User Behavior Analysis for Commercial Search Engines

Yiqun Liu Information Retrieval Group Department of Computer Science and Technology Tsinghua University







- * Tsinghua National Laboratory for Information Science and Technology
 - * One of the five national laboratories, only one in IT field
- * THUIR: our group
 - * Focused on IR researches since 2001
 - * <u>http://www.thuir.org/</u>





- * Research Interests
 - * Information retrieval models and algorithms
 - * Web search technologies
 - Computational social science
- * Members
 - * Leader: Prof. Shaoping Ma;
 - Professors: Min Zhang, Yijiang Jin, Yiqun Liu;
 - * Students: 11 Ph. D. students, 11 master students and 6 undergraduate students.











* Cooperation with industries

- * Tsinghua-Sohu joint lab on search engine technology
- * Tsinghua-Baidu joint course for undergraduate students: Fundamentals of Search Engine Technology
- * Tsinghua-Google joint course for graduate students: Search Engine Product Design and Implementation









- * For search engine: how to attract more users?
 - * To help users to meet their information needs
- * Key challenges (Google's viewpoint)
 - Challenges proposed by Henzinger et.al. (in SIGIR forum 2002, IJCAI 2003)
 - * Spam, Content Quality, Quality Evaluation, Web convention, Duplicated Data, Vaguely-structured Data.
 - * Challenges proposed by Amit Singhal (in SIGIR 2005, ECIR 2008)
 - * Search Engine Spam, Evaluation





* Research issues (our viewpoint)





* Research issues (our viewpoint)

- * Analysis on user's information need Research basics
- * Web Spam fighting
- * Search performance evaluation

Similar with google's challenges

- * How to meet the challenges
 - * With the help of "wisdom of the crowd"
 - * The "Ten thousand cent" project
- * Information sources
 - * user behavior information: search log, Web access log, input log, ...





- * User behavior & information need
- * Web spam fighting
- * Search performance evaluation





- * An important interaction function for search users
 - * Organize a better query
 - * Recommend related information

*	information retrieval的 retrieval retrieval翻译 3d model retrieval	相关搜索				they juery
*	information ratriaval	GOOOOO 1 2 3 4 5	000000gle 6 7 8 9 10	▶ <u>下一页</u>		ks on
1 [2] [3]	[4] [5] [6] [7] [8] [9] [1	10] 下一页	找到相关结果约9	,930,000个	q	
相关搜索	<u>信息检索与利用</u> 未检索到投档信息	<u>计算机信息检索</u> <u>电大信息检索与利用</u>	<u>word 信息检索</u> <u>信息检索试题</u>	<u>excel 信息检索</u> <u>信息检索教程</u>	<u>信息检索论文</u> <u>信息检索心得</u>	37
信息检索	Ę			百度一下	结果中找 帮助 举报	高级搜索



- * Previous solutions
 - * Recommending *similar* queries which were *previously proposed by users*.
 - * How to define "similarity"?
 - * Content based method (Fonseca, 2003; Baeza-Yates, 2004, 2007)
 - * Click-context based method (Wen et.al, 2001; Zaiane et.al, 2002; Cucerzan, 2007; Liu, 2008)
 - * *Problem:* We cannot suppose the recommended queries are better at representing information need. They are even not expressing a same information need.



* Query recommendation for "WWW 2010"

#	Baidu	Google China	Sogou
1	pes2010 (a popular computer game)	2010国家公务员职位表 (National civil service positions for 2010)	2010年国家公务员 (National civil service exam in 2010)
2	qq2010 (a software)	2010年国家公务员报名 (National civil service exam registration in 2010)	2010发型 (fashion hair styles in 2010)
3	实况2010 (a popular computer game)	2010国家公务员报名 (National civil service exam registration in 2010)	2010年考研报名 (Graduate entrance exam in 2010)
4	实况足球2010 (a popular computer game)		2010公务员报名 (civil service exam registration in 2010)
5	卡巴斯基2010 (Kaparsky 2010)	V	2010公务员考试 (civil service exam 2010)



- * How users describe their information needs?
 - * In their queries? May or may not...
 - * In the document they clicked? May or may not
 - * In the snippets they clicked? **Probably!**





* The probability of clicking a certain document is decided by both whether user views the snippet and whether user is interested in it.

