

S249/S249P/419C

DUBLIN INSTITUTE OF TECHNOLOGY
KEVIN STREET, DUBLIN 8

**BSc. (Honours)
Degree in Information Systems /
Information Technology
(Part-time)**

Stage 4

SUMMER EXAMINATIONS 2016

***** SOLUTIONS *****

ARTIFICIAL INTELLIGENCE II [CMPU4011]

Dr. John Kelleher
Dr. Deirdre. Lillis
Dr. Rem Collier

Monday 9th May 2016
4:00 pm to 6:00 pm

Question 1 is **compulsory**

Answer Question 1 (40 marks) **and**

any 2 Other Questions (30 marks each).

***** SOLUTIONS *****

***** SOLUTIONS *****

SOLUTIONS

1. (a) What is **supervised machine learning**?

(5 marks)

Supervised machine learning techniques automatically learn the relationship between a set of **descriptive features** and a **target feature** from a set of historical **instances**. Supervised machine learning is a subfield of machine learning. Machine learning is defined as an automated process that extracts patterns from data. In predictive data analytics applications, we use **supervised machine learning** to build models that can make predictions based on patterns extracted from historical data.

- (b) Explain what can go wrong when a machine learning classifier uses the wrong **inductive bias**.

(5 marks)

- If the inductive bias of the learning algorithm constrains the search to only consider simple hypotheses we may have excluded the real function from the hypothesis space. In other words, the true function is **unrealizable** in the chosen hypothesis space, (i.e., we are **underfitting**).
- If the inductive bias of the learning algorithm allows the search to consider complex hypotheses, the model may hone in on irrelevant factors in the training set. In other words the model with **overfit** the training data.

- (c) Table 1, on the next page, shows the predictions made for a categorical target feature by a model for a test dataset. Based on this test set, calculate the evaluation measures listed below.

- (i) A **confusion matrix**

(6 marks)

The confusion matrix can be written as		Prediction	
		'true'	'false'
Target	'true'	8	1
	'false'	0	11

- (ii) The **misclassification rate**

(4 marks)

$$\text{misclassification rate} = \frac{(FP + FN)}{(TP + TN + FP + FN)}$$

Misclassification rate can be calculated as

$$\begin{aligned} \text{misclassification rate} &= \frac{(FP + FN)}{(TP + TN + FP + FN)} \\ &= \frac{(0 + 1)}{(8 + 11 + 0 + 1)} \\ &= 0.05 \end{aligned}$$

(iii) The **precision, recall, and F₁ measure**

(12 marks)

$$\begin{aligned} \text{precision} &= \frac{TP}{(TP + FP)} \\ \text{recall} &= \frac{TP}{(TP + FN)} \\ F_1 \text{ measure} &= 2 \times \frac{(\text{precision} \times \text{recall})}{(\text{precision} + \text{recall})} \end{aligned}$$

We can calculate precision and recall as follows (assuming that the 'true' target level is the positive level):

$$\begin{aligned} \text{precision} &= \frac{TP}{(TP + FP)} \\ &= \frac{8}{(8 + 0)} \\ &= 1.000 \\ \text{recall} &= \frac{TP}{(TP + FN)} \\ &= \frac{8}{(8 + 1)} \\ &= 0.889 \end{aligned}$$

Using these figures, we can calculate the F₁ measure as

$$\begin{aligned} F_1 \text{ measure} &= 2 \times \frac{(\text{precision} \times \text{recall})}{(\text{precision} + \text{recall})} \\ &= 2 \times \frac{(1.000 \times 0.889)}{(1.000 + 0.889)} \\ &= 0.941 \end{aligned}$$

(iv) The **average class accuracy (harmonic mean)**. (During this calculation you should round all long floats to 3 places of decimal.)

