

Introduction to
Information Retrieval

CS276
Information Retrieval and Web Search
Pandu Nayak and Prabhakar Raghavan
Lecture 8: Evaluation

Introduction to Information Retrieval Sec. 6.2

This lecture

- How do we know if our results are any good?
 - Evaluating a search engine
 - Benchmarks
 - Precision and recall
- Results summaries:
 - Making our good results usable to a user

2

Introduction to Information Retrieval

EVALUATING SEARCH ENGINES

Introduction to Information Retrieval Sec. 8.6

Measures for a search engine

- How fast does it index
 - Number of documents/hour
 - (Average document size)
- How fast does it search
 - Latency as a function of index size
- Expressiveness of query language
 - Ability to express complex information needs
 - Speed on complex queries
- Uncluttered UI
- Is it free?

4

Introduction to Information Retrieval Sec. 8.6

Measures for a search engine

- All of the preceding criteria are *measurable*: we can quantify speed/size
 - we can make expressiveness precise
- The key measure: user happiness
 - What is this?
 - Speed of response/size of index are factors
 - But blindingly fast, useless answers won't make a user happy
- Need a way of quantifying user happiness

5

Introduction to Information Retrieval Sec. 8.6.2

Measuring user happiness

- Issue: who is the user we are trying to make happy?
 - Depends on the setting
- **Web engine**:
 - User finds what s/he wants and returns to the engine
 - Can measure rate of return users
 - User completes task – search as a means, not end
 - See Russell <http://dmrussell.googlepages.com/JCDL-talk-June-2007-short.pdf>
- **eCommerce site**: user finds what s/he wants and buys
 - Is it the end-user, or the eCommerce site, whose happiness we measure?
 - Measure time to purchase, or fraction of searchers who become buyers?

6

Introduction to Information Retrieval Sec. 8.6.2

Measuring user happiness

- **Enterprise** (company/govt/academic): Care about “user productivity”
 - How much time do my users save when looking for information?
 - Many other criteria having to do with breadth of access, secure access, etc.

7

Introduction to Information Retrieval Sec. 8.1

Happiness: elusive to measure

- Most common proxy: *relevance* of search results
- But how do you measure relevance?
- We will detail a methodology here, then examine its issues
- Relevance measurement requires 3 elements:
 1. A benchmark document collection
 2. A benchmark suite of queries
 3. A usually binary assessment of either **Relevant** or **Nonrelevant** for each query and each document
 - Some work on more-than-binary, but not the standard

8

Introduction to Information Retrieval Sec. 8.1

Evaluating an IR system

- Note: the **information need** is translated into a **query**
- Relevance is assessed relative to the **information need** *not* the **query**
- E.g., **Information need**: *I'm looking for information on whether drinking red wine is more effective at reducing your risk of heart attacks than white wine.*
- **Query**: **wine red white heart attack effective**
- Evaluate whether the doc addresses the information need, not whether it has these words

9

Introduction to Information Retrieval Sec. 8.2

Standard relevance benchmarks

- TREC - National Institute of Standards and Technology (NIST) has run a large IR test bed for many years
- Reuters and other benchmark doc collections used
- “Retrieval tasks” specified
 - sometimes as queries
- Human experts mark, for each query and for each doc, **Relevant** or **Nonrelevant**
 - or at least for subset of docs that some system returned for that query

10

Introduction to Information Retrieval Sec. 8.3

Unranked retrieval evaluation: Precision and Recall

- **Precision**: fraction of retrieved docs that are relevant = $P(\text{relevant} | \text{retrieved})$
- **Recall**: fraction of relevant docs that are retrieved = $P(\text{retrieved} | \text{relevant})$

	Relevant	Nonrelevant
Retrieved	tp	fp
Not Retrieved	fn	tn

- Precision $P = tp / (tp + fp)$
- Recall $R = tp / (tp + fn)$

11

Introduction to Information Retrieval Sec. 8.3

Should we instead use the accuracy measure for evaluation?


- Given a query, an engine classifies each doc as “Relevant” or “Nonrelevant”
- The **accuracy** of an engine: the fraction of these classifications that are correct
 - $(tp + tn) / (tp + fp + fn + tn)$
- **Accuracy** is a commonly used evaluation measure in machine learning classification work
- Why is this not a very useful evaluation measure in IR?

12

Introduction to Information Retrieval Sec. 8.3

Why not just use accuracy?

- How to build a 99.9999% accurate search engine on a low budget....



