

# 17s1: COMP9417 Machine Learning and Data Mining

---

**Lectures:** Introduction to Machine Learning and Data Mining

**Topic:** Questions from lecture topics

**Last revision:** Fri Mar 3 11:34:12 AEDT 2017

## Introduction

Some questions and exercises from the first course lecture, an “Introduction to Machine Learning and Data Mining”, focusing on reviewing some basic concepts and terminology, and building some intuitions about machine learning.

### Question 1

- a) What is the function that Linear Regression is trying to minimize ?
- b) Under what conditions would the value of this function be zero ?
- c) Can you suggest any other properties of this function ?

### Answer

- a) Linear regression aims to minimize the sum of squared errors<sup>1</sup>. This is defined as  $\sum_{i=1}^N (y_i - \hat{y}_i)^2$  where  $y_i$  is the actual value of the target (dependent) variable and  $\hat{y}_i$  is the value predicted by the learned regression function for example  $i$  in the training set.
- b) If the “line” (hyper-plane in general) passes through all of the values of the output in the training data, i.e.,  $\forall i \ y_i = \hat{y}_i$ .
- c) It is non-negative (because the error term is squared) and has a unique minimum (because the derivatives of the squared error term are linear).

**Question 2** Machine learning has a fair amount of terminology which it is important to get to know.

- a) Why do we need features ?
- b) What is the difference between a “task”, a “model” and a “learning problem” ?
- c) Can different learning algorithms be applied to the same tasks and features ?

### Answer

---

<sup>1</sup>Slides 16 and 17 of the lecture muddy the waters a bit here; the slides talk about Mean Squared Error (MSE), which is what is minimised in the typical setup, but what is shown on the slides is the sum of squared errors, also known as the residual sum of squares (RSS).

- a) Essentially, features are the “interface” between the raw data and the model which is to be learned. For example, you might be given raw data on income which has been split into a number of age ranges and you want to aggregate this into a single number, giving you a new feature, say “expected income”.
- b) A task defines a mapping from input features to output, e.g., mapping from demographic features to income, whereas a model is a specific form of that mapping that can be learned by an algorithm, such as a linear regression equation, thereby defining a learning problem.
- c) Yes, for example, you could use a linear classifier learning algorithm or a decision tree learner for the same classification task using the same features.

**Question 3** Suppose you run a learning algorithm that returns a basic linear classifier using the *homogeneous coordinates* representation on a set of training data and obtain the follow weight vector  $w = (-0.4, 0.3, 0.2)$ . For each of the following examples, what would the classification be using the weight vector  $w$  ?

- a)  $x_1 = (0.9, 1.1)$  ?
- b)  $x_2 = (0.3, 1.2)$  ?
- c)  $x_3 = (0.0, 2.0)$  ?

**Answer**

- a)  $\hat{y} = w \cdot x_1^\circ = (-0.4, 0.3, 0.2) \cdot (1, 0.9, 1.1) = 0.09$ , so the instance  $x_1$  is classified *positive*.
- b)  $\hat{y} = w \cdot x_2^\circ = (-0.4, 0.3, 0.2) \cdot (1, 0.3, 1.2) = -0.07$ , so the instance  $x_2$  is classified *negative*.
- c)  $\hat{y} = w \cdot x_3^\circ = (-0.4, 0.3, 0.2) \cdot (1, 0.0, 2.0) = 0$ , so the instance  $x_3$  is on the threshold. How it should be classified is an implementation issue, since the learned model does not give any clear guidance (although positive classification would be a typical choice).

**Question 4** You want to use a probabilistic model to learn to classify text files as containing either ‘business’ or ‘general’ news articles. To illustrate, we will only consider the presence or absence of two *keywords*, ‘valuation’ and ‘manufacturing’ in the text files. To simplify we assume the two classes are mutually exclusive, i.e., text files can only have one class, either ‘business’ or ‘general’.

Shown in Table 1 are the probabilities of the classes given the presence (1) or absence (0) of the keywords in the text.

Table 2 shows the marginal likelihoods of independently observing each of the keywords given each class.

- a) using the data from Table 1, what two patterns of occurrence of keywords in a text file lead to a prediction of ‘business’ ?
- b) what prediction should be made if we have an occurrence of ‘manufacturing’ but NOT ‘valuation’ in a text file ?

