



# **CS834 - Introduction to Information Retrieval**

## **Presentation #4**

Hussam Aldeen Hallak

# The Two Papers:

- A semantic approach to contextual advertising

SIGIR '07 Proceedings of the 30th annual international ACM SIGIR conference on Research and development in information retrieval  
Pages 559-566

- How much can behavioral targeting help online advertising?

WWW '09 Proceedings of the 18th international conference on World wide web  
Pages 261-270

# The First Paper:

- A semantic approach to contextual advertising

SIGIR '07 Proceedings of the 30th annual international ACM SIGIR conference on Research and development in information retrieval  
Pages 559-566

<http://clair.si.umich.edu/~radev/767w10/papers/Week12/ca/semantic.pdf>

# 1<sup>st</sup> Type of Advertising: Sponsored Search (SS)

- Placing ads on the result pages from a search engine

The image shows a Google search interface for '1997 ford f150 fuel pump'. The search bar is at the top with the Google logo on the left and a microphone icon on the right. Below the search bar are tabs for 'All', 'Shopping', 'Videos', 'Images', 'News', and 'More'. To the right of these tabs are 'Settings' and 'Tools'. A red arrow points from the 'Tools' link down to a red-bordered box containing a grid of five sponsored product listings. Each listing includes an image of a fuel pump, a title, a price, and the retailer's name. Below this grid are two more red-bordered boxes, each containing a sponsored text ad. A red arrow points from the 'Settings' link down to the first text ad, and another red arrow points from the 'Tools' link down to the second text ad.

Google






1997 ford f150 fuel pump

All Shopping Videos Images News More Settings Tools

About 1,180,000 results (0.85 seconds)

Shop for 1997 ford f150 fuel pump on Google

Sponsored ⓘ

|   |  |   |  |   |
|---|--|---|--|---|
| <br>1997 Ford F-150 Replacement Fu...<br>\$89.21<br>CarParts.com | <br>1997 Ford F150 Truck Fuel Pum...<br>\$69.95<br>1A Auto.com<br>Free shipping | <br>1997 Ford F-250 Replacement Fu...<br>\$89.21<br>CarParts.com | <br>1997 Ford Thunderbird...<br>\$54.68<br>CarParts.com | <br>1997 Ford E-150 Econoline Club...<br>\$66.54<br>CarParts.com |
|---|--|---|--|---|

Shop 1997 Ford F150 Fuel Pump - Free 2-day Shipping w/ Prime

Ad [www.amazon.com/automotive/parts](http://www.amazon.com/automotive/parts)

★★★★★ Rating for amazon.com: 4.7

Find Deals on 1997 Ford F150 Fuel Pump in Car Parts on Amazon.

Fuel Pumps - Advance Auto Parts® Official - [advanceautoparts.com](http://advanceautoparts.com)

Ad [shop.advanceautoparts.com/Advance-Auto/Spectra](http://shop.advanceautoparts.com/Advance-Auto/Spectra)

Shop For Fuel Pumps From Advance Auto. Buy Online, Pick Up In-Store Today!

Ratings: Prices 10/10 - Selection 10/10 - Quality 9.5/10 - Service 9/10 - Shipping 9/10 - Returns 9/10

## 2<sup>nd</sup> Type of Advertising: Context Match (CM)

- Commercial ads within the content of a web page.

https://www.kbb.com

Kelley Blue Book®  
The Trusted Resource

Home Car Values Cars for Sale Car Reviews Awards & Top 10s Research Tools Sign In ZIP code 0

New Cars by

Category Make Best Sellers  
Sponsored

SUV Crossover Sedan Truck Hatchback Convertible

Luxury Coupe Electric Hybrid Van/Minivan Wagon

2017 CHEVROLET SILVERADO  
2017 CHEVY CLOSEOUT  
GET SPECIAL CLOSEOUT PRICING ON OUR MOST POPULAR CHEVY TRUCKS<sup>1</sup>  
1. Not available with special financing, lease and some other offers. Take delivery by 11/30/17. See participating dealer for details.  
Learn More  
Advertisement

2017 Kia Sorento vs. 4 other midsize SUV's  
Compare now  
Presented by KIA  
Advertisement

5 / 46

# Why ads and page content should be related?


Improve user experience → increase clicks → increase revenue

<https://www.kbb.com/ford/excursion/2004/eddie-bauer-sport-utility-4d/?vehicleid=196210&intent=buy-used&mileage=141626&condition=excellent&pricet...>

## 2004 Ford Excursion

[photos](#)[videos](#)[360° view](#)[colors](#)

[all](#)[exterior](#)[interior](#)



1 of 3  
<previous | next>

Select another vehicle


Year ▼

Make ▼


Model ▼

Update

Offer disclosure

Go Further 

2017 EDGE  
SEL AWD





\$329 <sup>A</sup> MONTH / 39 <sup>MONTH</sup> RED CARPET LEASE  
\$2,369 CASH DUE AT SIGNING\*


[BUILD & PRICE](#)[VIEW OFFERS](#)

Local Ford Dealer


Advertisement



 To Find Your Next Vehicle!



Cavalier Ford Lincoln - Greenbrier

 866-362-1535

[LEARN MORE](#)

Advertisement

# Matching ads with pages

## **Syntactic Approach:**

Match words found in the page with words in ads.

## **Problems: Leads to irrelevant ads**

A page about the golfer “John **Maytag**” might trigger an ad for “**Maytag** dishwashers”

