

You are owing a supermarket mall and through membership cards, you have some basic data about your customers and their spending score. Spending Score is something you assign to the customer based on your defined parameters.

```
import numpy as np #linear algebra library of Python
import pandas as pd # build on top of numpy for data analysis, data manipulation and data I/O
import matplotlib.pyplot as plt #plotting library of Python
```

Now let's mount Google drive so that we can upload the diabetes.csv file. You can find the code in the 'Code' tab.

```
from google.colab import drive
drive.mount('/content/gdrive')
```



First thing that we do is take a look at the shape of the dataframe (df.shape) and take a look at first 5 lines of the dataframe (df.head()).

```
df=pd.read_csv('/content/gdrive/My Drive/Colab Notebooks/Mall_Customers.csv') #import data
df.head() #shows first 5 lines including column namesdf.shape # number of rows and columns
```



```
df.shape # provides # rows and # columns of the dataframe df - 200 rows and 5 columns
```

```
!rm -rf /content/gdrive # Drive already mounted at /content/gdrive; to attempt to forcibly remount, call drive.mount()
```

```
X=df.iloc[:, [3,4]].values # selects columns 3 and 4 of the dataframe (2 last columns)
```

Now we are ready for K-means and we will start by finding the best k. This is the so-called elbow method. Basically, as k increases the cost will decrease. The for loop will fit the kmeans algorithm to our data and will compute the cost, which is then appended to our wcss list.

Parameters kmeans:

init: 'random': random initialization method - choose k observations (rows) at random from data for the initial cluster centers

max\_iter: maximum number of iterations of the k-means algorithm for a single run to converge

n\_init: number of time the k-means algorithm will be run with different centroid seeds. The final results will be the lowest inertia.

```
from sklearn.cluster import KMeans
wcss=[] # wcss is 'within cluster sum of squares'
for i in range(1,11):
    kmeans=KMeans(n_clusters=i, init='random', max_iter=300, n_init=10, random_state=42)
    kmeans.fit(X) # fits the modelcompute k-means clustering
    wcss.append(kmeans.inertia_) # sum of squared distances of samples to their closes

plt.plot(range(1,11),wcss)
plt.title('The Elbow Method Graph')
plt.xlabel('Number of clusters')
plt.ylabel('WCSS')
plt.show()
```



Based on above graph (k vs. Cost) we suggest to take k=5 (5 centroids meaning 5 clusters)

```
model=KMeans(n_clusters=5,init='random', max_iter=300, n_init=10, random_state=42)
y_kmeans=model.fit_predict(X) # output is a vector with all samples and their cluster
```

```
plt.figure(figsize=(8,5))
plt.scatter(X[y_kmeans == 0, 0], X[y_kmeans == 0, 1], s = 100, c = 'red', label = 'Low Income & Not Spender')
plt.scatter(X[y_kmeans == 1, 0], X[y_kmeans == 1, 1], s = 100, c = 'blue', label = 'High Income & High Spender')
plt.scatter(X[y_kmeans == 2, 0], X[y_kmeans == 2, 1], s = 100, c = 'green', label = 'Low Income & High Spender')
plt.scatter(X[y_kmeans == 3, 0], X[y_kmeans == 3, 1], s = 100, c = 'cyan', label = 'High Income & Low Spender')
plt.scatter(X[y_kmeans == 4, 0], X[y_kmeans == 4, 1], s = 100, c = 'magenta', label = 'Low Income & Low Spender')

plt.title('Clusters of customers')
plt.xlabel('Annual Income (k$)')
plt.ylabel('Spending Score (1-100)')
plt.legend(loc=0)
```

```
plt.show()
```

`X[y_kmeans == 0, 0]` means outcome is cluster index 0 for column 0