



Big Data Analytic project

Google Play Store Apps Analysis

411021206 莊芳儒、411021209 張佩慈、411021231 林芷安

Table of contents



01 Introduction

Introduce the Dataset

02 Analysis

Data Preprocessing and Chart Analyzing

03 Forecasting

Correlation Analysis and Downloads Prediction

04 Demonstration

Run the Code









01

Introduction

Introduce the dataset



Data selection






LAVANYA · UPDATED 5 YEARS AGO

Google Play Store Apps

Web scraped data of 10k Play Store apps for analysing the Android market.

Data Explorer

Version 6 (9.03 MB)

-  googleplaystore.csv
-  googleplaystore_user_reviews.csv
-  license.txt

Data feature

Googleplaystore.csv

<u>App</u>	The name of the app.	<u>Price</u>	The price of the app.
<u>Category</u>	The category of the app.	<u>Content Rating</u>	The appropriate target audience of the app.
<u>Rating</u>	The rating of the app.	<u>Genres</u>	The genre of the app.
<u>Reviews</u>	The number of reviews of the app.	<u>Last Updated</u>	The date when the app was last updated.
<u>Size</u>	The size of the app.	<u>Current Ver</u>	The current version of the app.
<u>Install</u>	The number of installs of the app.	<u>Android Ver</u>	The minimum Android version required to run the app.
<u>Type</u>	Free or Paid.		

1	App	Category	Rating	Reviews	Size	Installs	Type	Price	Content Rating	Genres	Last Updated	Current Ver	Android Ver
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

Download

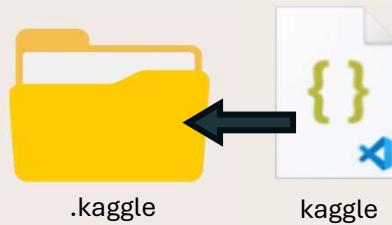
```
PS E:\code\python> pip install kaggle
```

API

Using Kaggle's beta API, you can interact with Competitions and Datasets to download data, make submissions, and more via the command line. [Read the docs](#)

Create New Token

Expire Token



Download (2 MB)

+ New notebook

Bookmark


Copy API command

Social share

Report issue

Download

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



02

Analysis

data preprocessing and chart analyzing.



Python Package Index

```
1 import numpy as np
2 import pandas as pd
3 import seaborn as sn
4 import matplotlib.pyplot as plt
5 from sklearn.cluster import KMeans
6 from sklearn import linear_model
7 import zipfile
8 import os
```

Data Preprocessing

Cleaning

Remove the missing value

1	App	Category
10474	Life Made WI-Fi Touchscreen Photo Frame	1.9
10475	osmino Wi-Fi: free WiFi	TOOLS

```
1 data.drop(index = 10472, inplace = True)
2 a = data.dropna()
```

Feature 'Installs'

Remove '+' and ',' then convert string to number

Installs
10,000+
10,000+

```
1 a.Installs = a.Installs.apply(lambda x: x.replace('+', ''))
2 a.Installs = a.Installs.apply(lambda x: x.replace(',', ''))
3 a.Installs = pd.to_numeric(a.Installs)
```

Data Preprocessing

Feature 'Price'

Remove '\$' then convert string to number

Price
\$1.49
0

```
1 a['Price'] = a['Price'].apply(lambda x: x.replace('$', ''))  
2 a['Price'] = pd.to_numeric(a['Price'])
```

Feature 'Reviews'

Convert string to number

```
1 a['Reviews'] = pd.to_numeric(a['Reviews'])
```

Data Preprocessing

Feature 'Size'

Change the unit to MB

Size
19M
14M

```
1 def convert(size):
2     if type(size) == str:
3         if 'k' in size:
4             return format(float(str(size).replace('k', '')) / 1024, '.3f')
5         if 'M' in size:
6             return float(str(size).replace('M', ''))
7     else:
8         return(size)
9 a['Size'] = a['Size'].str.replace('Varies with device', '0.00')
10 a['Size'] = a['Size'].apply(lambda x: convert(x))
11 a['Size'] = pd.to_numeric(a['Size'])
12 a['Size'].fillna(format(a['Size'].mean(), '.3f'), inplace = True )
13 a['Size'] = pd.to_numeric(a['Size'])
```

Chart Analysis

Category-Count bar

Display the number of apps in each category

```
1 value_count = a['Category'].value_counts()
2 value_count.plot(kind = 'bar',
3                   color = 'lightslategray')
4 plt.ylabel('Count')
5 plt.show()
```

