Information Retrieval WS 2016 / 2017

Lecture 2, Tuesday October 25th, 2016 (Ranking, Evaluation)

> Prof. Dr. Hannah Bast Chair of Algorithms and Data Structures Department of Computer Science University of Freiburg

REIBURG

Overview of this lecture

- Organizational
 - Your experiences with ES1
 - Requirements for ES2
- Contents
 - Ranking tf.idf and BM25
 - Evaluation Ground truth, Precision, ...

Inverted index

Test Cases

Exercise Sheet #2: implement BM25, tune your ranking, and then evaluate on a small benchmark
 There will be a small competition (table on the Wiki)

UNI FREIBURG

Experiences with ES1 1/3

Summary / excerpts

- Nice and interesting exercise ... most time spent on SVN, Checkstyle, etc. or due to lack of programming practice
 Don't worry, this will get much better soon
- Some problems understanding Exercise 2

In case of doubts, always ask on the Forum!

- Some battles fought with Jenkins / Checkstyle
 Don't worry, you will get used to it quickly
- Problems with the encoding on some platforms
 One solution: open(file, encoding=utf8)

Experiences with ES1 2/3

Results

– Queries that work well:

harry pottersufficiently specific keywordscolumbia pictures 2011sufficiently specific keywords

- Queries that don't work well

men in blackwords frequent in other docsthe kings speechapostrophe in doc (king's)spidermantwo words in doc (spider man)

Experiences with ES1 3/3

- Linear construction time?
 - Quite a few of you implemented something like this:
 for line in file:

if record_id not in self.inverted_lists[word]:
 self.inverted_lists[word].append(record_id)

 Then index construction on movies.txt takes very long: the "not in" takes linear time, not constant time which means the whole loop take quadratic time
 Super-important piece of advice: never use built-in functions without understanding their time-complexity

Requirements for ES2

Test Cases

- For ES1 you had to write the test cases yourself
- From now on, we will provide test cases, at least for the central functions of each exercise sheet
- This should save **you** some work, but also **our tutors**

Code that does not pass a basic unit test can be very hard and cumbersome to correct (similar to code that doesn't compile or has checkstyle errors)

– In return, you **have to** implement these unit tests

You can also extend them, but not restrict them

 As before, everything must run through on Jenkins without errors, otherwise no correction and no points

Ranking 1/14

Motivation

- Queries often return many hits
- Typically more than one wants to (or even can) look at
 For web search: often millions of documents
 - But even for less hits a proper ranking is **key** to usability

REN

- So we want to have the most "relevant" hits first
- Problem: how to measure what is how "relevant"

Ranking 2/14

Basic Idea

- In the inverted lists, for each doc id also have a score university 17 0.5 , 53 0.2 , 97 0.3 , 127 0.8 freiburg 23 0.1 , 34 0.8 , 53 0.1 , 127 0.7
- While merging, **aggregate** the scores, then **sort** by score
 - MERGED 17 0.5, 23 0.1, 34 0.8, 53 0.3, 97 0.3, 127 1.5 SORTED 127 1.5, 34 0.8, 17 0.5, 53 0.3, 97 0.3, 23 0.1
- The entries in the list are referred to as **postings**

Above, it's only doc id and score, but a posting can also contain more information, e.g. the position of a word

Ranking 3/14

Getting the top-k results

- A full sort of the result list takes time $\Theta(n \cdot \log n)$, where n is the number of postings in the list ZW

- Typically only the top-k hits need to be displayed
- Then a **partial sort** is sufficient: get the k largest elements, for a given k

Can be computed in time $\Theta(n + k \cdot \log k)$

k rounds of HeapSort yield time $\Theta(n + k \cdot \log n)$

For constant k these are both **O(n)**

For ES2, you can ignore this issue

Ranking 4/14

Meaningful scores

- How do we precompute good scores
 university
 17 0.5 , 53 0.2 , 97 0.3 , 127 0.8
 freiburg
 23 0.1 , 34 0.8 , 53 0.1 , 127 0.7
- Goal: the score for the posting for doc D_i in the inverted list for word w should reflect the relevance of w in D_i

In particular, the larger the score, the more relevant

Problem: relevance is somewhat subjective

But it has to be done somehow anyway !

Ranking 5/14

Term frequency (tf)

- The number of times a word occurs in a document
- Problem: some words are frequent in many documents, regardless of the content

REI

university	, 57 5 , , 123 2 ,
of	, 57 14 , , 123 23 ,
freiburg	, 57 3 ,, 123 1 ,
SCORE SUM	, 57 22 , , 123 26 ,

A word like "of" should not count much for relevance Some of you observed that already while trying out queries for ES1

Ranking 6/14

UNI FREIBURG

Document frequency (df)

- The number of documents containing a particular word $df_{university} = 16.384$, $df_{of} = 524.288$, $df_{freiburg} = 1.024$ For simplicity, number are powers of 2, see below why
- Inverse document frequency (idf)

 $idf = log_2 (N / df)$ N = total number of documents

For the example df scores above and N = $1.048.576 = 2^{20}$

 $idf_{university} = 6$, $idf_{of} = 1$, $idf_{freiburg} = 10$

Understand: without the \log_2 , small differences in df would have too much of an effect ; why exactly $\log_2 \rightarrow$ later slide

Ranking 7/14

djuminiersity = 6 dj_{og} = 1 dj_{freelung} = 10

Combining the two (tf.idf)

Reconsider our earlier tf only example

 university
 ..., 57
 5, ..., 123
 2, ...

 of
 ..., 57
 14, ..., 123
 23, ...

 freiburg
 ..., 57
 3, ..., 123
 1, ...

