

```
In [1]: 1 import numpy as np # Linear algebra
        2 import pandas as pd # data processing, CSV file I/O (e.g. pd.read_csv)
        3 import matplotlib.pyplot as plt
        4 import seaborn as sns
        5 import plotly as py
        6 import plotly.graph_objs as go
        7 from sklearn.cluster import KMeans
        8 import warnings
        9 import os
       10 warnings.filterwarnings("ignore")
       11 py.offline.init_notebook_mode(connected = True)
       12 #print(os.listdir("../input"))
```

```
In [2]: 1 df = pd.read_csv(r'Mall_Customers.csv')
        2 df.head()
```

```
Out[2]:
```

	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

```
In [3]: 1 df.shape
```

```
Out[3]: (200, 5)
```

```
In [4]: 1 df.describe()
```

```
Out[4]:
```

	CustomerID	Age	Annual Income (k\$)	Spending Score (1-100)
count	200.000000	200.000000	200.000000	200.000000
mean	100.500000	38.850000	60.560000	50.200000
std	57.879185	13.969007	26.264721	25.823522
min	1.000000	18.000000	15.000000	1.000000
25%	50.750000	28.750000	41.500000	34.750000
50%	100.500000	36.000000	61.500000	50.000000
75%	150.250000	49.000000	78.000000	73.000000
max	200.000000	70.000000	137.000000	99.000000

```
In [5]: 1 df.dtypes
```

```
Out[5]:
```

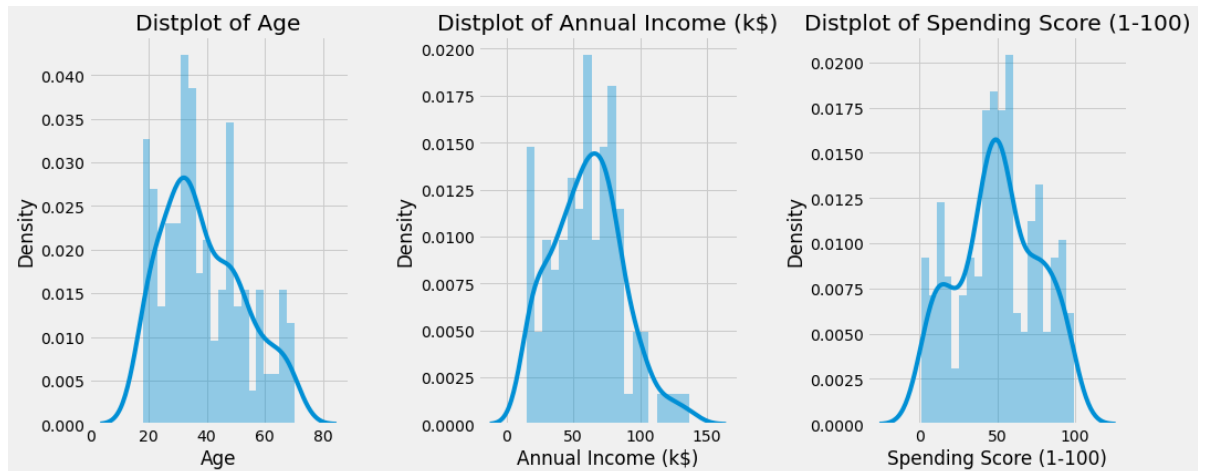
```
CustomerID      int64
Gender           object
```

```
In [6]: 1 df.isnull().sum()
```