 $P(click_i) = P(click_i, view_i) = P(view_i)P(click_i | view_i)$

* Users can only view the snippet while clicking, so $P(click_i | view_i) = P(snippet_i | Need)$

* Therefore,

 $P(click_i) = P(view_i)P(snippet_i | Need)$ $\Rightarrow P(snippet_i | Need) = \frac{P(click_i)}{P(view_i)}$



- Query recommendation performance
 - * Click-through data from September, 2009
 - * 9000 queries were randomly sampled as the test set (each was queried at least 20 times)





- * Find related queries for a given search topic
 - * e.g. find Epidemic related queries
- * Application: seasonal epidemic tendency tracing and predicting
 - * HFMD (hand foot mouth disease) prediction for Beijing in 2010
 - Varicella prediction for Beijing in 2009





* Find related queries for a given search topic * e.g. Find out whether users will buy a car



User behavior & information need

* Selected publications

iever @ Tsinghua University

- * Yiqun Liu, Junwei Miao, Min Zhang, Shaoping Ma, Liyun Ru. How Do Users Describe Their Information Need: Query Recommendation based on Snippet Click Model. Expert Systems With Applications. 38(11): 13847-13856, 2011.
- Danqing Xu, Yiqun Liu, Min Zhang, Liyun Ru, Shaoping Ma. Predicting Epidemic Tendency through Search Behavior Analysis. In Proceedings of the 22nd International Joint Conference on Artificial Intelligence (IJCAI-11) (Barcelona, Spain). 2361-2366.
- * Weize Kong, Yiqun Liu, Shaoping Ma and Liyun Ru. 2010. Detecting epidemic tendency by mining search logs. In Proceedings of the 19th WWW Conference. WWW '10. ACM, New York, NY, 1133-1134.
- Rongwei Cen, Yiqun Liu, Min Zhang, Liyun Ru, Shaoping Ma. 2010. Study language models with specific user goals. In Proceedings of the 19th WWW Conference . WWW '10. ACM, New York, NY, 1073-1074.



- * User behavior & information need
- * Web spam fighting
- * Search performance evaluation





Web spam fighting

* Spam pages are everywhere

金碑碑文范例 - 搜狗搜索 - Windows Internet Explorer ロ ロ × ロ ロ マ ロ マ ロ マ	
	P -
後 墓碑碑文范例 - Google 搜索 - Windows Internet Explorer	
	• م
🖕 收藏夹 🔡 🗸 🚼 臺磚磚文范例 - Googl X 🏈 臺磚磚文范例 🛛 👔 🤹 收藏夹 🔡 🗸 🚼 臺磚磚文范例 - Google 🏈 臺磚磚文范例 🗙	
网页图片视频地图资讯音乐问答»来吧»更多▼	×
Iuyiqun03@gmail.com 网络历史记录 设置 ▼ 退出 Google 臺碑碑文范例 Google 臺碑碑文范例 Google 臺級	团申请书注 加 加 面 面 面 五 位 の の 二 の の の 一 の の 二 の の の 一 の の こ の の の の つ 二 の の の の の の の の の の の の の
◎ 所有网页 ◎ 中文网页 ◎ 简体中文网页	<u>通作文范後</u> 觀信范例
	<u> </u>
算確確文范例 全国线到付款、高保障! 经一点、不然会叫的好: 日 china.Alibaba.com 全球千万网商丰富供求信息,交易便捷上阿里巴巴1688.com,丧葬用品销售批发 開部温暖湿润,性爱结合? 開部温暖湿润,性爱结合? 開部温暖湿润,性爱结合? 開部温暖湿润,性爱结合? 開設温暖湿润,性爱结合? 開設温暖湿润,性爱结合? 日	作计划范码 适作文范例
<u>摹碑碑文的写法规范-曾祥裕-职业日志-价值中国网,网络就是社会…</u> ☆ 2009年3月30日 范例:王洛宾薹碑碑文公元一九九六年三月十四日凌晨洛宾仙逝,二十 日向遗体告别,二十二日送骨灰进京,边城乌鲁木齐三降大雪。 www.chinavalue.net/Blog/142756.aspx - <u>网页快照</u>	论文缩过。 <u>出论文格式</u> <u>间墓碑碑文</u> <u>间墓碑碑文</u> 主荐站点 主荐站点 <u>打装工具</u> 巨名 <u>遗传</u>
<u>碑文范例(碑后文 墓志铭)殡葬原著选读</u> ☆	n tmotion 生优化大
Re:確文范例(確后文) 基本 (1) 2009-11-23 18:27:57 我父亲不幸遇车祸去世,准备立 差 (1) 算法 (1) [1] [1] [1] [1] [1] [1] [1] [1] [1] [1]	indows9) × = 司系统
臺碑碑文范例 章梁碑碑文范例 章子派母崔堂派教育 章子派子王子 章子派子 章子派母崔堂派教育 章子派子 章子派子 章子派子 章子》 章子派子 章子》 章子派子 章子 章子 章子 章子 章子 章子 章子 章子 章子 章	幸运用 讯公司 星Q40
	۲ ۱00% - 100%



