

R249/R249P/419C

DUBLIN INSTITUTE OF TECHNOLOGY
KEVIN STREET, DUBLIN 8

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

Stage 4

SUPPLEMENTAL EXAMINATIONS 2016

***** *SOLUTIONS* *****

ARTIFICIAL INTELLIGENCE II [CMPU4011]

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

Duration: 2 Hours

Question 1 is **compulsory**

Answer Question 1 (40 marks) **and**
any 2 Other Questions (30 marks each).

***** SOLUTIONS *****

***** SOLUTIONS *****

SOLUTIONS

1. (a) Explain what is meant by **inductive learning**.

(5 marks)

Inductive Learning involves the process of learning by example where a system tries to induce a general rule from a set of observed instances

- (b) In the context of machine learning, explain what is meant by the term **inductive bias** and illustrate your explanation using examples of inductive biases used by machine learning algorithms.

(15 marks)

- The inductive bias of a learning algorithm:
 - (i) is a set of assumption about what the true function we are trying to model looks like.
 - (ii) defines the set of hypotheses that a learning algorithm considers when it is learning.
 - (iii) guides the learning algorithm to prefer one hypothesis (i.e. the hypothesis that best fits with the assumptions) over the others.
 - (iv) is a necessary prerequisite for learning to happen because inductive learning is an ill posed problem.
- An example of the specific inductive bias introduced by particular machine learning algorithms would be good here. E.g.:
 - Maximum margin: when drawing a boundary between two classes, attempt to maximize the width of the boundary. This is the bias used in Support Vector Machines. The assumption is that distinct classes tend to be separated by wide boundaries.
 - Minimum cross-validation error: when trying to choose among hypotheses, select the hypothesis with the lowest cross-validation error.

- (c) Table 1 shows the predictions made for a categorical target feature by a model for a test dataset.

- (i) Create the **confusion matrix** for the results listed in Table 1.

(5 marks)

		Prediction	
		<i>true</i>	<i>false</i>
Target	<i>true</i>	1	3
	<i>false</i>	2	14

- (ii) Calculate the **classification accuracy** for the results listed in Table 1.

$$\text{classification accuracy} = \frac{(TP + TN)}{(TP + TN + FP + FN)}$$

(5 marks)

Classification accuracy can be calculated as

$$\begin{aligned} \text{classification rate} &= \frac{(TP + TN)}{(TP + TN + FP + FN)} \\ &= \frac{(1 + 14)}{(1 + 14 + 3 + 2)} \\ &= 0.75 \end{aligned}$$

- (iii) Calculate the **average class accuracy (harmonic mean)** for the results listed in Table 1. (During this calculation you should round all long floats to 4 places of decimal.)

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

(10 marks)

Note, in this solution we round all figures to four places of decimal. First, we calculate the recall for each target level:

$$\begin{aligned} recall_{true} &= \frac{1}{4} = 0.25 \\ recall_{false} &= \frac{14}{16} = 0.875 \end{aligned}$$

Then we can calculate a harmonic mean as

$$\begin{aligned} \text{average class accuracy}_{HM} &= \frac{1}{\frac{1}{|levels(t)|} \sum_{l \in levels(t)} \frac{1}{recall_l}} \\ &= \frac{1}{\frac{1}{2} \left(\frac{1}{0.25} + \frac{1}{0.875} \right)} \\ &= \frac{1}{\frac{1}{2} (4 + 1.1429)} \\ &= 0.38889 \end{aligned}$$

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	false	true
3	false	false	13	false	false
4	false	false	14	false	false
5	false	true	15	false	false
6	false	false	16	false	false
7	false	false	17	true	false
8	false	false	18	true	false
9	false	false	19	true	false
10	false	false	20	true	true

2. (a) You are building a recommender system for an large online shop that has a stock of over 100,000 items. In this domain the behaviour of individuals is captured in terms of what items they have bought or not bought.
- (i) Table 2 (below) lists 3 different models of similarity that work on binary data, similar to the data in this domain (**Russell-Rao**, **Sokal-Michener**, and **Jaccard**). Given that there are over 100,000 items available in the store which of these models of similarity (**Russell-Rao**, **Sokal-Michener**, or **Jaccard**) is most appropriate for this domain. Give an explanation for your choice.

