# Introduction to
# **Information Retrieval**

CS276
Information Retrieval and Web Search
Pandu Nayak and Prabhakar Raghavan

Evaluation
IIR Ch. 8 – Slides modified from Stanford CS276,
Spring 2015 (Manning and Nayak)
http://nlp.stanford.edu/IR-book/

---

## This lecture

- How do we know if our results are any good?
  - Evaluating a search engine
    - Benchmarks
    - Boolean retrieval metrics
    - Ranked retrieval metrics

- Results summaries:
  - Making our good results usable to a user

2

---

## **EVALUATING SEARCH ENGINES**

---

## 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

---

## 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

---

## 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

## Measuring user happiness

- <u>Enterprise</u> (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

## 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 <u>Relevant</u> or <u>Nonrelevant</u> for each query and each document
     - Some work on more-than-binary, but not the standard

8

## 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., <u>Information need</u>: *I'm looking for information on whether drinking red wine is more effective at reducing your risk of heart attacks than white wine.*
- <u>Query</u>: *wine red white heart attack effective*
- Evaluate whether the doc addresses the information need, not whether it has these words

9

## 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, <u>Relevant</u> or <u>Nonrelevant</u>
  - or at least for subset of docs that some system returned for that query

10

## Unranked retrieval evaluation: Precision and Recall

- **Precision**: fraction of retrieved docs that are relevant = P(relevant|retrieved)
- **Recall**: fraction of relevant docs that are retrieved = P(retrieved|relevant)

|  | Relevant | Nonrelevant |
|---|---|---|
| Retrieved | tp | fp |
| Not Retrieved | fn | tn |

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

11

## 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

## Why not just use accuracy?

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

snoogle.com

**Search for:**

*0 matching results found.*

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

13

## 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

## 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
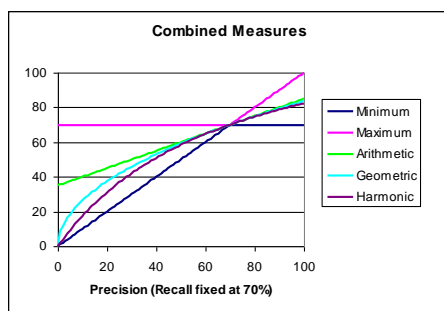
15

## 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

## $F_1$ and other averages

**Combined Measures**

- Minimum
- Maximum
- Arithmetic
- Geometric
- Harmonic

Precision (Recall fixed at 70%)

17

## Evaluating ranked results

- Up until now we've been considering metrics for boolean (set-based) retrieval
  - Precision, Recall, $F_1$

- But users don't really care about all results

- Users care about getting results near top of ranking…

18

3

## Metrics: Ranking

- Ranking results matters for human consumption of data

> Precision at k: P@10 does not distinguish between the two results

> Average Precision: prefers Result 2 to Result 1

1. **Precision @ k (P@k)**
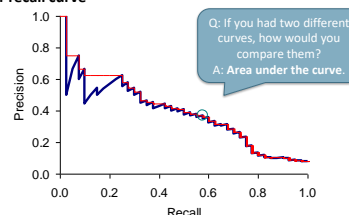   - Percent of relevant results (out of top k)

2. **Average Precision (AP or AveP)**
   - Weights higher ranks more
   - More on the exact definition shortly…

| Rank | Result 1 $P@k=0.5$ $AP=0.68$ | Result 2 $P@k=0.5$ $AP=1.00$ |
|------|------|------|
| 1 | ✓ | ✓ |
| 2 |   | ✓ |
| 3 | ✓ | ✓ |
| 4 |   | ✓ |
| 5 | ✓ | ✓ |
| 6 |   |   |
| 7 | ✓ |   |
| 8 |   |   |
| 9 | ✓ |   |
| 10 |   |   |

---

## Evaluating ranked results

- Sometimes we don't want to fix top "k"
  - The system can return any number of results
  - We can evaluate performance for a range of k by looking at the **precision-recall curve**



> Q: If you had two different curves, how would you compare them?
> A: **Area under the curve**.

One way to generate is to vary the length k of a (ranked) results list. 20

---

## Evaluation

- Graphs are good, but people want summary measures!
  - P@k good for most of web search… why? what k?

  - But P@k averages badly
    - If only 10 relevant docs, max P@100 is 0.1
    - Also has an arbitrary parameter of *k*
    - Sometimes **R-Prec** is better
      - R-Prec definition: P@k with k=#relevant docs (for query)
      - max R-Prec is 1.0, why?

