Department of Computing & Information Systems

COMP90049 Knowledge Technologies Final exam, Semester 1, 2017

Date: 9 June, 2017

Time: 08:30am

Reading Time allowed: 15 minutes

Writing Time allowed: 2 hours

Number of pages: 7 including this page, and the blank page overleaf

Instructions to candidates:

This paper counts for 50% of your final grade.

Answer all questions on the ruled pages in the script book(s) provided.

There are 85 marks in total, or 1 mark per 1.4 minutes. Note that questions are not of equal value. All questions should be interpretted 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.

Instructions to invigilators:

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.

Examiner's use only:

Q1	Q2	Q3	Q4	Q5	Q6	Q7	Q8	Q9	Q10
8	21	5	4	4	11	13	7	4	8
						j-	-		
						=			

(This page has been intentionally left blank.)

1. (a). crawling: Inoling and downloading downerts from the web,
gather information.

Canonical

Parsing: transform the data into canonical form. ofter some steps like stemming, (ase folder, turning obsuments into a train of topens.

Indexing: create a inverted that of each tokens for each downal index.

Week query: When we got the query attornise (in a similar manner to our downal we can apply query and or model like TF-1DT to apply a downer vanlery.

When we got the query attorner vanlery. (c) We only focus on the grey term Go we tighove the term which doesn't 2. (a) - the -> 1(2) -> 2(2) Spain > 1(1) -> Z(3) (at > 1(3) - 1(3) - 1(3) - 1(3) - 1(3)Wa, in = ==== in 7 1 (1) 72 (3) Wq,t= <0,1,0,1,0,0,0) hat > 1 (1) PaA: 22, 1, 3, 1, 1, 0, 0, 0) rain -> 2(1) Doc B: (2,3,0,3,0,1,1) mainly -> Z(1) :. COS (DOLA, Webt) = (2,1,3,1,1,0,0,0). < 0,1,0,1,0,0) plans -> 2(1) (6) | Spain | in | hat = 1+1+2 122+12+32+12+12. 12+12+12 DocA 1 1 F = 45 = 16 Doe A has three words, votur Doe A 65 (POB, wait) = (> 1 < > 1 < > 1 ... Du Builte ligher trank. 3+3+0 J3+4+9+9 J36 DOCA: LSpain> 1/2 10 1 10 -1 10 -1 10 = 12.525 - 56 > COSA DOCB= --= (17 nc17 nc07 = 0.

a).

Part I: Information Retrieval

[38 marks in total]

- 1. Consider building a "ranked Information Retrieval engine" for Web documents:
 - (a) Such an engine typically contains 4 main components: name them, and briefly (in a sentence or two) explain what the purpose of each one is. [6 marks]
 - Instead of the constantly-changing Web, consider building a ranked Information Retrieval engine for a mostly static database of newspaper articles. Give an example of a component that would be different.

not necessary to crawlette obcuments which have been searched, not the crawley, a list a V to,

- 2. For this question, consider the (very small) collection of documents, labelled A) and B) below (the label is not part of the document text):
 - A) the Spain cat cat in the cat hat
 - B) the rain in Spain mainly in the Spain plains in Spain and a query Q) Spain in hat
 - (a) For a standard inverted index consistent with the lecture or workshop notation; give a representation (in words or as a diagram) of the "inverted lists" for the 8 terms in this collection. (There is no need to explicitly indicate the "search structure" or "mapping table".)

 [6 marks]
 - (b) If we wish to apply the method of "Boolean querying" assuming that the query is implicitly a conjunction of terms describe the procedure by which the query engine would arrive at the result set {A}.
 - (c) Determine which of the documents above would be returned higher (nearer to the top of the ranking) for a "ranked query engine", based on the following "TF-IDF model", suitably interpretted in the context of this subject:

$$w_{d,t} = f_{d,t}$$

$$w_{q,t} = \frac{N}{f_t}$$

(Remember to show your work; there should be no need to simplify irrational square roots to solve this problem.) [7 marks]

(d) If we had instead used the following TF-IDF model, how do you expect the ranking to change compared to the model from (c)? Why is this the case?

