

Practical No: 01

* Aim :-

Familiarization with NumPy, Pandas and matplotlib by loading dataset in Python.

* Theory :-

1] NumPy Library :-

- NumPy is a general-purpose array-processing package. It provides a high-performance multidimensional array object, and tools for working with arrays. It is the fundamental package for scientific computing with python.
- One of the most popular open-source python packages. NumPy focuses on scientific and mathematical computation. It makes it easy to work with large matrices & multidimensional data.

2] Pandas Library :-

- It is a powerful and versatile library that simplifies the tasks of data manipulation in python.
- Pandas is an open-source library authorized under BSD. This well-known library utilized in the field of data science.

- Pandas is well suited for working with tabular data, such as spreadsheets or SQL tables.

3) Matplotlib library:-

- This library is responsible for the plotting of numerical data. It is utilized in data analysis for this reason.
- Matplotlib is a plotting library for python that provides a comprehensive set of tools that provides a tool for creating high quality 2D & 3D plots.

4) Scipy library:-

- Scipy is a library for scientific computing in python.
- This program demonstrates linear regression using 'unregress()' function.
- It provides a function for scientific and engineering application.
- It is for scientific computing, information processing, and high-level computing are primary uses for this open source library.



* Programs:-

Code 1 :- Numpy:-

```
import numpy as np.
```

```
arr = np.array ([20, 30, 40])
print ("Array with Rank 1:\n", arr)
```

```
arr = np.array ([[10, 20, 30],
                 [40, 50, 60]])
print ("Array with Rank 2:\n", arr)
```

```
arr = np.array ((10, 30, 20))
print ("An Array created using"
      " passed tuple :\n", arr)
```

Output:-

Array with Rank 1:

[20 30 40]

Array with Rank 2:

[[10 20 30]

[40 50 60]]

Array Created using passed tuple:

(10 30 20)

Code 2 :- panda

```
import pandas as pd
```

```
import numpy as np.
```

```
ser = pd.Series()
```



```

print ("Pandas Series: ", ser)
data = np.array(['g', 'e', 'e', 'K', 'S'])
ser = pd.Series(data)
print ("Pandas Series: \n", ser)

Output:-
Pandas Series: Series([], dtype: float64)
Pandas Series:
0    g
1    e
2    e
3    K
4    S
dtype: object.

```

Code 3: matplotlib.

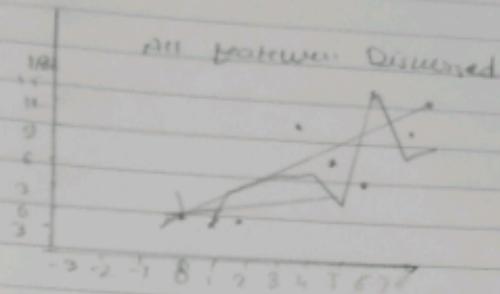
```

import matplotlib.pyplot as plt
a = [1, 2, 3, 4, 5]
b = [0, 0.6, 0.2, 15, 10, 8, 16, 21]
plt.plot(a)
plt.plot(b, "or")
plt.plot(list(range(0, 22, 3)))
plt.xlabel('Day →')
plt.ylabel('Temp →')
c = [4, 2, 6, 8, 3, 20, 13, 15]
plt.plot(c, label='14th Rep')
ax = plt.gca()
ax.legend(['1st Rep', '2nd Rep', '3rd Rep'])

```

```
plt.annotate('Temperature v/s Day',)  
plt.title('All features Discussed')  
plt.show()
```

O/P:-



code 4:- scipy

```
from scipy import sparse  
A = sparse lil_matrix((1000,1000))  
print(A).  
A[0,:100] = np.random.rand(100)  
A[1,100:200] = A[0,:100]  
A.setdiag(np.random.rand(1000))  
print A
```

O/P:-

(1, 120)	0.6315
(1, 131)	0.44
(1, 132)	0.57
(1, 133)	0.33
(1, 134)	0.22
(1, 135)	0.79

Result :- Hence, In python, the implementation
of Numpy, panda, matplotlib by Loading
dataset- are executed successfully.

