



Predicting Consumer Behavior using Clickstream Analysis

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

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

WebShop 2003, Univ of Maryland
13 June 2003

© 2003 by Alan Montgomery, All rights reserved.



Outline

- What is clickstream data?
- User Profiling
 - What does 'what you view' say about 'who you are?'
- Path Analysis
 - What does 'what you view' say about 'what you want?'
- Analyzing Textual Information in Clickstream Data
- Conclusions



Defining Clickstream Data

The raw input for web mining



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



Sources of clickstream data

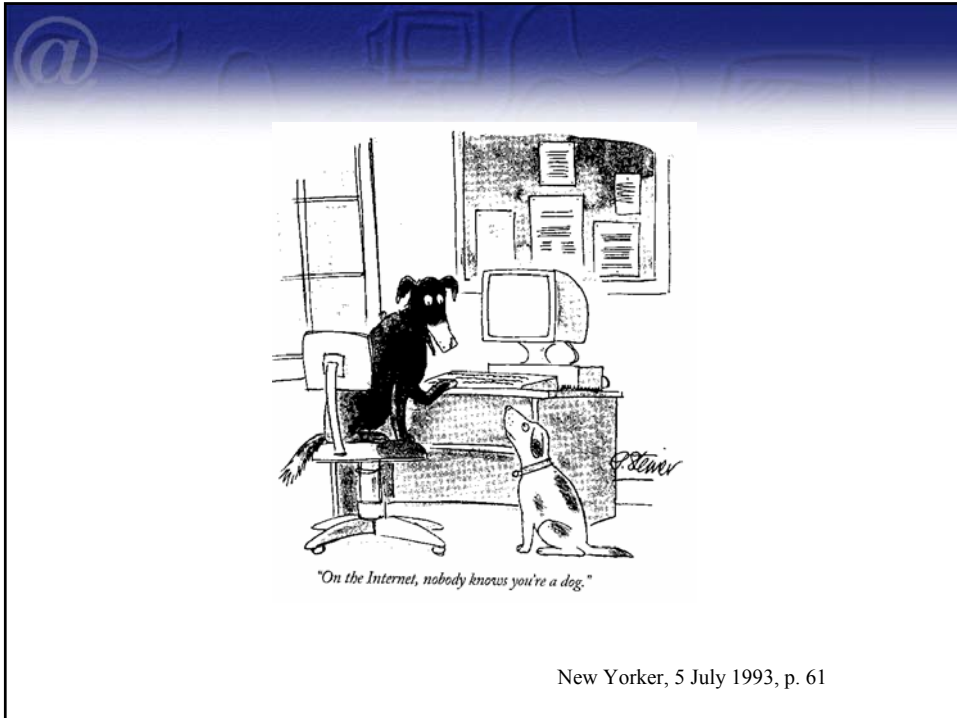
- Web Servers
 - Each hit is recorded in the web server log
- Media Service Providers
 - DoubleClick, Flycast
- ISP/Hosting Services
 - AOL, Juno, Bluelight.com
- Marketing Research Companies
 - ComScore Media Metrix and NetRatings

5



User Profiling

What does 'where you go' say
about 'who you are'?



Is this user male or female?

User visits the following five sites in the Doubleclick network

95% probability that user is female

8

Bayesian updating formula

Test the hypothesis that a user is female by updating the current guess using new information

$$\bar{p} = \frac{p \cdot \bar{p}}{p \cdot \bar{p} + (1 - p)(1 - \bar{p})}$$

New information
Old probability

New probability

Female
Male

9

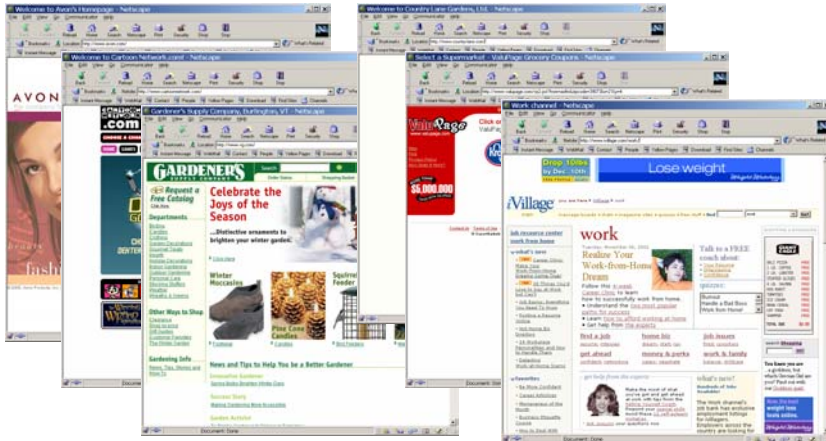
Probability user is female

	Probability a Female Visits the site	Probability visitor is Female given visits to
Overall Internet	45%	45.0%
cbs.com	54%	49.0%
ivillage.com	66%	65.1%
libertynet.org	63%	76.0%
nick.com	57%	80.8%
onlinepsych.com	83%	95.4%

Best Guess

10

What can we learn?



