



NAME: Hansika, P STD: X SEC: D ROLL NO: 46 SUB PODS

| S.No. | Date | Title | Page No. DM | Teacher's Sign / Remarks |
|-------|---------|---|----------------|--------------------------------|
| 1. | 19.7.25 | Numpy Commands | 10 | J P |
| 2. | 19.7.25 | Pandas Commands | 10 | J P |
| 3. | 22.7.25 | Train Dataset & Convert into DataFrame | 16 | J P |
| 4. | 22.7.25 | Data Preprocessing | 10 | J P |
| 5. | 26.7.25 | Sales prediction & Customer Segmentation using Linear Regression & K-means Clustering | 10 | J P |
| 6. | 29.7.25 | Text Preprocessing & Analysis Pipeline. | 10 | A Dr |
| 7. | 5.8.25 | NLP | 10 | A Dr |
| 8. | 9.8.25 | Exploratory Data Analysis | 10 | A Dr |
| 9. | 16.8.25 | Clustering | 10 | A Dr |
| 10. | 20.8.25 | Dashboards with Tableau using Google Sheets | 10 | A Dr |

A simple program to implements one dimensional array using numpy.

1. Import numpy

```
a = numpy.array([10, 20, 30, 40, 50])
print(a)
```

O/P: [10, 20, 30, 40, 50]

2. from numpy import *

```
a = array([10, 20, 30, 40, 50])
print(a)
```

O/P: [10, 20, 30, 40, 50]

3. from numpy import *

```
a = array([10, 20, 40.5, 50, 100])
print(a)
```

O/P: [10.0, 20.0, 40.5, 50.0, 100.0]

4. Import numpy as np

```
a = np.linspace(1, 10, 10)
print(a)
```

O/P: [1. 2. 3. 4. 5. 6. 7. 8. 9. 10.]

5. Import numpy as np

```
a = np.arange(10)
b = np.arange(5, 10)
c = np.arange(10, 1, -1)
```

Print(a)

Print(b)

Print(c)

O/P: [0, 1, 2, 3, 4, 5, 6, 7, 8, 9]

[5, 6, 7, 8, 9]

[10, 9, 8, 7, 6, 5, 4, 3, 2]

OIP: [0., 0., 0., 0., 0.]
[1., 1., 1., 1., 1.]

OIP: [15, 25, 35, 45, 55]
[5, 15, 25, 35, 45]
[50, 100, 150, 200, 250]
[2. 4. 6. 8. 10.]

OIP: [45, 5, 6, 7, 8] (Aliased)
[45, 5, 6, 7, 8] (Independent)

AP: [[2, 4, 6]
[6, 8, 10]]

OIP: 1

2

(2, 3)

6

Usually 1 (bytes)

6. Import numpy as np

$$k = np, T \text{ (max)} (5)$$

R = np · σ₀κε (5)

Prints (K)

Prints (P)

~~RULES MATHEMATICAL OPERATIONS~~

4. Import numpy as np

Keep away & cross

```
IC = np.array([10, 20, 30, 40, 50])
```

print CK+5)

print (*-5)

print (C#5)

print (K15)

ALIASING vs COPY

`k = np.array ([3, 5, 6, 7, 8])`

h = k

$$kT_0 J = 45$$

present (h)

~~h = k.copy()~~

$$K_{CO} = 100$$

~~print (n)~~

1 D Array

```
x = np.array([[[2, 4, 6], [6, 8, 10]]])
```

$\text{out}(x)$

array([5, 6, 7, 8])

$\lambda = \text{np. array}([4, 5, 6, 7, 8, 9])$

Fig. 1 (Continued)

1. *P. angustine*

prints (e.g.,
- B-snap)

paint (Ras-
pēnt)

spine (R -tensile)

11 P: 6 rows, 2 columns
1 row, 12 columns

Q1/P: Random Garbage Values
[E 0.03 T 0.03 E 0.03]
T E 1.3 T 1.3 E 1.3

014: [2 3 4 8 9 10] (8-4) entry
[2 3 4 8 9 10] (8-4) entry
[[2 3 4 5 6] (8-4) entry
[8 9 10] 3 (8-4) entry

101P:

[0 1 23 23 43 25 6 73] = N
[0 1 23 23 43 25 6 73] = N
[0 1 23 23 43 25 6 73] = N
[0 1 23 23 43 25 6 73] = N
[0 1 23 23 43 25 6 73] = N

o 1 P: $\sum_{k=1}^{\infty} (-1)^{k+1} \frac{(-1)^k}{k}$ converges as
 $\sum_{k=1}^{\infty} (-1)^k \frac{1}{k}$ converges by comparison test.

Reshaping

```
a = np.array([4, 5, 6, 7], [5, 6, 7, 8, 9, 6, 3])  
print(a = reshape(6, 12))  
print(a = reshape(1, 12))
```

Special Arrays

```
np.eye(3)  
np.empty([3, 23, 9])  
np.zeros([3, 23, 9])  
np.ones([3, 23, 9])
```

Joining Arrays

```
a = np.array([2, 3, 4])  
b = np.array([8, 9, 10])  
print(np.concatenate((a, b)))  
print(np.vstack((a, b)))  
print(np.vstack((a, b)))
```

Slicing Arrays

~~```
x = np.arange(8)
x1, x2, x3 = np.split(x, [3, 5])
print(x1, x2, x3)
a = np.arange(16).reshape((4, 4))
left, right = np.hsplit(a, 2)
top, bottom = np.vsplit(a, 2)
```~~

## Covariance

```
x = np.array([0, 1, 2])
y = np.array([2, 1, 0])
print(np.cov(x, y))
```

O/P: Cars percentage

|   |      |   |
|---|------|---|
| 0 | BMW  | 3 |
| 1 | VW   | 4 |
| 2 | Ford | 2 |

O/P: 0 10  
1 20  
2 30  
3 40  
4 50

dtype: int64

O/P: Series Representation

|   |    |
|---|----|
| 0 | 10 |
| 1 | 20 |
| 2 | 30 |
| 3 | 40 |
| 4 | 50 |

dtype: int64

Array Representation: [10 20 30 40 50]

Index Object: RangeIndex (start=0, stop=5, step=1)

O/P:

|   |   |
|---|---|
| x | 1 |
| y | 2 |
| z | 3 |

dtype: int64

O/P: a 10

b 20

c 30

d 40

e 50

dtype: int64

1. Import pandas  
`mydataset = {  
 'cars': ["BMW", "Volvo", "Ford"],  
 'passings': [3, 7, 2]}  
 myvar = pandas.DataFrame(mydataset)  
 print(myvar)`
2. Import pandas as pd  
`data = [10, 20, 30, 40, 50]  
series = pd.Series(data)  
print(series)`
3. Import pandas as pd  
`data = [10, 20, 30, 40, 50]  
series = pd.Series(data)  
print("Series Representation is", series)  
array_representation = series.values  
index_object = series.index  
print("Array Representation", array_representation)  
print("Index Object", index_object)`
4. a = [1, 7, 2]  
~~`myvar = pd.Series(a, index = ["x", "y", "z"])`~~  
~~print(myvar)~~
5. data = [10, 20, 30, 40, 50]  
~~`index_labels = ['a', 'b', 'c', 'd', 'e']`~~  
~~series = pd.Series(data, index = index\_labels)~~  
~~print(series)~~