- People doing information retrieval *want to find something* and have a certain tolerance for junk.

13

Introduction to Information Retrieval Sec. 8.3

Precision/Recall

- You can get high recall (but low precision) by retrieving all docs for all queries!
- Recall is a non-decreasing function of the number of docs retrieved
- In a good system, precision decreases as either the number of docs retrieved or recall increases
 - This is not a theorem, but a result with strong empirical confirmation

14

Introduction to Information Retrieval Sec. 8.3

Difficulties in using precision/recall

- Should average over large document collection/query ensembles
- Need human relevance assessments
 - People aren't reliable assessors
- Assessments have to be binary
 - Nuanced assessments?
- Heavily skewed by collection/authorship
 - Results may not translate from one domain to another

15

Introduction to Information Retrieval Sec. 8.3

A combined measure: F

- Combined measure that assesses precision/recall tradeoff is **F measure** (weighted harmonic mean):

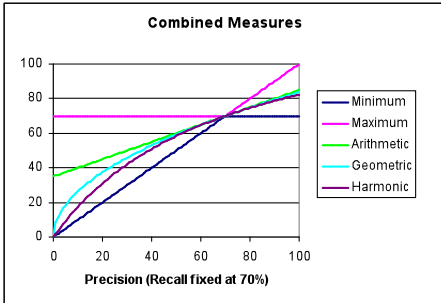
$$F = \frac{1}{\alpha \frac{1}{P} + (1-\alpha) \frac{1}{R}} = \frac{(\beta^2 + 1)PR}{\beta^2 P + R}$$

- People usually use balanced F_1 measure
 - i.e., with $\beta = 1$ or $\alpha = \frac{1}{2}$
- Harmonic mean is a conservative average
 - See CJ van Rijsbergen, *Information Retrieval*

16

Introduction to Information Retrieval Sec. 8.3

F_1 and other averages



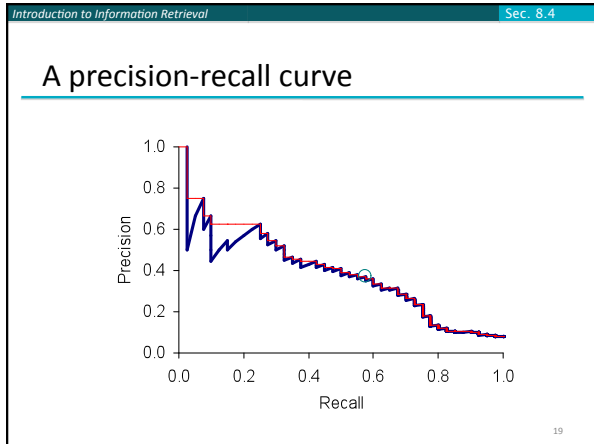
17

Introduction to Information Retrieval Sec. 8.4

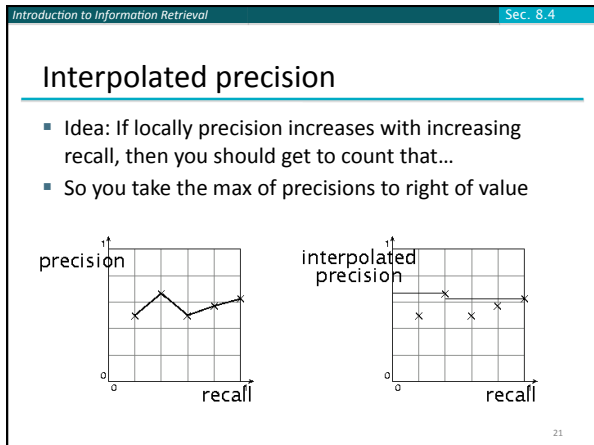
Evaluating ranked results

- Evaluation of ranked results:
 - The system can return any number of results
 - By taking various numbers of the top returned documents (levels of recall), the evaluator can produce a *precision-recall curve*

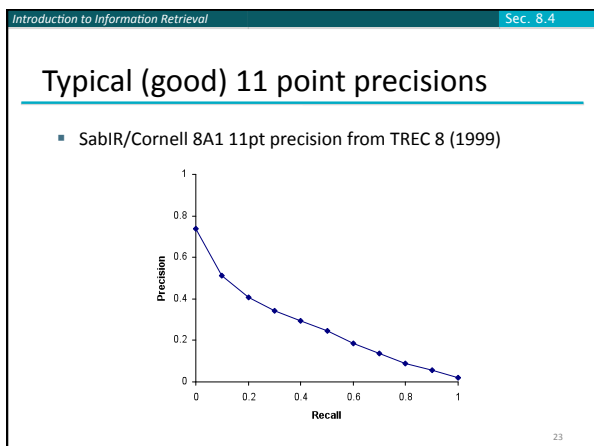
18



- Introduction to Information Retrieval Sec. 8.4
- ## Averaging over queries
- A precision-recall graph for one query isn't a very sensible thing to look at
 - You need to average performance over a whole bunch of queries.
 - But there's a technical issue:
 - Precision-recall calculations place some points on the graph
 - How do you determine a value (interpolate) between the points?
- 20



