

Department of Computing and Information Systems
The University of Melbourne
COMP30018/COMP90049
Knowledge Technologies (Semester 2, 2016)
Workshop exercises: Week 4

1. What is the difference between “data retrieval” and “information retrieval”? Why is the latter a knowledge task, but the former is not?
2. [EXTENSION] How many books are there in an average library? How many words are there in an average library? How many documents are there on the World Wide Web? How many words?
3. Identify some different types of “informational needs.” Often they are given in 3 categories — which categories might subsume other categories?
 - (a) Give examples of queries which might indicate a particular type of informational need. Are some of them ambiguous for the type of need?
 - (b) [EXTENSION] Input some queries of different types of informational needs into a web search engine like Google. Are search engines better at responding to some types of informational need than others? Is there any indication that the search engine is identifying the type of need and tailoring the results toward that?
4. Identify some differences between Boolean querying and ranked querying.
 - (a) [EXTENSION] Do search engines like Google use Boolean querying or ranked querying? How can we tell?
 - (b) [EXTENSION] With respect to the contents of a single web page, issues queries of successively greater length until only a single document is returned. Why does this occur? Pay attention to the time taken to resolve each query: what happens as you change the number of terms? Why?

DocID	apple	ibm	lemon	sun
Doc ₁	4	0	0	1
Doc ₂	5	0	5	0
Doc ₃	2	5	0	0
Doc ₄	1	2	1	7
Doc ₅	1	1	3	0

calculate the document ranking for the (conjunctive) queries: (a) apple and (b) apple lemon, based on the following TF-IDF term weighting model:

$$w_{d,t} = \begin{cases} 1 + \log_2 f_{d,t} & \text{if } f_{d,t} > 0 \\ 0 & \text{otherwise} \end{cases}$$
$$w_{q,t} = \begin{cases} \log(1 + \frac{N}{f_t}) & \text{if } f_{q,t} > 0 \\ 0 & \text{otherwise} \end{cases}$$

7. Revise “language model” approaches toward querying; in which ways is it similar to the “vector space model”? In which ways is it different?

Data Retrieval
• well-formatted
↳ database

Info Retrieval
• semi-formatted
↳ documents
• user → problem (Info Need)
↑ relevance
diff user judge diff

Info Need
user → problem
incomplete list, and overlapping
- Navigational
- Informational
- Geospatial
- Factoid
- Topic Tracking
- Transactional

Boolean
query
↓ this AND that OR NOT
result T/F

Ranked
query - this that
↓
result ranked by estimated most relevant
↓
est. - - least relevant

doc length
$$\text{sim}(q, d) = \frac{\sum_{t \in q} \text{TF-IDF}_{d,t}}{\sqrt{\sum_{t \in q} \text{IDF}_t^2}}$$

related to IDF

TF, apple: $(1 + \log_2 4) \cdot \log(1 + \frac{5}{8}) = 3$

IDF, lemon: $0 \times \log_2(1 + \frac{5}{3}) = 0$

$$\text{sim}(q, d_1) = \frac{3+0}{\sqrt{8^2 + 0^2 + 0^2 + [1 \cdot \log_2(1 + \frac{5}{8})]^2}}$$