

# Modeling Online Browsing and Path Analysis Using Clickstream Data

Alan Montgomery  
*Associate Professor*  
**Carnegie Mellon University**  
Graduate School of Industrial Administration

*With Shibo Li, Kannan Srinivasan, John Liechty*

**e-mail:** alan.montgomery@cmu.edu  
**web:** <http://www.andrew.cmu.edu/user/alm3>

*JSM Meetings, 4 August 2003*

© 2003 by Alan Montgomery, All rights reserved.

## Problem Overview

- Path Analysis
  - Infer consumer's purchase orientation from the navigation path
  - Impact of web design (i.e. hypertext links)
  - Impact of ads, price information, and free shipping promotion
- Purchase Conversion
  - Predict conversion dynamically
  - When and how to target online consumers



## **Outline**

- Clickstream
- Data
- Model
- Results
- Purchase Conversion
- Conclusions

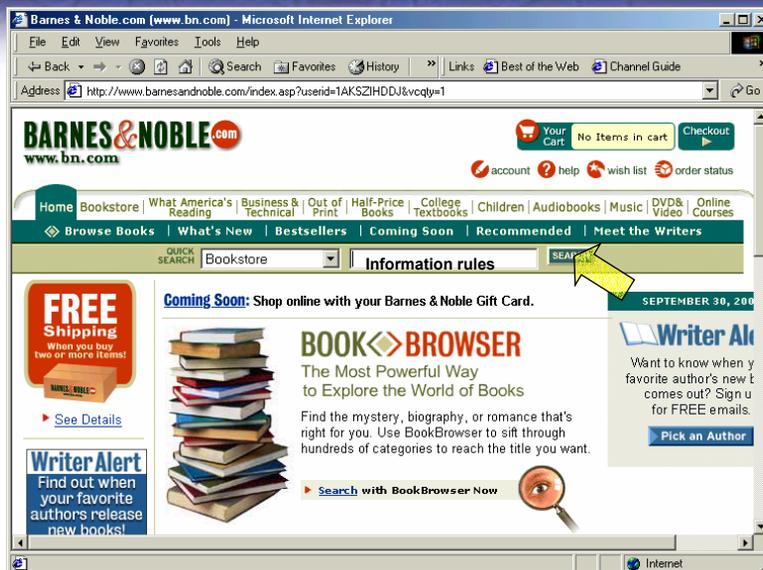


## **Clickstream**

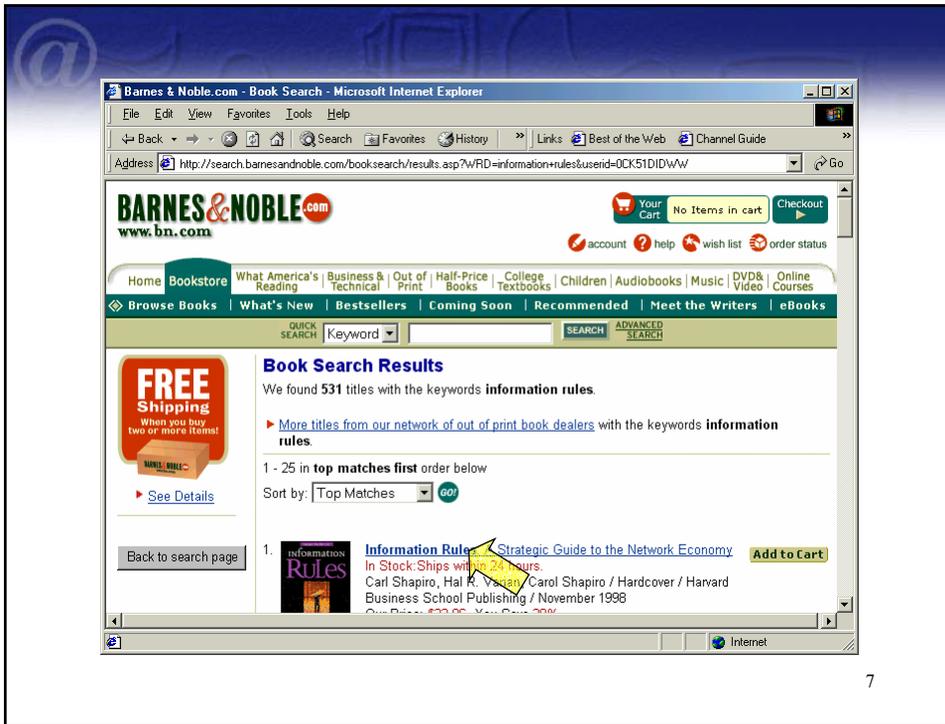
# What is clickstream data?

