


SIT330-770: Natural Language Processing

Week 2 - Information Retrieval Part 2

Probabilistic IR: the binary independence model, BM25, BM25F
Evaluation methods & NDCG
Dr. Mohamed Reda Bouadjeneq

School of Information Technology, Faculty of
Sci Eng & Built Env

reda.bouadjeneq@deakin.edu.au



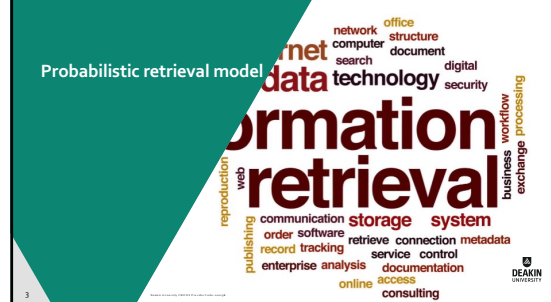
1

Outline

- Probabilistic IR
 - Introduction to the probability ranking principle
 - The Binary Independence Model: BIM
 - The BIM Ranking formula
 - The BM25 (Best Match 25) Model
- Evaluating search engines
 - Boolean Evaluating Metrics
 - Ranked evaluation metrics
 - Test collection for IR evaluation
 - Results presentation

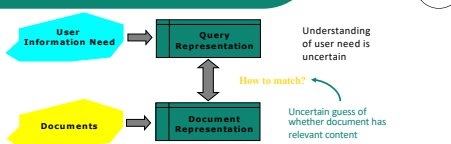
2

Probabilistic retrieval model



3

Why probabilities in IR?

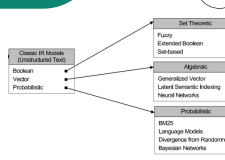


- In traditional IR systems, matching between each document and query is attempted in a semantically imprecise space of index terms
- Probabilities provide a principled foundation for uncertain reasoning.
- Can we use probabilities to quantify our search uncertainties?

4

Probabilistic IR topics

1. Classical probabilistic retrieval model
 - Probability ranking principle, etc.
 - Binary independence model (= Naive Bayes text cat)
 - (Okapi) BM25
2. Bayesian networks for text retrieval
3. Language model approach to IR
 - An important development



- Probabilistic methods are one of the oldest but also one of the currently hot topics in IR
 - Traditionally: neat ideas, but didn't win on performance
 - It seems to be different now

5

The document ranking problem

- We have a collection of documents
- User issues a query
- A list of documents needs to be returned
- Ranking method is the core of modern IR systems:
 - In what order do we present documents to the user?
 - We want the "best" document to be first, second best second, etc.
- Idea: Rank by probability of relevance of the document w.r.t. information need
 - $P(R=1|\text{document}, \text{query})$

6

7

8

9

10

11

12

Binary Independence Model

- Traditionally used in conjunction with PRP
- "Binary" = Boolean: documents are represented as binary incidence vectors of terms:
 - $d = (t_1, \dots, t_n)$
 - t_i term i is present in document d
- "Independence": terms occur in documents independently
- Different documents can be modeled as the same vector

13

Binary Independence Model

- Queries: binary term incidence vectors
- Given query q ,
 - for each document d need to compute $p(R|q, d)$
 - Interested only in ranking
- Will use odds and Bayes' Rule:

$$O(R|q, d) = \frac{p(R=1|q, d)}{p(R=0|q, d)} = \frac{p(R=1|q)p(d|R=1, q)}{p(R=0|q)p(d|R=0, q)}$$

14

Binary Independence Model

$$O(R|q, d) = \frac{p(R=1|q, d)}{p(R=0|q, d)} = \frac{p(R=1|q)}{p(R=0|q)} \times \frac{p(d|R=1, q)}{p(d|R=0, q)}$$

