## ⌄ Ex 7

Use the given dataset and implement K-Means from scratch and use the sklearn K-Means implementation. Compare the results of both the implementations and write your inferences in the ipynb file itself.

```python
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
from sklearn.cluster import KMeans
from sklearn.metrics import accuracy_score
from itertools import permutations

import warnings
warnings.filterwarnings('ignore')
```

```python
df = pd.read_csv("data.csv")
```
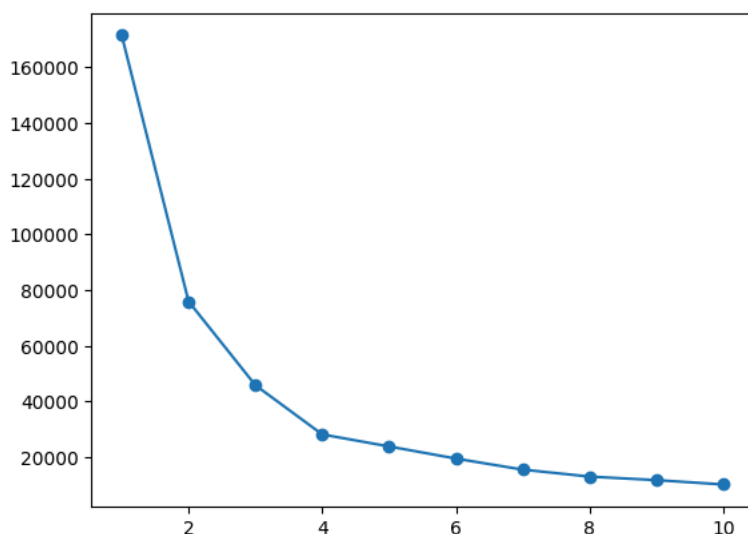
```python
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 200 entries, 0 to 199
Data columns (total 5 columns):
 #   Column                  Non-Null Count  Dtype
---  ------                  --------------  -----
 0   CustomerID              200 non-null    int64
 1   Gender                  200 non-null    object
 2   Age                     200 non-null    int64
 3   Annual Income (k$)      200 non-null    int64
 4   Spending Score (1-100)  200 non-null    int64
dtypes: int64(4), object(1)
memory usage: 7.9+ KB
```

```python
df.drop(['CustomerID'], axis = 1, inplace = True)

df['Gender'] = df['Gender'].map({'Male': 0, 'Female': 1})
```

```python
X = df[['Age', 'Spending Score (1-100)']].iloc[:, :].values
```

```python
wcss = []
for i in range(1, 11):
    kmeans = KMeans(n_clusters = i, init = 'k-means++', random_state = 42)
    kmeans.fit(X)
    wcss.append(kmeans.inertia_)
plt.plot(range(1, 11), wcss, marker = 'o')
plt.show()
```



## ⌄ K-Means - Sklearn

```python
kmeans = KMeans(n_clusters = 4, init = 'k-means++', random_state = 42)
y_pred = kmeans.fit_predict(X)
```

## ⌄ K-Means - Scratch

```python
class KMeansScratch:
    def __init__(self, n_clusters, max_iter=300):
        self.n_clusters = n_clusters
        self.max_iter = max_iter

    def fit(self, X):
        n_samples, n_features = X.shape
        self.centroids = X[np.random.choice(n_samples, self.n_clusters, replace=False)]
        for _ in range(self.max_iter):
            labels = self.assign_clusters(X)
            new_centroids = self.update_centroids(X, labels)
            if np.allclose(self.centroids, new_centroids):
                break
            self.centroids = new_centroids
        return labels

    def assign_clusters(self, X):
        distances = np.sqrt(((X - self.centroids[:, np.newaxis])**2).sum(axis=2))
        return np.argmin(distances, axis=0)

    def update_centroids(self, X, labels):
        new_centroids = np.zeros_like(self.centroids)
        for i in range(self.n_clusters):
            new_centroids[i] = np.mean(X[labels == i], axis=0)
        return new_centroids
```

```python
kmeans_scratch = KMeansScratch(n_clusters=4)
y_pred1 = kmeans_scratch.fit(X)
```

## ⌄ Inference

```python
print("Sklearn - labels: \n", y_pred)
print("Scratch - labels: \n", y_pred1)
```

```
    Sklearn - labels:
     [3 2 1 2 3 2 1 2 1 2 1 2 1 2 1 2 3 3 1 2 3 2 1 2 1 2 1 3 1 2 1 2 1 2 1 2 1
     2 1 2 0 2 0 3 1 3 0 3 3 3 0 3 3 0 0 0 0 0 3 0 0 3 0 0 0 3 0 0 3 3 0 0 0 0
     0 3 0 3 3 0 0 3 0 0 3 0 0 3 3 0 0 3 0 3 3 3 0 3 0 3 3 0 0 3 0 3 0 0 0 0 0
     3 3 3 3 3 0 0 0 0 3 3 3 2 3 2 0 2 1 2 1 2 3 2 1 2 1 2 1 2 1 2 3 2 1 2 0 2
     1 2 1 2 1 2 1 2 1 2 1 2 0 2 1 2 1 2 1 2 1 3 1 2 1 2 1 2 1 2 1 2 1 2 1 2 3
     2 1 2 1 2 1 2 1 2 1 2 1 2]
    Scratch - labels:
     [2 1 0 1 2 1 0 1 0 1 0 1 0 1 0 1 2 2 0 1 2 1 0 1 0 1 3 2 0 1 0 1 0 1 0 1 0
     1 0 1 3 1 3 2 0 2 3 2 2 2 3 2 2 3 3 3 3 3 2 3 3 3 2 3 3 3 2 3 3 3 2 3 3 3
     3 2 3 2 2 3 3 2 3 3 2 3 3 2 2 3 3 2 3 2 3 2 2 3 3 2 3 2 3 3 3 3 3
     2 2 2 2 2 3 3 3 3 3 2 2 2 1 2 1 3 1 0 1 0 1 2 1 0 1 0 1 0 1 0 1 2 1 0 1 3 1
     0 1 0 1 0 1 0 1 0 1 0 1 3 1 0 1 0 1 0 1 0 2 0 1 0 1 0 1 0 1 0 1 0 1 0 1 3
     1 0 1 0 1 0 1 0 1 0 1 0 1]
```

```python
best_acc = 0
for perm in permutations(range(4)):
    matched_labels = np.array([perm[label] for label in y_pred1])
    acc = accuracy_score(y_pred, matched_labels)
    if acc > best_acc:
        best_acc = acc

print("Accuracy: ", best_acc)
```

```
    Accuracy:  0.985
```