* Definition:

- * Web spam are designed to get "an unjustifiably favorable relevance or importance score" from search engines. (*Gyongyi et. al.* 2005)
- * How many spams are there on the Web?
 - * Over 10% Web pages are spams (*Fetterly et al.* 2004, *Gyöngyi et al.* 2004)
 - * Billions of spam pages...
- * How many can search engine index?
 - * Google: 8 billion@2004, Yahoo: 20 billion@2005



- * An important and difficult task
 - * Baidu.com: We banned over 30,000 spam sites each day on average. In the research field of Web spam fighting, we even spend more money than the whole Chinese search market value. (14 November, 2008)
 - * Why so difficult?
 - * Too many kinds of spamming techniques
 - keyword farm, link farm, weaving, cloaking, javascript/iframe redirecting, ...
 - * 道高一尺, 魔高一丈! (however persuasive good is, evil is still stronger)



- * Problems with existing methods
 - * Focus on existing spamming techniques, cannot deal with newly-appeared ones.
 - * How to identify spamming techniques you never see?
- * Our solution: spam v.s. users



- Containing no useful information
- Try to cheat search engines
- Try to attract more users



- Want to obtain useful information
- Rely on search engines
- Try to avoid visiting spam pages



- * Our solution (cont.)
 - * What do users do when they meet spams?
 - * What do users do when they visit ordinary pages?
- * User behavior features for spam fighting
 - * Search Engine Oriented Visit Rate
 - * Source Page Rate
 - * Short-time Navigation Rate
 - * Query Diversity
 - * Spam Query Number

*



* User behavior features for spam fighting (cont.)





- * Spam identification performance
 - * Better at identifying newly-appeared spam types
 - * Identified 1,000 spam sites on 2008/03/02; commercial search engines didn't recognize them until 2008/03/26
 - * Outperforms previous anti-spam algorithms

Algorithm	Recall = 25.00%	Recall = 50.00%	Recall = 75.00%	AUC
Content-based algorithm [Cormack et al. 2011]	81.63%	7.65%	4.08%	0.6414
Link-based algorithm [Gyöngyi et al. 2004]	74.43%	34.09%	18.75%	0.7512
User behavior algorithm	100.00%	76.14%	43.75%	0.9150



- * What if we cannot collect user browsing logs?
- * Search engine click-through logs may be enough...
- * Spam keywords are
 - * hot or reflect a heavy demand of search users
 - * lack of key recourses or authoritative results
- * Keyword Vampire
 - * Transform profitable keywords into affiliate links in a snap
 - * http://www.keywordvampire.com/





* A Label Propagation algorithm on query-URL bipartite graph

Query





- * Spam detection performance
 - Performs better than PageRank & TrustRank, works well together with PageRank & TrustRank
 - * A small seed set is enough to gain good performance





* Selected publications

- Yiqun Liu, Fei Chen, Weize Kong, Huijia Yu, Min Zhang, Shaoping Ma, Liyun Ru. Identifying Web Spam with the Wisdom of the Crowds. Accepted by ACM Transaction on the Web.
- * Chao Wei, Yiqun Liu, Min Zhang, Liyun Ru, Shaoping Ma, Kuo Zhang. Fighting against Web Spam: A Novel Propagation Method based on Clickthrough Data. Proceedings of the 35th Annual ACM SIGIR Conference (SIGIR 2012). ACM, New York, NY, 2012.
- * Data Cleansing for Web Information Retrieval using Query Independent Features. Yiqun Liu, Min Zhang, Liyun Ru, Shaoping Ma. Journal of the American Society for Information Science and Technology (JASIST), 58(12), Pages 1884-1898, 2007.
- Yiqun Liu, Yijiang Jin, Min Zhang, Shaoping Ma, Liyun Ru, User Browsing Graph: Structure, Evolution and Application. Late breaking result session in Second ACM International Conference on Web Search and Data Mining (WSDM 2009).



