

CSE 4/535

Information Retrieval

Sayantan Pal
PhD Student, Department of CSE
338Z Davis Hall



Department of CSE

Before we start

1. Midterm 2 Poll, Please answer by Wednesday (Oct 18)
2. Timeline will be updated based on Poll responses, check website on Friday (Oct 20th)
3. Today's lecture - Evaluation Methodology Result
Summaries (intuitive, less math)
 - a. Midterm 2 - Easy to score



Recap - Previous Class

1. Precision, Recall, F1
2. P@R, 11 point Precision



Today's Lecture...

- Evaluating a search engine
 - MAP score
 - Kappa Measure
 - DCG
 - A/B Testing




Single Value Measures

- Average precision at seen relevant documents
 - Precision figures after each new relevant document is observed are 1, .66, 0.5, 0.4, 0.3
 - Mean Avg precision is $(1+.66+.5+.4+.3)/5$ or 0.57
- R-precision
 - Generate a single value summary of ranking by computing precision at the R-th position in the ranking, where R is the total number of relevant documents
 - E.g. R=10, 4 relevant documents in first ten returned docs, R-precision is 0.4 (precision at 10)
- Precision histograms
 - Used to compare retrieval history of two algorithms



Mean Average Precision

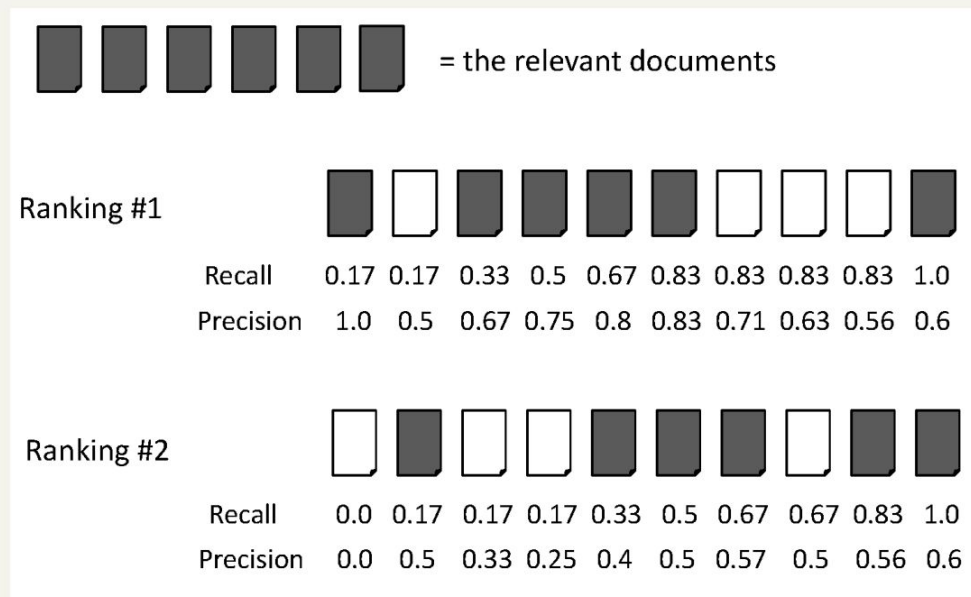
- Consider rank position of each **relevant** doc
 - $K_1, K_2, \dots K_R$
- Compute Precision@K for each $K_1, K_2, \dots K_R$
- Average precision = average of P@K
- Ex:



 has AvgPrec of $\frac{1}{3} \cdot \left(\frac{1}{1} + \frac{2}{3} + \frac{3}{5} \right) \approx 0.76$
- MAP is Average Precision across multiple queries/rankings