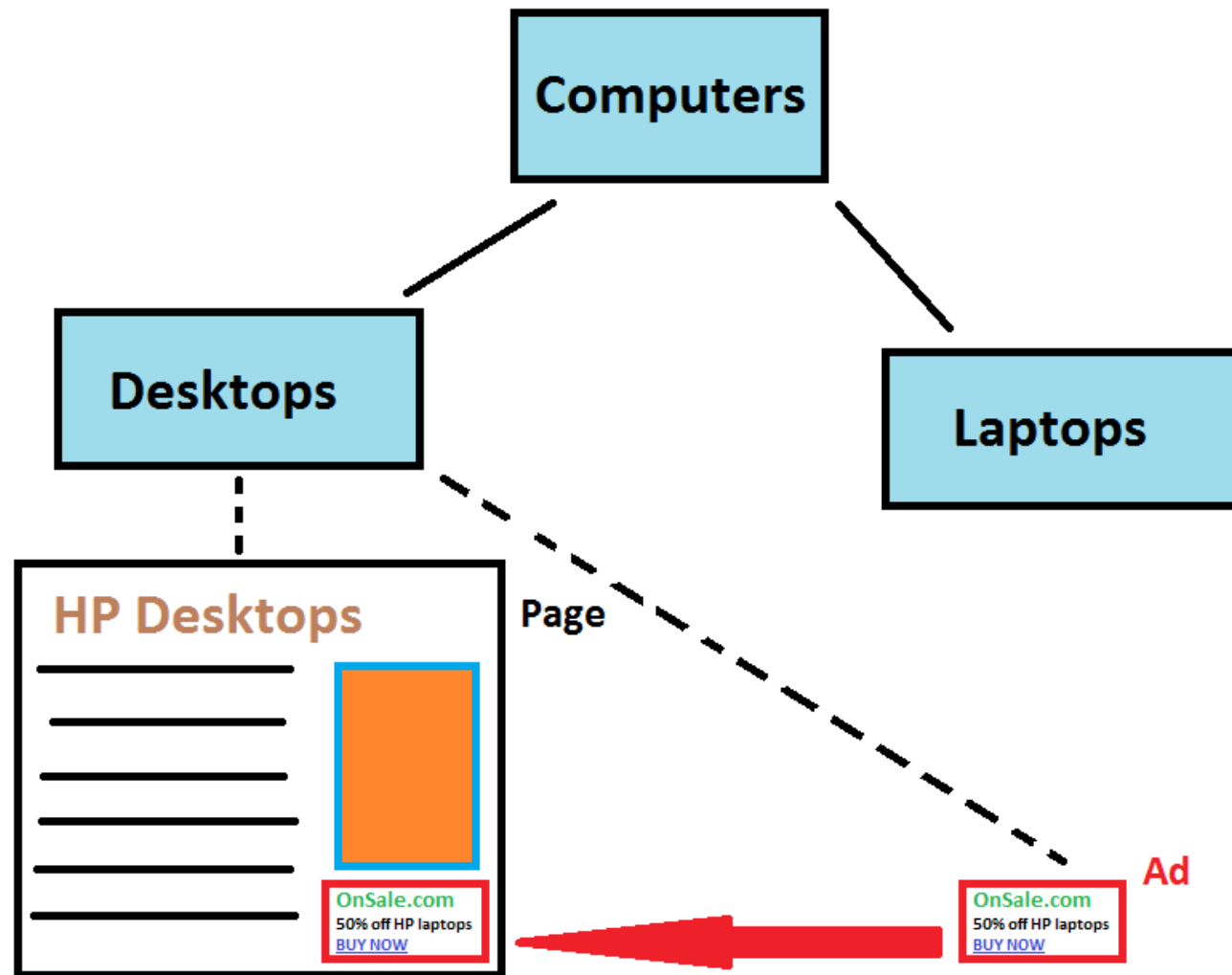
## **Solution:**

Combining Semantic (topical) and syntactic matching

## **The semantic phase:**

1. Classify the page and the ads into a taxonomy of topics
2. Use the proximity of the ad and page classes as a factor in the ad ranking formula

# Advantages of Using a hierarchical taxonomy



Classes —————> Set of applicable ads  
keywords —————> Narrow down the search.





# Taxonomy Choice

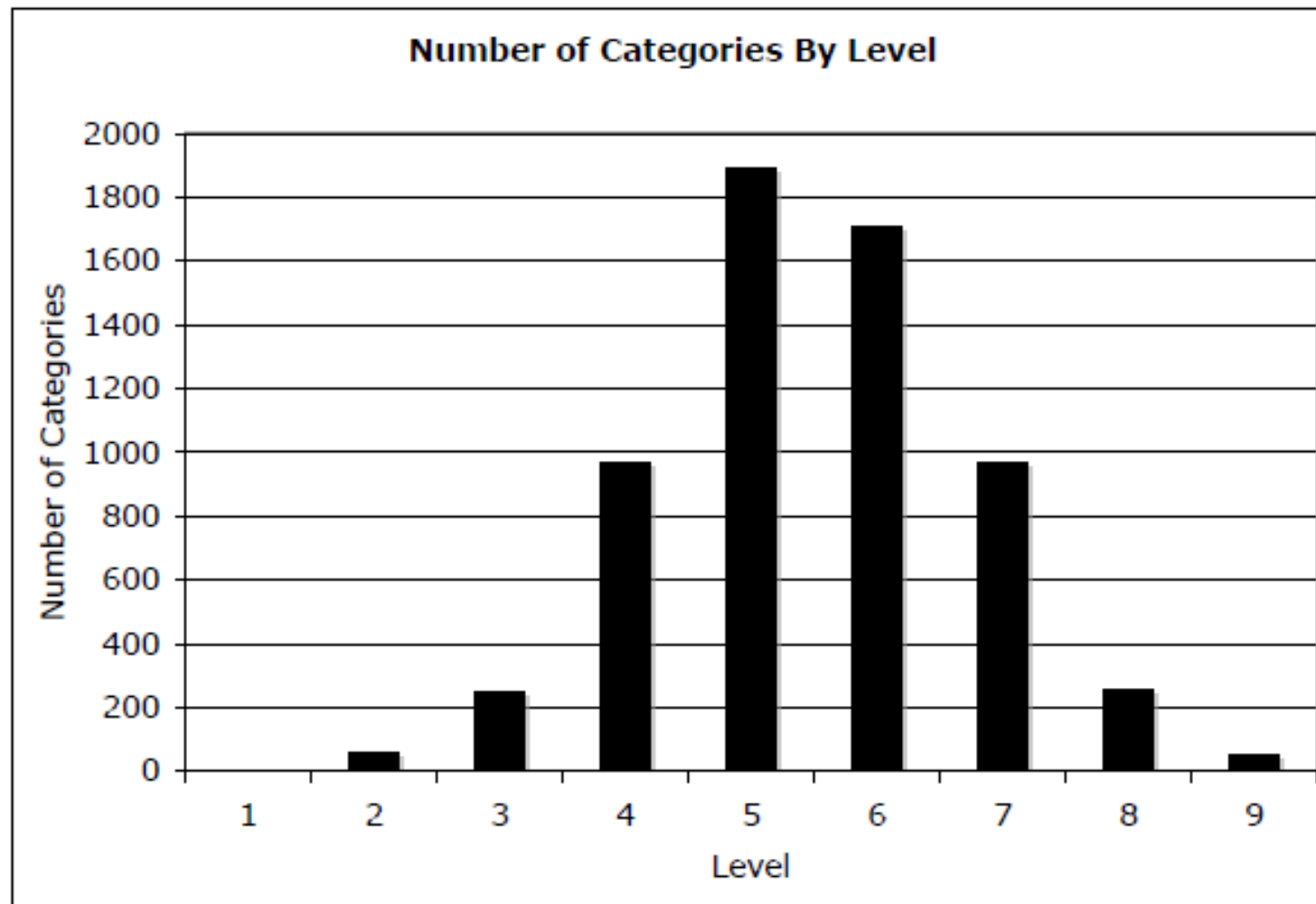
Built by a large web search engine in the US