- A record of an individual's movement through time at a web site
- Contains information about:
  - Time
  - URL content
  - User's machine
  - Previous URL viewed
  - Browser type

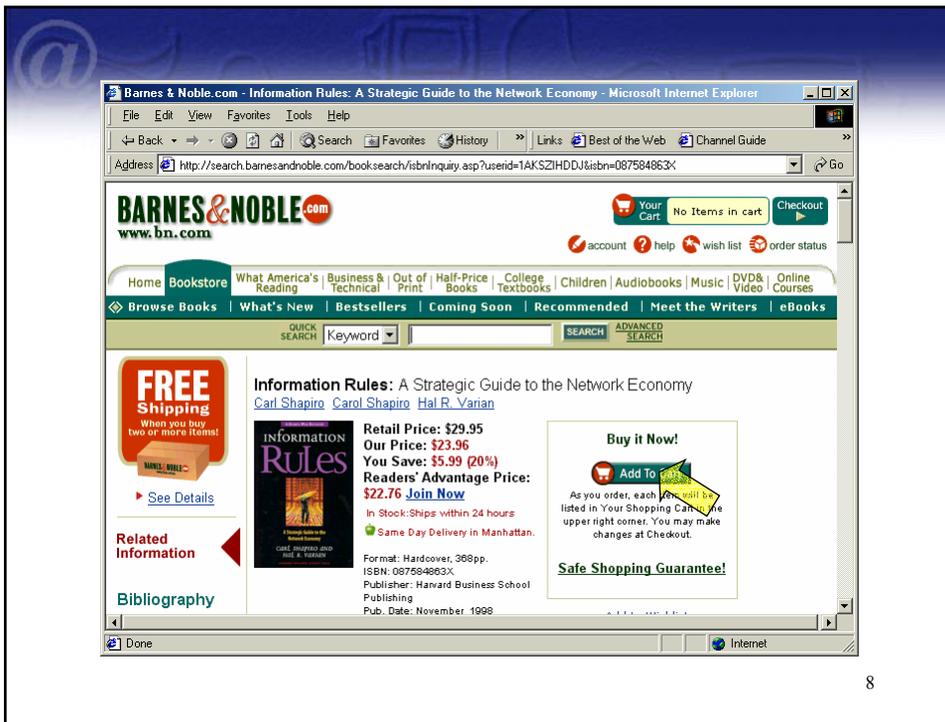
5



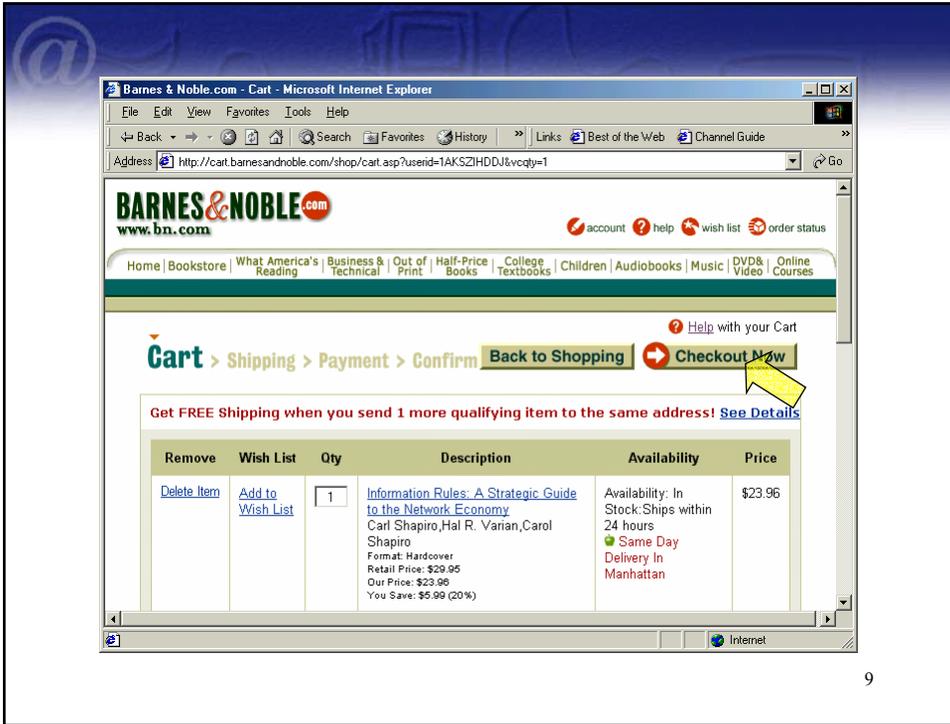
6



7



8



9

## Predicting Purchase Conversion

Home

Category

Product

Shopping Cart

What is the chance of this user making a purchase during this session?

}

1<sup>st</sup> viewing = 7%

2<sup>nd</sup> viewing = 14%

3<sup>rd</sup> viewing = 20%

4<sup>th</sup> viewing = 60%

10



## **Problems with Clickstream Data**

- Huge
  - Terabytes of data
- Collection Issues
  - Server versus client
  - Privacy and identifiability
- Unstructured
  - Textual
- Decision Problem
  - Redesign web site
  - Understand usage and implication on consumer behavior

11



## **Data**

# Data

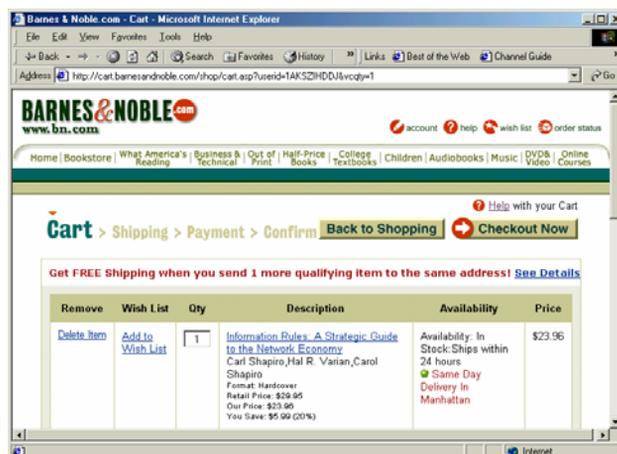
- Comscore Media Metrix Panelists from April 2002 who visit barnesandnoble.com
  - 1,160 users
  - 1,659 Sessions, 8.75 Viewings per session
  - 14,512 viewings
  - 114 purchases ⇒ 7% conversion rate
- Collect HTML content using Perl script
  - Parse the pages for content information, like the presence of a banner ad, price, ...

Path analysis can be done using server logs from bn.com directly

13

# Page Category Examples

- Home
- Category
- Product
- Information
- Account
- Shop. Cart



14

## Categorization of Pages

Abbr- eviation	Category Name	Description	% of Page Requests
H	Home	Home page, common starting page for B&N visitors	6.5%
A	Account	User account pages that allow a user to sign in, change address, and review order status	6.4%
C	Category	Pages that have list of items (links like cooking, fiction, etc. from the home page) or results of a book search	32.9%
P	Product	Pages with detailed product information, item description, price information, availability, and reviews	16.4%
I	Information	Pages with shipping, order status, and popup advertisements such as "Free Shipping"	29.1%
S	Shopping Cart	Pages related to the shopping cart, such as reviewing the cart, deleting items, entering purchase information	7.0%
O	Order	Confirmation page that denotes order has taken place	1.0%
E	Enter/Exit	Non B&N pages used to denote the beginning or end of a session	.8%

Table 3. Categorization scheme for B&N web pages for 9,180 unique pages requests.