```
OIP: 20
a 10
c 30
e 50
dtype: int64
```

```
OIP: series > 25 =
```

```
0 False
1 False
2 True
3 True
4 True
```

```
dtype: bool
```

Filtered Series:

```
2 30
3 40
4 50
```

```
dtype: int64
```

series after scalar Multiplication:

```
0 20
1 10
2 60
3 80
4 100
```

```
dtype: int64
```

```
OIP: series =
```

```
a 10
b 20
c 30
d 40
e 50
```

```
dtype: int64
```

6. `print (series['b'])`  
`print (series[['a', 'c', 'b']])`
7. `gt = series > 25`  
`Print ("Series > 25 = \n", gt)`  
~~filter~~ `series_filtered = series[series > 25]`  
`print ("Filtered Series \n", series_filtered)`  
`series_multiplied = series * 2`  
`print ("1 in series after scalar multiplication: \n",`  
~~series~~ `series_multiplied)`
8. `import numpy as np`  
`Print ("Series = \n", series)`  
`series_sqrt = np.sqrt (series)`  
`Print ("1 in series after Applying Square Root : \n",`  
~~series~~ `series_sqrt)`

Series after Applying Square Root:

|   |          |
|---|----------|
| a | 2.162278 |
| b | 4.472136 |
| c | 5.477226 |
| d | 6.324555 |
| e | 7.071068 |

dtype: float64

OIP: California 38332521  
Texas 26448193  
New York 19651127  
Florida 19552860  
Illinois 12882135

dtype: int64

OIP: True

False

OIP: a 10.0  
b NaN  
c 90.0  
d NaN  
e 50.0

dtype: float64

9. population - dict = { 'California' : 38332521,  
                          'Texas' : 26448193,  
                          'New York' : 19651127,  
                          'Florida' : 19552860,  
                          'Illinois' : 128821359}
- population = pd.Series(population - dict)
- print population
10. print ('California' in population)  
print ('Canada' in population)
11. data = [10, np.nan, 30, np.nan, 50]  
index\_labels = ['a', 'b', 'c', 'd', 'e']  
series\_with\_val = pd.Series(data, index=index\_labels)  
print series\_with\_val  
print ("In checking for NaN values %in",  
      series\_with\_val.isnull())  
print ("In checking for non-NaN values %in",  
      series\_with\_val.notnull())

Checking for Null values:

|   |       |
|---|-------|
| a | False |
| b | True  |
| c | False |
| d | True  |
| e | False |

dtype: bool

Checking for Non-Null Values:

|   |       |
|---|-------|
| a | True  |
| b | False |
| c | True  |
| d | False |
| e | True  |

dtype: bool

Df P: Population Series:

California 38332521

26448193

Texas 19651127

19552860

Florida

12882135

Illinois

dtype: int64

Area Series

California 423967

Texas 695662

New York 141297

Florida 140317

Illinois 149995

North Carolina 123456

dtype: int64

7. area-dict = { 'California': 423967, 'Texas': 695062,  
'New York': 141297, 'Florida': 170312, 'Illinois':  
149995, 'North Carolina': 1234567 }

area = pd.Series(area-dict)

Population-density = population / area

prev C1 in Population Density:  $10^4$ , population-density

population = area \* 1000

## Population Density:

California 90.413926

Florida 11.4806121

Illinois 185.893763

New York 139.076746

North Carolina NaN

Texas 38.018440

dtype: float64

## dp.. Total Population

California 38332521

Texas 26448193

New York 19552860

Illinois 12882135

dtype: float64

dp.. d

4.5

b 4.2

a -5.3

c 3.1

dtype: float64

a -5.3

b 7.2

c 3.6

d 4.5

e NaN

dtype: float64

13. male - population - dict = { 'California': 10000, 'Texas': 13000, 'New York': 95000, 'Florida': 96000, 'Illinois': 62500 }

female - population - dict = { 'California': 19222521, 'Texas': 13148193, 'New York': 10151124, 'Florida' }

male - population = pd \* scale(male - population - dict)

female - population = pd \* scale(female - population - dict)

total - population = male - population + female - population

print (total - population)

4. obj = pd \* scale ( [ 1.5, -1.2, -15.3, 3.6 ], 'order' = [ 'd', 'b', 'a', 'c' ] )

print (obj)

~~obj2 = obj \* render ( [ 'a', 'b', 'c', 'd', 'e' ] )~~

~~print (obj2)~~

EX-2013  
22.7.25

## LOAD THE TITANIC DATASET AND CONVERT INTO DATA FRAME

### AIM:

To perform basic preprocessing and exploratory data analysis (EDA) on the Titanic dataset using pandas, seaborn, and sklearn tools.

### PROCEDURE:

- \* Load Data into Data Frame.
- \* Preview Data: Use `.head()`.
- \* Check Info: `.info()` and `.isnull().sum()`
- \* Handle Missing Age: Apply `.ffill()` & `.bfill()` to 'Age'.
- \* Handle Missing Cabin: Fill with "unknown".
- \* Remove Duplicates: Use `.drop_duplicates()`.
- \* Encode 'Sex': Use Label Encoder.
- \* Scale 'Fare': Use Standard Scaler
- \* Pair Plot: For 'Pclass', 'Sex', 'Age', 'SibSp'
- \* Heatmap: For 'Pclass', 'Age', 'SibSp', 'Parch', 'Fare'.

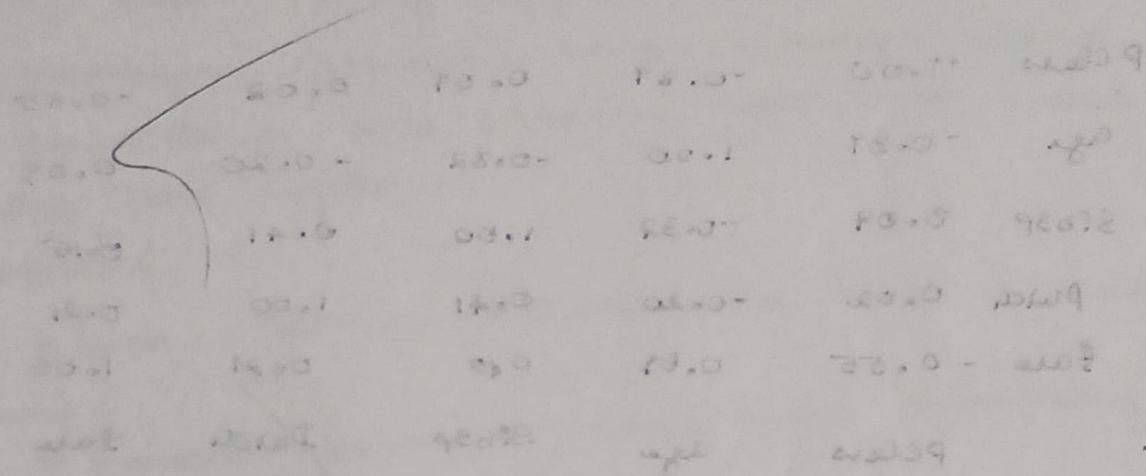
CODE:

```
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
from sklearn.preprocessing import LabelEncoder,
StandardScaler
df = sns.load_dataset('titanic')
print("First 5 rows:")
print(df.head())
print("In Dataset info:")
df.info()
print("In Missing values %n", df.isnull().sum())
df['Age'].ffill() = df['Age'].fillna()
df['Age'].bfill() = df['Age'].bfill()
df['deck'] = df['deck'].cat.add_categories('unkn')
df['deck'] = df['deck'].fillna('unknown', limit=5)
df.drop_duplicates(inplace=True)
le = LabelEncoder()
df['sex-encoded'] = le.fit_transform(df['sex'])
Scaler = StandardScaler()
df['fare-scaled'] = Scaler.fit_transform(df['fare'])
selected_features = ['Pclass', 'sex-encoded', 'age',
'sibsp']
sns.pairplot(df[selected_features].dropna())
plt.subtitle("Pair Plot", y=1.02)
plt.show()
corr_features = [pclass, 'age', 'sibsp', 'parch', 'fare']
corr_matrix = df[corr_features].corr()
```

```
plt.figure(figsize = (8, 6))
sns.heatmap(corr_matrix, annot=True, cmap=
 "coolwarm", fmt=".2f")
```

```
plt.title("Correlation Heatmap")
```

```
plt.show()
```



## RESULT:

Thus the basic preprocessing & EDA on the titanic dataset using pandas, Seaborn , sklearn tools.

Ex No: 9

## TITANIC DATASET ANALYSIS AND

Date: 28.7.25

## PREPROCESSING USING PANDAS

SIMPLE INPUTER

### AIM:

To load Titanic dataset from csv, handle missing values using simple inputer, analyse key passenger features, filter passenger based on conditions, and prepare data for model training and testing.

### PROCEDURE:

- \* Load titanic.csv handles into a data frame.
- \* Explore dataset shape & its summary statistics.
- \* Use simple inputer to fill missing Age.
- \* Fill missing cabin with "unknown"
- & embarked with mode.
- \* Visualise passenger class distribution with count plot.
- \* Filter passengers by genders, survival, class, age, family abroad and survival status.
- \* Identify top oldest survivors & zeros for passengers.
- \* Split training and testing sets.

### PROGRAM:

```
import pandas as pd
```

```
import seaborn as sns
```

```
import matplotlib.pyplot as plt.
```

### Correlation Heatmap

|        |       |       |       |       |       |
|--------|-------|-------|-------|-------|-------|
| PClass | +1.00 | -0.37 | 0.09  | 0.02  | -0.55 |
| age    | -0.37 | +1.00 | -0.32 | -0.20 | 0.09  |
| SibSp  | 0.09  | -0.32 | +1.00 | 0.41  | 0.15  |
| Parch  | 0.02  | -0.20 | 0.41  | +1.00 | 0.21  |
| fare   | -0.55 | 0.09  | 0.45  | 0.21  | +1.00 |
| PClass |       | age   | SibSp | Parch | fare  |

from sklearn import Imputer  
 from sklearn.model\_selection import train\_test\_split  
 df = sns.load\_dataset('titanic')  
 df['age'].fillna('mean')  
 df['age'].transform(lambda x: x.fillna('mean'))  
 df['deck'].cat.add\_categories(['Unknown'])  
 df['deck'].fillna('Unknown')  
 df['embarked'].fillna(df['embarked'].mode()[0])  
 sns.countplot(x='pclass', data=df).set\_title('Passenger class distribution')  
 plt.show()  
 print("Females who survived:", df[(df['sex'] == "female") & (df['survived'] == 1)].index.to\_list())  
 print("3rd class passengers under 18:", df[(df['pclass'] == 3) & (df['age'] < 18)].index.to\_list())  
 print("1st class passengers older than 40:", df[(df['pclass'] == 1) & (df['age'] > 40)].index.to\_list())  
 print("1st class passengers older than 40 who survived:", df[(df['pclass'] == 1) & (df['age'] > 40) & (df['survived'] == 1)].index.to\_list())  
 print("Male passengers who paid fare > 100:", df[(df['sex'] == "male") & (df['fare'] > 100)].index.to\_list())  
 print("Passengers embarked at 'C' and 2nd class:", df[(df['embarked'] == 'C') & (df['pclass'] == 2)].index.to\_list())  
 print("passengers with maximum 2 siblings / spouses abroad:", df[df['sibsp'] > 2].index.to\_list())

```
print("Top 5 oldest passengers who survived?")
print(df[df['survived'] == 1].head(5).values['age'])
according = False).mean() if age > 33
print("Passengers who paid zero fare:
df[df['fare'] == 0].index.to_list()")
train_df, test_df = train_test_split(df,
test_size=0.2, random_state=42)
print(f"Training set size: {len(train_df)}
Testing set size: {len(test_df)}")
```

### RESULTS

The program successfully identifies passengers with zero fare & efficiently splits the dataset into 80% training & 20% testing sets, ensuring reproducibility & readiness for machine learning tasks.

EX-N01 5 SALES PREDICTION AND CUSTOMER  
DATE: 28.7.25 SEGMENTATION USING LINEAR REGRESSION  
AND K MEANS CLUSTERING

AIMS:

To predict sales based on advertising budgets across TV, Radio & Newspaper using Linear Regression. Additionally, to identify distinct patterns in advertising spend and sales through k-means clustering.

ALGORITHMS:

- \* Read the advertising.csv file containing TV, Radio, Newspaper, budgets and sales.
- \* Prepare data: select features and target, then split into training & testing splits.
- \* Train and evaluate linear regression. Find the model on training data, predict sales on test data and calculate Mean Squared Error (MSE).
- \* Scale features & determine data into 3 groups to find patterns.
- \* Visualize features: plot actual vs predicted sales & visualize clusters to interpret model performance & segments.

## PROGRAM:

```
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.linear_model import LinearRegression
from sklearn.model_selection import train_test_split
from sklearn.metrics import mean_squared_error
from sklearn.cluster import KMeans
from sklearn.preprocessing import StandardScaler
df = pd.read_csv('advertising.csv')
print(df.head())
print(df.describe())
X = df[['TV', 'Radio', 'Newspaper']]
y = df['Sales']
X_train, X_test, y_train = train_test_split(X, y, test_size=0.2, random_state=0)
model = LinearRegression()
model.fit(X_train, y_train)
y_pred = model.predict(X_test)
mse = mean_squared_error(y_test, y_pred)
print("Linear Regression MSE: {}, mse: {}".format(mse))
plt.figure(figsize=(8, 5))
sns.scatterplot(x=X_test, y=y_pred)
plt.xlabel("Actual Sales")
plt.ylabel("Predicted Sales")
plt.title("Linear Regression: Actual vs Predicted")
plt.show()
```

scaled = StandardScaler()

scaled = scaled + fit - transform (df['W', 'Radio',  
'Newspaper'])

K-means = KMeans (n\_clusters = 3, random\_state =

df[['cluster']] = KMeans().fit(scaled).predict(scaled)

plt.figure (figsize = (8, 6))

sns.scatterplot (data = df, x = 'TV', y = 'Sales',  
hue = 'cluster', palette = "Set2")

plt.title ("K-Means (Clustered TV budget  
vs Sales")

plt.show()

## RESULT:

The Linear Regression model achieved

a Mean Squared Error (MSE) of 4052,  
showing a good fit between actual &

Predicted Sales. K-Means clustering  
grouped the data into 3 distinct  
clusters, revealing clear patterns in TV  
Advertising budgets & Sales.

EX NO: 6  
DATE: 29/07/25 TEXT PREPROCESSING AND  
ANALYTICS PIPELINE

AIM:

To clean & preprocess Amazon review texts for analysis, then identify the most frequent words.

ALGORITHM:

\* Load Dataset

Import Amazon reviews from amazon\_reviews.csv.

\* Clean Text:

Convert text to lowercase  
remove punctuation & special  
characters then tokenize using spacy.

\* Remove stopwords.

Filter out common  
stopwords & punctuable tokens.

\* Apply Cleaning:

Process the entire  
review text column to generate cleaned  
tokens.

\* Analyze frequency:

Flatten all tokens  
& count word frequencies using Counter.

Output:

we got this app for my band who  
is an OTR).

1. I'm a professional OTR truck  
driver, and I boll....
2. well, what can I say. I've had  
this unit in n....
3. Not going to write a long  
review, even megit....

Name: Review Text, dtype: Object

Review Text

0 we got this app for my  
band who is an OTR)

1. I'm a professional OTR truck  
driver & I boll...
2. well what can I say.
3. Not going to write a long  
review.
4. I've had mine for a  
year & here what we go.

## Cleaned tokens

- 0 [got, gps husband, off, road.]  
1 [me, professional, off, truck, driver]  
2 [ve, mit, truck, days, Poor.]  
3 [going, route, long, review]  
4 [ve, year, got, route, true]

## Top 15 frequent words in

### Amazon Reviews:

[('c', 3203, ('book', 1447), ('8', 962),  
('books', 005), ('Kindle', 561),  
('seen', 473), ('like', 452) ('lead', 13),  
(great, 122), ('use', 420)]

## PROGRAM:

```
Import Pandas as pd
Import re
Import spacy
nlp = spacy.load("en-core-web-sm")
df = pd.read_csv("amazon_reviews.csv")
print(df[["Review Text"]].head())
def clean_text_spacy(text):
 if pd.isnull(text):
 return []

 text = text.lower()
 text = re.sub(r'\W+', " ", text)
 text = text.encode("ascii", "ignore").decode("ascii")

 doc = nlp(text)
 tokens = [token for token in doc if
 not token.is_stop & not token.is_punct]

 return tokens

df['cleaned_tokens'] = df['Review Text'].
apply(clean_text_spacy)
print(df[['Review Text', 'cleaned_tokens']].
head(5))
```

## RESULT:

Thus cleaned tokens contain meaningful words without noise, and the top 5 frequent words reflect key terms from the review.

Ex: 5

## NLP

DATE: 5.9.25

AIM:

To analyze text using Part of Speech (POS) tagging (with Spacy) and for grammatical structure identification.

ALGORITHMS

- \* Load language model & load en-core in Spacy for POS tagging.
- \* Text Preprocessing Pass the input text to Spacy & extract POS tags.
- \* Prepare Corpus Collect sample documents & append the user query.
- \* Vectorization Convert documents into TF-IDF vectors.
- \* Similarity & Ranking - Compute cosine similarity between query & documents, then rank by score.

Q1 p:

## POS Tagging:

AI → PROPN

driven → VERB

Platforms → NOUN

Personalise → VERB

Learning → VERB

Paths → NOUN

and → COUNT

help → VERB

Students → NOUN

grasp → VERB

Concepts → NOUN

faster → ADV

→ PUNCT

## Top relevant documents:

Score 0.016 → AI helps automate  
grading & administrative tasks  
in schools.

Score 0.010 → Intelligent tutoring  
systems adapt to each student's  
learning style.

## PROGRAM

```
pip install spacy
```

```
import spacy
```

```
nlp = spacy.load("en-core-web-sm")
```

```
text = "AI - driven platforms personalize
learning paths + help students
```

```
group concepts together."
```

```
doc = nlp(text)
```

```
for token in doc:
```

```
print(f" {token.text} : {token.pos_}")
```

```
from sklearn.metrics.pairwise import
cosine_similarity
```

```
documents = [
```

" AI tools analyze student  
Performance & provide real time  
feedback " Intelligent tutoring system  
adapt to each student's learning  
style".

query = " How does AI support  
students in learning?"