100 queries for each node

Contains 6000 nodes

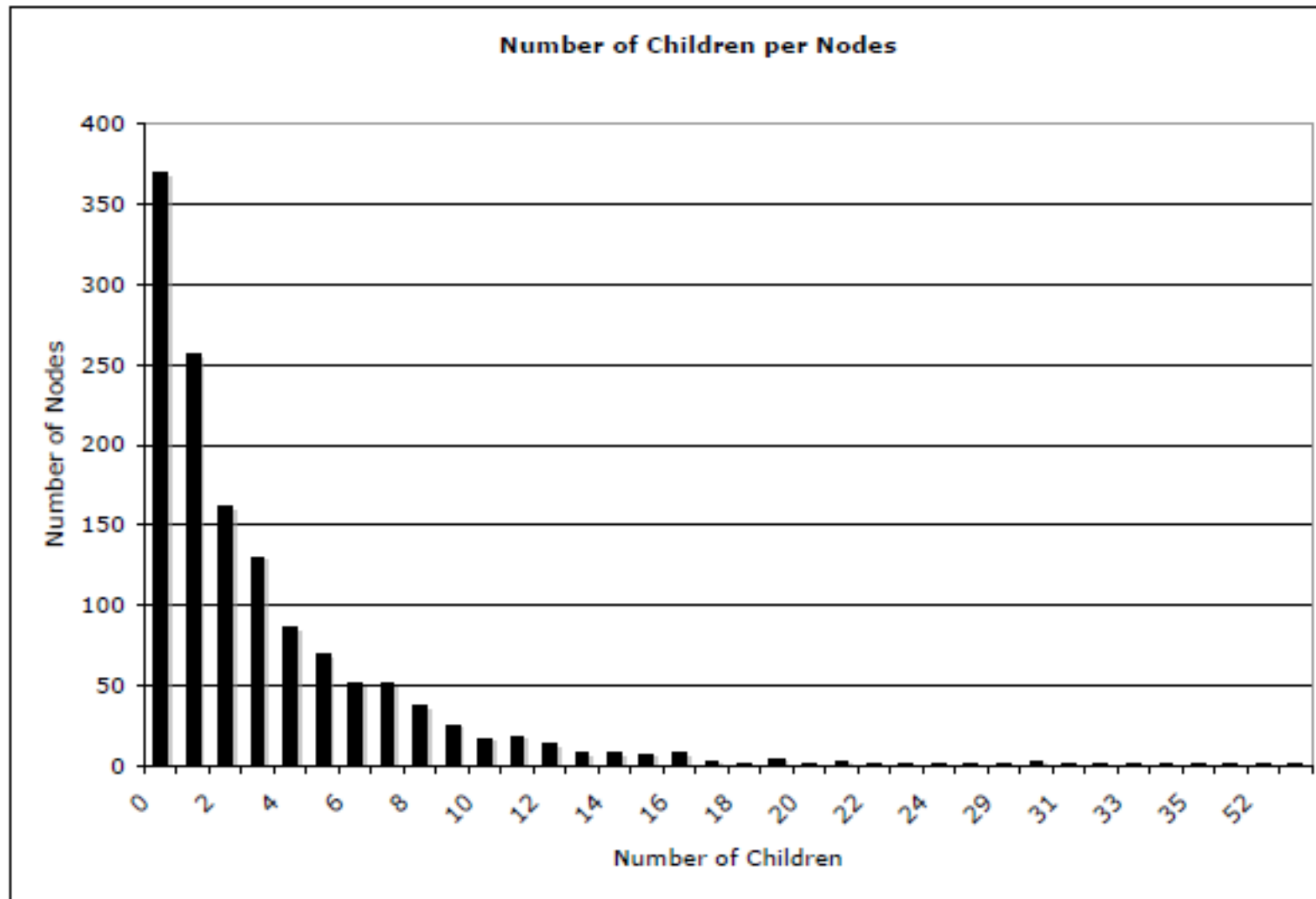
Used for classifying both pages and ads

# Categories per level



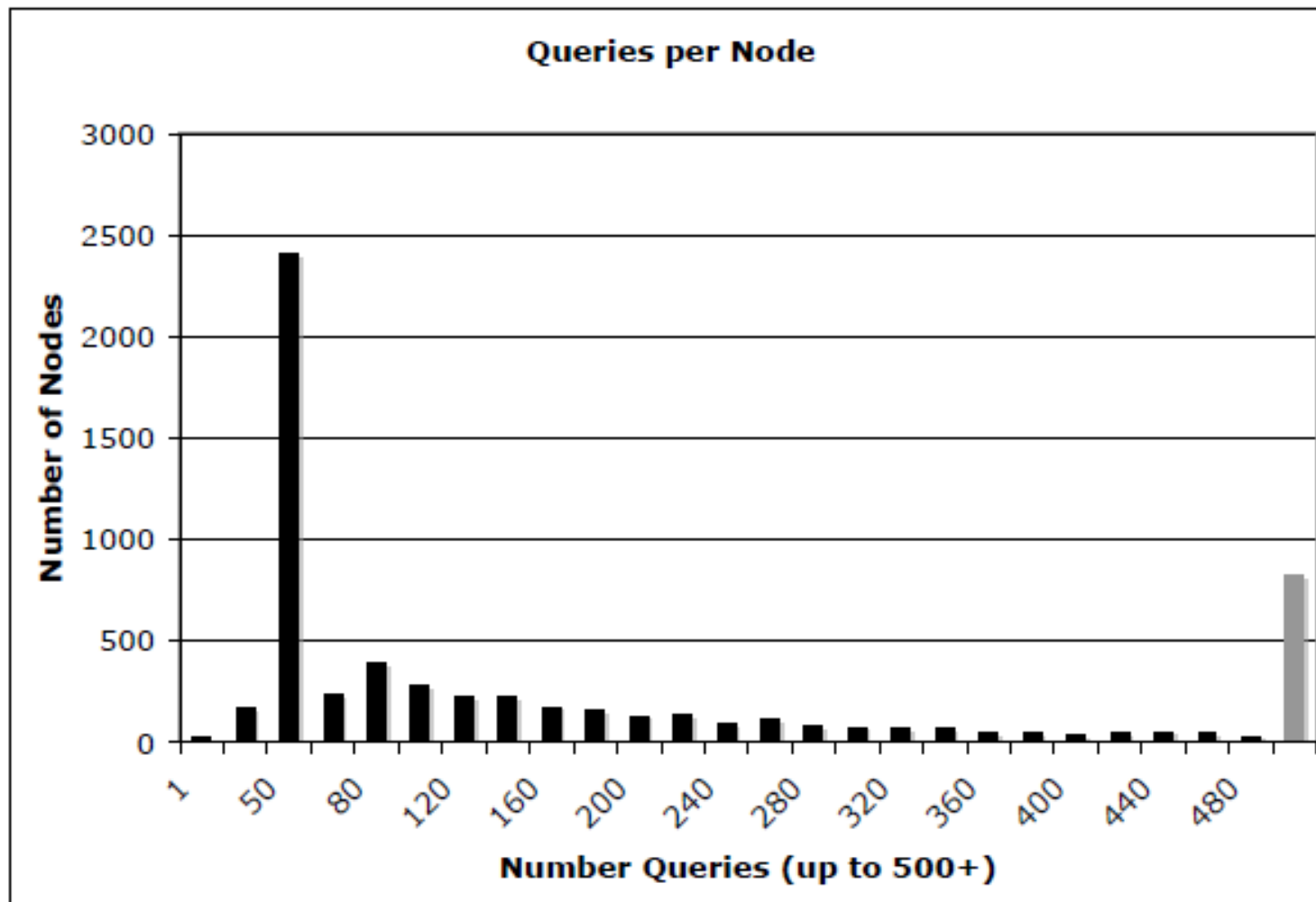
**Figure 1: Taxonomy statistics: categories per level**

# Number of children per node



**Figure 1: Taxonomy statistics: fanout for non-leaf nodes**

# Queries per node



**Figure 1: Taxonomy statistics: queries per node**



# Training data

Pages: Select top 10 results of Web search index for each class in the taxonomy

Ads: Select ads with a bid-phrase assigned to the class

# Classifiers:

- SVM and log-regression classifiers were slow!
- Rocchio's (nearest-neighbor) gave best performance!
- Each taxonomy node: a single meta-document (concatenation of all the example queries), represented as a centroid for the class.
- A centroid is the sum of the TF.IDF values of each term.

$$\vec{c}_j = \frac{1}{|C_j|} \sum_{\vec{q} \in C_j} \frac{\vec{q}}{\|\vec{q}\|}$$

where  $\vec{c}_j$  is the centroid for class  $C_j$   
 $q$  iterates over the queries in a particular class.

- Classification is based on the cosine of the angle between the document and the centroid.



# Semantic-Syntactic Matching

- Process the content of the page
- Extract features
- Search the ad space to find the best matching ads.