11

A Full Month of Browsing Example

% of female visitors during one month (Media Metrix):

48%	aol.com	63%	libertynet.org
64%	astronet.com	39%	lycos.com
75%	avon.com	27%	netradio.net
52%	blue-planet.com	57%	nick.com
56%	cartoonnetwork.com	59%	onhealth.com
54%	cbs.com	83%	onlinepsych.com
76%	country-lane.com	44%	simplenet.com
47%	eplay.com	76%	thriveonline.com
41%	halcyon.com	59%	valupage.com
70%	homearts.com	71%	virtualgarden.com
66%	ivillage.com	66%	womenswire.com

99.97% probability that user is female

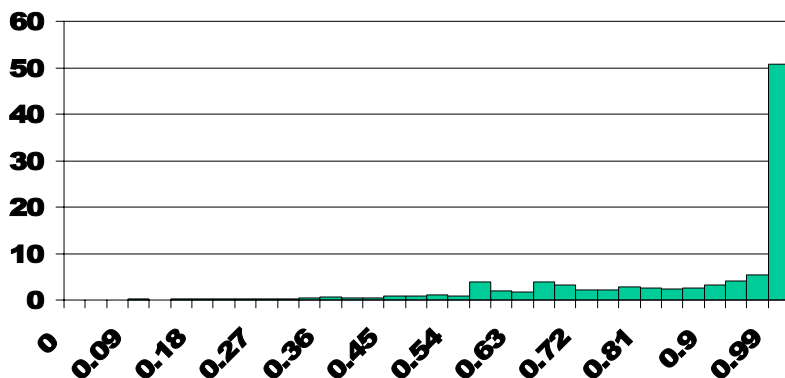
12

Results

- Analysis shows that there is a >99% probability this user is female.
 - Using only DoubleClick sites the probability is 95%.
- Using all user data for one month:
 - 90% of men are predicted with >80% confidence (81% accuracy)
 - 25% of women are predicted with >80% confidence (96% accuracy)

13

Probabilities of Predicting Male Users



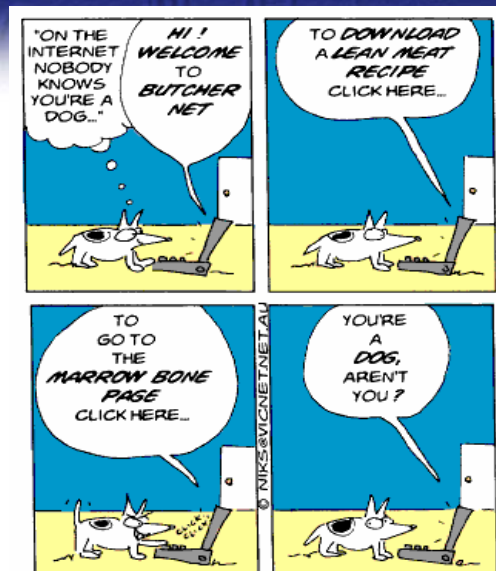
14

Key Points of User Profiling


We can identify 'who you are' from 'where you go'

- What the user views on the web reveals their interests and preferences
 - We can personalize the web experience without explicitly requiring customers to login and identify themselves
- Browsing and product choices can reveal key information about interest and price sensitivity
- Requires marketers to be smarter in designing their websites and analyzing their information. Big profitability gains if this is done correctly.

15




[Http://www.moreinfo.com/au.cranlerma/fo12.htm](http://www.moreinfo.com/au.cranlerma/fo12.htm)



Path Analysis

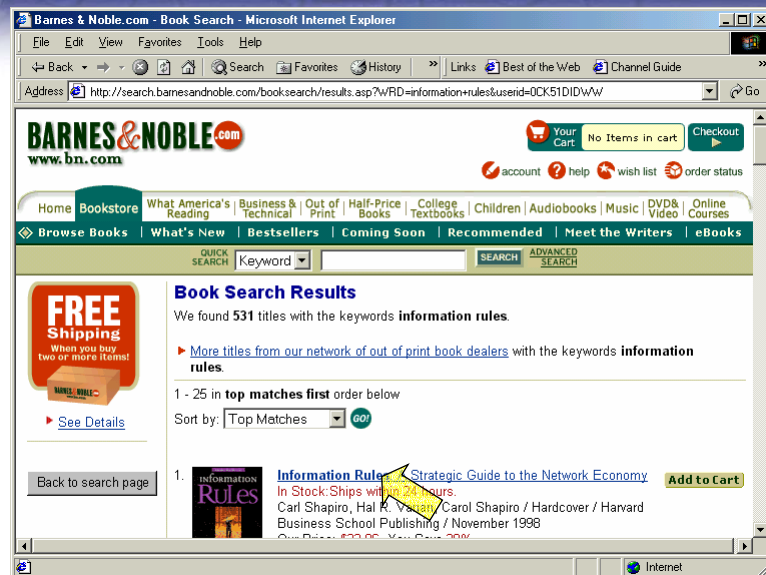
What does a user's web navigation path say about purchase conversion or a user's goals?



