

FIT2086 Assignment 3

Due Date: 11:59PM, Sunday, 27/10/2019

Introduction

There are total of two questions worth $20 + 16 = 34$ marks in this assignment.

This assignment is worth a total of 20% of your final mark, subject to hurdles and any other matters (e.g., late penalties, special consideration, etc.) as specified in the FIT2086 Unit Guide or elsewhere in the FIT2086 Moodle site (including Faculty of I.T. and Monash University policies).

Students are reminded of the Academic Integrity Awareness Training Tutorial Activity and, in particular, of Monash University's policies on academic integrity. In submitting this assignment, you acknowledge your awareness of Monash University's policies on academic integrity and that work is done and submitted in accordance with these policies.

Submission: No files are to be submitted via e-mail. Correct files are to be submitted to Moodle, as given above. You must submit the following three files:

1. One PDF file containing non-code answers to all the questions that require written answers. This file should also include all your plots.
2. The two required R script files containing R code answers as discussed in Question 1 and 2.

Please read these submission instructions carefully and take care to submit the correct files in the correct places.

Question 1 (20 marks)

In this question you will use R to analyse a dataset. The data is contained in the file `heart.train.ass3.2019.csv`. In this dataset, each observation represents a patient at a hospital that reported showing signs of possible heart disease. The outcome is presence of heart disease (HD), or not, so this is a classification problem. The predictors are summarised in Table 1 (overleaf). We are interested in learning a model that can predict heart disease from these measurements. To answer this question you must:

- Provide an R script containing all the code you used to answer the questions. Please use comments to ensure that the code used to identify each question is clearly identifiable. Call this `fn.sn.Q1.R`, where “fn.sn” is your first name followed by your family name.
- Provide appropriate written answers to the questions, along with any graphs, in a report document.

When answering this question, you must use the `rpart` package that we used in Studio 9. The wrapper function for learning a tree using cross-validation that we used in Studio 9 is contained in the file `wrappers.R`. Don't forget to source this file to get access to the function.

1. Using the techniques you learned in Studio 9, fit a decision tree to the data using the `tree` package. Use cross-validation with 10 folds and 1000 repetitions to select an appropriate size tree. What variables have been used in the best tree? How many leaves (terminal nodes) does the best tree have? **[2 marks]**
2. Plot the tree found by CV, and discuss what it tells you about the relationship between the predictors and heart disease. (*hint: use the `text(cv$best.tree,pretty=12)` function to add appropriate labels to the tree*). **[4 marks]**
3. For classification problems, the `rpart` package only labels the leaves with the most likely class. However, if you examine the tree structure in its textual representation on the console, you can determine the probabilities of having heart disease (see Question 2.3 from Studio 9 as a guide) in each leaf (terminal node). Take a screen-capture of the plot of the tree (don't forget to use the “zoom” button to get a larger image) or save it as an image using the “Export” button in R Studio.

Then, use the information from the textual representation of the tree available at the console and annotate the tree in your favourite image editing software; next to all the leaves in the tree, add text giving the probability of contracting heart disease. Include this annotated image in your report file. **[2 marks]**
4. According to your tree, which predictor combination results in the highest probability of having heart-disease? **[1 mark]**
5. We will also fit a logistic regression model to the data. Use the `glm()` function to fit a logistic regression model to the heart data, and use stepwise selection with the BIC score to prune the model. What variables does the final model include, and how do they compare with the variables used by the tree estimated by CV? Which predictor is the most important in the logistic regression? **[3 marks]**
6. Write down the regression equation for the logistic regression model you found using step-wise selection. **[1 mark]**

7. Using the `my.pred.stats()` function contained in the file `studio11.prediction.stats()`, compute the prediction statistics for both the tree and the step-wise logistic regression model on the heart data. Contrast and compare the two models in terms of the various prediction statistics? Would one potentially be preferable to the other as a diagnostic test? Justify your answer. [2 marks]
8. The file `heart.test.ass3.2019.csv` contains the data on a further $n' = 66$ individuals. Calculate the *odds* of having heart disease for the patient in the 45th row of this new dataset. The odds should be calculated for both:
 - (a) the tree model found using cross-validation; and
 - (b) the step-wise logistic regression model.

How do the predicted odds for the two models compare? [2 marks]

9. For the logistic regression model using the predictors selected by BIC in Question 1.6, use the bootstrap procedure (use at least 5,000 bootstrap replications) to find a confidence interval for the probability of having heart disease for patient in the 45th row in the test data. Use the `bca` option when computing this confidence interval. Discuss this confidence interval in comparison to the predicted probabilities of having heart disease for both the logistic regression model and the tree model. [2 marks]
10. Finally, use the bootstrap to compute a 95% confidence interval for the classification accuracy of the logistic regression model using the predictors selected by BIC. Plot the bootstrap results. Comment on the histogram of the bootstrapped classification accuracies. Compare the bootstrap 95% confidence interval for the classification accuracy against the actual classification accuracy you obtain on the testing data `heart.test.ass3.2019.csv` using the model you learned in Q1.6. [1 marks]

