# Big Data Analytics Project

## Google Play Store Apps Analysis

411021206　莊芳儒

411021209　張佩慈

411021231　林芷安

# 1. Dataset

來自網站 Kaggle，網址如下

https://www.kaggle.com/datasets/lava18/google-play-store-apps/data

選擇檔案為" googleplaystore.csv" 共 10842 筆

| A | B | C | D | E | F | G | H | I | J | K | L | M |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| App | Category | Rating | Reviews | Size | Installs | Type | Price | Content Rating | Genres | Last Updated | Current Ver | Android Ver |

使用到的項目有:

Category、Rating、Reviews、Size、Installs、Type、Price、Last Updated

# 2. 程式碼

主要分成 7 個部分:

①下載資料集

```
1   # kaggle API
2   os.system('kaggle datasets download -d lava18/google-play-store-apps')
3   # create directory
4   os.system('mkdir unzip')
5   # extract $zip_path and put it into $output_dir
6   def extract_files(zip_path, output_dir):
7       with zipfile.ZipFile(zip_path, 'r') as zipf:
8           zipf.extractall(output_dir)
9   extract_files('google-play-store-apps.zip', 'unzip')
10  path = "C:/Users/asus/.spyder-py3/unzip/googleplaystore.csv"
```

②資料預處理

原資料呈現:

| | App | Category | Rating | Reviews | Size | Installs | Type | Price | Content Rating | Genres | Last Updated | Current Ver | Android Ver |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | | | | | | | | | | | | | |
| 5 | Sketch - Draw & Paint | ART_AND_DESIGN | 4.5 | 215644 | 25M | 50,000,000+ | Free | 0 | Teen | Art & Design | 8-Jun-18 | Varies with device | 4.2 and up |

```
1   usecols = ['Category', 'Rating', 'Reviews', 'Size', 'Installs',
                                       'Type', 'Price', 'Last Updated']
2   data = pd.read_csv(path, usecols=usecols)
3   #cleaning
4   data.drop(index=10472, inplace=True)
5   a = data.dropna()
6   # installs
7   Installs = a.Installs.apply(lambda x: x.replace('+', ""))
8   Installs = a.Installs.apply(lambda x: x.replace(',', ""))
9   Installs = pd.to_numeric(a.Installs, errors="coerce")
```

```
10  # reviews
11  a['Reviews'] = pd.to_numeric(a['Reviews'])
12  # size
13  def convert(size):
14      if type(size) == str:
15          if 'k' in size:
16              return format(float(str(size).replace('k', "")) / 1024, '.3f')
17          if 'M' in size:
18              return float(str(size).replace('M', ""))
19      else: return(size)
20  a['Size'] = a['Size'].str.replace('Varies with device', '0.00')
21  a['Size'] = a['Size'].apply(lambda x: convert(x))
22  a['Size'] = pd.to_numeric(a['Size'])
23  a['Size'].fillna(format(a['Size'].mean(), '.3f'), inplace=True)
24  a['Size'] = pd.to_numeric(a['Size'])
25  # price
26  a['Price'] = a['Price'].apply(lambda x: x.replace('$', ""))
27  a['Price'] = pd.to_numeric(a['Price'])
```

③圖表分析

```
1   # category-count-bar
2   value_count = a['Category'].value_counts()
3   value_count.plot(kind='bar', color='lightslategray')
4   plt.ylabel('Count')
5   plt.show()
6   # category-Install-bar
7   S = a.groupby('Category')['Installs'].sum()
8   S.plot.bar(color = 'lightslategray')
9   plt.ylabel('Installs')
10  plt.show()
11  # rating-Count-bar(Free,Paid)
12  C_Free = a.groupby('Type').get_group('Free')
13  C_Paid = a.groupby('Type').get_group('Paid')
14  C_Freec = C_Free['Rating'].value_counts()
15  C_Paidc = C_Paid['Rating'].value_counts()
16  C_Freec.plot.bar(color='lightslategray')
17  plt.ylabel('Count')
18  plt.title('Free')
19  plt.show()
20  C_Paidc.plot.bar(color='lightslategray')
21  plt.ylabel('Count')
22  plt.title('Paid')
```

```
23  plt.show()
24  # free_vs_Paid-count-pie
25  C_Freecnt = C_Freec.sum()
26  C_Paidcnt = C_Paidc.sum()
27  x = [C_Freecnt, C_Paidcnt]
28  name = ["Free", "Paid"]
29  color = ['rosybrown', 'lightslategray']
30  plt.pie(x, autopct='%.3f%%', colors=color, labels=name)
31  plt.show()
```

④相關性分析

```
1  b = pd.DataFrame(a.iloc[:, [1, 2, 3, 4, 6]])
2  corr_matrix = b.corr()
3  sn.heatmap(corr_matrix, annot=True)
4  plt.show()
```

⑤聚類

rating-install

```
1   X = a.iloc[:, [1, 4]]
2   # elbow
3   wcss = []
4   for i in range(1, 11):
5       kmeans = KMeans(n_clusters=i, n_init='auto', random_state=42)
6       kmeans.fit(X)
7       wcss.append(kmeans.inertia_)
8   plt.plot(range(1, 11), wcss)
9   plt.show()
10  # cluster show
11  kmeans = KMeans(n_clusters=3, init='k-means++', random_state=42)
12  y_kmeans = kmeans.fit_predict(X)
13  X = np.array(X)
14  plt.scatter(X[y_kmeans == 0, 0], X[y_kmeans == 0, 1], s=10, c='red', label='low')
15  plt.scatter(X[y_kmeans == 1, 0], X[y_kmeans == 1, 1], s=10, c='blue', label='mid')
16  plt.scatter(X[y_kmeans == 2, 0], X[y_kmeans == 2, 1], s=10, c='green', label='high')
17  plt.scatter(kmeans.cluster_centers_[:, 0], kmeans.cluster_centers_[:, 1], s=100,
    c='black', label='Centroids')
18  plt.title('google play store apps')
19  plt.xlabel('Rating')
20  plt.ylabel('Installs')
21  plt.legend()
22  plt.show()
```

review-install

```
1   X = a.iloc[:, [2, 4]]
2   # elbow
3   wcss = []
4   for i in range(1, 11):
5       kmeans = KMeans(n_clusters=i, n_init='auto', random_state=42)
6       kmeans.fit(X)
7       wcss.append(kmeans.inertia_)
8   plt.plot(range(1, 11), wcss)
9   plt.show()
10  # cluster show
11  kmeans = KMeans(n_clusters=3, init='k-means++', random_state=42)
12  y_kmeans = kmeans.fit_predict(X)
13  X = np.array(X)
14  plt.scatter(X[y_kmeans == 0, 0], X[y_kmeans == 0, 1], s=10, c='red', label='low')
15  plt.scatter(X[y_kmeans == 1, 0], X[y_kmeans == 1, 1], s=10, c='blue', label='mid')
16  plt.scatter(X[y_kmeans == 2, 0], X[y_kmeans == 2, 1], s=10, c='green', label='high')
17  plt.scatter(kmeans.cluster_centers_[:, 0], kmeans.cluster_centers_[:, 1], s=100,
    c='black', label='Centroids')
18  plt.title('google play store apps')
19  plt.xlabel('Reviews')
20  plt.ylabel('Installs')
21  plt.legend()
22  plt.show()
```

size-install

```
1   X = a.iloc[:, [3, 4]]
2   # elbow
3   wcss = []
4   for i in range(1, 11):
5       kmeans = KMeans(n_clusters=i, n_init='auto', random_state=0)
6       kmeans.fit(X)
7       wcss.append(kmeans.inertia_)
8   plt.plot(range(1, 11), wcss)
9   plt.show()
10  # cluster show
11  kmeans = KMeans(n_clusters=3, init='k-means++', random_state=0)
12  y_kmeans = kmeans.fit_predict(X)
13  X = np.array(X)
14  plt.scatter(X[y_kmeans == 0, 0], X[y_kmeans == 0, 1], s=10, c='red', label='low')
15  plt.scatter(X[y_kmeans == 1, 0], X[y_kmeans == 1, 1], s=10, c='blue', label='mid')
16  plt.scatter(X[y_kmeans == 2, 0], X[y_kmeans == 2, 1], s=10, c='green', label='high')
```

```
17  plt.scatter(kmeans.cluster_centers_[:, 0], kmeans.cluster_centers_[:, 1], s=100,
    c='black', label='Centroids')
18  plt.title('google play store apps')
19  plt.xlabel('Size')
20  plt.ylabel('Installs')
21  plt.legend()
22  plt.show()
```

## ⑥預測

### 數據分類依日期分配

```
1  a['Last Updated'] = pd.to_datetime(a['Last Updated'])
2  before = a[a['Last Updated'] < '2018-07-01']
3  after = a[a['Last Updated'] >= '2018-07-01']
4  # divide data in 2 parts, traning set and test set
5  x_train = before[['Rating', 'Reviews', 'Size']]
6  y_train = before[['Installs']]
7  x_test = after[['Rating', 'Reviews', 'Size']]
8  y_test = after[['Installs']]
```

### 數據分類是電腦隨機分配

```
1  a['Last Updated'] = pd.to_datetime(a['Last Updated'])
2  x = a[['Rating', 'Reviews', 'Size']]
3  y = a[['Installs']]
4  x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.3,
    random_state=42)
```

### 線性迴歸

```
1   # train and construct model
2   regr1 = linear_model.LinearRegression()
3   regr1.fit(x_train,  y_train)
4   print("f(X1, X2, X3) = " + str(regr1.coef_[0][0]) + "X1 + " +  str(regr1.coef_[0][1]) +
    "X2 -", str(-regr1.coef_[0][2]) + "X3 + " + str(regr1.intercept_[0]))
5   # use the trained regression model to predict test dataset's target
6   y_predict = regr1.predict(x_test)
7   plt.plot(range(len(y_predict)), y_predict, 'b', label="predict")
8   plt.plot(range(len(y_predict)), y_test, 'r', label="test", alpha=0.5)
9   plt.legend(loc="upper right")
10  plt.show()
11  # accuracy
12  print("train_accuracy : ", regr1.score(x_train, y_train))
13  print("test_accuracy : ", regr1.score(x_test, y_test))
14  b = pd.DataFrame(y_predict, columns=['Predict'])
15  b.to_csv("test.csv", index=False)
```
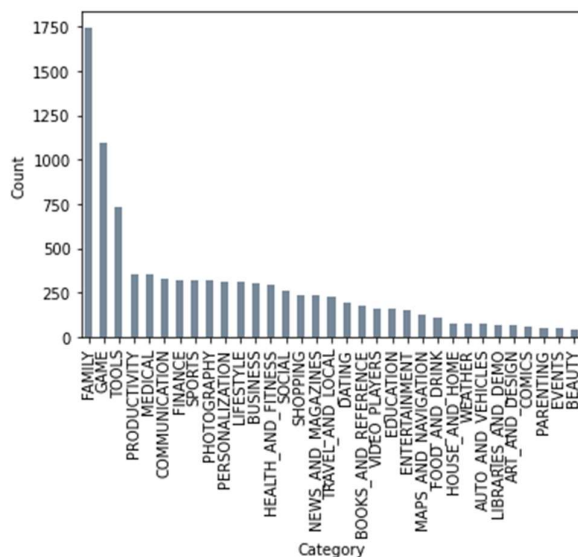
決策樹

```
1  dct = DecisionTreeRegressor()
2  dct.fit(x_train, y_train)
3  print("test_accuracy : ", dct.score(x_test, y_test))
```
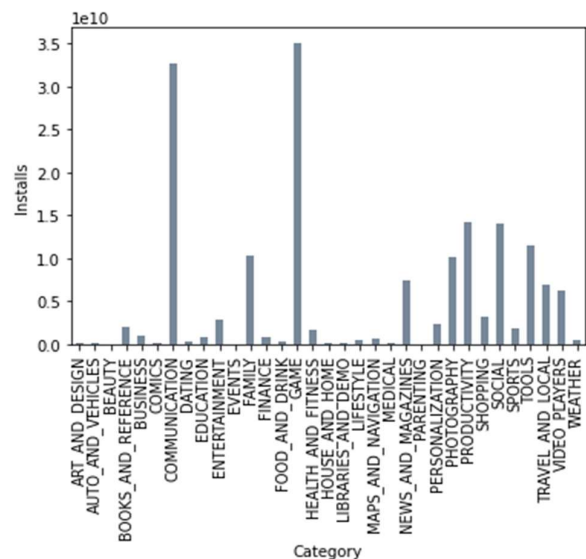
隨機森林迴歸器

```
1  Randomforestreg = RandomForestRegressor(n_estimators=100, n_jobs=-1, oob_score=True,
   bootstrap=True, random_state=42)
2  Randomforestreg.fit(x_train, y_train)
3  y_prediction_randomforest = Randomforestreg.predict(x_test)
4  randomforestscore = Randomforestreg.score(x_test, y_test)
5  print("test_accuracy : ", randomforestscore)
```
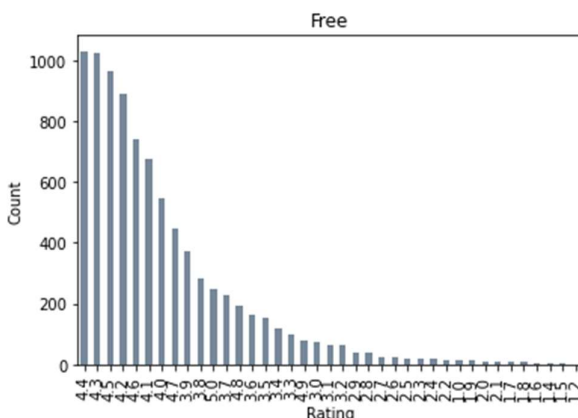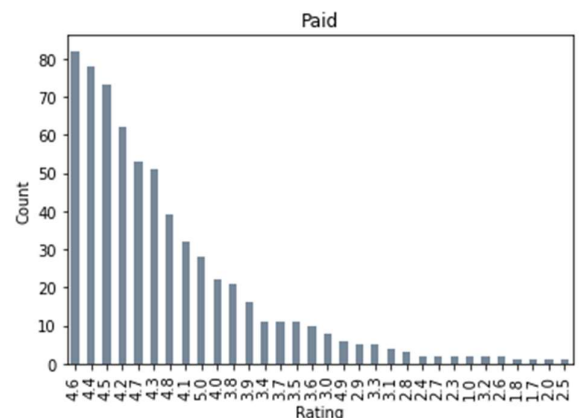
# 3. 測試結果展示

- Category-Count Bar



- Category-Install Bar



- Count-Rating in Free app



- Count-Rating in Paid app

- Free vs Paid app pie



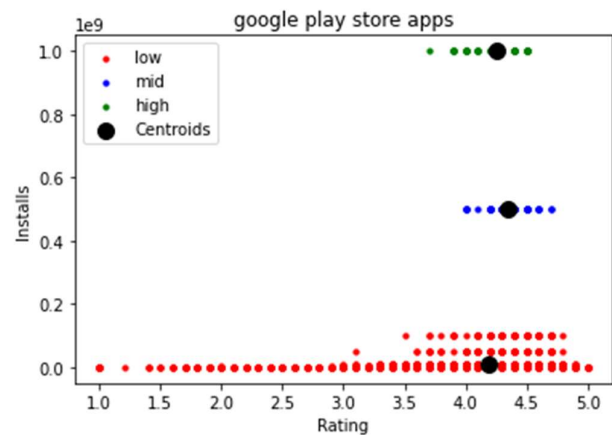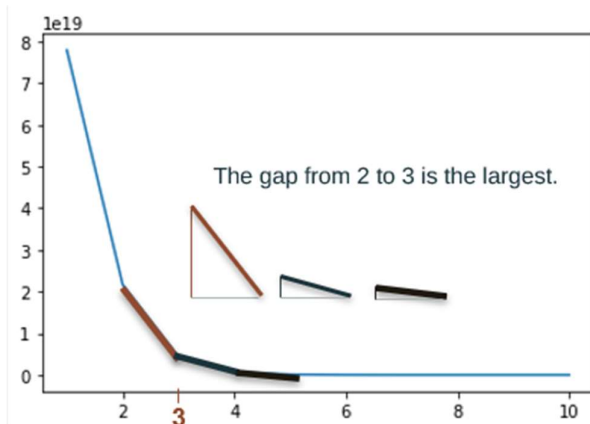## Correlation

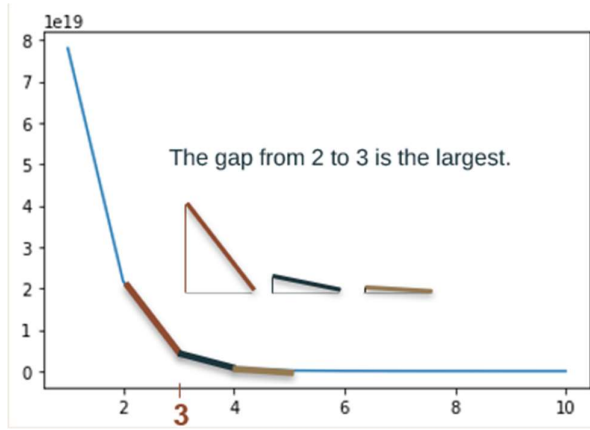我們選擇相關係數大於 0.05 的項目: Rating、Reviews、Size。



## Clustering

在每個 clustering 的圖中會發現，級距很明顯且彼此之間都沒有其他數據，這是因為數據 Installs 原本是以 '+' 代替後面的非整數，因此最後處理完的結果只會有 $5 * 10^n$ 和 $1 * 10^n$ 兩種類型的數據，如右圖。此外，圖中單位為 $10^9$，所以 low 為小於等於 $10^8$，mid 為 $5 * 10^8$，high 為大於等於 $10^9$。

```
Installs
1000000      1577
10000000     1252
100000       1150
10000        1010
5000000       752
1000          713
500000        538
50000         467
5000          432
100000000     409
100           309
50000000      289
500           201
500000000      72
10             69
1000000000     58
50             56
5               9
1               3
```
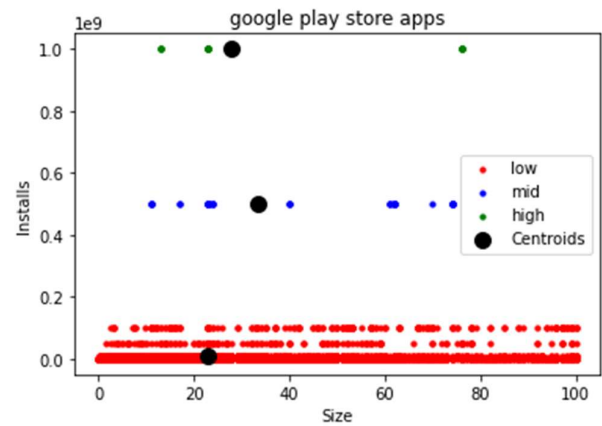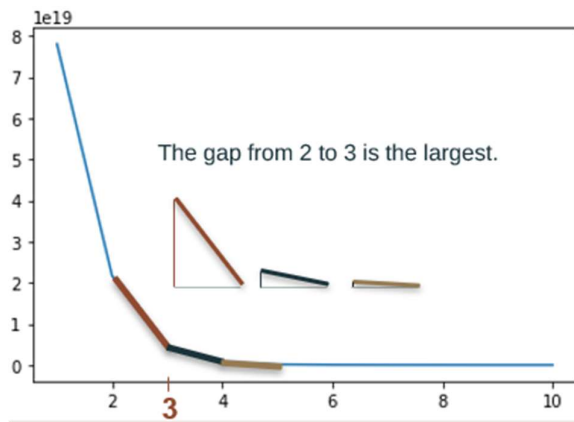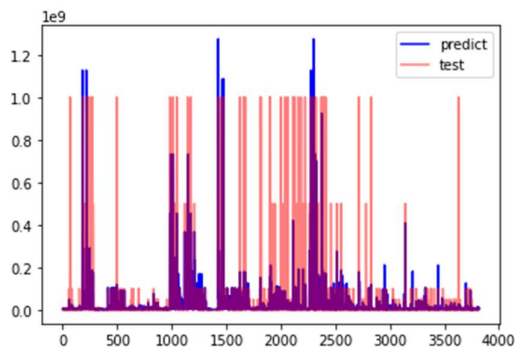
- Install-Rating

■ Install-Reviews



■ Install-Size



# Forecasting

■ Linear regression

由下圖可知 x 係數皆沒有過小的，因此沒有再進一步刪除項目而是直接預測

```
f(X1, X2, X3) = 447582.62405776535X1 + 16.268902752069355X2 - 4744.358168040289X3 + 699848.9572047135
```



✓ 依日期分類預測

```
train_accuracy :  0.34033503614315164
test_accuracy :   0.38172182613702543
```

✓ 電腦隨機分類預測(此準確率並非固定值，但平均大於 70%)

```
test_accuracy :   0.8499635555700153
```

- **Decision tree**
  - ✓ 依日期分類預測

    `test_accuracy :   0.2389810917196915`

  - ✓ 電腦隨機分類預測(此準確率並非固定值，但平均大於 70%)

    `test_accuracy :   0.857678394220075`

- **Random forest regressor**
  - ✓ 依日期分類預測

    `test_accuracy :   0.2521013309986284`

  - ✓ 電腦隨機分類預測(此準確率並非固定值，但平均大於 70%)

    `test_accuracy :   0.8622321162661485`

# 4. 心得與反思

　　做完此報告後才理解到老師所說的軍火展示跟真正分析的差別，原本我們的想法是偏於軍火展示，經過與老師的討論後才明白要確定目的並為此目的篩選所需項目，但那時已交完計畫報告，所以聚類分析即使後來覺得好像沒什麼用依然選擇做出來。預測的部分，由於準確率很低，所以我們又嘗試了其他種方法 (decision tree、random forest) 去提升它的準確率，但效果皆不顯著。除非測試集不分日期而是採電腦隨機分配，準確率才會超過 50%。對此，我們覺得數據集可能選錯了，我們到最後才發現 Install 這個數據並不是一個準確數據而是一個範圍數據，所以最後呈現的結果看起來不正常。經過這學期的實作練習，雖然最後專題結果不是很好，但對於大數據分析有比基本了解再更深入了，也學到一些網路爬蟲技巧以及分析數據相關技巧。