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

(Note that you do not need to calculate the steps of the model for this question.) [4 marks]

3. Consider the Information Retrieval evaluation metric "Average Precision" (AP): topk documents.

(a) Why is AP preferred over other metrics like Precision and Recall, in the context of evaluating a typical Information Retrieval engine?

users only prefer to bohat the top to or 20 docupe marks precision is

(b) Explain the procedure by which AP is calculated; you may give an example if you think it will help your arrival.

example if you think it will help your explanation. [3 marks] Recall is not used in IR, we don't hnow, based oh

4. One possible extension to a ranked Information Retrieval engine is "linksers analysis", of which "PageRank" is an example.

- (a) Define "link analysis", in the context of this subject. [1 marks]
- (b) What is the main idea behind incorporating this information into the ranking? [1 marks]

(c) What are the two main components of the PageRank algorithm, and

Weighty cless which alters the important of documents based anon.

(ink and web)

higher ranking downals will be more linked.

outgoing this.

outgoing this.

user can itself the link

for 90 to the webser.

to go to the webset

5. $d(1,7) = J(1-0)^2 + (2-0)^2 + (0-0)^2 = J5$ $d(2,7) = J0^2 + Z^2 + 0^2 = J4$ $d(3,7) = J3^2 = J9$ $d(4,7) = J3^2 = J9$ $d(5,7) = J7^2 + 1^2 = J2V$ $d(5,7) = J3^2 = J9$. $d(5,7) = J3^2 = J9$.

- 6. i). SVM will find a hyperplane to linearly separate two classes,
 If it isn't linearly seperatelles, then it would apply as bernal function to
 make it linearly separatell separable
 - Partition the training data based on the best line chyperplane).

 that divide the possible instance of the class the we're

 wokey for from regathe instances.

Part II: Data Mining/Machine Learning

[47 marks in total]

For Questions 5–8 in this section, we have a training dataset comprised of the following 6 instances, 3 attributes, and two classes YESTERDAY and TOMORROW, and a single test instance labelled with ?:

	tfw	ftw	wtf	CLASS
1	1	2	0	YESTERDAY
2	0	2	0	YESTERDAY
3	1	1	1	TOMORROW
X	0	0	3	TOMORROW
5	1	0	1	TOMORROW
6	0	3	0	TOMORROW
7	0	0	0	?

5. Classify the test instance according to the (simple) "3-Nearest Neighbour" method, given the following "distance" metric:

$$D(a,b) = \sqrt{\sum_i (a_i - b_i)^2}$$

(Show your work; it should not be necessary to simplify irrational square roots to solve this problem.) [4 marks]

- 6. Consider applying the method of Support Vector Machines to classify this test instance:
 - (a) Briefly explain (in a couple of sentences) the logic behind training a Support Vector Machine. You may use some logical simplifications (like "line" for "hyperplane"); you should aim to avoid the mathematical formulation unless you really understand it. [4 marks]

(You may include a diagram if you think it will help your explanation) For the given instances above, we will find that there is no solution.

Explain what this would look like, if we were to graph these instances.

(You may include a diagram if you think it will help your explanation) the linearly separate.

it can be approximate — please don't attempt to graph more than 2 dimensions simultaneously.) [3 marks]

(c) There are two main alternatives for building an SVM when we discover that there is no solution: briefly explain them (one sentence each).

(1) Soft margin. telax the notion of linear squadoility by allowing some number of points on the continued margin side of the

(2) Hernel function: use this function to map the convert points into higher dimensions, and find a hyperplane.

7-C) a sore the attribute on values

(2) Unearly scan these values, each time updating the count matrix and computing Gini index at points where class labor

charges

(3) choose the split position that has the Gini Index,

Few 0 0 1 2 2 3 class: tomorrow tomorrow tomorrow yeslady yeslady tomorrow. point B point B

> for the 41, 3 tomoron Gavi = 1- (3) - 0=0 Jew >1, [3] = 4 .. GINI sphl: \(\frac{1}{2}\times 0 + \frac{1}{2}\times \frac{1}{9} = \frac{2}{9}.

