

Data Science Portofolio

Mohamad Adhikasurya Haidar

Universitas Gadjah Mada

+62 82115529411

mohamadadhikasurya2019@mail.ugm.ac.id

Hello everyone, in this portofolio I would like to present some of my study from “Machine Learning Projects” books written by Muhammad Ardi. There are 2 machine learning project that i’d like to show which consist of :

- **Titanic Survival Analysis**
- **Customer Segmentation Using K-Means**

Titanic Survival Analysis

In this Analysis , i’m trying to make a model that can predict if any new passenger would survive or not. But first, let’s do an Exploratory Data Analysis over the data.

1. Let’s check if there’s any null values

```
In [8]: 1 df.isnull().sum()
```

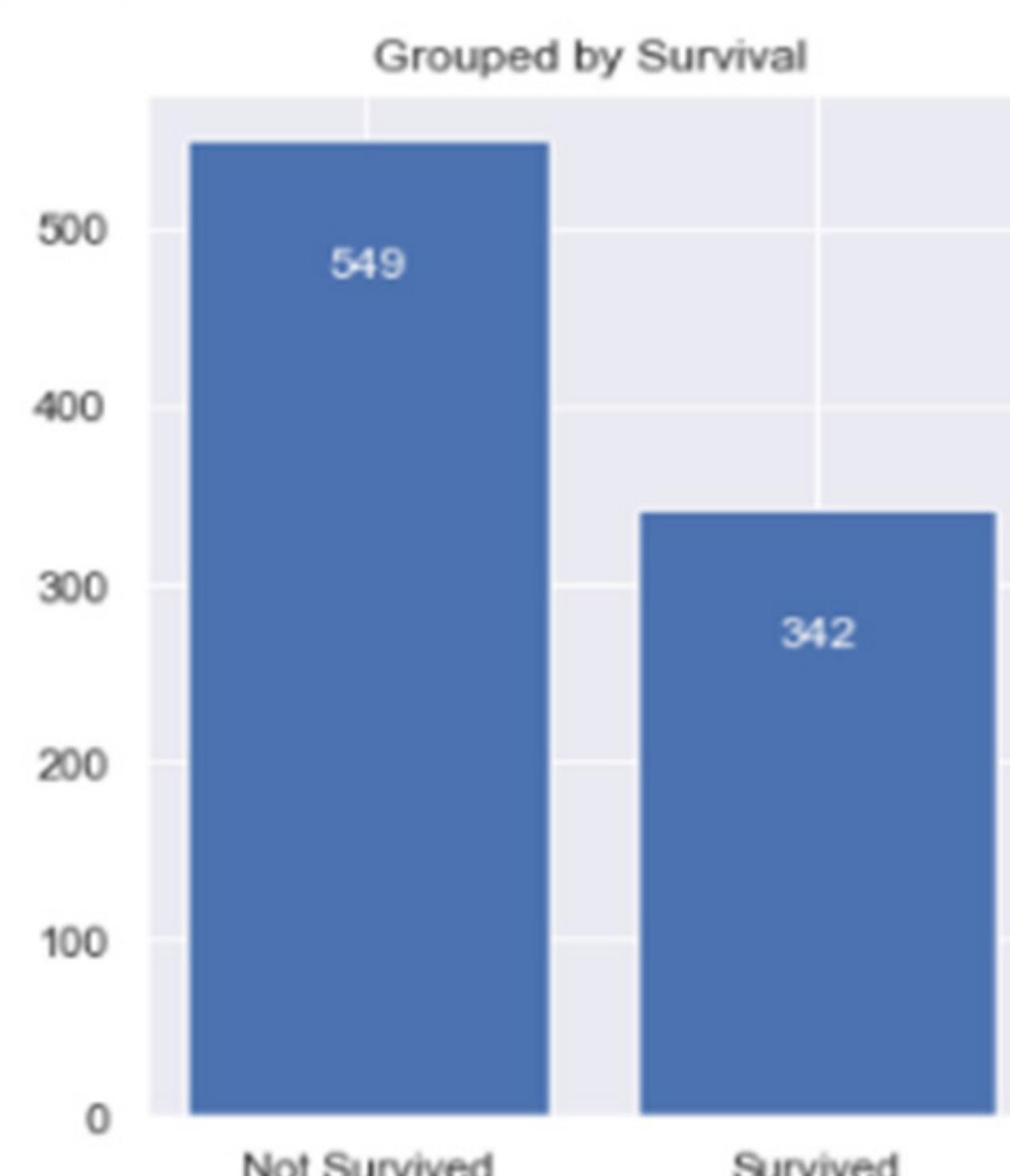
```
Out[8]: PassengerId      0
Survived        0
Pclass         0
Name          0
Sex           0
Age       177
SibSp         0
Parch         0
Ticket        0
Fare          0
Cabin       687
Embarked      2
dtype: int64
```

By using `.isnull().sum()` we can check the total number of null values in each of feature we have in the dataframe. As we can see, there’s a lot of null values in Age and Cabin values. We may have to replace those values for modelling purposes. On the other hand, since there’re only 2 null values in Embarked feature, we may just drop the rows which consist those null values in Embarked feature. We’ll process these null values in Data Preparation for modelling.

2. Make a bar plot to see the total of passengers that survived and doesn’t.

```
1 survived_count = df['Survived'].value_counts()
2 survived_count
3 plt.figure(figsize=(4,5))
4 plt.bar(survived_count.index,survived_count.values)
5 plt.title('Grouped by Survival')
6 plt.xticks([0,1],['Not Survived','Survived'])

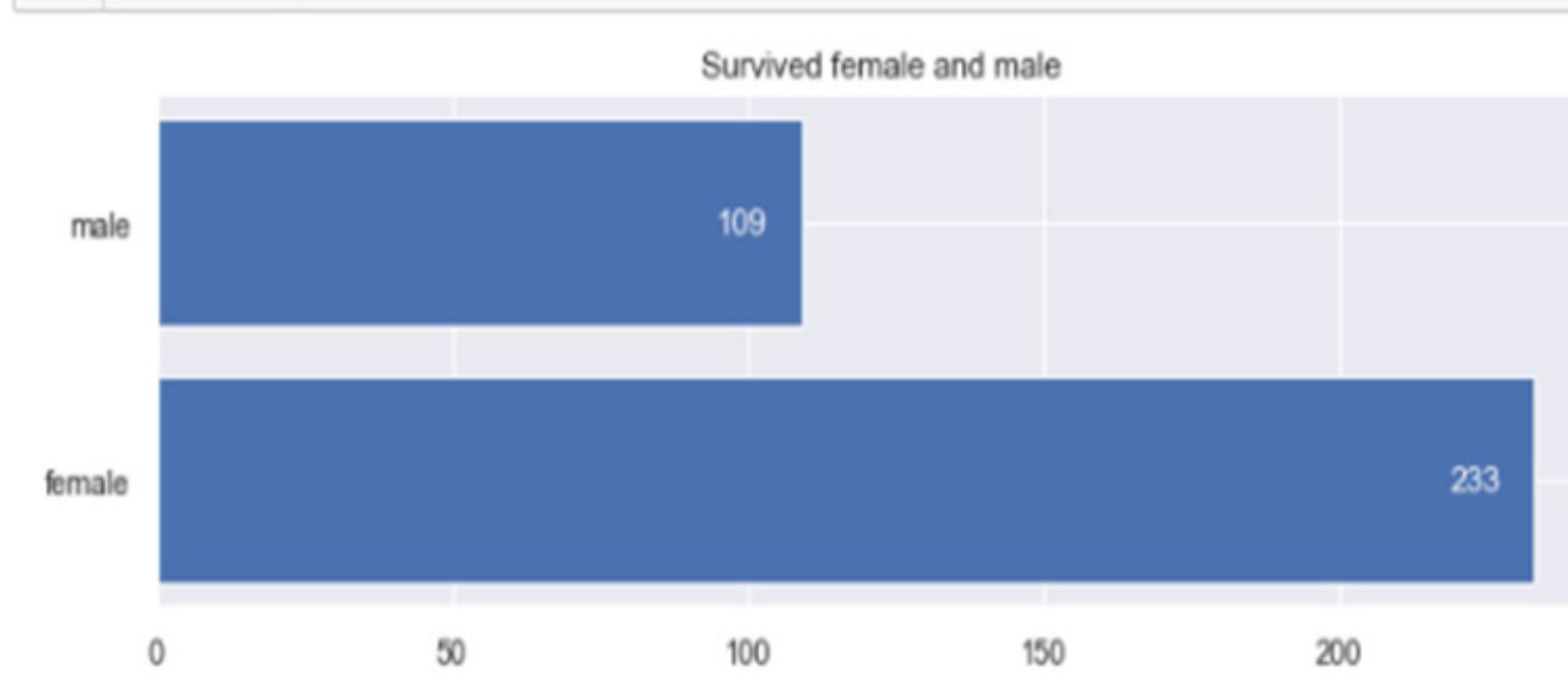
9 for i, value in enumerate(survived_count.values):
10     plt.text(i, value-70, s = str(value), fontsize=12, color='white',
11             horizontalalignment='center', verticalalignment='center')
12 plt.show()
```



We can tell from the bar plot that there are less people who survived than those who doesn’t.