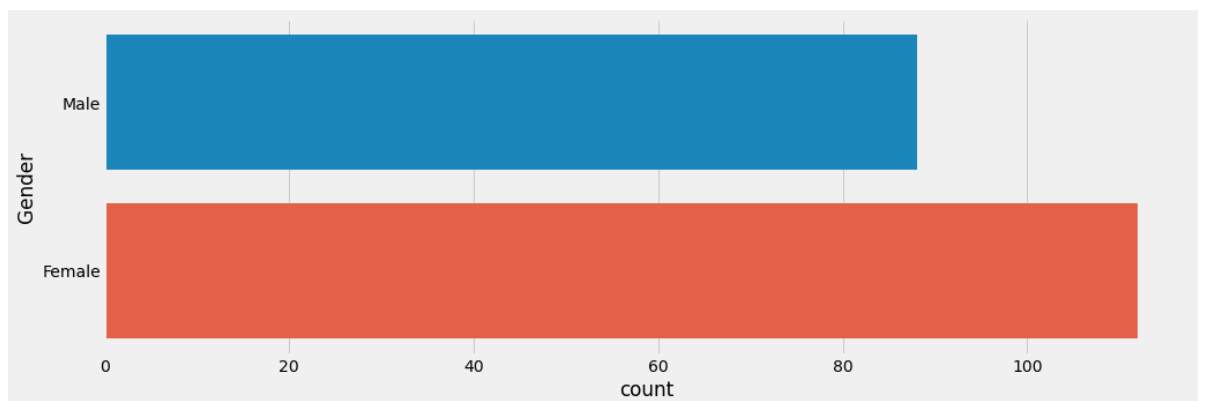
```
Out[6]: CustomerID      0
Gender           0
Age              0
Annual Income (k$)  0
Spending Score (1-100)  0
dtype: int64
```

```
In [7]: 1 plt.style.use('fivethirtyeight')
```

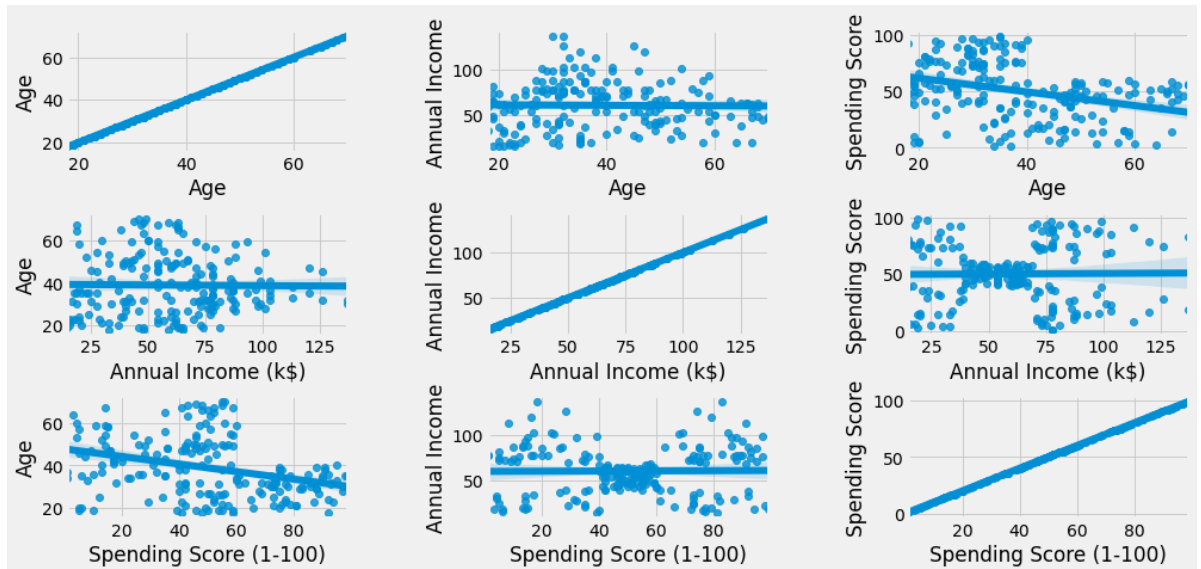
```
In [8]: 1 plt.figure(1 , figsize = (15 , 6))
2 n = 0
3 for x in ['Age' , 'Annual Income (k$)' , 'Spending Score (1-100)']:
4     n += 1
5     plt.subplot(1 , 3 , n)
6     plt.subplots_adjust(hspace = 0.5 , wspace = 0.5)
7     sns.distplot(df[x] , bins = 20)
8     plt.title('Distplot of {}'.format(x))
9 plt.show()
```