Average Precision

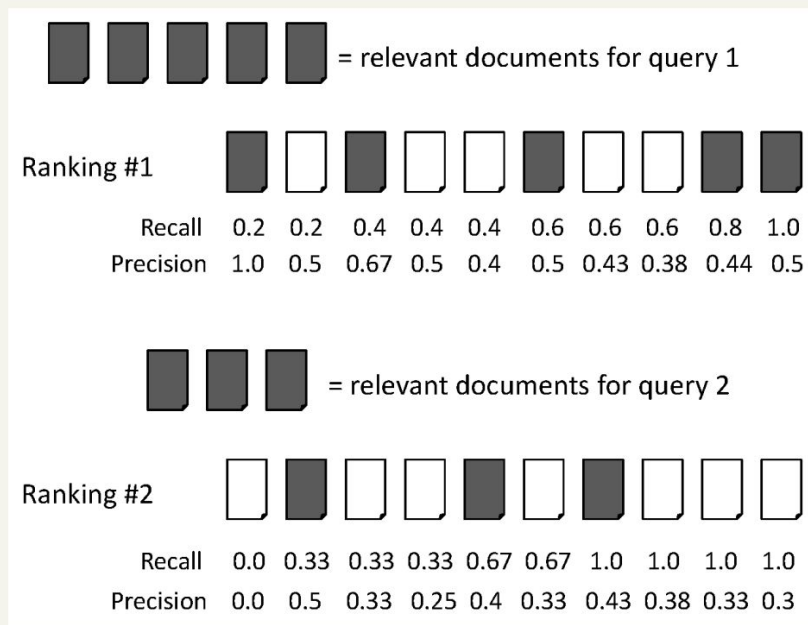


$$\text{Ranking \#1: } (1.0 + 0.67 + 0.75 + 0.8 + 0.83 + 0.6)/6 = 0.78$$

$$\text{Ranking \#2: } (0.5 + 0.4 + 0.5 + 0.57 + 0.56 + 0.6)/6 = 0.52$$



MAP



$$\text{average precision query 1} = (1.0 + 0.67 + 0.5 + 0.44 + 0.5)/5 = 0.62$$

$$\text{average precision query 2} = (0.5 + 0.4 + 0.43)/3 = 0.44$$

$$\text{mean average precision} = (0.62 + 0.44)/2 = 0.53$$



Mean average precision

- If a relevant document never gets retrieved, we assume the precision corresponding to that relevant doc to be zero
- MAP is macro-averaging: each query counts equally
- Now perhaps most commonly used measure in research papers
- Good for web search?
- MAP assumes user is interested in finding many relevant documents for each query
- MAP requires many relevance judgments in text collection

Variance

- For a test collection, it is usual that a system does crummily on some information needs (e.g., $\text{MAP} = 0.1$) and excellently on others (e.g., $\text{MAP} = 0.7$)
- Indeed, it is usually the case that the variance in performance of the same system across queries is much greater than the variance of different systems on the same query.
- That is, there are easy information needs and hard ones!



From document collections to test collections

- Still need
 - Test queries
 - Relevance assessments
- Test queries
 - Best designed by domain experts
 - Random query terms generally not a good idea
- Relevance assessments
 - Human judges, time-consuming
 - Are human panels perfect?

Unit of Evaluation

- We can compute precision, recall, F, and ROC curve for different units.
- Possible units
 - Documents (most common)
 - Facts (used in some TREC evaluations)
 - Entities (e.g., car companies)
 - May produce different results. Why?



Kappa measure for inter-judge (dis)agreement

- Kappa measure
 - Agreement measure among judges
 - Designed for categorical judgments
 - Corrects for chance agreement
- $\text{Kappa} = [P(A) - P(E)] / [1 - P(E)]$
- $P(A)$ – proportion of time judges agree
- $P(E)$ – what agreement would be by chance
- Kappa = 0 for chance agreement, 1 for total agreement.



