

import pandas as pd

from sklearn.preprocessing import LabelEncoder.

from sklearn.decomposition import PCA

import seaborn as sns.

from sklearn.preprocessing import StandardScaler.

import matplotlib.pyplot as plt.

import numpy as np

from sklearn.cluster import ~~K-Means~~ ^{K-Means}

from sklearn.decomposition import PCA

import warnings

warnings.filterwarnings('ignore')

import scipy.cluster.hierarchy as shc.

from sklearn.cluster import ~~KMeans~~ ^{KMeans}, AgglomerativeClustering

import statistics.

import numpy as np.

=> X = np.array([[1, 2], [1, 4], [1, 4], [1, 0], [10, 2], [10, 4], [10, 4], [10, 0]])

=> X = X[:, 0]

Y = X[:, 1]

sns.scatterplot(X, Y)

=> km = KMeans(n_clusters=2, random_state=0)

km.fit(X)

=> km.cluster_centers_

X = X[:, 0]

Y = X[:, 1]

sns.scatterplot(X, Y, hue=km.labels_)

=> km.predict([[2, 3]])

Step: 1

=> import pandas as pd.

data = pd.read_csv("HALL - customers.csv")

data.head()

=> data.shape

=> data.columns

=> data.dtypes.

=> data.info()

=> data.value_counts()

Lab11. Shopping Mall Customer Segmentation using Clustering

Objectives

In this lab, you will detect clusters from customer data and perform analytics to understand customers who visit malls, using KMeans and Agglomerative Clustering methods.

Learning Objectives

After completing this lab, you will be able to

- Perform Skew analysis and interpretation
- Check if data normalization is required
- Build KMeans model with the required no. of clusters
- Visualize clusters based on selected features
- Select the best value for K using inertia error values and Elbow method
- Perform Cluster Analysis to understand cluster statistics
- Reduce dimensions of data using PCA and build KMeans model
- Build MeanShift clustering model on dimensions reduced data
- Create hierarchical clusters using Agglomerative Clustering
- Visualize hierarchical clusters using Dendrogram

Business Use Case

You are given the data of customers who have visited your mall. The details of customers such as age, annual income, spending score and gender are collected for each customer. They have asked you to analyse this data and give insights. This will help them to design promotion strategies, marketing models and others so as to reach out specific group of customers. Also, those group of customers will be targeted via surveys to collect further details. Based on that feedback we can decide whether the new strategy is good for that customer segment or not, even before the strategy is released.

Step1. [Understand Data]

- Using Pandas, import "Mall_Customers.csv" file and print properties such as head, shape, columns, dtype, info and value_counts.
- For example: `customers_data = pd.read_csv("Mall_Customers.csv")`

Step2. [Label encode gender]

- Genre (ie., gender) is a string, so label encode into binary

Step3. [Check for variance]

- Use describe() on your data frame and check for variance. If variance is high for float columns, you need to normalize. Otherwise, ignore

Step4. [Check skewness]

- Check if float columns are skewed. Use skew() on your data frame. If skew value is greater than 0.75, then you can perform log transformation on those skew columns.

Step5. [Pair plot]

- Draw pair plot and observe correlations.

Step6. [Build KMeans]

- Create and fit KMeans (n_clusters variable can be set with any value)
- Print label_ and cluster_centers_ values

Step7. [Scatter plot]

- Draw scatter plot between any two features with hue as "labels_". This figure shows 5 clusters.

Step 2:

=> data.drop(['Customers ID'], axis=1, inplace=True)

=> label_encoder = LabelEncoder()

data['Genre'] = label_encoder.fit_transform(data['Genre'])

data

=> data.course.value_counts()

=> df = data.copy()

Step 3:

=> data.describe(include='all')

=> data.info()

=> data.var()

=> data.corr()

Step 4:

=> data.show()

=> data.sort_values(by=['Genre', 'Age', 'Annual Income (K\$)'],
encoding=False, inplace=True)

Step 5:

=> sns.pairplot(data)

Step 6:

=> model = KMeans(n_clusters=5)