(5 marks)

In a domain where there are 100,000's of items co-absences aren't that meaningful. For example, you may be in a domain where there are so many items most people haven't seen, listened to, bought or visited the vast majority of them and as a result the majority of features will be co-absences. The technical term to describe dataset where most of the features have zero values is **sparse data**. In these situations you should use a metric that ignore co-absences and if your features are binary then you should use the **Jaccard similarity** index.

- (ii) Table 3 (on the next page) lists the behaviour of two individuals in this domain for a subset of the items that at least one of the individuals has bought; and, Table 4 (also, on the next page) lists the behaviour of a customer **Q** that you want to generate recommendations for. Assuming that the recommender system uses the similarity metric you selected in Part (i) and that the system will recommend to person **Q** the items that the person most similar to person **Q** has already bought but that person **Q** has not bought, **which item or items will the system recommend to person Q?** Support your answer by showing your calculations and explaining your analysis of the results.

(10 marks)

Using a similarity metric the higher the value returned by the metric the more similar the two items are.

Assuming the student chose the **Jaccard** similarity metric then Person **A** is more similar to **Q** than Person **B**: $Jaccard(Q, A) = \frac{2}{2+1} = 0.6667$, $Jaccard(Q, B) = \frac{1}{4} = 0.25$. As a result the system will recommend item **498**.

If the student selected one of the other similarity metrics for part (a), the supporting calculations should be:

- $Russell-Rao(Q, A) = \frac{2}{5} = 0.4$
- $Russell-Rao(Q, B) = \frac{1}{5} = 0.2$
- $Sokal-Michener(Q, A) = \frac{4}{5} = 0.8$
- $Sokal-Michener(Q, B) = \frac{2}{5} = 0.4$

As is evident from these calculations regardless of which similarity metric is used Person **A** is more similar to **Q** than Person **B**. So the system will recommend item **498** regardless of which similarity metric is used.

Table 2: Similarity Metrics for Binary Data.

Russell-Rao(X,Y)	$= \frac{CP(X,Y)}{P}$
Sokal-Michener(X,Y)	$= \frac{CP(X,Y)+CA(X,Y)}{P}$
Jaccard(X,Y)	$= \frac{CP(X,Y)}{CP(X,Y)+PA(X,Y)+AP(X,Y)}$

- (b) Table 5 on the next page lists a sample of data from a census. There are four descriptive features in this dataset (AGE, EDUCATION, MARITAL STATUS, OCCUPATION) and the target feature ANNUAL INCOME has 3 levels (<25K, 25K–50K, >50K). Note, Table 6, also on the next page, lists some equations that you may find useful for this question.

- (i) Calculate the ENTROPY for this dataset.

(5 marks)

$$\begin{aligned}
 &H(\text{ANNUAL INCOME}, \mathcal{D}) \\
 &= - \sum_{l \in \left\{ \begin{array}{l} <25K, \\ 25K-50K, \\ >50K \end{array} \right\}} P(\text{AN. INC.} = l) \times \log_2(P(\text{AN. INC.} = l)) \\
 &= - \left(\left(\frac{2}{8} \times \log_2 \left(\frac{2}{8} \right) \right) + \left(\frac{5}{8} \times \log_2 \left(\frac{5}{8} \right) \right) + \left(\frac{1}{8} \times \log_2 \left(\frac{1}{8} \right) \right) \right) \\
 &= 1.2988 \text{ bits}
 \end{aligned}$$

- (ii) When building a decision tree, the easiest way to handle a continuous feature is to define a threshold around which splits will be made. What would be the optimal threshold to split the continuous AGE feature (use information gain based on entropy as the feature selection measure)?

(10 marks)

First sort the instances in the dataset according to the AGE feature,

as shown in the following table.

ID	AGE	ANNUAL INCOME
3	18	<25K
6	24	<25K
4	28	25K–50K
5	37	25K–50K
1	39	25K–50K
8	40	>50K
2	50	25K–50K
7	52	25K–50K

Based

on this ordering, the mid-points in the AGE values of instances that are adjacent in the new ordering but that have different target levels define the possible threshold points. These points are 26, 39.5, and 45.