Kappa Measure: Example

Number of docs	Judge 1	Judge 2
300	Relevant	Relevant
70	Nonrelevant	Nonrelevant
20	Relevant	Nonrelevant
10	Nonrelevant	relevant



Kappa Example

- $P(A) = 370/400 = 0.925$
- $P(\text{nonrelevant}) = (10+20+70+70)/800 = 0.2125$
- $P(\text{relevant}) = (10+20+300+300)/800 = 0.7878$
- $P(E) = 0.2125^2 + 0.7878^2 = 0.665$
- $\text{Kappa} = (0.925 - 0.665)/(1-0.665) = 0.776$

- $\text{Kappa} > 0.8$ = good agreement
- $0.67 < \text{Kappa} < 0.8 \rightarrow$ “tentative conclusions” (Carletta '96)
- Depends on purpose of study
- For >2 judges: average pairwise kappas

Standard relevance benchmarks: Others

- GOV2
 - Another TREC/NIST collection
 - 25 million web pages
 - Largest collection that is easily available
 - But still 3 orders of magnitude smaller than what Google/Yahoo/MSN index
- NTCIR
 - East Asian language and cross-language information retrieval
- Cross Language Evaluation Forum (CLEF)
 - This evaluation series has concentrated on European languages and cross-language information retrieval.
- Many others

Interjudge Agreement: TREC 3

information need	number of docs judged	disagreements	NR	R
51	211	6	4	2
62	400	157	149	8
67	400	68	37	31
95	400	110	108	2
127	400	106	12	94

Impact of Inter-judge Agreement

- Impact on **absolute** performance measure can be significant (0.32 vs 0.39)
- Little impact on ranking of different systems or **relative** performance
- Suppose we want to know if algorithm A is better than algorithm B
- A standard information retrieval experiment will give us a reliable answer to this question.



Critique of pure relevance

- Relevance vs Marginal Relevance
 - A document can be redundant even if it is highly relevant
 - Duplicates
 - The same information from different sources
 - Marginal relevance is a better measure of utility for the user.
- Using facts/entities as evaluation units more directly measures true relevance.
- But harder to create evaluation set



Can we avoid human judgment?

- No – actually, maybe we can use proxies
- Makes experimental work hard
 - Especially on a large scale
- In some very specific settings, can use proxies
 - E.g.: for approximate vector space retrieval, we can compare the cosine distance closeness of the closest docs to those found by an approximate retrieval algorithm
- But once we have test collections, we can reuse them (so long as we don't overtrain too badly)

BEYOND BINARY RELEVANCE



Evaluation at large search engines

- Search engines have test collections of queries and handranked results
- Recall is difficult to measure on the web
- Search engines often use precision at top k , e.g., $k = 10$
- . . . or measures that reward you more for getting rank 1 right than for getting rank 10 right.
 - NDCG (Normalized Discounted Cumulative Gain)
- Search engines also use non-relevance-based measures.
 - Click through on first result
 - Not very reliable if you look at a single click through ... but pretty reliable in the aggregate.
 - Studies of user behavior in the lab
 - A/B testing



YAHOO! Web Images Video Local Shopping More ▾

Toyota safety Search Options ▾

Also try: [toyota safety ratings](#), [toyota safety recall](#), [More...](#)

108,000,000 results for **Toyota safety**:

[Show All](#)

[Toyota](#)

[Motor Trend](#)

[CarsDirect](#)

[Shopping Sites](#)

Toyota Recall
Sponsored Results
Toyota Takes Care of its Customers. Read the FAQs at Toyota.com.
[www.Toyota.com/Recall](#)

Toyota Safety
Sponsored Results
& Latest Prices. Free Info. Toyota Research, Reviews.
[www.Toyota.Edmunds.com](#)

TOYOTA | Car Safety Innovation and Technology
Toyota home page for car safety and car technology Prius model.
[www.safetytoyota.com](#) - [Cached](#)

Toyota home page for car safety and car technology ...
We are presenting Toyota's safety technologies for cars. We clearly explain about car safety and car technology using movies and more.
[www.safetytoyota.com/en-gb](#) - [Cached](#)

Toyota Safety Ratings - Toyota Safety Features - Motor Trend ...
MotorTrend offers Toyota safety ratings, comprehensive auto safety reports, and more. View a all of the standard Toyota safety features. ...
[motortrend.com/new_cars/07/toyota/safety_ratings/index.html](#) - 149k - [Cached](#)

Toyota Motor Europe Corporate Site Safety
Our approach. Toyota believes that all stakeholders in the road safety equation share a responsibility to reduce the frequency of road accidents. ...
[www.toyota.eu/Safety](#) - [Cached](#)