Practical No. 2



* Aim :- Implement and demonstrate find-s Algorithm Read training data from a CSV file.

* Theory :-

- find-s Algorithm :-

- It is simple machine learning algorithm used for finding most specific hypothesis that fits training data.
- It works with binary classification problems where each instance is either a positive or negative example.

- Components :-

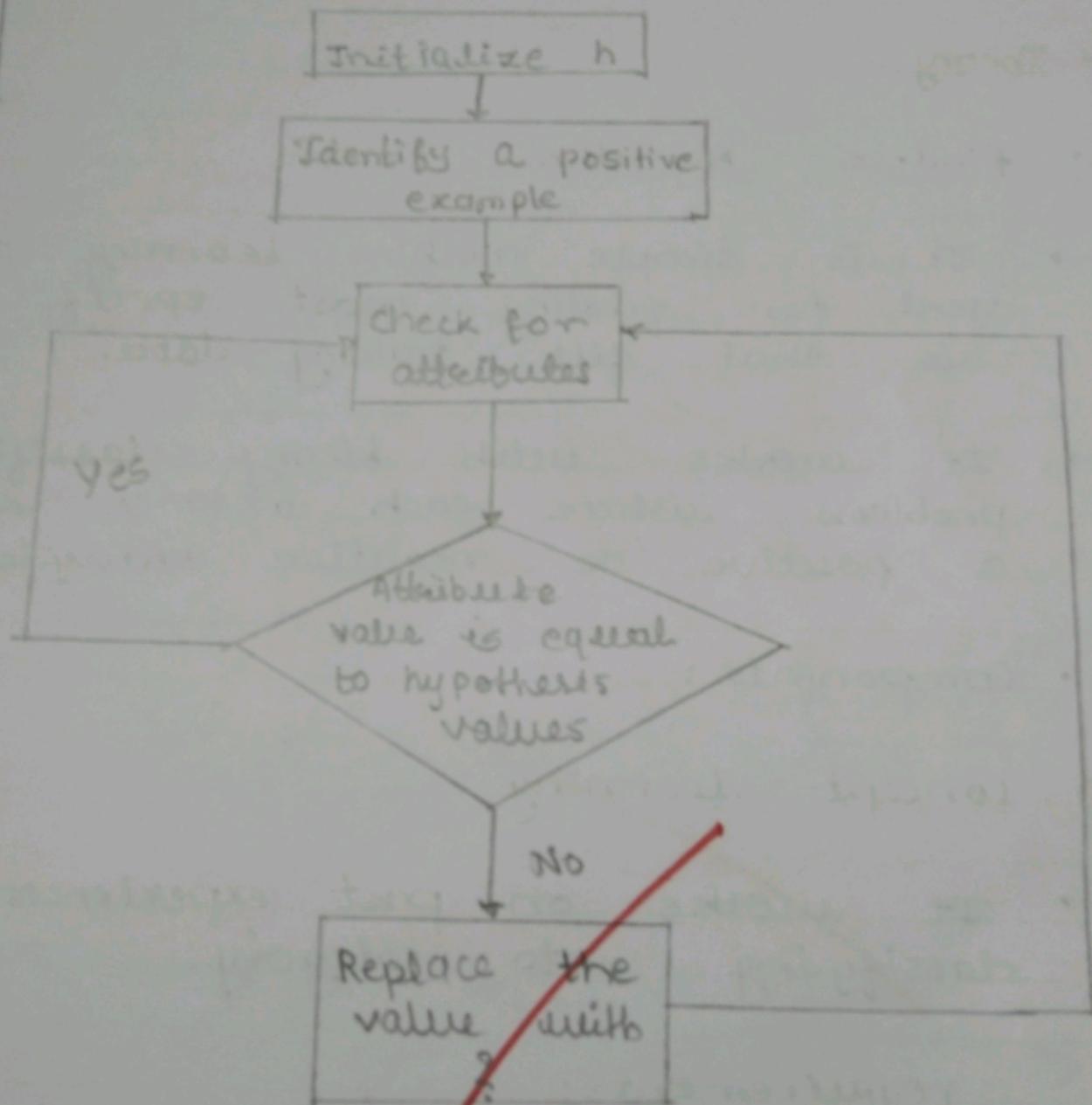
① Concept Learning :

- It works on past experiences and classifying into category.

requirements:-

- ① Training data
- ② Target concept.
- ③ Actual data objects

* Flowchart of finds Algorithm no:





② General Hypothesis :

- It is an explanation for something.
It states relationship betⁿ or among concepts.