3. Passengers survival based on genders

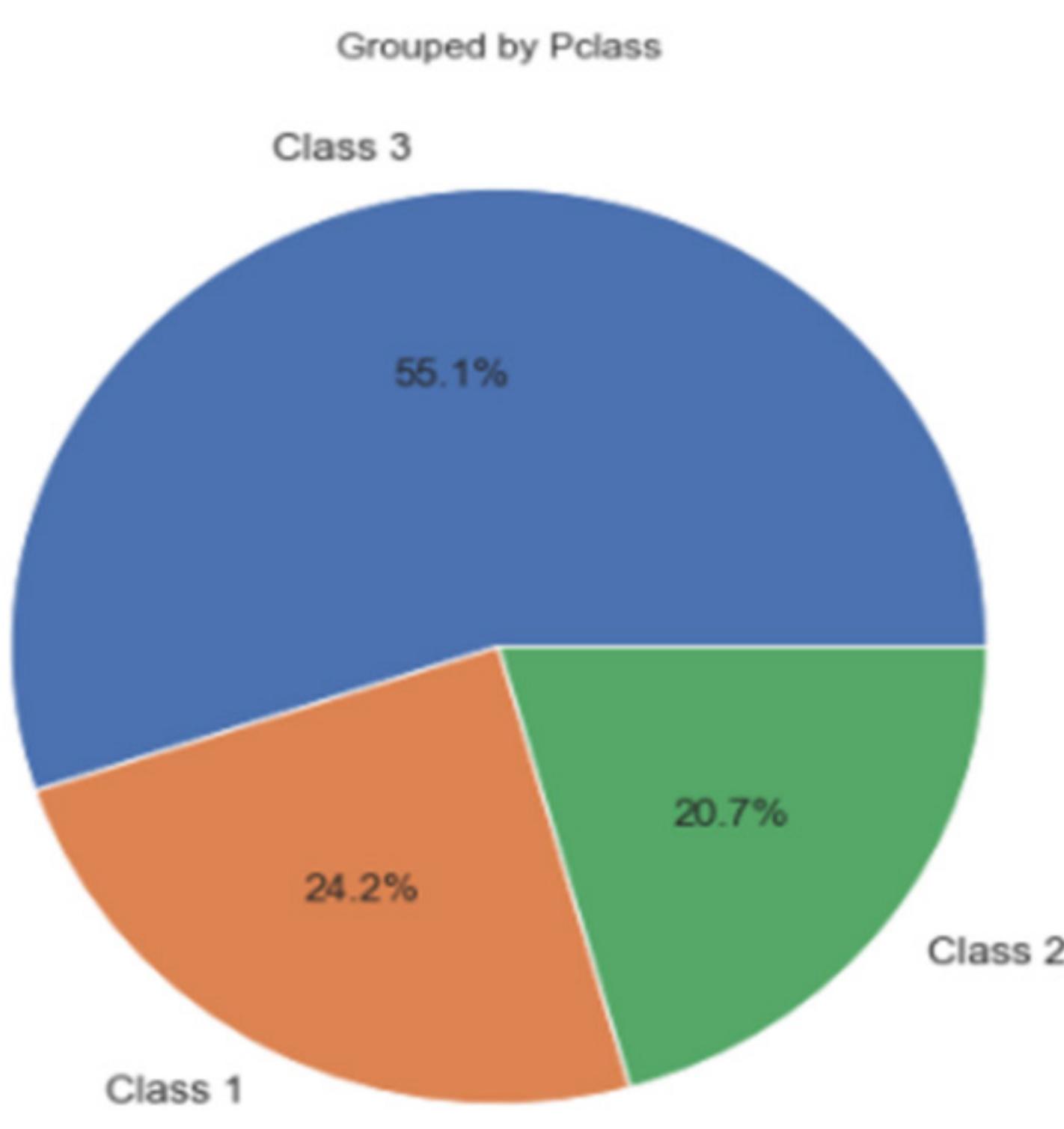
```
In [13]: 1 survived_sex = df[df['Survived']==1]['Sex'].value_counts()
2
3 plt.figure(figsize=(10,3))
4 plt.barh(survived_sex.index,survived_sex.values)
5 plt.title('Survived female and male')
6
7 for i, value in enumerate(survived_sex.values):
8     plt.text(value-10, i, s = str(value), fontsize=12, color='white',
9             horizontalalignment='center', verticalalignment='center')
10 plt.show()
```



We can tell from graph that there are more female survivors than man. It’s kinda obvious because women and children safety is top priority.

4. Ticket class distribution

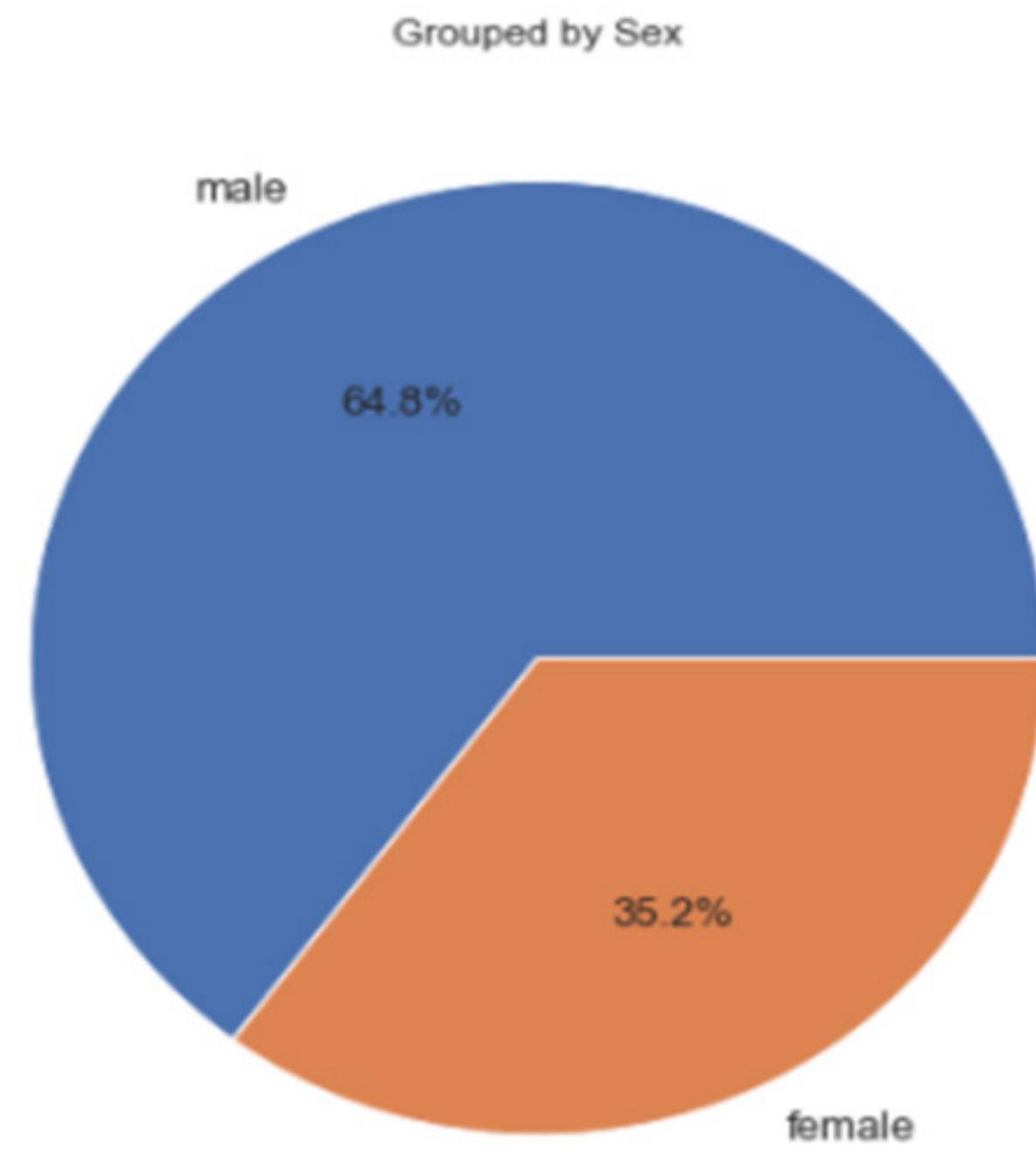
```
1 pclass_count = df['Pclass'].value_counts()
2 pclass_count
3 plt.figure(figsize=(7,7))
4 plt.title('Grouped by Pclass')
5 plt.pie(pclass_count.values,
6         labels = ['Class {}'.format(i) for i in pclass_count.index],
7         autopct='%.1f%%', textprops={'fontsize':13})
8 plt.show()
```



From the pie plot above, we can tell that most of the passengers in titanic ship bought class 3 tickets. It could be because people want to save their budget, or simply couldn't afford higher class ticket.

