entropy in Figure E1 based on 20 topics ( $T=20$ ). We also tried using 50 topics and 100 topics and the results are qualitatively consistent.

**Figure E1. Distribution of Topic Entropy ( $T=20$ )**



**References:**

- Steyvers, M., and Griffiths, T. 2007. Probabilistic Topic Models. Handbook of Latent Semantic Analysis (427:7), pp. 424-440.
- McCallum, A. K. 2002. MALLET: A Machine Learning for Language Toolkit.

<sup>3</sup> The number of topics  $T$  is pre-specified before estimating the LDA model. Entropy for hotel  $k$  is the smallest when there exists  $t \in \{1, \dots, T\}$  such that  $\tilde{\theta}_{kt} = 1$ ; Entropy is the largest when for all  $t \in \{1, \dots, T\}$ ,  $\tilde{\theta}_{kt} = 1/T$ .

## **Online Appendix F.**

### **Robustness Tests (1) – (3)**

To understand the robustness of the model estimation, and to analyze how unstructured social media and consumer heterogeneity (e.g., travel purposes) may affect the search cost and decisions of a consumer, we conduct three sets of robustness tests:

***Robustness Test 1: Exclude the unstructured social media data (i.e., no social media textual variables in utility or search-cost specifications).***

One of the main goals in our study is to examine how social media textual content affects consumer utility and search cost. Therefore, we are interested in comparing the differences in the search models with and without the set of social media textual variables. The results of this test are illustrated in Table F1, column 2 (“R1”). We find the estimated coefficients are qualitatively consistent with the main results. After computing the price elasticity, we notice the model that does not account for social media textual variables presents significantly higher price elasticity (2.128 vs. 1.619,  $p < 0.05$ ). This result indicates that the unstructured social media textual information plays an important role in consumer decision making, and that consumers’ cognitive costs to digest such information are non-negligible. Without accounting for such unstructured information during consumer product search can lead to an overestimation of price elasticity.

***Robustness Test 2: Include additional unstructured social media variable (i.e., “Topic Entropy” measurement derived from the online review textual content using topic modeling).***

To further understand the role of unstructured social media content during consumer search, in addition to the readability and subjectivity of review textual content, we also extracted an additional cognitive cost indicator based on the topic complexity of the customer reviews, “Topic Entropy.” In Robustness Test 2, we included this new social media variable into the search cost model and re-estimated our model. The results of this test are illustrated in Table F1, column 3 (“R2”). Overall, we find the estimated coefficients are qualitatively consistent with the main results. Moreover, we find the increase of Topic Entropy in the customer reviews for a hotel can lead to a significant increase in the search cost for that hotel. Intuitively, this result indicates that as the topics discussed in customer reviews become noisier, it can significantly increase the cognitive costs for consumers who are reading through them. Therefore, it might be more efficient for product search engines to provide a careful digest of the topics extracted from all the textual reviews, or to provide a guidance on the expected topics to be discussed from the reviewers (e.g., room service, friendliness of the staff, parking facilities, etc.).

