

# Optimizing Search Engines using Clickthrough Data

Thorsten Joachims

Presented by Boty Dimanov

# Overview

1. new algorithm for ranking
2. a way to **personalize** search engine queries

- Data
- Method
- Experiments

# Data Collection

- clickthrough
  - $(q, r, c)$  –  $q \in \text{query}$ ,  $r \in (N \times \text{links})$ ,  $c \equiv \text{clickedLinks} \subset \text{domain}(r)$
- easy acquisition
- information contained
  - relative relevance



# Data Collection

- clickthrough
  - $(q, r, c)$  –  $q \in \text{query}$ ,  $r \in \text{ranking}$ ,  $c = \text{linksClicked}(q, r)$
- easy acquisition
- information contained
  - relative relevance
  - partial relative relevance

**1. Kernel Machines**

*<http://svm.first.gmd.de/>*

**2. Support Vector Machine**

*<http://jbolivar.freesevers.com/>*

**3. SVM-Light Support Vector Machine**

*[http://ais.gmd.de/~thorsten/svm\\_light/](http://ais.gmd.de/~thorsten/svm_light/)*

**4. An Introduction to Support Vector Machines**

*<http://www.support-vector.net/>*

**5. Support Vector Machine and Kernel Methods References**

*<http://svm.research.bell-labs.com/SVMrefs.html>*

**6. Archives of SUPPORT-VECTOR-MACHINES@JISCMAIL.AC.UK**

*<http://www.jiscmail.ac.uk/lists/SUPPORT-VECTOR-MACHINES.html>*

**7. Lucent Technologies: SVM demo applet**

*<http://svm.research.bell-labs.com/SVT/SVMsvt.html>*

**8. Royal Holloway Support Vector Machine**

*<http://svm.dcs.rhbnc.ac.uk/>*

**9. Support Vector Machine - The Software**

*<http://www.support-vector.net/software.html>*

**10. Lagrangian Support Vector Machine Home Page**

*<http://www.cs.wisc.edu/dmi/lsvm>*

$$\text{link}_3 <_{\Gamma^*} \text{link}_2 \quad \text{link}_7 <_{\Gamma^*} \text{link}_2$$

$$\text{link}_7 <_{\Gamma^*} \text{link}_4$$

$$\text{link}_7 <_{\Gamma^*} \text{link}_5$$

$$\text{link}_7 <_{\Gamma^*} \text{link}_6$$

$$\text{link}_i <_{\Gamma^*} \text{link}_j$$

for all pairs  $1 \leq j < i$ , with  $i \in C$  and  $j \notin C$ .

$$\text{link}_3 <_{\Gamma^*} \text{link}_2 \quad \text{link}_7 <_{\Gamma^*} \text{link}_2$$

$$\text{link}_7 <_{\Gamma^*} \text{link}_4$$

$$\text{link}_7 <_{\Gamma^*} \text{link}_5$$

$$\text{link}_7 <_{\Gamma^*} \text{link}_6$$

$$\text{link}_i <_{\Gamma^*} \text{link}_j$$

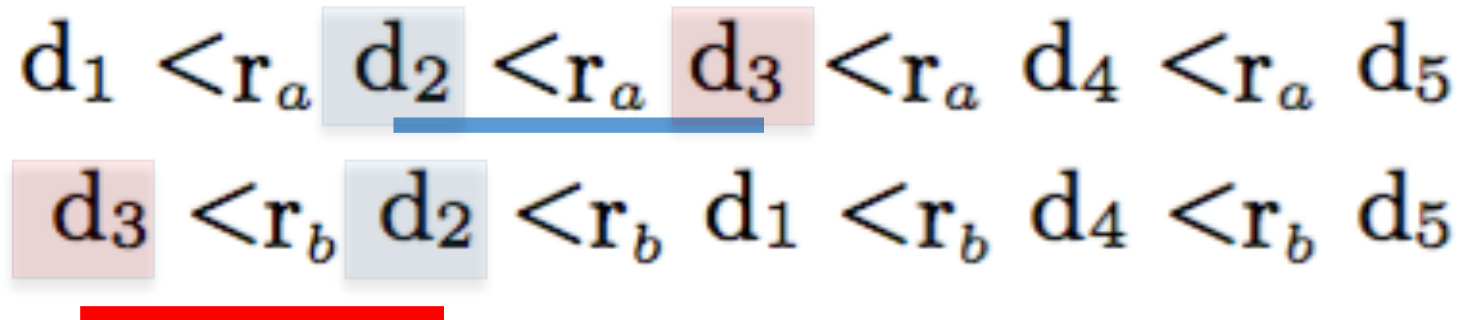
for all pairs  $1 \leq j < i$ , with  $i \in C$  and  $j \notin C$ .