```
Corpus = documents + [query]
```

vectorizer = TfidfVectorizer()

+ tfidft-matrix = vectorizer . fit\_transform  
(corpus)

similarities = cosine-similarity (+ tfidft-matrix  
T-1T, + tfidft-matrix T^T - 1T) . flattened  
ranked-docs = sorted (zip(similarities,  
documents), reverse=True)

print("Top relevant documents: ")

for score, doc in ranked\_docs:

print(f"score: {score}, doc: {doc}")

RESULT:

Thus, we have analyzed text  
using POS & for grammatical  
structure identification.

Ex no: 6

## Exploratory Data Analysis

DATE: 9.9.25

with Python

AIM:

To explore and analyze the Netflix dataset using Python in order to uncover insights about content types, release trends & countries of production & genre distributions.

ALGORITHM:

\* Load datasets & scale features as needed.

\* Apply Kmeans clustering & Mall customers data & use the ~~error~~ method to find the optimal clusters.

\* Perform multiple KMeans clusterings and see which dataset works varying cluster counts.

\* Bind a similarity matrix from base clusterings & apply spectral clustering for ensemble levels.

## OUTPUT:

<Class "pandas.core.frame.DataFrame">

RangeIndex: 8804 entries × 12 columns

Data columns (total 12 columns):

| # Column |            | Non-Null Count | D type          |
|----------|------------|----------------|-----------------|
| 0        | show_id    | 8804           | non-null object |
| 1        | type       | 8804           | non-null object |
| 2        | title      | 8804           | non-null object |
| 3        | director   | 6143           | non-null object |
| 4        | cast       | 7982           | non-null object |
| 5        | country    | 7982           | non-null object |
| 6        | date_added | 8794           | non-null object |

dtypes: float64(1), object(11)

memory usage: 82508+ KB

None

|    | show_id  | type     | title   | director | cast  | country |
|----|----------|----------|---------|----------|-------|---------|
| 1  | Ame      | Khosla   | Gail    | Thakur   | South |         |
|    | Asmada   | Nguma    | Mabaleu |          |       | Afghan  |
| 2  | Sani     | Tracy    | Samuel  | Nabi     | Nan   |         |
|    | Bonefile | Gotoas   | Tony    |          |       |         |
| 3. | Mayur    | Jitendra | Rajiv   | Alex K.  | Dulq  |         |
|    | More     | Kumar    | Rey     |          |       |         |

## PROGRAM:

```
import pandas as pd
import matplotlib.pyplot as plt
from sklearn.cluster import KMeans
from sklearn.preprocessing import StandardScaler
df = pd.read_csv("netflix_titles.csv")
print(df.info())
print(df.describe())
print(f"Number of unique countries: {df['Country'].nunique()}")
print(df['type'].value_counts())
print(df['release_year'].value_counts().head())
print(df.groupby('country')['type'].size().sort_values(ascending=False).head(10))
df['date_added'] = pd.to_datetime(df['date_added'], format="mixed", errors="coer")
df.set_index('date_added', inplace=True)
monthly_content = df.resample('M').size()
```

(date added release year writing initials)

|   |                    |      |       |           |
|---|--------------------|------|-------|-----------|
| 0 | September 25, 2021 | 2020 | PL-13 | Young     |
| 1 | September 24, 2021 | 2021 | TV-HA | 2 seasons |
| 2 | September 24, 2021 | 2021 | TV-HA | 1 season  |
| 3 | September 24, 2021 | 2021 | TV-HA | 1 season  |
| 4 | September 24, 2021 | 2021 | TV-HA | 2 seasons |

(Listed in)

0 Documentaries

- 1 International TV Shows, TV Dramas
- 2 Crime TV shows, TV Net
- 3 Documentaries, Reality TV
- 4 International TV Shows, Romantic TV Shows

description

- 0 As we follow near the end of  
    her life, film.
- 1 After crossing paths at a party,  
    a Cape Town ...
- 2 To protect her family from a  
    powerful drug lord ...
- 3 Feuds, flirtations & toilet talk go  
    down and ...
- 4 Be a city of wealthy citizens  
    busing to train? ...

```
plt.title ("Netflix Content Added Over Time")
plt.xlabel ("Date")
plt.ylabel ("Number of Titles Added")
plt.grid (True)
plt.show()

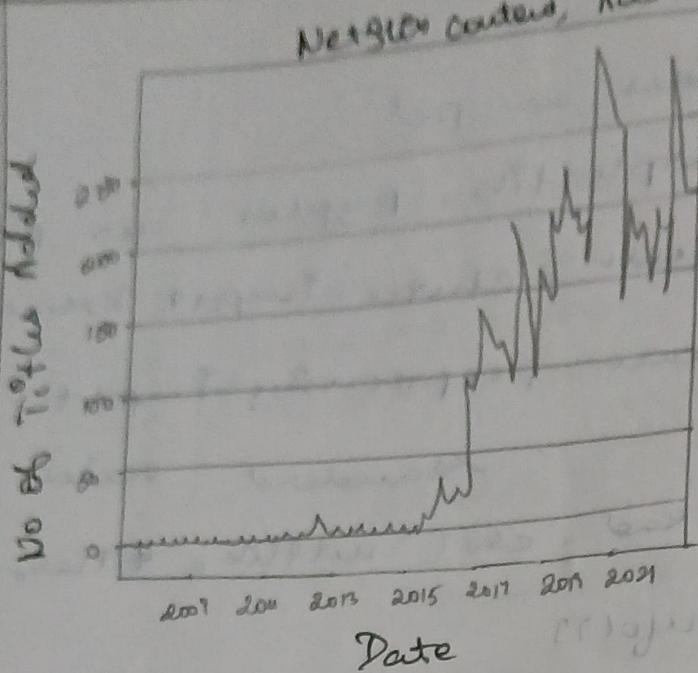
plt.figure (figsize = (10,5))
sns.countplot (df['release_year'], bins=30,
 rot=90, palette='Set2',
 edgecolor='black')
plt.title ("Distribution of Release Years")
plt.xlabel ("Release Year")
plt.ylabel ("Number of Titles")
plt.show()

sns.countplot (data=df[df['type'] == "Movie"], palette='Set2')
plt.title ("Count of Movies & TV Shows")
plt.show()