15

## Transformed Clickstream

	Time	URL	Category	Abbr.
1	8:36:11pm	/promo/coupon/popup/fs_usa_popup.asp?userid=xxx	Information	I
2	8:36:29pm	/booksearch/results.asp?wrld=70%2d215&userid=xxx	Category	C
3	8:36:48pm	/booksearch/results.asp?userid=xxx&mscssid=yyy&wrld=70%2d215&opr=a&sort=p	Category	C
4	8:37:14pm	/booksearch/isbninquiry.asp?userid=xxx&mscssid=yyy&isbn=0072134445	Product	P
5	8:38:10pm	/booksearch/results.asp?userid=xxx&mscssid=yyy&wrld=70%2d215&opr=a&sort=p	Category	C
6	8:44:32pm	/textbooks/booksearch/isbninquiry.asp?userid=xxx&mscssid=yyy&isbn=0619034971	Product	P
7	8:55:12pm	/promo/coupon/popup/fs_usa_popup.asp?userid=xxx	Information	I
8	8:55:24pm	/booksearch/results.asp?wrld=70%2d215&userid=xxx	Category	C
9	8:55:36pm	/booksearch/results.asp?userid=xxx&mscssid=yyy&wrld=70%2d215&opr=a&sort=p	Category	C
10	8:56:37pm	/shop/signin.asp?userid=xxx&mscssid=yyy	Account	A
11	8:58:16pm	/booksearch/results.asp?userid=xxx&mscssid=yyy&wrld=70%2d215&opr=a&sort=p	Category	C
12	8:58:40pm	/booksearch/isbninquiry.asp?userid=xxx&mscssid=yyy&isbn=0072224983	Product	P
13	8:59:21pm	/shop/cart.asp?userid=xxx&mscssid=yyy	Account	C
14	9:01:26pm	Exit	Exit	E

16

## Some Sample User Sessions

User	Path		
Browsers	1 2 3 4 5	ICCCCCCCCCPCCPCCCCCCCCCCCCCCCCCCCCCCCCCCCCCE IHHE IE IHICPPPCE IHIIICIE	
	Buyers	6 7 8 9 10	HIAAAAIHIIICICICICICIPPIPPPIPIICCSIIIPPPPIPSISISSOIIIIHE HCCPPPCPCCCCCCCCPSCSPCCPCPCCCCCSAAAAAAAAAASSOIIIIASCCCE IICICPCPPPCPCICPCCCPCPPPIPSIIAASSIIIOIE IISLASSIOIE IPPPSASSSSOIAAAHCCPCCCCCE

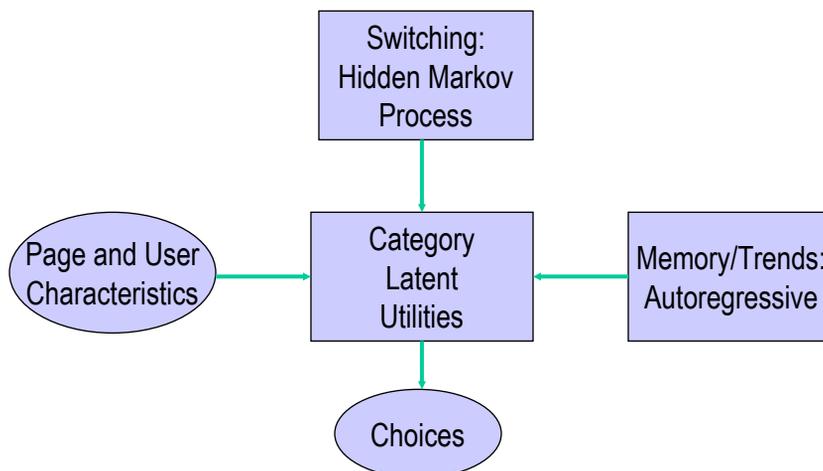
## Model

## Model Elements

- Dependent Variable is a Categorical Time Series
  - Choice modeled using Multinomial Probit (MNP Model)
- Two Types of Time Series Behavior:
  - Continuous trends: VAR(1) Process
  - Abrupt changes: Mixture Process with a Hidden Markov Model
- Covariates
  - Web design and marketing mix
  - Browsing behavior
  - Demographics
- Consumer Heterogeneity
  - Hierarchical Bayesian Approach