unsuitable format for machine learning algorithms

# Learn Ranking

- How good is a ranking?

- Kendall's  $\tau \equiv \tau(r_a, r_b) = \frac{P - Q}{P + Q} = \frac{1 - 2Q}{\binom{m}{2}}$

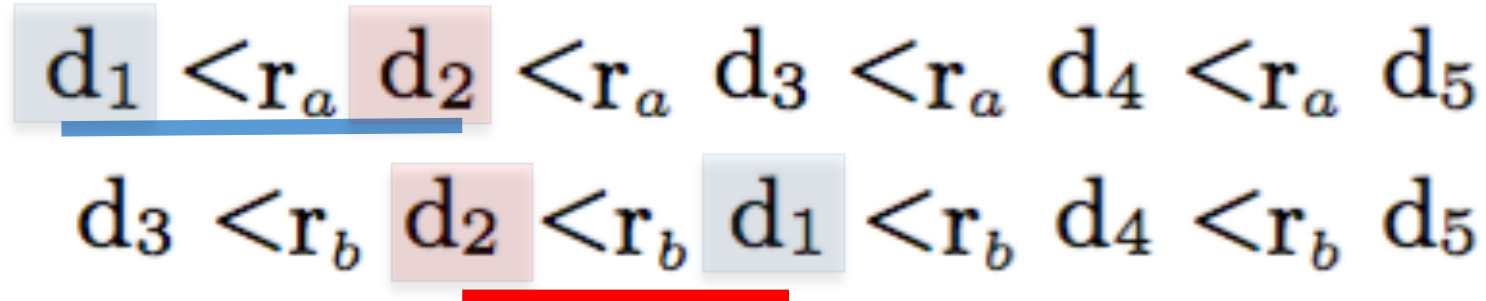




# Learn Ranking

- How good is a ranking?

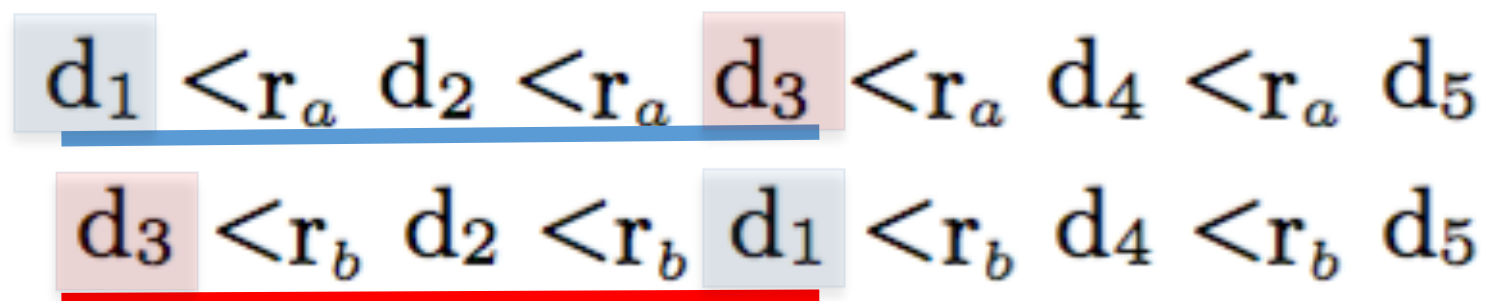
- Kendall's  $\tau \equiv \tau(r_a, r_b) = \frac{P - Q}{P + Q} = \frac{1 - 2Q}{\binom{m}{2}}$



# Learn Ranking

- How good is a ranking?

