<center>
School of Computing and Information Systems

The University of Melbourne

COMP90049 Knowledge Technologies (Semester 1, 2017)

Workshop exercises: Week 7
</center>

1. Assume that we have crawled the following "documents":

> (1) The South Australian Tourism Commission has defended a marketing strategy which pays celebrities to promote Kangaroo Island tourism to their followers on Twitter.
> (2) Mr O'Loughlin welcomed the attention the use of Twitter had now attracted.
> (3) Some of the tweeting refers to a current television advertisement promoting Kangaroo Island.
> (4) Those used by the Commission have included chef Matt Moran, TV performer Sophie Falkiner and singer Shannon Noll.
> (5) He said there was nothing secretive about the payments to celebrities to tweet the virtues of a tourism destination.
> (6) Marketing director of SA Tourism, David O'Loughlin, said there was no ethical problem with using such marketing and it might continue to be used.
> (7) Depending on their following, celebrities can be paid up to $750 for one tweet about the island.

- 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
- Using the vector space model and the cosine measure, rank the documents for the query commission to island on twitter
  - (a) Using the weighting functions $w_{d,t} = f_{d,t}$ and $w_{q,t} = \frac{N}{f_t}$
  - (b) Using the weighting functions $w_{d,t} = 1 + \log_2 f_{d,t}$ and $w_{q,t} = \log_2(1 + \frac{N}{f_t})$

2. What is the main problem of using **accumulators** when querying? What heuristics can we use to solve this problem?

3. What is a **phrase query**, and why is an inverted index — like the one from the question above — inadequate for phrase querying? How could the index be altered to support this style of querying?

4. What is **link analysis**? What aspects of user behaviour or the nature of data on the Web is it trying to model?