Using Independence Assumption:

$$\frac{p(d|R=1, q)}{p(d|R=0, q)} = \frac{\prod_{i=1}^n p(t_i|R=1, q)}{\prod_{i=1}^n p(t_i|R=0, q)}$$

$$O(R|q, d) = O(R|q) \times \prod_{i=1}^n \frac{p(t_i|R=1, q)}{p(t_i|R=0, q)}$$

Constant for a given query

Needs estimation

15

Binary Independence Model

$$O(R|q, d) = O(R|q) \times \prod_{i=1}^n \frac{p(t_i|R=1, q)}{p(t_i|R=0, q)}$$

- Since x_{ij} is either 0 or 1:

$$O(R|q, d) = O(R|q) \times \prod_{i=1}^n \frac{p(t_i|R=1, q)}{p(t_i|R=0, q)} \times \prod_{i=0}^n \frac{p(t_i|R=1, q)}{p(t_i|R=0, q)}$$

- Let $p_i = p(t_i|R=1, q)$ and $u_i = p(t_i|R=0, q)$
- Assume, for all terms not occurring in the query ($q_i = 0$) $p_i = u_i$:

$$O(R|q, d) = O(R|q) \times \prod_{i=1}^n \frac{p_i}{u_i} \times \prod_{q_i=0}^n \frac{(1-p_i)}{(1-u_i)}$$

16

Binary Independence Model

$$O(R|q, d) = O(R|q) \times \prod_{i=1}^n \frac{p_i}{u_i} \times \prod_{q_i=0}^n \frac{(1-p_i)}{(1-u_i)}$$

All matching terms

All query terms

Non-matching query terms

17

Binary Independence Model

$$O(R|q, d) = O(R|q) \times \prod_{i=1}^n \frac{p_i(1-u_i)}{u_i(1-p_i)} \times \prod_{q_i=0}^n \frac{(1-p_i)}{(1-u_i)}$$

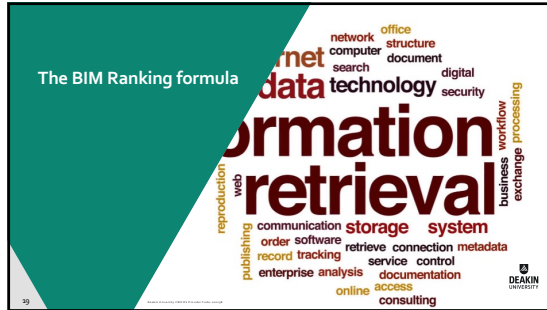
Constant for each query

Retrieval Status Value:

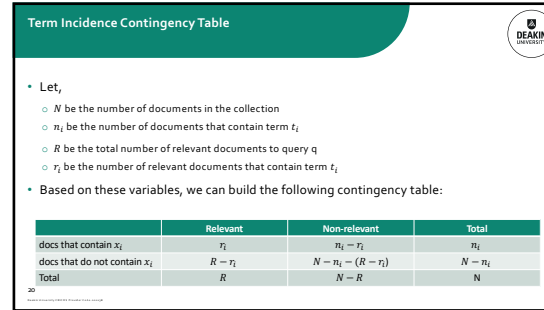
Only quantity to be estimated for rankings

$$RSV = \log \prod_{i=1}^n \frac{p_i(1-u_i)}{u_i(1-p_i)} = \sum_{i=1}^n \log \frac{p_i(1-u_i)}{u_i(1-p_i)} \approx Sim(q, d)$$