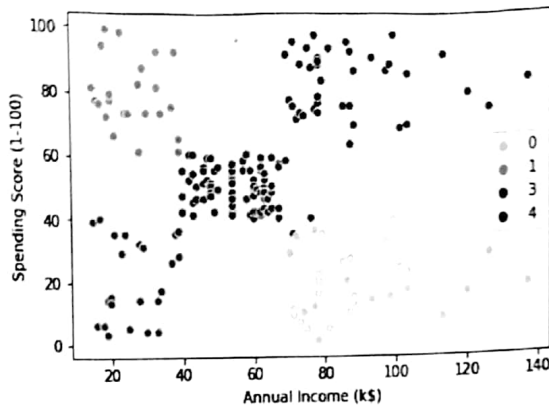
=> model.fit(data)

=> model.labels_

=> model.cluster_centers_

=> sns.scatterplot(data['Age'], data['Annual Income (K\$)'],
hue=model.labels_)

plt.show()

**Step8. [Cluster Analysis].**

Now, predict cluster labels for the same data. For example,

```
kmeans2 = KMeans(n_clusters = 5, init='k-means++')
kmeans2.fit(customers_data)
pred = kmeans2.predict(customers_data)
```

Now, add a new column for pred in a new dataframe, such as

```
frame = pd.DataFrame(customers_data)
frame['cluster'] = pred
```

This will create a new column to frame. That means, you have added a cluster prediction column, whose values say the cluster number to which the row belongs to. That is, that customer belongs to that cluster number.

Now, group customers based on cluster number. Remember, here we have 5 clusters from 0 to 4.

For each cluster group, print the following details.

```
Average age: 45.21739130434783
Average annual income: 26.304347826086957
Deviation of the mean for annual income: 7.893811054517766
No of customers ie shape: (23, 5)
From those customers we have 9 male and 14 female
```

Step9. [Find the best number of clusters]

Compute inertia value as shown below

```
SSE = []
for clust in range(1,20):
    km = KMeans(n_clusters=clust, init="k-means++")
    km = km.fit(customers_data)
    SSE.append(km.inertia_)
```

Plot a line chart between cluster number and its inertia value. You will get graph as below. Can you identify the best number of clusters from this graph?

Step: 7

=> sns.scatterplot (data['Annual Income (k\$)'], data['spending score (1-100)'], hue = model.labels').

plt.show()

Step: 8

=> data.head()

=> x = data.iloc[:, [2, 3]].values.

=> kmeans2 = KMeans(n_clusters=5, init = 'k-means++')

kmeans2.fit(x)

pred = kmeans2.predict(x)

=> frame = pd.DataFrame(x).

frame['cluster'] = pred.

data['cluster'] = pred

=> frame.cluster.values.tolist().

=> frame

=> d_frame0 = data[data['cluster'] == 0]

d_frame1 = data[data['cluster'] == 1]

d_frame2 = data[data['cluster'] == 2]

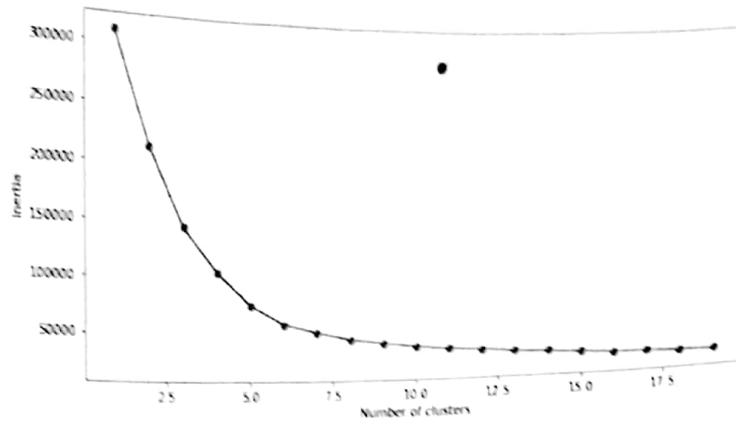
d_frame3 = data[data['cluster'] == 3]

d_frame4 = data[data['cluster'] == 4]

=> print("Average age for cluster 0:", d_frame0['Age'].mean())

print("Average Annual Income for cluster 0:", d_frame0['Annual Income (k\$)'].mean()).

print("Deviation of the mean for annual Income:", statistics.stdev(d_frame0['Annual Income (k\$)']))