5. Passengers distribution by gender

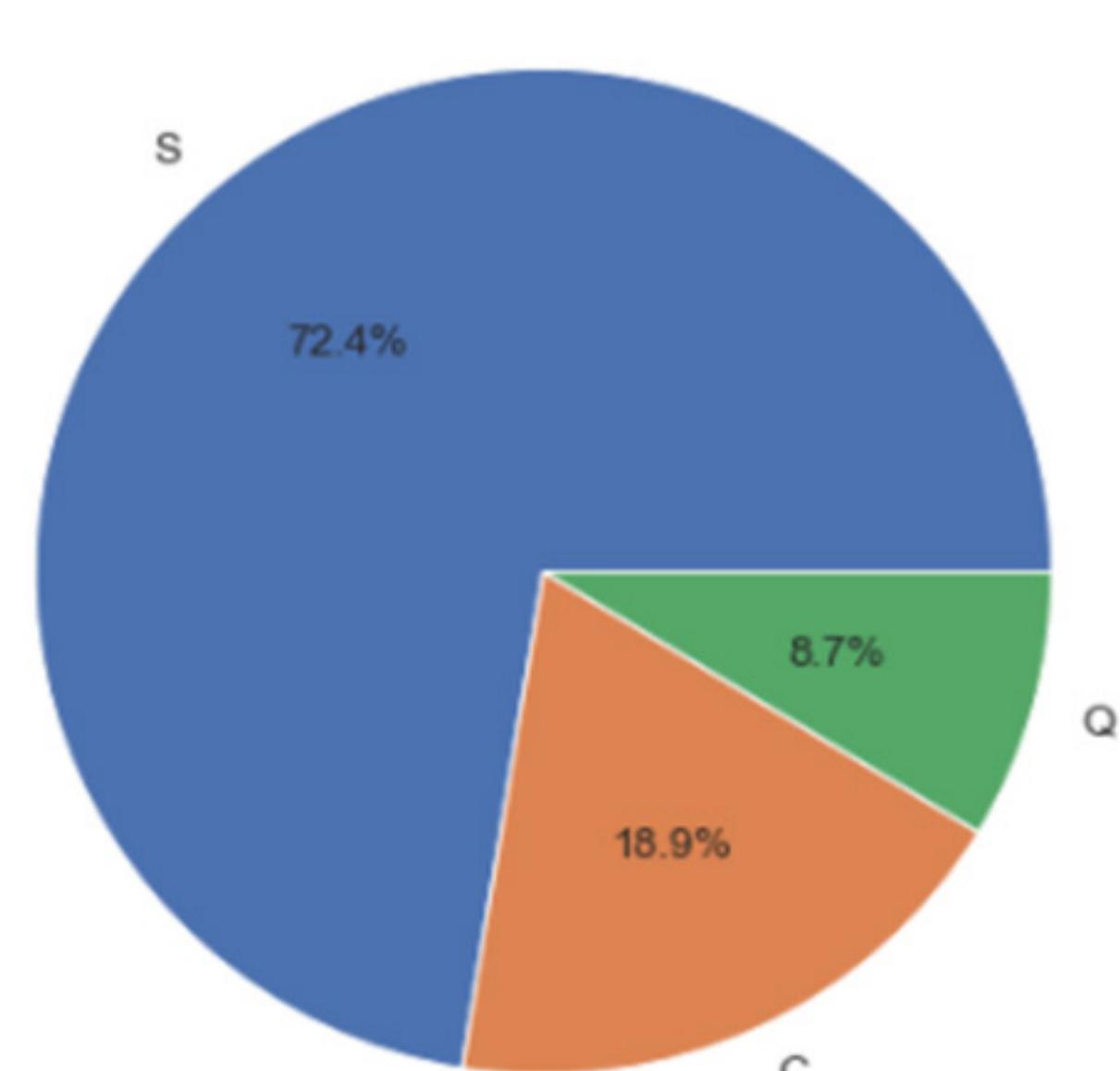
```
In [16]: 1 sex_count = df['Sex'].value_counts()
2 plt.figure(figsize=(7,7))
3 plt.title('Grouped by Sex')
4 plt.pie(sex_count.values,
5         labels = ['{}'.format(i) for i in sex_count.index],
6         autopct='%.1f%%', textprops={'fontsize':13})
7 plt.show()
```



There are more male passengers than female passengers. Probably most of the female passengers are wives or daughter

6. Passangers distribution by embarkation point

```
1 embarked_count = df['Embarked'].value_counts()
2 plt.figure(figsize=(7,7))
3 plt.title('Grouped by embarkation')
4 plt.pie(embarked_count.values,
5         labels = ['{}'.format(i) for i in embarked_count.index],
6         autopct='%.1f%%', textprops={'fontsize':13})
7 plt.show()
```

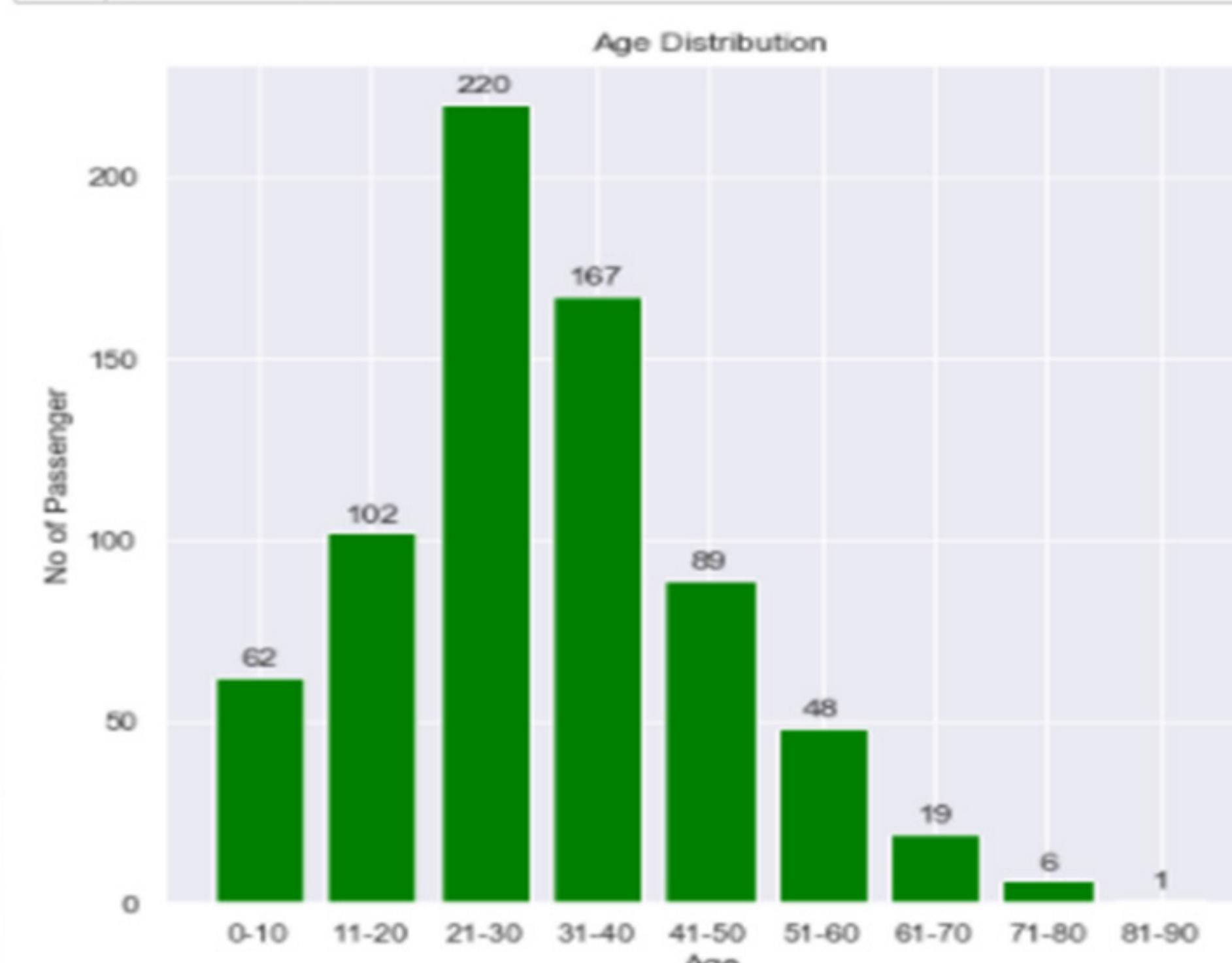


Most of the passengers embarked from southampton.