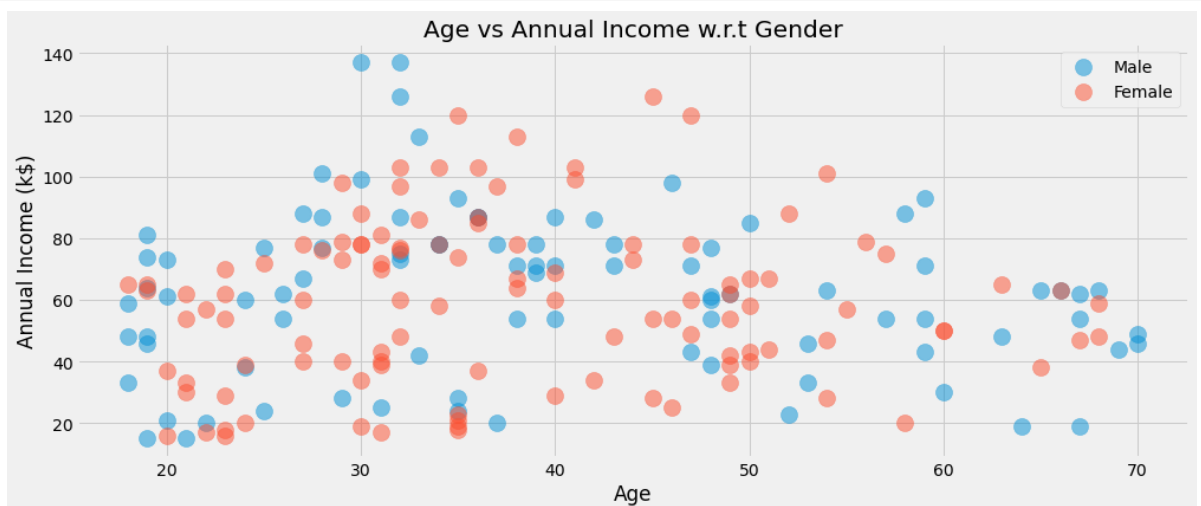
```
In [9]: 1 plt.figure(1 , figsize = (15 , 5))
2 sns.countplot(y = 'Gender' , data = df)
3 plt.show()
```



```
In [10]: 1 plt.figure(1 , figsize = (15 , 7))
2 n = 0
3 for x in ['Age' , 'Annual Income (k$)' , 'Spending Score (1-100)']:
4     for y in ['Age' , 'Annual Income (k$)' , 'Spending Score (1-100)']:
5         n += 1
6         plt.subplot(3 , 3 , n)
7         plt.subplots_adjust(hspace = 0.5 , wspace = 0.5)
8         sns.regplot(x = x , y = y , data = df)
9         plt.ylabel(y.split()[0]+' '+y.split()[1] if len(y.split()) > 1 else y)
10 plt.show()
```