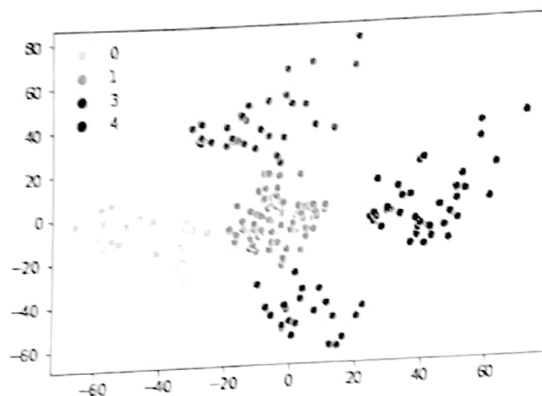


Step10. [Reduce Dimensions using PCA]

- Reduce 4 dimensions into 2 dimensions using PCA
- Create KMeans model, fit on the reduced dataset
- Print cluster_centers_ and labels_

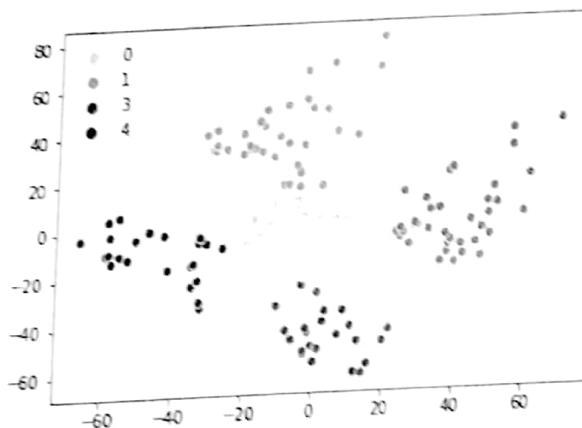
Step11. [Scatter plot]

- Draw a scatter plot between the 2 reduced dimensions, with hue as label_
- Your scatter plot may look like below



Step12. [MeanShift clustering]

- Create MeanShift clustering model and fit on the reduced data of PCA and visualize clusters on the reduced data, as shown below.



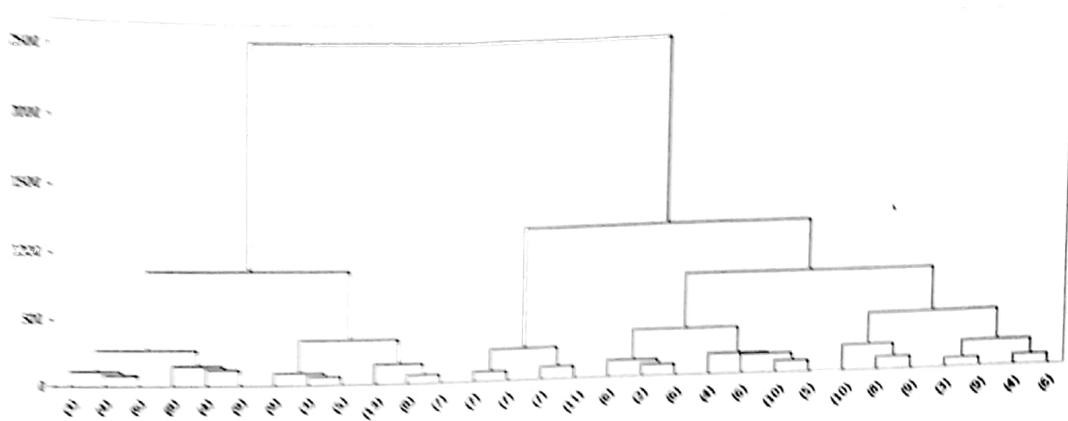
```

print("No. of customers i.e. shape of cluster 0:", dframe(0, slope)
print("From those customers we have", dframe(0, 'Age', 'value - counts'
                                             ($[0], "female"))
=> print("Average age for cluster 1:", dframe(1, 'Age', 'mean')).
print("Deviation of the Mean for annual income:", statistics.stdev -
      (dframe(1, 'Annual Income ($)', 'mean')).
print("No. of customers i.e. shape of cluster 1:", dframe(1, 'shape')
print("From those customers we have", dframe(1, 'Age', 'value - counts'
                                             ($[1], "male end", dframe(1, 'Age', 'value - counts'
                                             ($[1], "female"))
=> print("Average age for cluster 2:", dframe(2, 'Age', 'mean'))
print("Average Annual Income for cluster 2:", dframe(2, 'Annual Income'
                                                         ($[2], 'mean'))
print("Deviation of the Mean for annual income:", statistics.stdev -
      (dframe(2, 'Annual Income ($)', 'mean'))
=> print("Average age for cluster 3:", dframe(3, 'Age', 'mean'))
print("Average Annual Income for cluster 3:", dframe(3, 'Annual Income'
                                                         ($[3], 'mean'))
print("Deviation of the mean for annual income:", statistics.stdev -
      (dframe(3, 'Annual Income ($)', 'mean'))
print("No. of Customers i.e. shape of cluster 3:", dframe(3, 'shape')
print("From those customers we have", dframe(3, 'Age', 'value - counts'
                                             ($[3], "male end", dframe(3, 'Age', 'value - counts'
                                             ($[3], "female"))