- Kendall's  $\tau \equiv \tau(r_a, r_b) = \frac{P - Q}{P + Q} = \frac{1 - 2Q}{\binom{m}{2}}$



# Learn Ranking

- How good is a ranking?

- Kendall's  $\tau \equiv \tau(r_a, r_b) = \frac{P - Q}{P + Q} = \frac{1 - 2Q}{\binom{m}{2}}$

- higher  $\tau$  means higher similarity

- appropriate measure for IR

- learn retrieval function

$$\max \tau(f)$$

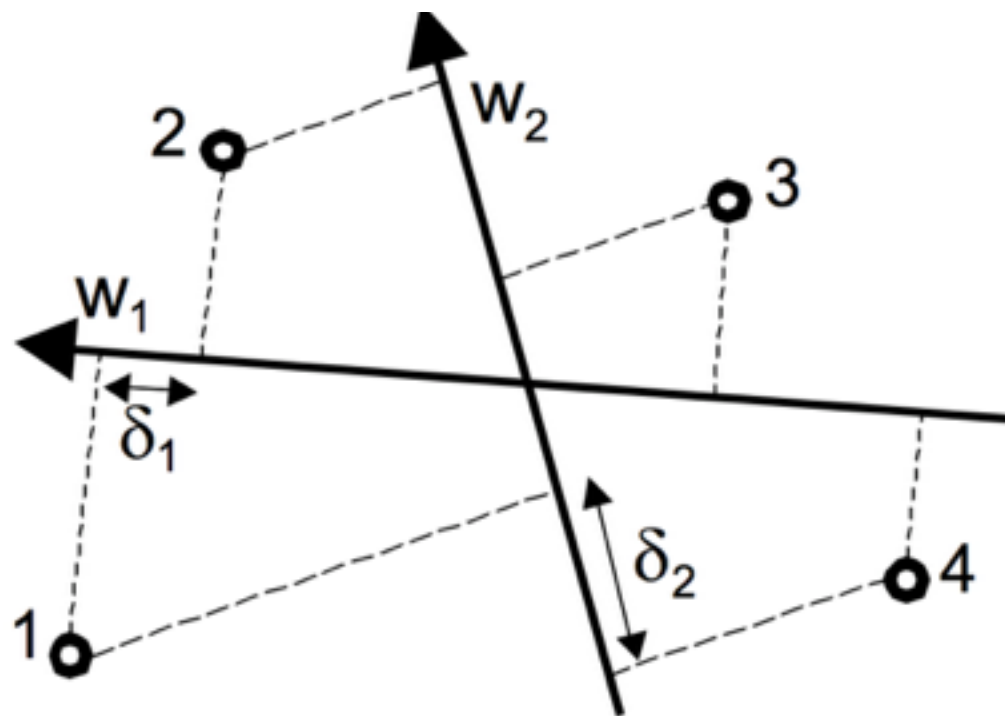
$$\tau(f) = \int \tau(r_{f(q)}, r^*) d\Pr(q, r^*)$$

$$\tau(f) = \frac{1}{n} \sum_{(q, r^*)} \tau(r_{f(q)}, r^*)$$

$$Q = |y - f(x)| \sim \tau(r_{f(q)}, r^*)$$

# Retrieval Function

$$(d_i, d_j) \in f_{\vec{w}}(q) \iff \vec{w}\Phi(q, d_i) > \vec{w}\Phi(q, d_j).$$



# Optimization problem

$$\bullet \max \tau(f) = \frac{1}{n} \sum_{(q, r^*) \in \text{train}} (r_{f(q)}, r^*) \quad \begin{matrix} \text{=} \\ \text{=} \end{matrix} \quad \min \#Q$$

$$\text{Kendall's } \tau \equiv \tau(r_a, r_b) = \frac{P - Q}{P + Q} = \frac{1 - 2Q}{\binom{m}{2}}$$

$$\forall (d_i, d_j) \in r_1^* : \vec{w}\Phi(q_1, d_i) > \vec{w}\Phi(q_1, d_j)$$

...