**Table F1. Estimation Results - Robustness Tests (1) & (2)**

Variable	Mean Effect (Std. Err) <sup>R1</sup>	Heterogeneity (Std. Err) <sup>R1</sup>	Mean Effect (Std. Err) <sup>R2</sup>	Heterogeneity (Std. Err) <sup>R2</sup>
(Preferences)	$\bar{\alpha}, \bar{\beta}, \bar{\lambda}$	$\sigma_{\alpha}, \Sigma_{\beta}, \Sigma_{\lambda}$	$\bar{\alpha}, \bar{\beta}, \bar{\lambda}$	$\sigma_{\alpha}, \Sigma_{\beta}, \Sigma_{\lambda}$
<i>PRICE</i> <sup>(L)</sup>	-1.706* (.029)	.490* (.079)	-1.249* (.023)	.423* (.071)
<i>PAGE</i>	-.187* (.003)	.099 (.158)	-.237* (.003)	.082 (.130)
<i>RANK</i>	-.250* (.008)	.189* (.057)	-.317* (.007)	.135* (.066)
<i>CLASS</i>	1.614* (.039)	.806* (.112)	1.511* (.021)	.934* (.180)
<i>AMENITYCNT</i> <sup>(L)</sup>	.156* (.034)	.039 (.080)	.142* (.032)	.065 (.072)
<i>ROOMS</i> <sup>(L)</sup>	.343* (.031)	.238 (.351)	.397* (.022)	.193 (.281)
<i>EXTAMENITY</i> <sup>(L)</sup>	.172* (.041)	.039* (.012)	.166* (.035)	.042 (.044)
<i>BEACH</i>	1.890* (.020)	.503* (.095)	1.541* (.028)	.560* (.097)
<i>LAKE</i>	-.784* (.118)	1.075* (.303)	-.667* (.115)	1.555* (.383)
<i>TRANS</i>	1.243* (.171)	.160 (.133)	1.339* (.141)	.197* (.065)
<i>HIGHWAY</i>	.399* (.112)	.055 (.042)	.443* (.091)	.071 (.063)
<i>DOWNTOWN</i>	.962* (.062)	.206 (.074)	1.195* (.063)	.475* (.094)
<i>CRIME</i>	-.159* (.033)	.020 (.053)	-.171* (.045)	.018 (.031)
<i>RATING</i>	2.898* (.020)	.983* (.082)	2.660* (.017)	1.309* (.089)
<i>REVIEWCNT</i> <sup>(L)</sup>	1.653* (.129)	.437* (.081)	1.228* (.106)	.366* (.062)
<i>STAFF</i>	----	----	.135* (.028)	.035 (.082)
<i>FOOD</i>	----	----	.223* (.034)	.138* (.002)
<i>BATHROOM</i>	----	----	.296 (.270)	.067 (.101)
<i>PARKING</i>	----	----	.097* (.005)	.079* (.014)
<i>BEDROOM</i>	----	----	-.179 (.236)	.251 (.271)
<i>FRONTDESK</i>	----	----	.066 (.108)	.021 (.070)
<i>BRAND</i>	Yes		Yes	
(Search Cost)	$\bar{\gamma}$	$\Sigma_{\gamma}$	$\bar{\gamma}$	$\Sigma_{\gamma}$
<i>Search Base Cost (Constant)</i>	-4.849* (.081)	1.303* (.124)	-7.976* (.095)	.904* (.168)
<b><i>TOPICENTROPY</i></b>	----	----	.298* (.036)	.334* (.097)
<i>COMPLEXITY</i>	----	----	.512* (.088)	.373* (.102)
<i>SYLLABLES</i> <sup>(L)</sup>	----	----	.633* (.121)	.705* (.099)
<i>SPELLERR</i> <sup>(L)</sup>	----	----	.286* (.085)	.039 (.111)
<i>SUB</i>	----	----	.187* (.042)	.063 (.232)
<i>SUBDEV</i>	----	----	.302* (.054)	.123 (.267)
<i>Maximum LL</i>	- 338,301		-409,742	
<i>Price Elasticity</i>	-2.128		-1.615	

(L) Logarithm of the variable. \* Statistically significant at 5% level.

R1: Robustness Test 1 (Main Model Without Social Media Cognitive Cost Variables).

R2: Robustness Test 2 (Main Model With Additional Topic Entropy Variable to Measure Topic Complexity)

***Robustness Test 3: Interaction effects between travel purposes and consumer preferences/search cost variables.***

To account for consumer heterogeneity during the search process, we focus on how consumers' heterogeneous travel purposes explain certain variation in search cost and in consumer preferences. To do so, we investigate the interaction effects between consumer travel purposes and consumer preferences/search cost variables. In particular, the variables on which we focus are the summary-page variables. To capture consumers' heterogeneous travel purposes, we define  $T_i$  as an indicator vector with identity components representing the travel purpose:

$$T_i' = [Family_i \ Business_i \ Romance_i \ Tourist_i \ Kids_i \ Senior_i \ Pets_i \ Disability_i]_{1 \times 8}. \quad (F1)$$

We acquire the empirical distribution of  $T_i$  from online consumer reviews and reviewers' profiles.<sup>4</sup>

Then we interact the summary-page variables with  $T_i$  in our search model. We estimate this model using a simulated maximum likelihood approach. We find that consumers' travel purposes can explain their heterogeneity toward specific search-cost and preferences variables. Interestingly, we find interesting interaction patterns between consumers' travel purposes and hotel characteristics. For example, we find that travelers on a romantic trip, relative to other types of travelers, tend to place more importance on online customer reviews (i.e., both the valance of online ratings and the volume of reviews). In addition, consistent with prior research (Ghose et al. 2012), we find that business travelers are the least price-sensitive, whereas tourists tend to be more sensitive to hotel price. We provide the detailed information on the estimated interaction effects in Table F2 ("R3").