- Introduction to Information Retrieval Sec. 8.4
- ## Evaluation
- Graphs are good, but people want summary measures!
 - Precision-at- k : Precision of top k results
 - Perhaps appropriate for most of web search: all people want are good matches on the first one or two results pages
 - But: averages badly and has an arbitrary parameter of k
 - 11-point interpolated average precision
 - The standard measure in the early TREC competitions: you take the precision at 11 levels of recall varying from 0 to 1 by tenths of the documents, using interpolation (the value for 0 is always interpolated!), and average them
 - Evaluates performance at all recall levels
- 22



- Introduction to Information Retrieval Sec. 8.4
- ## Yet more evaluation measures...
- Mean average precision (MAP)
 - Average of the precision value obtained for the top k documents, each time a relevant doc is retrieved
 - Avoids interpolation, use of fixed recall levels
 - MAP for query collection is arithmetic ave.
 - Macro-averaging: each query counts equally
 - R-precision
 - If we have a known (though perhaps incomplete) set of relevant documents of size Rel , then calculate precision of the top Rel docs returned
 - Perfect system could score 1.0.
- 24

Introduction to Information Retrieval Sec. 8.4

Variance

- For a test collection, it is usual that a system does crummily on some information needs (e.g., MAP = 0.1) and excellently on others (e.g., 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!

25

Introduction to Information Retrieval

CREATING TEST COLLECTIONS FOR IR EVALUATION

Introduction to Information Retrieval Sec. 8.5

Test Collections

TABLE 4.3 Common Test Corpora

Collection	NDocs	NQrys	Size (MB)	Term/Doc	Q-D RelAss
ADI	82	35			
AIT	2109	14	2	400	>10,000
CACM	3204	64	2	24.5	
CISI	1460	112	2	46.5	
Cranfield	1400	225	2	53.1	
LISA	5872	35	3		
Medline	1033	30	1		
NFL	11,429	93	3		
OSHMED	34,8566	106	400	250	16,140
Reuters	21,578	672	28	131	
TREC	740,000	200	2000	89-3543	» 100,000

27

Introduction to Information Retrieval Sec. 8.5

From document collections to test collections

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

28

Introduction to Information Retrieval Sec. 8.5

Kappa measure for inter-judge (dis) agreement

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

29

Introduction to Information Retrieval Sec. 8.5

Kappa Measure: Example

Number of docs	P(A)? P(E)?	
	Judge 1	Judge 2
300	Relevant	Relevant
70	Nonrelevant	Nonrelevant
20	Relevant	Nonrelevant
10	Nonrelevant	Relevant

30

Introduction to Information Retrieval Sec. 8.5

Kappa Example

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

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

31

Introduction to Information Retrieval Sec. 8.2

TREC

- TREC Ad Hoc task from first 8 TRECs is standard IR task
 - 50 detailed information needs a year
 - Human evaluation of pooled results returned
 - More recently other related things: Web track, HARD
- A TREC query (TREC 5)


```
<top>
<num> Number: 225
<desc> Description:
What is the main function of the Federal Emergency Management
Agency (FEMA) and the funding level provided to meet emergencies?
Also, what resources are available to FEMA such as people,
equipment, facilities?
</top>
```

32

Introduction to Information Retrieval Sec. 8.2

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

33

Introduction to Information Retrieval Sec. 8.5

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.

34

Introduction to Information Retrieval Sec. 8.5.1

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
- See Carbonell reference

35

Introduction to Information Retrieval Sec. 8.6.3

Can we avoid human judgment?