We calculate the information gain for each of these possible threshold points using the entropy value we calculated in part (a) of this question (1.2988 bits) as follows:

Split by Feature	Partition	Instances	Partition Entropy	Rem.	Info. Gain
>26	\mathcal{D}_1	d_3, d_6	0	0.4875	0.8113
	\mathcal{D}_2	$d_1, d_2, d_4, d_5, d_7, d_8$	0.6500		
>39.5	\mathcal{D}_3	d_1, d_3, d_4, d_5, d_6	0.9710	0.9456	0.3532
	\mathcal{D}_4	d_2, d_7, d_8	0.9033		
>45	\mathcal{D}_5	$d_1, d_3, d_4, d_5, d_6, d_8$	1.4591	1.0944	0.2044
	\mathcal{D}_6	d_2, d_7	0		

The threshold AGE > 26 has the highest information gain, and consequently, it is the best threshold to use if we are splitting the dataset using the AGE feature.

Table 3: A dataset showing the behaviour of two individuals in an online shop. A 1 indicates that the person bought the item a 0 indicates that they did not.

Person ID	Item 107	Item 498	Item 7256	Item 28063	Item 75328
A	1	1	1	0	0
B	1	0	0	1	1

Table 4: A query instance from the same domain as the examples listed in Table 3. A 1 indicates that the person bought the item a 0 indicates that they did not.

Person ID	Item 107	Item 498	Item 7256	Item 28063	Item 75328
Q	1	0	1	0	0

Table 5: Census data for the ID3 Algorithm Question

ID	AGE	EDUCATION	MARITAL STATUS	OCCUPATION	ANNUAL INCOME
1	39	bachelors	never married	transport	25K–50K
2	50	bachelors	married	professional	25K–50K
3	18	high school	never married	agriculture	<25K
4	28	bachelors	married	professional	25K–50K
5	37	high school	married	agriculture	25K–50K
6	24	high school	never married	armed forces	<25K
7	52	high school	divorced	transport	25K–50K
8	40	doctorate	married	professional	>50K

Table 6: Equations from information theory.

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

3. Table 7 lists a dataset of books and whether or not they were purchased by an individual (i.e., the feature PURCHASED is the target feature in this domain).

- (a) Calculate the probabilities (to four places of decimal) that a **naive Bayes** classifier would use to represent this domain.

(18 marks)

A naive Bayes classifier would require the prior probability for each level of the target feature and the conditional probability for each level of each descriptive feature given each level of the target feature:

$P(\text{Purchased} = \text{Yes}) = 0.4$	$P(\text{Purchased} = \text{No}) = 0.6$
$P(\text{2ndHand} = \text{True} \text{Purchased} = \text{Yes}) = 0.5$	$P(\text{2ndHand} = \text{True} \text{Purchased} = \text{No}) = 0.5$
$P(\text{2ndHand} = \text{False} \text{Purchased} = \text{Yes}) = 0.5$	$P(\text{2ndHand} = \text{False} \text{Purchased} = \text{No}) = 0.5$
$P(\text{Genre} = \text{Literature} \text{Purchased} = \text{Yes}) = 0.25$	$P(\text{Genre} = \text{Literature} \text{Purchased} = \text{No}) = 0.1667$
$P(\text{Genre} = \text{Romance} \text{Purchased} = \text{Yes}) = 0.5$	$P(\text{Genre} = \text{Romance} \text{Purchased} = \text{No}) = 0.3333$
$P(\text{Genre} = \text{Science} \text{Purchased} = \text{Yes}) = 0.25$	$P(\text{Genre} = \text{Science} \text{Purchased} = \text{No}) = 0.5$
$P(\text{Price} = \text{Cheap} \text{Purchased} = \text{Yes}) = 0.5$	$P(\text{Price} = \text{Cheap} \text{Purchased} = \text{No}) = 0.5$
$P(\text{Price} = \text{Reasonable} \text{Purchased} = \text{Yes}) = 0.25$	$P(\text{Price} = \text{Reasonable} \text{Purchased} = \text{No}) = 0.3333$
$P(\text{Price} = \text{Expensive} \text{Purchased} = \text{Yes}) = 0.25$	$P(\text{Price} = \text{Expensive} \text{Purchased} = \text{No}) = 0.1667$

- (b) Assuming conditional independence between features given the target feature value, calculate the **probability** of each outcome (PURCHASED=Yes, and PURCHASED=No) for the following book (marks will be deducted if workings are not shown, round your results to four places of decimal)

2ND HAND=False, GENRE=Literature, COST=Expensive

(10 marks)