**Table F2. Estimation Results - Robustness Test (3)**  
**Search Model with Interaction Effects Between Travel Purposes and Summary-Page Variables**

Variable	Mean Effect (Std. Err) <sup>R3</sup>	Family	Business	Romance	Tourist	Kids	Senior	Pets
<i>PRICE</i> <sup>(L)</sup>	-1.169* (.024)	<b>-.118*(.033)</b>	<b>.331*(.025)</b>	.123(.081)	<b>-.219*(.049)</b>	----	-.314(.257)	----
<i>PAGE</i>	-.212* (.021)	.005(.020)	<b>-.041*(.010)</b>	<b>.035*(.003)</b>	.021(.024)	----	----	----
<i>RANK</i>	-.288* (.034)	<b>-.037* (.008)</b>	<b>-.025* (.003)</b>	.022 (.021)	.154(.166)	-.011(.028)	----	----
<i>CLASS</i>	1.432* (.036)	<b>.067*(.021)</b>	<b>.092*(.022)</b>	<b>.065*(.016)</b>	<b>-.179*(.023)</b>	.200 (.363)	----	----
<i>RATING</i>	2.487* (.021)	.183(.226)	-.368(.399)	<b>.395*(.033)</b>	.040 (.057)	<b>.291*(.026)</b>	<b>.202*(.053)</b>	----
<i>REVIEWCNT</i>	1.315* (.043)	<b>.177*(.023)</b>	-.256(.219)	<b>.301*(.042)</b>	<b>.123*(.026)</b>	----	----	----

(L) Logarithm of the variable.

\* Statistically significant at 5% level.

R3: Robustness Test 3 (Search Model with Interaction Effects).

Note: Some interaction effects are dropped in the estimation due to practical reasons (e.g., collinearity or very low significance).

<sup>4</sup> After writing an online review for a hotel, a reviewer is asked to provide additional demographic and trip information—e.g., “What was the main purpose of this trip? (Select one from the eight choices.)” We derive the distribution of  $T_i$  based on reviewers' responses to this question.

## Online Appendix G.

### Model Comparisons – Details on the Alternative Models

Furthermore, to understand how the type and scale of data or modeling mechanisms may affect the performance of our analysis, we conducted model comparison analyses. In particular, we considered a set of alternative benchmark models using different data sets or modeling mechanisms. We discuss the details of the alternative modeling mechanisms in this appendix.

#### ***(1) Alternative Model I: Use the purchase data only (Mixed Logit Model).***

In reality, due to the unavailability of the individual-level click stream information, we may have access to only the purchase information. Classical static demand estimation models (e.g., Mixed Logit) are used to infer consumer preferences from the purchase data only. However, static demand estimation models do not consider the endogenous and limited nature of search-generated choice sets. With the recent growing pervasiveness of Internet, Web 2.0 and storage technologies, businesses have started tracking individual-level data beyond the final purchase to analyze a consumer's "path-to-purchase." However, individual-level click stream data are often at a much larger scale than the traditional purchase data and hence require more advanced methodologies and computing power for analysis.<sup>5</sup>

To examine how well our proposed search model performs by incorporating the additional click information in understanding consumer preferences, we consider a model that is widely used in the static demand estimation: the Mixed Logit model (e.g., McFadden and Train 2000).<sup>6</sup> To account for the variation in choice sets, we model the consumer decision process under the actual searched (limited) choice set, rather than under the universal choice set available in the market. Note the major difference between a static Mixed Logit model with actual choice sets and our proposed model is that our model captures not only the limited nature of the choice sets, but also the sequential and endogenous formation process of the choice sets. A static model typically takes the choice set as exogenously given.

Interestingly, we find that using a static model without accounting for consumers' search behavior can lead to an overestimation of the price elasticity (2.973 vs. 1.619,  $p < 0.05$ ). The interpretation of this finding can be attributed to the nature of the hotel search market. A model that captures consumers' actual search behaviors finds lower price elasticity, implying consumers in the hotel search market tend to highly evaluate the quality of hotels and put weight on non-price factors during search (e.g., class, amenities, or reviews). Our finding on price elasticity is consistent with prior findings by Koulayev (2014) and Brynjolfsson et al. (2010). Both studies

---

<sup>5</sup> For example, our original click stream data set contains approximately a total of seven million observations, whereas focusing only on the purchase data will reduce our total number of observations to about only eight thousand.