```

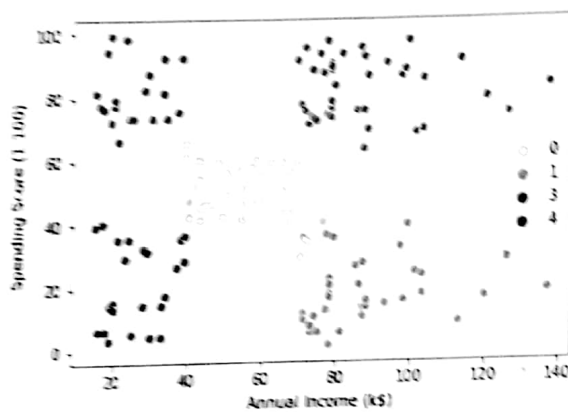
Step 13. [Predict hierarchical clusters using AgglomerativeClustering]

- Create 5 clusters using AgglomerativeClustering class. Your dendrogram will look like as below.



Step 14. [Visualize scatter plot with hue as agglomerativeclustering labels_]

- Visualize agglomerative clusters using the predicted label. Select any two features for X and Y with hue as labels_. Your scatter plot will look like below



```

print("Average age for cluster d:", d.frame.C4.Age.mean())
print("Average Annual Income for cluster d:", d.frame.C4["Annual Income (k$)"].mean())
print("Decision of the Means for Annual Income: 'Statistics',
      added(d.frame.C4["Annual Income (k$)"]))
print("No. of customers in shape of cluster d:", d.frame.C4.shape)
print("From these customers etc how", d.frame.C4.Curve.value_counts(X)
      "male and", d.frame.C4.Curve.value_counts - counts()[0], "female")

```


Step 9:-

$\Rightarrow SSE = []$

for cluster in range(1, 20):

km = KMeans(n_clusters = cluster, init = "k-means++")

km.fit(data)

SSE.append(km.inertia_)

$\Rightarrow plt.figure(figsize=(10, 6))$

$\Rightarrow plt.plot(np.arange(1, 20), SSE, "o-")$

$\Rightarrow plt.xlabel("Number of clusters")$

$\Rightarrow plt.ylabel("Inertia")$

plt.show()

Step 10:-

$\Rightarrow V = data$

$\Rightarrow pca = PCA(n_components=2)$

principal_components = pca.fit_transform(V)

PCA_components = pd.DataFrame(principal_components)

$\Rightarrow PCA_Components$

$\Rightarrow model1 = KMeans(n_clusters=5)$

$\Rightarrow model1.fit(PCA_Components)$

$\Rightarrow model1.cluster_centers$

$\Rightarrow model1.labels$

Step 11:-

$\Rightarrow sns.scatterplot(PCA_Components[0], PCA_Components[1],$

hue = model1.labels)

Step 12:-

$\Rightarrow model2 = MeanShift(bandwidth=25)$

model2.fit(PCA_components)

NOTE

⇒ cluster_centers = model2.cluster_centers.

⇒ Ans. Scatter plot (PCA-components [0], PCA-components [1], hue = model labels).

Step: 13:

⇒ model3 Agglomerative clustering (n_clusters = 3, linkage = "ward", compute_full_tree = True)

model3.fit(df)

⇒ model3.labels_

⇒ frame = df.copy()

frame['cluster'] = model3.labels_

⇒ frame.value_counts()

⇒ plt.figure(figsize=(10, 8))

plt.title('Dendrogram').

dend = sh.dendrogram(sh.linkage(frame[:, 30], method = 'ward'))

Step: 14

⇒ Ans. Scatter plot (frame['Annual Income (\$)'], frame['Spending Score (1-100)'], hue = model3.labels_)