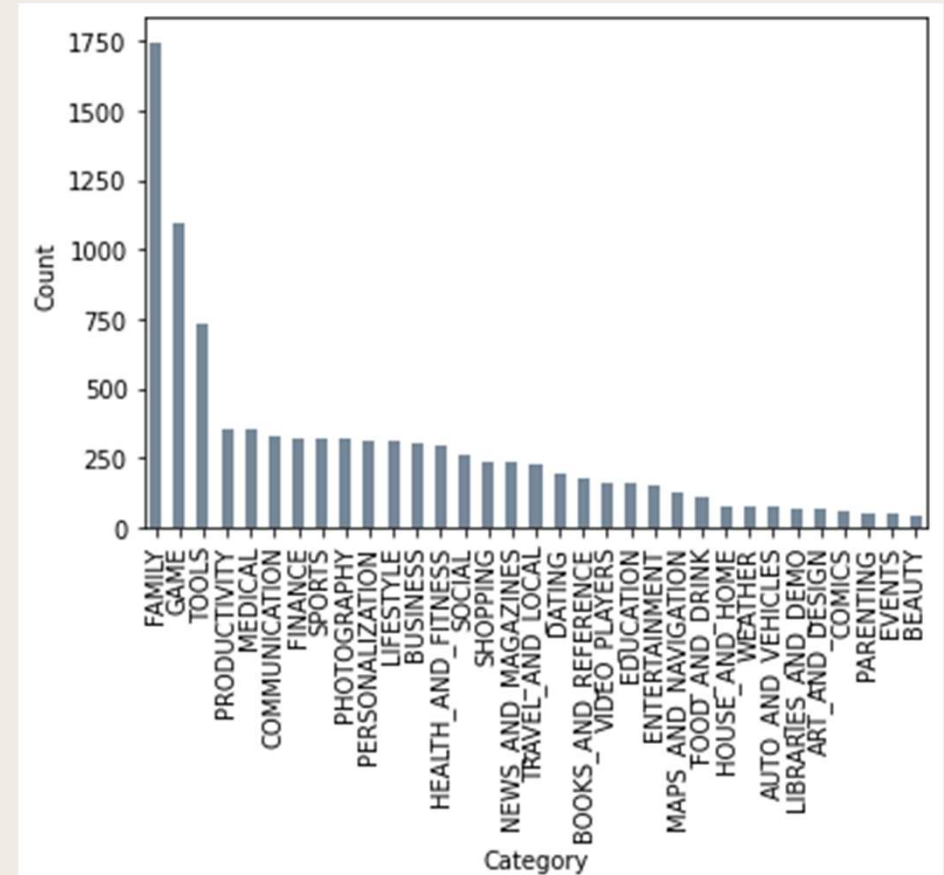


Chart Analysis

Category-Install bar

Display the downloads in each category

```
1 S = a.groupby('Category')['Installs'].sum()
2 S.plot.bar(color = 'lightslategray')
3 plt.ylabel('Installs')
4 plt.show()
```

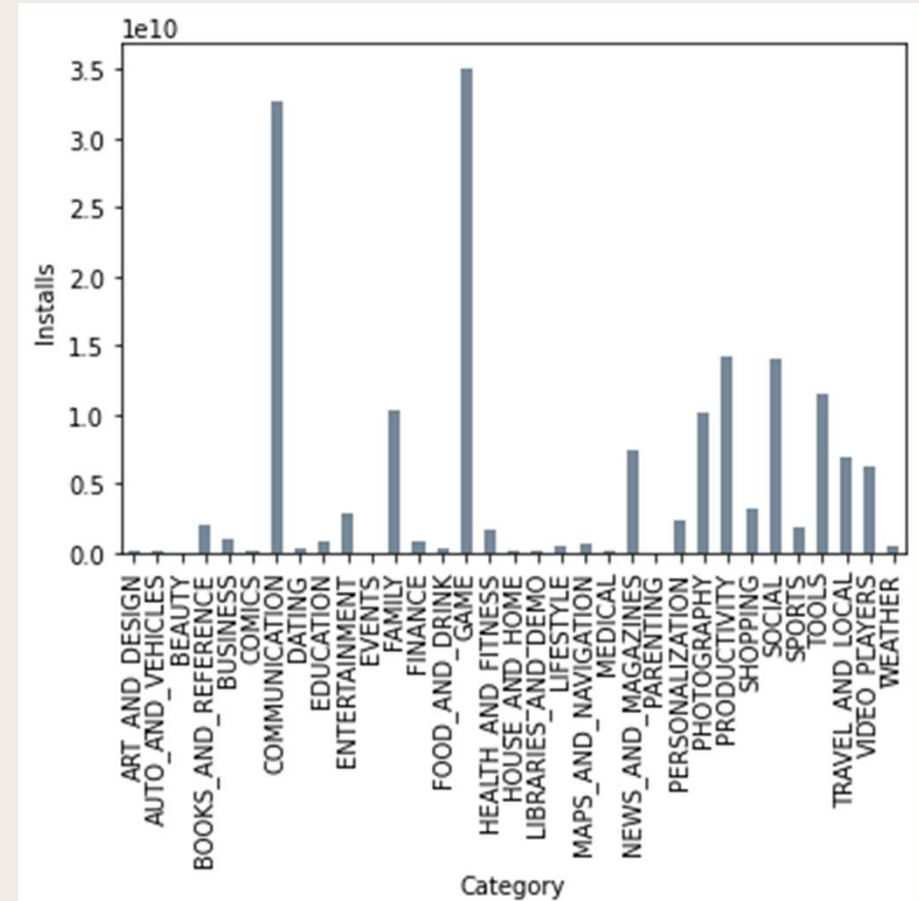


Chart Analysis

Free vs Paid app pie

Display the number of apps in free/paid.

```
1 C_Freecnt = C_Freec.sum()
2 C_Paidcnt = C_Paidc.sum()
3 x = [C_Freecnt, C_Paidcnt]
4 name = ["Free", "Paid"]
5 color = ['lightslategray', 'rosybrown']
6 plt.pie(x, autopct = '%.3f%%', colors = color,
7         labels = name)
8 plt.show()
```