<sup>6</sup> Note that the conditional purchase probability  $\eta_{i,j}$  here has a close form as defined in the Mixed Logit model, because the consideration set in this case is not endogenously generated, but exogenously given. Therefore, there is no selection issue and the error term follows its initial Type I EV distribution.

show that when consumers face a highly differentiated market (e.g., product differentiation or retailer differentiation), they are more likely to focus on non-price factors during search. Hence the estimated price elasticity is lower when incorporating consumers' search behaviors into the model. On the contrary, when a market is less differentiated, consumers become more price-sensitive and tend to focus on price search. Thus a search model that incorporates consumers' search behaviors may find a higher price elasticity of demand than a static model (e.g., de los Santos et al. 2012). The estimation results of this model are shown in Table G1, column 3 ("A1"). For easy comparison with the main model, we also provide the estimation results from our main model in Table G1, column 2 ("M").

***(2) Alternative Model II: Use the purchase data only (Mixed Logit Model + Additional Search Cost Variables).***

One concern towards the validity of *Alternative Model I* is that the static Mixed Logit model does not consider the search cost, so any difference in the estimation is likely to be caused by the missing variables that appear in the search cost from the search model. For example, the estimated parameter increase in price coefficient may be due to the correlation between search and prices—consumers search more for high-priced goods. Hence, inferences regarding price effects are likely to be biased when one does not control for quality (which may be related to the attributes that show up in the search-cost function). Therefore, we consider an alternative Mixed Logit model by incorporating the additional search cost variables. We provide the corresponding results in Table G2, column 2 ("A2").

We find the estimates from *Alternative Model II* are qualitatively consistent with our main model estimation results. Moreover, we also see a similar trend that the price elasticity from *Alternative Model II* is larger than the one from the main search model. This additional model provides further support that in a highly differentiated market (e.g., hotel market), ignoring consumers' search information in the demand model can lead to an overestimation of the price elasticity.

***(3) Alternative Model III: Use the click data only (Click Model).***

On the other hand, we sometimes have access to only the publicly available click-through data (but no purchase information). Therefore, finding out, given only the publicly available data, how well our search model can predict consumers' click behavior is important. For this purpose, we consider a "click model" in the robustness test. In particular, we include only the click-sequence-related information in the likelihood function using the click data only, as shown below in Equation (G1). We estimate this click model using a similar simulated maximum likelihood approach based on only the click probability. We find that the estimated coefficients are qualitatively consistent with the main results. For more details on the estimates from the click model, we provide the results in Table G1, column 4 ("A3").

$$\begin{aligned}
\text{Likelihood} &= \prod_i \Pr(\text{all\_clicks}_i, \text{all\_noclicks}_i) \\
&= \prod_i \left\{ \prod_{r(n) \in S^{i, \text{clicked}}} \pi_{i, r(n)} * \prod_{r(m) \in S^{i, \text{unclicked}}} (1 - \pi_{i, r(m)}) \right\}.
\end{aligned} \tag{G1}$$

**(4) Alternative Model IV: Use both the click and the purchase data (Joint Probabilistic Model of Click and Purchase + Additional Search Cost Variables, But No Click Sequence Information).**

From *Alternative Models (I) – (III)*, we find that using only the click data or only the purchase data are likely to overestimate the price elasticity, and therefore it is important to consider both click and purchase decisions when modeling consumer preferences. However, it is not clear the improvement in the model performance is attributed to the advantage of our holistic search model or simply to the use of more data. To examine this issue, we consider another alternative model – a joint probabilistic model of both click and purchase. The major difference between this join probabilistic model and our main search model is that instead of capturing the sequence of clicks and allowing clicks to be interdependent, the join model assumes each click decision to be independent. Correspondingly, it models the click decisions independently as following a discrete choice process (e.g., Logit model). The results of this model are illustrated in Table G2, column 3 (“A4”).

