

Programme Code: DT249, DT255
Module Code: CMPU4011
CRN: 22564, 32371

TECHNOLOGICAL UNIVERSITY DUBLIN
KEVIN STREET CAMPUS

BSc (Hons) Information Systems/Information
Technology (Part-Time)

BSc (Hons) Information Systems/Information
Technology (Full-Time)

Year 4

SEMESTER 2 EXAMINATIONS 2019/20

Machine Learning for Data Analytics

Dr. Bojan Božić
Dr. Deirdre Lillis
Professor Eleni Mangina

SOLUTIONS

***** SOLUTIONS *****

***** SOLUTIONS *****

SOLUTIONS

1. (a) Why is machine learning an **ill-posed problem**?

(5 marks)

Machine learning is an **ill-posed problem** for two reasons: first, when dealing with large datasets, it is likely there will be **noise** in the data and models consistent with a noisy dataset will make incorrect predictions; and second, for the majority of machine learning problems, the training set represents only a small sample of possible set of instances in the domain.

- (b) What is the **inductive bias** of a machine learning algorithm?

(5 marks)

The **inductive bias** is the set of assumptions that define the model selection criteria of a machine learning algorithm. There are two types of inductive bias that a machine learning algorithm can use: the **restriction bias** constrains the set of models the algorithm will consider during the learning process; and the **preference bias** guides the learning algorithm to prefer certain models over the others.

- (c) 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.

- (d) 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)

		Prediction	
		'true'	'false'
Target	'true'	8	1
	'false'	0	11

- (ii) The **classification accuracy**

(4 marks)

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

Classification accuracy can be calculated as

$$\begin{aligned} \text{misclassification rate} &= \frac{(TP + TN)}{(TP + TN + FP + FN)} \\ &= \frac{(8 + 11)}{(8 + 11 + 0 + 1)} \\ &= 0.95 \end{aligned}$$

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

(15 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

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

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	true	true	12	false	false
3	false	false	13	false	false
4	true	true	14	false	false
5	false	false	15	true	true
6	false	false	16	false	false
7	true	false	17	true	true
8	true	true	18	true	true
9	true	true	19	false	false
10	true	true	20	false	false

2. (a) Table 2, on the next page, lists a dataset containing examples described by two descriptive features, **Feature 1** and **Feature 2**, and labelled with a target class **Target**. Table 3, also on the next page, lists the details of a query for which we want to predict the target label. We have decided to use a **3-Nearest Neighbor** model for this prediction and we will use Euclidean distance as our distance metric:

$$d(x_1, x_2) = \sqrt{\sum_{i=1}^n ((x_1.f_i - x_2.f_i)^2)}$$

- (i) With which target class (**TypeA** or **TypeB**) will our **3-Nearest Neighbor** model label the query? Provide an explanation for your answer.

(8 marks)

The first stage is to calculate the Euclidean distance between each of the examples and the query:

ID	Euclidean Distance
101	60000
102	120000
103	120000
104	180000
105	240000

From this table we can see that the three closest examples to the query are examples 101, 102, and 103. Example 101 has a target label of TypeA and both 102 and 103 have target labels TypeB. Consequently TypeB is the majority label in local model constructed by the 3-Nearest Neighbor classifier for this query instance and the query will be labelled with class TypeB.

- (ii) There is a large variation in range between **Feature 1** and **Feature 2**. To account for this we decide to normalize the data. Compute the normalized versions of Feature 1 and Feature 2 to four decimal places of precision using range normalization

$$x_i.f' = \frac{x_i.f - \min(f)}{\max(f) - \min(f)}$$

(4 marks)

ID	Feature 1	Feature 2	Target
101	0.1667	0.2	TypeA
102	0.0000	0	TypeB
103	0.6667	0.8	TypeB
104	0.8333	0.4	TypeA
105	1.0000	1	TypeB

- (iii) Assuming we use the normalized dataset as input, with which target class (**TypeA** or **TypeB**) will our **3-Nearest Neighbor** model label the query? Provide an explanation for your answer.

(8 marks)

The normalize query instance is:

ID	Feature 1	Feature 2	Target
250	0.2	0.3333	?

The Euclidean distances between the normalized data and normalized query are:

ID	Euclidean Distance
101	0.1667
102	0.3887
103	0.6864
104	0.5385
105	1.0414

From this table we can see that the 3 closest neighbors are: 101, 102 and 104. 101 and 104 are both labelled as class **TypeA**. So **TypeA** is the majority class in the neighborhood and the query will be labelled as belonging to it.

- (b) A dataset showing the decisions made by a professional basketball team on whether to draft college players based on 4 features (1 continuous and 3 categorical) as listed in Table 4 on the next page. (Note that Table 5, also on the next page, lists some equations that you may find useful for this question.)

- (i) Given that the DRAFT 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:

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

- (ii) What the entropy tells us about a dataset? What happens with the distribution of a dataset if the entropy is increased?

(5 marks)

- A. The entropy is a measure of impurity of the elements in a dataset. It can be thought as the uncertainty associated with guessing the result if we are to make a random selection from the set.
- B. If the entropy of the dataset is increased, the dataset becomes less homogeneous. Therefore, the uncertainty associated with the dataset if we are to make a random selection is also increased.

Table 2: Dataset for the 3-Nearest Neighbor question

ID	Feature 1	Feature 2	Target
101	180000	4	TypeA
102	120000	3	TypeB
103	360000	7	TypeB
104	420000	5	TypeA
105	480000	8	TypeB

Table 3: Query instance for the 3-Nearest Neighbor question.

ID	Feature 1	Feature 2	Target
250	240000	4	?

Table 4: A dataset showing the decisions made by a professional basketball team on whether to draft college players.

ID	AGE	SPEED	AGILITY	ABILITY	DRAFT
1	20	1	1	3	<i>F</i>
2	21	2	2	1	<i>F</i>
3	20	2	1	2	<i>F</i>
4	22	2	1	1	<i>F</i>
5	22	4	4	4	<i>T</i>
6	21	5	4	5	<i>T</i>
7	23	5	5	4	<i>T</i>
8	19	4	5	5	<i>T</i>
9	22	5	5	5	<i>T</i>
10	21	1	1	2	<i>F</i>
11	20	5	5	4	<i>T</i>
12	21	3	1	1	<i>F</i>

Table 5: Equations for entropy.

$$H(\mathbf{f}, \mathcal{D}) = - \sum_{l \in \text{levels}(f)} P(f = l) \times \log_2(P(f = l))$$

3. Table 6 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 6: 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 7 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 8 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 7: 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 8: 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