(8 marks)

$$\text{average class accuracy}_{HM} = \frac{1}{\frac{1}{|levels(t)|} \sum_{l \in levels(t)} \frac{1}{recall_l}}$$

First, we calculate the recall for each target level:

$$recall_{true} = \frac{8}{9} = 0.889$$

$$recall_{false} = \frac{11}{11} = 1.000$$

Then we can calculate a harmonic mean as

$$\begin{aligned} \text{average class accuracy}_{HM} &= \frac{1}{\frac{1}{|levels(t)|} \sum_{l \in levels(t)} \frac{1}{recall_l}} \\ &= \frac{1}{\frac{1}{2} \left(\frac{1}{0.889} + \frac{1}{1} \right)} \\ &= 0.941 \end{aligned}$$

Table 1: The predictions made by a model for a categorical target on a test set of 20 instances

ID	Target	Prediction	ID	Target	Prediction
1	false	false	11	false	false
2	false	false	12	true	true
3	false	false	13	false	false
4	false	false	14	true	true
5	true	true	15	false	false
6	false	false	16	false	false
7	true	true	17	true	false
8	true	true	18	true	true
9	false	false	19	true	true
10	false	false	20	true	true

2. (a) A data analyst building a k -nearest neighbour model for a continuous prediction problem is considering appropriate values to use for k .

- (i) Initially the analyst uses a simple average of the target variables for the k nearest neighbours in order to make a new prediction. After experimenting with values for k in the range $0 - 10$ it occurs to the analyst that they might get very good results if they set k to the total number of instances in the training set. Do you think the analyst is likely to get good results using this value for k ?

(5 marks)

In answering this question students should realise that if the analyst set k to the number of training examples all predictions would essentially be the average target value across the whole dataset. To score very well students should realise that this is an example of massive underfitting.

- (ii) If the analyst was using a distance weighted averaging function rather than a simple average for their predictions would this have made their idea any more useful?

(5 marks)

Students should realise that yes, if distance weighted voting is used (particularly if a $\frac{1}{d^2}$ type distance weight is used) then examples that are far away from the query will have very little impact on the result. Again to score well students should mention that when distance weighted voting is used the value of k in k -NN classifiers is much less important.

- (b) A dataset showing the decisions made by an individual about whether to wait for a table at a restaurant is listed in Table 1 on the next page. (Note that Table 3, also on the next page, lists some equations that you may find useful for this question.)

- (i) Given that the WILLWAIT column lists the values of the target variable, compute the entropy for this dataset.

(5 marks)

There are 6 positive and 6 negative examples in this dataset. This means that the entropy for the dataset is:

$$\begin{aligned} I\left(\frac{6}{12}, \frac{6}{12}\right) &= -\frac{6}{12} \log_2 \frac{6}{12} - \frac{6}{12} \log_2 \frac{6}{12} \\ &= \left(-\frac{1}{2} \log_2 \frac{1}{2}\right) + \left(-\frac{1}{2} \log_2 \frac{1}{2}\right) \\ &= -\frac{1}{2}(-1) + -\frac{1}{2}(-1) \\ &= 1\text{bit} \end{aligned}$$

- (ii) What is the information gain for the PATRONS feature?

(5 marks)

$$\text{Gain}(\text{Patrons}) = 1 - \left(\frac{2}{12} I(0, 1) + \frac{4}{12} I(1, 0) + \frac{6}{12} I\left(\frac{2}{6}, \frac{4}{6}\right)\right) \approx 0.541 \text{ bits}$$

- (iii) What is the information gain for the TYPE feature?

(5 marks)

$$\text{Gain}(\text{Type}) = 1 - \left(\frac{2}{12} I\left(\frac{1}{2}, \frac{1}{2}\right) + \frac{1}{2} I\left(\frac{1}{2}, \frac{1}{2}\right) + \frac{4}{12} I\left(\frac{2}{4}, \frac{2}{4}\right) + \frac{4}{12} I\left(\frac{2}{4}, \frac{2}{4}\right) \right) = 0 \text{ bits}$$

- (iv) Given a choice between the PATRONS and TYPE feature, which feature would the ID3 algorithm choose as the root node for a decision tree?

(5 marks)

The ID3 algorithm would choose the Patrons feature as the root node for the decision tree because it has the higher information gain.

ID	BAR	PATRONS	PRICE	RAIN	TYPE	WILLWAIT
1	F	Some	€€€	F	French	T
2	F	Full	€	F	Thai	F
3	T	Some	€	F	Burger	T
4	F	Full	€	F	Thai	T
5	F	Full	€€€	F	French	F
6	T	Some	€€	T	Italian	T
7	T	None	€	T	Burger	F
8	F	Some	€€	T	Thai	T
9	T	Full	€	T	Burger	F
10	T	Full	€€€	F	Italian	F
11	F	None	€	F	Thai	F
12	T	Full	€	F	Burger	T

Table 2: A dataset describing the previous decisions made by an individual about whether to wait for a table at a restaurant.