We find that the estimation results are qualitatively consistent with our main findings. However, interestingly we find that although incorporating both click and purchase decisions information can improve the model estimation, the joint probabilistic model without considering the click sequence information can still lead to an overestimation of price elasticity (1.954 vs. 1.619,  $p < 0.05$ ). This result indicates that not only the final click or purchase decisions matter, but also the sequential click path of consumer search is critical in revealing consumer preferences. Failing to capture consumers’ search paths can lead to an overestimation of price elasticity in the online search market.



**Table G1. Estimation Results - Main Model and Alternative Models (I) & (III)**

Variable	Mean Effect (Std. Err) <sup>M</sup>	Heterogeneity (Std. Err) <sup>M</sup>	Mean Effect (Std. Err) <sup>A1</sup>	Heterogeneity (Std. Err) <sup>A1</sup>	Mean Effect (Std. Err) <sup>A3</sup>	Heterogeneity (Std. Err) <sup>A3</sup>
(Preferences)	$\bar{\alpha}, \bar{\beta}, \bar{\lambda}$	$\sigma_{\alpha}, \Sigma_{\beta}, \Sigma_{\lambda}$	$\bar{\alpha}, \bar{\beta}, \bar{\lambda}$	$\sigma_{\alpha}, \Sigma_{\beta}, \Sigma_{\lambda}$	$\bar{\alpha}, \bar{\beta}, \bar{\lambda}$	$\sigma_{\alpha}, \Sigma_{\beta}, \Sigma_{\lambda}$
PRICE <sup>(L)</sup>	-1.252* (.022)	.417* (.074)	-2.391* (.038)	1.064* (.082)	-1.744* (.029)	.472* (.088)
PAGE	-.239* (.003)	.080 (.133)	-.283* (.002)	.142 (.261)	-.206* (.003)	.117* (.140)
RANK	-.314* (.008)	.132* (.067)	-.341* (.008)	.138 (.211)	-.268* (.009)	.190* (.037)
CLASS	1.516* (.023)	.935* (.181)	1.882* (.012)	1.060* (.271)	2.057* (.020)	.574* (.128)
AMENITYCNT <sup>(L)</sup>	.146* (.034)	.066 (.070)	.212* (.051)	.059 (.126)	.137* (.015)	.056* (.022)
ROOMS <sup>(L)</sup>	.394* (.024)	.195 (.287)	.449* (.060)	.240 (.333)	.410* (.075)	.251 (.467)
EXTAMENITY	.165* (.036)	.041 (.046)	.207* (.049)	.041 (.107)	.199* (.051)	.046 (.039)
BEACH	1.539* (.028)	.561* (.099)	1.924* (.033)	.492* (.191)	1.227* (.077)	.388* (.104)
LAKE	-.663* (.116)	1.560* (.389)	-.745* (.081)	.974* (.267)	-.712* (.082)	1.568* (.235)
TRANS	1.336* (.140)	.192* (.064)	1.359* (.116)	.198 (.170)	1.503* (.182)	.193 (.216)
HIGHWAY	.447* (.093)	.068 (.061)	.464* (.080)	.057 (.109)	.374* (.092)	.053* (.011)
DOWNTOWN	1.198* (.061)	.471* (.093)	1.051* (.088)	.283 (.076)	.943* (.053)	.331 (.078)
CRIME	-.173* (.043)	.015 (.034)	-.189* (.041)	.036 (.067)	-.178* (.032)	.018 (.017)
RATING	2.661* (.015)	1.308* (.091)	2.361* (.017)	.926* (.083)	2.017* (.020)	1.334* (.092)
REVIEWCNT <sup>(L)</sup>	1.230* (.107)	.369* (.069)	1.102* (.128)	.405* (.057)	1.182* (.158)	.438* (.059)
STAFF	.139* (.027)	.034 (.088)	.142* (.021)	.031 (.081)	.147* (.022)	.033 (.095)
FOOD	.225* (.038)	.136* (.002)	.234* (.043)	.141* (.009)	.251* (.039)	.146* (.005)
BATHROOM	.290 (.271)	.060 (.103)	.278 (.259)	.082 (.122)	.242 (.277)	.091 (.118)
PARKING	.097* (.008)	.075* (.011)	.092* (.005)	.071* (.009)	.088* (.008)	.079* (.013)
BEDROOM	-.175 (.232)	.253 (.269)	-.164 (.241)	.277 (.256)	-.189 (.237)	.270 (.244)
FRONTDESK	.065 (.103)	.021 (.076)	.077 (.112)	.016 (.068)	.073 (.125)	.028 (.071)
BRAND	Yes					
(Search Cost)	$\bar{\gamma}$	$\Sigma_{\gamma}$	$\bar{\gamma}$	$\Sigma_{\gamma}$	$\bar{\gamma}$	$\Sigma_{\gamma}$
Search Base Cost	-7.511* (.089)	.971* (.176)	----	----	-4.041* (.092)	.932* (.241)
COMPLEXITY	.541* (.094)	.398* (.115)	----	----	.219 (.104)	.525 (.781)
SYLLABLES <sup>(L)</sup>	.678* (.115)	.721* (.106)	----	----	.582* (.165)	.633 (.958)
SPELLERR <sup>(L)</sup>	.329* (.082)	.033 (.101)	----	----	.192 (.226)	.053 (.283)
SUB	.196* (.045)	.057 (.229)	----	----	.141 (.123)	.070* (.011)
SUBDEV	.342* (.056)	.119 (.273)	----	----	.284* (.084)	.169* (.030)
Maximum LL	-405,418		-114,003		-352,359	
Price Elasticity	-1.619		-2.973		-2.183	