Clickstream Example #1



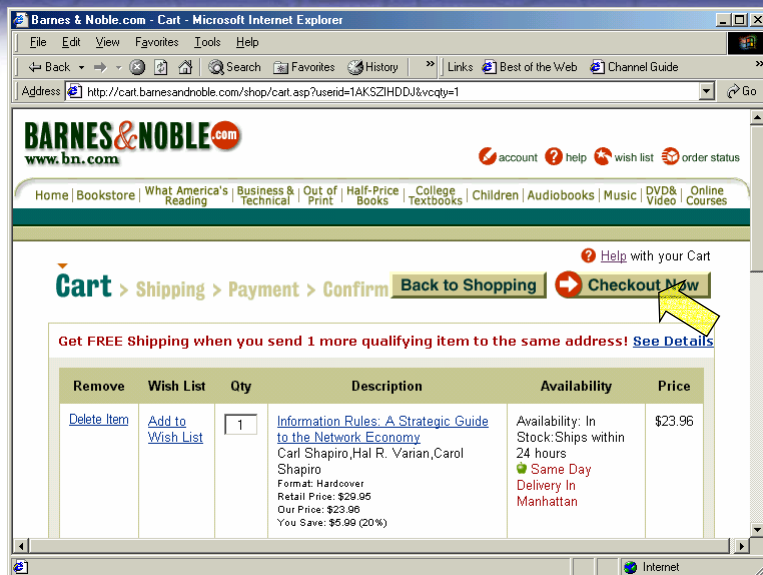
19



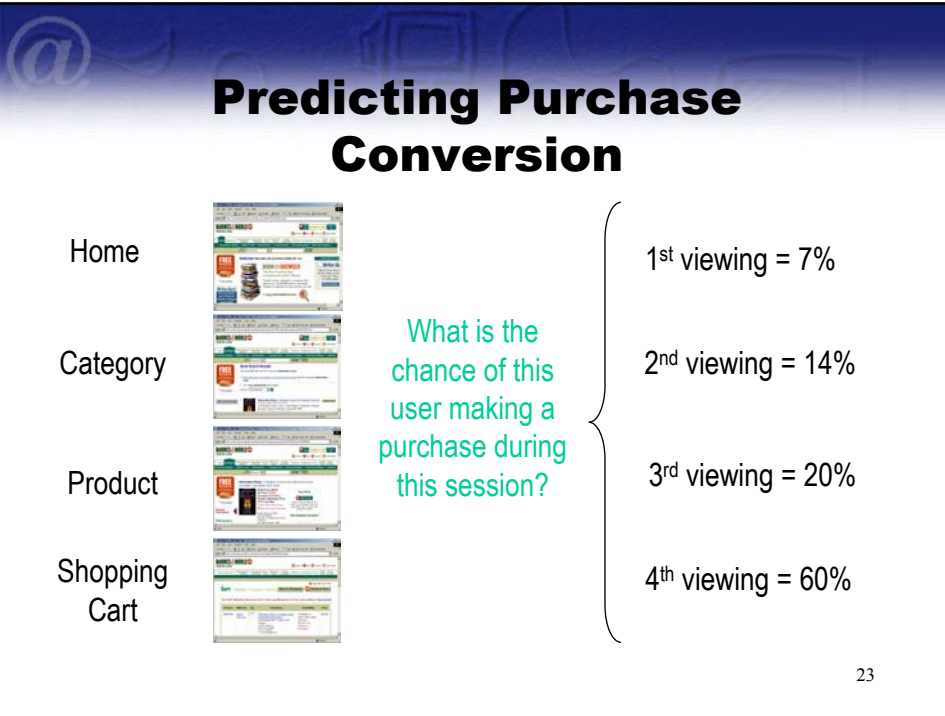
20



21



22



Clickstream Example #2

Will this user buy?

