

CITS4009 Project — Exploratory Data Analytics and Predictive Modelling

Worth 25% of the total assessment

Due: Friday, October 11th, 2024, 11:59pm

1. Project Preamble

This is an individual effort project. In this project, you are to demonstrate your understanding of how exploratory data analysis and modelling is carried out using R functions and visualisation tools

The submission should be an [HTML generated from your R Markdown file and a Shiny App](#).

2. Data

To demonstrate a full data science life cycle while fulfilling the educational goals, we would like to provide the following options for sourcing data for the project.

- A specified dataset
- Data from public repositories

Specified data

If you do not have any datasets or domain of interests in mind, we suggest you to use the **Countries and Death causes dataset** (`./Countries and death causes.csv`) which is collected from World Health Organisation (WHO) data. Monitoring the yearly number of deaths helps to address their causes and adapt health systems to react effectively, triggering responses of multiple sectors. Understanding the reasons why people die, can help comprehend the ways people live to improve health services and reduce preventable deaths in every country.

Data from public repositories

If you have specific domain of interests, for example, energy consumption, health sciences, transportation, sales, or sports, etc., you can browse some public dataset repositories to find a dataset that interests you. Note that the data needs to be in tabular form, i.e. not multi-media data, as we won't be able to deal with texts, images, videos, or audios. Also, the dataset should not be too simple. It should have continuous and discrete variables (columns) of various types (numerical, logical, character, etc). Your chosen dataset should have a similar level of complexity as the YouTube dataset.

A few well known public data repos are:

- <http://kaggle.com> (<http://kaggle.com/>)
- <https://data.gov/> (<https://data.gov/>)
- <https://data.gov.au/> (<https://data.gov.au/>)
- <https://data.world/> (<https://data.world/>)
- <https://archive.ics.uci.edu/ml/datasets.html> (<https://archive.ics.uci.edu/ml/datasets.html>)

You need to contact the unit coordinator as it is important to check if all learning outcomes can be demonstrated by your chosen dataset.

3. Exploratory Study of the Data

As a first step for this project, you are required to present the data in two different ways to stakeholders.

1. You are expected to produce an **HTML** file from your **R notebook** (a `.Rmd` file) that documents both the process and the R code that you used for your exploratory data analysis.
2. A **Shiny App** is also expected to coordinate the various plots with explanations to support interactive exploration of the data.

Using R functions to explore the data

Use R functions such as `str()`, `summary()` and `head()` to have a glance at the data. Document your interpretation of the data in the notebook.

Visualisation

Generate different types of plots and charts from the dataset for both single and multiple variable data exploration. Document any intuitions and observations you have in the notebook.

Data cleaning and transformation

Is there any missing values and data anomalies? Do you think it is important to do any data transformation? If so, document these in the notebook. Use visualisation to help justify the cleaning and transformation.

4. Modelling

In the second step, you need to work on modelling.

Classification

Firstly, **study your dataset and choose the response** (i.e., target) variable **suitable for a classification task**. The remaining variables can be your feature variables. You may need to **discard some character string and categorical columns**. Columns that have a unique value for each row should be discarded. If possible, **formulate it as a binary classification problem**, as multi-class classification is difficult and is not covered in the lectures.

Your next step is to split the data into a training set and a test set. You can use any meaningful split ratio (**90/10, 80/20, etc.**). You should implement R code for a **Decision Tree Classifier** and **choose one different classification techniques to compare with** (e.g., **logistic regression classifier, Naïve Bayes classifier, K nearest neighbours classifier**, etc.) and compare the performance of the two models. Your report should include discussions about attribute and feature selection and their impacts on the models. For example, you may have two attribute selection techniques, then you will build two models for each classifiers.

For example, using the provided dataset, we can predict the death rate/toll (high/low) for each country. You may make your own assumptions of high or low death rate. For example, death rate against the social economic status of the country (e.g. normalised using GDP, or GDP per capita).

The downloadable data from Word Bank is available at: <https://data.worldbank.org/indicator/SP.POP.TOTL>

([https://aus01.safelinks.protection.outlook.com/?](https://aus01.safelinks.protection.outlook.com/?url=https%3A%2F%2Fdata.worldbank.org%2Findicator%2FSP.POP.TOTL&data=05%7C02%7Cghulam.hassan%40uwa.edu.au%7Cc20b22f931c04562a9fb0&https://data.worldbank.org/indicator/NY.GDP.PCAP.CD)

[url=https%3A%2F%2Fdata.worldbank.org%2Findicator%2FSP.POP.TOTL&data=05%7C02%7Cghulam.hassan%40uwa.edu.au%7Cc20b22f931c04562a9fb0](https://aus01.safelinks.protection.outlook.com/?url=https%3A%2F%2Fdata.worldbank.org%2Findicator%2FSP.POP.TOTL&data=05%7C02%7Cghulam.hassan%40uwa.edu.au%7Cc20b22f931c04562a9fb0&https://data.worldbank.org/indicator/NY.GDP.PCAP.CD) and <https://data.worldbank.org/indicator/NY.GDP.PCAP.CD> (<https://data.worldbank.org/indicator/NY.GDP.PCAP.CD>)

Note: depending on what predication question you'd like to be answered, you may need to tidy the data into the right shape or merging multiple datasets. You can start with **single-variate models and then multi-variate models**, following the same process demonstrated in the lecture slides. **Highly correlated columns and semantic related columns will need to be removed.**

A good demonstration of the following investigations is expected:

1. Understanding of what a null model would look like in this context.
2. Aggregating, sub-setting, sampling or reshaping the data for better data preparation if necessary.
3. Transforming the data for single-variate model selection, if needed.
4. Using various measures to select a good combination of variables for multi-variate models.
5. Using **LIME** (see Labs) to find the determining feature(s) for the classification of a few instances from the test set. Discuss about your experimental findings in the report.
6. Evaluation of models. This step involves comparing your multi-variate models for the two implemented machine learning techniques on the training set and the test set using various measures (such as ROC plots, confusion matrices, deviances, etc).
7. Different techniques to generate combinations of variables. For example, information gain, principal component analysis, and forward selection. For bonus marks, you can also implement an attribute combination technique from reviewing literature.

Clustering

Choose or compute a set of feature variables; apply a clustering algorithm to these variables; visualise the clustering results and explain the rules discovered.

- Explain how and why you choose the distance measures and how the choices affect your clustering outcome.
- An investigation on the selection of k – the number of clusters.

5. Marking Criteria

- **Data preparation and exploration (10%):** Proficient use of data handling functions (e.g. pipes, merge/joins) or packages to construct clean and tidy training and testing datasets for classification. Sensible aggregation, transformation to obtain the right data for modelling, and good handling of missing or invalid data.
- **Classification (40%):** Good and thorough comparison of Decision Tree Classifiers against a different type of classification models with sensible interpretations of the performance measures, contextually with the right domain intuition. Good comparison of attribute or feature selection strategies.
- **Clustering (20%):** Good implementation of clustering with adequate investigations, demonstrations and explanation on the effect of hyper-parameters (e.g., k for the number of clusters) in these unsupervised techniques.
- **Shiny App (20%):** Sensible UI design to allow users to interactively viewing single variable model performance, two classification model performance and the clustering results.
- **Overall Report Quality (10%):** The report quality and the professionalism of video presentation will be considered. Treat this as presenting your final product to stakeholders and your client.
- To obtain marks in the HD range (80%-100%), you should aim at having the following elements in your project:
 - Well demonstrated exceptional understandings of the principles for model comparison and selection (NOTE: simply building a few more models will not automatically warrant marks in the HD range).
 - Thorough and effective treatment of data to improve model performance, e.g., imbalanced dataset treatment.
 - Exceptional data provenance for reproducible data science.
 - Exceptional understanding of feature variable treatment and selection.
 - Appropriate use of short text descriptions in diagrams (e.g., annotations).

6. Submission

1. Store the **Shiny App in an app.R** file;
2. Create a **short 2-3 minutes video demonstrating the Shiny App**, and provide a **Youtube link to the video in the .Rmd file**. Please show your face in your presentation. You will also need to use your own voice to explain the best features of your App, no AI voices accepted.
3. Generate an html file with the name **project.html** from your **.Rmd file**

4. Zip up the Shiny App file and .html file and submit it to LMS.

You should keep a record of your progressive work towards the final submission and also a copy of your latest work. You are encouraged to submit to LMS as often as you like (or need). The latest submission will overwrite the previous version.

If you like version controlling your work, then you can keep your working copies on GitHub (<https://github.com/>) before the final submission; however, do make sure you keep your GitHub repo **private**.

Submission Check List:

- Make sure that you have your name and student number clearly written at the beginning of your R markdown file. To make it easier for us in the marking process,
 - please ensure that you **have your student number written correctly**.
 - please ensure that you **use exactly the same name as shown on LMS**.
 - please put your **surname in uppercase letters**, e.g., Michael CHEN, John SMITH, Xiaolian HUANG.
- Submit the generated **.html** file (**not the .Rmd file or the .nb.html file**) and the Shiny App file. Before submission, check that you can view your **.html** file in a web browser outside Rstudio and that no diagrams are missing.

7. Special Considerations

Please make yourself familiar with the UWA special consideration process. All special considerations are handled centrally at UWA. Unit Coordinators are no longer processing individual request of project extensions. Please fill in the forms if needed:

<https://www.uwa.edu.au/students/My-course/Exams-assessments-and-results/Special-consideration> (<https://www.uwa.edu.au/students/My-course/Exams-assessments-and-results/Special-consideration>)

8. Penalty on Late Submissions

See the URL below about late submissions of assignments: https://ipoint.uwa.edu.au/app/answers/detail/a_id/2711/~/-consequences-for-late-assignment-submission (https://ipoint.uwa.edu.au/app/answers/detail/a_id/2711/~/-consequences-for-late-assignment-submission)

For example, if you get 70 marks (out of the total 100 marks) before applying the late penalty and if your submission is *two days* late, then you get 60 as your final mark (i.e., 10% of the **total mark** - NOTE not your mark - is deducted).

9. Use of ChatGPT (or other Generative AI tools)

You are permitted to use ChatGPT or other Generative AI tools to generate code and help you learn. However, you must cite and explain how you used it. Generated code and documentations can be easily detected by human markers. To avoid plagiarism and academic misconduct investigations, please do provide citation and explain what changes you have to made to the code to demonstrate your own understanding. Here are a few UWA articles about how to cite and use generative AI in your assessment.

- Can I use ChatGPT and other AI tools in my assessments? (https://ipoint.uwa.edu.au/app/answers/detail/a_id/3432/related/1)
- How do I cite and reference ChatGPT (https://ipoint.uwa.edu.au/app/answers/detail/a_id/3434/~/-how-to-cite-and-reference-chatgpt-and-other-generative-ai-tools-in-assessments)
- Guide for using AI Tools at UWA (<https://www.uwa.edu.au/students/-/media/Project/UWA/UWA/Students/Docs/STUDYSmarter/Using-AI-Tools-at-UWA.pdf>)

The library guide on writing academic papers (<https://guides.library.uwa.edu.au/strategicpublishing/draftpapercheck>) might also be relevant:

For example, where relevant, Mathematics and Physical Sciences papers should discuss and reference the following as a minimum, to assign attribution and enable readers to re-create the outcome:

- Source (e.g. OpenAI/Microsoft, Anthropic/Google, NVIDIA etc.)
- Model (e.g. GPT 3)
- Implementation (e.g. davinci-003)
- Fine-tuning (Where the user has fine-tuned the 'inbuilt' knowledge of LLMs based on their own libraries of content via APIs or other processes).