$$\forall (d_i, d_j) \in r_n^* : \vec{w}\Phi(q_n, d_i) > \vec{w}\Phi(q_n, d_j)$$

## OPTIMIZATION PROBLEM 1. (RANKING SVM)

$$\text{minimize:} \quad V(\vec{w}, \vec{\xi}) = \frac{1}{2} \vec{w} \cdot \vec{w} + C \sum \xi_{i,j,k} \quad (12)$$

subject to:

$$\begin{aligned} \forall (d_i, d_j) \in r_1^* : \vec{w}\Phi(q_1, d_i) &\geq \vec{w}\Phi(q_1, d_j) + 1 - \xi_{i,j,1} \\ &\dots \end{aligned} \quad (13)$$

$$\begin{aligned} \forall (d_i, d_j) \in r_n^* : \vec{w}\Phi(q_n, d_i) &\geq \vec{w}\Phi(q_n, d_j) + 1 - \xi_{i,j,n} \\ \forall i \forall j \forall k : \xi_{i,j,k} &\geq 0 \end{aligned} \quad (14)$$

$$\vec{w} (\Phi(q_k, d_i) - \Phi(q_k, d_j)) \geq 1 - \xi_{i,j,k},$$

$$(d_i, d_j) \in f_{\vec{w}^*}(q)$$

$$\iff \vec{w}^* \Phi(q, d_i) > \vec{w}^* \Phi(q, d_j)$$

$$\iff \sum \alpha_{k,l}^* \Phi(q_k, d_l) \Phi(q, d_i) > \sum \alpha_{k,l}^* \Phi(q_k, d_l) \Phi(q, d_j)$$

$$rsv(q, d_i) = \vec{w}^* \Phi(q, d_i) = \sum \alpha_{k,l}^* \Phi(q_k, d_l) \Phi(q, d_j)$$

# Experiments

- Offline – SVM learn a retrieval function  $f$  . max Kendel's  $\tau$
- Online - max Kendel's  $\tau$  improves retrieval quality
- “Striver” - combination method



**Ranking A:**

1. Kernel Machines  
<http://svm.first.gmd.de/>
2. SVM-Light Support Vector Machine  
<http://ais.gmd.de/~thorsten/svm-light/>
3. Support Vector Machine and Kernel ... References  
<http://svm.....com/SVMrefs.html>
4. Lucent Technologies: SVM demo applet  
<http://svm.....com/SVT/SVMsvt.html>
5. Royal Holloway Support Vector Machine  
<http://svm.dcs.rhbnc.ac.uk/>
6. Support Vector Machine - The Software  
<http://www.support-vector.net/software.html>
7. Support Vector Machine - Tutorial  
<http://www.support-vector.net/tutorial.html>
8. Support Vector Machine  
<http://jbolivar.freesevers.com/>

**Ranking B:**

1. Kernel Machines  
<http://svm.first.gmd.de/>
2. Support Vector Machine  
<http://jbolivar.freesevers.com/>
3. An Introduction to Support Vector Machines  
<http://www.support-vector.net/>
4. Archives of SUPPORT-VECTOR-MACHINES ...  
<http://www.jiscmail.ac.uk/lists/SUPPORT...>
5. SVM-Light Support Vector Machine  
<http://ais.gmd.de/~thorsten/svm-light/>
6. Support Vector Machine - The Software  
<http://www.support-vector.net/software.html>
7. Lagrangian Support Vector Machine Home Page  
<http://www.cs.wisc.edu/dmi/lsvm>
8. A Support ... - Bennett, Blue (ResearchIndex)  
<http://citeseer.../bennett97support.html>