19

## Choice Model of Browsing



20

## Dynamic Multinomial Probit Model with A Mixture Process

- Continuous-Time Hidden Markov Chain  $D_i$  with  $\lambda_s$ , the Waiting Time Intensity Parameter

$$P_i = \begin{bmatrix} 0 & P_{12i} & \cdots & P_{1Si} \\ P_{21i} & 0 & \cdots & P_{2Si} \\ \vdots & \vdots & \ddots & \vdots \\ P_{S1i} & P_{S2i} & \cdots & 0 \end{bmatrix}, \mathbf{v}_i = \begin{bmatrix} v_{1i} \\ v_{2i} \\ \vdots \\ v_{Si} \end{bmatrix}$$

- “Time” can refer to either per viewing, per session, or per user (which is a latent mixture model)

21

## Dynamic Multinomial Probit Model with A Mixture Process

- Observation Equation:

$$Y_{ict} \begin{cases} = 1 & \text{if } \max(U_{ikt}) < 0 \\ = c & \text{if } U_{ict} \geq \max(U_{ikt}) \end{cases}$$

- Model Setup:

Memory, Covariates, Hidden Markov States, Latent Utility

$$U_{ict} = \Phi_{ic\mathcal{S}} U_{i(t-1)} + \gamma_{ic\mathcal{S}}' X_{it} + \varepsilon_{ic\mathcal{S}}$$

with probability  $p_s$  such that  $\sum_{s=1}^S p_s = 1$

22

## Understanding State Model



- Continuous Time Model

$$f(D_i | P_i, \lambda_i, v_i) = \left( \prod_s v_i^{I(D_{i0}=s)} \right) \left( \prod_{s \neq 1} P_{is}^{N_{is}} \right) \left( \prod_s \lambda_{is}^{M_{is}} \right) \exp \left[ - \sum_s \left( \int_0^T \lambda_{is} I(D_{it} = s) dt \right) \right]$$

- Per viewing approximates First-Order Markov Model

$$\mathbf{P} = \begin{bmatrix} P_{11} & P_{12} \\ P_{21} & P_{22} \end{bmatrix} = \begin{bmatrix} 1 - \exp\{-\lambda_1\} & \exp\{-\lambda_1\} \\ \exp\{-\lambda_2\} & 1 - \exp\{-\lambda_2\} \end{bmatrix}$$

23

## Other Model Elements

- Hierarchical Bayesian Model
  - Parameters are “regressed” upon user characteristics
- Model includes many special cases:
  - Multinomial Probit Model (1 state)
  - Latent Mixture Model (Multiple states, zero-order markov)
  - Zero-Order Markov Process (only intercepts)
  - Approximates First-Order Markov Process with VAR(1)

24

## Predictive Results

## Page Transition Models

State Time	State Process	Number of States	Log Marginal Density	In-Sample Hit Rate (%)	Out-of-Sample Hit Rate (%)
Page	Zero-Order	1	-9378.1	72.05	65.40
		2	-9016.9	79.44	71.42
		3	-9064.0	80.34	70.56
	First-Order	1	-8545.4	83.23	79.95
		2	-8428.3	89.71	83.15
		3	-8474.0	89.97	81.14

**Best**

## Session Transition Models

State Time	State Process	Number of States	Log Marginal Density	In-Sample Hit Rate (%)	Out-of-Sample Hit Rate (%)
Session	Zero-Order	1	-9376.1	73.17	61.56
		2	-9051.0	77.90	70.48
		3	-9097.7	78.76	66.14
	First-Order	1	-8573.5	83.05	73.57
		2	-8464.9	88.44	81.48
		3	-8487.0	88.73	78.42

27

## User Transition Models

State Time	State Process	Number of States	Log Marginal Density	In-Sample Hit Rate (%)	Out-of-Sample Hit Rate (%)
User	Zero-Order	1	-9411.1	64.38	61.50
		2	-9124.2	70.04	64.12
		3	-9193.8	70.85	63.99

Table 7. Measures of fit for various dynamic probit models.

