



THE UNIVERSITY OF
MELBOURNE

Student Number: _____

Faculty: Computing and Information Systems

Subject Number: COMP90042

Subject Name: Web Search and Text Analysis

Writing Time: 2 hours

Reading Time: 15 minutes

Open Book Status: Closed book

Number of Pages: This paper has 5 pages including this cover page.

Authorized Materials: None

Instructions to Students: Every question and every part of each question should be attempted. A suggested amount of time for each question is shown next to the question.

Instructions to Invigilators: Invigilate vigilantly!

Paper to be held by Baillieu Library: Yes

Extra Materials Required: None

Question 1:

a). Consider the following Python function skeleton:

```
def tfidf_norm(dfdt, df, ndocs):
    """
```

Calculate unit-length-normalized TFIDF value for document.

Arguments:

- dfdt: a dictionary of { term : fdt }, giving the fdt for each term in the document
- df: a dictionary of { term : ndocs }, giving the number of documents that each term occurs in
- ndocs: the total number of documents in the collection (an integer)

Return value:

Returns a { term : w_dt } dictionary, giving the unit-length normalized TFIDF weight of each term in the document.

```
    """
```

```
pass
```

Implement this function. Use the TF and IDF formulae specified for Project 1.

[15 minutes]

b). Let \vec{q}_o be the original query vector, and $D_r = \{\vec{d}_1, \dots, \vec{d}_i, \dots, \vec{d}_r\}$ the document vectors of the top r documents returned for \vec{q}_o . What is the formula for the expanded pseudo-query-document, \vec{q}_e , in Rocchio pseudo-relevance feedback? (You may use $+$ to represent vector addition; that is, $\vec{a} + \vec{b} = \{a_1 + b_1, \dots, a_i + b_i, \dots, a_n + b_n\}$.)

[10 minutes]

[25 minutes in total]

Question 2:

Consider the following two formulae:

$$\mathbf{X}_{t \times d} = \mathbf{T}_{t \times t} \mathbf{\Sigma}_{t \times d} (\mathbf{D}_{d \times d})^T \quad (1)$$

$$\hat{\mathbf{X}}_{t \times d} = \hat{\mathbf{T}}_{t \times k} \hat{\mathbf{\Sigma}}_{k \times k} (\hat{\mathbf{D}}_{d \times k})^T \quad (2)$$

Equation 1 indicates the matrix operation known as singular value decomposition; and Equation 2 is a reduced-rank representation of the SVD.

a). What is the name of the text analysis technique that is built upon Equation 2?

[5 minutes]

end of page

b). Under the text analysis technique mentioned in Question 2.a, what does the dimension k in Equation 2 represent? [5 minutes]

c). Under the text analysis technique mentioned in Question 2.a, what information does the i 'th row of $\hat{\mathbf{T}}$ give us? [5 minutes]

[15 minutes in total]

Question 3:

We run a binary classifier against a test set of 1000 examples and get the following confusion matrix:

Predicted class	True class	
	1	0
1	6	14
0	24	956

where “1” is the positive class, “0” is the negative class.

a). What recall has the classifier achieved on the test set? [5 minutes]

b). What precision has the classifier achieved on the test set? [5 minutes]

c). What accuracy has the classifier achieved on the test set? [5 minutes]

d). The accuracy metric gives a very different picture of classifier performance from precision and recall. Why? Which metric score (taken at face value) is a better indicator of how well the classifier has done? [10 minutes]

[25 minutes in total]

Question 4:

In probabilistic IR, we aim to rank documents by decreasing probability of relevance to the query, $P(R|d, q)$. For general use, we only care about the ranking, not the precise probability of relevance. Therefore, any monotonic transformation of $P(R|d, q)$ will do.

end of page

a). One such transformation is the log odds ratio:

$$\log O(R|d, q) \propto \log P(d|R, q) - \log P(d|\bar{R}, q) \quad (3)$$

Show the working to get the above expression from:

$$O(R|d, q) = \frac{P(R|d, q)}{P(\bar{R}|d, q)} \quad (4)$$

[10 minutes]

b). To calculate Equation 3, we need a model for $P(d|R, q)$. One such model is the binary independence model (BIM). There are two main assumptions of the BIM; what are they? [5 minutes]

[15 minutes in total]

Question 5:

a). In the language model approach to information retrieval, we try to estimate $P(d|q)$, using the monotonic inversion:

$$P(d|q) \propto P(q|d) \quad (5)$$

Regarding the n -word query q as a list of word occurrences $\{q_1, \dots, q_i, \dots, q_n\}$, write an expression for $P(q|d)$ under the unigram language model. [5 minutes]

b). Jelinek-Mercer smoothing defines the quantity:

$$P_{JM}(w|d) = (1 - \lambda)P_{mle}(w|d) + \lambda P(w|C) \quad (6)$$

How is $P_{mle}(w|d)$ usually estimated?

[5 minutes]

c). In Equation 6, how is $P(w|C)$ usually estimated?

[5 minutes]

d). Smoothing performs two main functions in language models for IR. What are they?

[5 minutes]

[20 minutes in total]

end of page

Question 6:

a). In Naive Bayes classifiers, we estimate the probability that the n -word document $d = \{w_1, \dots, w_i, \dots, w_n\}$ belongs to class c as:

$$P(c|d) \propto P(c) \prod_i P(w_i|c) \quad (7)$$

How is $P(c)$ estimated? [5 minutes]

b). The Laplace-smoothed estimate for $P(w_i|c)$ in Equation 7 can be written:

$$P(w_i|c) = \frac{a + 1}{b + |V|} \quad (8)$$

where $|V|$ is the size of the vocabulary, and a and b are standing in for more meaningful variable names. What is a and what is b ? [5 minutes]

c). Naive Bayes tends to give too-extreme estimates of the probability of relevance. For instance, if we examined a set of documents each of which a Naive Bayes classifier estimated were 99% likely to be relevant, we might find that only 80% (say) were. Why does Naive Bayes tend to exaggerate probabilities like this? [10 minutes]

[20 minutes in total]

[120 MINUTES FOR EXAM]

end of exam



THE UNIVERSITY OF

MELBOURNE

Library Course Work Collections

Author/s:

Computing and Information Systems

Title:

Web Search And Text Analysis, 2014 Semester 1, Comp90042

Date:

2014

Persistent Link:

<http://hdl.handle.net/11343/52187>