# Relevance score

- **Convex combination of the keyword (syntactic) and classification (semantic) score:**

$$\begin{aligned} \text{Score}(p_i, a_i) = & \alpha \cdot \text{TaxScore}(\text{Tax}(p_i), \text{Tax}(a_i)) \\ & + (1 - \alpha) \cdot \text{KeywordScore}(p_i, a_i) \end{aligned}$$

- **$\alpha$  determines the relative weight of the taxonomy score and the keyword score.**
- **$\alpha = \text{TaxScore} / \text{KeywordScore}$**



# KeywordScore (syntactic relevance score)

- Uses Vector Space Model
- Pages and ads are vectors in n-dimensional space (one dimension for each distinct term)
- KeywordScore is the cosine of the angle between the page and the ad vectors

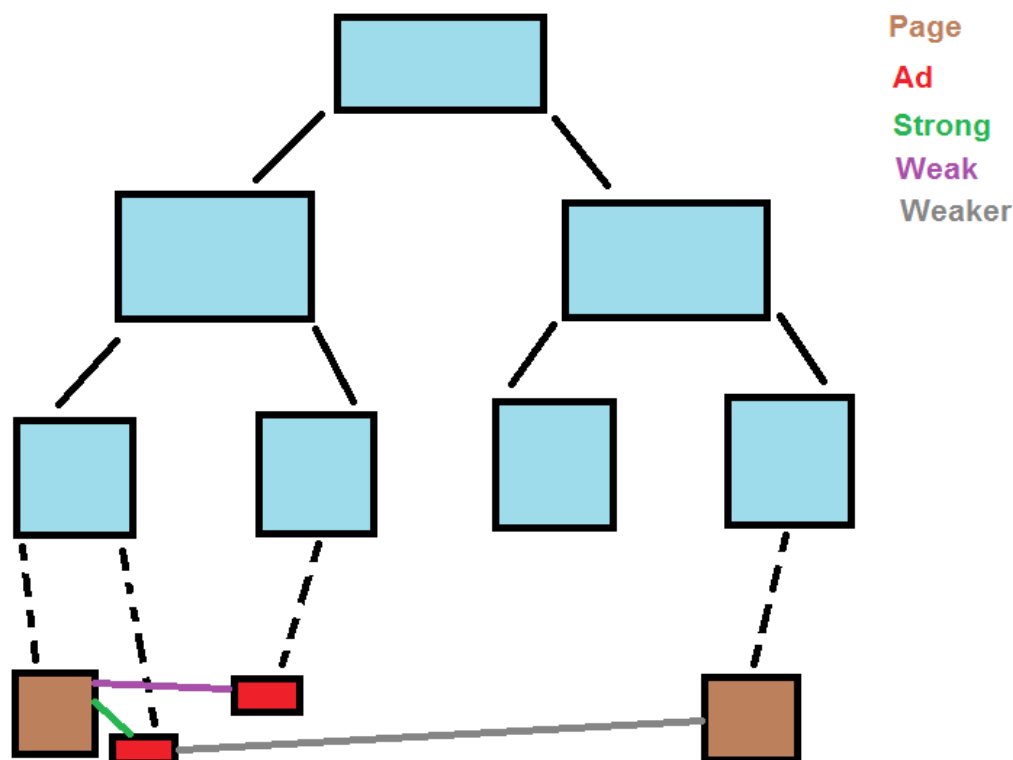
# TaxonomyScore (semantic relevance score) (page-ad topical match score)

- **Purpose:**

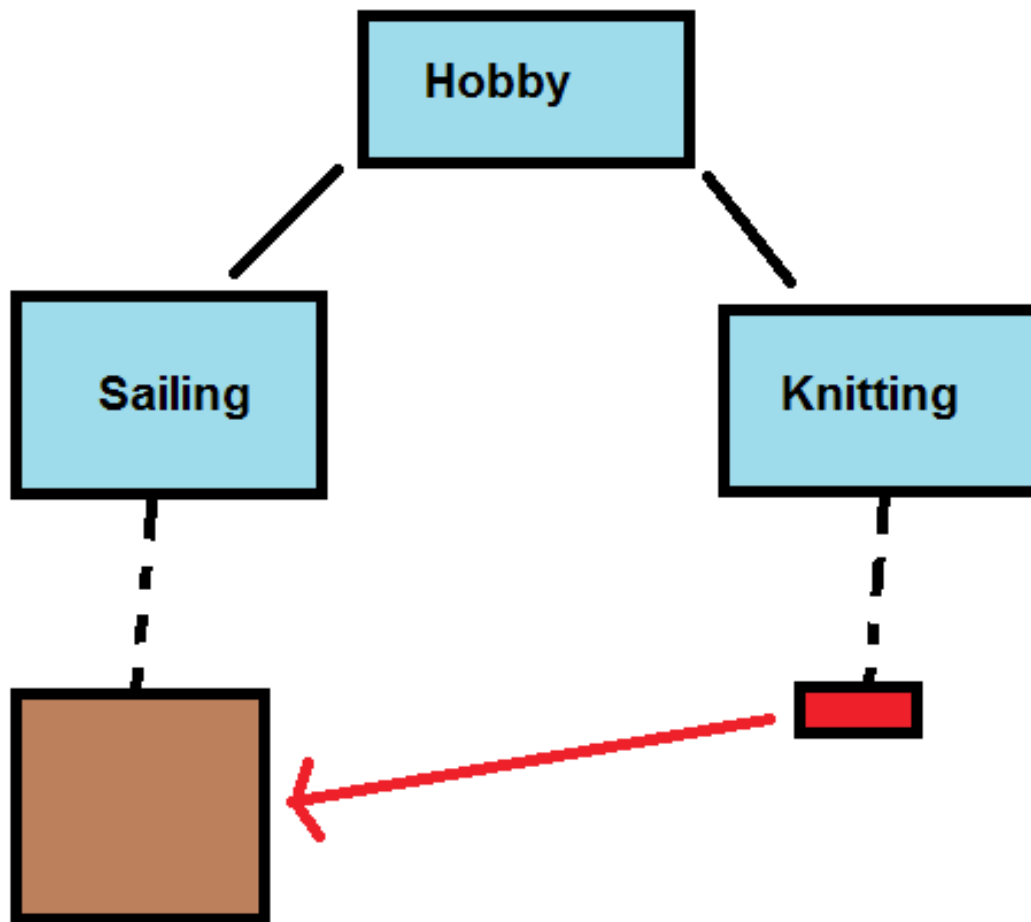
1. Match ads and pages based on the topic
2. Generalization within a taxonomy
3. Efficient search of the ad space (user is waiting)

- **Ideally:**

The match is stronger when both the ad and the page are classified into the same node and weaker when the distance between the nodes in the taxonomy gets larger.