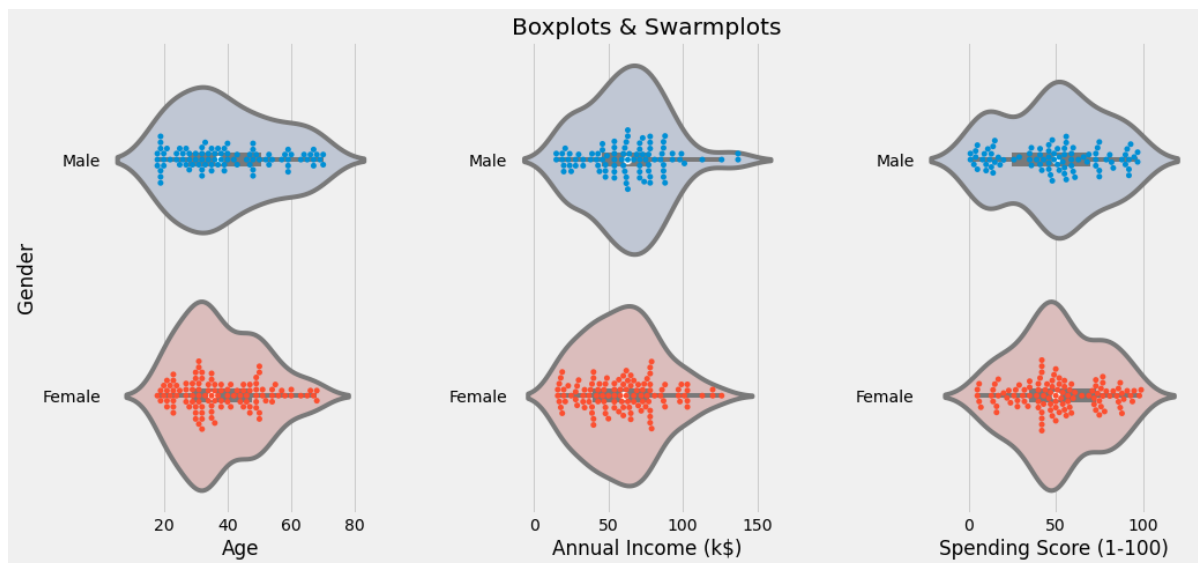


```
In [11]: 1 plt.figure(1 , figsize = (15 , 6))
2 for gender in ['Male' , 'Female']:
3     plt.scatter(x = 'Age' , y = 'Annual Income (k$)' , data = df[df['Gender'] == gender],
4                 s = 200 , alpha = 0.5 , label = gender)
5 plt.xlabel('Age') , plt.ylabel('Annual Income (k$)')
6 plt.title('Age vs Annual Income w.r.t Gender')
7 plt.legend()
8 plt.show()
```



Distribution of values in Age , Annual Income and Spending Score according to Gender

```
In [12]: 1 plt.figure(1 , figsize = (15 , 7))
2         n = 0
3         for cols in ['Age' , 'Annual Income (k$)' , 'Spending Score (1-100)']:
4             n += 1
5             plt.subplot(1 , 3 , n)
6             plt.subplots_adjust(hspace = 0.5 , wspace = 0.5)
7             sns.violinplot(x = cols , y = 'Gender' , data = df , palette = 'vlag')
8             sns.swarmplot(x = cols , y = 'Gender' , data = df)
9             plt.ylabel('Gender' if n == 1 else '')
10            plt.title('Boxplots & Swarmplots' if n == 2 else '')
11            plt.show()
```

