

R228/419C

DUBLIN INSTITUTE OF TECHNOLOGY  
KEVIN STREET, DUBLIN 8

---

**BSc. (Honours)**  
**Degree in Computer Science**

**Year 4**

---

**SUPPLEMENTAL EXAMINATIONS 2015**

**\*\*\* SOLUTIONS \*\*\***

---

**ARTIFICIAL INTELLIGENCE II [CMPU4011]**

Dr. John Kelleher  
Dr. Deirdre. Lillis  
Mr. P. Collins

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 4 (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 5 (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 6 (on the next page) lists a data set with of 6 examples described in terms of 3 binary descriptive features (**A**, **B**, and **C**) and a target feature (**Target**). You are asked to create a decision tree model using this data. **Which of the descriptive features will the ID3 decision tree induction algorithm choose as the feature for the root node of the decision tree?** Support you answer with appropriate calculations and dicussions of your results. Note that Table 3 (below) lists some equations that you may find useful for this question.

(15 marks)

Table 3: Equations from information theory.

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

The ID3 decision tree induction algorithm selects the descriptive feature with the highest information gain as the feature for the root node of the decision tree. The first step in calculating information gain is to calculate the entropy for the entire dataset:

$$\begin{aligned} H(DS) &= - \sum_{v \in \{C1, C2\}} p_v \log_2 p_v \\ &= - \left( \frac{3}{6} \log_2 \frac{3}{6} \right) + - \left( \frac{3}{6} \log_2 \frac{3}{6} \right) \\ &= 1.00 \text{ bits} \end{aligned}$$

The table below shows the calculation of the information gain for each of the descriptive features in the dataset:

| Split By | Feature Value | Partition | Examples  | Entropy of Partition | Remainder | Info. Gain |
|----------|---------------|-----------|-----------|----------------------|-----------|------------|
| A        | 1             | $DS_1$    | 1,2,3     | 0.9183               | 0.9183    | 0.0817     |
|          | 0             | $DS_2$    | 4,5,6     | 0.9183               |           |            |
| B        | 1             | $DS_3$    | 2,4,5,6   | 0.8113               | 0.5409    | 0.4591     |
|          | 0             | $DS_4$    | 1,3       | 0                    |           |            |
| C        | 1             | $DS_5$    | 1,2,3,4,6 | 0.9709               | 0.8091    | 0.1909     |
|          | 0             | $DS_6$    | 5         | 0                    |           |            |

From

this table we can see the feature **B** has the highest information gain and consequently the ID3 algorithm will chose this feature as the feature tested at the root node of the tree.



Table 4: 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 5: A query instance from the same domain as the examples listed in Table 4. 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 6: Dataset for the ID3 Algorithm Question

| ID | A | B | C | Target |
|----|---|---|---|--------|
| 1  | 1 | 0 | 1 | C1     |
| 2  | 1 | 1 | 1 | C2     |
| 3  | 1 | 0 | 1 | C1     |
| 4  | 0 | 1 | 1 | C2     |
| 5  | 0 | 1 | 0 | C1     |
| 6  | 0 | 1 | 1 | C2     |

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 restaurant?

(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 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 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. 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 9 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*.

And for the next query  $\langle -0.47, -0.50 \rangle$ :

$$\begin{aligned}\phi_0(\langle -0.47, -0.50 \rangle) &= 1 & \phi_4(\langle -0.47, -0.50 \rangle) &= 0.2500 \\ \phi_1(\langle -0.47, -0.50 \rangle) &= -0.47 & \phi_5(\langle -0.47, -0.50 \rangle) &= -0.1038 \\ \phi_2(\langle -0.47, -0.50 \rangle) &= -0.50 & \phi_6(\langle -0.47, -0.50 \rangle) &= -0.1250 \\ \phi_3(\langle -0.47, -0.50 \rangle) &= 0.2209 & \phi_7(\langle -0.47, -0.50 \rangle) &= 0.2350\end{aligned}$$

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

$$\begin{aligned}P(\text{TYPE} = \textit{dangerous}) &= \textit{Logistic}(-0.848 \times 1 + 1.545 \times -0.47 - 1.942 \times -0.50 + 1.973 \times 0.2209 \\ &\quad + 2.495 \times 0.25 + 0.104 \times -0.1038 + 0.095 \times -0.1250 + 3.009 \times 0.2350) \\ &= \textit{Logistic}(1.1404) \\ &= 0.7577\end{aligned}$$

This means that the probability of the query document causing a *dangerous* interaction is 0.7577, so we would return a *dangerous* prediction.

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:

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

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 8: 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 9: 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  |