7. Passengers distribution by Age

Since there exist a null value, we will only plot the non-null value data. Histogram will be used to show the age distribution.

```
1 ages = df[df['Age'].notnull()]['Age'].values
2 ages_hist = np.histogram(ages, bins=[0,10,20,30,40,50,60,70,80,90])
3 ages_hist_label = ['0-10','11-20','21-30',
4                     '31-40','41-50','51-60',
5                     '61-70','71-80','81-90']
6
7
8 plt.figure(figsize=(7,7))
9 plt.title('Age Distribution')
10 plt.bar(ages_hist_label, ages_hist[0], color='Green')
11 plt.xlabel('Age')
12 plt.ylabel('No of Passenger')
13 for i, bin in zip(ages_hist[0], range(9)) :
14     plt.text(bin, 1+5, str(int(i)), fontsize=12,
15               horizontalalignment='center', verticalalignment='center')
16 plt.show()
```



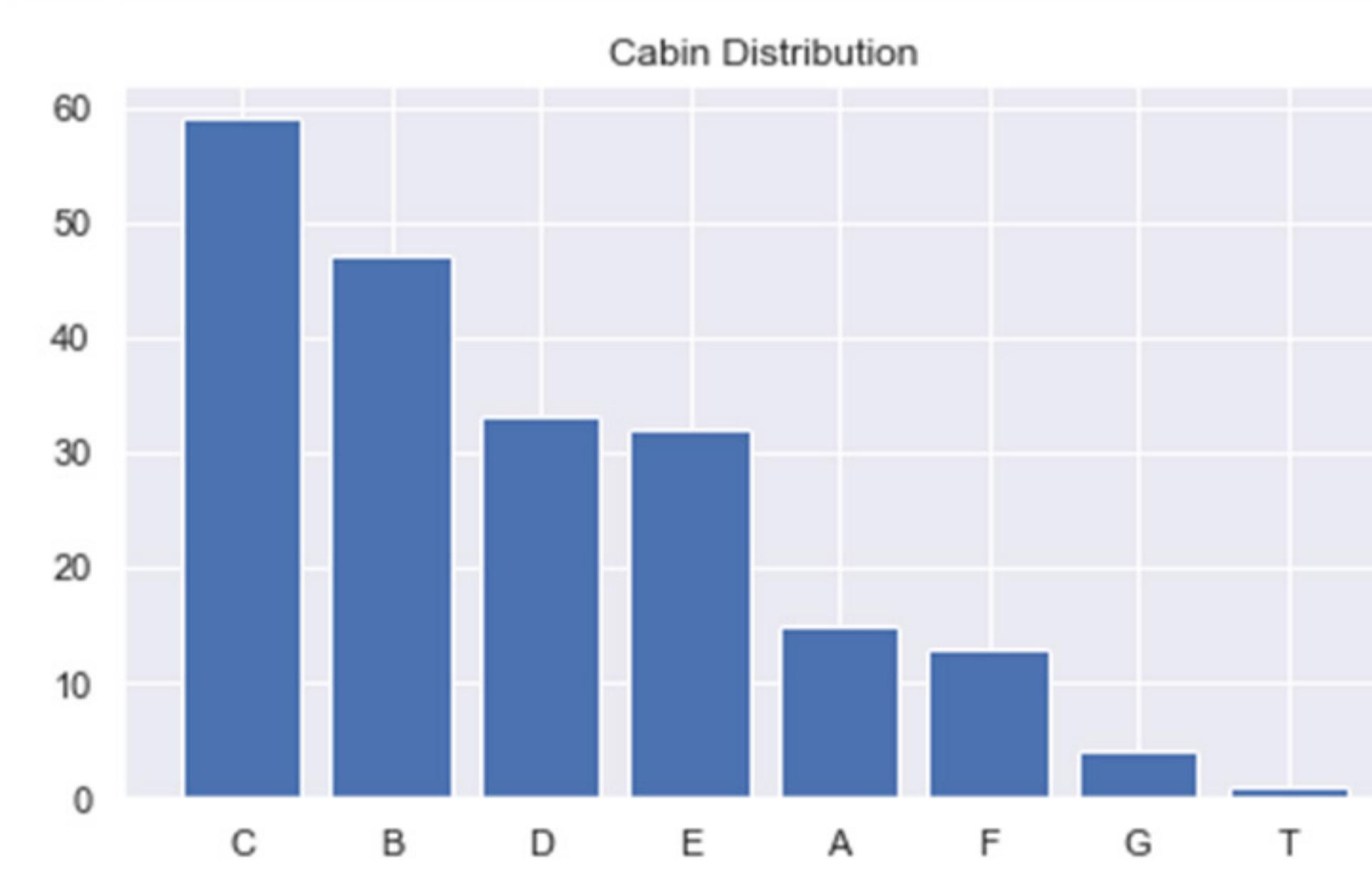
We can infer that most of the passengers are from ages 20s to 40s.

8. Cabins distribution

Just like the age distribution, since there exist a null value in cabin's feature, we need to use only the non-null value. Also, we're going to distribute the cabin's feature only by its initial name.

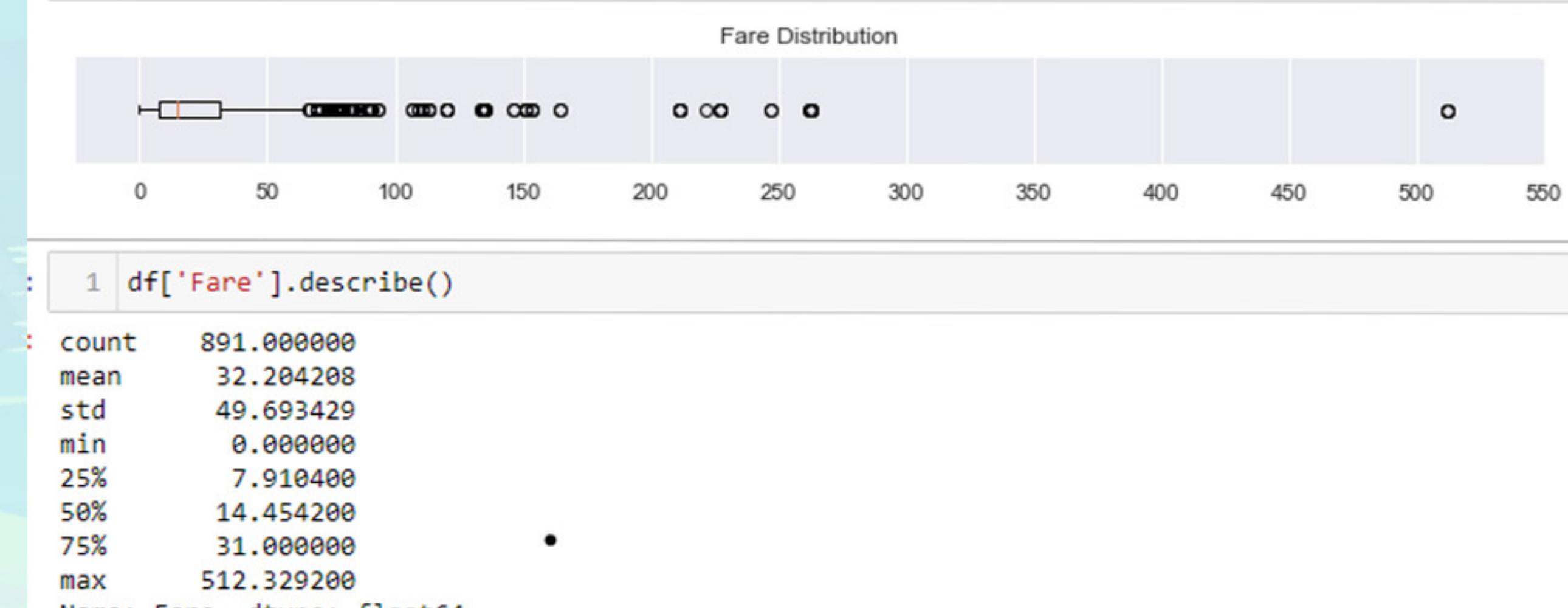
```
1 cabins = df['Cabin'].dropna()
2 def take_initial(x):
3     return x[0]
4 cabins = cabins.apply(take_initial)
```

```
1 cabins_count = cabins.value_counts()
2 plt.figure(figsize=(7,4))
3 plt.title("Cabin Distribution")
4 plt.bar(cabins_count.index, cabins_count.values)
5 plt.show()
```



9. Tickets Price distribution

```
1 plt.figure(figsize=(13,1))
2 plt.title("Fare Distribution")
3 plt.boxplot(df['Fare'], vert=False)
4 plt.xticks(range(0,600,50))
5 plt.yticks([])
6 plt.show()
```



Now that we've finished doing the ETL, we're going to do some feature engineering to make our data is ready to use for modelling using the Logistic Regression Algorithm.

Feature Engineering

1. SibSp and Parch Feature

We're going to combine the SibSp + Parch + 1(the passenger itself) as new column "family size" to measure the family size of each passengers. This feature would be usefull to be included in later modelling.

```

1 df["SibSp"].unique()
array([1, 0, 3, 4, 2, 5, 8], dtype=int64)

1 df["FamilySize"] = df['SibSp'] + df['Parch'] + 1
2 df[['Name', 'SibSp', 'Parch', 'FamilySize']]

      Name  SibSp  Parch  FamilySize
0    Braund, Mr. Owen Harris       1      0          2
1  Cumings, Mrs. John Bradley (Florence Briggs Th...       1      0          2
2   Heikkinen, Miss. Laina       0      0          1
3  Futrelle, Mrs. Jacques Heath (Lily May Peel)       1      0          2
4    Allen, Mr. William Henry       0      0          1
...
886   Montvila, Rev. Juozas       0      0          1
887  Graham, Miss. Margaret Edith       0      0          1
888  Johnston, Miss. Catherine Helen "Carrie"       1      2          4
889    Behr, Mr. Karl Howell       0      0          1
890   Dooley, Mr. Patrick       0      0          1
891 rows × 4 columns

```

2. Embarked feature

Since there only 2 missing data in this feature, we may drop the rows that include this missing data. The type of data of this features is categorical data, therefore we need to change it to a numerical value using one-hot encoding method.

```

1 print("Sebelum drop :", df.shape)
2 df = df.dropna(subset=['Embarked'])
3 print("Setelah drop :", df.shape)

Sebelum drop : (891, 13)
Setelah drop : (889, 13)

1 embarked_one_hot = pd.get_dummies(df['Embarked'], prefix='Embarked')
2 df = pd.concat([df, embarked_one_hot], axis=1)