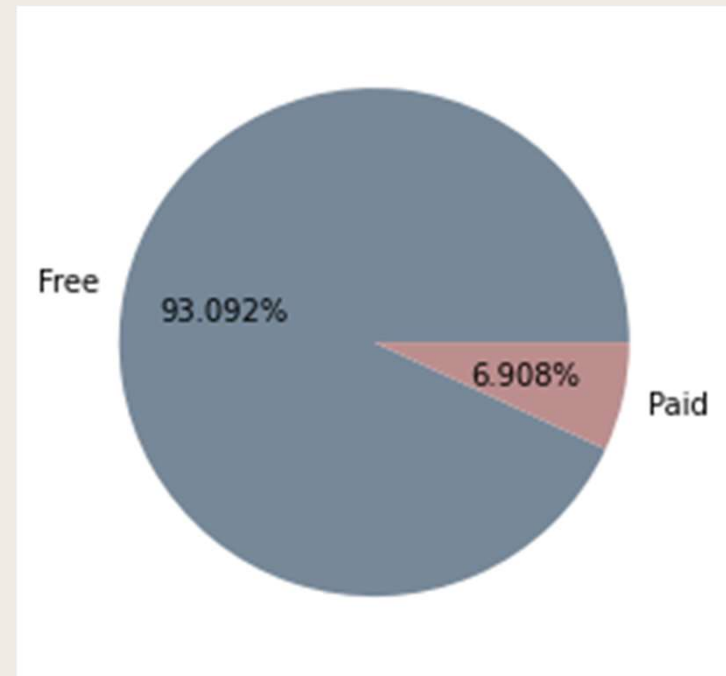


Chart Analysis

Count-Rating in Free app

Display the number of rating in free app.

```
1 C_Free = a.groupby('Type').get_group('Free')
2 C_Freec = C_Free['Rating'].value_counts()
3 C_Freec.plot.bar(color = 'lightslategray')
4 plt.ylabel('Count')
5 plt.title('Free')
6 plt.show()
```

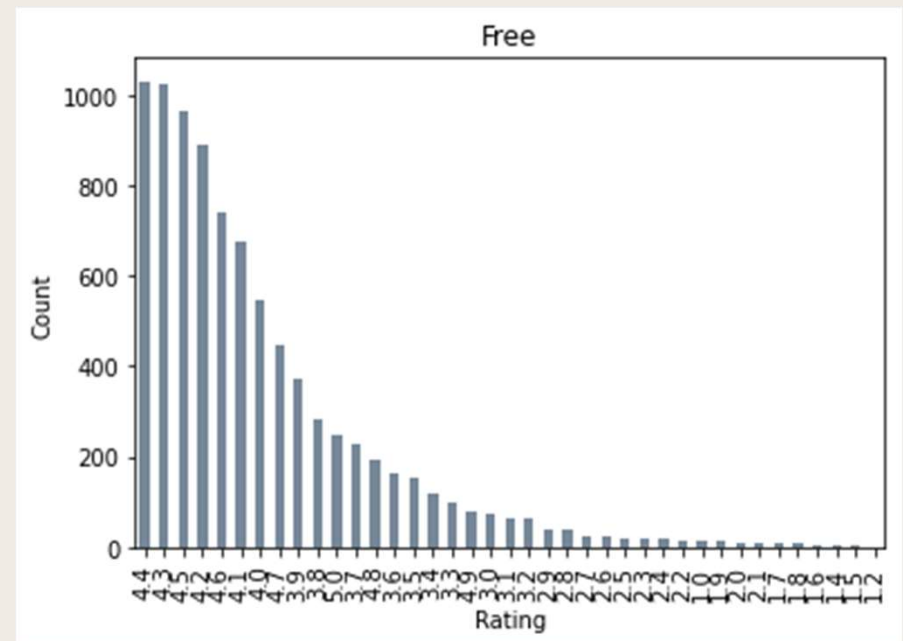


Chart Analysis

Count-Rating in Paid app

Display the number of rating in paid app.

```
1 C_Paid = a.groupby('Type').get_group('Paid')
2 C_Paidc = C_Paid['Rating'].value_counts()
3 C_Paidc.plot.bar(color = 'lightslategray')
4 plt.ylabel('Count')
5 plt.title('Free')
6 plt.show()
```