GW1=1-1-1-1-1-1 · GINZANI = 13 x \$+ 0x 6 = 329 : split of A. A

a) measure ele impurity of a mode. Maximum value of GIM when records are equally distributed among all classes most impure.

Minimum when records belong to I class most pure.

Chose ele linest GIM.

COMP90049, S1 2017

page 6 of 7

cmall.

7. If we wished to build a "Decision Tree" to classify the given test instance:
a) Explain the logic behind the GINI coefficient $(G = 1 - \sum P(j)^2)$, and gree is
how we can use it to build the tree. [2 marks]
(b) If we treated these attributes as "categorical", which attribute would of GZMZ
be placed at the root of the tree? Consequently, what would be the predicted label for the test instance? (You do not have to show your
work; an explanation which refers to the data is sufficient.) [4 marks]
(c) If we treated these attributes as "continuous", it is more difficult. Briefly demonstrate how we could determine the root attribute, working through the process for either (not both) of ftw or wtf? [4 marks]
(d) The method of "Bagging" usually wouldn't help with the procedure from (c) — but in this case, it could. Explain what might happen to vastly simplify the process, referring to the dataset where necessary. [3 marks]
Associated (d) => the >2, cz, the formore
8. The method of "Naive Bayes" that we discussed in this subject could only
8. The method of "Naive Bayes" that we discussed in this subject could only be applied to "cotogorical" attributes. Treat each number as a distinct the thank set

(a) To "train" a Naive Bayes model, we must estimate two types of probabilities. Name the two types. Give one example calculation for each, based on the training data above. [3 marks]

value when answering this question.

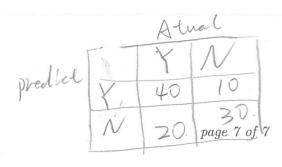
(b) There is a single crucial factor which causes this Naive Bayes model to label the given test instance as TOMORROW. Explain what this is, referring to the data as necessary. (You do not need to work through the entire procedure, but if you are unsure, we will accept the working as partial credit.)

[4 marks]

6) For Esmonow: PCT). P(tfw=0|T) (PCtfw=0|T) - P(w=0|T)

= \frac{3}{3} \times \frac{1}{2} \times \frac{1}{2}

1 - 0



COMP90049, S1 2017

- 9. Consider evaluating a Machine Learning classifier, such that:
 We have 100 instances in our development data, of which 60 are actually labelled as \$\Omega\$, and the rest as N. Our classifier predicted that 40 of the actual Y instances were indeed Y and 10 of the actual N instances were Y.
 - (a) Show the "Confusion Matrix" that summarises all of the classifications.

[2 marks]

(b) Find the "Accuracy" of this classifier.

[1 marks]

(c) Find the "Precision" (of class Y) for this classifier.

[1 marks]

(5) Accorn = 40+30 = 0.7=70%

10. Recall the definition of "knowledge tasks" from this subject: $\frac{\xi}{100} = \frac{100}{100} = \frac{100}$

- (a) The process of "Clustering" seems to clearly fit the most of criteria of things that are knowledge tasks. Give an example of this. [3 marks]
- (b) "Classification", on the other hand, mostly doesn't fit these criteria. Explain why, and explain how we might end up with "knowledge" nonetheless. [3 marks]
- (c) An example of an "Association Rule" is $\{tfw\}\rightarrow \{ftw\}$. Explain how such a rule encodes "knowledge". [2 marks]

a) knowledge tash: onleame are not well-deflect; clistley class group. Someon based on similarity, bosed on users.

b) concrete task.



Library Course Work Collections

Author/s:

Computing and Information Systems

Title:

Knowledge Technologies, 2017, Semester 1, COMP90049

Date:

2017

Persistent Link:

http://hdl.handle.net/11343/216574