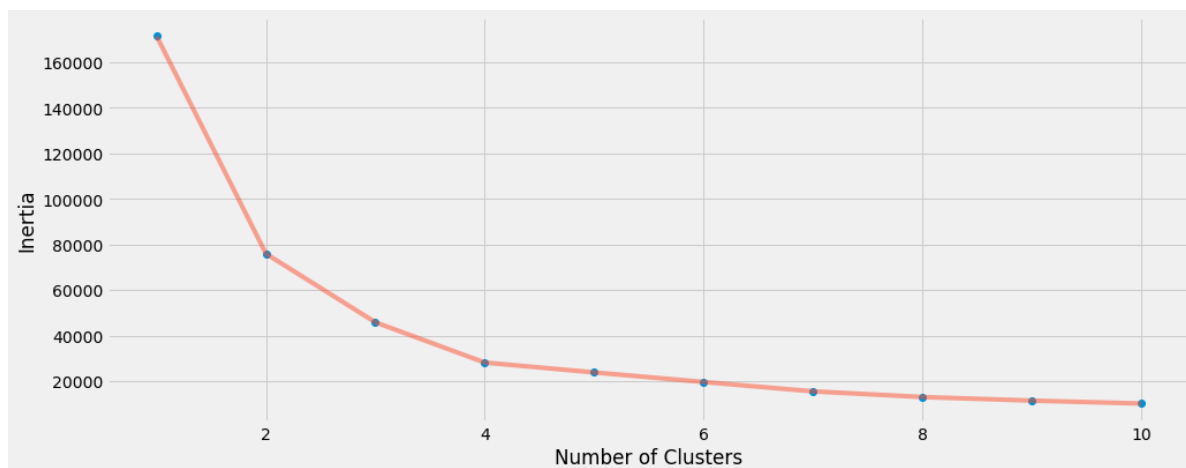


Clustering using K- means 1.Segmentation using Age and Spending Score

```
In [13]: 1 X1 = df[['Age' , 'Spending Score (1-100)']].iloc[:, :].values
2         inertia = []
3         for n in range(1 , 11):
4             algorithm = (KMeans(n_clusters = n , init='k-means++' , n_init = 10 , max
5                                tol=0.0001, random_state= 111 , algorithm='elkar
6             algorithm.fit(X1)
7             inertia.append(algorithm.inertia_)
```

In [14]:

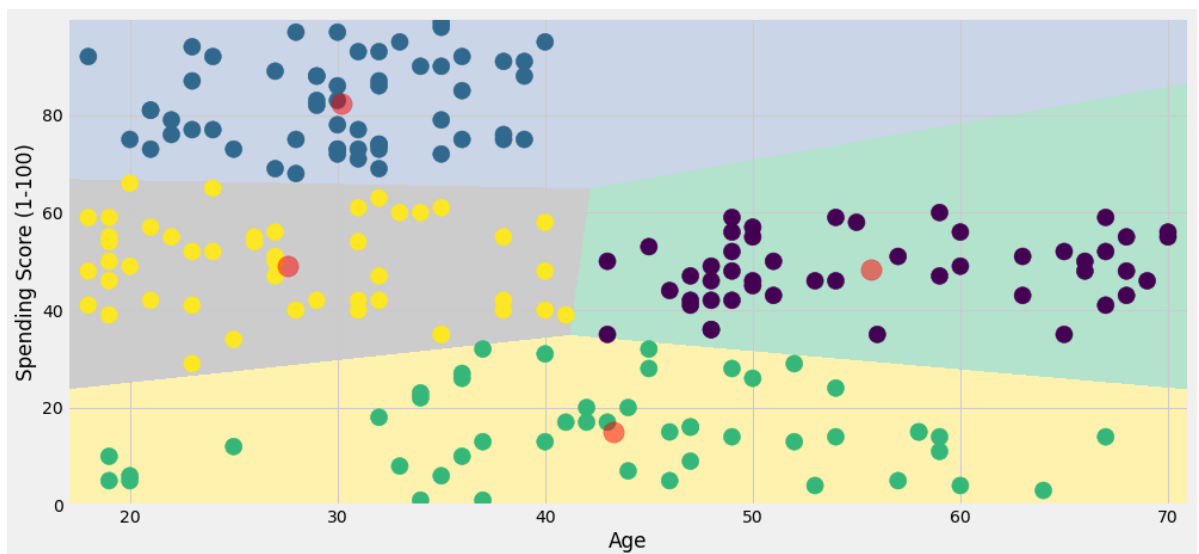
```
1 plt.figure(1 , figsize = (15 ,6))
2 plt.plot(np.arange(1 , 11) , inertia , 'o')
3 plt.plot(np.arange(1 , 11) , inertia , '-' , alpha = 0.5)
4 plt.xlabel('Number of Clusters') , plt.ylabel('Inertia')
5 plt.show()
```



```

In [15]: 1 algorithm = (KMeans(n_clusters = 4 ,init='k-means++', n_init = 10 ,max_ite
2           tol=0.0001, random_state= 111 , algorithm='elkar
3 algorithm.fit(X1)
4 labels1 = algorithm.labels_
5 centroids1 = algorithm.cluster_centers_
6 h = 0.02
7 x_min, x_max = X1[:, 0].min() - 1, X1[:, 0].max() + 1
8 y_min, y_max = X1[:, 1].min() - 1, X1[:, 1].max() + 1
9 xx, yy = np.meshgrid(np.arange(x_min, x_max, h), np.arange(y_min, y_max, h))
10 Z = algorithm.predict(np.c_[xx.ravel(), yy.ravel()])
11
12 plt.figure(1 , figsize = (15 , 7) )
13 plt.clf()
14 Z = Z.reshape(xx.shape)
15 plt.imshow(Z , interpolation='nearest',
16           extent=(xx.min(), xx.max(), yy.min(), yy.max()),
17           cmap = plt.cm.Pastel2, aspect = 'auto', origin='lower')
18
19 plt.scatter( x = 'Age' , y = 'Spending Score (1-100)' , data = df , c = labels1,
20             s = 200 )
21 plt.scatter(x = centroids1[:, 0] , y = centroids1[:, 1] , s = 300 , c = labels1)
22 plt.ylabel('Spending Score (1-100)') , plt.xlabel('Age')
23 plt.show()