Table 3: Equations from information theory.

$$\begin{aligned}
 H(\mathbf{f}, \mathcal{D}) &= - \sum_{l \in \text{levels}(f)} P(f=l) \times \log_2(P(f=l)) \\
 \text{rem}(\mathbf{f}, \mathcal{D}) &= \sum_{l \in \text{levels}(f)} \frac{|\mathcal{D}_{f=l}|}{|\mathcal{D}|} \times H(\mathbf{t}, \mathcal{D}) \\
 IG(\mathbf{d}, \mathcal{D}) &= H(\mathbf{t}, \mathcal{D}) - \text{rem}(\mathbf{d}, \mathcal{D})
 \end{aligned}$$

Table 4: Query Document
Machine learning is fun

Table 5: Document counts from the corpus for the words in the query.

Document counts for the DISLIKE data set				Document counts for the LIKE data set			
fun	is	machine	learning	fun	is	machine	learning
415	695	35	70	200	295	120	105

3. Lets assume we are given a set of **700** training documents that a friend has classified as DISLIKE and another **300** documents that they have classified as LIKE. We are now given a new document and asked to classify it. Table 4 lists the content of this query document and Table 5 gives the number of documents from each class (DISLIKE and LIKE) that the words in the query document occurred in. **What class will a Naive Bayes prediction model label the query document as belonging to?** (You must support your answer by showing the calculations that a Naive Bayes model will make.)

(30 marks)

A Naive Bayes model will label the query with the class that has the highest probability under the assumption of conditional independence between the evidence features. So to answer this question we need to calculate the probability of each class given the evidence and assuming conditional independence across the evidence.

To carry out these calculation we need to convert the raw documents counts into conditional probabilities by dividing each count by the total number of documents occurring in class:

w_k	Count	$P(w_k C = dislike)$
fun	415	$\frac{415}{700} = .593$
is	695	$\frac{695}{700} = .99$
learning	35	$\frac{35}{700} = .05$
machine	70	$\frac{70}{700} = .10$

w_k	Count	$P(w_k C = like)$
fun	200	$\frac{200}{300} = .667$
is	295	$\frac{295}{300} = .983$
learning	120	$\frac{120}{300} = .40$
machine	105	$\frac{105}{300} = .35$

We can now compute the probabilities of each class:

$P(dislike|Query\ Document)$

$$\begin{aligned}
 &= P(dislike) \times P(Machine|dislike) \times P(learning|dislike) \times P(is|dislike) \times p(fun|dislike). \\
 &= 0.7 \times 0.593 \times 0.99 \times 0.5 \times 0.1 \\
 &= 0.00205
 \end{aligned}$$

$P(like|Query\ Document)$

$$\begin{aligned}
 &= P(like) \times P(Machine|like) \times P(learning|like) \times P(is|like) \times p(fun|like). \\
 &= 0.3 \times 0.667 \times 0.983 \times 0.4 \times 0.35 \\
 &= 0.00275
 \end{aligned}$$

As $P(like|Query\ Document) > P(dislike|Query\ Document)$ the naive bayes classifier will return a label of LIKE.

4. (a) A multivariate logistic regression model has been built to predict the propensity of shoppers to perform a repeat purchase of a free gift that they are given. The descriptive features used by the model are the age of the customer, the average amount of money the customer spends on each visit to the shop, and the average number of visits the customer makes to the shop per week. This model is being used by the marketing department to determine who should be given the free gift. The trained model is

$$\begin{aligned}\text{REPEAT PURCHASE} = & -3.82398 - 0.02990 \times \text{AGE} \\ & + 0.74572 \times \text{SHOP FREQUENCY} \\ & + 0.02999 \times \text{SHOP VALUE}\end{aligned}$$

And, the logistic function is defined as:

$$\text{logistic}(x) = \frac{1}{1 + e^{-x}}$$

Assuming that the *yes* level is the positive level and the classification threshold is 0.5, use this model to make predictions for each of the query instances shown in Table 6, on the next page.