[PDF] pdf European Safety Brochure 2005
4047k - Adobe PDF - [View as html](#)
not guarantee that all accidents or injuries will be avoided when driving a Toyota and/or Lexus brand motor vehicle equipped with the safety systems ...
[www.toyota.no/Images/Safety_Brochure_tcm308-344461.pdf](#)

Toyota - Star Safety System
Star Safety System ... Toyota Mobility Program. Careers. Contact Us. Home, contact us, site map, your privacy rights, legal terms. Toyota Newsroom, sign up for info ...
[www.toyota.com/vehicles/demos/star-safety.html](#) - 68k - [Cached](#)

Toyota Prius Safety Ratings - CarsDirect
Get overall safety ratings and NHTSA crash test results for the Toyota Prius at CarsDirect.

Safety for a Toyota
Sponsored Results
Research Safety Ratings and Reviews For New Car at Kelley Blue Book.
[www.kbb.com](#)

Toyota Safety
Find Toyota Safety dealers, new cars, prices, and photos.
[www.NewCars.org](#)

Toyota Safety
Toyota safety Discount Prices Save Money Shopping Online Today.
[www.smarter.com](#)

Safety Toyota
Explore 5,000+ Pro Sports Choices. Save On Safety Toyota.
[BaseballGear.Shopzilla.com](#)

[See your message here...](#)

fair

fair

Good



DCG : Graded (Non-Binary) Relevance

- DCG: Two assumptions are made in using DCG
 - Highly relevant documents are more useful when appearing earlier in a search engine result list (have higher ranks)
 - Highly relevant documents are more useful than marginally relevant documents, which are in turn more useful than irrelevant documents.
 - highly relevant documents appearing lower in search result should be penalized as the graded relevance value is reduced logarithmically proportional to the position of the result. Discounted CG accumulated at rank position p , where rel_i is the graded relevance (0,1,2,3,4) of the result at position i .

$$DCG_p = rel_1 + \sum_{i=2}^p \frac{rel_i}{\log_2 i}$$



DCG Example

- 10 ranked documents judged on 0–3 relevance scale:
3, 2, 3, 0, 0, 1, 2, 2, 3, 0
- discounted gain:
 $3, 2/1, 3/1.59, 0, 0, 1/2.59, 2/2.81, 2/3, 3/3.17, 0$
 $= 3, 2, 1.89, 0, 0, 0.39, 0.71, 0.67, 0.95, 0$
- DCG:
3, 5, 6.89, 6.89, 6.89, 7.28, 7.99, 8.66, 9.61, 9.61



NDCG

- Search result lists vary in length depending on the query. Comparing a search engine's performance from one query to the next cannot be consistently achieved using DCG alone, so the cumulative gain at each position for a chosen value of p should be normalized across queries.
- Sort documents of a result list by relevance, producing an ideal DCG (IDCG) at position p . For a query, the *normalized discounted cumulative gain*, or nDCG, is computed as:

$$\text{nDCG}_p = \frac{DCG_p}{IDCG_p}$$



NDCG - Example

4 documents: d_1, d_2, d_3, d_4

i	Ground Truth		Ranking Function ₁		Ranking Function ₂	
	Document Order	r_i	Document Order	r_i	Document Order	r_i
1	d4	2	d3	2	d3	2
2	d3	2	d4	2	d2	1
3	d2	1	d2	1	d4	2
4	d1	0	d1	0	d1	0
	NDCG _{GT} =1.00		NDCG _{RF1} =1.00		NDCG _{RF2} =0.9203	

$$DCG_{GT} = 2 + \left(\frac{2}{\log_2 2} + \frac{1}{\log_2 3} + \frac{0}{\log_2 4} \right) = 4.6309$$

$$DCG_{RF1} = 2 + \left(\frac{2}{\log_2 2} + \frac{1}{\log_2 3} + \frac{0}{\log_2 4} \right) = 4.6309$$

$$DCG_{RF2} = 2 + \left(\frac{1}{\log_2 2} + \frac{2}{\log_2 3} + \frac{0}{\log_2 4} \right) = 4.2619$$

$$MaxDCG = DCG_{GT} = 4.6309$$



Mean Reciprocal Rank

- Consider rank position, K , of first relevant doc
 - Could be – only clicked doc
- Reciprocal Rank score = $\frac{1}{K}$
- MRR is the mean RR across multiple queries