**Combined Results:**

1. **Kernel Machines**  
<http://svm.first.gmd.de/>
2. Support Vector Machine  
<http://jbolivar.freesevers.com/>
3. **SVM-Light Support Vector Machine**  
<http://ais.gmd.de/~thorsten/svm-light/>
4. An Introduction to Support Vector Machines  
<http://www.support-vector.net/>
5. Support Vector Machine and Kernel Methods References  
<http://svm.research.bell-labs.com/SVMrefs.html>
6. Archives of SUPPORT-VECTOR-MACHINES@JISCMail.AC.UK  
<http://www.jiscmail.ac.uk/lists/SUPPORT-VECTOR-MACHINES.html>
7. **Lucent Technologies: SVM demo applet**  
<http://svm.research.bell-labs.com/SVT/SVMsvt.html>
8. Royal Holloway Support Vector Machine  
<http://svm.dcs.rhbnc.ac.uk/>
9. Support Vector Machine - The Software  
<http://www.support-vector.net/software.html>
10. Lagrangian Support Vector Machine Home Page  
<http://www.cs.wisc.edu/dmi/lsvm>

**rank<sub>X</sub>**: 100 minus rank in  $X \in \{\text{Google, MSN-Search, Altavista, Hotbot, Excite}\}$  divided by 100 (minimum 0)

**top1<sub>X</sub>**: ranked #1 in  $X \in \{\text{Google, MSNSearch, Altavista, Hotbot, Excite}\}$  (binary  $\{0, 1\}$ )

**top10<sub>X</sub>**: ranked in top 10 in  $X \in \{\text{Google, MSN-Search, Altavista, Hotbot, Excite}\}$  (binary  $\{0, 1\}$ )

**top50<sub>X</sub>**: ranked in top 50 in  $X \in \{\text{Google, MSN-Search, Altavista, Hotbot, Excite}\}$  (binary  $\{0, 1\}$ )