```

3. Cabins feature

Since there're too many missing value in this feature, we can't just simply drop the rows because we're going to lose too many data. And it's also not good to fill the missing value with the most used cabins, so we will just fill the missing value as unknown 'U'. Because the data type is categorical data, we have to do a one-hot encoding to change the data type into numerical value.

```

1 df['Cabin'] = df['Cabin'].fillna('U')

1 def take_initial(x):
2     return x[0]
3 df['Cabin'] = df['Cabin'].apply(take_initial)

1 cabin_one_hot = pd.get_dummies(df['Cabin'], prefix='Cabin')
2 df = pd.concat([df, cabin_one_hot], axis=1)
3
4 #Hapus kolom cabin yang aslinya
5 df = df.drop('Cabin', axis=1)
6
7 #Tampilkan nama tiap kolom
8 df.columns

Index(['PassengerId', 'Survived', 'Pclass', 'Name', 'Sex', 'Age', 'SibSp',
       'Parch', 'Ticket', 'Fare', 'Embarked', 'Cabin_A', 'Cabin_B', 'Cabin_C',
       'Cabin_D', 'Cabin_E', 'Cabin_F', 'Cabin_G', 'Cabin_T', 'Cabin_U'],
      dtype='object')

```

4. Name feature

We can make a categorical feature out of Name feature. It's that we can categorize the passengers name by its title, and then we make it as a numerical data types by one-hot encoding.

```

1 def get_title(x):
2     return x.split(',')[1].split('.')[0].strip()
3
4 df['Title'] = df['Name'].apply(get_title)

1 df['Title'].unique()

array(['Mr', 'Mrs', 'Miss', 'Master', 'Don', 'Rev', 'Dr', 'Mme', 'Ms',
       'Major', 'Lady', 'Sir', 'Mlle', 'Col', 'Capt', 'the Countess',
       'Jonkheer'], dtype=object)

1 title_one_hot = pd.get_dummies(df['Title'], prefix='Title')
2 df = pd.concat([df, title_one_hot], axis=1)
3
4 # Tampilkan semua kolom di dataframe
5 df.columns

Index(['PassengerId', 'Survived', 'Pclass', 'Name', 'Sex', 'Age', 'SibSp',
       'Parch', 'Ticket', 'Fare', 'FamilySize', 'Embarked_C', 'Embarked_Q',
       'Embarked_S', 'Cabin_A', 'Cabin_B', 'Cabin_C', 'Cabin_D', 'Cabin_E',
       'Cabin_F', 'Cabin_G', 'Cabin_T', 'Cabin_U', 'Title', 'Title_Capt',
       'Title_Col', 'Title_Don', 'Title_Dr', 'Title_Jonkheer', 'Title_Lady',
       'Title_Major', 'Title_Master', 'Title_Miss', 'Title_Mlle', 'Title_Mme',
       'Title_Mr', 'Title_Mrs', 'Title_Ms', 'Title_Rev', 'Title_Sir',
       'Title_the Countess'],
      dtype='object')

```

5. Sex feature

In this feature, we obviously just need to simply change the data types from categorical to numerical value through one-hot-encoding

```
1 sex_one_hot = pd.get_dummies(df['Sex'], prefix='Sex')
2 df = pd.concat([df, sex_one_hot], axis=1)
3
4 # Hapus kolom sex yang aslinya
5 df = df.drop('Sex', axis=1)
6
7 df.columns

Index(['PassengerId', 'Survived', 'Pclass', 'Name', 'Age', 'SibSp', 'Parch',
       'Ticket', 'Fare', 'FamilySize', 'Embarked_C', 'Embarked_Q',
       'Embarked_S', 'Cabin_A', 'Cabin_B', 'Cabin_C', 'Cabin_D', 'Cabin_E',
       'Cabin_F', 'Cabin_G', 'Cabin_T', 'Cabin_U', 'Title', 'Title_Capt',
       'Title_Col', 'Title_Don', 'Title_Dr', 'Title_Jonkheer', 'Title_Lady',
       'Title_Major', 'Title_Master', 'Title_Miss', 'Title_Mlle', 'Title_Mme',
       'Title_Mr', 'Title_Mrs', 'Title_Ms', 'Title_Rev', 'Title_Sir',
       'Title_the Countess', 'Sex_female', 'Sex_male'],
      dtype='object')
```

6. Age feature

We've noticed that there are lots of missing value in this feature. We're going to fill the missing value by the mean of ages of each titles just like in the following code :

```
1 age_median = df.groupby('Title')['Age'].median()
2 age_median

Title
Capt      70.0
Col       58.0
Don       40.0
Dr        46.5
Jonkheer   38.0
Lady      48.0
Major     48.5
Master    3.5
Miss      21.0
Mlle      24.0
Mme      24.0
Mr       30.0
Mrs      35.0
Ms       28.0
Rev       46.5
Sir       49.0
the Countess 33.0
Name: Age, dtype: float64

1 def fill_age(x) :
2     for index, age in zip(age_median.index, age_median.values):
3         if x['Title'] == index:
4             return age
5 df['Age'] = df.apply(lambda x : fill_age(x) if np.isnan(x['Age']) else x['Age'], axis=1)
```

7. We're going to drop feature that we're not going to use in the modeling.

```
1 df = df.drop(['PassengerId', 'Ticket', 'Title', 'Name', 'Cabin', 'embarked'], axis=1)
```

Now that the data is ready, we're going to start to make the model using Logistic Regression to classify/predict whether any new passenger will survive. or not.

1. We split the data into train data and test data. We use the train data for training the model, and the test data to evaluate the accuracy of the model

```
1 from sklearn.model_selection import train_test_split
2 X_train, X_test, y_train, y_test = train_test_split(X, y, random_state=21, test_size=0.2)
3
4 print('X_train.shape\t:', X_train.shape)
5 print('y_train.shape\t:', y_train.shape)
6 print('X_test.shape\t:', X_test.shape)
7 print('y_test.shape\t:', y_test.shape)

X_train.shape : (711, 37)
y_train.shape : (711,)
X_test.shape  : (178, 37)
y_test.shape  : (178,)
```

We use 80% of the data to train the model, and the rest for testing the model.

Train the Model using Logistic Regression Algorithm

1. We split the data into train data and test data. We use the train data for training the model, and the test data to evaluate the accuracy of the model

```
1 from sklearn.model_selection import train_test_split
2 X_train, X_test, y_train, y_test = train_test_split(X, y, random_state=21, test_size=0.2)
3
4 print('X_train.shape\t:', X_train.shape)
5 print('y_train.shape\t:', y_train.shape)
6 print('X_test.shape\t:', X_test.shape)
7 print('y_test.shape\t:', y_test.shape)

X_train.shape : (711, 37)
y_train.shape : (711,)
X_test.shape  : (178, 37)
y_test.shape  : (178,)
```

We use 80% of the data to train the model, and the rest for testing the model.

2. Train the model using *scikit learn* Library

```
1 from sklearn.linear_model import LogisticRegression
2 clf = LogisticRegression()
3 clf.fit(X_train, y_train)
```

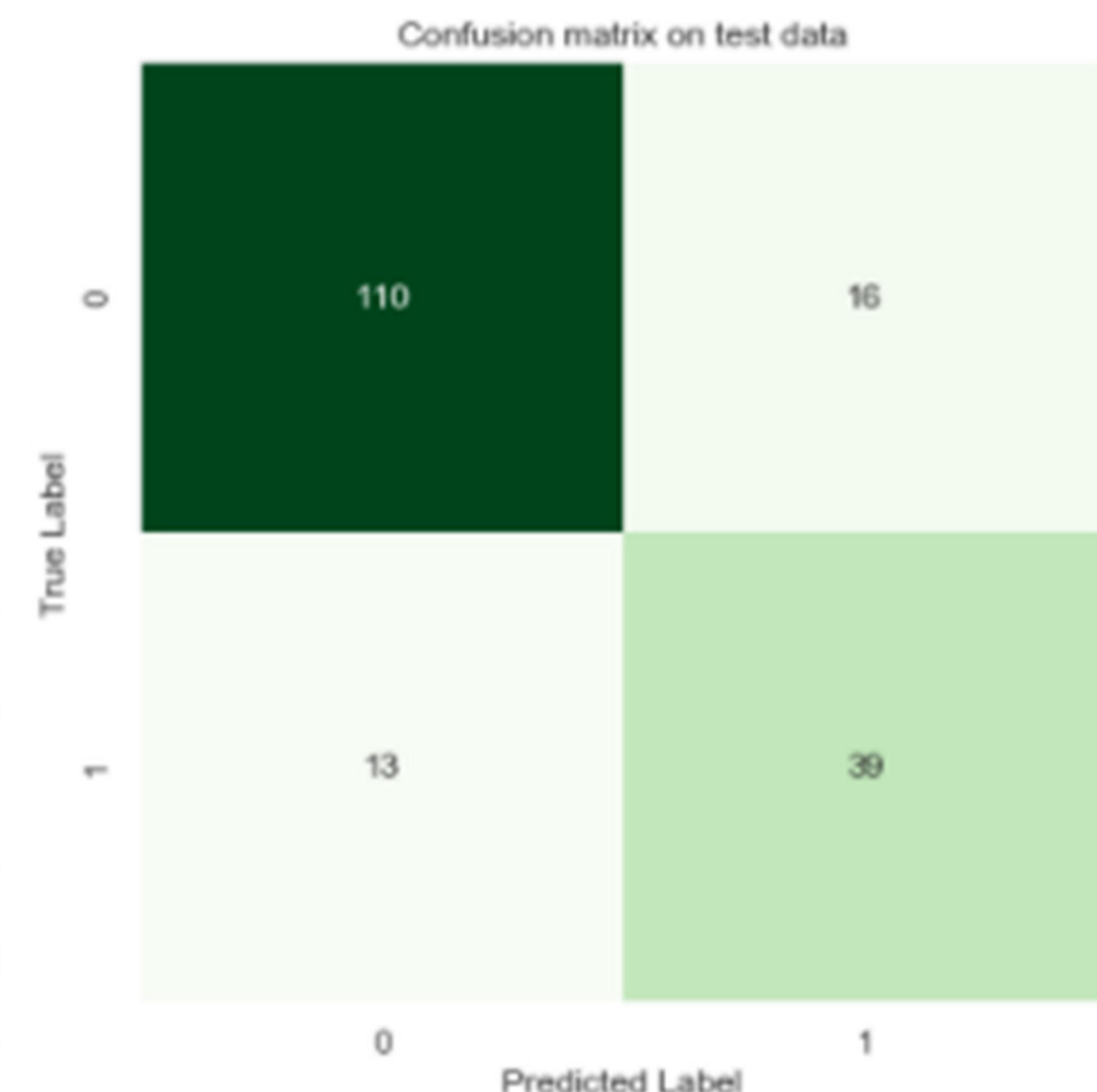
3. Evaluate the model using confusion matrix

```
1 test_preds = clf.predict(X_test)

1 cm = confusion_matrix(y_test, test_preds)
2 cm

array([[110,  16],
       [ 13,  39]], dtype=int64)

1 plt.figure(figsize=(6,6))
2 plt.title('Confusion matrix on test data')
3 sns.heatmap(cm, annot=True, fmt='d', cmap=plt.cm.Greens, cbar=False)
4 plt.xlabel('Predicted Label')
5 plt.ylabel('True Label')
6 plt.show()
```



