RMIT

Computer Science & IT, School of Science COSC 2670 — Practical Data Science

Assignment 2: Data Modelling and Presentation

Due: 23:59, Tuesday 22 May (week 12)

This assignment is worth 35% of your overall mark.

Assignment Teams

This assignment should be carried out in groups of two.

It is up to you to form a team.

Once you have formed your team, you should register using the form at:

https://goo.gl/forms/KQixgtipMcTX5Wxa2.

Important: you must register your team by 15 May at the latest (one week before the assignment is due).

If you are a postgraduate student and have strong reasons for needing to complete the assignment individually, you may apply to do so by sending an email to the lecturer, explaining your reasons. However, bear in mind that the requirements and available marks will be the same as for pairs, so you are strongly advised to work in a team.

Introduction

This assignment focuses on *data modelling*, a core step in the data science process. You will need to develop and implement appropriate steps, in IPython, to complete the corresponding tasks.

This assignment is intended to give you practical experience with the typical 5th and 6th steps of the data science process: data modelling, and presentation and automation.

The "Practical Data Science" Canvas contains further announcements and a discussion board for this assignment. Please be sure to check these on a regular basis – it is your responsibility to stay informed with regards to any announcements or changes. Login through https://rmit.instructure.com/.

Where to Develop Your Code: Jupyter Notebook: Anaconda (5.1.0) with Python 2.7

The testing environment for this assignment is Anaconda (5.1.0) with Python 2.7.

You can download it from

https://www.anaconda.com/download/#windows

on your own computer;

Or you can use the same version installed on Lab computers.

For Lab Computers, you can find Jupyter Notebook via:

 $Start \rightarrow All \ Programs \rightarrow Anaconda2 \ (64-bit) \rightarrow Jupyter \ Notebook$

Then,

- Select New \rightarrow Python 2
- The new created '*.ipynd' is created at the following location:
 - C:\Users\sXXXXXXX
 - where sXXXXXXX should be replaced with a string consisting of the letter "s" followed by your student number.

Plagiarism

RMIT University takes plagiarism very seriously. All assignments will be checked with plagiarism-detection software; any student found to have plagiarised will be subject to disciplinary action as described in the course guide. Plagiarism includes submitting code that is not your own or submitting text that is not your own. Allowing others to copy your work is also plagiarism. All plagiarism will be penalised; there are no exceptions and no excuses. For further information, please see the *Academic Integrity* information at http://www1.rmit.edu.au/academicintegrity.

General Requirements

This section contains information about the general requirements that your assignment must meet. Please read all requirements carefully before you start.

• You must do all modelling in IPython.

- You must include a plain text file called "readme.txt" with your submission. This file should include your name(s) (if you are a group of two) and student ID(s), and instructions for how to execute your submitted script files. This is important as automation is part of the 6th step of data science process, and will be assessed strictly.
- Parts of this assignment will include a written report, this must be in PDF format.
- Please ensure that your submission follows the file naming rules specified in the tasks below. File names are case sensitive, i.e. if it is specified that the file name is gryphon, then that is exactly the file name you should submit; Gryphon, GRYPHON, griffin, and anything else but gryphon will be rejected.

Task 1: Data Retrieving (3%)

This assignment will focus on data modelling, and you can choose to focus on one approach: Classification or Clustering.

For this assignment, you need to select **one** suitable dataset, from the following options:

- 1. Find and then analyse your own data set, in a domain that is of interest to you. If you choose this option, you will need to:
 - include a detailed description of the data, and each attribute of it, including the type, the range of possible values, whether it contains any missing values/errors
 - submit a copy of the dataset, to allow the assessment of your modelling result.
- 2. Select one data set from the UCI Repository: http://archive.ics.uci.edu/ml/. Choose one dataset from either the *Classification* or *Clustering* task.

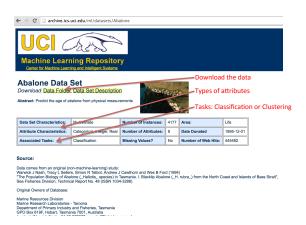


Figure 1: Example of the Abalone data set with instructions: where to download the data, the types of the attributes in the data, and the suitable task(s) of this data.

Being a careful data scientist, you know that it is vital to thoroughly check any available data before starting to analyse it. Please ensure you understand the data you selected, including the meaning of each attribute. For datasets from the UCI repository, you can obtain this information from the corresponding web page under the sections *Data Set Information* and *Attribute Information*.

Task 2: Data Exploration (5%)

Explore the selected data, carrying out the following tasks:

- Explore each column, using appropriate descriptive statistics and graphs (if appropriate), e.g. the distribution of a numerical attribute, the proportion of each value of a categorical attribute. Format each graph carefully, and use it in your final report. You need to include appropriate labels on the x-axis and y-axis, a title, and a legend. The fonts should be sized for good readability. Components of the graphs should be coloured appropriately, if applicable.
- Explore the relationship between all pairs of attributes (or at least 10 pairs of attributes, if there are more in the data), and show their relationship in an appropriate graph.

Task 2: Data Modelling (10%)

Model the data by treating it as **either** a *Classification* or *Clustering* Task, depending on which dataset you previously selected.

You must choose **two** models within the particular Task category (i.e. two Classification models, or two Clustering models), and carry out the following steps for *each* model:

- Select the appropriate model (e.g. DecisionTree for classification) from sklearn.
- If you choose to do a Classification Task,
 - Split the data into training set and the test set.
 - Train the model by selecting appropriate values for each parameter in the model.
 - * You need to show how do you choose this value, and justify why you choose it (for example, k in the KNearestNeighbor model).
 - Test the accuracy of the model on the *test* set, and report the performance of the model in the following terms:
 - * Confusion Matrix
 - * Classification Error Rate

- * Precision
- * Recall
- * F1-Score
- If you choose to do a *Clustering* Task,
 - Train the model by selecting appropriate values for each parameter in the model.
 - * Show how do you choose this value, and justify why you choose it (for example, k in the k-means model).
 - Evaluate the performance of the clustering model by:
 - * Checking the clustering results against the true observation labels
 - * Constructing a "confusion matrix" to analyse the meaning of each cluster by looking at the majority of observations in the cluster. (You can do this by using a pen and a piece of paper, as we did in Practical Exercise 3 in Tute/Lab 06 (week7); if you prefer, you can also explore how to do this step directly in IPython.)

After you have built two Classification models, or two Clustering models, on your data, the next step is to *compare* the models. You need to include the results of this comparison, including a recommendation of which model should be used, in your report (see next section).

Task 3: Report (12%)

Your report should be saved in a file called report.pdf, and must be in PDF format. Remember to clearly cite any sources (including books, research papers, course notes, etc.) that you referred to while designing aspects of your programs.

Your report must have the following structure:

- A cover page, including
 - Title
 - Author (your name(s))
 - Affiliations
 - Contact details
 - Date of report
- Table of Content
- An abstract/executive summary
- Introduction

- Methodology
- Results
- Discussion
- Conclusion
- References

Please revisit p. 51-57 in the Week1 lecture slides.

Task 4: Presentation (5%)

You will be required to do a presentation for your assignment 2 in Week 12's Tute/Prac:

- The presentation should
 - briefly describe your chosen data set,
 - state the hypotheses/questions that you were investigating,
 - then explain what the analysis and results were.
- The presentations are a maximum of 3 minutes per group, and we suggest each group to have at most 3 slides, and print them out on a4 paper, to put on the document camera for presentation (to save time connecting computers between presentations).

If you have your teammates are in different Tute/Prac sessions, you can choose to attend one of the sessions and do the presentation together. But, if you prefer to do the presentation separately in each of your sessions, which is also acceptable.

Please note: the students in 19:30-21:30 slot will present in the lecture, and the rest present in their own tute/pracs.

Optional Extension: Additional Model (Up to 3.5% bonus marks)

ONLY attempt this section if you have completed all previous sections of the assignment.

In the previous sections you chose to analyse your data using either two Clustering or two Classification models. In this bonus section, you should learn about and then apply a *third* model within the area same grouping.

Include your analysis of the data using this third model, and also compare the performance of this model with your previous two. Explain your results in your report. Also include your additional IPython code as part of your assignment submission.

What to Submit, When, and How

The assignment is due at

23:59, Tuesday 22 May (week 12).

Assignments submitted after this time will be subject to standard late submission penalties. There are three files you need to submit:

- Notebook file containing your python commands, 'Assignment2.ipynb'.
- # For the notebook files, please make sure to clean them and remove any unnecessary lines of code (cells). Follow these steps before submission:
 - 1. Main menu \rightarrow Kernel \rightarrow Restart & Run All
 - 2. Wait till you see the output displayed properly. You should see all the data printed and graphs displayed.
- Your report.pdf file.
- Your presentation slides.

They must be submitted as ONE single zip file, named as your student numbers (for example, 1234567_7321283.zip if your student ID are s1234567 and s7321283). The zip file must be submitted in Canvas:

Assignments/Assignment 2.

Please do NOT submit other unnecessary files.