

# LGMVIP March Data Science Internship Task 1

Iris flower KMeans Clustering ML Project

In [1]:

```
1 #import packages
2 import numpy as np
3 import pandas as pd
4 import matplotlib.pyplot as plt
5 from sklearn.cluster import KMeans
6 from sklearn.decomposition import PCA
7 from sklearn.preprocessing import StandardScaler
```

In [2]:

```
1 #Import data
2 df=pd.read_csv("E:\Data files\IRIS.csv")
3 df
```

Out[2]:

	sepal_length	sepal_width	petal_length	petal_width	species
0	5.1	3.5	1.4	0.2	Iris-setosa
1	4.9	3.0	1.4	0.2	Iris-setosa
2	4.7	3.2	1.3	0.2	Iris-setosa
3	4.6	3.1	1.5	0.2	Iris-setosa
4	5.0	3.6	1.4	0.2	Iris-setosa
...	...	...	...	...	...
145	6.7	3.0	5.2	2.3	Iris-virginica
146	6.3	2.5	5.0	1.9	Iris-virginica
147	6.5	3.0	5.2	2.0	Iris-virginica
148	6.2	3.4	5.4	2.3	Iris-virginica

In [3]:

```
1 x = df.iloc[:, :-1]
2 y = df.iloc[:, -1]
```

In [4]:

```
1 x.head(5)
```

Out[4]:

	sepal_length	sepal_width	petal_length	petal_width
0	5.1	3.5	1.4	0.2
1	4.9	3.0	1.4	0.2
2	4.7	3.2	1.3	0.2
3	4.6	3.1	1.5	0.2
4	5.0	3.6	1.4	0.2

In [5]:

```
1 y.unique()
```

Out[5]:

```
array(['Iris-setosa', 'Iris-versicolor', 'Iris-virginica'], dtype=object)
```

In [6]:

```
1 #Standardize data
2 SS = StandardScaler()
3 X = SS.fit_transform(x)
4 X = pd.DataFrame(X, columns=df.columns[:-1])
5 X.head()
```

Out[6]:

	sepal_length	sepal_width	petal_length	petal_width
0	-0.900681	1.032057	-1.341272	-1.312977
1	-1.143017	-0.124958	-1.341272	-1.312977
2	-1.385353	0.337848	-1.398138	-1.312977
3	-1.506521	0.106445	-1.284407	-1.312977
4	-1.021849	1.263460	-1.341272	-1.312977

In [7]:

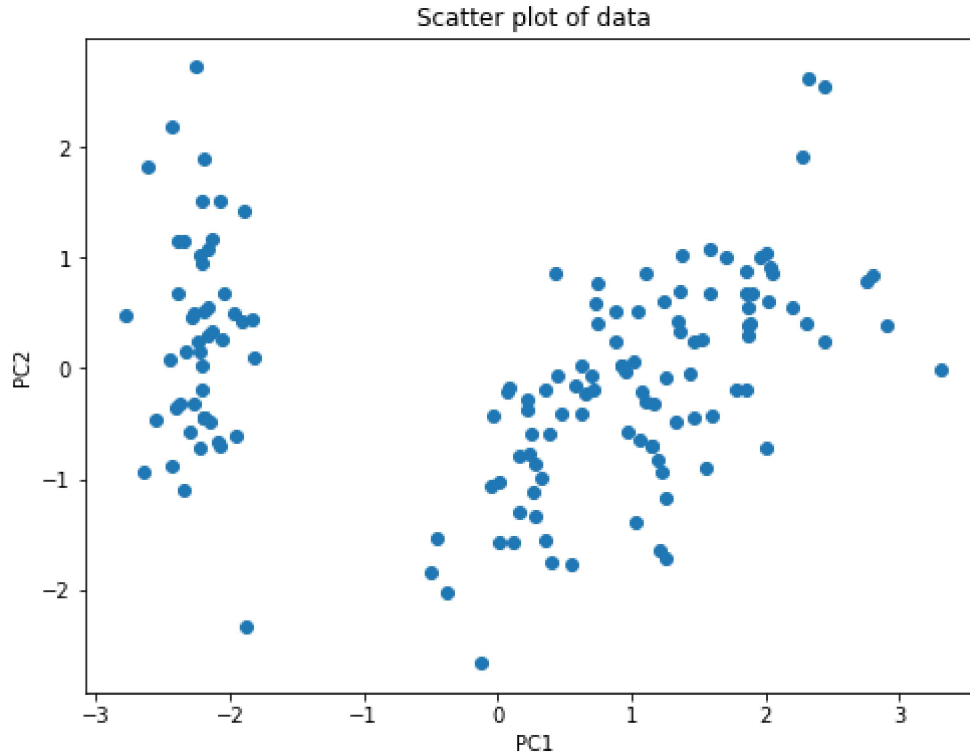
```
1 #Reduce dimension (PCA)
2 pca = PCA(n_components=2)
3 X = pca.fit_transform(X)
4 X = pd.DataFrame(X, columns=['PC1', 'PC2'])
5 X.head()
```

