

Unsupervised learning

Unsupervised learning is a machine learning technique that involves finding patterns or structure in data without the use of explicit labels or a predefined target variable. In unsupervised learning, the algorithm is left to discover the underlying structure of the data on its own, without any prior knowledge or guidance. One of the most common unsupervised learning techniques is clustering.

Clustering is the process of grouping similar data points together based on their intrinsic properties or characteristics. In other words, clustering algorithms identify groups of data points that are similar to each other and different from the rest of the data. Clustering algorithms can be used to identify natural groupings or clusters within a dataset, which can be useful for a wide range of applications, such as customer segmentation, image segmentation, and anomaly detection.

There are many different clustering algorithms, each with its own strengths and weaknesses. Some of the most popular clustering algorithms include k-means clustering, hierarchical clustering, and density-based clustering. These algorithms differ in terms of their approach to grouping data points and their ability to handle different types of data and clustering structures.

K Means

K-means clustering is a popular unsupervised machine learning algorithm used for clustering data points into groups or clusters based on their similarities. It is an iterative algorithm that partitions a dataset into k clusters where k is a predefined number of clusters. The algorithm aims to minimize the sum of the squared distances between data points and their assigned cluster centers. The squared distance between a data point and a cluster center is also known as the within-cluster sum of squares (WCSS).

Here's a overview of the K-means algorithm:

Randomly select k data points from the dataset as the initial cluster centers.

Assign each data point in the dataset to the nearest cluster center based on its distance.

Recalculate the new cluster centers by taking the mean of all the data points assigned to each cluster.

Repeat steps 2 and 3 until convergence (i.e., until the assignment of data points to clusters no longer changes or a maximum number of iterations is reached).

K-means clustering can be used for a variety of applications, such as customer segmentation, image segmentation, and anomaly detection. However, the algorithm has some limitations. One limitation is that it requires the number of clusters to be specified beforehand, which can be challenging when the number of clusters is not known in advance. Another limitation is that it assumes clusters are spherical and equally sized, which may not always be the case in real-world datasets.

Despite these limitations, K-means clustering remains a widely used algorithm in unsupervised learning due to its simplicity and effectiveness for many clustering tasks.

```
In [4]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sn
from sklearn.cluster import KMeans
```

```
In [5]: customer_Data = pd.read_csv('Mall_Customers.csv')
```

```
In [6]: customer_Data.head()
```

```
Out[6]:
```

	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 [7]: customer_Data.shape
```

```
Out[7]: (200, 5)
```

```
In [8]: customer_Data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 200 entries, 0 to 199
Data columns (total 5 columns):
#   Column                Non-Null Count  Dtype
---  -
0   CustomerID            200 non-null   int64
1   Gender                200 non-null   object
2   Age                   200 non-null   int64
3   Annual Income (k$)    200 non-null   int64
4   Spending Score (1-100) 200 non-null   int64
dtypes: int64(4), object(1)
memory usage: 7.9+ KB
```

```
In [9]: customer_Data.isnull().sum()
```

```
Out[9]: CustomerID          0  
Gender          0  
Age             0  
Annual Income (k$)  0  
Spending Score (1-100)  0  
dtype: int64
```