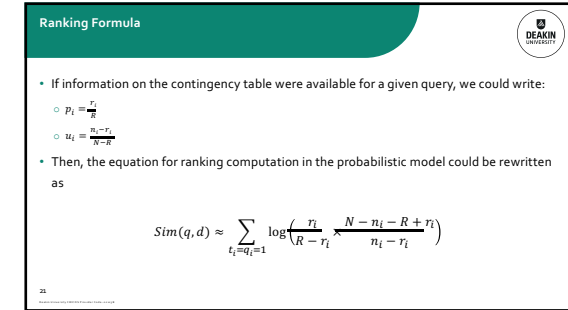
18



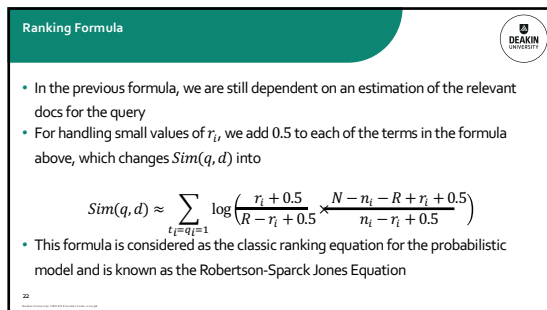
19



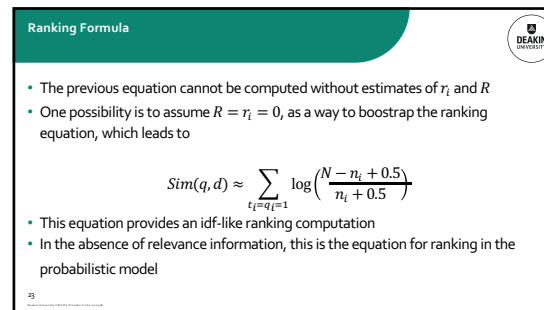
20



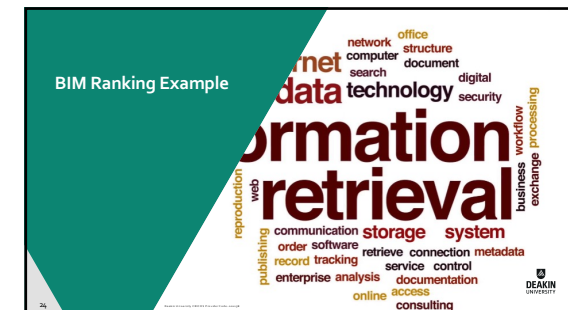
21



22



23



24

Ranking Example

- Document ranks computed by the previous probabilistic ranking equation for the query **"to do"**

d₁: To do is to be. To be is to do.

d₂: To be or not to be. I am what I am.

d₃: I think therefore I am. Do be do be do.

d₄: Do do do, da da da. Let it be, let it be.

doc	rank computation	rank
<i>d₁</i>	$\log \frac{4-2+0.5}{2+0.5} + \log \frac{4-3+0.5}{3+0.5}$	- 1.222
<i>d₂</i>	$\log \frac{4-2+0.5}{2+0.5}$	0
<i>d₃</i>	$\log \frac{4-3+0.5}{3+0.5}$	- 1.222
<i>d₄</i>	$\log \frac{4-3+0.5}{3+0.5}$	- 1.222

25

Ranking Example

- The ranking computation led to negative weights because of the term **"do"**
- Actually, the probabilistic ranking equation produces negative terms whenever $n_i > N/2$
- One possible artifact to contain the effect of negative weights is to change the previous equation to:

$$Sim(q, d) \approx \sum_{t_i=q_i=1} \log \left(\frac{N+0.5}{n_i+0.5} \right)$$

- By doing so, a term that occurs in all documents ($n_i = N$) produces a weight equal to zero

26

Ranking Example

- Using this latest formulation, we redo the ranking computation for our example collection for the query **"to do"** and obtain

d₁: To do is to be. To be is to do.

d₂: To be or not to be. I am what I am.

d₃: I think therefore I am. Do be do be do.

d₄: Do do do, da da da. Let it be, let it be.

doc	rank computation	rank
<i>d₁</i>	$\log \frac{4+0.5}{2+0.5} + \log \frac{4+0.5}{3+0.5}$	1.210
<i>d₂</i>	$\log \frac{4+0.5}{2+0.5}$	0.847
<i>d₃</i>	$\log \frac{4+0.5}{3+0.5}$	0.362
<i>d₄</i>	$\log \frac{4+0.5}{3+0.5}$	0.362

