

School of Computing and Information Systems
The University of Melbourne
COMP90049

Knowledge Technologies (Semester 1, 2017)

Workshop exercises: Week 4

Following on from last week, suppose that we have observed the token **lended**, and we have a dictionary as follows:

```
addendum
blenders
commodity
deaden
end
leader
leant
lent
lemonade
pleading
```

1. Finish last week's problems, where necessary.
2. Assuming that the "correct" (intended) dictionary entry was **lent**, calculate the precision of the following methods of finding approximate entries from the dictionary.
 - For each method, we will consider how many dictionary entries it returns as a result (predicts as a good match), as well as how many it got correct — in this case, there is only a single possible correct answer, so the value will be 0 or 1.
 - (a) Neighbourhood search, with a neighbourhood of 1
 - There were any results returned from the dictionary, so precision isn't well-defined ($\frac{0}{0}$)
 - (b) Neighbourhood search, with a neighbourhood of 2
 - There was one entry returned from the dictionary (**leader**), but it wasn't **lent**, so the precision is $\frac{0}{1} = 0$.
 - (c) Neighbourhood search, with a neighbourhood of 3
 - There were five entries returned from the dictionary, and **lent** was one of them. The precision of this system is the number of correct responses (1) out of the total number of attempted responses (5), $\frac{1}{5} = 20\%$
 - (d) Global Edit Distance, with a parameter $[m, i, d, r] = [1, -1, -1, -1]$
 - There were two (tied) results from the dictionary (**blenders** and **leader**), but no **lent**, so the precision is $\frac{0}{2} = 0$
 - (e) Local Edit Distance, with a parameter $[m, i, d, r] = [1, -1, -1, -1]$
 - There was just a single result (**blenders**) which wasn't **lent**, so the precision is 0
 - (f) N-gram Distance, where n is 2 (and padding with terminals)
 - There was a single result which was **lent**, so the precision is $\frac{1}{1} = 100\%$
 - (g) Using the Soundex transformation, and then looking for exact matches
 - (h) Using the Soundex transformation, and then permitting a 1-neighbourhood
 - There weren't any exact matches with the Soundex code of **lended**, so precision isn't well defined
 - Allowing approximate matches of the Soundex code meant that there were four results, including **lent**, so the precision is $\frac{1}{4} = 25\%$

3. What is the difference between “data retrieval” and “information retrieval”? Why is the latter a knowledge task, but the former is not?
 - The main difference here is the existence of people — users. Because people are wildly divergent, the notion of a relevant result in information retrieval depends on contextualising the data to the particular user (which may be very difficult, because we have an imperfect model of the user and indeed the user has an imperfect model of their needs!). Whereas with data retrieval, there is a particular unit of data (bitstream) that we need to access in memory or on a hard drive, and there is generally no ambiguity.
4. **[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?
 - These aren’t straightforward questions to answer. But, to an order of magnitude, a small city library might have about 10K books; a larger one, maybe 30K. I might estimate the word count of a typical book to be about 50K (many are longer; many are shorter; there are varying definitions of “word”), which would situate a library as carrying roughly 1G words. The US Library of Congress catalogues about 2.3M books, so perhaps 100G words.
 - As of 2008, Google claimed to index 1T unique urls (<http://googleblog.blogspot.com.au/2008/07/we-knew-web-was-big.html>); by 2012, this had supposedly risen to 30T (<http://www.google.com/insidesearch/howsearchworks/thestory/>). But maybe take all of this with a grain of salt! :-) Estimating the number of words on the Web is even harder — Google tells me that the mean document size is about 400KB, but much of that isn’t going to be text. I might ballpark about 1000 words (roughly 6KB of the 400KB) per document, which might be upwards of 10000T words (or more!, but probably less)!
5. Identify some different types of “informational needs.”
 - Rehashing the lecture slides:
 - Requests for informations, e.g. “global warming”
 - Factoid questions, e.g. “melting point of lead”
 - Topic tracking, e.g. “Dutch elections” [as in, the most recent ones]
 - Navigational, e.g. “University of Melbourne home page”
 - Service or transactional, e.g. “Mac powerbook”
 - Geospatial, e.g. “Carlton restaurant”
 - This isn’t an exhaustive list. Nor is it non-overlapping: for example, most queries can be construed as being navigational in nature (as the user is likely to click through to a relevant document), and many are informational as well.