

```
import numpy as np
import pandas as pd
import matplotlib as mpl
import matplotlib.pyplot as plt
import seaborn as sns
sns.set(style="white", color_codes=True)
import warnings
warnings.filterwarnings("ignore")
%matplotlib inline
```

```
data = pd.read_csv('/content/Iris.csv')
data = data.drop('Id', axis=1)
data.head()
```



|   | SepalLengthCm | SepalWidthCm | PetalLengthCm | PetalWidthCm | 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 |

```
X = data.iloc[:,0:4]
y = data.iloc[:, -1]
print(X.sample(5))
print(y.sample(5))
```

```

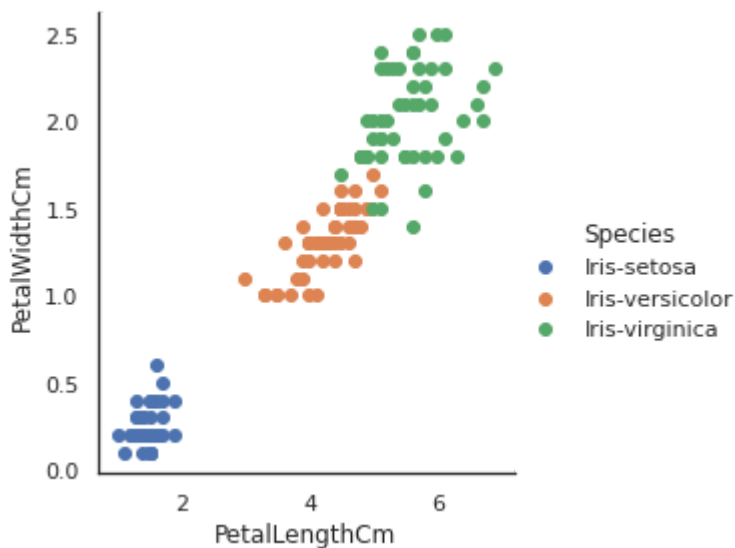
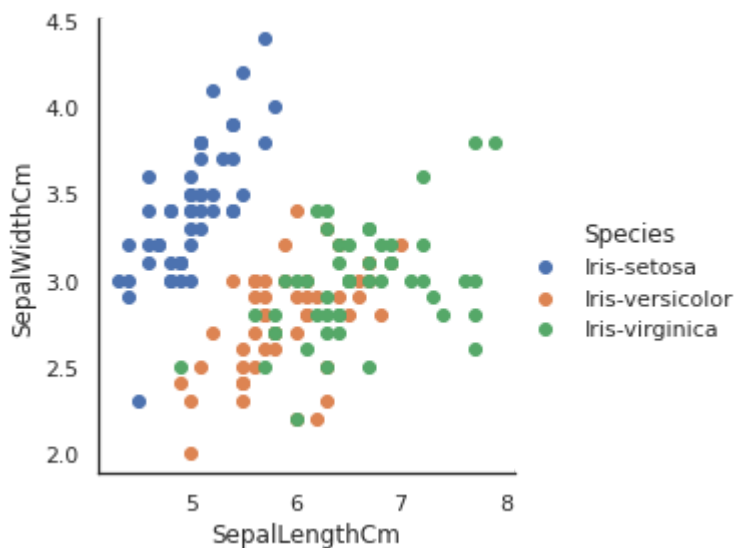
      SepalLengthCm  SepalWidthCm  PetalLengthCm  PetalWidthCm
29              4.7             3.2            1.6           0.2
140             6.7             3.1            5.6           2.4
116             6.5             3.0            5.5           1.8
118             7.7             2.6            6.9           2.3
109             7.2             3.6            6.1           2.5
108  Iris-virginica
61   Iris-versicolor
46   Iris-setosa
3    Iris-setosa
100  Iris-virginica
Name: Species, dtype: object
```

```
data["Species"].value_counts()
```

```

Iris-setosa      50
Iris-versicolor  50
Iris-virginica   50
Name: Species, dtype: int64
```

```
# use seaborn to make scatter plot showing species for each sample
sns.FacetGrid(data, hue="Species", size=4) \
    .map(plt.scatter, "SepalLengthCm", "SepalWidthCm") \
    .add_legend();
sns.FacetGrid(data, hue="Species", size=4) \
    .map(plt.scatter, "PetalLengthCm", "PetalWidthCm") \
    .add_legend();
```



```
from sklearn import preprocessing

scaler = preprocessing.StandardScaler()

scaler.fit(X)
X_scaled_array = scaler.transform(X)
X_scaled = pd.DataFrame(X_scaled_array, columns = X.columns)

X_scaled.sample(5)
```

|           | SepalLengthCm | SepalWidthCm | PetalLengthCm | PetalWidthCm |
|-----------|---------------|--------------|---------------|--------------|
| <b>15</b> | -0.173674     | 3.114684     | -1.284407     | -1.050031    |
| <b>65</b> | 1.038005      | 0.106445     | 0.364699      | 0.264699     |
| <b>41</b> | -1.627688     | -1.744778    | -1.398138     | -1.181504    |
| <b>64</b> | -0.294842     | -0.356361    | -0.090227     | 0.133226     |
| <b>39</b> | -0.900681     | 0.800654     | -1.284407     | -1.312977    |

## K-means clustering

```
from sklearn.cluster import KMeans

nclusters = 3
seed = 0

km = KMeans(n_clusters=nclusters, random_state=seed)
km.fit(X_scaled)

y_cluster_kmeans = km.predict(X_scaled)
y_cluster_kmeans

array([1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
       1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
       1, 1, 1, 1, 1, 1, 2, 2, 2, 2, 0, 0, 0, 2, 0, 0, 0, 0, 0, 0, 0, 0, 2,
       0, 0, 0, 0, 2, 0, 0, 0, 0, 2, 2, 2, 0, 0, 0, 0, 0, 0, 0, 2, 2, 0,
       0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 2, 0, 2, 2, 2, 2, 0, 2, 2, 2,
       2, 2, 2, 0, 0, 2, 2, 2, 2, 0, 2, 0, 2, 0, 2, 2, 0, 2, 2, 2, 2,
       2, 0, 0, 2, 2, 2, 0, 2, 2, 2, 0, 2, 2, 2, 0, 2, 2, 0], dtype=int32)
```

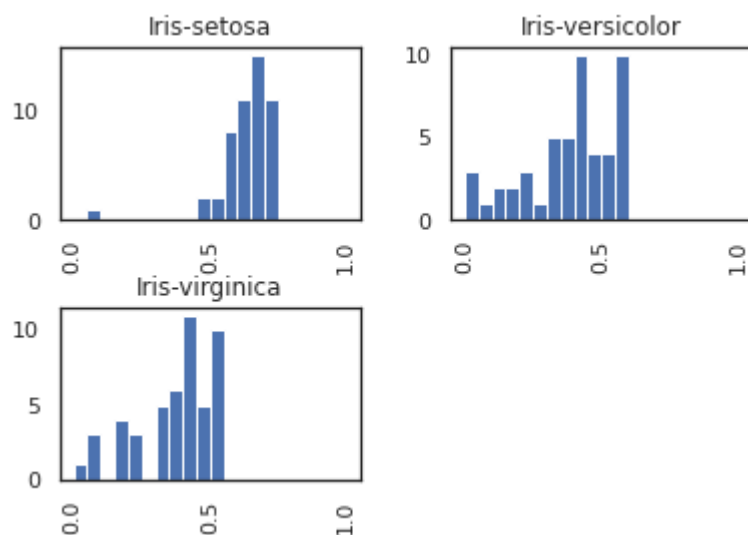
```
from sklearn import metrics
score = metrics.silhouette_score(X_scaled, y_cluster_kmeans)
score
```

```
0.4589717867018717
```

```
scores = metrics.silhouette_samples(X_scaled, y_cluster_kmeans)
sns.distplot(scores);
```



```
df_scores = pd.DataFrame()
df_scores['SilhouetteScore'] = scores
df_scores['Species'] = data['Species']
df_scores.hist(by='Species', column='SilhouetteScore', range=(0,1.0), bins=20);
```



## PCA

```
from sklearn.decomposition import PCA

ndimensions = 2

pca = PCA(n_components=ndimensions, random_state=seed)
pca.fit(X_scaled)
X_pca_array = pca.transform(X_scaled)
X_pca = pd.DataFrame(X_pca_array, columns=['PC1','PC2']) # PC=principal component
X_pca.sample(5)
```

|            | PC1       | PC2       |
|------------|-----------|-----------|
| <b>26</b>  | -2.052063 | 0.266014  |
| <b>130</b> | 2.435497  | 0.246654  |
| <b>60</b>  | -0.124697 | -2.658063 |
| <b>114</b> | 1.464062  | -0.444148 |
| <b>14</b>  | -2.192292 | 1.889979  |

```

y_id_array = pd.Categorical(data['Species']).codes
#Categorical.from_array(data['Species']).codes

df_plot = X_pca.copy()
df_plot['ClusterKmeans'] = y_cluster_kmeans
df_plot['SpeciesId'] = y_id_array # also add actual labels so we can use it in later plots
df_plot.head(5)

```

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

```

def plotData(df, groupby):
    "make a scatterplot of the first two principal components of the data, colored by the group"
    fig, ax = plt.subplots(figsize = (7,7))
    cmap = mpl.cm.get_cmap('prism')
    for i, cluster in df.groupby(groupby):
        cluster.plot(ax = ax,
                     kind = 'scatter',
                     x = 'PC1', y = 'PC2',
                     color = cmap(i/(nclusters-1)),
                     label = "%s %i" % (groupby, i),
                     s=30)

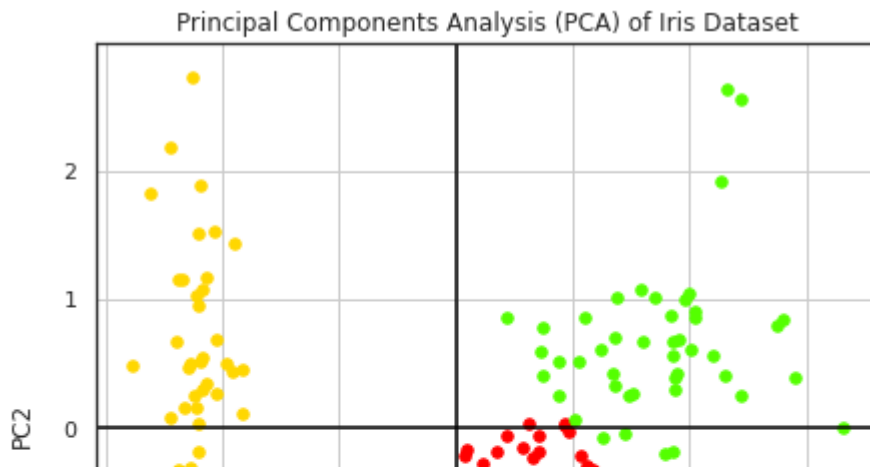
    ax.grid()
    ax.axhline(0, color='black')
    ax.axvline(0, color='black')
    ax.set_title("Principal Components Analysis (PCA) of Iris Dataset");

```

```

# plot the clusters each datapoint was assigned to
plotData(df_plot, 'ClusterKmeans')

```



## GMM



```
from sklearn.mixture import GaussianMixture
```

```
gmm = GaussianMixture(n_components=nclusters)
gmm.fit(X_scaled)
```

```
# predict the cluster for each data point
y_cluster_gmm = gmm.predict(X_scaled)
y_cluster_gmm
```

```
array([1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
       1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
       1, 1, 1, 1, 1, 1, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2,
       2, 2, 0, 2, 0, 2, 0, 2, 2, 2, 2, 0, 2, 2, 2, 2, 2, 0, 2, 2, 2, 2,
       2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
       0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
       0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0])
```

```
# add the GMM clusters to our data table and plot them
df_plot['ClusterGMM'] = y_cluster_gmm
plotData(df_plot, 'ClusterGMM')
```