```
In [10]: X = customer_Data.iloc[:, [3, 4]].values
```

```
In [11]: X
```

```
Out[11]: array([[ 15,  39],  
               [ 15,  81],  
               [ 16,   6],  
               [ 16,  77],  
               [ 17,  40],  
               [ 17,  76],  
               [ 18,   6],  
               [ 18,  94],  
               [ 19,   3],  
               [ 19,  72],  
               [ 19,  14],  
               [ 19,  99],  
               [ 20,  15],  
               [ 20,  77],  
               [ 20,  13],  
               [ 20,  79],  
               [ 21,  35],  
               [ 21,  66],  
               [ 23,  29],  
               [ 23,  98],  
               [ 24,  35],  
               [ 24,  73],  
               [ 25,   5],  
               [ 25,  73],  
               [ 28,  14],  
               [ 28,  82],  
               [ 28,  32],  
               [ 28,  61],  
               [ 29,  31],  
               [ 29,  87],  
               [ 30,   4],  
               [ 30,  73],  
               [ 33,   4],  
               [ 33,  92],  
               [ 33,  14],  
               [ 33,  81],  
               [ 34,  17],  
               [ 34,  73],  
               [ 37,  26],  
               [ 37,  75],  
               [ 38,  35],  
               [ 38,  92],  
               [ 39,  36],  
               [ 39,  61],  
               [ 39,  28],  
               [ 39,  65],  
               [ 40,  55],  
               [ 40,  47],  
               [ 40,  42],  
               [ 40,  42],  
               [ 42,  52],  
               [ 42,  60],  
               [ 43,  54],  
               [ 43,  60],  
               [ 43,  45],  
               [ 43,  41],  
               [ 44,  50],  
               [ 44,  46],  
               [ 46,  51],  
               [ 46,  46],  
               [ 46,  56],  
               [ 46,  55],  
               [ 47,  52],  
               [ 47,  59],  
               [ 48,  51],  
               [ 48,  59],  
               [ 48,  50],  
               [ 48,  48],  
               [ 48,  59],  
               [ 48,  47],  
               [ 49,  55],
```

[49, 42],
[50, 49],
[50, 56],
[54, 47],
[54, 54],
[54, 53],
[54, 48],
[54, 52],
[54, 42],
[54, 51],
[54, 55],
[54, 41],
[54, 44],
[54, 57],
[54, 46],
[57, 58],
[57, 55],
[58, 60],
[58, 46],
[59, 55],
[59, 41],
[60, 49],
[60, 40],
[60, 42],
[60, 52],
[60, 47],
[60, 50],
[61, 42],
[61, 49],
[62, 41],
[62, 48],
[62, 59],
[62, 55],
[62, 56],
[62, 42],
[63, 50],
[63, 46],
[63, 43],
[63, 48],
[63, 52],
[63, 54],
[64, 42],
[64, 46],
[65, 48],
[65, 50],
[65, 43],
[65, 59],
[67, 43],
[67, 57],
[67, 56],
[67, 40],
[69, 58],
[69, 91],
[70, 29],
[70, 77],
[71, 35],
[71, 95],
[71, 11],
[71, 75],
[71, 9],
[71, 75],
[72, 34],
[72, 71],
[73, 5],
[73, 88],
[73, 7],
[73, 73],
[74, 10],
[74, 72],
[75, 5],
[75, 93],
[76, 40],
[76, 87],
[77, 12],
[77, 97],
[77, 36],
[77, 74],
[78, 22],
[78, 90],
[78, 17],
[78, 88],
[78, 20],
[78, 76],
[78, 16],
[78, 89],
[78, 1],
[78, 78],
[78, 1],
[78, 73],

```
[ 79, 35],
[ 79, 83],
[ 81, 5],
[ 81, 93],
[ 85, 26],
[ 85, 75],
[ 86, 20],
[ 86, 95],
[ 87, 27],
[ 87, 63],
[ 87, 13],
[ 87, 75],
[ 87, 10],
[ 87, 92],
[ 88, 13],
[ 88, 86],
[ 88, 15],
[ 88, 69],
[ 93, 14],
[ 93, 90],
[ 97, 32],
[ 97, 86],
[ 98, 15],
[ 98, 88],
[ 99, 39],
[ 99, 97],
[101, 24],
[101, 68],
[103, 17],
[103, 85],
[103, 23],
[103, 69],
[113, 8],
[113, 91],
[120, 16],
[120, 79],
[126, 28],
[126, 74],
[137, 18],
[137, 83]], dtype=int64)
```

```
In [13]: wcss = []
         for i in range(1, 11):
             kmeans = KMeans(n_clusters=i, init='k-means++', random_state=42)
             kmeans.fit(X)
             wcss.append(kmeans.inertia_)
```

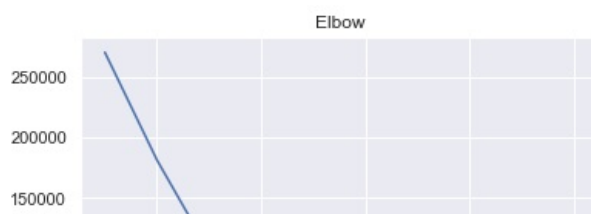
C:\Users\User\anaconda3\lib\site-packages\sklearn\cluster_kmeans.py:881: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=1.