```



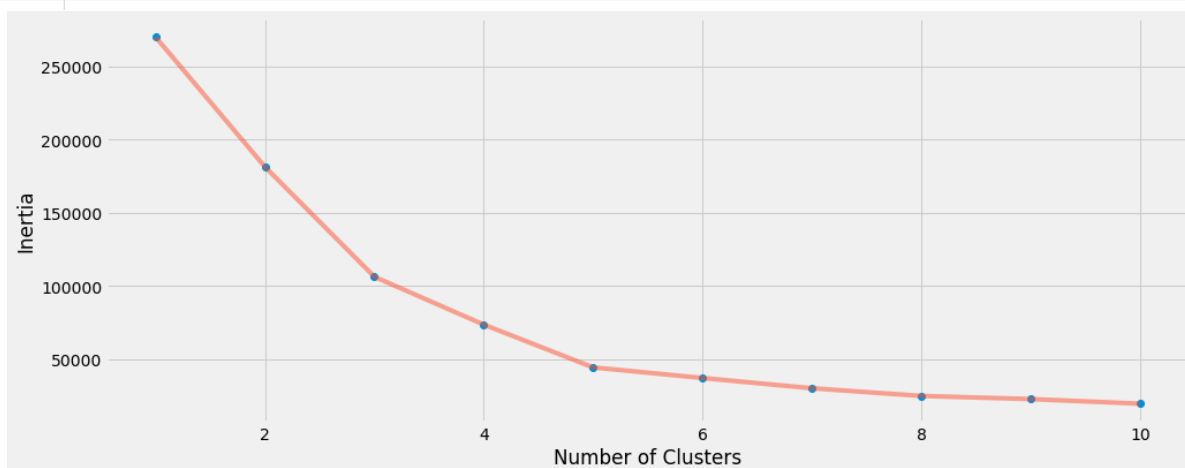
2. Segmentation using Annual Income and Spending Score

```

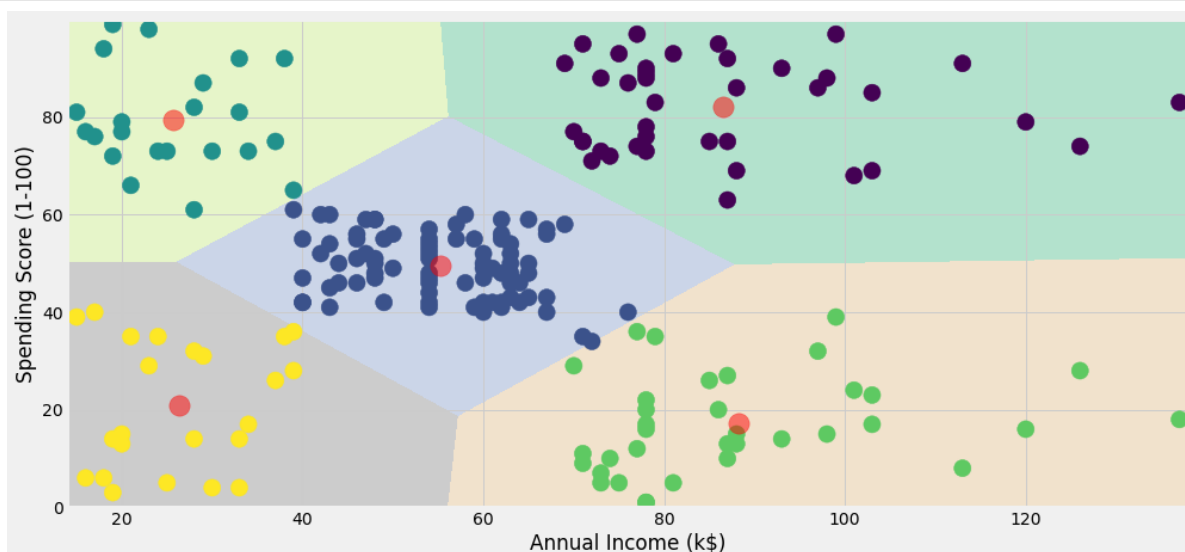
In [16]: 1 X2 = df[['Annual Income (k$)' , 'Spending Score (1-100)']].iloc[:, :].values
2 inertia = []
3 for n in range(1 , 11):
4     algorithm = (KMeans(n_clusters = n ,init='k-means++', n_init = 10 ,max_ite
5           tol=0.0001, random_state= 111 , algorithm='elkar
6     algorithm.fit(X2)
7     inertia.append(algorithm.inertia_)
8 plt.figure(1 , figsize = (15 , 6))
9 plt.plot(np.arange(1 , 11) , inertia , 'o')
10 plt.plot(np.arange(1 , 11) , inertia , '-', alpha = 0.5)
11 plt.xlabel('Number of Clusters') , plt.ylabel('Inertia')