# Generalization Challenge



Page  
Ad

# Generalization Advantage

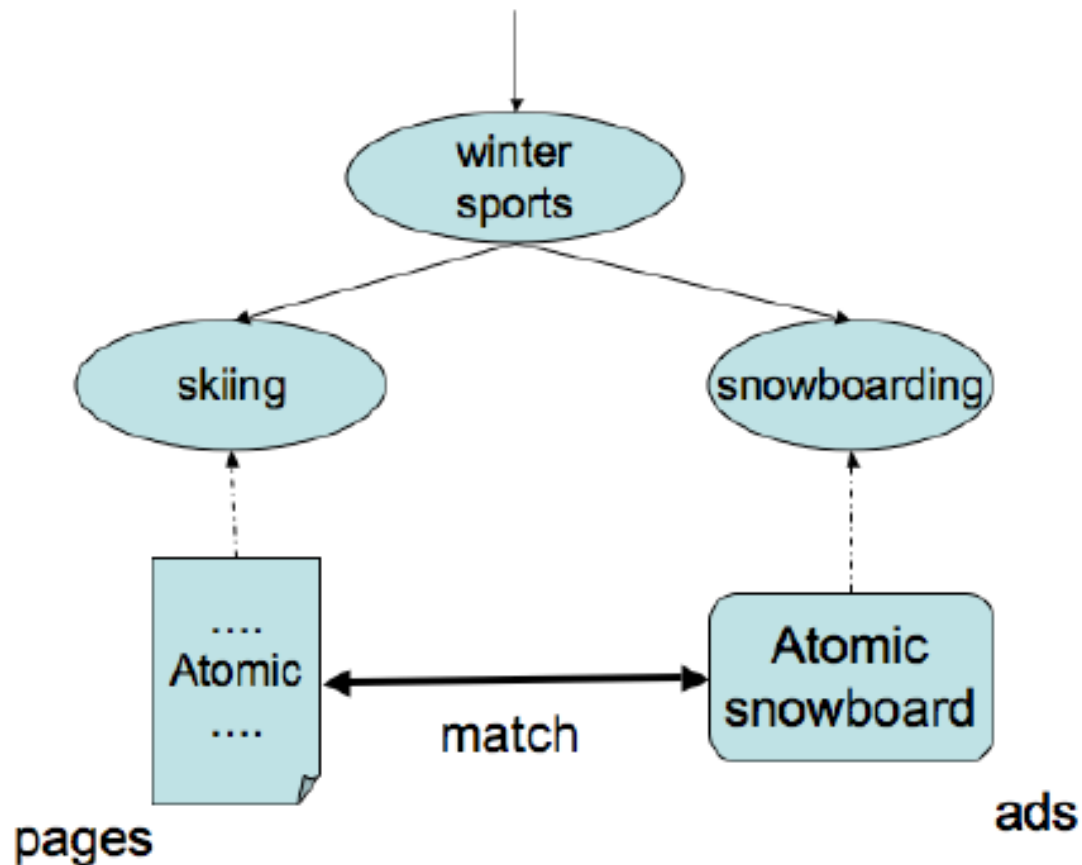


Figure 2: Two generalization paths

# Generalization cost

- Density: The probability of an ad belonging to the parent topic being suitable for the child topic.

$$idist(c, p) = \frac{n_c}{n_p}$$

Where:

$c$  is the child node.

$p$  is the parent node.

$n_c$  is the number of document classified into the subtree rooted at  $c$ .

$n_p$  is the number of document classified into the subtree rooted at  $p$ .

- $0 \leq idist \leq 1$   
 $idist = 1$  (the page and ad belong to the same class/node)  
 $idist = 0$  (the least common ancestor of page and ad is the root)

# Searching the ad space

- Using inverted index
- The ads are parsed into terms
- Each term has a weight based on a section where it appears
- Challenge: Preserving class information in the index
- Solution: Annotate ads with a unique meta-term for each class
- Cons: Generalization is lost!
- Instead: Also annotate each ad with one meta-term for each ancestor of the assigned class; utilize weights!
- Weights of the meta-terms: the value of `idist()` function



Example:

# Data and Methodology

- Data: 105 pages randomly selected from 20 million pages with contextual advertising
- Tens of millions of ads from advertising network in the US
- Human judges for each page-ad pair on a 1 to 3 scale:

1. Relevant

Page: The National Football League

Ad: Tickets for NFL games

2. Somewhat Relevant

Page: The National Football League

Ad: NFL branded products

3. Irrelevant:

Page: The National Football League

Ad: NFL player “John Maytag” triggers “Maytag” dishwasher ads on NFL page.



# Results

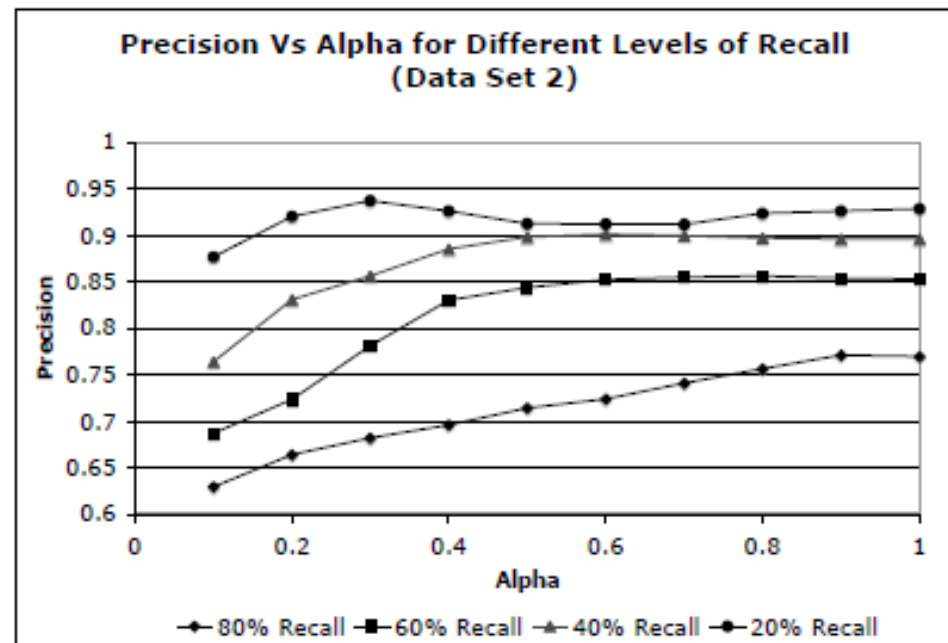


Figure 4: Impact of  $\alpha$  on precision for different levels of recall

|                               |      |
|-------------------------------|------|
| pages                         | 105  |
| words per page                | 868  |
| judgments                     | 2946 |
| judg. inter-editor agreement  | 84%  |
| unique ads                    | 2680 |
| unique ads per page           | 25.5 |
| page classification precision | 70%  |
| ad classification precision   | 86%  |

Table 1: Dataset statistics

In most cases, precision grows or is flat when Alpha is increased.

# What to take away?

- Syntactic (keyword) matching between pages and ads leads to irrelevant ads.
- Semantic (topical) matching relies on the matching between pages and ads in topic
- Semantic matching complemented with syntactic matching leads to better results.

# The Second Paper:

- How much can behavioral targeting help online advertising?

WWW '09 Proceedings of the 18th international conference on World wide web

Pages 261-270

<http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.215.1473&rep=rep1&type=pdf>

# What is Behavioral Targeting BT Ads?

- The delivery of ads to targeted users based on their web search and browsing history.
- How much can BT help online advertising?
- Strategies to represent the users' behavior:
  1. Web browsing history
  2. Search queries
- Which BT strategy is better for user segmentation?
- The performance of online advertising is measured by ads Click-Through Rate (CTR)
- Calculate & compare ads click entropy, precision, recall, and F-measure for different BT strategies.
- Short window (1 day) vs long window (7 days)

# Modeling User Browsing History

- Users are represented by a matrix  $U_{g \times l}$   
g: number of users  
l: number of urls
- TF.IDF:  
Users are the documents  
URLs are the terms
- Each entry in the matrix is given the value:

|                  | Lnk <sub>1</sub> | Lnk <sub>2</sub> | ... | Lnk <sub>l</sub> |
|------------------|------------------|------------------|-----|------------------|
| usr <sub>1</sub> |                  |                  |     |                  |
| usr <sub>2</sub> |                  |                  |     |                  |
| ...              |                  |                  |     |                  |
| usr <sub>g</sub> |                  |                  |     |                  |

$$u_{ij} = \log[\text{count}(\text{clicks on } URL_j \text{ by } User_i) + 1] \times \log\left[\frac{l}{\text{count}(\text{users who have clicked on } URL_j)}\right]$$