 SCORE SUM
 ..., 57
 22, ..., 123
 26, ...

Now combined with idf scores from previous slide

university	, 57	30 , , 123 12 ,
of	, 57	14 , , 123 23 ,
freiburg	, 57	30 , , 123 10 ,
SCORE SUM	, 57	74 , , 123 45 ,

Ranking 8/14

Problems with tf.idf in practice

– The idf part is fine, but the tf part has several problems

- Let w be a word, and D_1 and D_2 be two documents
- **Problem 1:** assume that D_1 is longer than D_2

Then tf for w in D_1 tends to be larger then tf for w in D_2 , because D_1 is longer, not because it's more "about" w

 Problem 2: assume that D₁ and D₂ have the same length, and that the tf of w in D₁ is twice the tf of w in D₂

Then it is reasonable to assume that D_1 is more "about" w than D_2 , but just a little more, and not twice more

UNI FREIBURG

h = O

The BM25 (best match) formula

 This tf.idf style formula has consistently outperformed other formulas in standard benchmarks over the years

BM25 score = $tf^* \cdot \log_2 (N / df)$, where

 $\mathbf{tf^*} = \mathbf{tf} \cdot (\mathbf{k} + 1) / (\mathbf{k} \cdot (1 - \mathbf{b} + \mathbf{b} \cdot \mathbf{DL} / \mathbf{AVDL}) + \mathbf{tf})$

tf = term frequency, DL = document length, AVDL = average document length

- Standard setting for **BM25**: k = 1.75 and $b = 0.75 \implies a = 1.4$

Ranking 10/14

Plausibility argument for BM25, part 1

- Start with the simple formula $tf \cdot idf$
- Replace tf by tf* such that the following properties hold:
 - tf* = 0 if and only if tf = 0
 - tf* increases as tf increases
 - $tf^* \rightarrow fixed limit as tf \rightarrow \infty$

 $B = 0 = 0 B^{*} = 0 B$ $B^{*} = \frac{2+1}{1+2}B B$ $B^{*} = \frac{2+1}{1+2}B - 2 + 1$ $B^{*} = \frac{2+1}{1+2}B - 2 + 1$

UNI FREIBURG

- The "simplest" formula with these properties is • $tf^* = tf \cdot (k + 1) / (k + tf) = \frac{\mathcal{B} \cdot (\mathcal{B} + 1)}{\mathcal{B} + \mathcal{B}} = \frac{\mathcal{B} + 1}{\mathcal{B} / (\mathcal{A} + 1)}$

Ranking 11/14

Plausibility argument for BM25, part 2

- So far, we have $tf^* = tf \cdot (k + 1) / (k + tf)$

- Normalize by the length of the document
 - Replace tf by tf / α
 - Full normalization: $\alpha = DL / AVDL \dots$ too extreme

REI

b = 0 = 0 = 2 = 1 $ab = 1 = 0 = 2 = \frac{DL}{AVDI}$

- Some normalization: $\alpha = (1 b) + b \cdot DL / AVDL$
- This gives us tf* = tf / $\alpha \cdot (k + 1) / (k + tf / \alpha)$

- And hence $tf^* = tf \cdot (k + 1) / (k \cdot \alpha + tf)$

Lots of "theory" behind this formula, but to me not really more convincing than these simple plausibility arguments

Ranking 12/14

Implementation advice

– First compute the inverted lists with **tf** scores

You already did that (implicitly or explicitly) for ES1

 Along with that compute the document length (DL) for each document, and the average document length (AVDL) ZW

You can measure DL (and AVDL) via the number of words

 Make a second pass over the inverted lists and replace the tf scores by tf* • idf scores

 $tf \cdot (k + 1) / (k \cdot (1 - b + b \cdot DL / AVDL) + tf) \cdot \log_2 (N / df)$

Note that the **df** of a word is just the length (number of postings) in its inverted list

Ranking 13/14

Further refinements

- For ES2, play around with the BM25 parameters ${\bf k}$ and ${\bf b}$
- Boost results that match each query word at least once

Warning: when you **restrict** your results to such matches, you might miss some relevant results

For example: steven spielberg **movies**

- Somehow take the popularity of a movie into account
 In the file on the Wiki, movies are sorted by popularity
 Popularity scores also have a Zipf distribution, so you might take ~ N^{-α} as popularity score for the N-th movie in the list
- Anything else that comes to your mind and might help ...