- * User behavior & information need
- * Web spam fighting
- * Search performance evaluation



Search performance evaluation

* Evaluation is important for search engines

Tsinahua Ilai

- * **Research**: Evaluation became central to R&D in IR to such an extent that new designs and proposals and their evaluation became one. (*Saracevic*, SIGIR 1995)
- * Advertising: Search advertisers choose the most profitable platform.
- * **Engineering:** Search engineers has to decide whether proposed algorithms are effective.
- * Cranfield-like evaluation approaches
 - * A set of query topics, their corresponding answers (usually called qrels) and evaluation metrics.

Search performance evaluation

- * Problems with previous Cranfield-like method
 - * Labor intensive: 9 people months are required to judge 1 topic for a collection of 8M documents. (Voorhees, 2001)
 - * Objective: Assessors disagree on 58% documents for a query topic in a task of TREC 2008.
- * Our solution
 - * Annotate answers with the help of wisdom of the crowd.
 - Construction of user click models
 - Satisfaction instead of relevance

* For navigational type queries (e.g. yahoo mail)

Tsinahua Ilni

* Basic assumption: The result clicked by more users should be more relevant than the one clicked by fewer users.



- * For informational/transactional type queries
 - * The basic assumption fails for non-navigational type queries. e.g. the query "电影" (movie)



* For informational/transactional type queries (cont.)

Tsinahus Ilni

* Improved assumption: click-through data from multiple search engines are more informational and less biased than that from a single engine.



* For long-tail queries (cont.)

* Only a few clicks for a long-tail query: each click should make difference in answer annotation process.

No.	Feature Name	Description		
1	NumQueries	Num of unique queries in current session		
2	ClickEntropy	Entropy of click distribution in current session		
3	ClickSelfInformation	Self-information of current doc in current session		
4	4 ClickDocNumInSession Click number of current doc in current session			
5	ClickDocNumInQuery	Click number of current doc in current query		
6	IsFirstClickInSession	=1 if the first click in current session, =0 otherwise		
7	7 IsLastClickInSession =1 if the last click in current session,=0 other			
8	IsFirstClickInQuery	=1 if the first click in current query, =0 otherwise		
9	IsLastClickInQuery	=1 if the last click in current query, =0 otherwise		
10	ClickOrderInSession	Order of current click in current session		
11	ClickOrderInQuery	Order of current click in current query		
12	ClickDocRank	Result rank of current click		

* For long-tail queries (cont.)

@ Tsinghua University

- * Clicks with reliable relevance feedback information is different from unreliable ones.
- * A learning based framework can be adopted to separate reliable clicks.



* For long-tail queries (cont.)

tierevial cuthnist



nks

搜狗搜索

- * Problems with Cranfield-like approaches
 - * Time consuming, objective
 - * Relevance annotation of "query-result" pairs
 - * Ignore the representation of results

<u>新闻 **网页** 音乐 图片 视频 地图 知识 更多>></u>

Sogou 搜狗

inll cuthnisT @

找到约 189,277 条结果(用时 0.095 秒)

C	> 网页结果	<u>清华大学建校</u>
	音乐	清华学堂旧照历史悠久的 <mark>清华大学</mark> 是我国一所引人瞩目、令人向往的著名重点大学。 <mark>清华大学</mark> 的
	図出	前身叫清华学堂,是1911年清朝政府用美国"退还"的一部分庚子赔款办的一所
_		www.todayonhistory.com/eJianXiao.html - 2011-8-16 - <u>快照</u>
> 🔀	视频	
音	知识	胡锦涛在庆祝清华大学建校100周年大会上的讲话
5	新闻	2011年4月24日 新华社北京4月24日电 在庆祝清华大学建校100周年大会上的讲话 胡锦涛(2
Ĕ	9811+41	011年4月24日) 4月24日,庆祝清华大学建校100周年大会在北京人民
初	〉全部时间	www.gov.cn/04/24/content_1851436.htm - 2011-4-24 - <u>快照</u>
矨	一天内	清华大学建校 互动百科
新	一周内	
	一月内	往的著名重点大学。清华大学的前身叫清华学堂,是1911年清朝政府用美国
> 全	一年内	www.hudong.com/wiki/ 清华大学建校 - <u>快照</u>

- * User satisfaction evaluation instead of relevance judgment
 - * What's a satisfied user session?
 - * Navigational: top result should be the target.
 - * Informational: top ranked results answer user's question with a different aspect.
 - * Transactional: user can accomplish task with the top few results.
 - * Behavior patterns in satisfied/unsatisfied search sessions should be different.