③ Specific Hypothesis :

- If It fills in all the important details about variables given in general hypothesis

$$S = \{ ' \phi ', ' \phi ', ' \phi ', ' \phi ', \dots, ' \phi ' \}$$

* Algorithm :-

① Initialize Hypothesis :

- start with most specific hypothesis,
e.g., $h = (? , ? , \dots, ?)$, where ? represents any h-value.

② Iterate through examples :

- for each positive example (x, y) :
 - update 'h' to match positive example by replacing '?' in 'h' with corresponding attribute value from 'x'.
- Negative examples are ignored.



⑤ Output: - The final hypotheses 'h' represents most specific hypothesis consistent with all positive examples.

* Outputs Data frame:

	Name	branch	result	mood	beh.	color	med.	abse
0	rohit	AI	pass	happy	good	fit	white	yes
1	Komal	CS	pass	happy	good	fit	black	yes
2	Palak	EE	fail	unhappy	bad	unfit	white	no
3	dhan.	ETC	fails	unhappy	bad	unfit	black	yes
4	Vinay	CS	fail	happy	good	fit	white	no
5	aditya	CE	fail	unhappy	bad	unfit	black	yes
6	Harshal	AI	pass	unhappy	bad	fit	white	no
7	Piyush	EE	fail	happy	good	unfit	white	yes
8	Tushar	ETC	pass	happy	bad	unfit	black	no
9	aditya	AI	fail	unhappy	good	fit	white	yes
10	Nikhil	ETC	pass	happy	bad	fit	black	no

Attribute:-

```
[['Komal', 'CS', 'Pass', 'Happy', 'Good', 'White', 'Fit'],
 ['Palak', 'CI', 'Fail', 'Unhappy', 'Bad', 'Black', 'Unfit'],
 ['Aditya', 'EE', 'Pass', 'Unhappy', 'Good', 'White', 'Unfit'],
 ['Dhananjay', 'ETC', 'Fail', 'Happy', 'Bad', 'Black', 'Fit'],
 ['Vinay', 'CS', 'Pass', 'Unhappy', 'Bad', 'Black', 'Fit'],
 ['Aditya', 'CI', 'Fail', 'Unhappy', 'Good', 'White', 'Unfit'],
 ['Harshal', 'AI', 'Fail', 'Happy', 'Bad', 'White', 'Unfit'],
 ['Piyush', 'EE', 'Pass', 'Unhappy', 'Good', 'Black', 'Unfit'],
 ['Nikhil', 'ETC', 'Pass', 'Happy', 'Bad', 'White', 'Fit']]
```

The target is ['Yes', 'No', 'Yes', 'No', 'Yes', 'No', 'Yes', 'No']

* Result:- Hence, the find-s algorithm implemented and demonstrated successfully

Practical no: 08



* Aim:- Demonstrate various data pre-processing techniques and plot density plot, box plot using suitable datasets.

* Theory:-

① Creating database:

```
import pandas as pd
stu = ["abc", "xyz", "rst", "lmn"]
ye = [1, 1, 3, 3]
dict = {'student': stu, "year": ye}
```

```
df = pd.DataFrame(dict)
print(df)
```

	student	year
0	abc	1
1	xyz	1
2	rst	3
3	lmn	3

② Connecting into CSV files :-

→ To connect a database frames into a CSV file using python's pandas library, we can use to_csv method.