As we mentioned before, we use the test set to evaluate our model. First, we use the model to predict the passengers in the test data whether each of them will survive or not. And then we're comparing the predicted result (test_preds) with the true result (y_test) using the confusion matrix.

As we can see, from the confusion matrix we have a 110 true positive, 13 false positive (the predicted says it doesn't survive, but the true result is survive), 16 False positive, and 39 false negative.

Since it has the accuracy of the test data is 83 %, i'd say this model is quite good to use.

```
1 print("Train accuracy\t:", clf.score(X_train, y_train))
2 print("Test accuracy\t:", clf.score(X_test, y_test))

Train accuracy : 0.8368495077355836
Test accuracy  : 0.8370786516853933
```

Customer Segmentation Using K-Means Algorithm

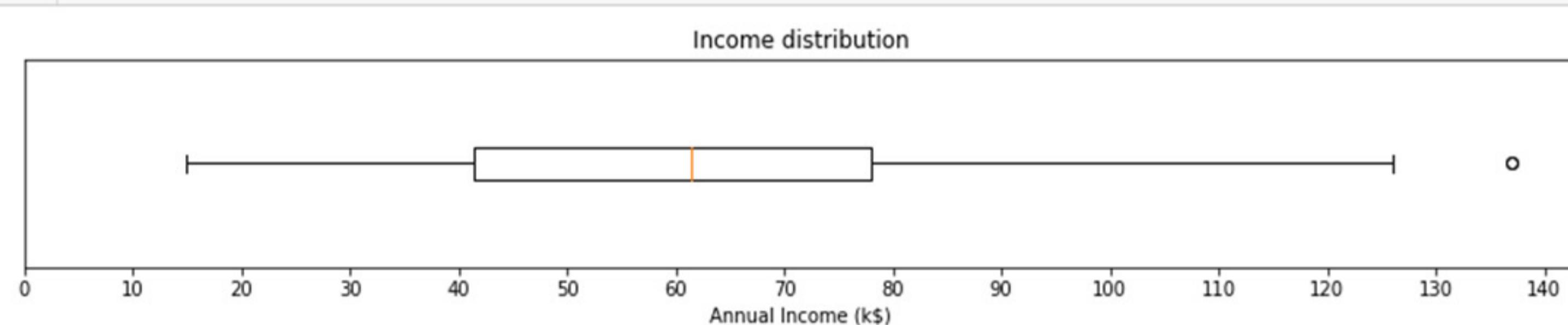
In this project, we will try to clusterize customers at Mall. The purpose of this clusterization is to help the marketing team for targeting the mall's future customers.

First let's take a look at our data

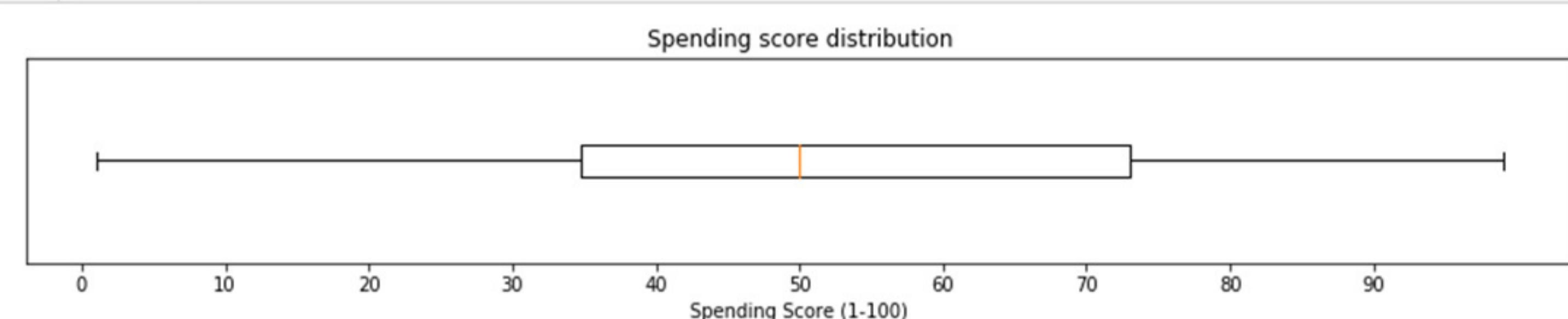
```
1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 from sklearn.cluster import KMeans
6
7 df = pd.read_csv("Mall_Customers.csv")
8 df.head(10)
```

	CustomerID	Gender	Age	Annual Income (k\$)	Spending Score (1-100)
0	1	Male	19	15	39
1	2	Male	21	15	81
2	3	Female	20	16	6
3	4	Female	23	16	77
4	5	Female	31	17	40
5	6	Female	22	17	76
6	7	Female	35	18	6
7	8	Female	23	18	94
8	9	Male	64	19	3
9	10	Female	30	19	72

```
1 plt.figure(figsize=(15,2))
2 plt.boxplot(df['Annual Income (k$)'], vert = False)
3 plt.title('Income distribution')
4 plt.xticks([])
5 plt.yticks(range(0,150,10))
6 plt.xlabel('Annual Income (k$)')
7 plt.show()
```



```
1 plt.figure(figsize=(15,2))
2 plt.boxplot(df['Spending Score (1-100)'], vert = False)
3 plt.title('Spending score distribution')
4 plt.xticks([])
5 plt.yticks(range(0,100,10))
6 plt.xlabel('Spending Score (1-100)')
7 plt.show()
```



We're going to use only these two feature to clusterize our customers, since the other features wouldn't really help at all with our analysis.

Here in the tabel/data we have information about CustomerID, gender, age, Annual income, and spending score of each customer who's spend some of their money at the mall.

To understand more of the data, we're going to do some Exploratory Data Analysis.