# Modeling User Search History

- Query history uses Bag-of-Words (BOW) model to populate TF.IDF matrix.
- Stop words are removed and Porter stemming is used
- Terms that only appear once are removed.
- Number of terms is reduced from 765k to 294k

# Examined BT Strategies:

- Two BT strategies (pages visited, queries searched)
- Two window sizes (long term: 7 days, short term: 1 day)
- Lead to four possible BT strategies to assess:
  1. Long term behavior based on page views (LP)
  2. Long term behavior based on query terms (LQ)
  3. Short term behavior based on page views (SP)
  4. Short term behavior based on query terms (SQ)

# Raw Data Clean-Up!

- Source: A commercial search engine (Bing?)
- log data tracked:
  1. Web page clicking
  2. Ad clicking
- Period: 7 days' click-through log data ranging from June 1<sup>st</sup> to 7<sup>th</sup> 2008
- Users removed if they clicked on more than 100 ads in one day (robots)
- Ads removed if they have fewer than 30 clicks over seven days (cannot be used to draw a reliable statistical conclusions)
- Only English queries are considered!



# Click-Through Log Format

**Table 1. Format of click-through log used in our study.**

|                  |  |   |
|------------------|--|---|
| <b>UserID</b>    | UID030608473X  | A user ID for each unique user.   |
| <b>QueryText</b> | xbox   | The detailed query text used by the user  |
| <b>QueryTime</b> | 08-06-03 21:15:47  | The time when the query was issued  |
| <b>ClickTime</b> | 08-06-03 21:16:02  | The time when the click occurred after the query was issued   |
| <b>ClickURL</b>  | http://www.xbox365.com   | The URL which has been clicked by the user  |
| <b>IsAd</b>      | 0  | A Boolean value to show the clicked URL is an ad or not   |
| <b>NumberAd</b>  | 3  | The number of ads displayed in the search results   |
| <b>DisplayAd</b> | http://video-games.half.ebay.com/<br>http://accessories.us.dell.com/<br>http://www.gamefly.com | The URL list of all the ads that displayed by the query. (To save space, we only reserve top domain of the ad URL in this example.) |

# Definitions:

- $A = \{a_1, a_2, \dots, a_n\}$  is the set of ads
- $Q_i = \{q_{i1}, q_{i2}, \dots, q_{imi}\}$  queries which have displayed or clicked  $a_i$
- $U_i = \{U_{i1}, U_{i2}, \dots, U_{imi}\}$  users who have displayed or clicked  $a_i$
- $\delta(u_{ij})$  is used to show whether the user  $u_{ij}$  has clicked the ad  $a_i$

$$\delta(u_{ij}) = \begin{cases} 1 & \text{if } u_{ij} \text{ clicked } a_i \\ 0 & \text{otherwise} \end{cases}$$

$l_i = \sum_j \delta(u_{ij})$ : the number of users clicked ad  $a_i$

- K-means and CLUTO (a clustering software package) cluster users into groups:  $g_k(U_i)$  is all users in  $U_i$

$$G(U_i) = \{g_1(U_i), g_2(U_i), \dots, g_K(U_i)\}, i=1,2,\dots,n$$

# Calculating Similarity Between Users

- Cosine similarity is used:

$$Sim(u_{ij}, u_{st}) = \frac{\langle u_{ij}, u_{st} \rangle}{||u_{ij}|| ||u_{st}||}$$

- Within-ad similarity: Users who clicked the same ad

$$S_w(a_i) = \frac{2}{l_i(l_i - 1)} \sum_{\delta(u_{ij})=1} \sum_{\substack{\delta(u_{it})=1 \\ t \neq j}} Sim(u_{ij}, u_{it})$$

- Between-ad similarity: Users who clicked different ads

$$S_b(a_i, a_s) = \frac{1}{l_i l_s} \sum_{\delta(u_{ij})=1} \sum_{\delta(u_{st})=1} Sim(u_{ij}, u_{st})$$

# Similarity Ratio

- Ratio between within-ad and between-ad similarity

$$R(a_i, a_s) = \frac{S_w(a_i) + S_w(a_s)}{2S_b(a_i, a_s)}$$

- The larger the ratio, the more confident we are on the basic assumption of BT for a pair of ads  $a_i$  and  $a_s$

# Ad Click-Through Rate (CTR)

- The CTR of ad  $a_i$  is defined as the number of users who clicked it over the number of users who either clicked it or only displayed it.

$$CTR(a_i) = \frac{1}{m_i} \sum_{j=1}^{m_i} \delta(u_{ij})$$

# F-measure

- CTR can be used to calculate precision and recall
- Positive instance: Users displayed and clicked  $a_i$
- Negative instance: Users displayed  $a_i$  but didn't click it

- Precision: CTR of segment

$$Pre(a_i|g_k) = CTR(a_i|g_k)$$

- Recall: clicks of segment/total clicks

$$Rec(a_i|g_k) = \frac{\sum_{u_{ij} \in g_k(u_i)} \delta(u_{ij})}{\sum_{j=1}^{m_i} \delta(u_{ij})}$$

- F-measure:

$$F(a_i|g_k) = \frac{2Pre(a_i|g_k)Rec(a_i|g_k)}{Pre(a_i|g_k) + Rec(a_i|g_k)}$$

- The larger the F measure is, the better the achieved performance is by user segmentation for BT

## Ads-Click Entropy

- For ad  $a_i$ , the probability of users in segment  $g_k$ , who will click this ad, is estimated by:

$$P(g_k|a_i) = \frac{1}{m_i} \sum_{u_{ij} \in g_k(U_i)} \delta(u_{ij})$$



- Ads-Click Entropy (mathematically):

$$Enp(a_i) = - \sum_{k=1}^K P(g_k|a_i) \log P(g_k|a_i)$$

- Smaller Entropy => Better user segmentation

# Within- and between- ads user similarity

**Table 2. Within- and between- ads user similarity.**

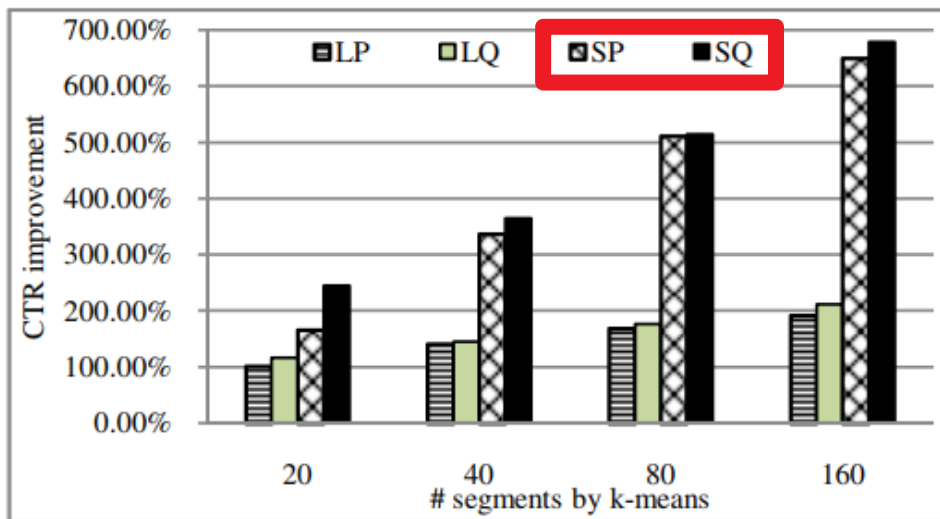
|           |  $S_w$ |  $S_b$ | $R$     |
|-----------|---|--|---------|
| <b>LP</b> | 0.1417  | 0.0252   | 28.9217 |
| <b>LQ</b> | 0.2239  | 0.0196   | 44.2908 |
| <b>SP</b> | 0.1532  | 0.0281   | 24.5086 |
| <b>SQ</b> | 0.2594  | 0.0161   | 91.1890 |