- {Home}
- {Category}
- {Category}
- {Category}
- {Shop Cart}
- {Account}

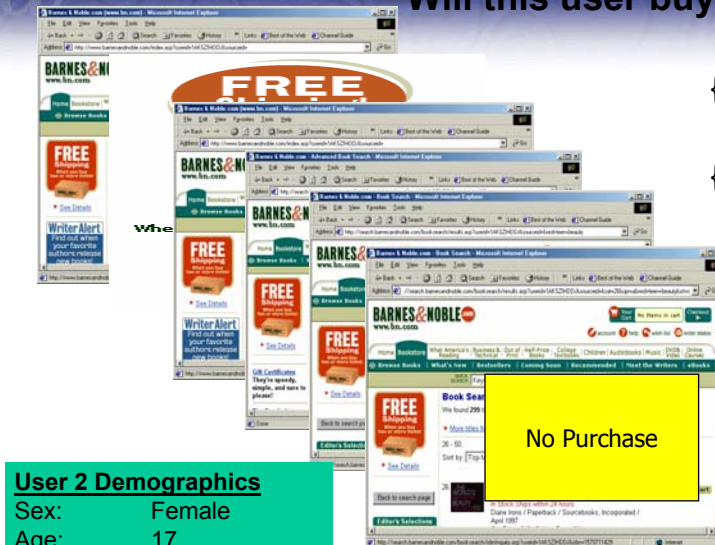
User 1 Demographics

Sex: Male
 Age: 55
 Occupation: Service Worker
 State: Washington

Clickstream Example #3

Will this user buy?

{Home}
{Information}
{Home}
{Information}
{Category}
{Category}



User 2 Demographics

Sex: Female
Age: 17
Occupation: Student
State: Virginia

27

Identifying Browsing Patterns

Categorizing Pages

Abbr	Category	Description
H	Home	Home page
A	Account	User account pages
C	Category	Page with list of products
P	Product	Product information pages
I	Information	Shipping, order status, etc
S	ShoppingCart	Pre-order pages
O	Order	Confirmation/purchase page
E	Enter/Exit	Non B&N pages

29

Some Sample User Sessions

User	Path
Browsers	1 ICCCCCCCCCPCGCCCCCCCCCCCCCCCCCCCCCCCCCCCE
	2 IHHE
	3 IE
	4 IHICPPCE
	5 IHIIICIE
Buyers	6 HIAAAIAIIIIICICICICICICIPPIPPPIPIICCSIIIPPPPIPIPSISISSOIIIIHE
	7 HCCPPPCPCCCCCCCCPSCSPCCPCPCCCCCSAAAAAAAAAASSOIIIIISACCCE
	8 IIICPCPPPCICICIPCCPCPPPIPSIIAASSIIIOIE
	9 IISASSIOIE
	10 IPPPPASSSSOIAAAHCCPCCCCCE

30

Probability of Viewing a Page

Category	Purchaser	Browser	Odd Ratio
Home	1%	9%	1/9
Account	13%	4%	3/1
Category	27%	35%	.8/1
Product	17%	17%	1/1
Information	24%	33%	.8/1
Shopping Cart	15%	2%	7/1
Purchase	3%	0%	Inf

31

Transition Matrix

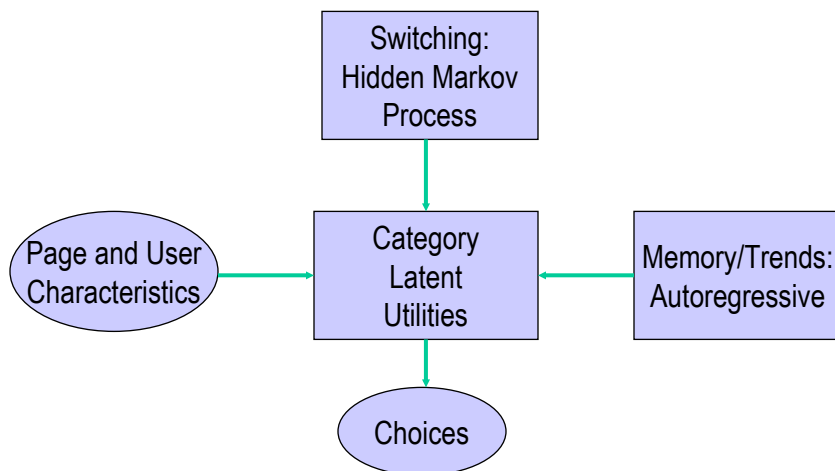
		Category of Previous Viewing							
Category of Current Viewing	Category	Home	Account	Category	Product	Inform.	ShopCart	Order	Exit
	Home	.23	.01	.01	.01	.10	.02	0	.16
	Account	.01	.69	.01	.01	.02	.15	0	.01
	Category	.17	.02	.60	.31	.15	.05	0	.16
	Product	.01	0	.20	.43	.10	.05	0	.05
	Information	.25	.06	.08	.12	.46	.15	.87	.61
	Shop. Cart	.01	.16	.01	.03	.02	.45	.13	.01
	Order	0	0	0	0	0	.10	0	0
	Exit	.32	.06	.09	.09	.14	.02	0	0
	Marginal	.06	.05	.32	.17	.23	.05	.01	.11
Initial Prob.		.16	.02	.16	.06	.60	.01	0	0

Table 6. Sample transition matrix for categories of viewings. (Notice that the columns sum to one, and there are a total of 14,512 observations.)