Ranking 14/14

Advanced methods

 There is a multitude of other sources / approaches for improving the quality of the ranking, for example: REI

Using query logs and click-through data

Who searches what and clicked on what ... main pillar for the result quality of big search engines like Google

Learning to rank

Using machine learning (more about that in a later lecture) to find the best parameter setting

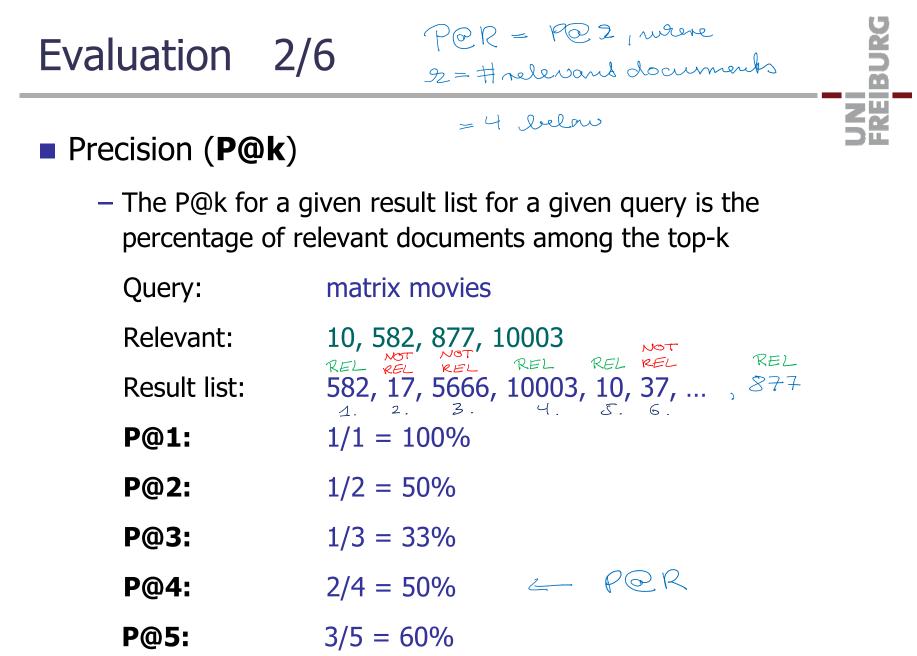
Evaluation 1/6

UNI FREIBUR

Ground truth

- For a given query, the ids of all documents relevant for that query
 - Query: matrix movies
 - Relevant: 10, 582, 877, 10003
- For ES2, we have built a ground truth for 10 queries

Building a good and large enough ground truth is a common (and time-consuming) part in research in IR



Evaluation 3/6

Average Precision (AP)

- Let R_1 , ..., R_k be the sorted list of positions of the relevant document in the result list of a given query

REI

– Then AP is the average of the k P@Ri values

Query:	matrix movies
Relevant:	10, 582, 877, 10003 REL REL REL REL REL , REL
Result list:	582, 17, 5666, 10003, 10,, 877
R ₁ ,, R ₄ :	1, 4, 5, 40
P@R ₁ ,, P@R ₄ :	100%, 50%, 60%, 10%
AP:	55%

Note: for docs not in result list, just take $P@R_i = 0$

Evaluation 4/6

UNI FREIBURG

Mean Precisions (MP@k, MP@R, MAP)

- Given a benchmark with several queries + ground truth
- Then one can capture the quality of a system by taking the mean (average) of a given measure over all queries
 MP@k = mean of the P@k values over all queries
 MP@R = mean of the P@R values over all queries
 MAP = mean of the AP values over all queries
 These are very common measures, which you will find in a lot of research papers on information retrieval

Evaluation 5/6

Other measures

- There is a BIG variety of other evaluation measures, e.g.

– **nDCG** = normalized discounted cumulative growth

Takes into account that documents can have varying degrees of relevance, e.g. "primary" and "secondary" Gives credit if "primary" is ranked before "secondary"

erearen prinary is rankea berere secondar

BPref = binary relevance ... preference relation
 Takes into accounts that some documents are unjudged
 This is a frequent problem in benchmarks for huge text corpora, where complete judgment is impossible

E.g. all relevant documents for "tom hanks" on the web

Evaluation 6/6

UNI FREIBURG

Overfitting

 Tuning parameters / methods to achieve good results on a given benchmark is called **overfitting**

In an extreme case: for each query from the benchmark, output the list of relevant docs from the ground truth

In a realistic environment (real search engine or competition), one is given a **training** set for development

The actual evaluation is on a **test** set, which must not be used / was not available during development

For ES2, do the development / tuning on some queries of your choice, then evaluate without further changes

References

UNI FREIBURG

In the Manning/Raghavan/Schütze textbook

Section 6: Scoring and Term Weighting

Section 8: Evaluation in Information Retrieval

Relevant Papers

Probabilistic Relevance: BM25 and Beyond FnTIR 2009

Test Collection Based Evaluation of IR Systems FnTIR 2010

Relevant Wikipedia articles

http://en.wikipedia.org/wiki/Okapi BM25

https://en.wikipedia.org/wiki/Information_retrieval #Performance_and_correctness_measures