28

## Alternative Model Specifications

Model	Log Marginal Density	In-Sample Hit Rate (%)	Out-of-Sample (%)
Zero-Order Markov Model (1 State)	-20410.4	20.48	12.62
Zero-Order Markov Model (2 States)	-19458.3	28.18	19.02
First-Order Markov Model (1 State)	-16444.5	56.06	51.59
First-Order Markov Model (2 States)	-16076.0	58.61	52.08
Latent Class Model (1 State)	-17849.2	35.47	30.78
Latent Class Model (2 States)	-17673.9	44.29	40.21
Latent Class Model (3 States)	-17722.3	45.29	36.14
Independent	-19086.4	33.23	30.35
Only-Intercept	-19335.9	29.37	23.12
Intercept + VAR	-13768.4	71.13	64.38

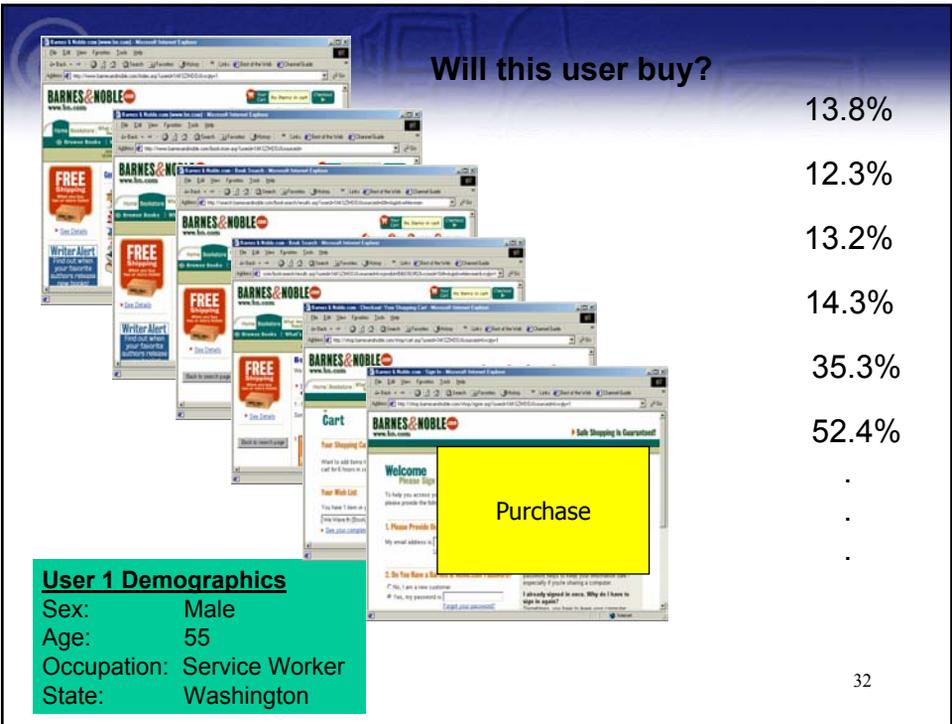
29

## Summary of Findings

- Banner ads tend to encourage browsing oriented individuals but discourage continued browsing from deliberation oriented users
- Including price information tends to discourage browsing oriented users but encourage deliberation oriented users
- More home links leads to less use of home page by browsers and more by deliberation oriented users
- Weekend users more likely to surf
- Browsers who have purchased before are more likely to order
- Rich set of relationships in dynamic browsing behavior