  - *But P@k and R-Prec still use a fixed k. Does any ranking metric approximate area under precision-recall curve?*
    - Well yes, average precision does just that…

22

---

## Definition of (Mean) Average Precision

- **Average Precision (AveP or AP) and Mean AP (MAP)**

$$\text{MAP} = \frac{\sum_{q=1}^{Q} \text{AveP(q)}}{Q} \qquad \text{AveP} = \frac{\sum_{k=1}^{n} (P(k) \times \text{rel}(k))}{\text{number of relevant documents}}$$

- AP = higher ranked docs are counted more often
  - Unlike P@k, ordering matters!

- AP $\approx$ area under precision-recall curve when n$\rightarrow$#all docs!
  - Good discussion in IIR book and on Wikipedia
    https://en.wikipedia.org/wiki/Information_retrieval#Performance_and_correctness_measures

- Mean AP (MAP) = mean over queries
  - Note: this is **macro-averaging**: queries weighted equally
  - Empirically correlates with human evaluation of retrieval systems

23

---

## 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.

- **There are easy information needs and hard ones!**

28

---

## CREATING TEST COLLECTIONS FOR IR EVALUATION

## 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 | | |
| NPL | 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 |

30

---

## 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?

31

---

## 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.

32

---

P(A)? P(E)?

## Kappa Measure: Example

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

33

---

## Kappa Example

- P(A) = 370/400 = 0.925
- P(nonrelevant) = (10+20+70+70)/800 = 0.2125
- P(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 -> "tentative conclusions" (Carletta '96)
- Depends on purpose of study
- For >2 judges: average pairwise kappas

34

---

## 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>

35

## 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

36

## 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.

37

## 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

38

## 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)

39

## 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

40

## 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

41

**RESULTS PRESENTATION**

42

## 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/Informing/Issues · 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/16438320 · 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

43

## 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

44

## 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

45

## Dynamic summaries

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

Google  christopher manning    **Christopher Manning**, Stanford NLP
**Christopher Manning**, Associate Professor of Computer Science and Linguistics, Stanford University.
nlp.stanford.edu/~manning/ - 12k - Cached - Similar pages

Google  christopher manning machine translation    **Christopher Manning**, Stanford NLP
**Christopher Manning**, Associate Professor of Computer Science and Linguistics, ... computational semantics, **machine translation**, grammar induction, ...
nlp.stanford.edu/~manning/ - 12k - Cached - Similar pages

YAHOO!  christopher manning    **Christopher Manning**, Stanford NLP
**Christopher Manning**, Associate Professor of Computer Science and Linguistics, Stanford University ... **Chris Manning** works on systems and formalisms that can ...
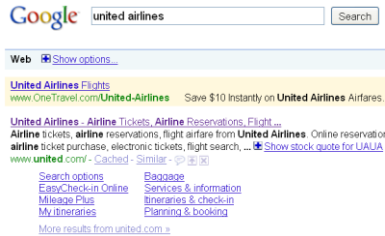nlp.stanford.edu/~manning - Cached

## 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 12[th]
- Challenges in evaluation: judging summaries
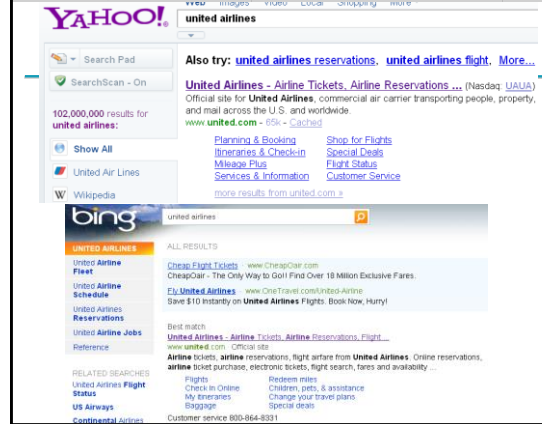  - Easier to do pairwise comparisons rather than binary relevance assessments

47

7

## Quicklinks

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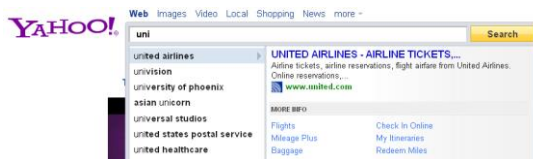- Quicklinks provide navigational cues on that home page



48

49

## Alternative results presentations?



50

## 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.

51