

--

COMP90049 Knowledge Technologies

Reading Time allowed: 15 minutes

Number of pages: 8 including this page

This paper counts for 50% of your final grade.

There are 85 marks in total, or 1 mark per 1.5 minutes. Note that questions are not of equal value. All questions should be interpreted as referring to concepts given in this subject, whether or not it is explicitly stated.

No external materials may be used for this exam, but calculators are permitted (although not necessary). You may leave square roots and logarithms without integer solutions (like $\sqrt{2}$) unsimplified.

Unless otherwise indicated, you must show your working for each problem. Please indicate your final answers clearly for problems where you show intermediate steps.

The students require script books.

Calculators are permitted; other materials are not authorised.

The examination paper should not leave the examination hall; this exam is to be held on record in the Baillieu Library.

[illegible]

Part I : Text Processing**[28 marks in total]**

1. Describe two (or more) steps that we would typically perform in the “tokenisation” process for an Information Retrieval collection, according to the description from this subject. [4 marks]
 - Folding case — making everything lowercase
 - Stripping punctuation — removing some nonalphabetic characters like “,”
 - Stemming — removing suffixes (in English) to change a word to a base form
 - Splitting into tokens — splitting the document based on whitespace (in English) and maybe some punctuation
 - And other possible answers

2. It has been claimed that there are three primary types of “information need” in a web search context: “informational”, “navigational”, and “transactional”. Briefly describe each of these, optionally with the aid of an example. [4 marks]
 - Informational: tell me more about this topic, e.g. **history of Australia**
 - Navigational: take me to the URL corresponding to this topic, e.g. **Unimelb homepage**
 - Transactional: interface with a database, so that I can perform some service (like buying a product), e.g. **iphone ebay**

continued ...

3. In the context of Information Retrieval:

(a) Explain how “data retrieval” is different to “information retrieval”.

[2 marks]

- Data retrieval: getting some variable value out of memory, or a record out of a database, etc.
- Information retrieval: trying to find some document(s) which meet the users information need expressed by the query
- Information retrieval doesnt have an exact answer; whether the results are useful depends on the user issuing the query

(b) Give an example of a method or source of information that we might incorporate into our engine, that is specific to Web-scale information retrieval.

[1 marks]

- Link analysis
- Clickthrough data
- And other possible answers

4. ...and other questions to add up to the marks as stated above. :-)

Part II: Data Mining/Machine Learning [57 marks in total]

For these questions, we have a training dataset comprised of the following 6 instances, 3 attributes, and two classes F and T, with a single test instance labelled with “?”:

ele	fed	aus	CLASS
1	1	1	F
1	0	0	F
1	1	0	T
1	1	0	T
1	1	1	T
1	1	1	T
0	0	0	?

5. Classify the test instance according to the method of “Naive Bayes”, as described in this subject. [4 marks]

- We need to pre-calculate a bunch of probabilities: $P(f) = \frac{2}{6}$, $P(t) = \frac{4}{6}$; $P(e = 0|f) = 0$, $P(e = 0|t) = 0$, $P(fed = 0|f) = \frac{1}{2}$, $P(fed = 0|t) = 0$, $P(a = 0|f) = \frac{1}{2}$, $P(a = 0|t) = \frac{2}{4}$
- When we substitute, we need to replace 0 values with ϵ , a small positive constant value.
- We calculate the scores for the two classes F and T:

$$\begin{aligned}
 \text{F} &: P(f)P(e = 0|f)P(fed = 0|f)P(a = 0|f) \\
 &= \left(\frac{1}{3}\right)(\epsilon)\left(\frac{1}{2}\right)\left(\frac{1}{2}\right) = \frac{\epsilon}{12} \\
 \text{T} &: P(t)P(e = 0|t)P(fed = 0|t)P(a = 0|t) \\
 &= \left(\frac{2}{3}\right)(\epsilon)(\epsilon)\left(\frac{1}{2}\right) = \frac{\epsilon^2}{3}
 \end{aligned}$$

- ϵ is less than $\frac{1}{4}$, so F has the larger value — so that is the class we choose.

6. Explain why “1-Nearest Neighbour” will give a different prediction to “3-Nearest Neighbour” on the given test instance. (It is not necessary to show your work for this question; an explanation which refers to the data should suffice.) [2 marks]

- Regardless of the distance metric were using, clearly the second instance (1,0,0:F) has the smallest distance; so, 1-NN will predict F.
- The next best instance(s) are (1,1,0:T), of which there are two.
- So, for 3-NN, we will observe 2 T instances and 1 F instance among the 3 nearest neighbours; there are more T than F, so we classify it as T.
- (Since the question doesnt ask for working, it is possible to explain this more compactly.)

