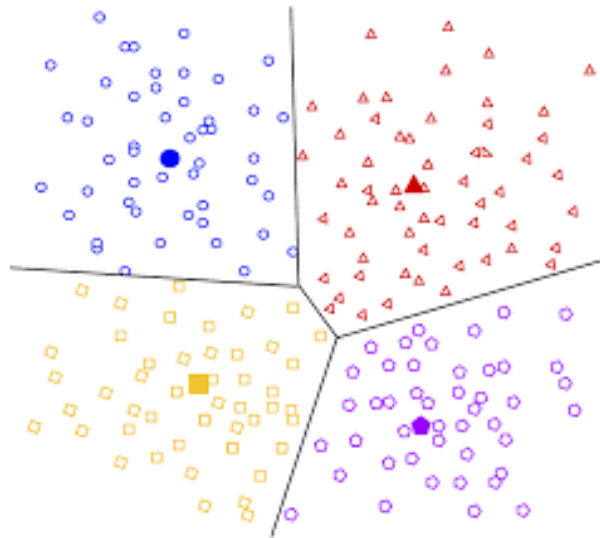


Kmeans Clustering

Using Mall Customers Data



Packages Used

```
In [1]: 1 import numpy as np
        2 import matplotlib.pyplot as plt
        3 import pandas as pd
        4 import sklearn as sk
```

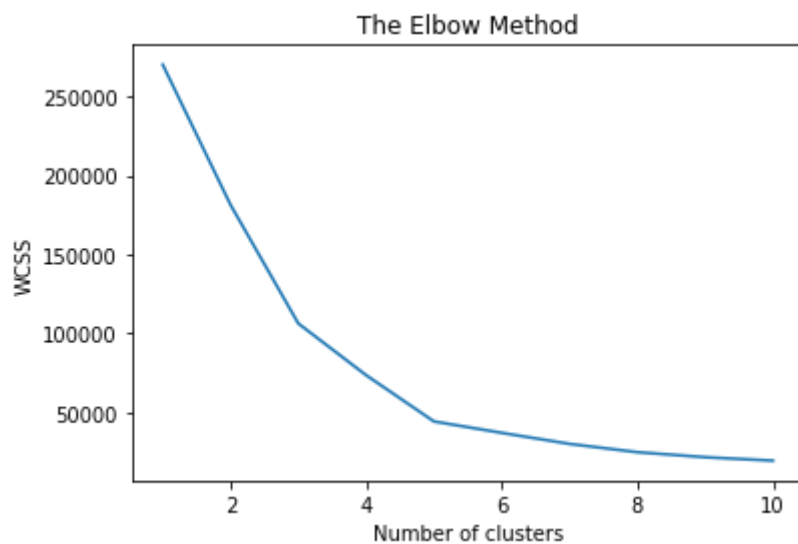
Load and Preprocess Data

```
In [2]: 1 # Importing the dataset
        2 dataset = pd.read_csv('../_resources/data/Mall_Customers.csv')
        3 X = dataset.iloc[:, [3, 4]].values
        4
```

Create Kmeans Class With all Methods

```
In [3]: 1 class Kmeans:
2
3
4     def __init__(self, k, seed = None, max_iter = 200):
5         self.k = k
6         self.seed = seed
7         self.centroids = []
8         if self.seed is not None:
9             np.random.seed(self.seed())
10        self.max_iter = max_iter
11
12    def initialise_centroids(self, data):
13        initial_centroids = np.random.permutation(data.shape[0])[:self.k]
14        self.centroids = data[initial_centroids]
15        return self.centroids
16
17    def assign_clusters(self, data):
18        if data.ndim == 1:
19            data = data.reshape(-1,1)
20        dist_to_centroid = sk.metrics.pairwise_distances(data, self.centroids,
21                                                         metric = 'euclidean')
22        self.cluster_labels = np.argmin(dist_to_centroid, axis = 1)
23        return self.cluster_labels
24
25    def update_centroids(self, data):
26        self.centroids = np.array([data[self.cluster_labels == i].mean(
27                                     for i in range(self.k))]
28        return self.centroids
29
30
31    def predict(self, data):
32        return self.assign_clusters(data)
33
34    def fit_kmeans(self, data):
35        self.centroids = self.initialise_centroids(data)
36        for iter in range(self.max_iter):
37            self.cluster_labels = self.assign_clusters(data)
38            self.centroids = self.update_centroids(data)
39            if iter % 100 == 0:
40                print('Running Model Iteration %d' %iter)
41        print('Model finished running')
42        return self
43
44
45
```

```
In [4]: 1 # Using the elbow method to find the optimal number of clusters
2 from sklearn.cluster import KMeans
3 wcss = []
4 for i in range(1, 11):
5     kmeans = KMeans(n_clusters = i, init = 'k-means++', random_state =
6     kmeans.fit(X)
7     wcss.append(kmeans.inertia_)
8 plt.plot(range(1, 11), wcss)
9 plt.title('The Elbow Method')
10 plt.xlabel('Number of clusters')
11 plt.ylabel('WCSS')
12 plt.show()
13
```



```
In [ ]: 1 # Fitting K-Means to the dataset
2 kmeans = Kmeans(k = 5)
3 y_kmeans = kmeans.predict(X)
4
5 # Visualising the clusters
6 plt.scatter(X[y_kmeans == 0, 0], X[y_kmeans == 0, 1], s = 100, c = 'red')
7 plt.scatter(X[y_kmeans == 1, 0], X[y_kmeans == 1, 1], s = 100, c = 'blue')
8 plt.scatter(X[y_kmeans == 2, 0], X[y_kmeans == 2, 1], s = 100, c = 'green')
9 plt.scatter(X[y_kmeans == 3, 0], X[y_kmeans == 3, 1], s = 100, c = 'cyan')
10 plt.scatter(X[y_kmeans == 4, 0], X[y_kmeans == 4, 1], s = 100, c = 'magenta')
11 plt.scatter(kmeans.cluster_centers_[0, 0], kmeans.cluster_centers_[0, 1], s = 100, c = 'black')
12 plt.title('Clusters of customers')
13 plt.xlabel('Annual Income (k$)')
14 plt.ylabel('Spending Score (1-100)')
15 plt.legend()
16 plt.show()
```

```
In [ ]: 1
```