# HW1 Report (Group 29)

Group Members: 張柏豐(111651070)、王睿綸(110511290)

## A. Summary of K-Nearest Neighbors (KNN) Processing Steps

1. **Import Necessary Functions and Download the Specified Kaggle Dataset**

   Begin by importing the required functions and downloading the designated Kaggle dataset.

2. **Remove Data with a Mix of English and Numerals (Having Minimal Impact on Classification Features)**

   Exclude data instances that combine English and numerals, as these have minimal influence on classification features.

3. **Calculate the Feature Correlation Matrix Using the Feature Correlation Matrix**

   Utilize the feature correlation matrix to assess the relationships between each feature and the classification labels.

4. **Select the Most Relevant Features from the Training and Testing Sets**

   Identify the most highly correlated features with the training and testing datasets.

5. **Evaluate Accuracy Using a KNN Model**

   Assess accuracy by feeding these features into a KNN model for testing.

6. **Enhance Precision with Normalization and Weight**

   Improve precision by introducing normalization and weighting techniques.

## B. Detailed Code Explanation

### 1. Import Function

```
# import function
!pip install kaggle
from google.colab import files
import pandas as pd
import numpy as np
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import accuracy_score

# Z-Score Normalization
from sklearn.preprocessing import StandardScaler

# Robust Normalization
from sklearn.preprocessing import RobustScaler
```

First, import all the libraries that will be used in the following code. The 'files' library is imported to allow the upload of the Kaggle.json file, which enables us to use Kaggle's API to download the dataset for this assignment. The 'KNeighborsClassifier' library is used for our k-nearest neighbors (k-NN) model, and the 'accuracy_score' function is imported to evaluate accuracy. Finally, the last two libraries are used for normalizing the training and test datasets.

## 2. Download data in Kaggle

```
# Download data in Kaggle
uploaded = files.upload() # kaggle.json
!mkdir -p ~/.kaggle/
!cp kaggle.json ~/.kaggle/
!chmod 600 ~/.kaggle/kaggle.json # Setting Permissions of the File (kaggle.json), only readable by the current user
!kaggle competitions download -c groupchallenge-knn-competition-pds
!unzip groupchallenge-knn-competition-pds.zip
```

Next, we start by uploading the 'kaggle.json' file to the Google Colab workspace. We then move this file to a location where we can access it and set the permissions to be accessible only by the current user to prevent others from using the API to download datasets from this account. After that, we proceed to download and unzip the dataset for this assignment, preparing for the next steps.

## 3. Data formatting

```
# For Train
df = pd.read_csv("train.csv", sep=',', header=None, dtype=str, skipinitialspace=True)
df = df.fillna(0)
train_array = df.values

# For Test
df = pd.read_csv("test.csv", sep=',', header=None, dtype=str, skipinitialspace=True)
df = df.fillna(0)
test_array = df.values
```

To facilitate the upcoming data processing, we first read the original CSV file and convert it into an array, recording the data in each cell. Each entry in the CSV file is separated by a comma (",") and stored in a matrix. The "dtype" parameter specifies the data type of the stored data. Finally, the "skipinitialspace" option is used to skip any whitespace before and after the commas.

Next, we use the "df.fillna(0)" function to replace missing values ("nan") in the data table with the numerical value 0. Finally, we convert the DataFrame into a Numpy array format.

## 4. Data Cleaning

```
# Change Y, N to 1, 0 in columns 7, 8
train_array[:, 7] = np.where(train_array[:, 7] == "Y", 1, 0)
train_array[:, 8] = np.where(train_array[:, 8] == "Y", 1, 0)
test_array[:, 7] = np.where(test_array[:, 7] == "Y", 1, 0)
test_array[:, 8] = np.where(test_array[:, 8] == "Y", 1, 0)

# Segment IDX
train_array = train_array[1:, :]
test_array = test_array[1:, :]

# Remove mixed alphanumeric data (low-relevance)
# [IDX, id, full_name, pdes, name, prefix, orbit_id, equinox, class]
delete_columns = [0, 1, 3, 4, 5, 6, 13, 17, 44]
train_array = np.delete(train_array, delete_columns, axis=1)
test_array = np.delete(test_array, delete_columns, axis=1)
```

Next, to prepare the data, we'll start by addressing the columns 7 and 8 in both the train_array and test_array. These columns contain only three possible values: "N," "Y," and "nan" (which we treat as 0). We will use One-Hot Encoding to encode these values as follows: "N" becomes 0, "Y" becomes 1, and "nan" becomes 0.