```

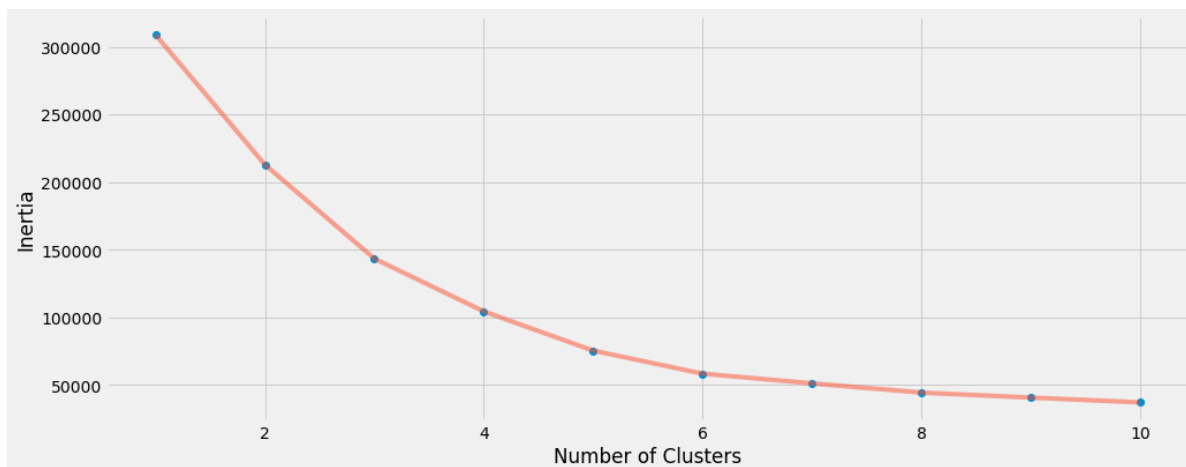
12 plt.show()



```
In [17]: 1 algorithm = (KMeans(n_clusters = 5 ,init='k-means++', n_init = 10 ,max_iter
2               tol=0.0001, random_state= 111 , algorithm='elkar
3 algorithm.fit(X2)
4 labels2 = algorithm.labels_
5 centroids2 = algorithm.cluster_centers_
6 h = 0.02
7 x_min, x_max = X2[:, 0].min() - 1, X2[:, 0].max() + 1
8 y_min, y_max = X2[:, 1].min() - 1, X2[:, 1].max() + 1
9 xx, yy = np.meshgrid(np.arange(x_min, x_max, h), np.arange(y_min, y_max, h))
10 Z2 = algorithm.predict(np.c_[xx.ravel(), yy.ravel()])
11
12 plt.figure(1 , figsize = (15 , 7) )
13 plt.clf()
14 Z2 = Z2.reshape(xx.shape)
15 plt.imshow(Z2 , interpolation='nearest',
16           extent=(xx.min(), xx.max(), yy.min(), yy.max()),
17           cmap = plt.cm.Pastel2, aspect = 'auto', origin='lower')
18
19 plt.scatter( x = 'Annual Income (k$)' ,y = 'Spending Score (1-100)' , data
20             s = 200 )
21 plt.scatter(x = centroids2[:, 0] , y = centroids2[:, 1] , s = 300 , c =
22 plt.ylabel('Spending Score (1-100)' ) , plt.xlabel('Annual Income (k$)')
23 plt.show()
```



```
In [18]: 1 X3 = df[['Age' , 'Annual Income (k$)' , 'Spending Score (1-100)']].iloc[: ,
2         inertia = []
3         for n in range(1 , 11):
4             algorithm = (KMeans(n_clusters = n ,init='k-means++', n_init = 10 ,max
5                                 tol=0.0001, random_state= 111 , algorithm='elkar
6             algorithm.fit(X3)
7             inertia.append(algorithm.inertia_)
8         plt.figure(1 , figsize = (15 ,6))
9         plt.plot(np.arange(1 , 11) , inertia , 'o')
10        plt.plot(np.arange(1 , 11) , inertia , '-' , alpha = 0.5)
11        plt.xlabel('Number of Clusters') , plt.ylabel('Inertia')
12        plt.show()
```



```
In [19]: 1 algorithm = (KMeans(n_clusters = 6 ,init='k-means++', n_init = 10 ,max_ite
2                                 tol=0.0001, random_state= 111 , algorithm='elkar
3         algorithm.fit(X3)
4         labels3 = algorithm.labels_
5         centroids3 = algorithm.cluster_centers_
6         df['label3'] = labels3
7         trace1 = go.Scatter3d(
8             x= df['Age'],
9             y= df['Spending Score (1-100)'],
10            z= df['Annual Income (k$)'],
11            mode='markers',
12            marker=dict(
13                color = df['label3'],
14                size= 20,
15                line=dict(
16                    color= df['label3'],
17                    width= 12
18                ),
19                opacity=0.8
20            )
21        )
22
23        data = [trace1]
24        layout = go.Layout(
25            # margin=dict(
26            #         l=0,
27            #         r=0,
```



```
28 #         b=0,
29 #         t=0
30 #     )
31     title= 'Clusters',
32     scene = dict(
33         xaxis = dict(title = 'Age'),
34         yaxis = dict(title = 'Spending Score'),
35         zaxis = dict(title = 'Annual Income')
36     )
37 )
38 fig = go.Figure(data=data, layout=layout)
39 py.offline.iplot(fig)
```

In []: 1