(L) Logarithm of the variable.

\* Statistically significant at 5% level.

M: Main Model.

A1: Alternative Model I (Use Purchase Data Only – Mixed Logit Model, with Actual Limited Consideration Set).

A3: Alternative Model III (Use Click Data Only – Click Model).

**Table G2. Estimation Results – Alternative Models (II) & (IV)**

Variable	Mean Effect (Std. Err) <sup>A2</sup>	Heterogeneity (Std. Err) <sup>A2</sup>	Mean Effect (Std. Err) <sup>A4</sup>	Heterogeneity (Std. Err) <sup>A4</sup>
(Preferences)	$\bar{\alpha}, \bar{\beta}, \bar{\lambda}$	$\sigma_{\alpha}, \Sigma_{\beta}, \Sigma_{\lambda}$	$\bar{\alpha}, \bar{\beta}, \bar{\lambda}$	$\sigma_{\alpha}, \Sigma_{\beta}, \Sigma_{\lambda}$
PRICE <sup>(L)</sup>	-2.036* (.034)	.989* (.086)	-1.663* (.027)	.421* (.052)
PAGE	-.240* (.004)	.124 (.248)	-.190* (.009)	.080 (.163)
RANK	-.309* (.004)	.130 (.186)	-.245* (.006)	.177* (.055)
CLASS	1.767* (.023)	1.020* (.244)	1.606* (.025)	.824* (.155)
AMENITYCNT <sup>(L)</sup>	.201* (.045)	.050 (.105)	.144* (.038)	.059 (.075)
ROOMS <sup>(L)</sup>	.425* (.037)	.211 (.309)	.320* (.035)	.159 (.292)
EXTAMENITY <sup>(L)</sup>	.193* (.051)	.043 (.128)	.174* (.043)	.036 (.049)
BEACH	1.958* (.023)	.491* (.199)	1.868* (.012)	.491* (.088)
LAKE	-.886* (.076)	1.026* (.202)	-.797* (.109)	1.318* (.361)
TRANS	1.161* (.089)	.186 (.195)	1.021* (.183)	.172* (.072)
HIGHWAY	.412* (.072)	.056 (.110)	.335* (.099)	.030 (.086)
DOWNTOWN	.991* (.083)	.263 (.083)	.918* (.055)	.389* (.082)
CRIME	-.159* (.041)	.034 (.061)	-.166* (.023)	.014 (.041)
RATING	2.215* (.011)	.883* (.098)	2.794* (.019)	.972* (.097)
REVIEWCNT <sup>(L)</sup>	1.077* (.114)	.413* (.044)	1.208* (.192)	.408* (.079)
STAFF	.146* (.022)	.033 (.081)	.132* (.025)	.033 (.082)
FOOD	.241* (.045)	.145* (.007)	.230* (.040)	.133* (.005)
BATHROOM	.286 (.267)	.084 (.121)	.292 (.266)	.062 (.104)
PARKING	.091* (.004)	.074* (.007)	.090* (.003)	.079* (.012)
BEDROOM	-.160 (.240)	.282 (.273)	-.171 (.227)	.254 (.260)
FRONTDESK	.079 (.111)	.017 (.067)	.066 (.098)	.022 (.075)
COMPLEXITY	-.130* (.022)	.063 (.331)	-.149* (.025)	.093 (.443)
SYLLABLES <sup>(L)</sup>	-.175* (.041)	.271* (.048)	-.171* (.037)	.246* (.033)
SPELLERR <sup>(L)</sup>	-.083* (.003)	.019 (.032)	-.081* (.004)	.034 (.054)
SUB	-.109* (.008)	.038* (.006)	-.136* (.011)	.040* (.010)
SUBDEV	-.139* (.019)	.073 (.054)	-.165* (.012)	.067* (.033)
BRAND	Yes			
Maximum LL	-119,752		-388,106	
Price Elasticity	-2.393		1.954	