27



28

Estimating r_i and R

- Our examples above considered that $r_i = R = 0$
- An alternative is to estimate r_i and R performing an initial search:
 - select the top 10-20 ranked documents
 - inspect them to gather new estimates for r_i and R
 - remove the 10-20 documents used from the collection
 - rerun the query with the estimates obtained for r_i and R
- Unfortunately, procedures such as these require human intervention to initially select the relevant documents

29

Improving the Initial Ranking

- Consider the equation

$$Sim(q, d) \approx \sum_{t_i=q_i=1} \log \frac{p_i(1-u_i)}{u_i(1-p_i)}$$

- How obtain the probabilities p_i and u_i ?
- Estimates based on assumptions:
 - $p_i = 0.5$
 - $u_i = \frac{n_i}{N}$ where n_i is the number of docs that contain t_i
 - Use this initial guess to retrieve an initial ranking
 - Improve upon this initial ranking

30

Improving the Initial Ranking

- Substituting p_i and u_i into the previous Equation, we obtain:

$$Sim(q, d) \approx \sum_{t_i=q, i=1} \log\left(\frac{N - n_i}{n_i}\right)$$
- That is the equation used when no relevance information is provided, without the 0.5 correction factor
- Given this initial guess, we can provide an initial probabilistic ranking
- After that, we can attempt to improve this initial ranking as follows

31

Improving the Initial Ranking

- We can attempt to improve this initial ranking as follows
- Let
 - D_i : set of docs initially retrieved
 - D_i : subset of docs retrieved that contain t_i
- Reevaluate estimates:
 - $p_i = \frac{n_i}{D}$
 - $u_i = \frac{n_i - D_i}{N - D_i}$
- This process can then be repeated recursively
- To avoid $D = 0$ and $D_i = 0$:
 - $p_i = \frac{D_i + 0.5}{D + 1}$
 - $u_i = \frac{n_i - D_i + 0.5}{N - D + 1}$
- Also,
 - $p_i = \frac{D_i + \frac{1}{2}}{D + 1}$
 - $u_i = \frac{n_i - D_i + \frac{1}{2}}{N - D + 1}$

32

Pluses and Minuses

- Advantages:
 - Docs ranked in decreasing order of probability of relevance
- Disadvantages:
 - need to guess initial estimates for p_i
 - method does not take into account tf factors
 - the lack of document length normalization

33

Comparison of Classic Models

- Boolean model does not provide for partial matches and is considered to be the weakest classic model
- There is some controversy as to whether the probabilistic model outperforms the vector model
- Croft suggested that the probabilistic model provides a better retrieval performance
- However, Salton *et al* showed that the vector model outperforms it with general collections
- This also seems to be the dominant thought among researchers and practitioners of IR

34

The BM (Best Match) Models

35

BM25 (Best Match 25)

- BM25 was created as the result of a series of experiments on variations of the probabilistic model
- A good term weighting is based on three principles
 - inverse document frequency
 - term frequency
 - document length normalization
- The classic probabilistic model covers only the first of these principles
- This reasoning led to a series of experiments with the Okapi system, which led to the BM25 ranking formula

36

BM1, BM11 and BM15 Formulas

- At first, the Okapi system used the Equation below as ranking formula

$$Sim(q, d) \approx \sum_{t_i=q_i=1} \log \left(\frac{N - n_i + 0.5}{n_i + 0.5} \right)$$

which is the equation used in the probabilistic model, when no relevance information is provided

- It was referred to as the BM1 formula (**Best Match 1**)