32

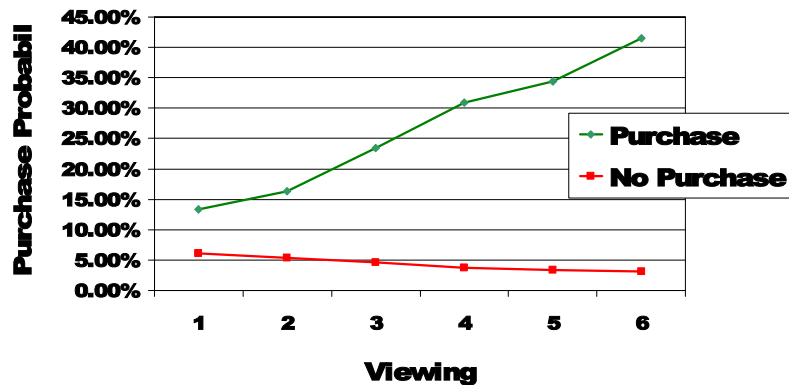
Methodology

Choice Model of Browsing





Predicting Purchase Conversion



37

Key Points of Path Analysis

We can infer 'what you want' from 'what you view'

- The path a user takes reveals goals and interests
 - We look at pages we are interested in
 - Avoid those pages that are irrelevant
- Path Analysis indicates we can intervene before a non-purchaser leaves the site
- Presenting promotional information to purchasers is distracting, but increases conversion for surfers
- Show the right information at the right time

38

Text Classification

Categorizing Web Viewership Using Statistical Models of Web Navigation and Text Classification



Information Available

Clickstream Data

- Panel of representative web users collected by Jupiter Media Metrix
- Sample of 30 randomly selected users who browsed during April 2002
 - 38k URLs viewings
 - 13k unique URLs visited
 - 1,550 domains
- Average user
 - Views 1300 URLs
 - Active for 9 hours/month

Classification Information

- Dmoz.org - Pages classified by human experts
- Page Content - Text classification algorithms from Comp. Sci./Inform. Retr.

41

Dmoz.org

- Largest, most comprehensive human-edited directory of the web
- Constructed and maintained by volunteers (open-source), and original set donated by Netscape
- Used by Netscape, AOL, Google, Lycos, Hotbot, DirectHit, etc.
- Over 3m+ sites classified, 438k categories, 43k editors (Dec 2001)

Categories

1. Arts
2. Business
3. Computers
4. Games
5. Health
6. Home
7. News
8. Recreation
9. Reference
10. Science
11. Shopping
12. Society
13. Sports
14. Adult

42



Problem

- Web is very large and dynamic and only a fraction of pages can be classified
 - 147m hosts (Jan 2002, Internet Domain Survey, isc.org)
 - 1b (?) web pages+
- Only a fraction of the web pages in our panel are categorized
 - 1.3% of web pages are exactly categorized
 - 7.3% categorized within one level
 - 10% categorized within two levels
 - 74% of pages have no classification information

43



Text Classification

Background

- Informational Retrieval
 - Overview (Baeza-Yates and Ribeiro-Neto 2000, Chakrabarti 2000)
 - Naïve Bayes (Joachims 1997)
 - Support Vector Machines (Vapnik 1995 and Joachims 1998)
 - Feature Selection (Mladenic and Grobelnik 1998, Yang Pederson 1998)
 - Latent Semantic Indexing
 - Support Vector Machines
 - Language Models (MacKey and Peto 1994)

45

True Class: Sports

ESPN.com - MLD - Recap - White Sox at Twins - 06/24/2002 - Microsoft Internet Explorer

Address: <http://sports.espn.go.com/mlb/recap?gameId=220624109>

ESPN Baseball **GAME TRACK** LIVE GAME UPDATES/BOX SCORES

GAME DAY RELAP Monday, June 24

Jones, Hunter cap late-inning comeback

RECAP | [BOX SCORE](#) | [GAME LOG](#)

MINNEAPOLIS (AP) -- After dominating the [Chicago White Sox](#) last year, Minnesota had to wait nearly three months this season before playing them for the first time Monday night.