Command: df.to_csv ('Pandas.csv')



③ Loading a database :-

→ Loading a database into a pandas data-frame typically involves using a data-base connector and the pandas library.

```
df = pd.read_csv('Pandas CSV')
print(df.head())
```

	unnamed:	0	student	year
0	0	abc	1	
1	1	xyz	1	
2	2	rst	3	
3	3	lmn	3	

④ Getting information of database in detail:-

→ It includes using a single command in python i.e. 'info'.

df.info()

```
<class 'pandas.core.frames.DataFrame'>
RangeIndex: 4 entities, 0 to 3
Data columns (Total 3 columns)
```

*	column	not-null count	Dtype
0	unnamed:0	4-not null	int 64
1	student	4-not null	object
2	year	4-not null	int 64.



dtypes : int 64(2), object(1)
memory usage: 224.04 bytes

⑥ To get summary of the database :-

→ Knowing the detailed description of the database including Using a single command "describe()".

dt.describe(),

	Unamed: 0	year
count	4.00000	4.00000
mean	1.50000	2.00000
std	1.290994	1.154701
min	0.000000	1.00000
25%	0.750000	1.00000
50%	1.50000	2.00000
75%	2.250000	3.00000
max	3.00000	3.00000

⑥ To know about null values :-

→ It present in the database includes using the command is "null().sum()".

command : df.isnull().sum()



df. is null(), sum()

unnamed : 0	0
student	0
year	0
d type : int 64	

- ⑦ To Know about training
2 testing data

→ Segregating data gives specifically the data present in the training part and the amount of data present in testing part following is the method.

$x_{\text{train}}, x_{\text{test}}, y_{\text{train}}, y_{\text{test}} = \text{train-test-}$
split(x, y, test_size = 0.05, random_state = 0)

$$\text{ratio} = 0.75$$

$$\text{total_rows} = \text{df. shape}[0]$$

$$\text{train_size} = \text{int}(\text{total_rows} * \text{ratio})$$

$$\text{train} = \text{df} (0 : \text{train_size})$$

$$\text{test} = \text{df} [\text{train_size} :]$$

~~print(train)~~

unnamed : 0	student	year
0	abc	1
1	xyz	1

Date :



2

2

rst

3

print (test)

3

unnamed : 0

3

student

unn

year
3

* Result :- Hence, various data pre-processing techniques & plot density plot in python is done successfully.

~~Result :-~~

Practical No. 4.

* Aim :- Demonstrate Dimensionality reduction using principle component analysis (PCA) method.

* Theory :

* Dimensionality Reduction :

- As a number of dimensions increased, the volume of the feature space grows exponentially. leading to sparse data, which makes it difficult for models to generalize.
- Dimensionality reduction is the transformation of data from a high-dimensional space into a low-dimensional space.

* Principle Component Analysis (PCA) :

- principle component analysis is a dimensionality reduction method that is often used to reduce the dimensionality of large dataset.

* Step-by-step explanation of PCA:

~~Step 1:- Standardization.~~

The aim of this step is to standardize the range of continuous initial variables so that each one of them contributes equally to analysis.

$$\text{Mathematically, } z = \frac{\text{value} - \text{mean}}{\text{s.d.}}$$

Step 2:- covariance Matrix computation.

The aim of this step is to understand how the variables of the input data set are very from mean with respect to each other.

$$cov(x, y) = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

Step 3:- calculate Eigenvalues & eigenvectors once you have the covariance matrix, the next step is to calculate the eigenvalues & eigenvectors. Eigenvectors represent directions (principle components) of the new feature space.

Step 4:- Sort eigenvalues & eigenvectors:

- Sort the eigenvalues in descending order
- and rearrange the corresponding eigenvectors. the eigenvectors with highest eigenvalue are principle components that capture the most variance in the data.

Step 5:- Select principle components.

- Decide the numbers of principle components you want to keep. Usually, the number of components is chosen based on the cumulative explained variance. which indicate how much of total variance is retained by selected components.



- * Cumulative Explained Variance: The proportion of the total variance captured by the first k principal components.

Step 5:- Transform the data:

- finally transform the original dataset into the new space defined by the selected principal components. This step involves multiplying the original data by the matrix of selected eigenvectors.

$$Z = X \cdot W$$

where,

Z is the transformed dataset

X is the original standardized data.

W is matrix of selected eigenvectors.

Step 7:- Interpret the Result:

- The resulting transformed dataset consists of new features that are uncorrelated & ordered by the amount of variance they explain.