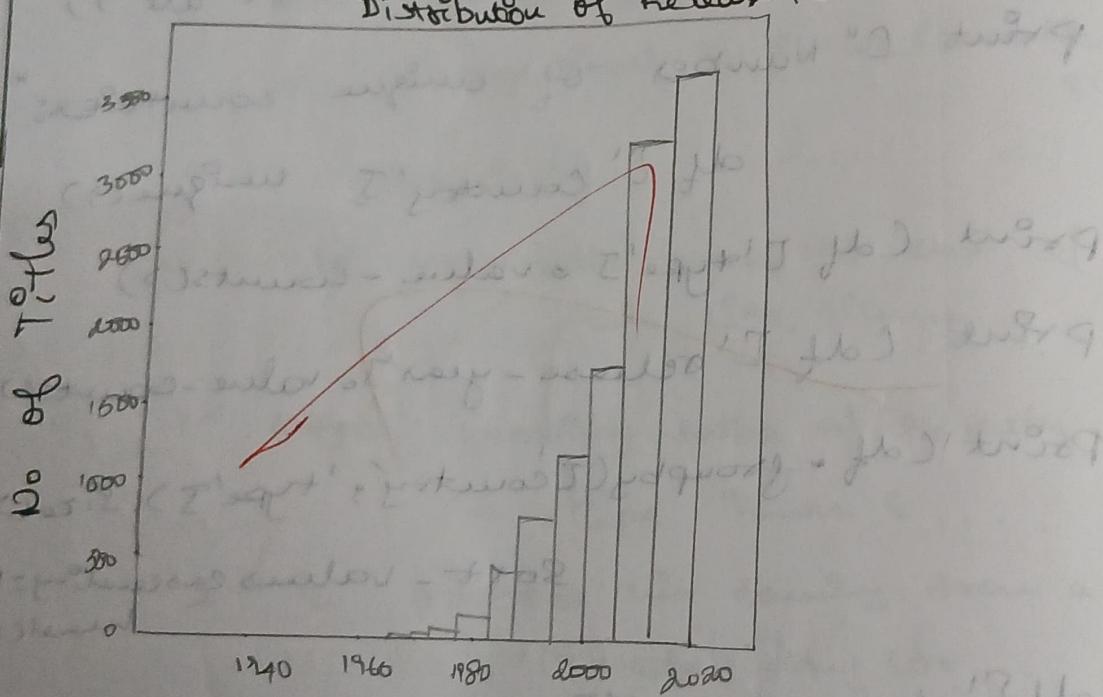
top_countries = df['country'].value_counts().head(10)
top_countries.plot (kind='bar', color='skyblue')
plt.title ("Top 10 countries by Number of Titles")
plt.ylabel ("Count")
plt.xticks (rotation=45)
plt.show()

genres = df['listed_in'].str.split (',', expand=True)
```

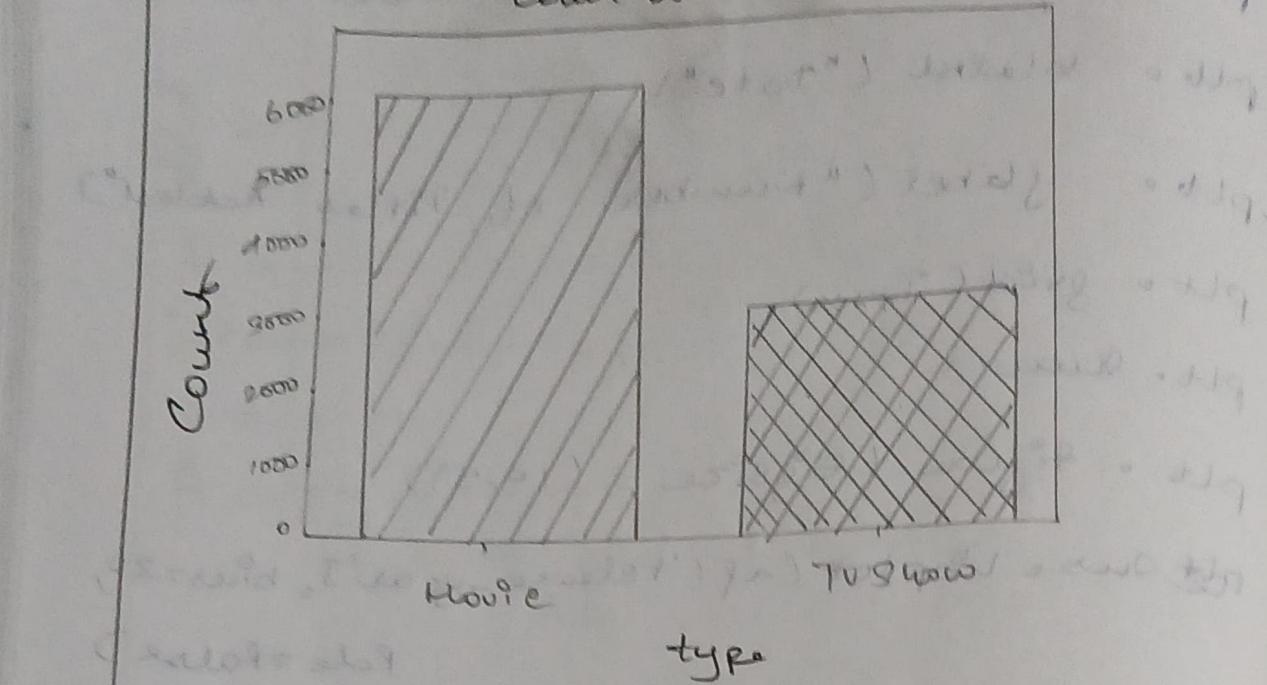
Net titles added over time



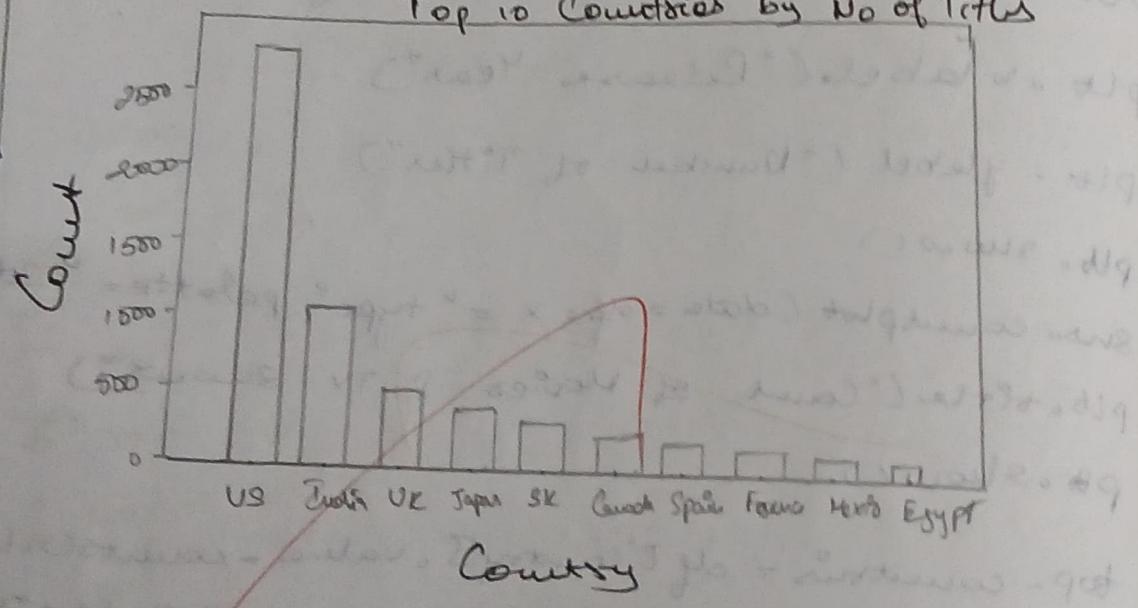
Distribution of Release Years



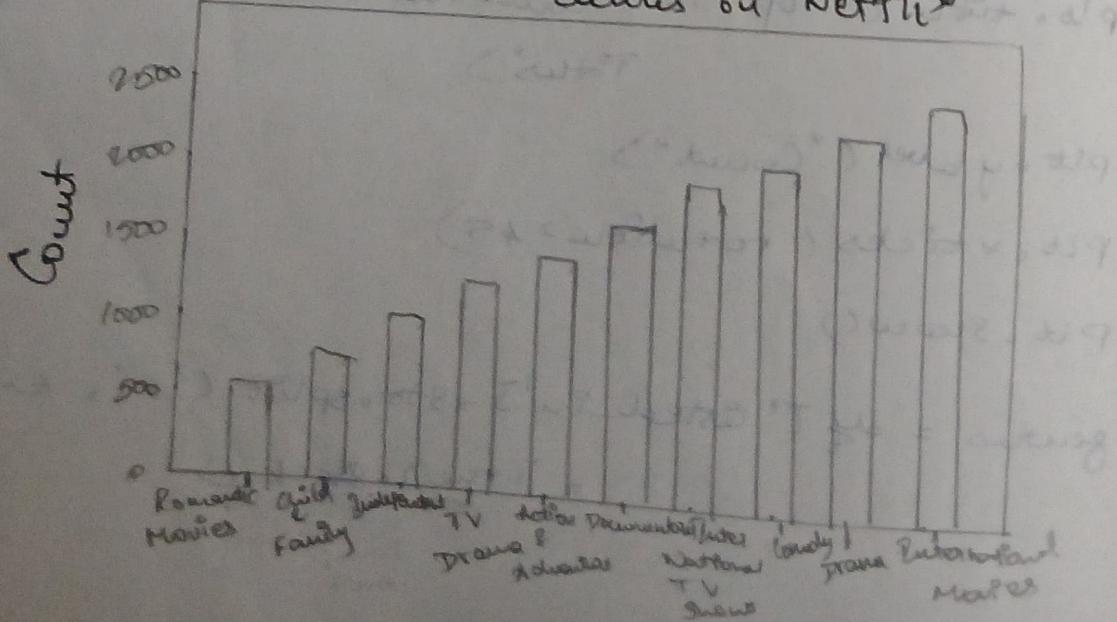
### Count of Movies & TV Shows



### Top 10 Countries by No of Titles



### TOP 10 Genres on Netflix