Variable name	Description	Values
AGE	Age of patient in years	29 – 77
SEX	Sex of patient	M = Male F = Female
CP	Chest pain type	Typical = Typical angina Atypical = Atypical angina NonAnginal = Non anginal pain Asymptomatic = Asymptomatic pain
TRESTBPS	Resting blood pressure (in <i>mmHg</i>)	94 – 200
CHOL	Serum cholesterol in <i>mg/dl</i>	126 – 564
FBS	Fasting blood sugar > 120 <i>mg/dl</i> ?	<120 = No >120 = Yes
RESTECG	Resting electrocardiographic results	Normal = Normal ST.T.Wave = ST wave abnormality Hypertrophy = showing probable hypertrophy
THALACH	Maximum heart rate achieved	71 – 202
EXANG	Exercise induced angina?	N = No Y = Yes
OLDPEAK	Exercise induced ST depression relative to rest	0 – 6.2
SLOPE	Slope of the peak exercise ST segment	Up = Up-sloping Flat = Flat Down = Down-sloping
CA	Number of major vessels colored by flourosopy	0 – 3
THAL	Thallium scanning results	Normal = Normal Fixed.Defect = Fixed fluid transfer defect Reversible.Defect = Reversible fluid transfer defect
HD	Presence of heart disease	N = No Y = Yes

Table 1: Heart Disease Data Dictionary. ST depression refers to a particular type of feature in an electrocardiograph (ECG) signal during periods of exercise. Thallium scanning refers to the use of radioactive Thallium to check the fluid transfer capability of the heart.

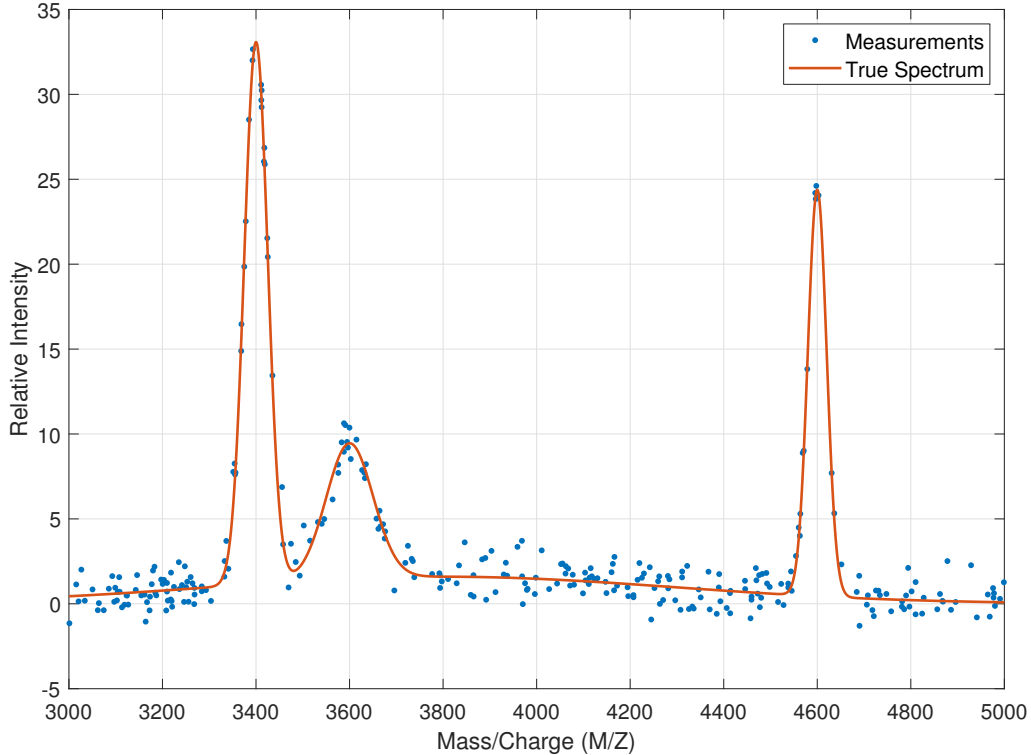


Figure 1: Noisy measurements from a (simulated) mass spectrometry reading. The “true” (unknown) measurements are shown in orange, and the noisy measurements are shown in blue.

Question 2 (16 marks)

Data Smoothing

Data “smoothing” is a very common problem in data science and statistics. We are often interested in examining the unknown relationship between a dependent variable (y) and an independent variable (x), under the assumption that the dependent variable has been imperfectly measured and has been contaminated by measurement noise. The model of reality that we use is

$$y = f(x) + \varepsilon$$

where $f(x)$ is some unknown, “true”, potentially non-linear function of x , and $\varepsilon \sim N(0, \sigma^2)$ is a random disturbance or error. This is called the problem of function estimation, and the process of estimating $f(x)$ from the noisy measurements y is sometimes called “smoothing the data” (even if the resulting curve is not “smooth” in a traditional sense, it is less rough than the original data). In this question you will use several tools to try and estimate an underlying function in a simulated, but realistic example of a data science problem.