* Resulting Program has been executed successfully.

Practical No. 5.



* Aims- Implement the multiple linear regressions algorithm in order to fit data points.

* Theory:

* Multiple linear regression:

In this, previous topic, simple linear regression, where a single independent / Predictor (x) variable is used to model the response variable is used to model the response variable (y). But there are may be varieties cases in which the response variable is affected by more than one predictor variable.

We assume that y is linearly dependent on the factors according to:

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k + \varepsilon$$

- The variable y is dependent or predicted
- The slope of y depends on the y -intercept. - that is, when x_1 & x_2 are both zero, y will be β_0 .

* Multiple Regression working process:

1. Data Preparation: Gather and preprocess the data, including handling missing values,



encoding categories variables, and scaling the features if necessary.

2. Splitting the Dataset: Divide the dataset into a training set & a test set to evaluate the model's performance on unseen data.
3. fitting the model: Use the training set to fit the multiple linear regression model by finding the best coefficients $B_0, B_1, B_2, \dots, B_n$ that minimize the sum of squared residuals (errors)
4. Making Predictions: Use the fitted model to predict the dependent variable on the test set.
5. Evaluating the model: Assess the model's performance using metrics such as R-squared, Mean Absolute Error (MAE), Mean squared Error (MSE), and Root Mean squared error (RMSE).

* Result: Hence, a implementation of multiple linear regression algorithm is successfully.

~~Q&A~~

Practical NO. 6.



* Aim :- Apply K-nearest Neighbour algorithm to classify the data set.

* Theory :- K-Nearest Neighbour (KNN) Algorithm for machine learning.

- K-Nearest Neighbour is one of the simplest Machine learning algorithms based on supervised learning technique.
- K-NN algorithm assumes the similarity between the new case / data into the category that is most similar to the available categories.

* How does K-NN work :

- The K-NN working can be explained on the basis of the below algorithm.
- step 1 : Select the number K of the neighbours.
- step - 3 : Take the K nearest neighbour
as per the calculated Euclidean distance.
- ~~step - 4 : As many these K neighbours, count the numbers of the data points in each category.~~

- step 6: Assign the new data points to that category for which the no. of the neighbor is maximum.
- step 6: Our model is ready.
- * How to select value of k.
 - A very low values for k such as $k=1$ or $k=r$, can be noisy & lead to the effects of outliers in the model.
 - Large values for k are good, but it may find some difficulties.

* Distance Matrix in KNN

1. Euclidean Distance: The most commonly used distance matrix that measures the "straight line" distance between two points p in Euclidean space.

$$d(x, y) = \sqrt{\sum (x_i - y_i)^2}$$

* Key features of KNN,

~~- Non-parametric: KNN does not make any assumptions about the underlying data distribution, making it flexible for various type of datasets.~~



- Lazy learning : KNN is a lazy algorithm because it does not learn a discriminative memorizer the training dataset. Predictions are made only when the query sample is provided.

* Result: Hence the implementation of k-nearest neighbour algorithm successfully.

Result A

Practical No. 7.

* Aim:- Demonstrate the working of the decision tree based ID3 algorithm.

* Theory:-

- The ID3 algorithm (Iterative Dichotomiser) is a popular decision tree algorithm used for classification tasks.
- The ID3 algorithm builds the decision tree by recursively partitioning the dataset based on the attribute that results in the maximum information gain.

* Steps in ID3 Algorithm :

1. Entropy calculation:

- Entropy measures the impurity or uncertainty in the dataset.

formula:-

$$\text{Entropy } (S) = - \sum p_i \log_2 (p_i)$$

- lower entropy means less randomness, and higher entropy means ate disorder in the data.

2. Information Gain:

- Information Gain (IG) helps select the attribute that best splits the data by reducing entropy the most.

$IG_I(S.A) = \text{Entropy}(S) - \sum_{i=1}^n (f_{Si} * \text{Entropy}(S|A_i))$

The attribute with the highest information gain is selected as the root or interested node.

3. Tree construction:

- Recursively partition the dataset based on the attribute with the highest information gain until:
 - All instances belong to the same class.
 - There are no more attributes to split.

