

# FIT5196-S2-2022 assessment 2

***This is an individual assessment and is worth 35% of your total mark for FIT5196.***

Due date: **Check Moodle**

## Data Cleansing (60%)

For this assessment, you are required to write Python code to analyze your dataset, find and fix the problems in the data. The input and output of this task are shown below:

**Table 1. The input and output of the task**

Input	Output	Jupyter notebook
<Student_id>_dirty_data.csv <Student_id>_outlier_data.csv <Student_id>_missing_data.csv branches.csv edges.csv nodes.csv	<Student_id>_dirty_data_solution.csv <Student_id>_outlier_data_solution.csv <Student_id>_missing_data_solution.csv	<Student_id>_ass2.ipynb <Student_id>_ass2.py

**Note1: Output file names must exactly match the above standard. All files must be zipped into a file named <Student\_id>\_ass2.zip (please use zip and not rar, 7z, tar, etc.)**

**Note2: Replace <Student\_id> with your student id (do not include <>)**

**Note3: Each student can find their input files on the shared google drive**

**Note4: <Student\_id>\_ass2.py is an export of your ipynb in py format and will be used for academic integrity checks.**

Exploring and understanding the data is one of the most important parts of the data wrangling process. You are required to perform graphical and/or non-graphical EDA methods to understand the data first and then find the data problems. You are required to:

- Detect and fix errors in <Student\_id>\_dirty\_data.csv
- Detect and **remove** outlier rows in <Student\_id>\_outlier\_data.csv (outliers are to be found w.r.t. *delivery\_fee* attribute)
- Impute the missing values in <Student\_id>\_missing\_data.csv

As a starting point, here is what we know about the dataset in hand:

The dataset contains Food Delivery data from a restaurant in Melbourne, Australia. The restaurant has three branches around CBD area. All three branches share the same menu but they have different management so they operate differently.

Each instance of the data represents a single order from said restaurant. The description of each data column is shown in Table 2.

**Table 2. Description of the columns**

COLUMN	DESCRIPTION
order_id	A unique id for each order
date	The date the order was made, given in YYYY-MM-DD format
time	The time the order was made, given in hh:mm:ss format
order_type	A categorical attribute representing the different types of orders namely: Breakfast, Lunch or Dinner
branch_code	A categorical attribute representing the branch code in which the order was made. Branch information is given in the <i>branches.csv</i> file.
order_items	A list of tuples representing the order items: first element of the tuple is the item ordered, and the second element is the quantity ordered for that item.
order_price	A float value representing the order total price
customer_lat	Latitude of the customer coming from the <i>nodes.csv</i> file
customer_lon	Longitude of the customer coming from the <i>nodes.csv</i> file
customerHasloyalty?	A logical variable denoting whether the customer has a loyalty card with the restaurant (1 if the customer has loyalty and 0 otherwise)
distance_to_customer_KM	A float representing the shortest distance, in kilometers, between the branch and the customer nodes with respect to the <i>nodes.csv</i> and the <i>edges.csv</i> files. <a href="#">Dijkstra algorithm</a> can be used to find the shortest path between two nodes in a graph. Reading materials can be found <a href="#">here</a> .
delivery_fee	A float representing the delivery fee of the order

### Notes:

1. The output csv files **must** have the exact same columns as the input.
2. There is **at least one anomaly in the dataset** from each category of the data anomalies (i.e., syntactic, semantic, and coverage), and each anomaly has only one possible fix.
3. In the file `<Student_id>_dirty_data.csv`, **any row can carry no more than one anomaly**. (i.e. there can only be one anomaly in a single row and all anomalies are fixable)
4. There are no data anomalies in the file `<Student_id>_outlier_data.csv`, only outliers. Similarly, there are no data anomalies other than missing values in the file `<Student_id>_missing_data.csv`
5. There are three types of meals:
  - **Breakfast - served during morning (8am - 12pm),**
  - **Lunch - served during afternoon (12:00:01pm - 4pm)**
  - **Dinner - served during evening (4:00:01pm - 8pm)**Each meal has a distinct set of items in the menu (ex: breakfast items can't be served during lunch or dinner and so on).
6. A useful python package to solve a linear system of equations is [numpy.linalg](#)
7. Delivery fee is calculated using a different method for each branch.  
The fee depends linearly (but in different ways for each branch) on:
  - a. **weekend or weekday (1 or 0)**
  - b. **time of the day (morning 0, afternoon 1, evening 2)**
  - c. **distance between branch and customer**

**Note: If a customer has loyalty, they get a 50% discount on delivery fee**

8. The restaurant uses Dijkstra algorithm to calculate the shortest distance between customer and restaurant. (explore **networkx** python package for this or alternatively find a way to implement the algorithm yourself)
9. We know that the below columns are error-free:
  - **order\_id**
  - **time**
  - **the numeric quantity in order\_items**
  - **delivery\_fee**
10. As EDA is part of this assessment, no further information will be given publicly regarding the data. However, you can brainstorm with the teaching team during tutorials and consultation sessions.

## Methodology (25%)

The report should demonstrate the methodology (including all steps) to achieve the correct results.

## Documentation (15%)

The cleaning task must be explained in a well-formatted report (with appropriate sections and subsections). Please remember that the report must explain the complete EDA to examine the data, your methodology to find the data anomalies and the suggested approach to fix those anomalies.