- * A number of behavior features
 - * Result click behavior: first click position, last click position, revisit click, non-click, ...
 - * Other click behavior: recommendation click, next page



* Compared with human assessors

	AUC	AUC d	lifference
Human assessor	0.87		/
Informational/Transactional	0.75	-16	5.00%
Navigational	0.80	-8	.75%
	A	B	System
Assessor A as ground truth	1.00	0.80	0.80
Assessor B as ground truth	0.80	1.00	0.76
System as ground truth	0.80	0.76	1.00

* Query frequency v.s. applicable queries

Query frequency	1	2~3	4~10	11~100	100~	top
Percentage of Applicable queries	2.59%	14.96%	36.11%	65.61%	89.56%	94.11%

Search performance evaluation

* Selected publications

triever @ Tsinghua University

- Bo Zhou, Yiqun Liu, Min Zhang, Yijiang Jin, Shaoping Ma, Incorporating Web Browsing Information into Anchor Texts for Web Search, Information Retrieval. Volume 14, Issue 3: 290-314, 2011.
- * Danqing Xu, Yiqun Liu, Min Zhang and Shaoping Ma. Incorporating Revisiting Behaviors into Click Models. Accepted by **the 5th ACM International Conference on Web Search and Data Mining**. WSDM 2012.
- * Yiqun Liu, Yupeng Fu, Min Zhang, Shaoping Ma, Liyun Ru. Automatic Search Engine Performance Evaluation with Click-through Data Analysis. in Proceedings of the 16th international Conference on World Wide WWW '07. ACM, New York, 1133-1134.
- Rongwei Cen, Yiqun Liu, Min Zhang, Bo Zhou, Liyun Ru, Shaoping Ma. Exploring Relevance for Clicks. In Proceeding of the 18th ACM Conference on information and Knowledge Management. CIKM '09. 1847-1850.



- * Key challenges of search engines
- * Meet challenges with the help of "wisdom of the crowd"
 - * User behavior and information need
 - * Web spam fighting
 - * Search performance evaluation





Welcome to visit our homepage

http://www.thuir.cn/

On-line Demos Search Engine Evaluation Seasonal Epidemic Prediction Web Spam Page Identification Web News Event Clustering



- Problems with the automatic answer annotation process
 - * Each click is regarded as a relevance voting for the corresponding result.
 - * However, results aren't equally examined: **Position Bias**



 Model click behavior to solve the position bias problem

 $C_i = 1 \rightarrow \qquad R_i = 1, R_i = 1$

* How to estimate the examine probability?

- * Cascade model: $P(E_{i+1} = 1 | E_i = 1, C_i) = 1 C_i$
- * Dependent click model (DCM):

 $P(E_{i+1} = 1 | E_i = 1, C_i = 0) = 1$ $P(E_{i+1} = 1 | E_i = 1, C_i = 1) = \lambda_i$

* User browsing model (UBM):

$$P(E_i = 1 | C_{1\dots i-1}) = \lambda_{r_i, d_i}$$

* Lots of other models: DBM, CCM, ...

- * Problem with the existed models
 - * Results are not sequentially examined and clicked
 - * Eye-tracking experiments (Lorigo et.al, 2005) show that lots of users revisit
 - * Revisit behavior is popular (24.1% multi-click sessions)



* From ranking position to timing sequence



* The Temporal Hidden Click Model (THCM)

 $p(E_{t(i+1)} = 1 | E_{ti} = 1) = \alpha$ Forward examination

 $p(E_{t(i-1)} = 1 | E_{ti} = 1) = \gamma$ **Backward examination (revisit)**

 $0 < \alpha + \gamma < 1, \alpha > 0, \gamma > 0$

 $p(C_t|C_1, C_2, \dots, C_{t-1}) = p(C_t|C_{t-1})$ One-order model (to be simple) $p(C_t = i) = p(E_{ti} = 1) \cdot p(R_i = 1)$ similar with old models

Data requirement: the click sequence should be recorded



* THCM performance