The initial score for each outcome is calculated as follows:

$$(\text{Purchased} = \text{Yes}) = 0.5 \times 0.25 \times 0.25 \times 0.4 = 0.0125$$

$$(\text{Purchased} = \text{No}) = 0.5 \times 0.1667 \times 0.1667 \times 0.6 = 0.0083$$

However, these scores are not probabilities. To get real probabilities we must normalise these scores. The normalisation constant is calculated as follows:

$$\alpha = 0.0125 + 0.0083 = 0.0208$$

The actual probabilities of each outcome is then calculated as:

$$P(\text{Purchased} = \text{Yes}) = \frac{0.0125}{0.0208} = (0.600961...) = 0.6010$$

$$P(\text{Purchased} = \text{No}) = \frac{0.0083}{0.0208} = (0.399038...) = 0.3990$$

- (c) What prediction would a **naive Bayes** classifier return for the above book?

(2 marks)

A naive Bayes classifier returns outcome with the maximum a posteriori probability as its prediction. In this instance the outcome PURCHASED=Yes is the MAP prediction and will be the outcome returned by a naive Bayes model.

Table 7: A dataset describing the a set of books and whether or not they were purchased by an individual.

ID	2ND HAND	GENRE	COST	PURCHASED
1	False	Romance	Expensive	Yes
3	True	Romance	Cheap	Yes
4	False	Science	Cheap	Yes
10	True	Literature	Reasonable	Yes
2	False	Science	Cheap	No
5	False	Science	Expensive	No
6	True	Romance	Reasonable	No
7	True	Literature	Cheap	No
8	False	Romance	Reasonable	No
9	True	Science	Cheap	No

4. (a) A multivariate linear regression model has been built to predict the HEATING LOAD in a residential building based on a set of descriptive features describing the characteristics of the building. Heating load is the amount of heat energy required to keep a building at a specified temperature, usually 65° Fahrenheit, during the winter regardless of outside temperature. The descriptive features used are the overall surface area of the building, the height of the building, the area of the building's roof, and the percentage of wall area in the building that is glazed. This kind of model would be useful to architects or engineers when designing a new building. The trained model is

$$\begin{aligned}\text{HEATING LOAD} = & -26.030 + 0.0497 \times \text{SURFACE AREA} \\ & + 4.942 \times \text{HEIGHT} - 0.090 \times \text{ROOF AREA} \\ & + 20.523 \times \text{GLAZING AREA}\end{aligned}$$

Use this model to make predictions for each of the query instances shown in the Table 8 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.

$$\begin{aligned}1: & -26.030 + 0.0497 \times 784.0 + 4.942 \times 3.5 - 0.090 \times 220.5 + 20.523 \times 0.25 \\ & = 15.5\end{aligned}$$

$$\begin{aligned}2: & -26.030 + 0.0497 \times 710.5 + 4.942 \times 3.0 - 0.09 \times 210.5 + 20.523 \times 0.10 \\ & = 7.2\end{aligned}$$

- (b) 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 9 on the next page.

(18 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:

1: $Logistic(-3.82398 + -0.0299 \times 56 + 0.74572 \times 1.6 + 0.02999 \times 109.32)$
 $= Logistic(-1.02672) = \frac{1}{1 - e^{1.02672}}$
 $= 0.26372 \Rightarrow no$

2: $Logistic(-3.82398 + -0.0299 \times 21 + 0.74572 \times 4.92 + 0.02999 \times 11.28)$
 $= Logistic(-0.44465) = \frac{1}{1 - e^{0.44465}}$
 $= 0.390633 \Rightarrow no$

3: $Logistic(-3.82398 + -0.0299 \times 48 + 0.74572 \times 1.21 + 0.02999 \times 161.19)$
 $= Logistic(0.477229) = \frac{1}{1 - e^{-0.477229}}$
 $= 0.6205 \Rightarrow yes$

Table 8: The queries for the multivariate linear regression HEATING LOAD question

ID	SURFACE	HEIGHT	ROOF	GLAZING
	AREA		AREA	AREA
1	784.0	3.5	220.5	0.25
2	710.5	3.0	210.5	0.10

Table 9: The queries for the multivariate logistic regression question

ID	AGE	SHOP	SHOP
		FREQUENCY	VALUE
1	56	1.60	109.32
2	21	4.92	11.28
3	48	1.21	161.19