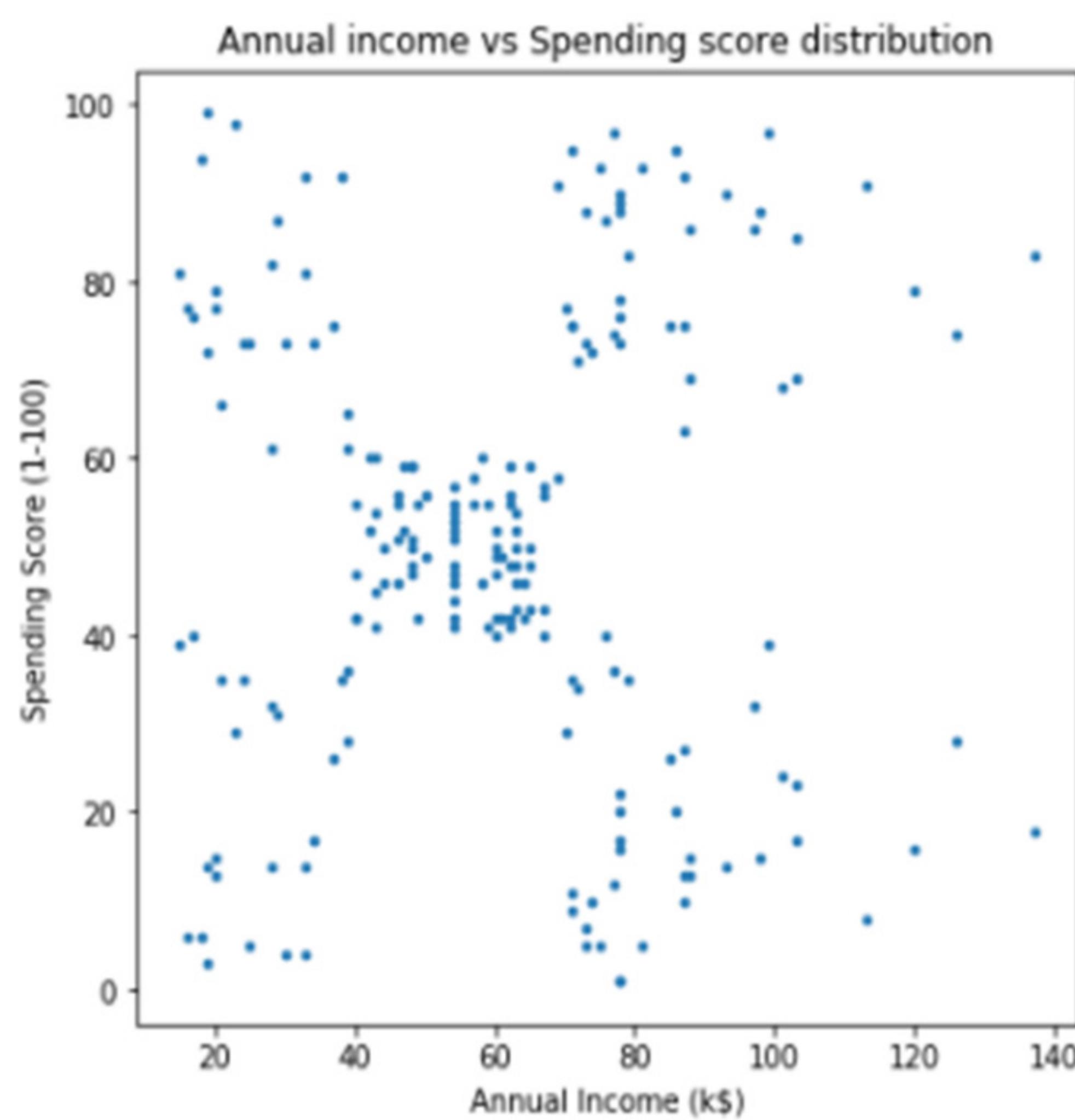
Let's create a boxplot to see the distribution of the Annual Income (k\$) and Spending Score (1-100) feature using matplotlib library.

From Annual income distribution boxplot we can tell that most of the costumers has an annual salary from approximately 40K to 80 K \$. In the boxplot we can see there's an outlier data above the max value, which means there's a costumer that has much more annual salary than the other customer.

From the spending score distribution boxplot we can tell that most of spending score lies around 35K to 73 K \$. And there's no outlier data at the spending score distribution. Which means the distribution is (strictly) normally distributed.

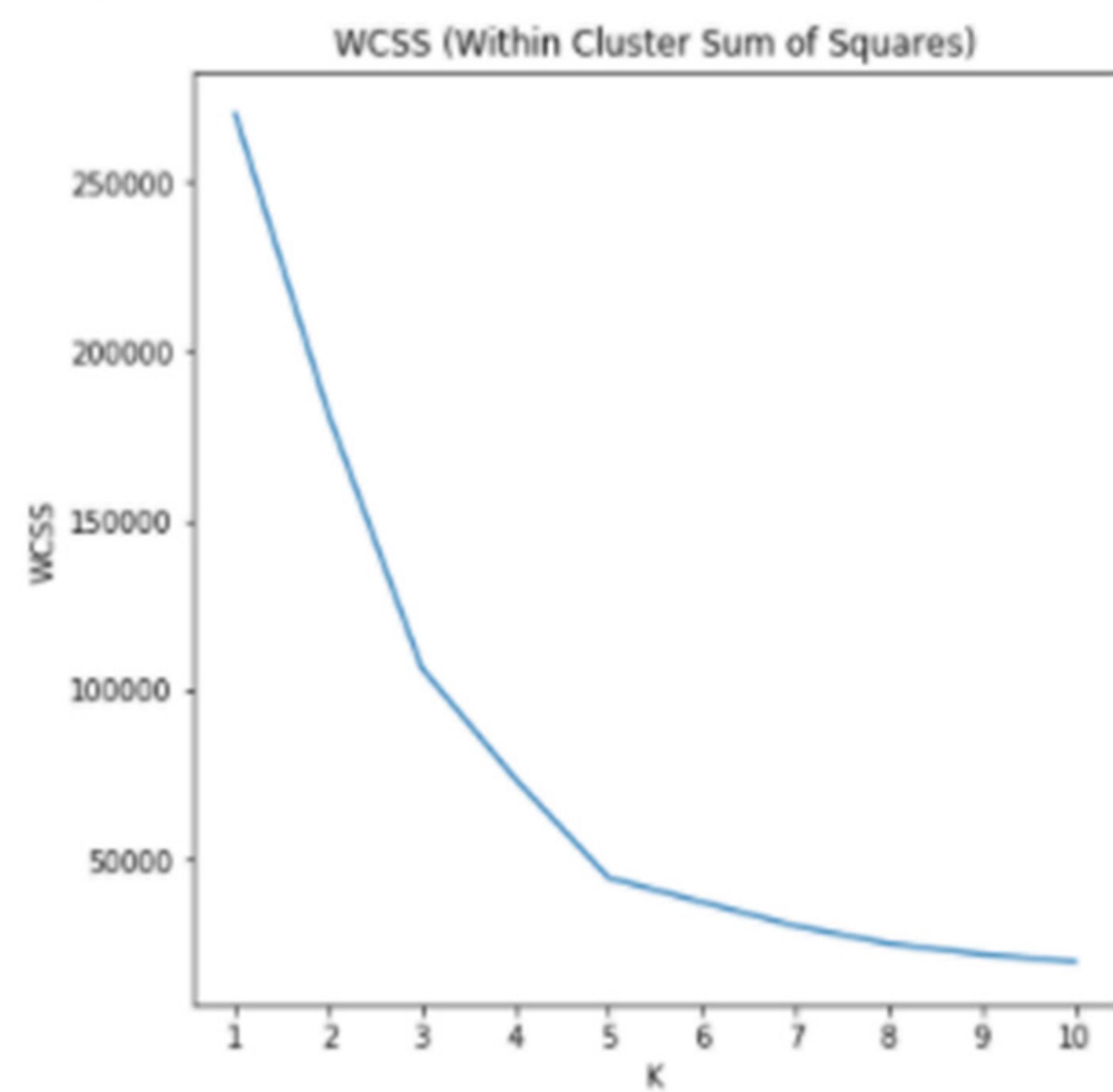
K-Means Clustering algorithm

```
1 X = df[['Annual Income (k$)', 'Spending Score (1-100)']].values  
  
1 plt.figure(figsize=(6,6))  
2 plt.scatter(X[:,0], X[:,1], s=7)  
3 plt.title('Annual income vs Spending score distribution')  
4 plt.xlabel('Annual Income (k$)')  
5 plt.ylabel('Spending Score (1-100)')  
6 plt.show()
```



```
1 # Ngitung WCSS
2 WCSS = []
3 for i in range(1,11):
4     kmeans = KMeans(n_clusters=i, random_state=42)
5     kmeans.fit(X)
6     WCSS.append(kmeans.inertia_)

1 # Menampilkan graph WCSS untuk melakukan elbow method
2 plt.figure(figsize=(6,6))
3 plt.title('WCSS (Within Cluster Sum of Squares)')
4 plt.plot(range(1,11), WCSS)
5 plt.xticks(range(1,11))
6 plt.xlabel('K')
7 plt.ylabel('WCSS')
8 plt.show()
```



```
1 # Inisialisasi object KMeans
2 kmeans = KMeans(n_clusters=5, random_state = 41)
3
4 # Memulai proses clustering
5 kmeans.fit(X)

KMeans(algorithm='auto', copy_x=True, init='k-means++', max_iter=300,
       n_clusters=5, n_init=10, n_jobs=None, precompute_distances='auto',
       random_state=41, tol=0.0001, verbose=0)

1 y_kmeans = kmeans.predict(X)
2 print(y_kmeans)

[2 4 9 4 2 4 9 4 9 4 9 4 9 4 2 4 2 4 2 4 9 4 9 4 2 4 2 4 9 4 9 4 9 4 9
 4 2 4 2 4 2 0 2 0 0 0 2 2 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
 0 0 0 0 0 7 0 0 7 7 0 0 0 0 0 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7
 7 7 7 7 7 7 7 7 7 7 7 7 3 7 5 7 3 1 5 1 5 7 5 1 3 1 5 1 5 1 3 1 5 1 3 7 3 1 3 7 5
 1 3 1 3 1 5 1 3 1 5 1 5 7 3 1 3 1 5 1 3 1 5 1 5 1 3 1 5 1 3 1 5 1 3 8 3 8 3 8
 3 8 6 8 6 8 6 8 6 8 6 8 6 8 6 8 6]
```



```
1 # koordinat center (titik pusat) dari tiap cluster
2 centroids = kmeans.cluster_centers_
3 print(centroids)

[[86.53846154 82.12820513]
 [55.2962963 49.51851852]
 [25.72727273 79.36363636]
 [88.2 17.11428571]
 [26.30434783 20.91304348]]
```

First thing we do to use this k-means algorithm is to get our features values just as we can see at our code.