The wait was worth it for the Twins and their fans.

[Jacque Jones](#) drove in the go-ahead run in the eighth after [Toni Hunter](#) tied it with a two-run homer an inning earlier, giving Minnesota a 5-4 victory.

"It's fun to play on a team where nobody quits," said starter [Matt Kinney](#), who left after six innings with a 4-2 deficit but watched Jones give the Twins the lead with a line drive to deep center against Bobby Howry (0-1).

[Luis Rojas](#) walked with two outs, took second on a wild pitch and scored the hit by Jones, who was thrown out to end the inning trying to stretch his hit to a triple.

The White Sox arrived two games under .500 and six games behind the Twins, so this four-game series is critical for Chicago, which has now lost four in a row and 19 of 27.

"At this point in the season, we're running out of time," Howry said. "We know

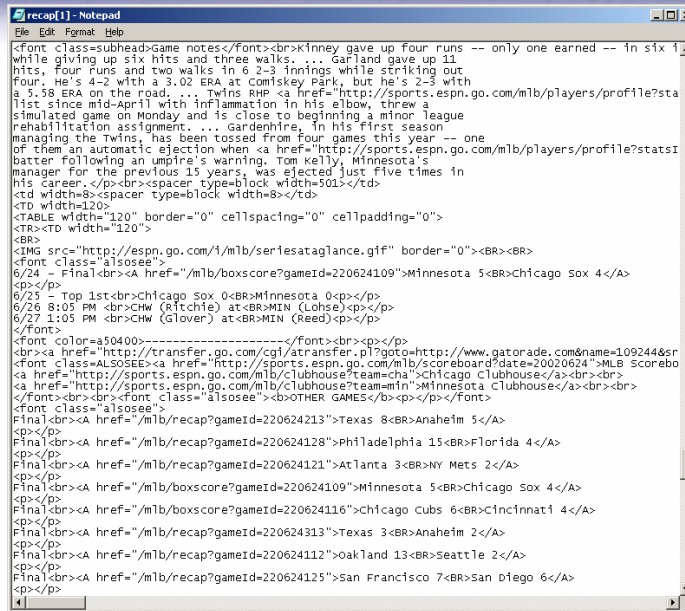
SERIES AT A GLANCE

6/24 - Final	Minnesota 5	Chicago Sox 4
6/25 - Top 1st	Chicago Sox 0	Minnesota 0
6/26 6:05 PM	CHW (Rothie) at	MIN (Lohse)
6/27 1:05 PM	CHW (Glover) at	MIN (Reed)

Gatorade first responders

46

Page Contents = HTML Code + Regular Text



```

recap(1) - Notepad
File Edit Format Help
<font class=subhead>Game notes</font><br>Kinney gave up four runs -- only one earned -- in six i
while giving up six hits and three walks. ... Garland gave up 11
hits, four runs and two walks in 6 2-3 innings while striking out
four. He's 4-2 with a 3.02 ERA at Comiskey Park, but he's 2-3 with
a 3.58 ERA on the road. ... Twins RHP <a href="http://sports.espn.go.com/mlb/players/profile?sta
list since mid-April with inflammation in his elbow, threw a
simulated game on Monday and is close to beginning a minor league
rehabilitation assignment. ... Gardenhire, in his first season
managing the Twins, has been tossed from four games this year -- one
of them an automatic ejection when <a href="http://sports.espn.go.com/mlb/players/profile?stats1
batter following an umpire's warning. Tom Kelly, Minnesota's
manager for the previous 15 years, was ejected just five times in
his career.</p><br><spacer type=block width=501></td>
<td width=8><spacer type=block width=8></td>
<td width=120>
<table width="120" border="0" cellspacing="0" cellpadding="0">
<tr><td width="120">
<br>
<br><br>
<font class=alsolee">
6/24 - Final<br><a href="/mlb/boxscore?gameId=220624109">Minnesota 5<br>Chicago Sox 4</a>
<p></p>
6/25 - Top 1st<br>Chicago Sox 0<br>Minnesota 0<p></p>
6/26 8:05 PM <br>CHW (Ritchie) at<br>MIN (Lohse)<p></p>
6/27 1:05 PM <br>CHW (Glover) at<br>MIN (Reed)<p></p>
</font>
<font color=#50400>-----</font><br><p></p>
<br><a href="http://transfer.go.com/cgi/atransfer.pl?gotow=http://www.gatorade.com&name=109244&er
<font class=alsolee"><a href="http://sports.espn.go.com/mlb/scoreboard?date=20020624">MLB Scorebo
<a href="http://sports.espn.go.com/mlb/clubhouse?team=cha">Chicago Clubhouse</a><br><br>
<a href="http://sports.espn.go.com/mlb/clubhouse?team=min">Minnesota Clubhouse</a><br><br>
</font><br><br><font class=alsolee"><b>OTHER GAMES</b><p></p></font>
<font class=alsolee">
Final<br><a href="/mlb/recap?gameId=220624213">Texas 8<br>Anaheim 5</a>
<p></p>
Final<br><a href="/mlb/recap?gameId=220624128">Philadelphia 15<br>Florida 4</a>
<p></p>
Final<br><a href="/mlb/recap?gameId=220624121">Atlanta 3<br>NY Mets 2</a>
<p></p>
Final<br><a href="/mlb/boxscore?gameId=220624109">Minnesota 5<br>Chicago Sox 4</a>
<p></p>
Final<br><a href="/mlb/boxscore?gameId=220624116">Chicago Cubs 6<br>Cincinnati 4</a>
<p></p>
Final<br><a href="/mlb/recap?gameId=220624313">Texas 3<br>Anaheim 2</a>
<p></p>
Final<br><a href="/mlb/recap?gameId=220624112">Oakland 13<br>Seattle 2</a>
<p></p>
Final<br><a href="/mlb/recap?gameId=220624125">San Francisco 7<br>San Diego 6</a>
<p></p>
</font>

