

# COMP20008 Project 2

V1.0: 16th April 2019

## Due Date

The assignment is worth 20 marks, worth (20% of subject grade) **and is due 11:59pm 8th May 2019**. Submission is via the LMS. Please ensure you get a submission receipt via email. If you don't receive a receipt via email, this means your submission has not been received and hence cannot be marked. Late penalty structure is described at the end of this document.

## Objectives

The objectives of this assignment are

- To practice a further selection of techniques discussed in lectures and workshops, including clustering, correlations and predictions.
- To practice using a widely used Python library for data processing - *pandas* and a widely used library for prediction and clustering *scikit-learn*.
- To gain experience using library functions which are unfamiliar and which require consultation of additional documentation from resources on the Web.

## Background

In this project you will further practice your Python wrangling skills with the same publicly available dataset used in Project 1. The dataset is the food nutrient database from Food Standards Australia and New Zealand<sup>1</sup> in the file *food\_nutrient\_2011\_13\_AHS.csv*. It contains a detailed breakdown of nutrient values for a wide range of food items. Each nutrient value is presented on a per 100 gram edible portion basis. It was used to support the 2011-13 Australian Health Survey. Table 1 shows a high level summary description of each feature.

There is one data file provided on the project page on the LMS:

- *food\_nutrient\_2011\_13\_AHS.csv*: Nutrient information for about 5.7k foods.

Key libraries to use are Pandas, Matplotlib, NumPy, SciPy, seaborn and sklearn. You will need to write Python 3 code (Jupyter notebook) and work with the topics discussed in workshops weeks (6-8). If you are using other packages, you must provide an explanation in your code about why it is necessary.

---

<sup>1</sup><http://www.foodstandards.gov.au/science/monitoringnutrients/ausnut/ausnutdatafiles/Pages/foodnutrient.aspx>

| Field Name            | Description   |
|-----------------------|---|
| Food ID               | 8 digit alpha numeric food identification code, based on FSANZ standard   |
| Survey ID             | 8 digit food identification code that was specific to Australian Health Survey  |
| Survey flag           | can be ignored  |
| Food Name             | Name of the food  |
| nutrient name (units) | Describes the full nutrient name of each of the nutrients<br>e.g. 'Protein' and includes the units the nutrient is presented in<br>e.g. '(g)' for grams. A value is then provided for each food and nutrient. |

Table 1: Summary of Features for *food\_nutrient\_2011\_13\_AHS.csv*

## Importing required Python Libraries and loading the data

Please begin your Jupyter notebook by listing all the Python libraries you will be using. Also load the dataset (.csv) in a dataframe object.

```
#import ....
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
from sklearn.preprocessing import StandardScaler
from sklearn.decomposition import PCA
from scipy.spatial.distance import pdist, squareform
import seaborn as sns
from sklearn.model_selection import train_test_split
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import accuracy_score
from sklearn.tree import DecisionTreeClassifier
from sklearn.cluster import KMeans
from scipy.spatial.distance import cdist
food = pd.read_csv("food_nutrient_2011_13_AHS.csv", header=0, low_memory=False)
```

There are also some helpful functions like VAT (visualization of clustering tendency) and MI (mutual information) that you can use and integrate into your code. These are provided in the example Jupyter notebook file, and they were covered in workshops for weeks 6 and 7.

## 1 Feature standardisation(1 mark)

In order for a number of types of analysis techniques to be effective, it is important for features to be scaled or standardised.

a) Create a new data frame that contains only the continuous features from *food\_nutrient\_2011\_13\_AHS.csv*. I.e. the features starting from *Energy, with dietary fibre (kJ)* and ranging to *Total trans fatty acids (mg)*. (53 features in total). This data frame is the input for part b).

b) Apply the *sklearn.preprocessing.StandardScaler()* function to standardise each feature separately, to have 0 mean and unit variance. Store the transformed features in *foodscaled*.

c) Print the number of rows and columns, as well as the minimum, maximum, mean and standard deviation (computed taking into account all values) of *foodscaled*, using the following format:

```
***
Q1.c: foodscaled matrix details
Number of rows: #
Number of columns: #
Min: #
Max: #
Mean: #
Standard Deviation: #
***
```

where each # is the appropriate value rounded to 1 decimal place.

Hint: Useful functions include [Standard Scaler](#)

## 2 Principal components analysis (2 marks)

It is difficult to visualise the data in *foodscaled*, due to the large number of features. One strategy to deal with this is to apply principal components analysis.

a) Create a feature *EnergyLevel* which has value "1" if (unstandardized) *Energy, with dietary fibre (kJ)* is greater than 1000 kJ and "0" otherwise.

b) Apply principal components analysis to your data in *foodscaled* to compute the first two principal components. Store the result in *foodreduced*.

c) Produce a scatter plot of the data in *foodreduced*. A food should have red color if it has 'High' energy (*EnergyLevel* of "1") and blue if it has 'Low' energy (*EnergyLevel* of "0"). There should be an accompanying legend mapping colours to *EnergyLevel*.

d) Comment on your scatter plot in c). What does it show? How can it be interpreted? What are the advantages and disadvantages of using PCA for visualising this food dataset?

Hint: Useful functions are [sklearn PCA](#)

## 3 Clustering visualisation (3 marks)

a) Similar to what you did in Project 1, create a new feature *Food category*, which contains the first two digits of the feature Survey ID.

b) From *foodscaled*, select only foods from categories 13, 20 and 24. Store these foods in *foodscaledsample*.

c) Use the `VAT()` function, discussed in Workshop Week 6, to compute the ordered dissimilarity matrix for the `foodscaledsample` matrix.

d) Plot the heatmap for the ordered dissimilarity matrix from c). For your plot, select a colormap that is effective in revealing the cluster structure.

e) How many clusters are apparent in your heatmap? Is this expected (why/why not)? Why does use of different colormaps produce visualisations of varying usefulness? Explain what properties an optimal colormap should have for this task.

Hint: Useful items include [Seaborn heatmap](#), [Colormap reference](#).

## 4 K-means and sum of squared errors (1.5 marks)

Recall question 5 from the Week 6 workshop. This question will take you through a practical example of computing sum of squared errors (SSE), which is a quality measure that can be used for evaluating a clustering produced by *k*-means.

a) Using the dataset `foodscaled`, produce a plot of the SSE value of the *k*-means clustering of the dataset (y-axis), versus *k* value (x axis). *k* should range across 2, 3, 4, ..., 25. This requires generation of 24 different *k*-means clusterings. When generating each clustering, use `random.state=100` and default values for all other parameters of the `KMeans` function, except `n_clusters`, which will need to be varied.

b) Comment on the shape and trend of your plot from a). Where is the elbow? Is this expected (why/why not)?

Hint: Useful functions are the VAT code provided in the example Jupyter notebook and `sklearn.cluster.KMeans` and `scipy.spatial.distance.cdist`.

## 5 Correlation and Mutual Information(4 marks)

a) **Visualise correlation matrix:** Plot a Pearson correlation matrix for the first 10 nutrients (shown below) in the `food` data frame. The calculated 10\*10 symmetric matrix should contain the correlation between every pair of attributes. For example, a value in row *i* and column *j* should contain the Pearson correlation  $r_{ij}$  between two features *i* and *j*. Plot a heatmap for your computed correlation matrix, including feature names on each axis of the heatmap.

```
'Energy, with dietary fibre (kJ)',  
'Energy, without dietary fibre (kJ)',  
'Moisture (g)',  
'Protein (g)',  
'Total fat (g)',  
'Available carbohydrates, with sugar alcohols (g)',
```

'Available carbohydrates, without sugar alcohol (g)',  
'Starch (g)',  
'Total sugars (g)',  
'Added sugars (g)'

**b) Analysing the effect of binning**

- For the pair of features in the *food* dataframe: '*Protein (g)*' and '*Energy, with dietary fibre (kJ)*': Calculate a series of Mutual Information (MI) values, according to increasing numbers of *equal-width* bins. Use [2, 10, 20, 30, 40, ... , 200] numbers of bins and calculate 21 MI values, one for each number of bins.
- Create a plot, where the x axis is number of bins and the y axis is the MI of the feature pair.
- Explain the trend in your plot.

**c) Comparison of Pearson correlation and Mutual Information:**

- Considering the full set of 53 continuous features from the *food* data frame, find the top-10 feature-pairs by Pearson correlation strength (highest to lowest)
- Considering the full set of 53 continuous features from the *food* data frame, find the top-10 feature-pairs by Mutual Information, using 20 equal-width bins
- Print out the top-10 feature-pairs for Pearson correlation (highest to lowest) and by MI (highest to lowest). Comment on the similarities and differences between the two lists.

Hint: Useful functions are `pandas.DataFrame.corr` and the provided mutual information function in the example Jupyter notebook. Note that executing the code to answer this question may be more computationally demanding, and so you may wish to use a smaller dataset during initial development and debugging.

## 6 Prediction models: decision trees (2.5 marks)

This question explores predicting the *Food category* label, using the values of the various nutrients. I.e. *Food category* is the class label.

a) Randomly split the *foodscaled* data into 80% training and 20% testing. You should output two matrices: *X\_train* and *X\_test*, and two vectors/columns: *y\_train* and *y\_test*.

- *X\_train* contains features of training instances (i.e. 80%) and *y\_train* contains labels for the training instances.
- *X\_test* contains features of test instances (i.e. 20%) and *y\_test* contains labels for test instances.

The output of this question should print the shape (number of rows and columns) of  $X_{train}$ ,  $X_{test}$ ,  $y_{train}$  and  $y_{test}$  in the following format:

```
***
Q6.a: Train Test Split Results
X_train matrix: #
y_train labels: #
X_test matrix: #
y_test labels: #
***
```

where  $\#$  are the calculated shape values.

b) Using the same training and testing data, you will next analyse the effect of varying the  $max\_depth$  parameter for the decision tree. Produce a plot showing accuracy (y-axis) versus  $max\_depth$  (x-axis), where  $max\_depth$  varies in the range  $1, 2, 3, \dots, 40$ . Each data point in the plot corresponds to the accuracy of a decision tree that is trained using the 80% training data and evaluated using the 20% test data, to predict the *Food category*.

c) Comment on the shape of the plot in Part c). Suggest reasons for the shape of the plot and any local peaks.

Hint: Useful functions are `sklearn.model_selection.train_test_split` and `sklearn.tree.DecisionTreeClassifier`

## 7 Prediction models: K-NN (2 marks)

a) Create an instance of the *KNeighborsClassifier* and set  $k$  (number of neighbors) to 1. Train/fit the model using  $X_{train}$  and  $y_{train}$ . Evaluate the model by calculating its accuracy when applied to both the train and the test set (i.e.  $X_{train}$  and  $X_{test}$ ).

Then print out the following:

```
***
Q7a: Food category prediction using k-NN (k=1)
Train accuracy: # %
Test accuracy: # %
***
```

Where  $\#$  is the calculated  $k$ -NN classifier accuracy rounded to 1 decimal place.

b) Repeat the steps of a), this time using  $k = 3$ . The output will be as follows

```
***
Q7.b: Food category prediction using k-NN (k=3)
Train accuracy: # %
Test accuracy: # %
***
```

c) Comment on the differences and similarities between your output for parts a) and b). Where they are different, suggest reasons why.

d) In parts b) and c), the reported accuracy may be over-optimistic, due to the way standardisation was applied in Question 1. Explain why this is the case and how this issue could be addressed.

Hint: Useful functions are `sklearn.neighbors.KNeighborsClassifier`

## 8 Feature generation (4 marks - harder)

In order to achieve higher prediction accuracy for  $k$ -NN, one can investigate the use of feature generation and selection. Again, *Food category* is the class label to be predicted.

*Feature generation* involves the creation of additional features, additional to the ones already present in *foodscaled*. Two methods are

- *Interaction term pairs*. Given a pair of features  $f_1$  and  $f_2$  in *foodscaled*, create a new feature  $f_{pair} = f_1 \times f_2$  or by  $f_{pair} = \frac{f_1}{f_2}$ . All possible pairs can be considered.
- *Clustering labels*: Apply  $k$ -means clustering to the data in *foodscaled* and then use the resulting cluster labels as the values for a new feature  $f_{clusterlabel}$ . At test time, a label for a testing instance can be created by assigning it to its nearest cluster.

Given a set of  $N$  features (the original features plus generated features), *feature selection* involves selecting a smaller set of  $n$  features ( $n < N$ ). Here one computes the mutual information between each feature and the class label (on the training data), sorts the features from highest to lowest mutual information value, and then retains the top  $n$  features from this ranking, to use for classification with  $k$ -NN.

Your task in this question is to evaluate whether the above methods for feature generation and selection, can deliver a boost in prediction accuracy compared to using  $k$ -NN just on the original features in *foodscaled*. You should:

- Implement the above two methods for feature generation. You may need to experiment with different parameter values, including  $k$  and different numbers of generated features.
- Implement feature selection using mutual information. Again you may need to experiment with different parameter values, such as how many features to select.

Your output for this question should include i) your implemented code that generates accuracies or plots, ii) a discussion (2-3 paragraphs) of whether feature selection+generation can deliver an accuracy boost, based on your evidence from i).

Note that you should use a training and testing data split, to evaluate accuracies. Feature generation for training data should not require any knowledge of the testing data. The feature selection step (on the training data) should not require any knowledge of the testing data.

Hint: This is a more open ended question. It must be clear from your answer that you have successfully implemented the proposed strategy, experimented with different parameter settings, and have a clear, evidenced-based explanation of the extent to which it succeeded.

## Marking scheme

*Correctness (20 marks):* For each of the questions, a mark will be allocated for level of correctness (does it provide the right answer, is the logic right), according to the number in parentheses next to each question. Note that your code should work for any data input formatted in the same way as *food\_nutrient\_2011\_13\_AHS.csv*. E.g. If a random sample of 1000 records was taken from *food\_nutrient\_2011\_13\_AHS.csv*, your code should provide a correct answer if this was instead used as the input.

Correctness will also take into account the readability and labelling provided for any plots and figures (plots should include title of the plot, labels/scale on axes, names of axes, and legends for colours symbols where appropriate).

*Coding style:* A global deduction of 1 mark may be made if any of the following aspects are lacking:

- Formatting of code (e.g. use of indentation and overall readability for a human)
- Code modularity and flexibility. Use of functions or loops where appropriate, to avoid highly redundant or excessively verbose definitions of code.
- Use of Python library functions (you should avoid reinventing logic if a library function can be used instead)
- Code commenting and clarity of logic. You should provide comments about the logic of your code for each question, so that it can be easily understood by the marker.

## Resources

The following are some useful resources, for refreshing your knowledge of Python, and for learning about functionality of pandas.

- [Python tutorial](#)
- [Python beginner reference](#)
- [pandas 10min tutorial](#)
- [Official pandas documentation](#)
- [Official matplotlib tutorials](#)
- [Python pandas Tutorial by Tutorialspoint](#)
- [pandas: A Complete Introduction by Learn Data Sci](#)
- [pandas quick reference sheet](#)
- [Python Data Analytics by Fabio Nelli](#) (available via University of Melbourne sign on)



## Submission Instructions

Via the LMS, submit a Jupyter notebook (A template Jupyter notebook "example-notebook-project2.ipynb" is provided in the folder with the datasets) containing your Python 3 code.

## Other

*Extensions and Late Submission Penalties:* If requesting an extension due to illness, please submit a medical certificate to the lecturer. If there are any other exceptional circumstances, please contact the lecturer with plenty of notice. Late submissions without an approved extension will attract the following penalties

- $0 < \text{hourslate} \leq 24$  (2 marks deduction)
- $24 < \text{hourslate} \leq 48$  (4 marks deduction)
- $48 < \text{hourslate} \leq 72$ : (6 marks deduction)
- $72 < \text{hourslate} \leq 96$ : (8 marks deduction)
- $96 < \text{hourslate} \leq 120$ : (10 marks deduction)
- $120 < \text{hourslate} \leq 144$ : (12 marks deduction)
- $144 < \text{hourslate}$ : (20 marks deduction)

where *hourslate* is the elapsed time in hours (or fractions of hours).

This project is expected to require 20-25 hours work.

## Academic Honesty

You are expected to follow the academic honesty guidelines on the University website <https://academichonesty.unimelb.edu.au>

## Further Information

A project discussion forum has also been created on the subject LMS. Please use this in the first instance if you have questions, since it will allow discussion and responses to be seen by everyone. There will also be a list of frequently asked questions on the project page.