Next, we remove the first row, which contains labels, and retain only the data.

Finally, for the data that consists of a mix of English and numeric values, since most of it is primarily numeric and does not significantly affect the KNN prediction results, we can remove it. This is mainly because both the Feature Correlation Matrix and the KNN model cannot handle non-numeric data, so it must be excluded. The selected columns for removal include: [IDX, id, full_name, pdes, name, prefix, orbit_id, equinox, class]

As a result, the train_array and test_array matrices we obtain will consist entirely of numerical values.
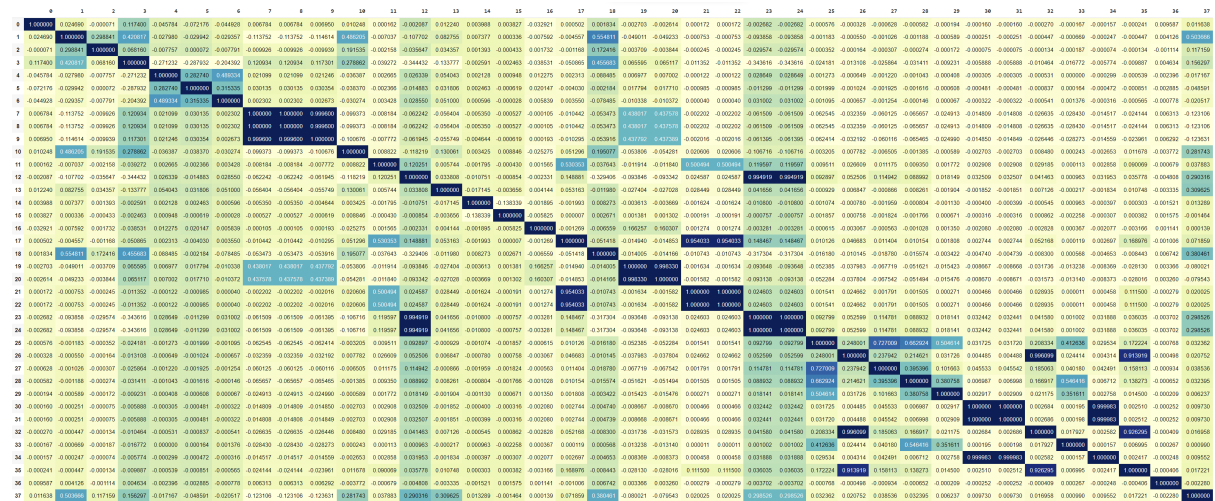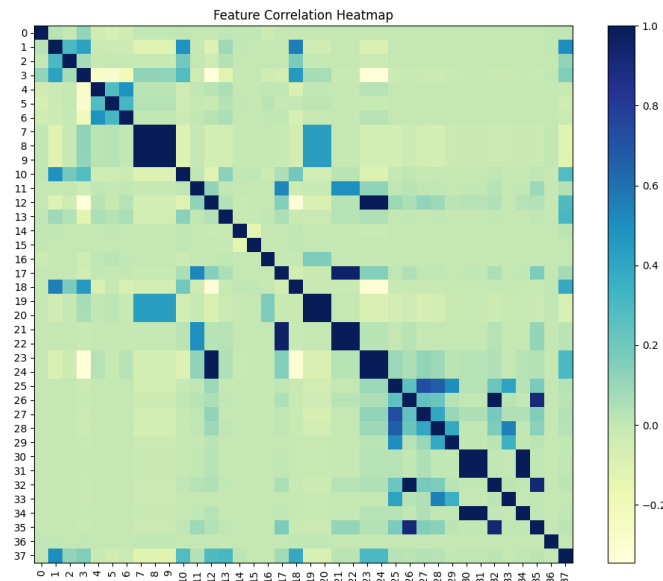
## 5. Feature Correlation Matrix

```
# Creat feature correlation matrix
train_df = pd.DataFrame(train_array.astype("float"))
correlation_matrix = train_df.corr(numeric_only=True)
display(correlation_matrix.style.background_gradient(cmap='YlGnBu').set_properties(2))

# Segment labels and answers
train_data = train_array[:, :-1]
train_label = train_array[:, -1]
test_data = test_array[:, :-1]
test_label = test_array[:, -1]
```

Next, we'll export a Feature Correlation Matrix from the previously obtained train_array to understand the relationships between each feature. It's essential to note that the data matrix we obtained earlier consists of both integers and floating-point numbers. Different data types can prevent us from conducting feature correlation tests. Therefore, we need to convert all values to floating-point numbers (float) first before generating the chart. We configure the chart's color gradient using the "YlGnBu" color format and set the font size for the correlation values of each feature to level 2 (.set_properties(2)).

After generating the chart, we proceed to separate the answers and features in both the train_array and test_array, storing them in two separate matrices.

Feature Correlation Heatmap

## 6. Reserve Correlated Column

```
# Select features with higher correlation
reserved_column = [1, 10, 12, 13, 18, 23, 24] #[neo, e, q, i, n, moid, moid_ld]
train_data = train_data[:, reserved_column]
test_data = test_data[:, reserved_column]
```

From the previously generated feature correlation matrix, we can observe that features in columns 1, 10, 12, 13, 18, 23, and 24 (neo, e, q, i, n, moid, moid_ld) exhibit the highest correlations with the last column (class_id). The correlation values range from 30% to 50%. Therefore, we will only select these specific features from train_data and test_data for training and prediction using k-nearest neighbors (KNN).
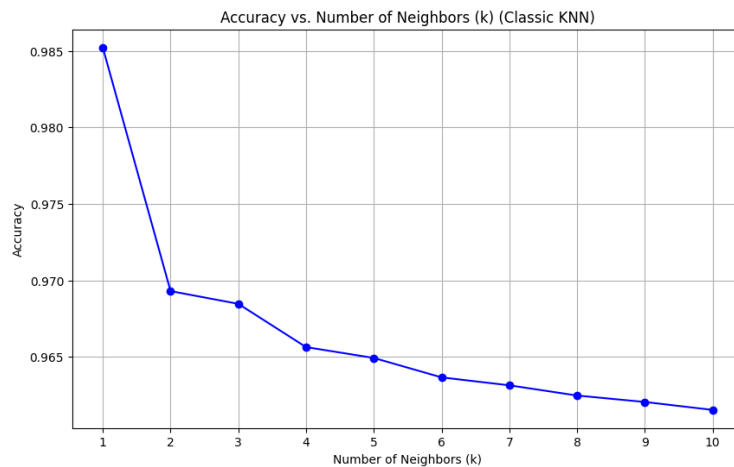
## 7. Classic KNN w/ k=1

```
# Train & Predict (None)
accuracies = []

for k in range(1, 11):
  knn = KNeighborsClassifier(n_neighbors = k)
  knn.fit(train_data, train_label)
  test_pred_none = knn.predict(test_data)
  accuracy = accuracy_score(test_label, test_pred_none)
  accuracies.append(accuracy)

plt.figure(figsize=(10, 6))
plt.plot(range(1, 11), accuracies, marker='o', linestyle='-', color='b')
plt.title('Accuracy vs. Number of Neighbors (k) (Classic KNN)')
plt.xlabel('Number of Neighbors (k)')
plt.ylabel('Accuracy')
plt.grid(True)
plt.xticks(range(1, 11))
plt.show()
```

```
Accuracy: [0.9852202701386861, 0.9693033704739251, 0.9684722280181948, 0.9656449526353128, 0.9649355256019697,
           0.9636662076256451, 0.9631480476846249, 0.9624803516532432, 0.9620595640501328, 0.9615414041091127]
```

Accuracy vs. Number of Neighbors (k) (Classic KNN)

Next, we move on to the evaluation of training and prediction. Here, we perform KNN training without adding weights or normalizing the data. We directly train the KNN model. We also vary the parameter n_neighbors from 1 to 10 and record the accuracy for each case to create the chart as shown above. We can observe that when n_neighbors = 1, the highest accuracy is achieved (approximately 0.9852). We speculate that this phenomenon occurs because, after feature correlation analysis, we have already selected the most relevant features. Therefore, the nearest neighbors are more likely to be the correct class, reducing the chances of misclassification due to less relevant features.
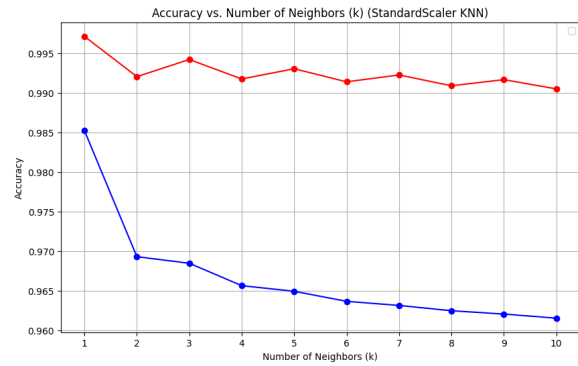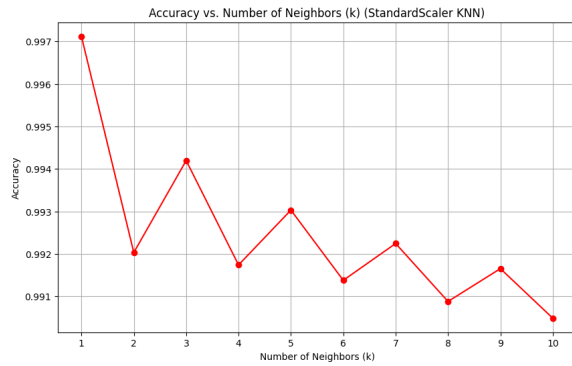
## 8. StandardScaler Method

```
# Train & Predict (StandardScaler)

scaler = StandardScaler()
scaled_train_data = scaler.fit_transform(train_data)
scaled_test_data = scaler.fit_transform(test_data)

accuracies_Standard = []
for k in range(1, 11):
  knn = KNeighborsClassifier(n_neighbors = k)
  knn.fit(scaled_train_data, train_label)
  test_pred_standard = knn.predict(scaled_test_data)
  accuracy = accuracy_score(test_label, test_pred_standard)
  accuracies_Standard.append(accuracy)

plt.figure(figsize=(10, 6))
plt.plot(range(1, 11), accuracies_Standard, marker='o', linestyle='-', color='r')
plt.title('Accuracy vs. Number of Neighbors (k) (StandardScaler KNN)')
plt.xlabel('Number of Neighbors (k)')
plt.ylabel('Accuracy')
plt.grid(True)
plt.xticks(range(1, 11))
plt.show()
```

```
Accuracy: [0.9971170832811695, 0.9920328562088776, 0.9941959131438746, 0.991740739195148, 0.9930309226724534,
           0.991379070511483, 0.9922449888021811, 0.9908817760714435, 0.9916537996077286, 0.9904888091363073]
```

Next, we introduce a standardization technique called StandardScaler, which is one of the most common methods for standardization. It involves subtracting the mean and dividing by the standard deviation. After applying this technique to both the train_data and test_data, we conduct the same evaluation and create the chart as shown above. It's evident that data accuracy significantly improves after applying StandardScaler, with the highest accuracy still observed at k=1. Notably, the accuracy line exhibits a sawtooth pattern as k increases, with odd values achieving higher accuracy. However, the overall trend still shows a gradual decrease as k increases.
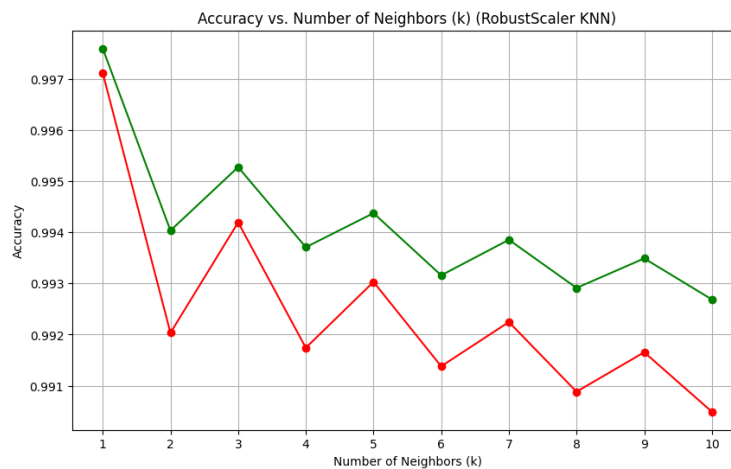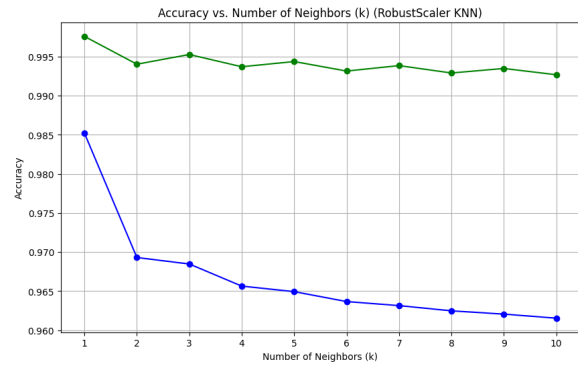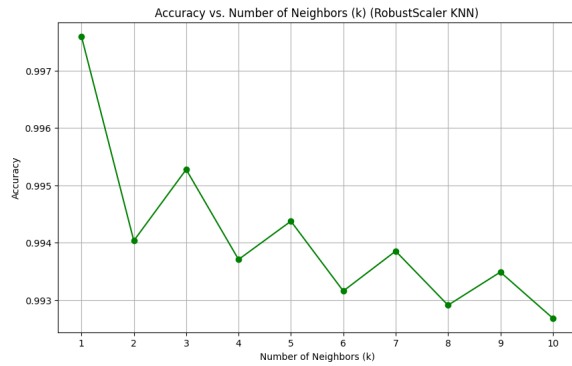
## 9. RobustScaler Method

```
# Train & Predict (RobustScaler)

scaler = RobustScaler()
scaled_train_data = scaler.fit_transform(train_data)
scaled_test_data = scaler.fit_transform(test_data)

accuracies_Robust = []
for k in range(1, 11):
  knn = KNeighborsClassifier(n_neighbors = k)
  knn.fit(scaled_train_data, train_label)
  test_pred_Robust = knn.predict(scaled_test_data)
  accuracy = accuracy_score(test_label, test_pred_Robust)
  accuracies_Robust.append(accuracy)

plt.figure(figsize=(10, 6))
plt.plot(range(1, 11), accuracies_Robust, marker='o', linestyle='-', color='g')
plt.title('Accuracy vs. Number of Neighbors (k) (RobustScaler KNN)')
plt.xlabel('Number of Neighbors (k)')
plt.ylabel('Accuracy')
plt.grid(True)
plt.xticks(range(1, 11))
plt.show()
```

```
Accuracy: [0.9975969898037252, 0.9940359443030227, 0.9952774416113731, 0.9937090514543254, 0.9943767474857071,
           0.9931595932618342, 0.9938551099611902, 0.9929126848335629, 0.9934899636940283, 0.9926866419062722]
```

Accuracy vs. Number of Neighbors (k) (RobustScaler KNN)



Accuracy vs. Number of Neighbors (k) (RobustScaler KNN)



Accuracy vs. Number of Neighbors (k) (RobustScaler KNN)

Next, we introduce another standardization technique called RobustScaler, which computes the center and range of the data using the median and quartiles. The key difference compared to the previous method lies in RobustScaler's ability to automatically exclude outliers to prevent them from affecting feature interpretation erroneously.

After applying this technique to both the train_data and test_data, we generate the chart as shown above. We can observe that this method yields significantly higher accuracy compared to KNN without standardization, and it also outperforms StandardScaler. The accuracy trend still exhibits a sawtooth pattern with decreasing accuracy as k increases.

Since this data preprocessing approach results in the most significant accuracy improvement, we choose to use RobustScaler along with weighted k-NN for prediction.

## 10. RobustScaler + Weight Method

```
# Train & Predict (RobustScaler + Weight)

scaler = RobustScaler()
scaled_train_data = scaler.fit_transform(train_data)
scaled_test_data = scaler.fit_transform(test_data)

accuracies_Robust_weight = []
for k in range(1, 11):
  knn = KNeighborsClassifier(n_neighbors = k, weights='distance')
  knn.fit(scaled_train_data, train_label)
  test_pred_Robust_weight = knn.predict(scaled_test_data)
```

```
    accuracy = accuracy_score(test_label, test_pred_Robust_weight)
    accuracies_Robust_weight.append(accuracy)

plt.figure(figsize=(10, 6))
plt.plot(range(1, 11), accuracies_Robust, marker='o', linestyle='-', color='g')
plt.plot(range(1, 11), accuracies_Robust_weight, marker='s', linestyle='--', color='g')
plt.title('Accuracy vs. Number of Neighbors (k) (RobustScaler KNN)')
plt.xlabel('Number of Neighbors (k)')
plt.ylabel('Accuracy')
plt.grid(True)
plt.xticks(range(1, 11))
plt.legend()
plt.show()
```
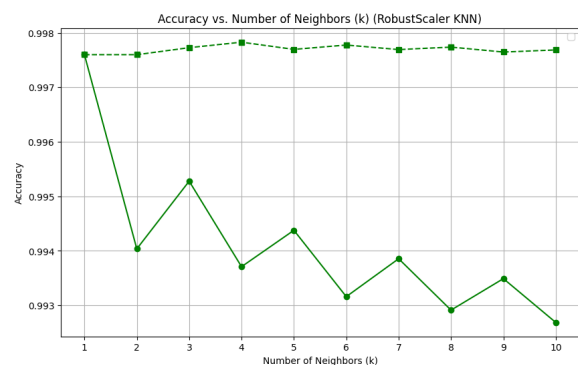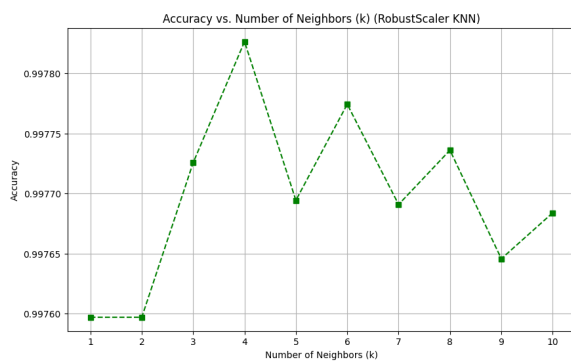
```
Accuracy: [0.9975969898037252, 0.9975969898037252, 0.9977256603931061, 0.9978265103145126, 0.997694362141635,
           0.997774346562061, 0.9976908845581383, 0.9977360931435963, 0.9976456759726801, 0.9976839293911447]
```



We set the weight to "distance," meaning that closer data points have higher weights. After introducing weighted k-NN on top of RobustScaler, we observe a further increase in accuracy compared to using RobustScaler alone. Interestingly, the highest accuracy is achieved at k=4, not k=1. We believe that this phenomenon occurs because when each data point in space is weighted based on proximity, the model considers not only the nearest neighbor but also nearby data points with a degree of proximity. This helps prevent overfitting that might occur when considering only the nearest neighbor.

After reaching the highest accuracy at k=4, the accuracy gradually decreases as k increases, still exhibiting a sawtooth pattern. Notably, the accuracy remains higher than when considering only 1 nearest neighbor.

## 11. Organize The Output

```
# Save as .csv file

pred = test_pred_Robust_weight
idx_values = np.arange(len(pred))
result = np.column_stack((idx_values, pred))
result = pd.DataFrame(result, columns=['IDX', 'Target'])
result.to_csv('test_pred_Robust_weight.csv', index=False)
files.download('test_pred_Robust_weight.csv')
```

Based on the experiments conducted, we have observed that the accuracy significantly improves when the data is standardized and weighted between data points in the k-NN model. Ultimately, we choose the RobustScaler with distance-weighted k-NN, as it achieved the highest accuracy. We can now finalize the predictions and export them as a CSV file.