And then we might want to check how our data scattered (using the matplotlib library) to roughly predict how many cluster do we might have from our data.

We can tell that (intuitively) there are about 5 cluster.

Just to make sure, we're going to use the WCSS (within cluster sum of squares) method or the elbow method to determine the total of our clusters.

Using the elbow method, we can see that the plot starting to fashionly linear (from the graph) when the value of k is 5. which means our previous prediction was true that the total of the cluster for our analysis is 5 cluster.

Now that we have the total of our cluster, we can start to train our data to clusterize our customers.

```

1 # Inisialisasi object KMeans
2 kmeans = KMeans(n_clusters=5, random_state = 41)
3
4 # Memulai proses clustering
5 kmeans.fit(X)

KMeans(algorithm='auto', copy_x=True, init='k-means++', max_iter=300,
       n_clusters=5, n_init=10, n_jobs=None, precompute_distances='auto',
       random_state=41, tol=0.0001, verbose=0)

```

```

1 y_kmeans = kmeans.predict(X)
2 print(y_kmeans)

[2 4 9 4 2 4 9 4 9 4 9 4 9 4 9 4 2 4 2 4 2 4 9 4 9 4 2 4 2 4 9 4 9 4 9 4 9
4 2 4 2 4 2 0 2 0 0 0 2 2 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 7 0 0 7 7 0 0 0 0 0 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7
7 7 7 7 7 7 7 7 7 7 7 3 7 5 7 3 1 5 1 5 7 5 1 3 1 5 1 5 1 3 7 3 1 3 7 5
1 3 1 3 1 5 1 3 1 5 1 5 7 3 1 3 1 5 1 3 1 5 1 3 1 3 1 5 1 3 8 3 8 3 8
3 8 6 8 6 8 6 8 6 8 6 8 6 8 6]
```

```

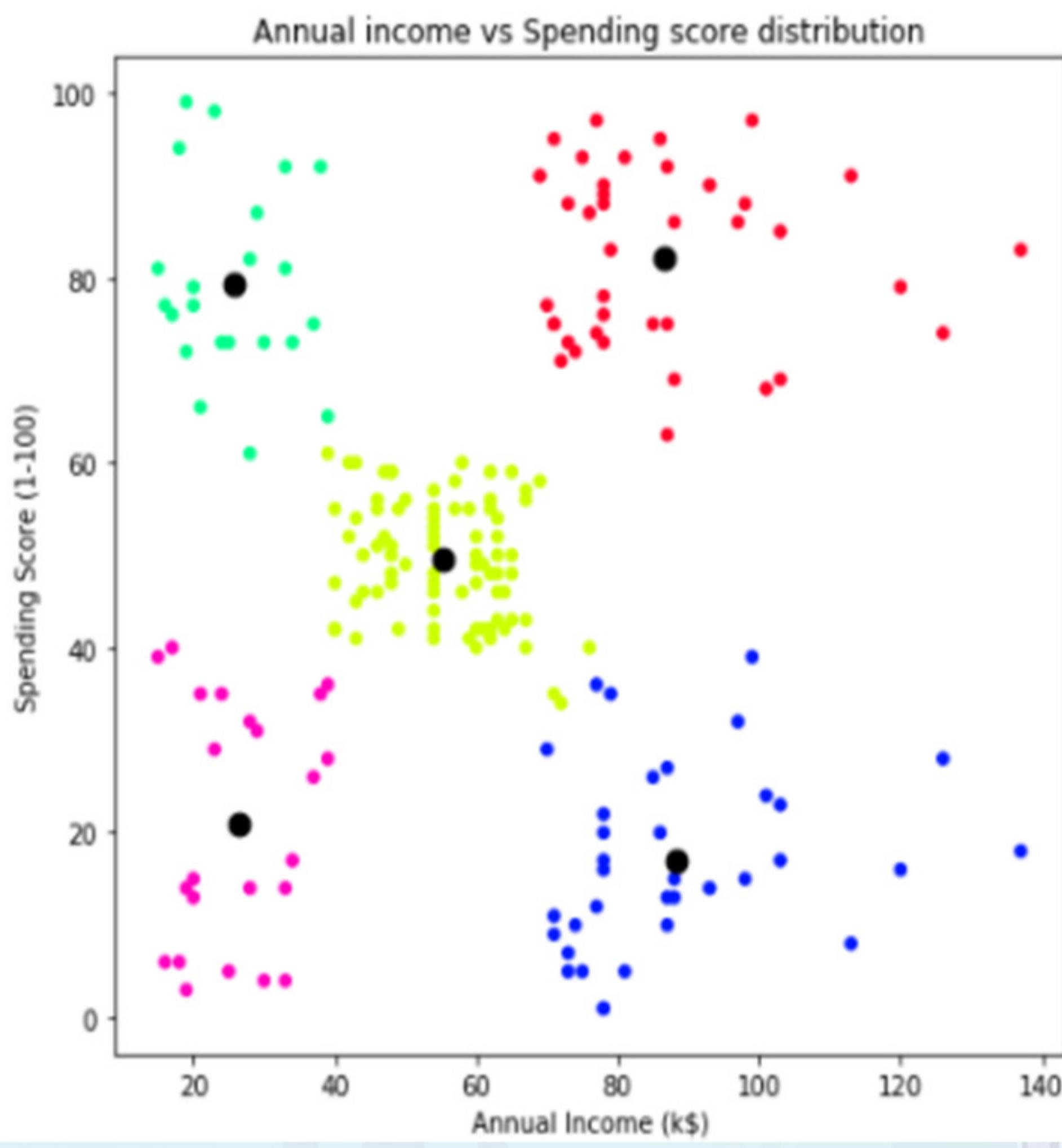
1 # koordinat center (titik pusat) dari tiap cluster
2 centroids = kmeans.cluster_centers_
3 print(centroids)

[[86.53846154 82.12820513]
 [55.2962963 49.51851852]
 [25.72727273 79.36363636]
 [88.2 17.11428571]
 [26.30434783 20.91304348]]
```

```

1 plt.figure(figsize=(7,7))
2 plt.scatter(X[:,0], X[:,1], s=20, c=y_kmeans,cmap='gist_rainbow')
3 plt.scatter(centroids[:,0],centroids[:,1], s=75, c='black')
4 plt.title('Annual income vs Spending score distribution')
5 plt.xlabel('Annual Income (k$)')
6 plt.ylabel('Spending Score (1-100)')
7 plt.show()

```



To make our model, first we have to input number of cluster into the *KMeans* parameter. For the random state I specify 42, it actually doesn't really matter what value we want to specify to the *random_state* parameter. And then, we train our model with our data to starting clusterize our data and getting the *centroids of each cluster*. Now that we have our centroids and clusterized each of the customer, we can visualize it in a graph to see how our data is clusterize. But what actually can we tell from this clusterization ?

1. The **red** cluster represents the customers who both have high annual income and spending score. The mall need keep advertize to these customers so that tey will keep shopping at the mall. We can advertize through email or maybe through text massages to reach these customers.
2. The **blue** cluster represents the customers who has high annual income but doesn't spend so much. The mall need to advertize more to these customer, or maybe have a questioner to ask their opinion to improve the mall so that these customers would shop more ofter.
3. The **green** clusters represents customers who have medium annual income and also medium spending score. I believe there's no need any special treatment for these customer.
4. The **pink** clusters represents customers who have low annual income and also low spending score. I honestly don't think we need to put more ad to these customer since they don't spend so much and also probably these customers try to save their money as much as they can.
5. The last cluster (i'm not sure what color they are) is kind of weird. It's represents customers who have low income but they spend much money at the mall. It could be they already have a big saving (from previous jobs or from parents) or could be they always in debt to someone

Our table with added cluster columns

```
: 1 df.head(20)
```

	CustomerID	Gender	Age	Annual Income (k\$)	Spending Score (1-100)	Cluster
0	1	Male	19	15	39	low income, low spend
1	2	Male	21	15	81	Low Income, spend much
2	3	Female	20	16	6	low income, low spend
3	4	Female	23	16	77	Low Income, spend much
4	5	Female	31	17	40	low income, low spend
5	6	Female	22	17	76	Low Income, spend much
6	7	Female	35	18	6	low income, low spend
7	8	Female	23	18	94	Low Income, spend much
8	9	Male	64	19	3	low income, low spend
9	10	Female	30	19	72	Low Income, spend much
10	11	Male	67	19	14	low income, low spend
11	12	Female	35	19	99	Low Income, spend much
12	13	Female	58	20	15	low income, low spend
13	14	Female	24	20	77	Low Income, spend much
14	15	Male	37	20	13	low income, low spend
15	16	Male	22	20	79	Low Income, spend much
16	17	Female	35	21	35	low income, low spend
17	18	Male	20	21	66	Low Income, spend much
18	19	Male	52	23	29	low income, low spend
19	20	Female	35	23	98	Low Income, spend much

On the table above, i've added new column which is cluster to represents each of the customer's cluster. I've write the *Red* cluster as 'High income, spend much', the *blue* cluster as 'High Income, low spend', the *green* cluster as 'medium income, medium spend', the *pink* cluster as 'low income, low spend', and the last cluster as 'Low Income, spend much'.

This way the marketing team can read the data easier to target their marketing area.

*Background pictures are taken from freepik.com

Thank you for your time!