- No
- 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)

36

Evaluation at large search engines

- Search engines have test collections of queries and hand-ranked 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 Cumulative Discounted Gain)
- Search engines also use non-relevance-based measures.
 - Clickthrough on first result
 - Not very reliable if you look at a single clickthrough ... but pretty reliable in the aggregate.
 - Studies of user behavior in the lab
 - A/B testing

37

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

38

RESULTS PRESENTATION

39

Result Summaries

- Having ranked the documents matching a query, we wish to present a results list
- Most commonly, a list of the document titles plus a short summary, aka “10 blue links”

John McCain

John McCain 2008 - The Official Website of John McCain's 2008 Campaign for President ... African American Coalition; Americans of Faith; American Indians for McCain; Americans with ...
www.johnmccain.com - Cached page

JohnMcCain.com - McCain-Palin 2008

John McCain 2008 - The Official Website of John McCain's 2008 Campaign for President ... African American Coalition; Americans of Faith; American Indians for McCain; Americans with ...
www.johnmccain.com/InformingIssues - Cached page

John McCain News- msnbc.com

Complete political coverage of John McCain. ... Republican leaders said Saturday that they were worried that Sen. John McCain was heading for defeat unless he brought stability to ...
www.msnbc.msn.com/id/1643500 - Cached page

John McCain | Facebook

Welcome to the official Facebook Page of John McCain. Get exclusive content and interact with John McCain right from Facebook. Join Facebook to create your own Page or to start ...
www.facebook.com/johnmccain - Cached page

40

Summaries

- The title is often automatically extracted from document metadata. What about the summaries?
 - This description is crucial.
 - User can identify good/relevant hits based on description.
- Two basic kinds:
 - Static
 - Dynamic
- A **static summary** of a document is always the same, regardless of the query that hit the doc
- A **dynamic summary** is a *query-dependent* attempt to explain why the document was retrieved for the query at hand

41

Static summaries

- In typical systems, the static summary is a subset of the document
- Simplest heuristic: the first 50 (or so – this can be varied) words of the document
 - Summary cached at indexing time
- More sophisticated: extract from each document a set of “key” sentences
 - Simple NLP heuristics to score each sentence
 - Summary is made up of top-scoring sentences.
- Most sophisticated: NLP used to synthesize a summary
 - Seldom used in IR; cf. text summarization work

42

Introduction to Information Retrieval Sec. 8.7

Dynamic summaries

- Present one or more “windows” within the document that contain several of the query terms
 - “KWIC” snippets: Keyword in Context presentation

The image shows two search engines, Google and Yahoo!, displaying search results for the query 'christopher manning'. Each result includes a snippet of text from the source document, with the search terms highlighted. For example, the Google result shows: 'Christopher Manning, Stanford NLP. Christopher Manning, Associate Professor of Computer Science and Linguistics, Stanford University. nlp.stanford.edu/~manning/ - 12k - Cached - Similar pages'.

Introduction to Information Retrieval Sec. 8.7

Techniques for dynamic summaries

- Find small windows in doc that contain query terms
 - Requires fast window lookup in a document cache
- Score each window wrt query
 - Use various features such as window width, position in document, etc.
 - Combine features through a scoring function – methodology to be covered Nov 12th
- Challenges in evaluation: judging summaries
 - Easier to do pairwise comparisons rather than binary relevance assessments

44

Introduction to Information Retrieval

Quicklinks

- For a *navigational query* such as **united airlines** user’s need likely satisfied on www.united.com
- Quicklinks provide navigational cues on that home page

The image shows a Google search for 'united airlines'. Below the search bar, there are several quicklinks provided by Google, such as 'United Airlines Flights', 'United Airlines - Airline Tickets, Airline Reservations, Flight...', and 'Airline tickets, airline reservations, flight airfare from United Airlines'. These links are designed to help users quickly find the information they need on the United Airlines website.

Introduction to Information Retrieval

The image shows a Yahoo! search for 'united airlines'. The search results page displays the official site for United Airlines, along with various quicklinks such as 'Cheap Flight Tickets', 'United Airline Schedule', and 'United Airlines Reservations'. The page also includes a 'Best match' section and 'RELATED SEARCHES'.

Introduction to Information Retrieval

Alternative results presentations?

The image shows a Yahoo! search for the query 'uni'. The search results are presented in a list format, with the entry for 'united airlines' highlighted. The highlighted entry includes a snippet of text and a link to 'www.united.com'. Other results include 'univision', 'university of phoenix', 'asian unicorn', 'universal studios', 'united states postal service', and 'united healthcare'.

Introduction to Information Retrieval

Resources for this lecture

- IIR 8
- MIR Chapter 3
- MG 4.5
- Carbonell and Goldstein 1998. The use of MMR, diversity-based reranking for reordering documents and producing summaries. SIGIR 21.

48