(12 marks)

Calculating the predictions made by the model simply involves inserting the descriptive features from each query instance into the prediction model. With this information, the predictions can be made as follows:

$$\begin{aligned}\text{A: } & \text{Logistic}(-3.82398 + -0.0299 \times 37 + 0.74572 \times 0.72 + 0.02999 \times 170.65) \\ & = \text{Logistic}(0.724432) = \frac{1}{1 - e^{-0.724432}} \\ & = 0.673582 \Rightarrow \text{yes}\end{aligned}$$

$$\begin{aligned}\text{B: } & \text{Logistic}(-3.82398 + -0.0299 \times 32 + 0.74572 \times 1.08 + 0.02999 \times 165.39) \\ & = \text{Logistic}(0.984644) = \frac{1}{1 - e^{-0.984644}} \\ & = 0.728029 \Rightarrow \text{yes}\end{aligned}$$

- (b) The effects that can occur when different drugs are taken together can be difficult for doctors to predict. A machine learning has been trained to distinguish between dosages of two drugs that cause a dangerous interaction and those that cause a safe interaction. There are just two continuous features in this dataset, DOSE1 and DOSE2, and two target levels, *dangerous* and *safe*. There is a non-linear decision boundary between dangerous and safe interactions and, consequently, the following set of basis functions were defined:

$$\begin{aligned}\phi_0(\langle \text{DOSE1}, \text{DOSE2} \rangle) &= 1 & \phi_1(\langle \text{DOSE1}, \text{DOSE2} \rangle) &= \text{DOSE1} \\ \phi_2(\langle \text{DOSE1}, \text{DOSE2} \rangle) &= \text{DOSE2} & \phi_3(\langle \text{DOSE1}, \text{DOSE2} \rangle) &= \text{DOSE1}^2 \\ \phi_4(\langle \text{DOSE1}, \text{DOSE2} \rangle) &= \text{DOSE2}^2 & \phi_5(\langle \text{DOSE1}, \text{DOSE2} \rangle) &= \text{DOSE1}^3 \\ \phi_6(\langle \text{DOSE1}, \text{DOSE2} \rangle) &= \text{DOSE2}^3 & \phi_7(\langle \text{DOSE1}, \text{DOSE2} \rangle) &= \text{DOSE1} \times \text{DOSE2}\end{aligned}$$

Training a logistic regression model using this set of basis functions leads to the

following model:

$$\begin{aligned}
 P(\text{TYPE} = \textit{dangerous}) = \\
 \textit{Logistic} & \left(-0.848 \times \phi_0(\langle \text{DOSE1}, \text{DOSE2} \rangle) + 1.545 \times \phi_1(\langle \text{DOSE1}, \text{DOSE2} \rangle) \right. \\
 & - 1.942 \times \phi_2(\langle \text{DOSE1}, \text{DOSE2} \rangle) + 1.973 \times \phi_3(\langle \text{DOSE1}, \text{DOSE2} \rangle) \\
 & + 2.495 \times \phi_4(\langle \text{DOSE1}, \text{DOSE2} \rangle) + 0.104 \times \phi_5(\langle \text{DOSE1}, \text{DOSE2} \rangle) \\
 & \left. + 0.095 \times \phi_6(\langle \text{DOSE1}, \text{DOSE2} \rangle) + 3.009 \times \phi_7(\langle \text{DOSE1}, \text{DOSE2} \rangle) \right)
 \end{aligned}$$

Use this model to make predictions for the query instances in Table 7 and using these prediction explain whether or not the dosage combinations are likely to lead to a dangerous or safe interaction.

