

Student Number

--

The University of Melbourne
School of Computing and Information Systems

Final Examination, Semester 1, 2017
COMP30027 Machine Learning

Reading Time: 15 minutes. **Writing Time:** 2 hours.

This paper has 7 pages including this cover page.

Instructions to Invigilators:

Students should be provided with script books, and should answer all questions in the provided script book. Students may not remove any part of the examination paper from the examination room.

Instructions to Students:

There are 9 questions in the exam worth a total of 120 marks, making up 60% of the total assessment for the subject.

- Please answer all questions in the script book provided, starting each question on a new page. Please write your student ID in the space above and also on the front of each script book you use. When you are finished, place the exam paper inside the front cover of the script book.
- Your writing should be clear; illegible answers will not be marked.

Authorised Materials: No materials are authorised.

Calculators: Students are permitted to use calculators.

Library: This paper may be held by the Baillieu Library.

<i>Examiners' use only</i>									
1	2	3	4	5	6	7	8	9	Total

Section A: Short Answer Questions [28 marks]

Answer each of the questions in this section as briefly as possible. Expect to answer each sub-question in a couple of lines.

Question 1: Short Answer Questions [28 marks]

1. What is the primary difference between “supervised” and “unsupervised” learning? [1.5 marks]
2. For each of the following features, indicate which of “numeric”, “ordinal”, and “categorical” best captures its type: [4 marks]
 - (a) blood pressure level, with possible values {**low**, **medium**, **high**}
 - (b) age, with possible values [0,120]
 - (c) weather, with possible values {**clear**, **rain**, **snow**}
 - (d) abalone sex, with possible values {**male**, **female**, **infant**}
3. Describe a strategy for measuring the distance between two data points comprising of “categorical” features. [1.5 marks]
4. What is the relationship between “accuracy” and “error rate” in evaluation? [1.5 marks]
5. With the aid of a diagram, describe what is meant by “maximal marginal” in the context of training a “support vector machine”. [1.5 marks]
6. What makes a feature “good”, i.e. worth keeping in a feature representation? How might we measure that “goodness”? [3 marks]
7. For each of the following models, state whether it is canonically applied in a “classification”, “regression” or “clustering” setting: [6 marks]
 - (a) multi-layer perceptron with a softmax final layer
 - (b) soft k -means
 - (c) multi-response linear regression
 - (d) logistic regression
 - (e) model tree
 - (f) support vector regression
8. With the aid of an example, briefly describe what a “hyperparameter” is. [3 marks]
9. With the use of an example, outline what “stacking” is. [1.5 marks]
10. What is the convergence criterion for the “EM algorithm”? [1.5 marks]
11. Outline the basis of “purity” as a form of cluster evaluation. [1.5 marks]
12. What is the underlying assumption behind active learning based on “query-by-committee”? [1.5 marks]

Section B: Methodological Questions [30 marks]

In this section you are asked to demonstrate your conceptual understanding of a subset of the methods that we have studied in this subject.

Question 2: Random Forests [7 marks]

“Random forests” are based on decision trees under different dimensions of “randomisation”. With reference to the following toy training dataset, provide a brief outline of two (2) such “random processes” used in training a random forest. (You should give examples as necessary; it is not necessary to draw the resulting trees, although you may do so if you wish.)

<i>Instance ID</i>	<i>Feature 1</i>	<i>Feature 2</i>	<i>Feature 3</i>	<i>Class</i>
A:	1	1	1	True
B:	0	2	0	True
C:	3	0	0	False
D:	1	1	2	False
E:	2	0	2	False

Question 3: HMMs [9 marks]

This question is on “hidden Markov models”.

1. In the “forward algorithm”, $\alpha_t(j)$ is used to “memoise” a particular value for each state j and observation t . Describe what each $\alpha_t(j)$ represents. [3 marks]
2. In the “Viterbi algorithm”, two memoisation variables are used describe for each combination of state j and observation t : (1) $\beta_t(j)$ (which plays a similar role to $\alpha_t(j)$ in the forward algorithm); and (2) $\psi_t(j)$. Describe what each $\psi_t(j)$ represents. [3 marks]
3. Why do we tend to use “log probabilities” in the Viterbi algorithm but not the forward algorithm? [3 marks]

Question 4: Model learning [14 marks]

We have discussed that the objective in much of supervised machine learning is to derive a model that explains (“fits”) a set of (labelled) data. In a basic “curve fitting” scenario, we will assume that all of the needed data is available to us at the time of building the model; the objective is to fit a model to that data. In machine learning, in contrast, we have only a subset of possible data available to us when we build the model. This contrast has certain implications for how we approach machine learning, which we explore in this question.

Discuss the implications of knowing that we do not have all possible data for a given problem, in terms of the following aspects:

1. Is our primary objective in machine learning to derive a model that fits the subset of the data that we do have? Why or why not? [3 marks]
2. Explain how we can use our limited data, in a machine learning context, to demonstrate whether or not our objective has been met. [3 marks]
3. Identify and explain one important problem that can emerge with respect to this primary objective, even if we are successful in deriving a good model for the data that we have. Name one specific technique discussed in class that can be applied to mitigate this problem, and explain how it does so. [3 marks]
4. Define “bias” and “variance”, indicating how we might detect each one. Discuss how bias and variance relate to each other in the context of our primary objective. [5 marks]

Section C: Numeric Questions [42 marks]

In this section you are asked to demonstrate your understanding of a subset of the methods that we have studied in this subject, in being able to perform numeric calculations. Questions 5 and 6 both make use of the following training data set:

early	tie	label
N	Y	dinner
Y	N	tea
Y	N	dinner
N	N	dinner
Y	N	dinner
N	N	tea

Table 1: Training Dataset for Questions 5 and 6

Question 5: Naive Bayes [9 marks]

Given the training dataset from Table 1:

1. Using the method of “maximum likelihood estimation” (without smoothing), compute $P(\text{dinner})$ and $P(\text{dinner}|\text{tie} = \text{Y})$. [1.5 marks]
2. Apply the method of “Naive Bayes” as it was discussed in the lectures, to predict the label of the test instance $\{\text{early}=\text{N}, \text{tie}=\text{N}\}$; show your workings. [7.5 marks]

Question 6: Decision Trees [15 marks]

1. Briefly explain — in at most two sentences — the basic logic behind the “ID3” algorithmic approach toward building decision trees. This should be focussed on labelling the nodes and leaves; you do not have to explain edge-cases. [3 marks]
2. The criterion for labelling a node, as explained in the lectures, was based around the idea of “entropy” — what does entropy tell us about a node of a decision tree? [3 marks]
3. A very similar alternative to entropy is the so-called GINI coefficient, defined as follows:

$$GINI = 1 - \sum_{j \in C} [p(j)]^2$$

Calculate the GINI coefficient based on the labels in Table 1. [3 marks]

4. Extend the notion of “Information Gain” to “GINI Gain”, and demonstrate why `tie` would be chosen as the root of the decision tree on the given data. [3 marks]
5. Why will the resulting decision tree have difficulty classifying test instances like $\{\text{early}=\text{N}, \text{tie}=\text{N}\}$? [3 marks]

Question 7: Nearest Prototype and k -Nearest Neighbour [10.5 marks]

Given the following training dataset:

abv	opacity	label
4.8	0.20	ale
5.2	0.10	ale
5.0	0.33	ale
4.7	0.02	lager
5.1	0.23	lager
4.6	0.05	lager

1. Generate the “prototype” for each class in the training data. [3 marks]
2. For the test instance $\{\text{abv}=5.1, \text{opacity}=0.23\}$, determine which of the prototypes is more similar. (You will need to choose an appropriate metric, and show your work; do not just use inspection.) [3 marks]
3. What should happen if we instead use the “1-Nearest Neighbour” method to predict the label of that test instance? (You do not need to show your work.) [1.5 marks]
4. Give one possible reason why each of these two methods could reasonably claim to be making a better prediction for this instance. [3 marks]

Question 8: Evaluation [7.5 marks]

Assume that our development set contains 100 instances (truly) labelled as one of three classes as follows: 50 `acro` instances, 30 `base` instances, and 20 `claw` instances. We then build a classifier, apply it to this dataset, and observe the following confusion matrix:

		Actual		
		acro	base	claw
Predicted	acro	28	5	7
	base	10	10	0
	claw	0	8	12

1. Determine the “classification accuracy” of the system described above. [1.5 marks]
2. Calculate the “micro-averaged precision” and “macro-averaged precision” of this system. (Show your workings; you should simplify these to a single fraction or decimal value.) [3 marks]
3. In some contexts these three values are equal; here they are not. Briefly explain why. [3 marks]



THE UNIVERSITY OF

MELBOURNE

Library Course Work Collections

Author/s:

Computing and Information Systems

Title:

Machine Learning, 2017 Semester1, COMP30027

Date:

2017

Persistent Link:

<http://hdl.handle.net/11343/190833>