

School of Computing and Information Systems
The University of Melbourne
COMP90049

Knowledge Technologies (Semester 1, 2019)
Workshop sample solutions: Week 7

1. What are the four primary components of a **Web-scale Information Retrieval engine**? Briefly describe our goal in each of them.

- **Crawling**: finding and downloading as many documents as we can from the web (hopefully all of them, although this isn't possible in practice).
- **Parsing**: turning each document into a list of tokens (or terms), probably by removing page metadata, case folding, stemming, etc.
- **Indexing**: building an inverted index out of all of the tokens in our downloaded document collection. (We stop worrying about the original documents at this point.)
- **Querying**: after the previous three steps have been completed (off-line), we are ready to accept user queries (on-line), where we aim to produce a document ranking. The queries take the form of keywords, that we tokenise (in a similar manner to our document collection) and then apply our querying model (e.g. TF-IDF) based on the information in the inverted index.
- (Optionally) **Additional things**: change the above ranking, based on ad-hoc application of certain factors, for example, PageRank, HITS, click-through data, zones, anchor text, etc.

2. Recall the (hypothetical) method of **crawling** given in the lectures:

(a) Would this method be *effective* at solving the problem of crawling? Why or why not?

- Somewhat, although it depends on having a large, random set of seeds, which isn't really possible in practice. The method will miss large numbers of pages that aren't linked to from pages in the "core" of the World Wide Web, as well as all sorts of rich data encoded in databases, etc.

(b) Would this method be *efficient* at solving the problem of crawling? Why or why not?

- An efficient approach depends on having a good model of **web page duplication**, so that we can decide (quickly) whether a given page has already been crawled. For example, having a hash function with respect to the page's contents (although consider how large the set is, and how difficult it would be to avoid collisions!).

3. When we **tokenise** text, we often **canonicalise** it. What are these generally accepted as referring to?

(Note the terminology is not used consistently in the literature; for example "tokenisation" occasionally refers to all three ideas.)

- Tokenisation, here, means decomposing the larger document into smaller information bearing units ("tokens") than we can compare against (the keywords present in) our query.
- Canonicalisation, here, means having a single representation of (each token within) a text, which we can use to sensibly compare one text with another (in our case, a document on the Web, and a query).

(a) What are some issues that arise when canonicalising text written in English?

- There are various point suggested in the lectures, including:
 - Stopwords
 - Stemming
 - Date/number formatting
 - Dialect variation
 - Spelling errors

Web Search | Information Retrieval

Four components {
crawling: gather data
parsing: translate data into canonical form
indexing: need smart algorithm to narrow the subset of query
querying: efficiently.

add-on technique: {
(HP2) snippet generation: same doc generate different snippet depends on what queries \Rightarrow different conclusion.
As-you-type query
Query correction
Answer consolidation: (合并)
info boxes

Crawling: Basic challenge: there is no central index of URLs of interest.

- ① Prioritised list L of URL to visit, List V that have been visited.
- ② Repeat
1) choose a URL u from L , fetch the page $p(u)$ at location u .
2) parse and index $p(u)$, extract URLs $\{u'\}$ from $p(u)$
3) add u to V , remove from L . Add $\{u'\}$ to L .
4) Process V or 'old' URLs to L .

Page processing is much faster than URL resolution

Parsing: first step determining the format of the page.

{ tokenisation
stemming
zoning
character encoding { ASCII \leftarrow
UTF-8, 8-32 bits, first 128 of 8 bits

c) avoid indexing invisible content \Rightarrow misleads users and allows spoofing (white text in white background.)

(2) scraping: ignore the ads in blog website

Tokenisation: If tokenisation is successful, tokens will match without approximate matching.

{
Hyphenation: 连字符 'under-coating' \rightarrow one or two
compounding: 'football' \rightarrow one or two
Possessives: 所有格 'Smith's'

- Note that English is **easy** compared to many languages!

4. Assume that we have crawled the following “documents”:

- Parse each document into terms.

- Construct an inverted index over the documents, for (at least) the terms `and`, `australia`, `celebrity`, `commission`, `island`, `on`, `the`, `to`, `tweet`, `twitter`

- | | | | | | | | | | |
|------------|---|------|---|------|---|------|---|------|---------------|
| and | → | 4(1) | → | 6(1) | | | | | |
| australia | → | 1(1) | | | | | | | |
| celebrity | → | 1(1) | → | 5(1) | → | 7(1) | | | |
| commission | → | 1(1) | → | 4(1) | | | | | |
| island | → | 1(1) | → | 3(1) | → | 7(1) | | | |
| on | → | 1(1) | → | 7(1) | | | | | |
| the | → | 1(1) | → | 2(2) | → | 3(1) | → | 4(1) | → 5(2) → 7(1) |
| to | → | 1(2) | → | 3(1) | → | 5(2) | → | 6(1) | → 7(1) |
| tweet | → | 3(1) | → | 5(1) | → | 7(1) | | | |
| twitter | → | 1(1) | → | 2(1) | | | | | |

- (a) Using the weighting functions $w_{d,t} = f_{d,t}$ and $w_{q,t} = \frac{N}{f_t}$

- 2

Canonicalisation : \Rightarrow Canonical form.