```
top_genres = genres.value_counts()[:10]
top_genres.plot(kind="bar", color="coral")
plt.title("Top 10 Genres on Netflix")
plt.xlabel("Count")
plt.yticks().invert_yaxis()
plt.show()
```

RESULT:

~~This~~ ~~about~~ ~~is~~ ~~to~~ ~~the~~ ~~dataset~~ ~~implementation~~ & analyze the  
dataset using Python is.  
Implementation & executed successfully.

EX NO: 4

DATE: 16.9.25

## CLUSTERING

AIM:

\* Mall Customers: To group customers using K Means clustering & find the optimal number of clusters with the Elbow Method.

\* Wine Dataset: To apply CSPA Ensemble clustering with spectral clustering & visualise clusters using PCA.

ALGORITHM:

\* MALL CUSTOMERS (K Means):

\* Load dataset & select features.

\* Applying K Means Clustering.

\* Use Elbow Method to determine best K.

\* Visualise clusters with Scatter plot.

\* WINE DATASET:

\* Load & Standardise data.

- \* Run k-Means with different  $k$  values.
- \* Build similarity matrix.
- \* Apply spectral clustering.
- \* Evaluate with silhouette score.
- \* Reduce dimensions with PCA and plot clusters.

CODE:

```

import pandas as pd
import matplotlib.pyplot as plt
from sklearn.cluster import KMeans
from sklearn.preprocessing import StandardScaler
import seaborn as sns
df = pd.read_csv("Imall - customers.csv")
kmeans = KMeans(n_clusters=5, random_state=42)
df['cluster'] = kmeans.fit_predict()
distortions = []
for i in range(1, 11):
 km = KMeans(n_clusters=i)
 km.fit(df[['Annual Income (k$)', 'Spending Score (1-100)']])
 distortions.append(km.inertia_)

distortions = np.array(distortions)

```

```
plt = plt + plot(orange[1, 11], distortions, marker="o")
plt = title("Elbow Method")
plt = xlabel("Number of clusters")
plt = ylabel("Distortion")
plt = show()
sns = scatterplot(data=df, x = "Annual Income($)",
y = "Spending Score (1-100)", hue="Cluster",
palette="Set2")
```