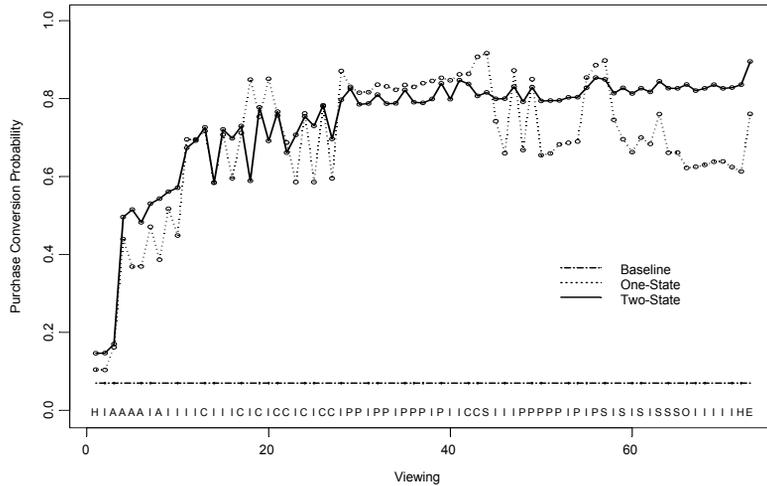
30

# Purchase Conversion



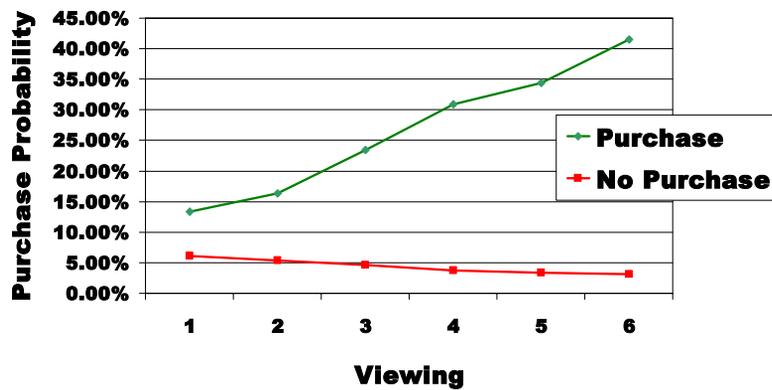


## User 2: Purchase Conversion Predictions



35

## Predicting Purchase Conversion

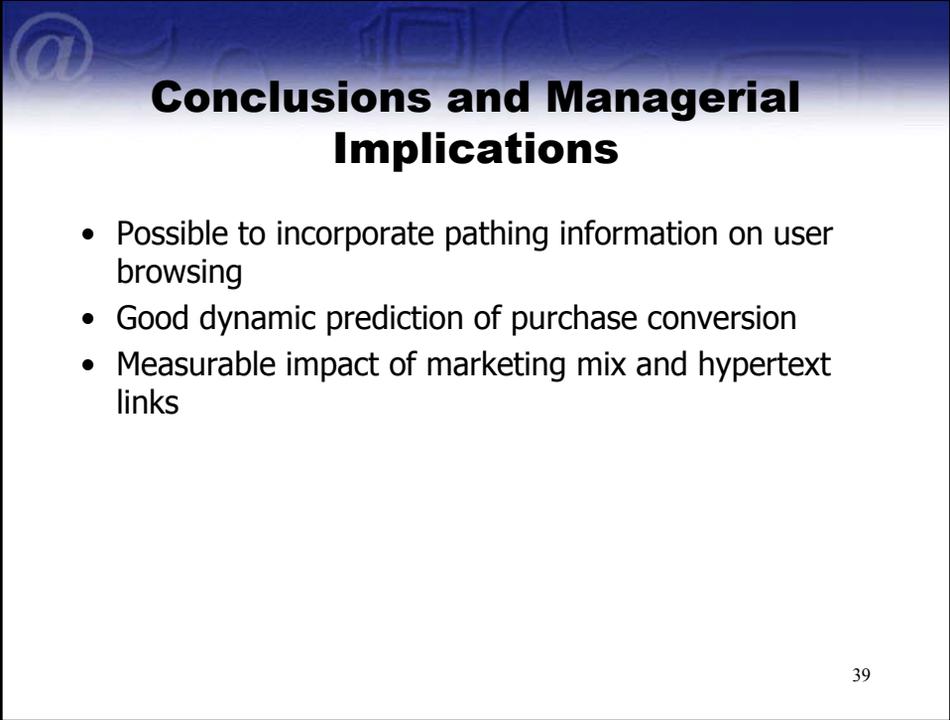


36

## Predicted Purchase Conversions

Sample	Session Type	Number of Sessions	Forecast Origin/Number of viewings during session					
			1	2	3	4	5	6
Estimation	Purchase	83	13.3% (0.48)	16.3% (0.52)	23.4% (0.60)	30.9% (0.65)	34.4% (0.67)	41.5% (0.69)
	No Purchase	1129	6.1% (0.33)	5.4% (0.32)	4.6% (0.30)	3.7% (0.27)	3.4% (0.26)	3.1% (0.25)
	All	1212	6.6% (0.35)	6.1% (0.34)	5.9% (0.33)	5.6% (0.33)	5.5% (0.32)	5.7% (0.33)