Scores are average across all ads or query terms

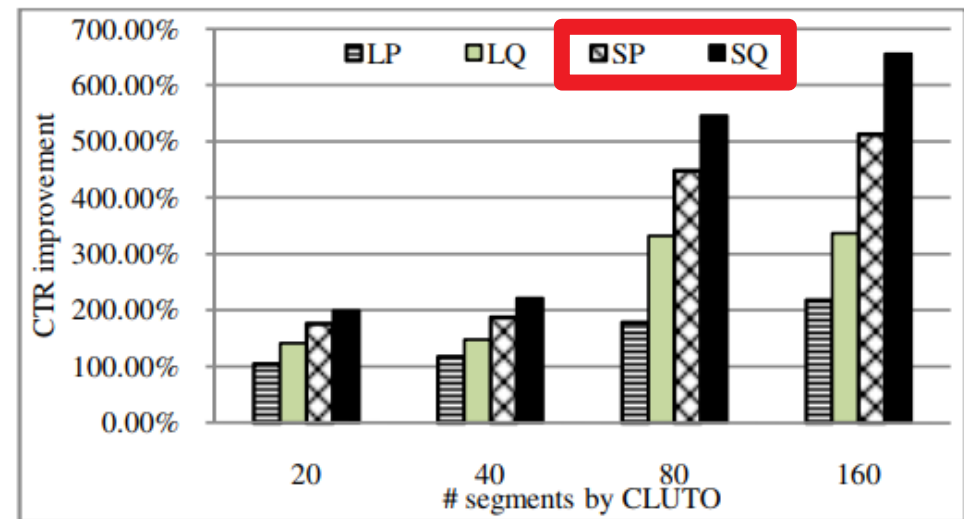
- Users who clicked the same ad are up to 91 times more similar
- For all ad pairs, 99.37% had a higher within-ad user similarity than between-ad similarity
- Search queries are more effective than pages clicked in BT



# User segmentation improves CTR by 670%:



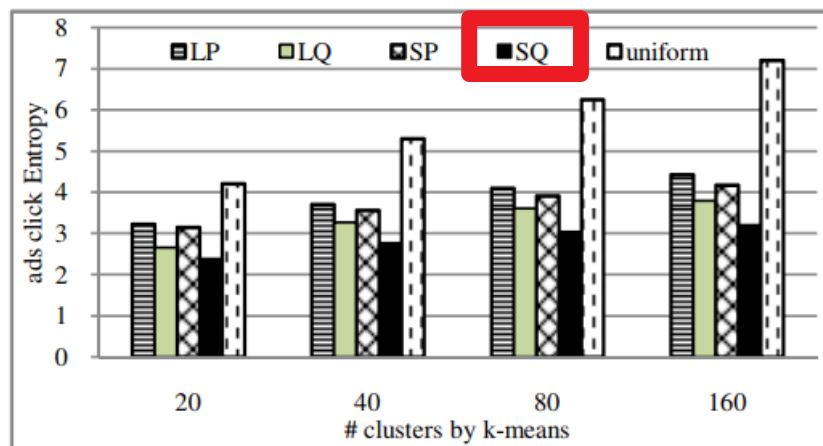
(a) User clustering by k-means



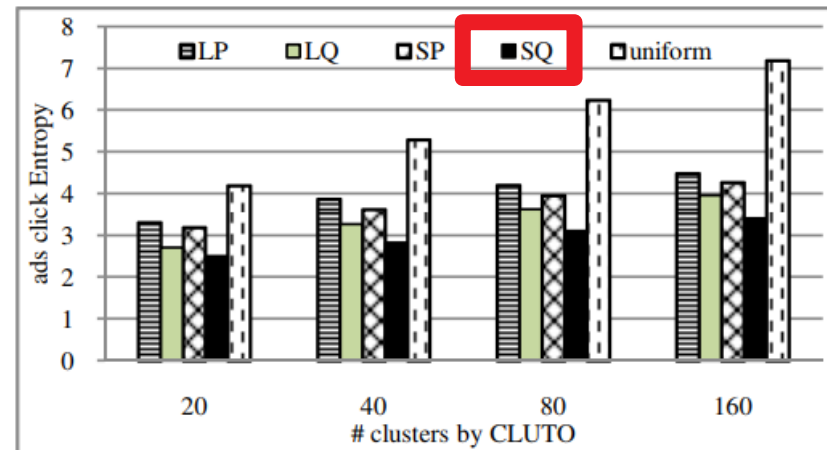
(b) User clustering by CLUTO

Figure 1. CTR improvements by user segmentation for BT.

- CTR was improved by up to 670% off of the non segmented CTR



(a) Cluster by k-means



(b) Cluster by CLUTO

Figure 2. Ads click Entropy of user segmentation for BT.

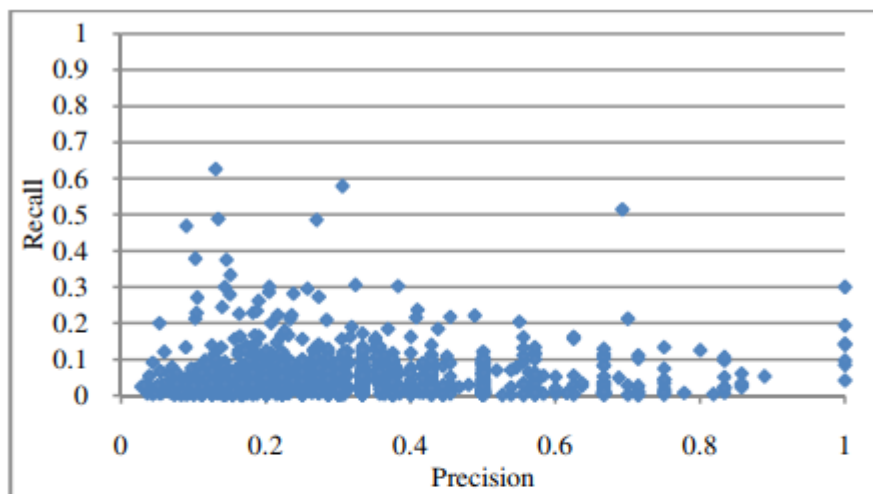
# Overall Performance by different methods:

Table 5. F measure of different BT strategies

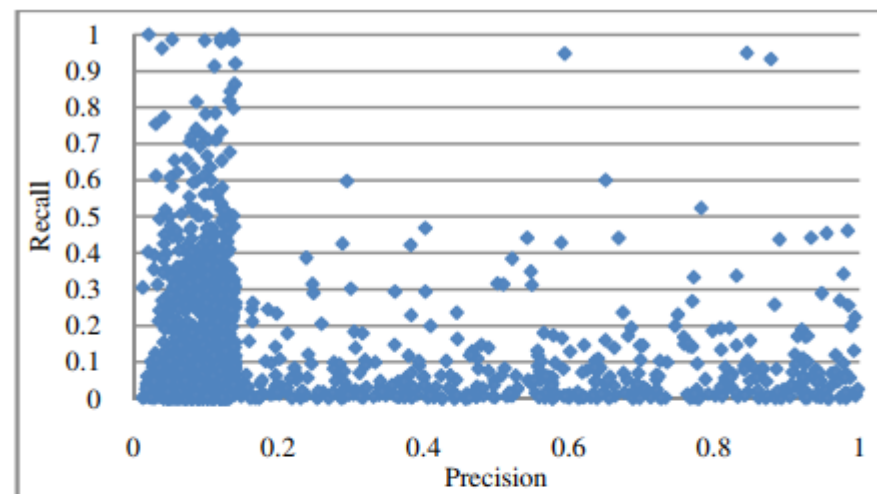


|                           |            | LP     | LQ            | SP     | SQ            |
|---------------------------|------------|--------|---------------|--------|---------------|
| K-means<br>(20 segments)  | <i>Pre</i> | 8.67%  | 8.60%         | 13.35% | <b>17.08%</b> |
|                           | <i>Rec</i> | 10.20% | 22.34%        | 7.63%  | <b>25.58%</b> |
|                           | <i>F</i>   | 0.08   | 0.10          | 0.08   | <b>0.16</b>   |
| CLUTO<br>(20 segments)    | <i>Pre</i> | 8.62%  | 8.56%         | 14.61% | <b>19.13%</b> |
|                           | <i>Rec</i> | 10.01% | 20.51%        | 7.86%  | <b>21.43%</b> |
|                           | <i>F</i>   | 0.08   | 0.10          | 0.07   | <b>0.15</b>   |
| K-means<br>(40 segments)  | <i>Pre</i> | 8.84%  | 9.23%         | 19.76% | <b>20.53%</b> |
|                           | <i>Rec</i> | 9.48%  | 18.20%        | 4.83%  | <b>20.75%</b> |
|                           | <i>F</i>   | 0.08   | 0.10          | 0.06   | <b>0.16</b>   |
| CLUTO<br>(40 segments)    | <i>Pre</i> | 8.76%  | 9.14%         | 19.38% | <b>22.80%</b> |
|                           | <i>Rec</i> | 8.44%  | <b>17.88%</b> | 4.52%  | 17.78%        |
|                           | <i>F</i>   | 0.08   | 0.10          | 0.06   | <b>0.14</b>   |
| K-means<br>(80 segments)  | <i>Pre</i> | 9.02%  | 9.63%         | 23.47% | <b>23.49%</b> |
|                           | <i>Rec</i> | 8.93%  | 17.62%        | 4.06%  | <b>19.35%</b> |
|                           | <i>F</i>   | 0.08   | 0.10          | 0.06   | <b>0.16</b>   |
| CLUTO<br>(80 segments)    | <i>Pre</i> | 8.85%  | 9.51%         | 23.09% | <b>27.00%</b> |
|                           | <i>Rec</i> | 7.82%  | <b>16.65%</b> | 4.00%  | 15.55%        |
|                           | <i>F</i>   | 0.07   | 0.10          | 0.06   | <b>0.15</b>   |
| K-means<br>(160 segments) | <i>Pre</i> | 9.09%  | 9.93%         | 25.68% | <b>25.81%</b> |
|                           | <i>Rec</i> | 8.54%  | 17.98%        | 3.92%  | <b>19.78%</b> |
|                           | <i>F</i>   | 0.074  | 0.10          | 0.06   | <b>0.17</b>   |
| CLUTO<br>(160 segments)   | <i>Pre</i> | 8.87%  | 9.84%         | 25.43% | <b>31.02%</b> |
|                           | <i>Rec</i> | 7.24%  | <b>15.58%</b> | 3.78%  | 14.52%        |
|                           | <i>F</i>   | 0.07   | 0.10          | 0.06   | <b>0.15</b>   |

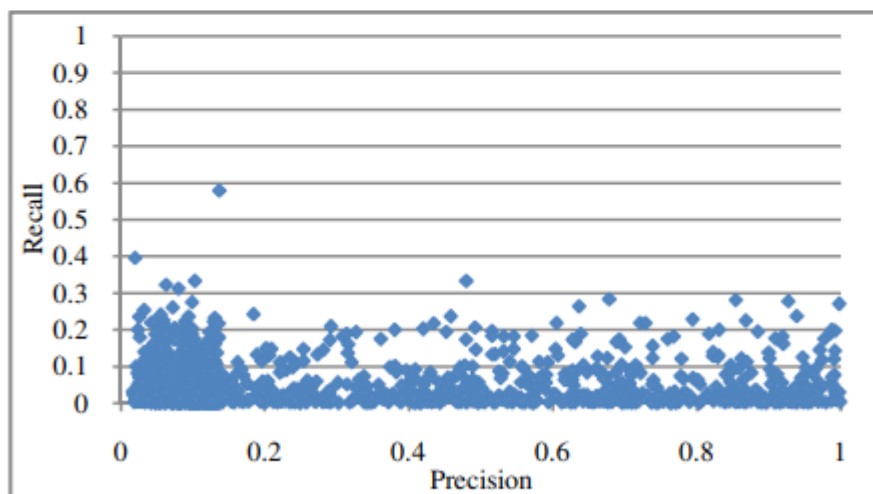
# Precision vs Recall for 160 user segments



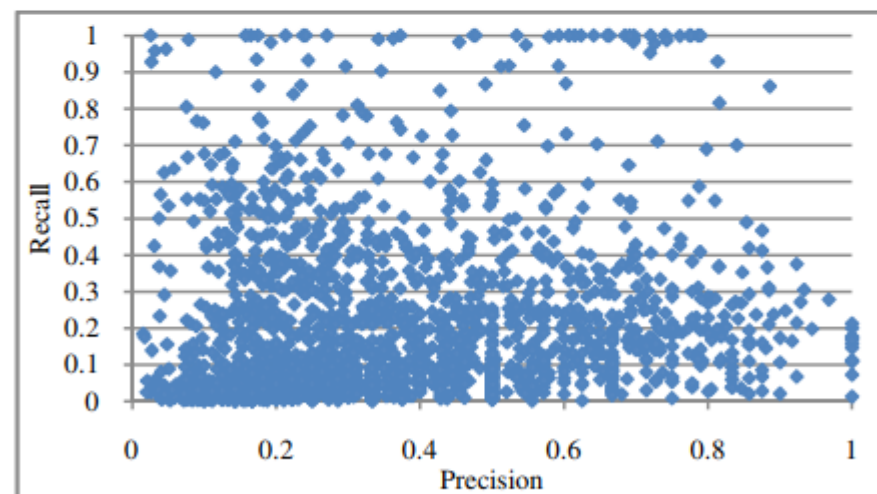
(a)LP



(b)LQ



(c)SP



(d)SQ

**Figure 3. Scatter plot of Precision and Recall over all the ads (CLUTO-160 user segments).**

## Take away message:

- Behavioral targeting can help online advertisement
- Users who click the same ads have similar behavior and they are 91 times more similar than those who don't
- User segmentation based on search queries gives better results than pages visited
- It is better to use short window to segment users
- Increasing the number of segments provides better targeting.
- CTR can be improved by 670% using user segmentation

# References:

- A semantic approach to contextual advertising

SIGIR '07 Proceedings of the 30th annual international ACM SIGIR conference on Research and development in information retrieval  
Pages 559-566

<http://clair.si.umich.edu/~radev/767w10/papers/Week12/ca/semantic.pdf>

- How much can behavioral targeting help online advertising?

WWW '09 Proceedings of the 18th international conference on World wide web

Pages 261-270

<http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.215.1473&rep=rep1&type=pdf>

# Why did the 2<sup>nd</sup> paper cite the 1<sup>st</sup> paper?

- Both papers introduce methods to deliver ads related to users' interests