37

37

BM1, BM11 and BM15 Formulas

- The first idea for improving the ranking was to introduce a term-frequency factor $\mathcal{F}_{t,d}$ in the BM1 formula
- This factor, after some changes, evolved to become

$$\mathcal{F}_{t,d} = S_1 \times \frac{tf_{t,d}}{K_1 + tf_{t,d}}$$

Where

- $tf_{t,d}$ is the frequency of term t within document d
- K_1 is a constant setup experimentally for each collection
- S_1 is a scaling constant, normally set to $S_1 = (K_1 + 1)$
- If $K_1 = 0$ this whole factor becomes equal to 1 and bears no effect in the ranking

38

38

BM1, BM11 and BM15 Formulas

- The next step was to modify the $\mathcal{F}_{t,d}$ factor by adding document length normalization to it, as follows:

$$\hat{\mathcal{F}}_{t,d} = S_1 \times \frac{tf_{t,d}}{\frac{K_1 \times len(d)}{avg_doclen} + tf_{t,d}}$$

Where

- $len(d)$ is the length of document d (computed, for instance, as the number of terms in the document)
- avg_doclen is the average document length for the collection

39

39

BM1, BM11 and BM15 Formulas

- Next, a correction factor \mathcal{G}_q dependent on the document and query lengths was added

$$\mathcal{G}_q = K_2 \times len(q) \times \frac{avg_doclen - len(d)}{ave_doclen + len(d)}$$

Where

- $len(q)$ is the query length (number of terms in the query)
- K_2 is a constant

40

40

BM1, BM11 and BM15 Formulas

- A third additional factor, aimed at taking into account term frequencies within queries, was defined as

$$\mathcal{F}_{t,q} = S_3 \times \frac{tf_{t,q}}{K_3 + tf_{t,q}}$$

Where

- $tf_{t,q}$ is the frequency of term t within query q
- K_3 is a constant
- S_3 is a scaling constant related to K_3 , normally set to $S_3 = (K_3 + 1)$

41

41

BM1, BM11 and BM15 Formulas

- Introduction of these three factors led to various BM (Best Matching) formulas, as follows:

$$Sim_{BM1}(q, d) \approx \sum_{t_i=q_i=1} \log \left(\frac{N - n_i + 0.5}{n_i + 0.5} \right)$$

$$Sim_{BM15}(q, d) \approx \mathcal{G}_q + \sum_{t_i=q_i=1} \mathcal{F}_{t,d} \times \mathcal{F}_{t,q} \times \log \left(\frac{N - n_i + 0.5}{n_i + 0.5} \right)$$

$$Sim_{BM11}(q, d) \approx \mathcal{G}_q + \sum_{t_i=q_i=1} \hat{\mathcal{F}}_{t,d} \times \mathcal{F}_{t,q} \times \log \left(\frac{N - n_i + 0.5}{n_i + 0.5} \right)$$

42

42

BM1, BM11 and BM15 Formulas

- Experiments using TREC data have shown that BM11 outperforms BM15
- Further, empirical considerations can be used to simplify the previous equations, as follows:
 - Empirical evidence suggests that a best value of K_2 is 0, which eliminates the g_q factor from these equations
 - Further, good estimates for the scaling constants S_1 and S_3 are $K_1 + 1$ and $K_3 + 1$, respectively
 - Empirical evidence also suggests that making K_3 very large is better. As a result, the $F_{t,q}$ factor is reduced simply to $f_{t,d}$
 - For short queries, we can assume that $f_{t,d}$ is 1 for all terms

43

BM1, BM11 and BM15 Formulas

- These considerations lead to simpler equations as follows:

$$Sim_{BM1}(q, d) \approx \sum_{t_i=q_i=1} \log \left(\frac{N - n_i + 0.5}{n_i + 0.5} \right)$$

$$Sim_{BM15}(q, d) \approx g_q + \sum_{t_i=q_i=1} F_{t,d} \times F_{t,q} \times \log \left(\frac{N - n_i + 0.5}{n_i + 0.5} \right)$$

$$Sim_{BM11}(q, d) \approx g_q + \sum_{t_i=q_i=1} F_{t,d} \times F_{t,q} \times \log \left(\frac{N - n_i + 0.5}{n_i + 0.5} \right)$$