```

47

Tokenization & Lexical Parsing

- HTML code is removed
- Punctuation is removed
- All words are converted to lowercase
- Stopwords are removed
 - Common, non-informative words such as 'the', 'and', 'with', 'an', etc...

Determine the term frequency (TF) of each remaining unique word

48

Result: Document Vector



home	2
game	8
hit	4
runs	6
threw	2
ejected	1
baseball	5
major	2
league	2
bat	2

49

Classifying Document Vectors

Test Document

home	2
game	8
hit	4
runs	6
threw	2
ejected	1
baseball	5
major	2
league	2
bat	2

?

?

?

bush	58
congress	92
tax	48
cynic	16
politician	23
forest	9
major	3
world	29
summit	31
federal	64

{News Class}

game	97
football	32
hit	45
goal	84
umpire	23
won	12
league	58
baseball	39
soccer	21
runs	26

{Sports Class}

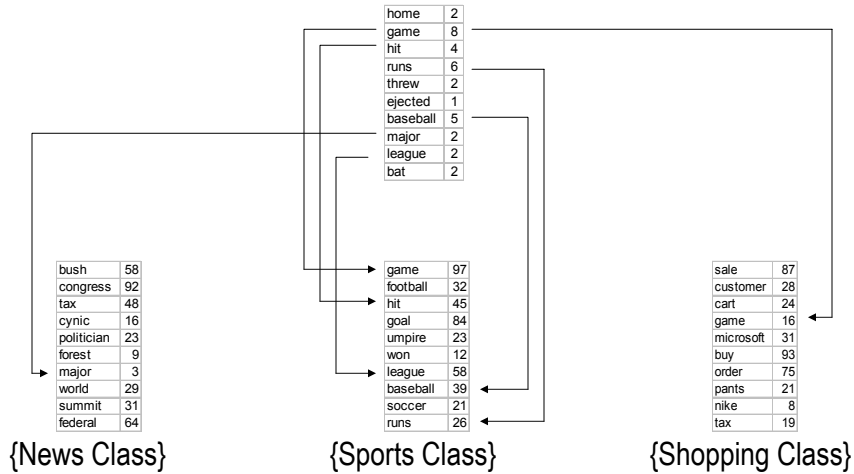
sale	87
customer	28
cart	24
game	16
microsoft	31
buy	93
order	75
pants	21
nike	8
tax	19

{Shopping Class}

50

Classifying Document Vectors

Test Document



51

Classifying Document Vectors

Test Document

home	2
game	8
hit	4
runs	6
threw	2
ejected	1
baseball	5
major	2
league	2
bat	2

$P(\{\text{News}\} \mid \text{Test Doc}) = 0.02$

bush	58
congress	92
tax	48
cynic	16
politician	23
forest	9
major	3
world	29
summit	31
federal	64

{News Class}

$P(\{\text{Sports}\} \mid \text{Test Doc}) = 0.91$

game	97
football	32
hit	45
goal	84
umpire	23
won	12
league	58
baseball	39
soccer	21
runs	26

{Sports Class}

$P(\{\text{Shopping}\} \mid \text{Test Doc}) = 0.07$

sale	87
customer	28
cart	24
game	16
microsoft	31
buy	93
order	75
pants	21
nike	8
tax	19