(L) Logarithm of the variable.

\* Statistically significant at 5% level.

A2: Alternative Model II (Mixed Logit Model, with Actual Limited Consideration Set + Additional Search Cost Variables).

A4: Alternative Model IV (Joint Probabilistic Model of Click and Purchase, Using Both Click and Purchase Data + Additional Search Cost Variables, But No Click Sequence Information).

## Online Appendix H.

### Breakdown Analysis for Predicted Search Engine Revenues

As shown in Table 4 from the policy experiments, under various different scenarios product search engines will experience a revenue increase when providing different sets of information on the search summary page. To examine where the revenue increase comes from (i.e., existing consumers or better market coverage), we conducted an additional analysis on the breakdown of the revenue in the simulation.

More specifically, the predicted total search engine revenues can be computed as follows:

$$\text{Predicted Total Revenues} = \sum_{\text{All Sessions}} \sum_{\text{All Hotels}} (\text{Price} * \text{Predicted Purchase Probability} * \text{Commission Rate } 20\%). \quad (\text{H1})$$

Based on Equation (G1), we were able to separately compute the predicted revenues from all hotels for each session  $i$  as follows:

$$\sum_{\text{All Hotels in Session } i} (\text{Price} * \text{Predicted Purchase Probability} * \text{Commission Rate } 20\%). \quad (\text{H2})$$

Then, we categorized all sessions from our observed data into two types: 1) sessions without any purchase, and 2) sessions with purchase. We separately computed the total predicted revenues from each of the two categories, and then compare the difference in the predicted revenues and the observed revenues.

We found that the revenue increase (i.e., increase under all scenarios except for “Existing - Price”) came from both types of sessions. In particular, for sessions without any purchase, the predicted revenue increase indicates a potential increase in market coverage (i.e., consumers who did not purchase in the past become likely to purchase). This is consistent with our model intuition that providing additional product information on the search summary page can reduce the potential error in consumers’ expectation towards product utility and search costs before click. As a consequence, consumers are more likely to click on the best set of products that will provide them the highest utility. Hence, the maximum utility discovered from this click-generated consideration set is more likely to exceed the utility of the outside good. As a result, consumers are less likely to miss a good-value deal (i.e., leave without purchase).

For sessions with purchase, the predicted revenue increase indicates a potential increase in spending from the existing consumers (consumers who were less likely to purchase in the past become more likely to purchase; or consumers who were likely to spend less in the past become likely to spend more). We provide more details on the breakdown of the revenue increase in Table H1 below for your convenience.

**Table H1. Breakdown Analysis Results on Search Engine Revenue Increase**

	Overall Search Engine Revenue	Revenue Increase		
		All Sessions	Sessions With Purchase	Sessions Without Purchase
<i>Existing</i>	\$452,781	\$0	\$0	\$0
<i>Existing + Location Information</i>	\$553,136	\$100,355	\$91,323	\$9,032
<i>Existing + Service Information</i>	\$467,369	\$14,588	\$11,962	\$2,626
<i>Existing – Price Information</i>	\$420,132	\$-32,649	\$-28,731	\$-3,918
<i>Existing + Review Information (Text Features)</i>	\$507,160	\$54,379	\$47,853	\$6,526
<i>Existing + Review Information (Topic Entropy)</i>	\$490,063	\$37,282	\$30,572	\$6,711

In addition, we also found that the revenue increase occurs for both hotels that have been purchased in the past (i.e., existing hotels) and hotels that have not been purchased in the past (i.e., new hotels). This finding provides additional supports that with carefully designed information on search summary page, search engine can improve the market coverage of consumers as well as the diversity of products consumed, which can lead to a potential increase in consumer surplus.