44

The BM25 Model

45

BM25 Ranking Formula

- BM25: combination of the BM11 and BM15
- The motivation was to combine the BM11 and BM25 term frequency factors as follows:

$$B_{t,d} = \frac{(K_1 + 1) f_{t,q}}{K_1 \left[(1 - b) + b \frac{\text{len}(d)}{\text{avg_doclen}} \right] + f_{t,q}}$$

Where b is a constant with values in the interval $[0, 1]$

- If $b = 0$, it reduces to the BM15 term frequency factor
- If $b = 1$, it reduces to the BM11 term frequency factor
- For values of b between 0 and 1, the equation provides a combination of BM11 with BM15

46

BM25 Ranking Formula

- The ranking equation for the BM25 model can then be written as:

$$Sim_{BM25}(q, d) \approx \sum_{t_i=q_i=1} B_{t,d} \times \log \left(\frac{N - n_i + 0.5}{n_i + 0.5} \right)$$

where K_1 and b are empirical constants

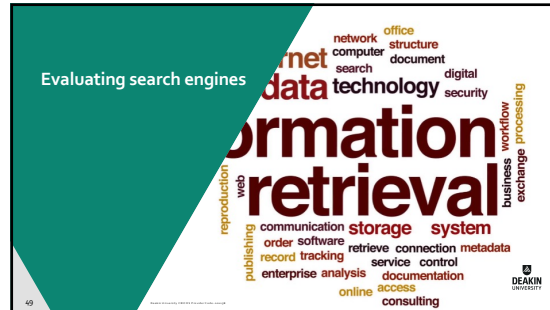
- $K_1 = 1$ works well with real collections
- b should be kept closer to 1 to emphasize the document length normalization effect present in the BM11 formula
- For instance, $b = 0.75$ is a reasonable assumption
- Constants values can be fine tuned for particular collections through proper experimentation

47

BM25 Ranking Formula

- Unlike the probabilistic model, the BM25 formula can be computed without relevance information
- There is consensus that BM25 outperforms the classic vector model for general collections
- Thus, it has been used as a baseline for evaluating new ranking functions, in substitution to the classic vector model

48



49

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?

50

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

51

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?

52

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.

53

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

54

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

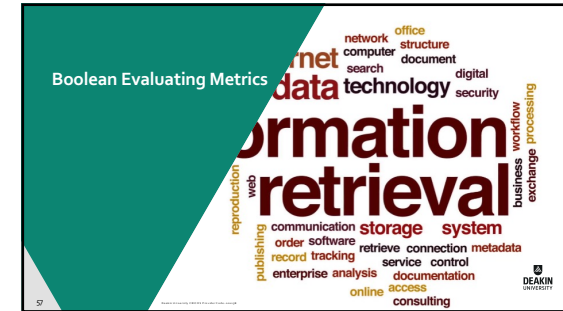
55

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

56

Boolean Evaluating Metrics



57

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 $R = tp / (tp + fn)$

58

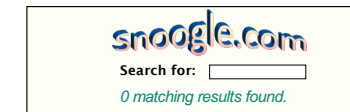
Should we instead use the accuracy measure for evaluation?

- 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?

59

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.