Human judgments are

- Expensive
- Inconsistent
 - Between raters
 - Over time
- Decay in value as documents/query mix evolves
- Not always representative of “real users”
 - Rating vis-à-vis query, don’t know underlying need
 - May not understand meaning of terms, etc.
- So – what alternatives do we have?



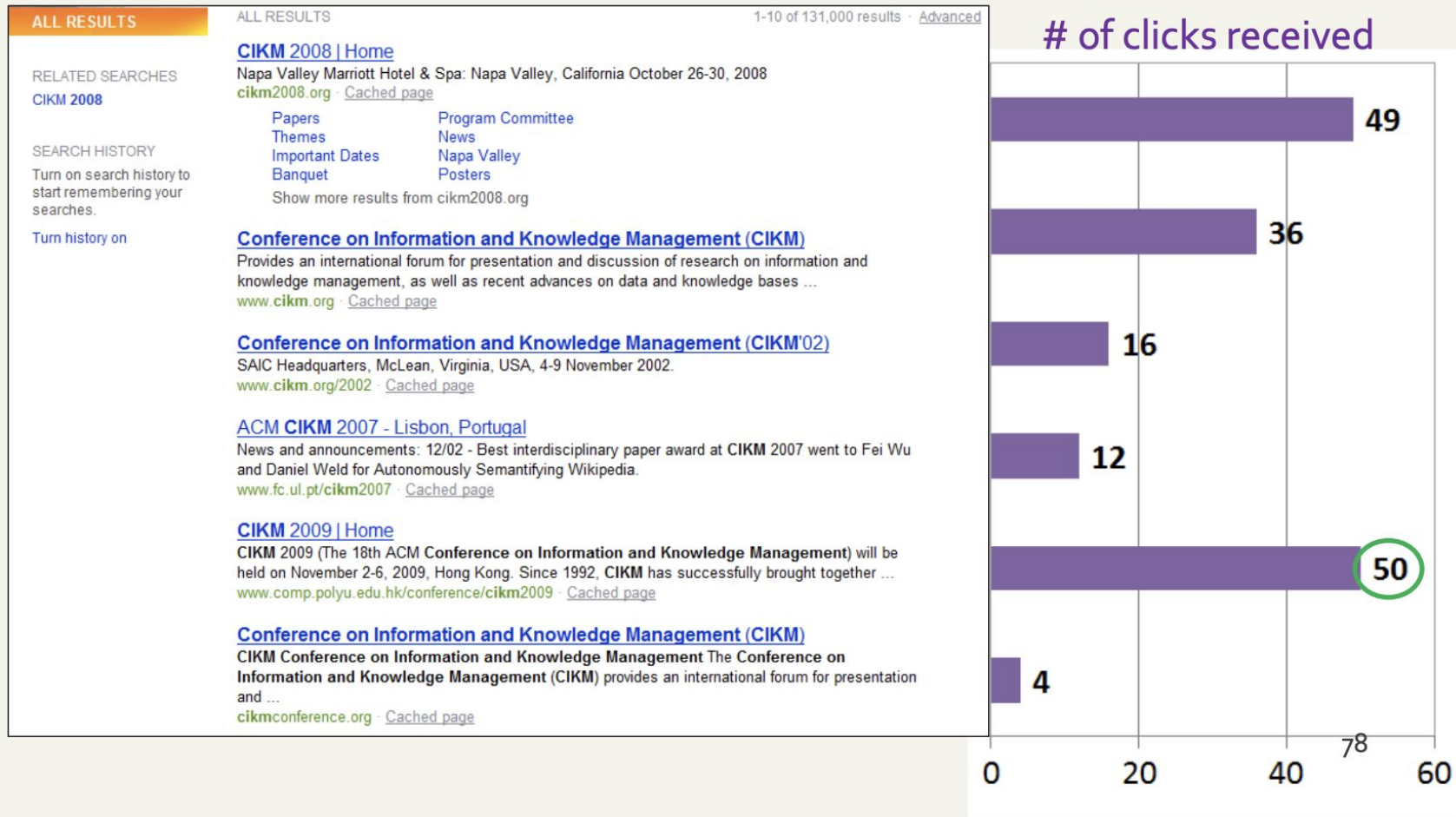
A/B testing

- Purpose: Test a single innovation
- Prerequisite: You have a large search engine up and running.
- Have most users use old system
- Divert a small proportion of traffic (e.g., 1%) to the new system that includes the innovation
- Evaluate with an “automatic” measure like clickthrough on first result
- Now we can directly see if the innovation does improve user happiness.
- Probably the evaluation methodology that large search engines trust most
- In principle less powerful than doing a multivariate regression analysis, but easier to understand