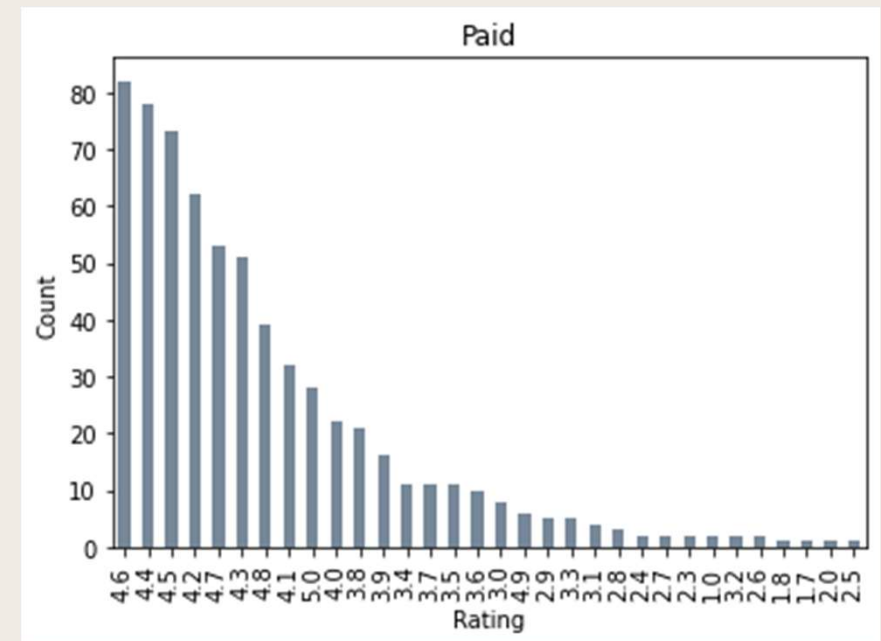
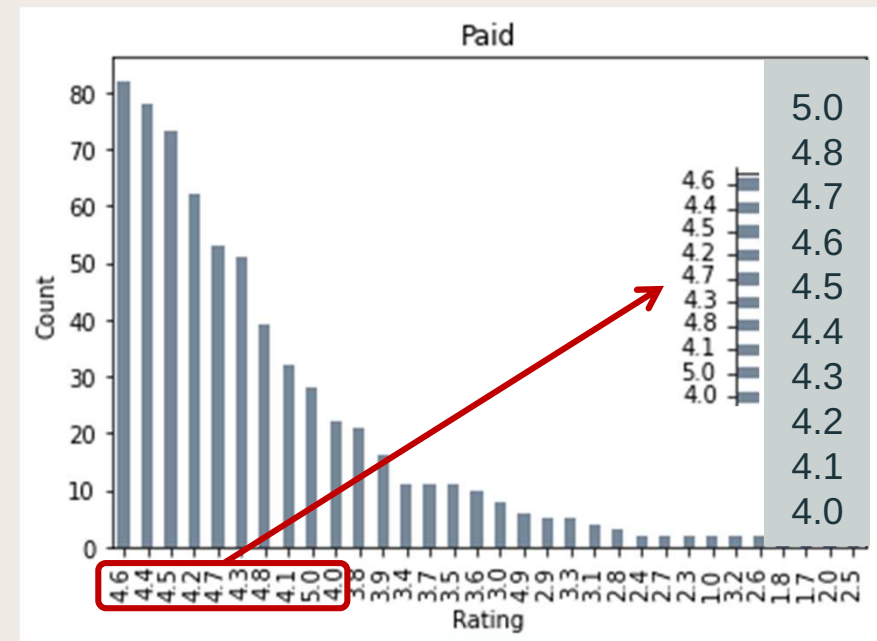
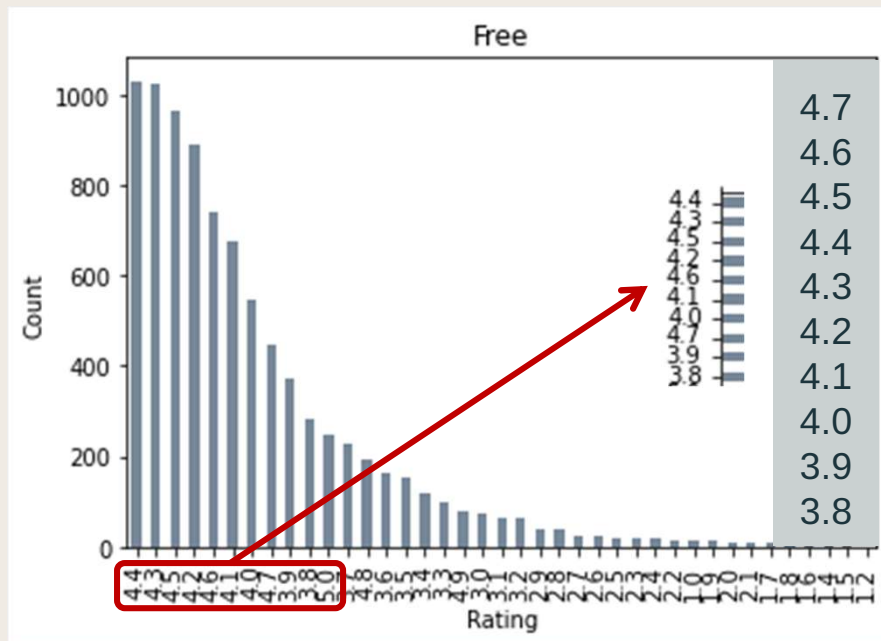


Chart Analysis

Count-Rating: Free vs Paid app





03

Forecasting

Correlation analysis and downloads prediction.



Correlation analysis

Filtering feature

We select features with correlation coefficients greater than 0.05.

```
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()
```

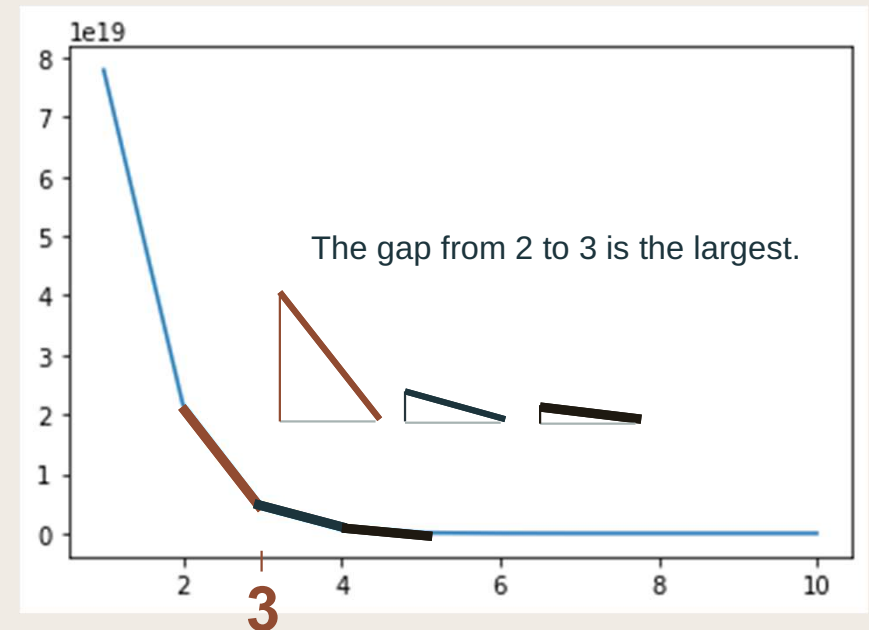


Cluster

Rating-Installs

Determine the number of clusters.

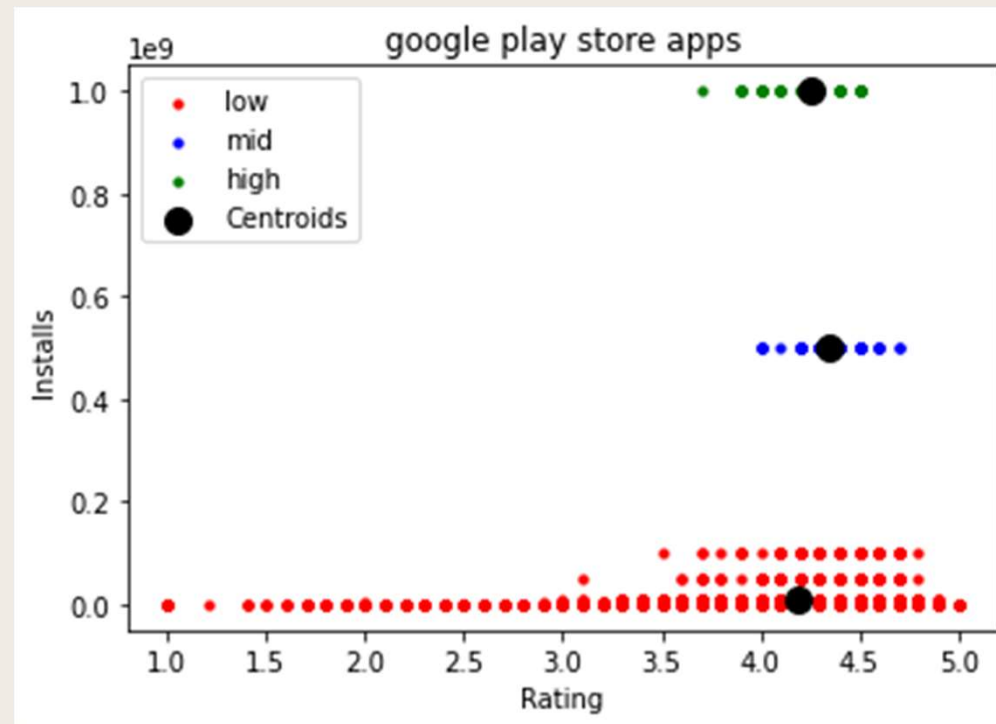
```
1 X = a.iloc[:, [1, 4]]
2 wcss = []
3 for i in range(1, 11):
4     kmeans = KMeans(n_clusters = i,
5                     n_init = 'auto', random_state = 0)
6     kmeans.fit(X)
7     wcss.append(kmeans.inertia_)
8 plt.plot(range(1, 11), wcss)
9 plt.show()
```



Cluster

```
1 kmeans = KMeans(n_clusters = 3, init = 'k-means++', random_state = 0)
2 y_kmeans = kmeans.fit_predict(X)
3 X = np.array(X)
4 plt.scatter(X[y_kmeans == 0, 0], X[y_kmeans == 0, 1], s = 10, c = 'red', label = 'low')
5 plt.scatter(X[y_kmeans == 1, 0], X[y_kmeans == 1, 1], s = 10, c = 'blue', label = 'mid')
6 plt.scatter(X[y_kmeans == 2, 0], X[y_kmeans == 2, 1], s = 10, c = 'green', label = 'high')
7 plt.scatter(kmeans.cluster_centers_[0, 0], kmeans.cluster_centers_[0, 1],
8             s = 100, c = 'black', label = 'Centroids')
9 plt.title('google play store apps')
10 plt.xlabel('Rating')
11 plt.ylabel('Installs')
12 plt.legend()
13 plt.show()
```

Cluster

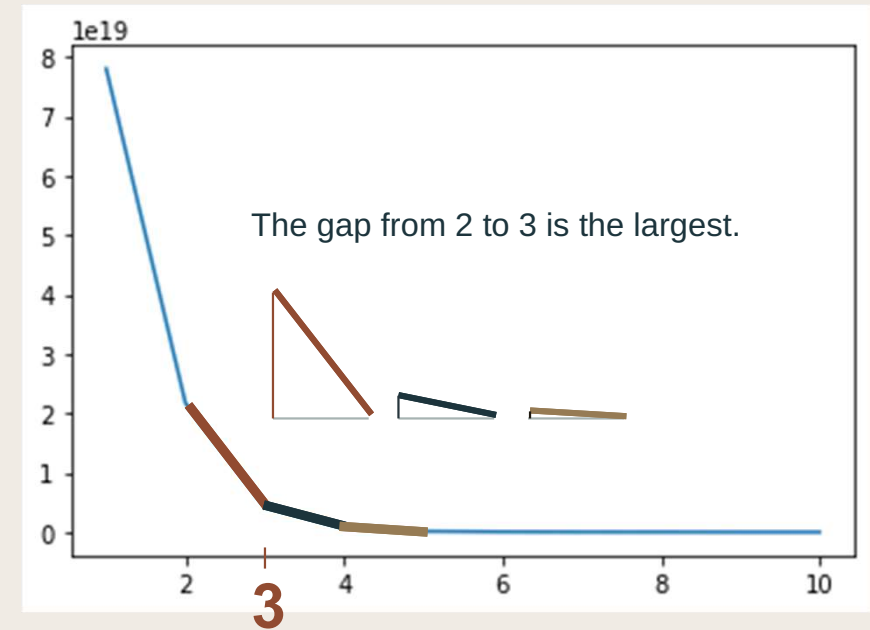


Cluster

Reviews-Installs

Determine the number of clusters.

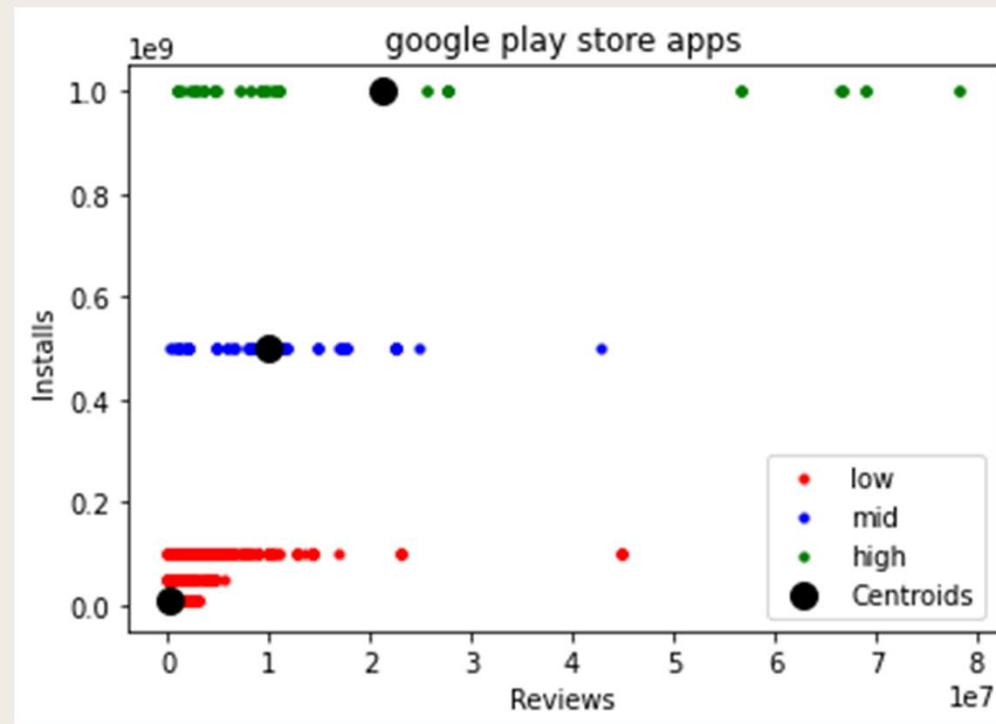
```
1 X = a.iloc[:, [2, 4]]
2 wcss = []
3 for i in range(1, 11):
4     kmeans = KMeans(n_clusters = i,
5                     n_init = 'auto', random_state = 0)
6     kmeans.fit(X)
7     wcss.append(kmeans.inertia_)
8 plt.plot(range(1, 11), wcss)
9 plt.show()
```



Cluster

```
1 kmeans = KMeans(n_clusters = 3, init = 'k-means++', random_state = 0)
2 y_kmeans = kmeans.fit_predict(X)
3 X = np.array(X)
4 plt.scatter(X[y_kmeans == 0, 0], X[y_kmeans == 0, 1], s = 10, c = 'red', label = 'low')
5 plt.scatter(X[y_kmeans == 1, 0], X[y_kmeans == 1, 1], s = 10, c = 'blue', label = 'mid')
6 plt.scatter(X[y_kmeans == 2, 0], X[y_kmeans == 2, 1], s = 10, c = 'green', label = 'high')
7 plt.scatter(kmeans.cluster_centers_[0, 0], kmeans.cluster_centers_[0, 1],
8             s = 100, c = 'black', label = 'Centroids')
9 plt.title('google play store apps')
10 plt.xlabel('Reviews')
11 plt.ylabel('Installs')
12 plt.legend()
13 plt.show()
```

Cluster

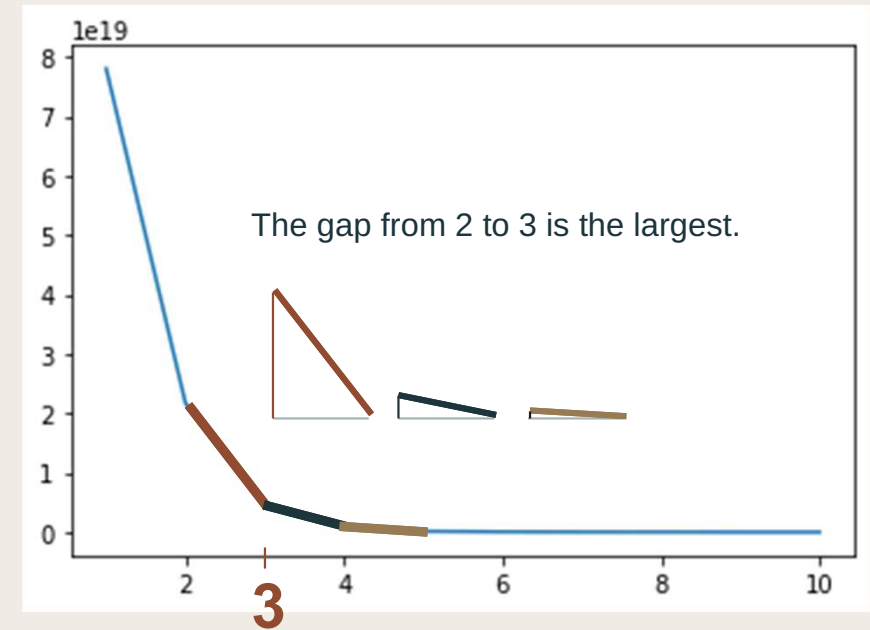


Cluster

Size-Installs

Determine the number of clusters.

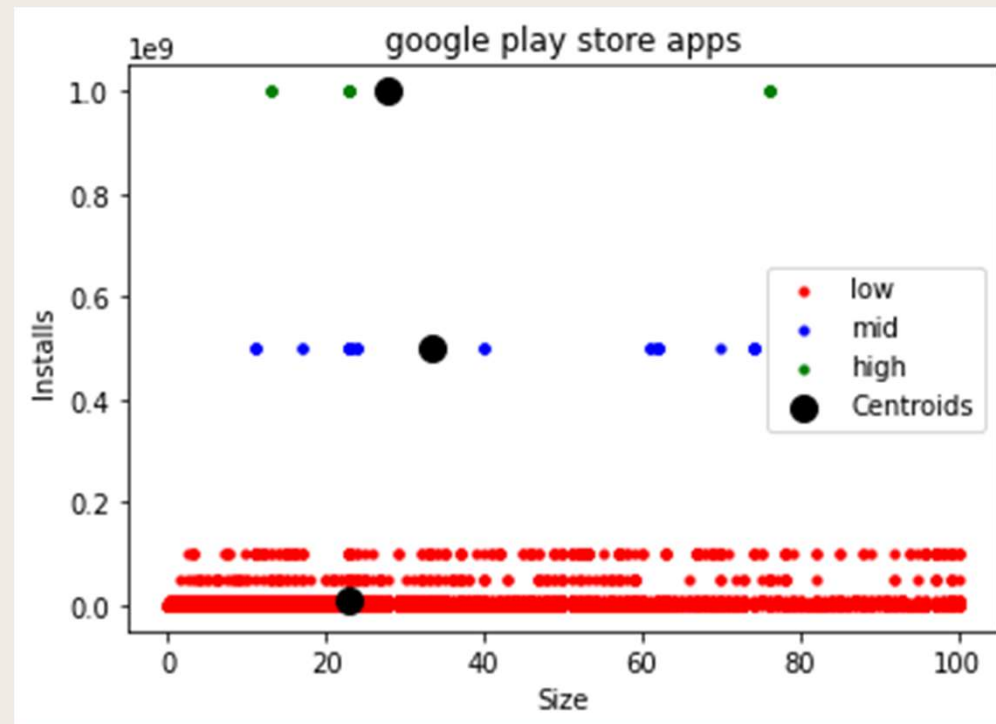
```
1 X = a.iloc[:, [2, 4]]
2 wcss = []
3 for i in range(1, 11):
4     kmeans = KMeans(n_clusters = i,
5                     n_init = 'auto', random_state = 0)
6     kmeans.fit(X)
7     wcss.append(kmeans.inertia_)
8 plt.plot(range(1, 11), wcss)
9 plt.show()
```



Cluster

```
1 kmeans = KMeans(n_clusters = 3, init = 'k-means++', random_state = 0)
2 y_kmeans = kmeans.fit_predict(X)
3 X = np.array(X)
4 plt.scatter(X[y_kmeans == 0, 0], X[y_kmeans == 0, 1], s = 10, c = 'red', label = 'low')
5 plt.scatter(X[y_kmeans == 1, 0], X[y_kmeans == 1, 1], s = 10, c = 'blue', label = 'mid')
6 plt.scatter(X[y_kmeans == 2, 0], X[y_kmeans == 2, 1], s = 10, c = 'green', label = 'high')
7 plt.scatter(kmeans.cluster_centers_[0, 0], kmeans.cluster_centers_[0, 1], s = 100,
8             c = 'black', label = 'Centroids')
9 plt.title('google play store apps')
10 plt.xlabel('Size')
11 plt.ylabel('Installs')
12 plt.legend()
13 plt.show()
```

