

3) i) Problem Statement :-

Develop a program to implement principal component analysis (PCA) for reducing the dimensionality of the iris dataset from 4 features to 2

ii) Objective :-

The purpose of this program is to reduce the dimensionality of the iris dataset using PCA from 4 features to 2, making it easier to visualize and analyze the data while retaining most of the important information

iii) Algorithm :-

- a) Import necessary libraries :- import numpy, pandas, matplotlib and sklearn libraries
- b) Load Dataset : load the iris dataset using the load-iris() function
- c) Prepare data : convert the dataset into a pandas DataFrame and add labels
- d) Apply PCA : Reduce the dimensionality from 4 features to 2 using PCA
- e) Create DataFrame : create a DataFrame for the reduced data with labels
- f) Visualize Data : Plot the reduced data points with different colors for each class
- g) Display plot : Show the plot with labels, grid and plot

iv) Source / Program Code

```
import numpy as np
import pandas as pd
from sklearn.datasets import load_iris
from sklearn.decomposition import PCA
import matplotlib.pyplot as plt
```

```
iris = load_iris()
data = iris.data
labels = iris.target
label_names = iris.target_names
```

```
iris_df = pd.DataFrame(data = iris.data, column = iris.  
                        - feature_names)
```

```
pca = PCA(n_components = 2)
```

```
data_reduced = pca.fit_transform(data)
```

```
reduced_df = pd.DataFrame(data = data_reduced,  
                           - columns = ['PC1', 'PC2'])
```

```
reduced_df['labels'] = labels
```

```
plt.figure(figsize = (10, 6))  
colors = ['r', 'g', 'b']
```

```
for i, label in enumerate(np.unique(labels)):  
    plt.scatter(
```


Name of the Experiment :

```
reduced_df[reduced_df['label'] == label]
                                ['PC1'],
reduced_df[reduced_df['label'] == label]
                                ['PC2'],
c = colors[i],
label = label - names[label]
)
```

```
plt.title('PCA on IRIS IRIS Dataset')
plt.xlabel('PC1')
plt.ylabel('PC2')
plt.legend()
plt.grid()
plt.show()
```

v) Compilation and Execution steps

- * Save the program in python file (eg, iris-pca.py)
- * Open a terminal or command prompt
- * Navigate to the directory containing the python file
- * Run the program using the command

```
python iris-pca.py
```

vi) Sample input and output

Sample input:

The program uses the built-in iris dataset, so no external input is required.

Sample output:

A scatter plot showing the reduced Iris dataset with 2 principal components (PC1 and PC2) plotted against each other. Different colors represent different classes.

vii) Explanation of Output:

The plot shows the Iris dataset reduced to two principal components using PCA. Each class (Setosa, Versicolour, Virginical) is represented by a different color. The plot helps in understanding how well the PCA has separated the different classes.

viii) Observation and Analysis:

- * PCA reduced the dataset from 4 dimension to 2 dimension effectively
- * Setosa is well separated from the other two classes
- * Versicolour and Virginica have some overlapping points, indicating similarities between them

ix) Conclusion:

The program successfully implemented PCA to reduce the dimensionality of the Iris dataset. Visualization of the reduced data helps in identifying the separation between different classes, which can be further used for classification tasks.

4) i) Problem Statement

For a given set of training data examples stored in a .CSV file, implement and demonstrate the find-S algorithm to output a description of the set of all hypotheses consistent with the training examples

ii) Objective :

The purpose of this program is to implement the Find-S algorithm for finding the most specific hypothesis from a given set of training data stored in a .CSV file

iii) Algorithm :

- * Import necessary libraries : Import the pandas library for handling the .CSV file
- * Define Find-S Algorithm Function : Implement the Find-S Algorithm as a function that reads data from a .CSV file
- * Initialize Hypothesis : Start with the most general hypothesis
- * Update Hypothesis : For each positive example, refine the hypothesis by comparing attribute values
- * Return Final Hypothesis : Output the most specific hypothesis consistent with all positive training example.

iv) Source / Program Code

```
import pandas as pd

def find-s-algorithm(file-path):
    data = pd.read_csv(file-path)

    print("Training data:")
    print(data)
    attributes = data.columns[:-1]
    class-label = data.columns[-1]
    hypothesis = ['?'] for i in attributes

    for index, row in data.iterrows():
        if row[class-label] == 'Yes':
            for i, value in enumerate(row
                                      - [attributes]):
                if hypothesis[i] == "?" or
                   hypothesis[i] == value:
                    hypothesis[i] = value
            else:
                hypothesis[i] = '?'

    return hypothesis

file-path = 'training-data.csv'
hypothesis = find-s-algorithm(file-path)
print("In the final hypothesis is:", hypothesis)
```


v) Compilation and Execution steps

- * prepare a csv file (e.g training-data.csv) with training data examples
- * Save the program in a python file (e.g find-s.py)
- * open a terminal or Command prompt
- * Navigate to the directory containing the python file and csv file
- * Run the program using the command

```
python find-s.py
```

vi) Sample Input and Output

Sample Input:

A s .csv file containing training data examples with various attributes and a class label (Yes/No)

Sample output:

The final hypothesis generated by the find-s algorithm, displayed as a list of attribute values or '?' indicating generalization

vii) Explanation of Output

The output shows the most specific hypothesis that covers all positive training examples. If a value differs between positive examples, the hypothesis will have a '?' at that position, indicating generalization.

viii) Observation and Analysis:

- * The Find-S algorithm produces a hypothesis that is as specific as possible, covering only positive example
- * It cannot handle noisy data or negative examples properly

xi) Conclusion:

The find-S algorithm was successfully implemented and tested with a .CSV file. The algorithm provides the most specific hypothesis that fits all positive training examples. However, it is sensitive to noise and incomplete data.