USING USER CLICKS

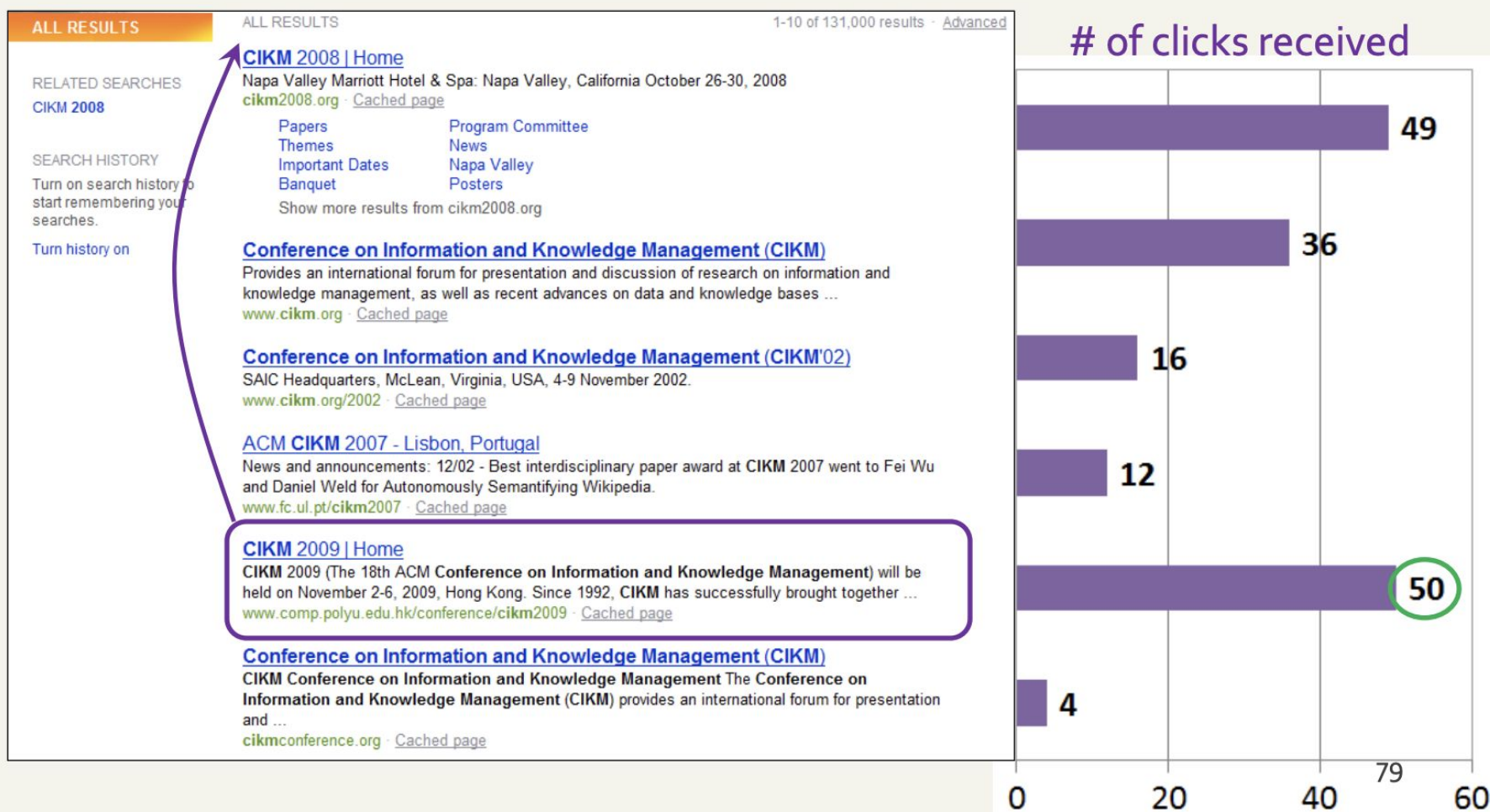


■ Search Results for "CIKM" (in 2009!)

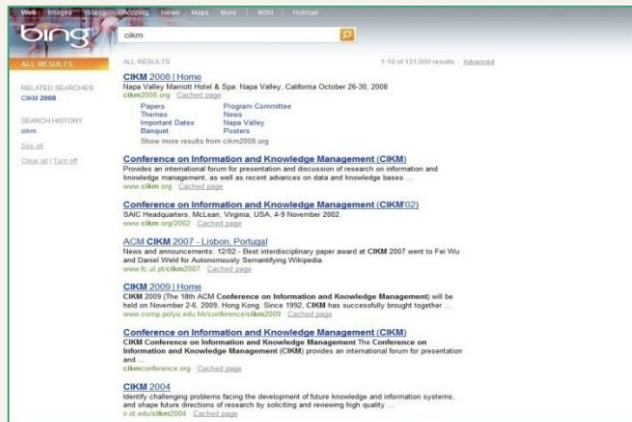




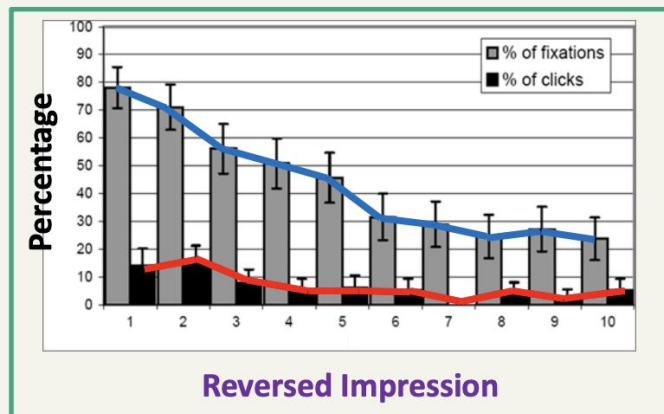
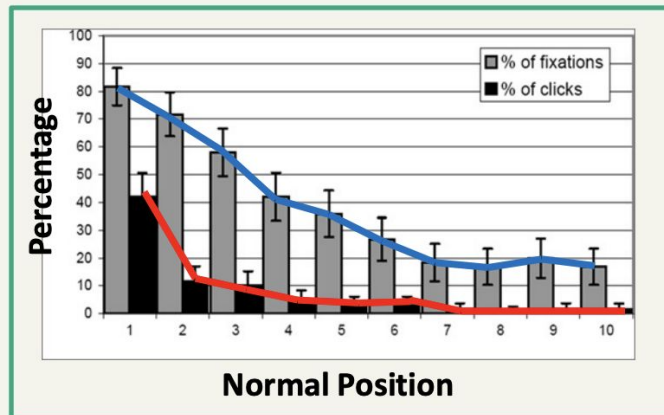
■ Adapt ranking to user clicks?



Eye Tracking User Study



Click Position-bias



- Higher positions receive more **user attention (eye fixation)** and **clicks** than lower positions.
- This is true even in the extreme setting where the order of positions is **reversed**.
- “Clicks are informative but biased”.