(18 marks)

The first step in making a prediction is to generate the outputs of the basis functions. This is done for the first query as follows:

$$\begin{aligned}\phi_0(\langle 0.50, 0.75 \rangle) &= 1 & \phi_4(\langle 0.50, 0.75 \rangle) &= 0.5625 \\ \phi_1(\langle 0.50, 0.75 \rangle) &= 0.50 & \phi_5(\langle 0.50, 0.75 \rangle) &= 0.1250 \\ \phi_2(\langle 0.50, 0.75 \rangle) &= 0.75 & \phi_6(\langle 0.50, 0.75 \rangle) &= 0.4219 \\ \phi_3(\langle 0.50, 0.75 \rangle) &= 0.25 & \phi_7(\langle 0.50, 0.75 \rangle) &= 0.3750\end{aligned}$$

We can now use the regression model to make a prediction:

$$\begin{aligned}P(\text{TYPE} = \textit{dangerous}) &= \textit{Logistic}(-0.848 \times 1 + 1.545 \times 0.50 - 1.942 \times 0.75 + 1.973 \times 0.25 \\ &\quad + 2.495 \times 0.5625 + 0.104 \times 0.1250 + 0.095 \times 0.4219 + 3.009 \times 0.3750) \\ &= \textit{Logistic}(1.5457) \\ &= 0.8243\end{aligned}$$

This means that the probability of the query dosages causing a *dangerous* interaction is 0.8243, so we would say that the result for this query is *dangerous*.

For the next query $\langle 0.10, 0.75 \rangle$:

$$\begin{aligned}\phi_0(\langle 0.10, 0.75 \rangle) &= 1 & \phi_4(\langle 0.10, 0.75 \rangle) &= 0.5625 \\ \phi_1(\langle 0.10, 0.75 \rangle) &= 0.10 & \phi_5(\langle 0.10, 0.75 \rangle) &= 0.0010 \\ \phi_2(\langle 0.10, 0.75 \rangle) &= 0.75 & \phi_6(\langle 0.10, 0.75 \rangle) &= 0.4219 \\ \phi_3(\langle 0.10, 0.75 \rangle) &= 0.01 & \phi_7(\langle 0.10, 0.75 \rangle) &= 0.0750\end{aligned}$$

We can now use the regression model to make a prediction:

$$\begin{aligned}P(\text{TYPE} = \textit{dangerous}) &= \textit{Logistic}(-0.848 \times 1 + 1.545 \times 0.10 - 1.942 \times 0.75 + 1.973 \times 0.01 \\ &\quad + 2.495 \times 0.5625 + 0.104 \times 0.0010 + 0.095 \times 0.4219 + 3.009 \times 0.0750) \\ &= \textit{Logistic}(-0.4613) \\ &= 0.3867\end{aligned}$$

This means that the probability of the query dosages causing a *dangerous* interaction is 0.3867, so we would say that these dosages are *safe* together.

And for the last query $\langle -0.47, 0.18 \rangle$:

$$\begin{aligned}\phi_0(\langle -0.47, 0.18 \rangle) &= 1 & \phi_4(\langle -0.47, 0.18 \rangle) &= 0.0324 \\ \phi_1(\langle -0.47, 0.18 \rangle) &= -0.47 & \phi_5(\langle -0.47, 0.18 \rangle) &= -0.1038 \\ \phi_2(\langle -0.47, 0.18 \rangle) &= 0.18 & \phi_6(\langle -0.47, 0.18 \rangle) &= 0.0058 \\ \phi_3(\langle -0.47, 0.18 \rangle) &= 0.2209 & \phi_7(\langle -0.47, 0.18 \rangle) &= -0.0846\end{aligned}$$

We can now use the regression model to make a prediction:

$$\begin{aligned}P(\text{TYPE} = \textit{dangerous}) &= \textit{Logistic}(-0.848 \times 1 + 1.545 \times -0.47 - 1.942 \times 0.18 + 1.973 \times 0.2209 \\ &\quad + 2.495 \times 0.0324 + 0.104 \times -0.1038 + 0.095 \times 0.0058 + 3.009 \times -0.0846) \\ &= \textit{Logistic}(-1.672106798) \\ &= 0.1581\end{aligned}$$

This means that the probability of the query dosages causing a *dangerous* interaction is 0.1581, so we would say that, instead, this is a *safe* dosage pair.

Table 6: The queries for the multivariate logistic regression question

ID	AGE	SHOP	SHOP
		FREQUENCY	VALUE
A	37	0.72	170.65
B	32	1.08	165.39

Table 7: The query instances for the dosage prediction problem

ID	DOSE1	DOSE2
1	0.50	0.75
2	0.10	0.75
3	-0.47	0.18