**top1count<sub>X</sub>**: ranked #1 in  $X$  of the 5 search engines

**top10count<sub>X</sub>**: ranked in top 10 in  $X$  of the 5 search engines

**top50count<sub>X</sub>**: ranked in top 50 in  $X$  of the 5 search engines

2. Query/Content Match (3 features total):

**query\_url\_cosine**: cosine between URL-words and query (range  $[0, 1]$ )

**query\_abstract\_cosine**: cosine between title-words and query (range  $[0, 1]$ )

**domain\_name\_in\_query**: query contains domain-name from URL (binary  $\{0, 1\}$ )

3. Popularity-Attributes ( $\sim 20.000$  features total):

**url\_length**: length of URL in characters divided by 30

**country<sub>X</sub>**: country code  $X$  of URL (binary attribute  $\{0, 1\}$  for each country code)

**domain<sub>X</sub>**: domain  $X$  of URL (binary attribute  $\{0, 1\}$  for each domain name)

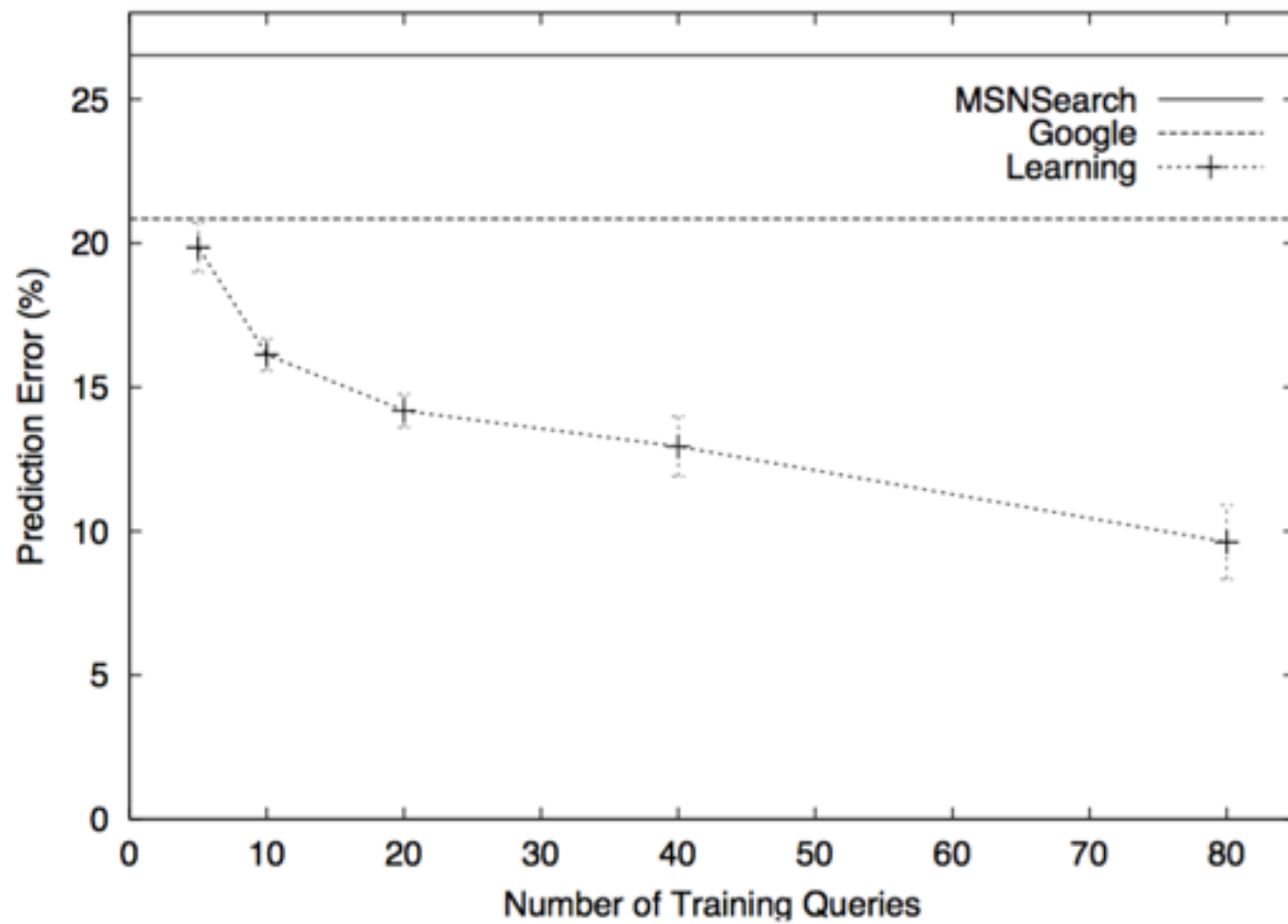
**abstract\_contains\_home**: word "home" appears in URL or title (binary attribute  $\{0, 1\}$ )

**url\_contains\_tilde**: URL contains "~" (binary attribute  $\{0, 1\}$ )

**url<sub>X</sub>**: URL  $X$  as an atom (binary attribute  $\{0, 1\}$ )

# Offline

- • Combination Google + MSN Search
- 112 queries
- ranked  $d_i <_r$  random  $d_j$



# Online

- training - 216 queries
- evaluation - Combination (Learned + random SE)

Comparison	more clicks on learned	less clicks on learned	tie (with clicks)	no clicks	total
Learned vs. Google	29	13	27	19	88
Learned vs. MSNSearch	18	4	7	11	40
Learned vs. Toprank	21	9	11	11	52

weight	feature
0.60	query_abstract_cosine
0.48	top10_google
0.24	query_url_cosine
0.24	top1count_1
0.24	top10_msnsearch
0.22	host_citeseer
0.21	domain_nec
0.19	top10count_3
0.17	top1_google
0.17	country_de
...	
0.16	abstract_contains_home
0.16	top1_hotbot
...	
0.14	domain_name_in_query
...	
-0.13	domain_tu-bs
-0.15	country_fi
-0.16	top50count_4
-0.17	url_length
-0.32	top10count_0
-0.38	top1count_0

**Table 3: Features with largest and smallest weights as learned from the training data in the online experiment.**

# Conclusion

- Ranking SVM can successfully learn an improved retrieval function
- function automatically adapts
- Good results - taking the best of all search engines