7. Consider the method of “Random Forests”:

- (a) Briefly explain how a Random Forest would be constructed on the training data above. [4 marks]
- For a Random Forest, we will construct a bunch of “Random Trees”, in this case, let's say 10 of them.
 - For each tree, we will use Bagging to come up with a different training dataset: we will re-sample the instances with replacement, until we have 6 (possibly repeated) training instances.
 - When building our tree, at each node, we only consider a random sample of the attributes. Because we have so few attributes, let's say: we randomly choose 2 of the 3 attributes for consideration at the root node; we consider both of the remaining attributes at the second layer; we consider the final attribute at the third layer.
- (b) Is there any evidence that a Random Forest would label the given test instance differently to a regular “Decision Tree”? [3 marks]
- (Aside: this semester, we didn't cover this topic in enough detail to actually construct a sensible response. On the other hand, it is worth noting that there will be some more difficult questions like this one. If you need to think about the problem, the harder questions might take longer to answer than the marks suggest!)
 - Probably not:
 - The regular decision tree will have **fed** at the root, as it is clearly the most useful attribute — and therefore classify the test instance as F.
 - When bagging, the chance of any individual instance being present in the training data is about 63%. If the second instance isn't present, we are going to predict T — this will happen for 37% of the trees.
 - If the second instance is present at least once, and **fed** is in the choices at the root ($0.63 \times \frac{2}{3} = 42\%$ of the trees), we will predict F.
 - (There will be a small number of trees where the instance distribution makes all of the attributes useless — the prediction will be the majority class T. For example, if all of the instances have the same attribute values. This will be substantially less than 50% of the remaining trees.)
 - For the other trees, **a** will be probably placed at the root (**e** is completely useless, regardless of the instance distribution). **fed** will be placed at the **a=0** branch, and then we will yet again predict F.

continued ...

- All in all, more than 50% of the trees will predict τ , so the prediction is likely to be the same.

8. Exclude the CLASS labels from the dataset, and cluster all 7 instances using the method of “ k -means”. Apply the Manhattan Distance as a similarity measure; use the second (1,0,0) and third (1,1,0) instances as seeds. [4 marks]

- Let’s say Cluster 1 C_1 begins at 1,0,0 and Cluster 2 C_2 begins at 1,1,0.
- For each instance, we calculate the Manhattan distance to the two clusters. I will show the workings for one instance; it is obviously crazy to try to write the whole formula 14 times in 5–6 minutes.
 - First instance to C_1 : $|1 - 1| + |1 - 0| + |1 - 0| = 2$; to C_2 : $|1 - 1| + |1 - 1| + |1 - 0| = 1$.
 - Second instance to C_1 : 0; to C_2 : 1.
 - Third instance to C_1 : 1; to C_2 : 0.
 - Fourth instance is the same as third instance; fifth and sixth instances are the same as first instance.
 - Seventh instance to C_1 : 1; to C_2 : 2.
- So, the first, third, fourth, fifth, and sixth instances are closer to C_2 ; the second and seventh are closer to C_1 . We now update our centroids:

$$C_1 : \frac{1}{2}[(1, 0, 0) + (0, 0, 0)] = (0.5, 0, 0)$$

$$C_2 : \frac{1}{5}[(1, 1, 1) + (1, 1, 0) + (1, 1, 0) + (1, 1, 1) + (1, 1, 1)] = (1, 1, 0.6)$$

- Now, we re-calculate the Manhattan distances:
 - First instance to C_1 : $|1 - 0.5| + |1 - 0| + |1 - 0| = 2.5$; to C_2 : $|1 - 1| + |1 - 1| + |1 - 0.6| = 0.4$.
 - Second instance to C_1 : 0.5; to C_2 : 1.6.
 - Third instance to C_1 : 1.5; to C_2 : 0.6.
 - Fourth instance is the same as third instance; fifth and sixth instances are the same as first instance.
 - Seventh instance to C_1 : 0.5; to C_2 : 2.6.
- So, the first, third, fourth, fifth, and sixth instances are closer to C_2 ; the second and seventh are closer to C_1 . This is the same as the previous iteration, so this is the clustering.

9. ...and other questions to add up to the marks as stated above. :-)