Table 1: Posterior probability distribution of classes given word occurrence (bold font indicates more probable class).

valuation	manufacturing	$P(Y = \text{business} \text{valuation}, \text{manufacturing})$	$P(Y = \text{general} \text{valuation}, \text{manufacturing})$
0	0	0.3	<b>0.7</b>
0	1	0.5	0.5
1	0	<b>0.6</b>	0.4
1	1	0.9	<b>0.1</b>

Table 2: Marginal likelihoods: think of these as probabilities of observing the data items (words) independently of any others, given the repetitive classes.

$Y$	$P(\text{valuation} = 1 Y)$	$P(\text{valuation} = 0 Y)$
business	0.3	0.7
general	0.1	0.9

$Y$	$P(\text{manufacturing} = 1 Y)$	$P(\text{manufacturing} = 0 Y)$
business	0.4	0.6
general	0.2	0.8

- c) suppose we are given a text file to classify, and we know that ‘manufacturing’ occurs in the text file, but we know some words are missing from the file for some reason, and we are uncertain if ‘valuation’ occurred or not. However, we do know that the probability of ‘valuation’ occurring in any text file is 0.05. Compute the probability of each class for the given text file.
- d) using the values from Table 2 compute the likelihood ratios for each of the four possible patterns of occurrence of the keywords.

### Answer

- a) if we see at least one occurrence each of ‘valuation’ and ‘manufacturing’, or just at least one occurrence of ‘valuation’, we should predict ‘business’. Note: these scenarios will give the same classification, but with different probabilities !
- b) both classes are equally probable, so without any further information it is irrational to make a prediction; however, if we are told that one class is more probable *a priori* then we could use that fact to make a prediction by default.
- c) we use *marginalisation* to average over the two possibilities, i.e., that ‘valuation’ did or did not occur. We compute the conditional probabilities for each class given this average evidence. We obtain the probabilities from Table 1. First, the formula to use is:  $P(Y|\text{valuation} = 0, \text{manufacturing})P(\text{valuation} = 0) + P(Y|\text{valuation} = 1, \text{manufacturing})P(\text{valuation} = 1)$  and we evaluate this for **each** of the classes  $Y = \text{business}$  and  $Y = \text{general}$ . For  $Y = \text{business}$  this evaluates to  $(0.5 * 0.95) + (0.9 * 0.05) = 0.52$  and for  $Y = \text{general}$  this evaluates to  $(0.5 * 0.95) + (0.1 * 0.05) = .48$ . Since ‘valuation’ mostly does NOT occur, we see that this is

pretty close to the posterior probabilities of each class (0.5) in the second row of Table 1 when ‘valuation’ is KNOWN not to occur.

- d) we need to multiply together the (independent) marginal likelihood ratios to obtain the overall likelihood ratio, for each instantiation of the two keywords denoting whether the word appears in the document, or not. Letting  $X_1$ ,  $X_2$  stand for the occurrence of the keywords, the formula is  $\frac{P(X_1|Y=\text{business})}{P(X_1|Y=\text{general})} \times \frac{P(X_2|Y=\text{business})}{P(X_2|Y=\text{general})}$ . Expanding this out for each of the combinations of keyword occurrences, this gives:

$$\frac{P(\text{valuation} = 0|Y = \text{business})}{P(\text{valuation} = 0|Y = \text{general})} \times \frac{P(\text{manufacturing} = 0|Y = \text{business})}{P(\text{manufacturing} = 0|Y = \text{general})} = \frac{0.7}{0.9} \frac{0.6}{0.8} = 0.58 \quad (0.3)$$

$$\frac{P(\text{valuation} = 0|Y = \text{business})}{P(\text{valuation} = 0|Y = \text{general})} \times \frac{P(\text{manufacturing} = 1|Y = \text{business})}{P(\text{manufacturing} = 1|Y = \text{general})} = \frac{0.7}{0.9} \frac{0.4}{0.2} = 1.55 \quad (0.5)$$

$$\frac{P(\text{valuation} = 1|Y = \text{business})}{P(\text{valuation} = 1|Y = \text{general})} \times \frac{P(\text{manufacturing} = 0|Y = \text{business})}{P(\text{manufacturing} = 0|Y = \text{general})} = \frac{0.3}{0.1} \frac{0.6}{0.8} = 2.25 \quad (0.6)$$

$$\frac{P(\text{valuation} = 1|Y = \text{business})}{P(\text{valuation} = 1|Y = \text{general})} \times \frac{P(\text{manufacturing} = 1|Y = \text{business})}{P(\text{manufacturing} = 1|Y = \text{general})} = \frac{0.3}{0.1} \frac{0.4}{0.2} = 6.00 \quad (0.9)$$

The decision in all but the first row is: predict ‘business’. This agrees with the decision rule using the posteriors in all but the second row (where both predictions are equally probable).