

Name of Candidate: .....

Student number: .....

Signature: .....

## COMP9417 Machine Learning and Data Mining

### Final Examination:

## SAMPLE QUESTIONS + ANSWERS

HERE ARE SEVEN QUESTIONS WHICH ARE *somewhat* REPRESENTATIVE OF THE TYPE THAT WILL BE IN THE FINAL EXAM. EACH QUESTION IS OF EQUAL VALUE, AND SHOULD NOT TAKE LONGER THAN 30 MINUTES TO COMPLETE. IN THE ACTUAL EXAM THERE WILL BE SOME DEGREE OF CHOICE AS TO WHICH QUESTIONS YOU CAN ANSWER. CANDIDATES MAY BRING AUTHORISED CALCULATORS TO THE EXAMINATION, BUT NO OTHER MATERIALS WILL BE PERMITTED.

This page intentionally left blank.

### Question 1 [20 marks]

#### Comparing Lazy and Eager Learning

The following truth table gives an “ $m$ -of- $n$  function” for three Boolean variables, where “1” denotes true and “0” denotes false. In this case the target function is: “exactly two out of three variables are true”.

X	Y	Z	Class
0	0	0	false
0	0	1	false
0	1	0	false
0	1	1	true
1	0	0	false
1	0	1	true
1	1	0	true
1	1	1	false

A) [4 marks]

Construct a decision tree which is complete and correct for the examples in the table. [Hint: draw a diagram.]

B) [4 marks]

Construct a set of classification rules from the tree which is complete and correct for the examples in the table. [Hint: use an *if-then-else* representation.]

C) [10 marks]

Suppose we define a simple measure of distance between two equal length strings of Boolean values, as follows. The distance between two such strings  $B_1$  and  $B_2$  is:

$$\text{distance}(B_1, B_2) = |(\sum B_1) - (\sum B_2)|$$

where  $\sum B_i$  is simply the number of variables with value 1 in string  $B_i$ . For example:

$$\text{distance}(\langle 0, 0, 0 \rangle, \langle 1, 1, 1 \rangle) = |0 - 3| = 3$$

and

$$\text{distance}(\langle 1, 0, 0 \rangle, \langle 0, 1, 0 \rangle) = |1 - 1| = 0$$

What is the LOOCV (“Leave-one-out cross-validation”) error of 2-Nearest Neighbour using our distance function on the examples in the table ? [Show your working.]

D) [2 marks]

Compare your three models. Which do you conclude provides a better representation for this particular problem ? Give your reasoning (one sentence).

### Question 1 ANSWER

A) Here is one complete and correct tree (there can be others):

```
X = 0
|   Y = 0: false (2)
|   Y = 1
|   |   Z = 0: false (1)
|   |   Z = 1: true (1)
X = 1
|   Y = 0
|   |   Z = 0: false (1)
|   |   Z = 1: true (1)
|   Y = 1
|   |   Z = 0: true (1)
|   |   Z = 1: false (1)
```

This is in “horizontal” format (numbers at leaves refer to numbers of examples at that leaf).

B) Here is one set of ordered rules:

```
IF (X = 0 AND Y = 0) THEN false ELSE
IF X = 0 THEN (IF Z = 0 THEN false ELSE true) ELSE
IF X = 1 THEN (IF Y = 0 THEN (IF Z = 0 THEN false ELSE true) ELSE
(IF Z = 0 THEN true ELSE false))
```

where the parentheses indicate nesting of compound conditions, or rules.

C) This distance function is unusual for  $k$ -NN, because it can result in multiple examples ( $> 2$ ) that are the same distance from the query. We need to make a design decision; here we simply take the majority vote of all neighbours that are the shortest distance from the query.

Example	Nearest neighbours	Majority Vote	Actual Class	Error
1	2,3,5	false	false	0
2	3,5	false	false	0
3	2, 5	false	false	0
4	6,7	true	true	0
5	2,3	false	false	0
6	4, 7	true	true	0
7	4,6	true	true	0
8	4,6,7	true	false	1

So the LOOCV error is  $\frac{1}{8}$ .