[Joachims+07]₈₂

Relative vs absolute ratings

ALL RESULTS

RELATED SEARCHES
 CIKM 2008

SEARCH HISTORY
 Turn on search history to start remembering your searches.
 Turn history on

ALL RESULTS
 1-10 of 131,000 results · [Advanced](#)

[CIKM 2008 | Home](#)
 Napa Valley Marriott Hotel & Spa: Napa Valley, California October 26-30, 2008
[cikm2008.org](#) · [Cached page](#)

Papers
 Themes
 Important Dates
 Banquet

Program Committee
 News
 Napa Valley
 Posters

Show more results from cikm2008.org

[Conference on Information and Knowledge Management \(CIKM\)](#)
 Provides an international forum for presentation and discussion of research on information and knowledge management, as well as recent advances on data and knowledge bases ...
[www.cikm.org](#) · [Cached page](#)


[Conference on Information and Knowledge Management \(CIKM'02\)](#)
 SAIC Headquarters, McLean, Virginia, USA, 4-9 November 2002.
[www.cikm.org/2002](#) · [Cached page](#)

[ACM CIKM 2007 - Lisbon, Portugal](#)
 News and announcements: 12/02 - Best interdisciplinary paper award at CIKM 2007 went to Fei Wu and Daniel Weld for Autonomously Semantifying Wikipedia.
[www.fc.ul.pt/cikm2007](#) · [Cached page](#)

[CIKM 2009 | Home](#)
 CIKM 2009 (The 18th ACM Conference on Information and Knowledge Management) will be held on November 2-6, 2009, Hong Kong. Since 1992, CIKM has successfully brought together ...
[www.comp.polyu.edu.hk/conference/cikm2009](#) · [Cached page](#)

[Conference on Information and Knowledge Management \(CIKM\)](#)
 CIKM Conference on Information and Knowledge Management The Conference on Information and Knowledge Management (CIKM) provides an international forum for presentation and ...
[cikmconference.org](#) · [Cached page](#)

User's click sequence



Hard to conclude Result1 > Result3
 Probably can conclude Result3 > Result2

Evaluating pairwise relative ratings

- Pairs of the form: Doc A better than Doc B for a query
- Doesn't mean that Doc A relevant to query
- Now, rather than assess a rank-ordering wrt per doc relevance assessments ...
- Assess in terms of conformance with historical pairwise preferences recorded from user clicks
- BUT!
 - Don't learn and test on the same ranking algorithm
 - I.e., if you learn historical clicks from nozama and compare Sergey vs nozama on this history ...



Facts/entities (what happens to clicks?)

+Prabhakar
Search
Images
Mail
Drive
Calendar
Sites
Groups
Contacts
More -

pragh@google.com 0 + Share

Web
Images
Maps
Shopping
News
More ▾
Search tools

About 1,300,000 results (0.39 seconds)

29,029' (8,848 m)
Mount Everest, Elevation

[Mount Everest - Wikipedia, the free encyclopedia](#)
https://en.wikipedia.org/wiki/Mount_Everest ▾

By the same measure of base to summit, **Mount McKinley**, in Alaska, is also taller than **Everest**. Despite its **height** above sea level of only 6,193.6 m (20,320 ft), ...

[List of deaths on eight - List of people who died ... - Timeline of climbing Mount](#)

[Facts About Mt. Everest - Scholastic](#)
teacher.scholastic.com/activities/hillary/archive/evfacts.htm ▾

Number of people to successfully climb **Mt. Everest**: 660. Number of

©2013 Google Map data ©2013 Google

Mount Everest

Mountain

Mount Everest is the Earth's highest mountain, with a peak at 8,848 metres above sea level and the 5th tallest mountain measured from the centre of the Earth. It is located in the Mahalangur section of the Himalayas.

Wikipedia

Elevation: 29,029' (8,848 m)
First ascent: May 29, 1953
Base camp: 29,029' (8,848 m)

References

1. Slides provided by Sougata Saha (Instructor, Fall 2022 - CSE 4/535)
2. Materials provided by Dr. Rohini K Srihari
3. <https://nlp.stanford.edu/IR-book/information-retrieval-book.html>