4. Stopping Criteria:

- When no. more attributes are available, or when all instances in the subset have the same label.

* Application:

- used in classification problems such as medical diagnosis, customer segmentation, and more.

Result :- Hence, a program executed successfully.

Naresh

Practical No. 08



* Aims- Implement a Bayesian network algorithm considering medical data, weather data.

* Theory:-

- Naive Bayes algorithm

It is the one of the simple and most effective classification algorithm which help in building part model that can make with predictions. It is the probitatic classifier which mean it predict on the basis of probability of an object.

- Naive of is called Naive because if assume that the occurrence of certain feature is independent of occurrence of other features.

- Bayes It is called Bayes because : It depends on principle of Bayes theorem.

formula :-
$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}$$

• $P(A)$ is prior of 'n' ; a probability of went before evidence seen.

• $P(B|A)$ is open probability the likelihood that hypothesis will come to based on evidence.



- $P(A|B)$ is Posterior probability the likelihood and that a hypothesis will come to based on outcomes.
- Types :-
 - 1) Gaussian Naive Bayes.
 - 2) Multinomial Naive Bayes.
 - 3) Bernoulli Naive Bayes.

* Advantages :-

- Simple & fast.
- Works well with high-dimensional data.
- Effective for that classification.

* Disadvantages .

- Independent Assumption.
- zero probability problem.

* Results - Hence we perform practical successfully.

Model A

Practical NO. 9.



* Aim:- Implement & demonstrate the K-means algorithm.

* Theory:-

K-mean clustering is an algorithm that is used to solve the clustering problems in machine learning & data science. It is an unsupervised learning algorithm, which groups the unlabeled dataset into different clusters.

Here K defines the no. of pre-defined clusters that need to be created in the process, as if $K=2$, there will be two clusters, & for $K=3$, there will be three clusters, & so on.

The algorithm takes the unlabeled dataset as input, divides the dataset into K-no. of clusters, & repeats the process until it does not find the best clusters.

• The K-means clustering algorithm mainly performs two tasks:-

i) ~~Determines the best value for k center points or centroids by an iterative process.~~

ii) Assigns each data point to its closest K-center. Those data points which are near to the particular K-center, create a cluster.



* How does the K-means algorithm work?

Step 1:- Select no. k to decide no. of clusters.

Step 2:- Select random k points as centroids.

Step 3:- Assign each data point to their closest centroid, which will form the predefined k-clusters.

Step 4:- calculate the variance and place a new centroid of each cluster.

Step 5:- Repeat the third steps, which means reassign each datapoint to the new closest centroid of each cluster.

Step 6:- If any reassignment occurs, then go to step 4 else go to FINISH.

Step 7:- The model is ready.

Its objective is to divide data into clusters, making similar data points part of the same group

* Result:- Hence, the implementing the K-mean algorithm successfully.

Ques A

Practical No. 10.

- * Aims:- Write a python program for to check whether email is spam or not.

- * Theory:-

- Spam detection is the process of identifying & filtering unsolicited and potentially harmful emails, typically referred to as "spam". Machine learning, natural language processing (NLP), and statistical techniques are commonly used to classify emails as "spam" or "not spam" based on patterns in the email content, metadata, or subject line.

- In a spam detection system, we use a dataset of labeled emails (either spam or not) to train a machine learning model. The model learns patterns into words, phrases, and structures of these emails to predict the probability that new, unseen emails are spam.

- Key techniques include text processing (removing stop words, punctuation, normalization), converting text to numerical format using methods like bag of words or TF-IDF, and classification algorithms like naive Bayes, support vector machines (SVM), or neural networks.

* Key Components:

① Data processing: Text is converted to lower case. Punctuation is removed.

② Vectorization:

- CountVectorizer converts the text into a matrix of token counts (Bag of words).

③ Model training:

- The multinomial NB classifier is trained on the vectorized email data.

④ Prediction:

- The model is used to predict whether new emails are spam or not.

⑤ Evaluation:

- The program calculates the accuracy of the model and provides a confusion matrix to evaluate the performance.

* Result: Hence, a program executed successfully.

Mill A