D) We can observe that to represent an “ $m$ -of- $n$  function” will result in very complex decision trees or rule sets due to the problem of replicating sub-trees or rules to express all the cases of the function, especially as  $m, n$  increase, but  $k$ -NN does not suffer from this representational issue, so on this criterion it is a better approach.

We can make the comment, however, that a linear threshold function learner like a perceptron would be a still better choice for this target concept.

## Question 2 [20 marks]

### Bayesian Learning

A) [2 marks]

Explain the difference between the *maximum a posteriori* hypothesis  $h_{\text{MAP}}$  and the *maximum likelihood* hypothesis  $h_{\text{ML}}$ .

B) [4 marks]

Would the Naive Bayes classifier be described as a generative or as a discriminative probabilistic model ? Explain your reasoning informally in terms of the conditional probabilities used in the Naive Bayes classifier.

C) [10 marks]

Using the multivariate Bernoulli distribution to model the probability of some type of weather occurring or not on a given day, from the following data calculate the probabilities required for a Naive Bayes classifier to be able to decide whether to play or not.

Use pseudo-counts (Laplace correction) to smooth the probability estimates.

Day	Cloudy	Windy	Play tennis
1	1	1	no
2	0	1	no
3	1	1	no
4	0	1	no
5	0	1	yes
6	1	0	yes

D) [4 marks]

To which class would your Naive Bayes classifier assign each of the following days ?

Day	Cloudy	Windy	Play tennis
7	0	0	?
8	0	1	?

## Question 2 ANSWER

A) In the Bayesian setting, both  $h_{\text{MAP}}$  and  $h_{\text{ML}}$  are single *most probable* hypotheses from the hypothesis space, given training data. However,  $h_{\text{MAP}}$  is found by taking the product of the likelihood and the prior, whereas  $h_{\text{ML}}$  is found simply by taking the likelihood.

B) The Naive Bayes classifier is used to determine the most probable class  $Y$  given the data  $X$ , so we can view this as computing the probability  $P(Y|X)$ . There are algorithms that learn this probability directly, such as logistic regression (see slide 120 in the lecture on Classification). Such models are characterised as *Discriminative*. On the other hand, the Naive Bayes classifier actually learns the joint probability  $P(X|Y)P(Y) = P(Y, X)$ . These models are termed *Generative* since they can be used to “generate” (sample) the examples for learning (see slide 67 in the lecture on Classification). Naive Bayes applies Bayes Theorem to actually do the classification.

C) Using the multivariate Bernoulli, treat examples as bit vectors (don't forget for the probability smoothing to add two “pseudo-examples” for each class, one with all bits set to 1 and the other with all 0). We have 2 examples of class ‘yes’ and 4 examples of class ‘no’ in the data. Adding bit vectors for each class results in (2, 4) for ‘no’ and (1, 1) for ‘yes’. We need to divide the counts by the number of examples in each class, but first add the counts for the probability smoothing, giving  $(\frac{3}{6}, \frac{5}{6}) = (0.5, 0.83)$  for ‘no’ and  $(\frac{2}{4}, \frac{2}{4}) = (0.5, 0.5)$  for ‘yes’.

D) To classify an example we need to compute the probabilities of the data given each class, then predict the class with the higher probability.

Day	Cloudy	Windy	Play tennis ?	Probability
7	0	0	no	$(1 - 0.5) \times (1 - 0.83) = 0.09$
			yes	$(1 - 0.5) \times (1 - 0.5) = 0.25$
8	0	1	no	$(1 - 0.5) \times 0.83 = 0.41$
			yes	$(1 - 0.5) \times 0.5 = 0.25$

So for example 7 the prediction is ‘yes’ and for example 8 it is ‘no’.

### Question 3 [20 marks]

#### Neural Networks

A) [4 marks]

A *linear unit* from neural networks is a linear model for numeric prediction that is fitted by gradient descent. Explain the differences between the *batch* and *incremental* (or *stochastic*) versions of gradient descent.

B) [4 marks]

Stochastic gradient descent would be expected to deal better with local minima during learning than batch gradient descent – true or false ? Explain your reasoning.