Holdout	Purchase	31	10.4% (0.97)	12.8% (1.06)	15.2% (1.14)	18.0% (1.21)	19.1% (1.24)	21.2% (1.29)
	No Purchase	416	6.9% (0.80)	5.5% (0.72)	5.1% (0.70)	4.2% (0.63)	3.5% (0.58)	3.2% (0.56)
	All	447	7.2% (0.82)	5.9% (0.75)	5.8% (0.74)	5.1% (0.70)	4.6% (0.66)	4.4% (0.65)

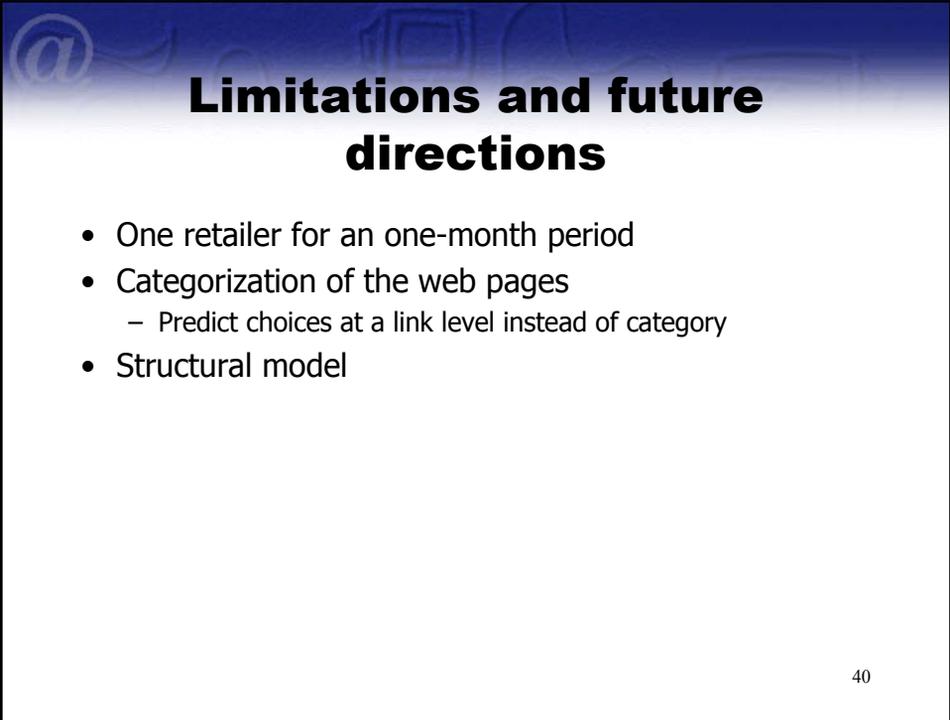
## Conclusions



## **Conclusions and Managerial Implications**

- Possible to incorporate pathing information on user browsing
- Good dynamic prediction of purchase conversion
- Measurable impact of marketing mix and hypertext links

39



## **Limitations and future directions**

- One retailer for an one-month period
- Categorization of the web pages
  - Predict choices at a link level instead of category
- Structural model

40