```
warnings.warn(
```

```
In [14]: wcss
```

```
Out[14]: [269981.280000000014,
          181363.59595959607,
          106348.37306211119,
          73679.78903948837,
          44448.45544793369,
          37233.81451071002,
          30259.657207285458,
          25011.839349156595,
          21850.16528258562,
          19672.07284901432]
```

```
In [16]: sn.set()
         plt.plot(range(1,11), wcss)
         plt.title("Elbow")
         plt.show()
```



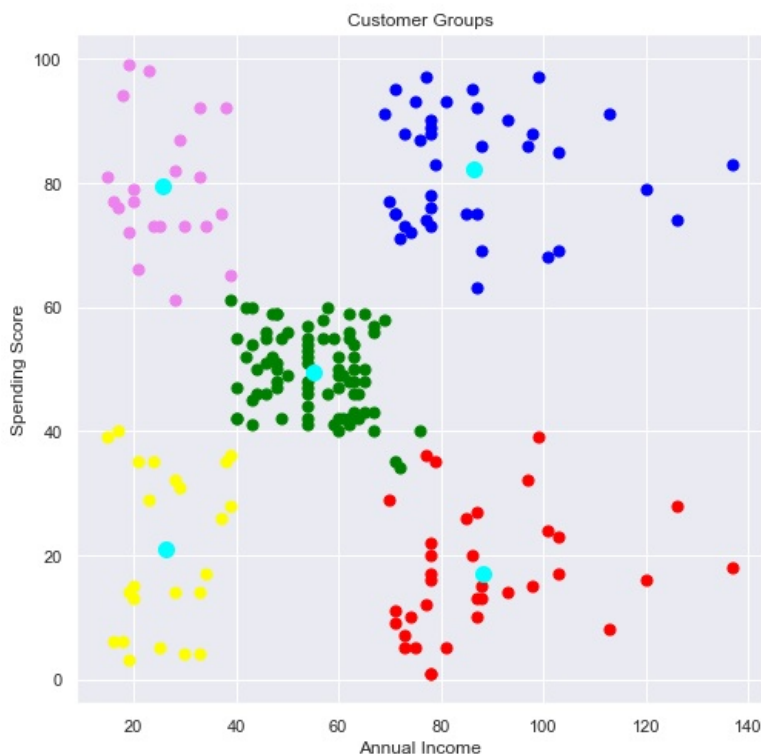


```
In [18]: kmeans = KMeans(n_clusters=5, init='k-means++', random_state=42)
Y = kmeans.fit_predict(X)
Y
```

```
Out[18]: array([2, 3, 2, 3, 2, 3, 2, 3, 2, 3, 2, 3, 2, 3, 2, 3, 2, 3, 2, 3,
2, 3, 2, 3, 2, 3, 2, 3, 2, 3, 2, 3, 2, 3, 2, 3, 2, 3, 2, 0,
2, 3, 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, 0, 0, 0, 0, 0, 0, 0,
0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 4, 1, 4, 0, 4, 1, 4, 1, 4,
0, 4, 1, 4, 1, 4, 1, 4, 1, 4, 0, 4, 1, 4, 1, 4, 1, 4, 1, 4,
1, 4, 1, 4, 1, 4, 1, 4, 1, 4, 1, 4, 1, 4, 1, 4, 1, 4, 1, 4,
1, 4, 1, 4, 1, 4, 1, 4, 1, 4, 1, 4, 1, 4, 1, 4, 1, 4, 1, 4,
1, 4])
```

```
In [19]: plt.figure(figsize=(8,8))
plt.scatter(X[Y==0,0], X[Y==0,1], s=50, c='green', label='Cluster 1')
plt.scatter(X[Y==1,0], X[Y==1,1], s=50, c='red', label='Cluster 2')
plt.scatter(X[Y==2,0], X[Y==2,1], s=50, c='yellow', label='Cluster 3')
plt.scatter(X[Y==3,0], X[Y==3,1], s=50, c='violet', label='Cluster 4')
plt.scatter(X[Y==4,0], X[Y==4,1], s=50, c='blue', label='Cluster 5')

plt.scatter(kmeans.cluster_centers_[0,0], kmeans.cluster_centers_[0,1], s=100, c='cyan', label='Centroids')
plt.title('Customer Groups')
plt.xlabel('Annual Income')
plt.ylabel('Spending Score')
plt.show()
```



In []:

In []:

In []:

In []:

In []:

In []:

In []:

In []:

In []:

In []:

In []:

In []:

In []:

In []: