

## KMeans

```
In [14]: 1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 from scipy.spatial.distance import cdist
5 import seaborn as sns
```

```
In [15]: 1 def kmean(k,x,random_state=0):
2     s = x.shape[0]
3     m = np.random.seed(random_state)
4     r = np.random.choice(s,k,replace = False)
5     centroids = x[r]
6     while True:
7         old_centroids = centroids.copy()
8         dist = cdist(x,centroids)
9         label = np.argmin(dist,axis=1)
10        for cluster_id in range(k):
11            cluster = x[label == cluster_id]
12            if len(cluster):
13                centroids[cluster_id] = cluster.mean(axis=0)
14            if np.all(old_centroids != centroids):
15                break
16        return (centroids,label)
```

## Evaluation Metrics

```
In [16]: 1 def distortion(x,label,centroids):
2     return np.linalg.norm(x-(centroids[label]),axis=-1).mean()
3
4 def silhouette_score(k,x,label):
5     y = np.array([])
6     for i in range(k):
7         c = x[label == i]
8         if len(c):
9             y = np.append(y,cdist(c,c).mean())
10    a=y.mean()
11    z = np.array([])
12    for i in range(k):
13        c = x[label == i]
14        v = np.array([])
15        for j in range(k):
16            if i != j:
17                d = x[label == j]
18                if len(d) and len(c):
19                    v = np.append(v,cdist(c,d).mean())
20        z = np.append(z,np.min(v))
21    b=z.mean()
22    return (b-a)/np.max(np.dstack([a,b]),axis=-1).mean()
```

## Loading the Iris dataset

```
In [17]: 1 names = ["s_len", "s_wid", "p_len", "p_wid", "species"]
```

```
In [18]: 1 x_data = pd.read_csv("C:\\Users\\ASUS\\Iris.csv", names = names)
2 x = x_data.drop("species", axis=1).values
```

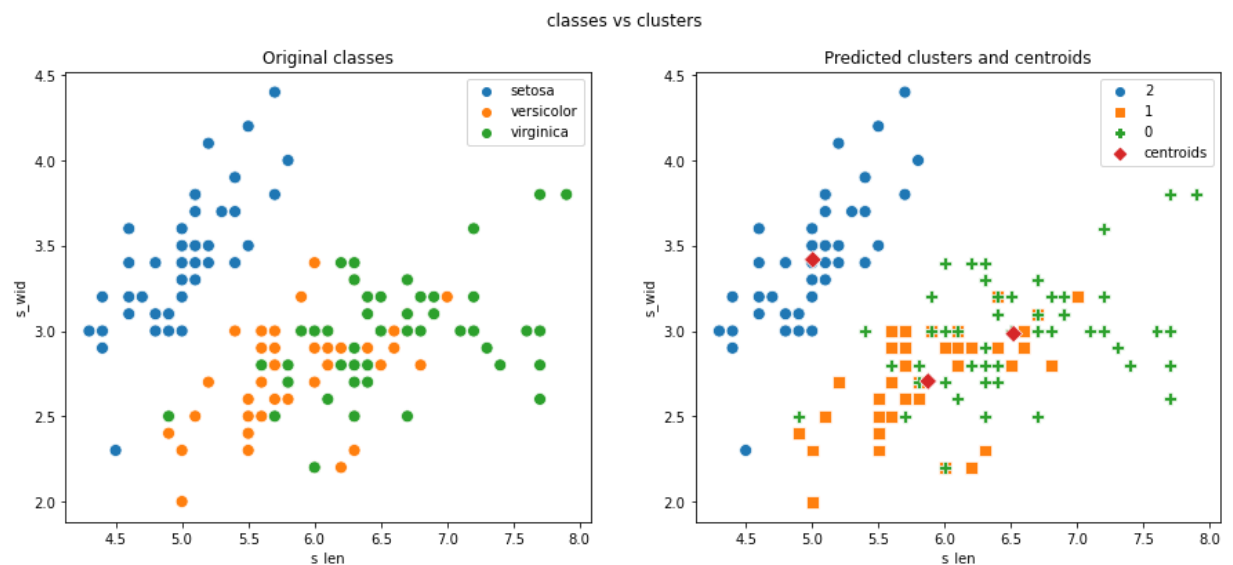
## Implementation of KMeans

```
In [19]: 1 centroids, label = kmean(3, x)
```

```
In [20]: 1 all_df = x_data.copy(deep=True)
2 clusters = pd.DataFrame(centroids, columns=names[:-1])
3 clusters["cluster"] = "centroids"
4 all_df["cluster"] = label.astype("str")
5 all_df = pd.concat([all_df, clusters])
```

## Comparing the clustering results with original classes

```
In [21]: 1 plt.figure(figsize=(15,6))
2 plt.suptitle("classes vs clusters")
3 plt.subplot(121)
4 plt.title("Original classes")
5 sns.scatterplot(data=all_df, x="s_len", y="s_wid", hue="species", s=80)
6 plt.legend(loc="upper right")
7 plt.subplot(122)
8 plt.title("Predicted clusters and centroids")
9 sns.scatterplot(data=all_df, x="s_len", y="s_wid", hue="cluster", style="cluster")
10 plt.legend(loc="upper right")
11 plt.show()
```



## Testing with multiple values of k (Hyperparameter tuning)

```
In [22]: 1 sil_coefs = []
2 distortions = []
3 K = np.arange(2,6)
4 for k in K:
5     centroids,label = kmean(k,x)
6     avg_sil_coef = silhouette_score(k,x,label)
7     dist = distortion(x,label,centroids)
8     print(f"For k={k:<4} Avg.Sil.Coeff: {avg_sil_coef:<10.5f} Distortion: {di
9     distortions.append(dist)
10    sil_coefs.append(avg_sil_coef)
```

For k=2	Avg.Sil.Coeff: 0.53721	Distortion: 1.30604
For k=3	Avg.Sil.Coeff: 0.59123	Distortion: 0.67576
For k=4	Avg.Sil.Coeff: 0.55830	Distortion: 0.59402
For k=5	Avg.Sil.Coeff: 0.43412	Distortion: 0.56721

### Using Elbow method to find the optimal value of k

```
In [23]: 1 plt.plot(K, distortions, 's-', markersize=8)
2 plt.xlabel('k')
3 plt.xticks(K)
4 plt.ylabel('Distortion')
5 plt.title('Elbow Method for Finding Optimal k')
6 plt.show()
```