{Shopping Class}

52

Classifying Document Vectors

Test Document

home	2
game	8
hit	4
runs	6
threw	2
ejected	1
baseball	5
major	2
league	2
bat	2

$P(\{\text{Sports}\} \mid \text{Test Doc}) = 0.91$

game	97
football	32
hit	45
goal	84
umpire	23
won	12
league	58
baseball	39
soccer	21
runs	26

{Sports Class}

53

Classification Model

- A document is a vector of term frequency (TF) values, each category has its own term distribution
- Words in a document are generated by a multinomial model of the term distribution in a given class:

$$d_c \sim M\{n, \vec{p} = (p_1^c, p_2^c, \dots, p_{|V|}^c)\}$$

- Classification: $\arg \max_{c \in C} \{P(c \mid d)\}$

$$\arg \max_{c \in C} \{P(c) \prod_{i=1}^{|V|} P(w_i \mid c)^{n_i^c}\}$$

$|V|$: vocabulary size

n_i^c : # of times word i appears in class c

54



Results

- 25% correct classification
- Compare with random guessing of 7%
- More advanced techniques perform slightly better:
 - Shrinkage of word term frequencies (McCallum et al 1998)
 - n-gram models
 - Support Vector Machines

55



User Browsing Model

User Browsing Model

- Web browsing is "sticky" or persistent: users tend to view a series of pages within the same category and then switch to another topic
- Example:



57

Markov Switching Model

	arts	business	computers	games	health	home	news	recreation	reference	science	shopping	society	sports	adult
arts	83%	4%	5%	2%	1%	2%	6%	3%	2%	6%	2%	3%	4%	1%
business	3%	73%	5%	3%	2%	3%	6%	2%	3%	3%	3%	2%	3%	2%
computers	5%	11%	79%	3%	3%	7%	5%	3%	4%	4%	5%	5%	2%	2%
games	1%	3%	2%	90%	1%	1%	1%	1%	0%	1%	1%	1%	1%	0%
health	0%	0%	0%	0%	84%	1%	1%	0%	0%	1%	0%	1%	0%	0%
home	0%	1%	1%	0%	1%	80%	1%	1%	0%	1%	1%	1%	0%	0%
news	1%	1%	1%	0%	1%	0%	69%	0%	0%	1%	0%	1%	1%	0%
recreation	1%	1%	1%	0%	1%	1%	1%	86%	1%	1%	1%	1%	1%	0%
reference	0%	1%	1%	0%	1%	0%	1%	0%	85%	2%	0%	1%	1%	0%
science	1%	0%	0%	0%	1%	1%	0%	1%	75%	0%	0%	1%	0%	0%
shopping	1%	3%	2%	1%	1%	2%	1%	1%	0%	1%	86%	1%	1%	0%
society	1%	1%	2%	0%	2%	1%	3%	1%	2%	2%	0%	82%	1%	1%
sports	2%	1%	1%	0%	0%	0%	3%	1%	1%	0%	0%	1%	85%	0%
adult	1%	1%	1%	0%	0%	0%	1%	0%	0%	0%	0%	1%	0%	93%
	16%	10%	19%	11%	2%	3%	2%	6%	3%	2%	7%	6%	5%	7%

Pooled transition matrix, heterogeneity across users

58

Implications

- Suppose we have the following sequence:



- Using Bayes Rule can determine that there is a 97% probability of news, unconditional=2%, conditional on last observation=69%

59

Results



Methodology

Bayesian setup to combine information from:

- Known categories based on exact matches
- Text classification
- Markov Model of User Browsing
 - Introduce heterogeneity by assuming that conditional transition probability vectors drawn from Dirichlet distribution
- Similarity of other pages in the same domain
 - Assume that category of each page within a domain follows a Dirichlet distribution, so if we are at a “news” site then pages more likely to be classified as “news”

61



Findings

Random guessing	7%
Text Classification	25%
+ Domain Model	41%
+ Browsing Model	78%

62



Findings about Text Classification



Key Points of Text Processing

Can turn text and qualitative data into quantitative data

- Each technique (text classification, browsing model, or domain model) performs only fairly well (~25% classification)
- Combining these techniques together results in very good (~80%) classification rates



Applications

- Newsgroups
 - Gather information from newsgroups and determine whether consumers are responding positively or negatively
- E-mail
 - Scan e-mail text for similarities to known problems/topics
- Better Search engines
 - Instead of experts classifying pages we can mine the information collected by ISPs and classify it automatically
- Adult filters
 - US Appeals Court struck down Children's Internet Protection Act on the grounds that technology was inadequate

65



Conclusions



Lessons about Behavior

- We reveal a wealth of information about ourselves through clicking, which can then be used to accurately predict about who we are and our interests.
- This works because we tend for information that is compatible with our interests and goals.