60

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

61

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

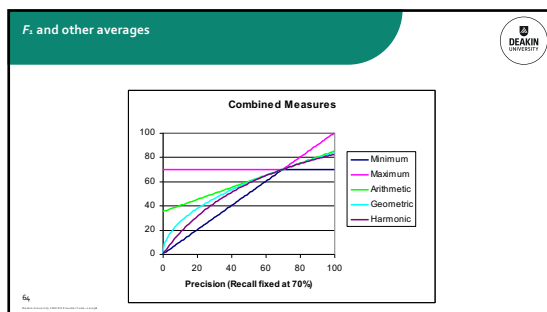
62

A combined measure: F

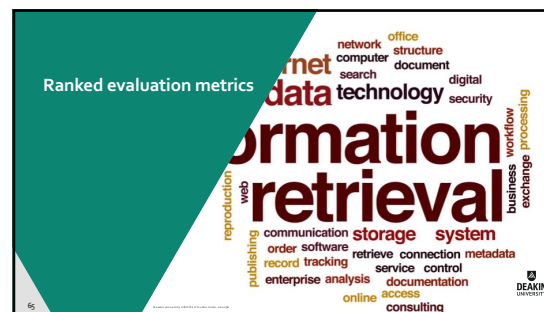
- Combined measure that assesses precision/recall tradeoff is **F measure** (weighted harmonic mean):

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

63



64



65

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...

66

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

Rank	Result 1	Result 2
1	✓	✓
2	✓	✓
3	✓	✓
4	✓	✓
5	✓	✓
6	✓	✓
7	✓	✓
8	✓	✓
9	✓	✓
10	✓	✓

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...

67

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**

68

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 20 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...

69

Definition of (Mean) Average Precision

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

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

- AP = higher ranked docs are counted more often
 - Unlike P@k, ordering matters!
- AP = area under precision-recall curve when n → #all docs!
 - Good discussion in IIR book and on Wikipedia: https://en.wikipedia.org/wiki/Average_precision
- Mean AP (MAP) = mean over queries
 - Note: this is **macro-averaging**: queries weighted equally
 - Empirically correlates with human evaluation of retrieval systems

70

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!

71

Test collection for IR evaluation

72

Test Collections

TABLE 4.3 Common Test Corpora

Collection	<i>N</i> Docs	<i>N</i> Qry	Size (MB)	Terms/Doc	Q-D Pairs
ADI	82	35			
ATT	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	5972	35	3		
Medline	1033	30	1		
NPL	11,429	93	3		
OSHMEM	34,8566	106	400	250	16,140
Reuters	21,578	672	28	131	
TREC	740,000	200	2000	89-3543	~ 100,000

73

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?

74

TREC

- TREC Ad Hoc task from first 8 TRECs is standard IR task
 - go 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>
```

75

Qrels example

```
225 q EP-1003976 2
225 q EP-0946545 2
225 q EP-0952483 2
226 q EP-0810279 2
226 q EP-1254844 1
227 q EP-0999933 2
227 q EP-0855293 1
228 q EP-0600930 1
228 q EP-0032015 2
228 q EP-1348868 2
229 q EP-0152364 2
229 q EP-1249379 1
229 q EP-0554488 1
229 q EP-0946539 2
230 q EP-0933578 2
230 q EP-1008229 2
```

76

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)

77

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

78

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

79

79

Results presentation

80

80

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"

81

81

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

82

82

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

83

83

Dynamic summaries

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

84

84

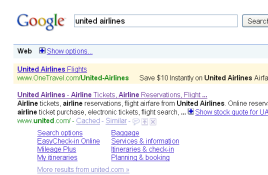
Techniques for dynamic summaries

- 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

85

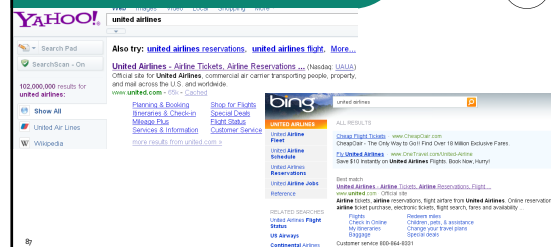
Quicklinks

- 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



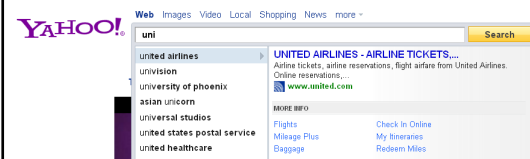
86

Quicklinks



87

Alternative results presentations?



88