Out[7]:

	PC1	PC2
0	-2.264542	0.505704
1	-2.086426	-0.655405
2	-2.367950	-0.318477
3	-2.304197	-0.575368
4	-2.388777	0.674767

In [8]:

```
1 plt.figure(figsize=(8,6))
2 plt.scatter(X['PC1'], X['PC2'])
3 plt.xlabel('PC1')
4 plt.ylabel('PC2')
5 plt.title('Scatter plot of data')
6 plt.show()
```



In [9]:

```
1 # Train the model
2 k_means = KMeans(n_clusters=3)
3 k_means_model = k_means.fit(X)
```

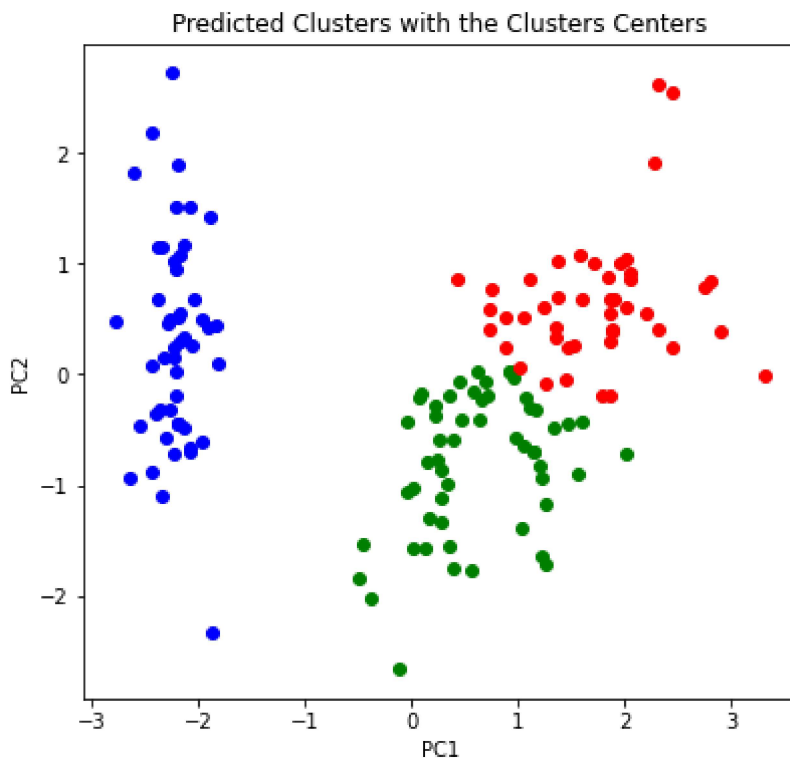
C:\Users\khush\anaconda3\lib\site-packages\sklearn\cluster\\_kmeans.py:133  
4: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP\_NUM\_THREADS=1.  
warnings.warn(

In [10]:

```
1 fig = plt.figure(figsize=(14,6))
2 #predicted
3 ax1 = fig.add_subplot(121)
4 pred = k_means_model.predict(X)
5 colours = ['blue', 'green', 'red']
6 for idx in range(3):
7     plt.scatter(X[pred == idx]['PC1'], X[pred == idx]['PC2'], c = colours[idx])
8 plt.xlabel('PC1')
9 plt.ylabel('PC2')
10 plt.title('Predicted Clusters with the Clusters Centers')
```

Out[10]:

Text(0.5, 1.0, 'Predicted Clusters with the Clusters Centers')



In [11]:

```

1 k =3  #3 clusters
2 df_clusters = [X[k_means_model.labels_==i] for i in range(k)]
3 len(df_clusters)

```

Out[11]:

3

In [12]:

```

1 stat_dict = {
2     'Cluster' : list(range(k)),
3     'Size' : [len(df_clusters[i]) for i in range(k)],
4     'Mean PC1' : [round(df_clusters[i]['PC1'].mean(), 2) for i in range(k)],
5     'Std PC1' : [round(df_clusters[i]['PC1'].std(), 2) for i in range(k)],
6     'Mean PC2' : [round(df_clusters[i]['PC2'].mean(), 2) for i in range(k)],
7     'Std PC2' : [round(df_clusters[i]['PC2'].std(), 2) for i in range(k)],
8 }

```

In [13]:

```

1 df_clusters_stats = pd.DataFrame(stat_dict)
2 df_clusters_stats

```

Out[13]:

	Cluster	Size	Mean PC1	Std PC1	Mean PC2	Std PC2
0	0	50	-2.22	0.20	0.29	0.94
1	1	55	0.61	0.56	-0.80	0.63
2	2	45	1.72	0.63	0.65	0.58

In [ ]:

1

In [ ]:

1