```
from sklearn.datasets import load_wine
from sklearn.preprocessing import StandardScaler
from sklearn.cluster import KMeans
from sklearn.decomposition import PCA
from sklearn.metrics import silhouette_score
import matplotlib.pyplot as plt
import numpy as np
import pandas as pd
```

wine = load\_wine()

x = pd.DataFrame(wine.data, columns=wine.feature\_names)

x\_scaled = StandardScaler().fit\_transform(x)

base\_clusterings = [ ]

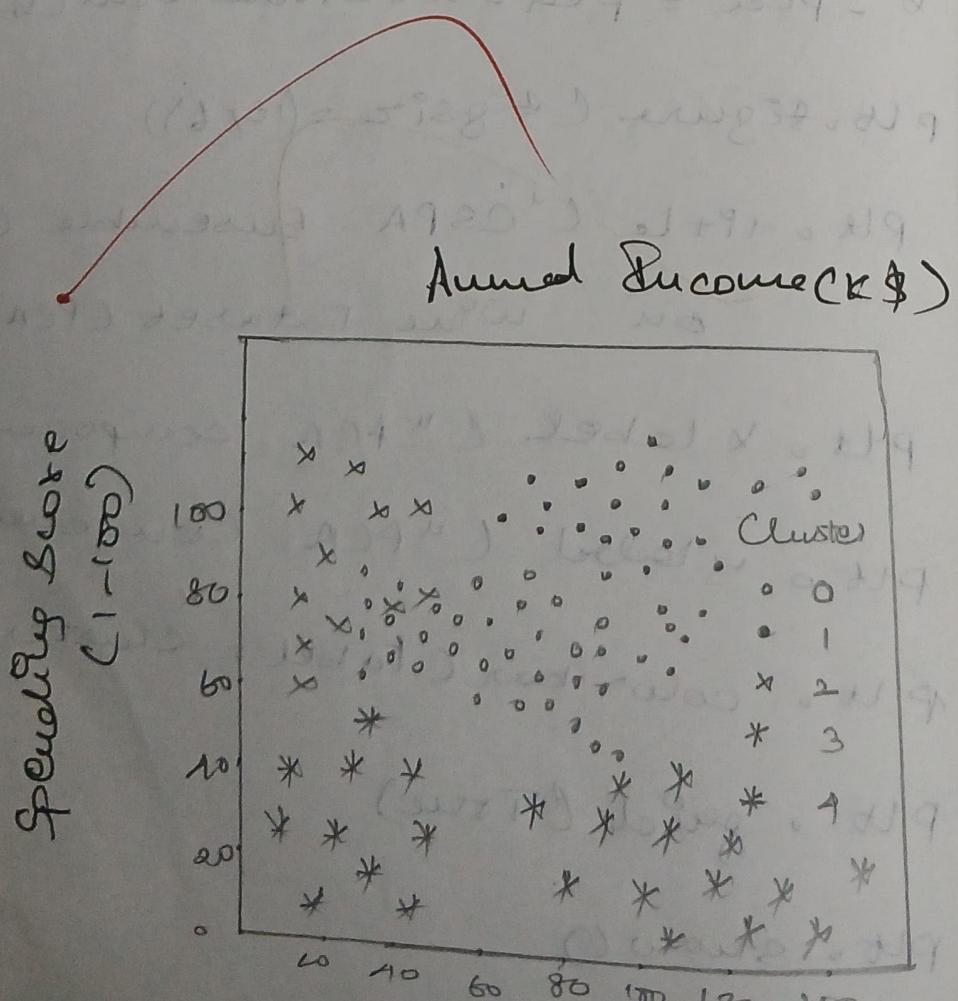
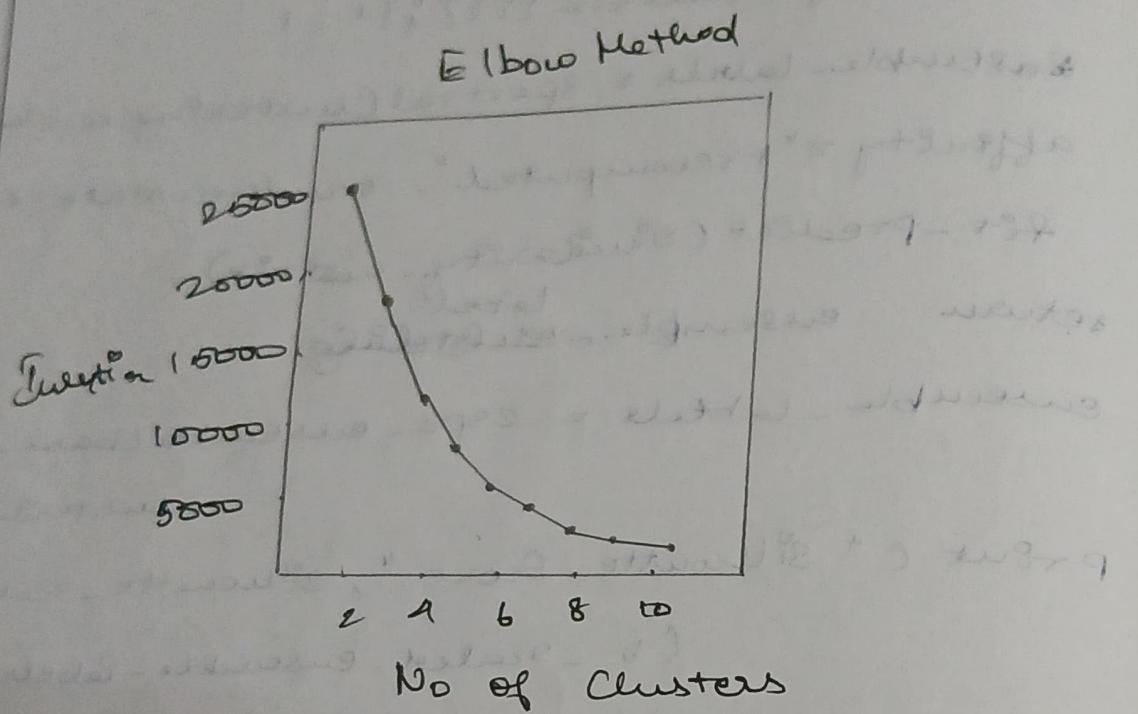
for k in [3, 4, 5]:

km = KMeans(n\_clusters=k, random\_state=42)  
base\_clusterings.append(km.fit\_predict(x\_scaled))

```

for clustering in clusterings:
 for p in range(n - samples):
 for q in range(n - samples):
 if clustering[i][j] == - clustering[j][i]:
 similarity_matrix[i][j] += 1
ensemble_labels = Spectral Clustering(n_clusters=3,
affinity="precomputed", random_state=42),
fit_predict(similarity_matrix)
return ensemble_labels
ensemble_labels = cspa.ensemble(base -
clusterings)
print("Silhouette Score:", silhouette_score
(b_scaled, ensemble_labels))
PCA = PCA(n_components=2)
b_pca = pca.fit_transform(b_scaled)
plt.figure(figsize=(10, 6))
plt.title("CSPA Ensemble Clustering
on Wine Dataset (PCA-reduced")
plt.xlabel("PCA component 1")
plt.ylabel("PCA component 2")
plt.colorbar(label="Cluster label")
plt.gcf().ax[0].set_label("True")
plt.show()

```



## RESULT:

- \* In the Mall customers dataset, customers are successfully grouped into distinct clusters showing spending behaviour patterns.
- \* In the wine dataset, LSPA ensemble clustering improved clustering quality confirmed by a good silhouette score & clear cluster separation in PCA visualisation.

Ex no: 8

## Interactive Dashboards

With Tableau using  
Google Sheets

AIM:

To create an interactive Tableau dashboard that integrates data from Google Sheets & visualizes it using bar charts, line plots and scatter plots, allowing dynamic exploration of sales or demographic data.

MATERIALS REQUIRED:

- \* Tableau Desktop or Tableau Public
- \* Google account with Google Sheets.
- \* Sample Sales or demographic dataset on Google Sheets.
- \* Internet connection.

PROCEDURE:

- \* Prepare the data.
- \* Connect Tableau to Google Sheets.
- \* Create individual visualizations.
  - ⇒ Bar charts (Sales by product or Region)

- ⇒ Line Plot (Sales over time)
  - ⇒ Scatter plot (Customer Demographics or profit vs sales)
- \* Build the Dashboard.
- \* Publish or Share.

~~RESULT~~ Created an Interactive Tableau dashboard linked to Google Sheets, allowing users to explore sales trends, compare product performance & analyze demographics.