A) [12 marks]

Suppose a single unit has output  $o$  the form:

$$o = w_0 + w_1x_1 + w_1x_1^2 + w_2x_2 + w_2x_2^2 + \cdots + w_nx_n + w_nx_n^2$$

The problem is to learn a set of weights  $w_i$  that minimize squared error. Derive a batch gradient descent training rule for this unit.



### Question 3 ANSWER

A) For a linear unit, batch and stochastic gradient descent differ as follows:

- in batch gradient descent the gradient is computed over *all* the examples in the training set before the weight update is applied
- in stochastic gradient descent the gradient is computed for a *single* example, and then the weight update is applied

B) With batch gradient descent, the direction of the gradient is computed over all the training data, so in some sense this is the true gradient. If the algorithm is located near some local minimum then it will move in the direction of the steepest descent and converge at that minimum. However, in stochastic gradient descent, where some example is selected at random, we would expect that the gradient computed for that example may not be the true gradient, so the algorithm may instead move in a different direction, thus avoiding the local minimum. So our answer is 'true'.

C) The approach to deriving the required training rule pretty much follows the method in the lecture slides. We have a single unit with output  $o$  of the form:

$$o = w_0 + w_1x_1 + w_1x_1^2 + w_2x_2 + w_2x_2^2 + \cdots + w_nx_n + w_nx_n^2$$

Using homogeneous coordinates we can write this as:

$$o = \sum_{i=0}^n w_i(x_i + x_i^2)$$

Assuming the same error function and gradient definition as before (slides 15-18 on the lecture on Neural Learning) we can derive the following:

$$\begin{aligned}\frac{\partial E}{\partial w_i} &= \frac{\partial}{\partial w_i} \frac{1}{2} \sum_d (t_d - o_d)^2 \\ &= \frac{1}{2} \sum_d \frac{\partial}{\partial w_i} (t_d - o_d)^2 \\ &= \frac{1}{2} \sum_d 2(t_d - o_d) \frac{\partial}{\partial w_i} (t_d - o_d) \\ &= \sum_d (t_d - o_d) \frac{\partial}{\partial w_i} (t_d - \sum_{i=0}^n w_i(x_i + x_i^2)) \\ \frac{\partial E}{\partial w_i} &= \sum_d (t_d - o_d)(-x_{i,d} - x_{i,d}^2)\end{aligned}$$

#### Question 4 [20 marks]

##### Ensemble Learning

A) [8 marks]

As model complexity increases from low to high, what effect does this have on:

- 1) Bias ?
- 2) Variance ?
- 3) Predictive accuracy on training data ?
- 4) Predictive accuracy on test data ?

B) [3 marks]

Is decision tree learning relatively stable ? Describe decision tree learning in terms of bias and variance in no more than two sentences.

C) [3 marks]

Is nearest neighbour relatively stable ? Describe nearest neighbour in terms of bias and variance in no more than two sentences.

D) [3 marks]

Bagging reduces bias. True or false ? Give a one sentence explanation of your answer.

E) [3 marks]

Boosting reduces variance. True or false ? Give a one sentence explanation of your answer.

#### Question 4 ANSWER

- A) These questions should be answered in the setting of sampling theory, which means all errors are taken over a large (possibly infinite) set of training sets of the same size drawn at random from the same target distribution, where the error can be broken down into bias and variance. As model complexity increases from low to high:
- 1) Bias will reduce, since the probability of systematic error due to a mismatch between the target model class and the learner's model class will be reduced;
  - 2) Variance will increase, since the amount of error due to variation over the training samples drawn repeatedly from the target distribution will increase due to increased flexibility in fitting any particular set of training samples;
  - 3) Predictive accuracy on training data will be reduced, since in general increasing model complexity will lead to better fitting of the training data;
  - 4) Predictive accuracy on test data will increase, since in general increasing model complexity will lead to overfitting of the training data, unless this is controlled by regularisation.
- B) Decision tree learning is relatively unstable. This is because a decision tree can be grown to completely fit a training set (except for noise in labelling examples). This can be understood as decision tree learning being in general low bias and high variance.
- C) Nearest neighbour tree learning is relatively stable. This is because (depending on the value of  $k$  for the number of nearest neighbours to be used for prediction) small changes in the examples appearing in the training set will not change the predictions that much. This can be understood as nearest neighbour learning being in general high bias and low variance.
- D) Bagging is designed to reduce variance, not reduce bias, since bagging uses a majority vote of base learning methods, which not changed by the ensemble, so bias is not changed.
- E) Boosting is designed to reduce bias, not reduce variance, since multiple instances of a boosted base learning method are combined using an error-based weighted voting scheme on which predictions are made, so most of the error reduction is due to this rather than variance reduction.

### Question 5 [20 Marks]

#### Computational Learning Theory

A) [8 marks]

An instance space  $X$  is defined using  $m$  Boolean attributes. Let the hypothesis space  $H$  be the set of decision trees defined on  $X$  (you can assume two classes). What is the largest set of instances in this setting which is shattered by  $H$  ? [Show your reasoning.]

B) [10 marks]

Suppose we have a consistent learner with a hypothesis space restricted to conjunctions of exactly 8 attributes, each with values  $\{\text{true}, \text{false}, \text{don't care}\}$ . What is the size of this learner's hypothesis space ? Give the formula for the number of examples sufficient to learn with probability at least 95% an approximation of any hypothesis in this space with error of at most 10%. [Note: you are *not* required to compute the solution.]

C) [2 marks]

Informally, which of the following are consequences of the No Free Lunch theorem:

- a) averaged over all possible training sets, the variance of a learning algorithm dominates its bias
- b) averaged over all possible training sets, no learning algorithm has a better off-training set error than any other
- c) averaged over all possible target concepts, the bias of a learning algorithm dominates its variance
- d) averaged over all possible target concepts, no learning algorithm has a better off-training set error than any other
- e) averaged over all possible target concepts and training sets, no learning algorithm is independent of the choice of representation in terms of its classification error

## Question 5 ANSWER

- A) There are  $n = 2^m$  possible instances in the instance space  $X$ . The number of possible dichotomies is the number of possible subsets of the instance space (each subset can be labelled as ‘positive’ and its complement labelled ‘negative’), which is  $2^n$ .

Each subset of instances can be uniquely defined by a conjunction of Boolean attributes, which can be represented in a decision tree by the path from root node to leaf node. So any dichotomy can be represented by a decision tree defined in this way.

So the size of the largest set of instances in this setting that can be shattered by  $H$  is the size of the instance space  $X$ , which is  $n$ .

- B) If there are 8 attributes, each with 3 values, and the hypothesis space  $H$  can be formed by conjunctions of attributes, then the size of hypothesis space  $H$  is  $3^8$ .

Recall the formula for the sample complexity for a consistent learner:

$$m \geq \frac{1}{\epsilon} (\log(|H|) + \log(\frac{1}{\delta}))$$

Here  $\epsilon = 0.1$  and  $\delta = 1 - 0.95 = 0.05$ , so the expression for the number of examples required

$$m \geq \frac{1}{0.1} (\log(3^8) + \log(\frac{1}{0.05}))$$

- C) Answer is d) averaged over all possible target concepts, no learning algorithm has a better off-training set error than any other

### Question 6 [20 Marks]

#### Mistake Bounds

Consider the following learning problem on an instance space which has only one feature, i.e., each instance is a *single integer*. Suppose instances are always in the range  $[1, 5]$ . The hypothesis space is one in which each hypothesis is an interval over the integers. More precisely, each hypothesis  $h$  in the hypothesis space  $H$  is an interval of the form  $a \leq x \leq b$ , where  $a$  and  $b$  are integer constants and  $x$  refers to the instance. For example, the hypothesis  $3 \leq x \leq 5$  classifies the integers 3, 4 and 5 as positive and all others as negative.

Instance	Class
1	Negative
2	Positive
3	Positive
4	Positive
5	Negative

A) [15 marks]

Apply the HALVING ALGORITHM to the five examples in the order in which they appear in the table above. Show each class prediction and whether or not it is a mistake, plus the initial hypothesis set and the hypothesis set at the end of each iteration.

B) [5 marks]

What is the worst-case mistake bound for the HALVING ALGORITHM given the hypothesis space described above ? Give an informal derivation of your bound.

## Question 6 ANSWER

A) Applying the HALVING ALGORITHM with hypotheses of the form  $a \leq x \leq b$  to the data. Start by enumerating the initial hypothesis space<sup>1</sup>.

$1 \leq x \leq 5$     $2 \leq x \leq 5$     $3 \leq x \leq 5$     $4 \leq x \leq 5$     $5 \leq x \leq 5$   
 $1 \leq x \leq 4$     $2 \leq x \leq 4$     $3 \leq x \leq 4$     $4 \leq x \leq 4$   
 $1 \leq x \leq 3$     $2 \leq x \leq 3$     $3 \leq x \leq 3$   
 $1 \leq x \leq 2$     $2 \leq x \leq 2$   
 $1 \leq x \leq 1$

Now we apply the algorithm to each instance in turn. The algorithm works by running the current instance against each hypothesis in the current hypothesis space. If a hypothesis is consistent with the instance it votes 'POSITIVE' otherwise 'NEGATIVE' (here an instance is consistent with the hypothesis if it is in the range defined by the hypothesis). The class with more votes is the prediction. Then the algorithm checks the actual class of the instance to see if it made a mistake, and it eliminates all of hypotheses that misclassified the instance.

Instance	POSITIVE votes	NEGATIVE votes	Prediction	Actual	Mistake ?
1	5	15	NEGATIVE	NEGATIVE	no

Updated hypothesis space:

$2 \leq x \leq 5$     $3 \leq x \leq 5$     $4 \leq x \leq 5$     $5 \leq x \leq 5$   
 $2 \leq x \leq 4$     $3 \leq x \leq 4$     $4 \leq x \leq 4$   
 $2 \leq x \leq 3$     $3 \leq x \leq 3$   
 $2 \leq x \leq 2$

Get next instance:

Instance	POSITIVE votes	NEGATIVE votes	Prediction	Actual	Mistake ?
2	4	6	NEGATIVE	POSITIVE	yes

Updated hypothesis space:

$2 \leq x \leq 5$   
 $2 \leq x \leq 4$   
 $2 \leq x \leq 3$   
 $2 \leq x \leq 2$

---

<sup>1</sup>Note: all hypotheses of the form  $b \leq x \leq a$  where  $a < b$ , such as  $4 \leq x \leq 2$ , have been omitted. Why? These are all equivalent in the sense that they are semantically the same, since they exclude all possible instances. We could have included one of these, but it would not change the results much.

Get next instance:

Instance	POSITIVE votes	NEGATIVE votes	Prediction	Actual	Mistake ?
3	3	1	POSITIVE	POSITIVE	no

Updated hypothesis space:

$$2 \leq x \leq 5$$

$$2 \leq x \leq 4$$

$$2 \leq x \leq 3$$

Get next instance:

Instance	POSITIVE votes	NEGATIVE votes	Prediction	Actual	Mistake ?
4	2	1	POSITIVE	POSITIVE	no

Updated hypothesis space:

$$2 \leq x \leq 5$$

$$2 \leq x \leq 4$$

Get next instance:

Instance	POSITIVE votes	NEGATIVE votes	Prediction	Actual	Mistake ?
5	1	1	POSITIVE	NEGATIVE	yes

The last vote required a tie-break. On this occasion it led to a mistake. Nevertheless the algorithm has converged to the correct hypothesis.

Final hypothesis space:

$$2 \leq x \leq 4$$

B) The worst-case mistake bound for the HALVING ALGORITHM is  $\lfloor \log_2 |H| \rfloor$  where  $H$  is the hypothesis space. Informally, this can be explained as a kind of “binary chop” procedure. For each instance, the algorithm makes a classification based on a majority vote of the hypotheses in the hypothesis space. If the predicted class is the same as the actual class, there is no mistake, otherwise there is. However, in *either* case all hypotheses that predicted incorrectly are eliminated. So on each mistake, at least half of the hypotheses will be eliminated (because of majority voting).

In this case, the algorithm converged with 2 mistakes, less than the worst-case  $\lfloor \log_2(15) \rfloor = 3$ .



END OF PAPER