Canonical

1) Dates: 5/4/2011, 4/5/2011.

2) Numbers: 18 million, 18,000,000.

3) Variant spelling: Color colour.

4) Variant usage: Dr. Doctor

5) Variant punctuation: e.g. eg

most significant form is stemming for English.

words from root and stem

\Rightarrow stemming is the process of stripping away affixes.

in + expense + ive \rightarrow inexpensive

glass.es \rightarrow glass \Rightarrow stop when we arrive at a dictionary entry.

Indexing: web documents can be segmented into discrete zones like title, anchor text, headings. \Rightarrow calculate weights of zones \Rightarrow highlight the content, rank higher. and compute similarities for documents.

Indexing: inverted index ★

the \rightarrow 1(2) \rightarrow 2(1) \rightarrow 4(1)

\uparrow \uparrow \uparrow

doc1 Frequency doc2 doc4

we \rightarrow 2(1) \rightarrow 3(2)

Querying { boolean query: TDM (term-document matrix)

ranked query: TF-IDF, cosine. $S(q,d) = \frac{q \cdot d}{|q| |d|}$

Weighting functions: { (1) $W_{d,t} = f_{d,t}$ $W_{a,t} = \frac{N}{f_t}$ term frequency

(2) $W_{d,t} = 1 + \log_2 f_{d,t}$ $W_{a,t} = \log_2 (1 + \frac{N}{f_t})$

documents in the collection, and 2 of them contain **commission**, and so on:

$$\begin{aligned}
 w_{q, \text{commission}} &= \frac{N}{f_{\text{commission}}} = \frac{7}{2} = 3.500 \\
 w_{q, \text{island}} &= \frac{7}{3} \approx 2.333 \\
 w_{q, \text{on}} &= \frac{7}{2} = 3.500 \\
 w_{q, \text{to}} &= \frac{7}{5} = 1.400 \\
 w_{q, \text{twitter}} &= \frac{7}{2} = 3.500
 \end{aligned}$$

Handwritten notes:

$$w_{q, \text{commission}} = \langle 3.5, 2.333, 3.5, 1.4, 3.5 \rangle$$

$$w_{d, 1} = \langle 1, 1, 1, 2, 1 \rangle$$

$$\cos(d_1, q) = \frac{\langle 1, 1, 1, 2, 1 \rangle \cdot \langle 3.5, 2.333, 3.5, 1.4, 3.5 \rangle}{\sqrt{11} \sqrt{44.15}} = \frac{15.633}{\sqrt{11} \sqrt{44.15}} \approx 0.709$$

- All of the other **term weights** in the **query** are 0, so we can **ignore them** (and the weight documents are unimportant for the dot-product). The **document weights** are just the frequency of the term in each document (which we have recorded in our inverted index above). For example, for **island** our representation would look like:

$$\begin{aligned}
 w_{d_1, i} &= 1 \\
 w_{d_2, i} &= 0 \\
 w_{d_3, i} &= 1 \\
 w_{d_4, i} &= 0 \\
 w_{d_5, i} &= 0 \\
 w_{d_6, i} &= 0 \\
 w_{d_7, i} &= 1
 \end{aligned}$$

- The cosine model gives us:

$$\begin{aligned}
 \cos(d_1, q) &= \frac{\langle 1, 1, 1, 2, 1 \rangle \cdot \langle 3.5, 2.333, 3.5, 1.4, 3.5 \rangle}{\sqrt{\langle 0, 1, 1, 1, 1, 1, 1, 2, 0, 1 \rangle} \cdot \sqrt{\langle 3.5, 2.333, 3.5, 1.4, 3.5 \rangle}} \\
 &\approx \frac{15.633}{\sqrt{11} \sqrt{44.15}} \approx 0.709 \\
 \cos(d_2, q) &= \frac{\langle 0, 0, 0, 0, 1 \rangle \cdot \langle 3.5, 2.333, 3.5, 1.4, 3.5 \rangle}{\sqrt{\langle 0, 0, 0, 0, 0, 2, 0, 0, 1 \rangle} \cdot \sqrt{\langle 3.5, 2.333, 3.5, 1.4, 3.5 \rangle}} \\
 &\approx \frac{3.5}{\sqrt{5} \sqrt{44.15}} \approx 0.236 \\
 \cos(d_3, q) &= \frac{\langle 0, 1, 0, 1, 0 \rangle \cdot \langle 3.5, 2.333, 3.5, 1.4, 3.5 \rangle}{\sqrt{\langle 0, 0, 0, 0, 1, 0, 1, 1, 1, 0 \rangle} \cdot \sqrt{\langle 3.5, 2.333, 3.5, 1.4, 3.5 \rangle}} \\
 &\approx \frac{3.733}{\sqrt{4} \sqrt{44.15}} \approx 0.281 \\
 \cos(d_4, q) &= \frac{\langle 1, 0, 0, 0, 0 \rangle \cdot \langle 3.5, 2.333, 3.5, 1.4, 3.5 \rangle}{\sqrt{\langle 1, 0, 0, 1, 0, 0, 1, 0, 0, 0 \rangle} \cdot \sqrt{\langle 3.5, 2.333, 3.5, 1.4, 3.5 \rangle}} \\
 &\approx \frac{3.5}{\sqrt{3} \sqrt{44.15}} \approx 0.304
 \end{aligned}$$

Query cost: { Disk space
Memory space:
CPU time
Disk traffic.

Phrase query: { 1) bag-of-words, terms can occur anywhere in matching documents, then post-process to eliminate false matches,
2) add word positions to the index entries
3) use phrase index or word-pair index.

Link analysis: { HITS: hyperlinked-induced topic search
PageRank:

Note
1) Approximations can be used to reduce querying costs
2) Link and anchor information can be dominant evidence.

$$\begin{aligned}
\cos(d_5, q) &= \frac{\langle 0, 0, 0, 2, 0 \rangle \cdot \langle 3.5, 2.333, 3.5, 1.4, 3.5 \rangle}{|\langle 0, 0, 1, 0, 0, 0, 2, 2, 1, 0 \rangle| \cdot |\langle 3.5, 2.333, 3.5, 1.4, 3.5 \rangle|} \\
&\approx \frac{2.8}{\sqrt{10}\sqrt{44.15}} \approx 0.133 \\
\cos(d_6, q) &= \frac{\langle 0, 0, 0, 1, 0 \rangle \cdot \langle 3.5, 2.333, 3.5, 1.4, 3.5 \rangle}{|\langle 1, 0, 0, 0, 0, 0, 0, 1, 0, 0 \rangle| \cdot |\langle 3.5, 2.333, 3.5, 1.4, 3.5 \rangle|} \\
&\approx \frac{1.4}{\sqrt{2}\sqrt{44.15}} \approx 0.149 \\
\cos(d_7, q) &= \frac{\langle 0, 1, 1, 1, 0 \rangle \cdot \langle 3.5, 2.333, 3.5, 1.4, 3.5 \rangle}{|\langle 0, 0, 1, 0, 1, 1, 1, 1, 1, 0 \rangle| \cdot |\langle 3.5, 2.333, 3.5, 1.4, 3.5 \rangle|} \\
&\approx \frac{7.233}{\sqrt{6}\sqrt{44.15}} \approx 0.444
\end{aligned}$$

- All of the documents have non-zero similarity, so they all get returned, and the order is: {1,7,4,3,2,6,5}.
 - Notice also that each cosine calculation involves the same division by the length of the query (in this case, $\sqrt{44.15}$) since all we care about is the order of the documents and not the actual score, we can safely leave this term out (but it isn't the cosine anymore, so I'll call it \cos' below).
- (b) Using the weighting functions $w_{d,t} = 1 + \log_2 f_{d,t}$ and $w_{q,t} = \log_2(1 + \frac{N}{f_t})$
- Let's observe first that nothing has changed with the document weights: frequencies of 0 still have a weight of 0; frequencies of 1 are now $1 + \log_2(1) = 1 + 0 = 1$; frequencies of 2 are now $1 + \log_2(2) = 1 + 1 = 2$. (Higher frequency terms would have their weights reduced with this model.) This also means that the document lengths are the same.
 - We can now calculate the query weights according to the new model:
q: $\langle 2.170, 1.737, 2.170, 1.263, 2.170 \rangle$, and the cosine calculations look similar:

$$\begin{aligned}
\cos'(d_1, q) &= \frac{\langle 1, 1, 1, 2, 1 \rangle \cdot \langle 2.170, 1.737, 2.170, 1.263, 2.170 \rangle}{|\langle 0, 1, 1, 1, 1, 1, 2, 0, 1 \rangle|} \quad \text{[< 2.17, 1.737, 2.17, 1.263, 2.17 >]} \\
&\approx \frac{10.773}{\sqrt{11}} \approx 3.248 \quad \text{be ignore} \\
\cos'(d_2, q) &= \frac{\langle 0, 0, 0, 0, 1 \rangle \cdot \langle 2.170, 1.737, 2.170, 1.263, 2.170 \rangle}{|\langle 0, 0, 0, 0, 0, 0, 2, 0, 0, 1 \rangle|} \\
&\approx \frac{2.170}{\sqrt{5}} \approx 0.970 \\
\cos'(d_3, q) &= \frac{\langle 0, 1, 0, 1, 0 \rangle \cdot \langle 2.170, 1.737, 2.170, 1.263, 2.170 \rangle}{|\langle 0, 0, 0, 0, 1, 0, 1, 1, 1, 0 \rangle|} \\
&\approx \frac{3.000}{\sqrt{4}} \approx 1.500 \\
\cos'(d_4, q) &= \frac{\langle 1, 0, 0, 0, 0 \rangle \cdot \langle 2.170, 1.737, 2.170, 1.263, 2.170 \rangle}{|\langle 1, 0, 0, 1, 0, 0, 1, 0, 0, 0 \rangle|} \\
&\approx \frac{2.170}{\sqrt{3}} \approx 1.253 \\
\cos'(d_5, q) &= \frac{\langle 0, 0, 0, 2, 0 \rangle \cdot \langle 2.170, 1.737, 2.170, 1.263, 2.170 \rangle}{|\langle 0, 0, 1, 0, 0, 0, 2, 2, 1, 0 \rangle|} \\
&\approx \frac{2.526}{\sqrt{10}} \approx 0.799
\end{aligned}$$

$$\begin{aligned}
\cos'(d_6, q) &= \frac{\langle 0, 0, 0, 1, 0 \rangle \cdot \langle 2.170, 1.737, 2.170, 1.263, 2.170 \rangle}{| \langle 1, 0, 0, 0, 0, 0, 1, 0, 0 \rangle |} \\
&\approx \frac{1.263}{\sqrt{2}} \approx 0.893 \\
\cos'(d_7, q) &= \frac{\langle 0, 1, 1, 1, 0 \rangle \cdot \langle 2.170, 1.737, 2.170, 1.263, 2.170 \rangle}{| \langle 0, 0, 1, 0, 1, 1, 1, 1, 0 \rangle |} \\
&\approx \frac{5.170}{\sqrt{6}} \approx 2.111
\end{aligned}$$

- The document ranking has changed a little bit this time (because the relative differences between common terms and rare terms is smaller due to the log): {1,7,3,4,2,6,5}
- Also note that some of the values are greater than 1, because I left out the division by the length of the query (\cos'), which gives a factor of $\sqrt{18.74}$ in the denominator it doesn't change the order though.

5. When querying, what is an **accumulator**? What is the main problem, if we wish to use them? What heuristics can we use to solve this problem?

- We use accumulators to keep track of the partial sum of the dot product (for the cosine similarity) as we process query terms one-by-one. This makes sense, because our inverted index keeps track of the documents where a given term appears: consequently, we can take a term from the query, read off each document from the index, and update (add) the TF-IDF value for that document to the corresponding accumulator. Once we have processed all of the query terms, we divide each accumulator by its corresponding document length, and then take the greatest values as the top results.
- The list of accumulators is typically an array (we need random-access when writing to make it worth our while in terms of time). We need a single accumulator for each document unfortunately, there are lots of documents. A collection with 10M documents, using 32-bit floats for each accumulator value, would require 80MB of memory: not a problem, but we're going to have plenty of cache misses. For a more realistic number of documents on the Web, we have no hope of fitting this data structure into memory.
- To solve this problem, we generally use only a few accumulators (perhaps a few hundred or a few thousand, depending on how many results we wish to return), and aggressively prune the possible result space note that most documents have a trivially low similarity anyway! There are two main pruning strategies that we use:
 - Limiting: Only create an accumulator for the first few documents; these documents are then the results set, and the further calculations are just for re-ordering. This depends on processing the most important terms first, namely, the ones with highest IDF: we wish to have accumulators for the documents that contain the important terms; if we run out of accumulators for documents which contain the less important terms (but not the important ones), we accept that they probably weren't going to be returned at the top of the ranking anyway.
 - Thresholding: Only create an accumulator when the TF-IDF value (for the query term we're processing) is above some threshold. Similar to the above, but documents with low-frequency rare terms might not get added to the result set; whereas documents with high-frequency moderately rare terms might get added instead.