Mass Spectrometry Data Smoothing

The file `ms.train.csv` contains $n = 300$ (simulated) measurements from a mass spectrometer. Mass spectrometry is a chemical analysis tool that provides a measure of the physical composition of a material. The outputs of a mass spectrometry reading are the intensities of various ions, indexed by their mass-to-charge ratio. The resulting spectrum usually consists of a number of relatively sharp peaks that indicate a concentration of particular ions, along with an overall background level. A standard problem is that the measurement process is generally affected by noise – that is, the sensor readings are imprecise and corrupted by measurement noise. Therefore, smoothing, or removing the noise is crucial as it allows us to get a more accurate idea of the true spectrum, as well as determine the relative quantity of the ions more accurately. However, we would ideally like for our smoothing procedure to not damage the important information contained in the spectrum (i.e., the heights of the peaks).

The file `ms.train.csv` contains the simulated measurements of our mass spectrometry reading; `ms.train$MZ` are the mass-to-charge ratios of various ions, and `ms.train$intensity` are the measured (noisy) intensities of these ions in our material. The file `ms.test.csv` contains $n = 10,000$ different values of MZ along with the “true” intensity values, stored in `ms.test$intensity`. These true values would be unknown in a real-life situation, but can be used here to see how close your estimated spectrum is to the truth. For reference, the samples `ms.train$intensity` and the value of the true spectrum `ms.test$intensity` are plotted in Figure 1 against their respective MZ values.

To answer this question you must:

- Provide an R script containing all the code you used to answer the questions. Please use comments to ensure that the code used to identify each question is clearly identifiable. Call this file `fn.sn.Q2.R`, where “fn.sn” is your first name followed by your family name.
- Provide appropriate written answers to the questions, along with any graphs, in a non-handwritten report document.

To answer this question, you must use the `kkn` and `rpart` packages that we used in Studio 9.

Questions

1. Use the k -nearest neighbours method (k -NN) to estimate the underlying spectrum from the training data. Use the `kkn` package we examined in Studio 9 to provide predictions for the MZ values in `ms.test`, using `ms.train` as the training data. You should use the `kernel = "optimal"` option when calling the `kkn()` function. This means that the predictions are formed by a weighted average of the k points nearest to the point we are trying to predict, the weights being determined by how far away the neighbours are from the point we are trying to predict.
 - (a) For each value of $k = 1, \dots, 25$, use k -NN to estimate the values of the spectrum for the MZ values in `ms.test$MZ`. Then, compute the mean-squared error between your estimates of the spectrum, and the true values in `ms.test$intensity`. Produce a plot of these errors against the various values of k . [1 mark]
 - (b) Produce four graphs, each one showing: (i) the training data points (`ms.train$intensity`), (ii) the true spectrum (`ms.test$intensity`) and (iii) the estimated spectrum (predicted `intensity` values for the MZ values in `ms.test.csv`) produced by the k -NN method for four different values of k ; do this for $k = 2$, $k = 5$, $k = 10$ and $k = 25$. Make sure the graphs have clearly labelled axes and a clear legend. Use a different colour for your estimated curve. [4 marks]

- (c) Discuss, qualitatively, and quantitatively (in terms of mean-squared error on the true spectrum) the four different estimates of the spectrum. **[2 marks]**
2. Use the cross-validation functionality in the `knn` package to select an estimate of the best value of k (make sure you still use the `optimal` kernel). What value of k does the method select? How does it compare to the (in practice, unknown) value of k that would minimise the actual mean-squared error (as computed in Q2.1a)? **[1 mark]**
 3. Using the estimates of the curve produced in the previous question, see if you can provide an estimate of the variance of the sensor/measurement noise that has corrupted our intensity measurements. **[1 mark]**
 4. Do any of the estimated spectra achieve our aim of providing a smooth, low-noise estimate of background level as well as accurate estimation of the peaks? Explain why you think the k -NN method is able to achieve, or not achieve, this aim. **[1 marks]**.
 5. We can also use a decision tree to smooth our data. Learn a decision tree on `ms.train` to predict `intensity` given `MZ` using CV. Produce a plot, as above, with the predicted intensity values for the `MZ` values in `ms.test.csv` using the final tree found by CV. Compare this qualitatively, and quantitatively (in terms of mean-squared error on the true spectrum) to the k -NN method. **[2 marks]**
 6. Compare and contrast the decision tree approach with the k -NN approach on this one dimensional problem. Thinking about the way in which both methods work, discuss how they are similar and they are different. **[1 marks]**
 7. Apply the bootstrap procedure to the spectrum estimated by k -NN using the value of k found in Q2.2. Using the outputs of the bootstrap function `boot()`, plot an approximate 95% confidence interval for the estimated spectrum for the `MZ` values in `ms.test` along with the true spectrum in `ms.test` and the average of the bootstrap estimates of the spectra (to keep things clean do not include the measurements on this plot). Run it for at least 100 bootstrap iterations. Discuss the confidence interval in comparison to the true spectrum. (*hint: the `$t` field of the object returned by `boot()` gives you access to the bootstrap estimates for each bootstrap iteration as a matrix*). **[3 marks]**