Cluster



Linear Regression

```
1 # split data into 2 parts, before and after 2018.07
2 a['Last Updated'] = pd.to_datetime(a['Last Updated'])
3 before = a[a['Last Updated'] < '2018-07-01']
4 after = a[a['Last Updated'] >= '2018-07-01']
5
6 # select feature and put into training set and test set
7 x_train = before[['Rating', 'Reviews', 'Size']]
8 y_train = before[['Installs']]
9 x_test = after[['Rating', 'Reviews', 'Size']]
10 y_test = after[['Installs']]
```

Linear Regression

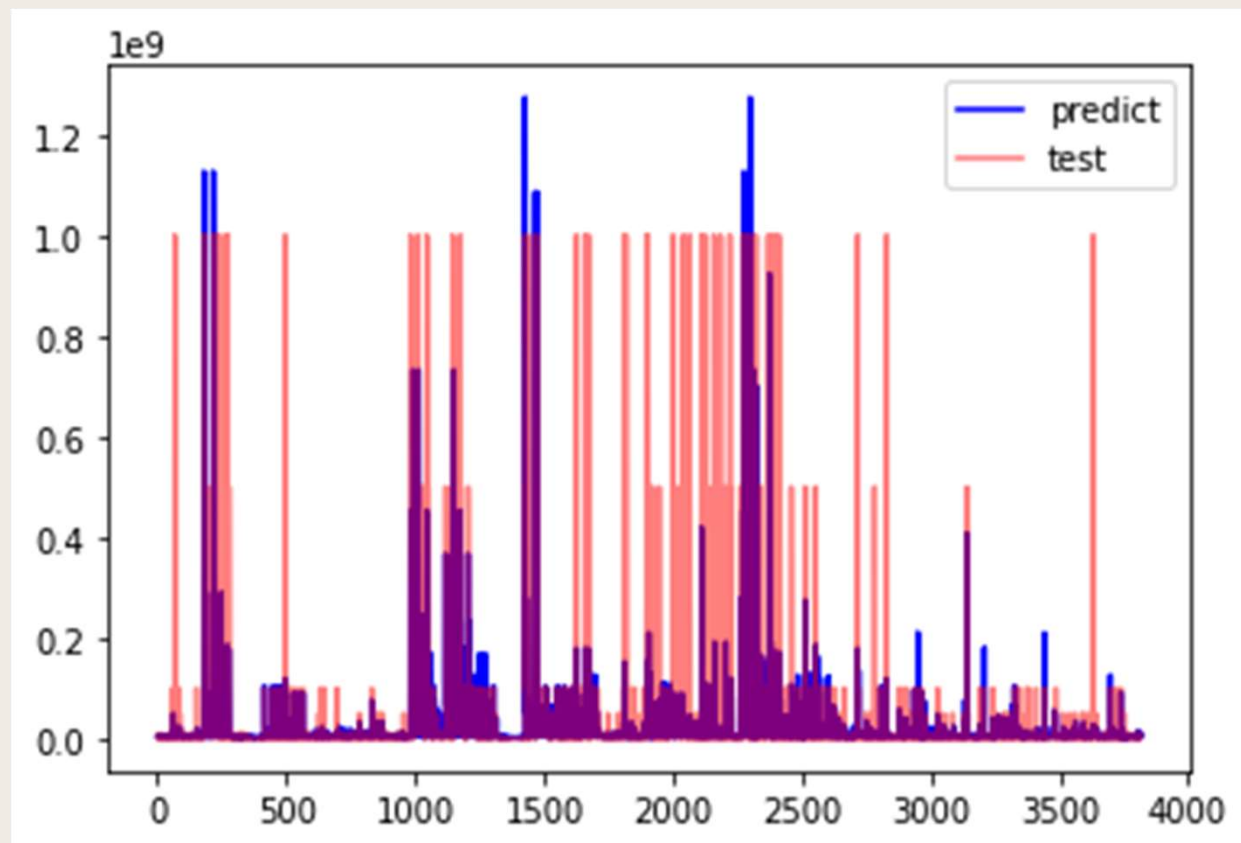
```
1 # train and construct model
2 regr1 = linear_model.LinearRegression()
3 regr1.fit(x_train,y_train)
4 print("f(X1, X2, X3) = " + str(regr1.coef_[0][0]) + "X1 + " + str(regr1.coef_[0][1]) + "X2 -",
5       str(-regr1.coef_[0][2]) + "X3 + " + str(regr1.intercept_[0]))
```

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

Linear Regression

```
1 # use the trained regression model to predict test dataset's target
2 y_predict = regr1.predict(x_test)
3 plt.plot(range(len(y_predict)), y_predict, 'b', label = "predict")
4 plt.plot(range(len(y_predict)), y_test, 'r', label = "test", alpha = 0.5)
5 plt.legend(loc = "upper right")
6 plt.show()
```


Prediction



Prediction

```
1 print("train_accuracy : ", regr1.score(x_train, y_train)) # accuracy
2 print("test_accuracy : ", regr1.score(x_test, y_test))
```


Output :

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

```
1 b = pd.DataFrame(y_predict, columns=['Predict'])
2 b.to_csv("test.csv", index = False)
```

Output :

```
[4185903.05404785]
[2790748.31718829]
[2755251.32242758]
...
[2650401.53661713]
[2920747.4636996 ]
[9103845.80874038]
```



04

Demonstration

Run the Code





Thanks!

