# IRIS DATASET ANALYSIS REPORT

**1.Introduction:**

The Iris dataset is one of the most famous datasets used in the field of machine learning and statistics. It contains 150 observations of iris flowers, categorized into three species: setosa, versicolor, and virginica. The dataset includes four features (sepal length, sepal width, petal length, and petal width) which are numerical attributes representing the physical dimensions of the flowers.

In this report, we will go through the steps taken to perform exploratory data analysis (EDA) on the Iris dataset, including the methodologies used, and provide explanations for the patterns identified.

**2.Methodology:**

2.1 Importing Modules

The first step in any data analysis project is to import the necessary libraries. For this analysis, we used the following Python libraries:

- pandas: For data manipulation and analysis.
- numpy: For numerical operations.
- Matplotlib and seaborn: For data visualization.
- scikit-learn: For machine learning tasks.

Code:

```
import pandas as pd

import numpy as np

import matplotlib.pyplot as plt

import os

import seaborn as sns

from sklearn.preprocessing import LabelEncoder

from sklearn.model_selection import train_test_split

from sklearn.linear_model import LogisticRegression

from sklearn.neighbors import KNeighborsClassifier

from sklearn.tree import DecisionTreeClassifier
```

- import pandas as pd: Imports the pandas library and aliases it as `pd`. Pandas is used for data manipulation and analysis, providing data structures like DataFrames.
- import numpy as np: Imports the numpy library and aliases it as `np`. Numpy is used for numerical operations and handling arrays.
- import matplotlib.pyplot as plt: Imports the `pyplot` module from the matplotlib library and aliases it as `plt`. It is used for creating static, interactive, and animated visualizations in Python.

- import os: Imports the os module, which provides a way of using operating system-dependent functionality like reading or writing to the file system.
- import seaborn as sns: Imports the seaborn library and aliases it as `sns`. Seaborn is used for making statistical graphics, providing a high-level interface for drawing attractive and informative statistical graphics.
- from sklearn.preprocessing import LabelEncoder: Imports the `LabelEncoder` class from the scikit-learn library's preprocessing module. It is used to convert categorical labels into numerical values.
- from sklearn.model_selection import train_test_split: Imports the `train_test_split` function from scikit-learn's model_selection module. It is used to split datasets into training and testing sets.
- from sklearn.linear_model import LogisticRegression: Imports the `LogisticRegression` class from scikit-learn's linear_model module. It is used to perform logistic regression, a linear model for binary classification tasks.
- from sklearn.neighbors import KNeighborsClassifier: Imports the `KNeighborsClassifier` class from scikit-learn's neighbors module. It is used for implementing the k-nearest neighbors algorithm for classification tasks.
- from sklearn.tree import DecisionTreeClassifier: Imports the `DecisionTreeClassifier` class from scikit-learn's tree module. It is used for implementing decision tree algorithms for classification tasks.

## 2.2 Loading the Dataset

The Iris dataset can be loaded directly from the seaborn library, which comes with a variety of datasets for practice and demonstration purposes.

Code:

```
df = pd.read_csv('Iris.csv')
```

## 2.3 Data Preprocessing

Data preprocessing is a crucial step in the data analysis and machine learning pipeline. It involves cleaning and transforming raw data into a format that is more suitable for analysis.

### 2.3.1 Handling Missing Values

The first step in data preprocessing is to check for and handle any missing values. For the Iris dataset, we need to verify if there are any missing values.

Code:

```
# Check for missing values

print(iris.isnull().sum())
```

If there are missing values, we can handle them by either removing the rows/columns with missing values or imputing them with appropriate values (e.g., mean, median)

**2.4 Exploratory Data Analysis (EDA)**

We started by examining the basic statistics of the dataset to understand the distribution and central tendencies of the features.

Code:

iris.describe()

df.head()

df.info()

2.4.1 Data Visualization

Visualizing the data helps in identifying patterns, correlations, and outliers.

Histograms:

Histograms are useful for understanding the distribution of individual features.

Scatter plots:

Scatter plots help in visualizing the relationship between two numerical features.


**2.5 Correlation Matrix**

A correlation matrix helps in identifying relationships between different features.

Code:

df_numeric = df.drop('Species', axis=1)

corr = df_numeric.corr()

fig, ax = plt.subplots(figsize=(5,4))

sns.heatmap(corr, annot=True, ax=ax, cmap='coolwarm')

plt.show()


**2.6 Label Encoding**

Since machine learning models typically require numerical inputs, we need to convert the categorical labels (species) into numerical values.

Code:

from sklearn.preprocessing import LabelEncoder

le = LabelEncoder()

df['Species'] = le.fit_transform(df['Species'])

**2.6 Model Training**

For model training, we used the Random Forest classifier due to its robustness and accuracy.

Splitting the Data: The dataset was split into training and testing sets.

Steps for Model Training

1. Import Required Libraries

2. Prepare Data for Training

3. Train Models

4. Evaluate Models

Code:

```
from sklearn.model_selection import train_test_split

X = df.drop(columns=['Species'])

Y = df['Species']

# Calling  the train_test_split function with the arguments

x_train, x_test, y_train, y_test = train_test_split(X, Y, test_size=0.3, random_state=42)
```

**Logistic Regression:**

- **Classification Report**: Provides precision, recall, F1-score, and support for each class.

- **Confusion Matrix**: Shows the counts of true positives, true negatives, false positives, and false negatives.

Code:

```
from sklearn.linear_model import LogisticRegression

model = LogisticRegression()

model.fit(x_train,y_train)

print("Accuracy:",model.score(x_test,y_test)*100)
```

**Decision Tree:**

- **Classification Report**: Provides precision, recall, F1-score, and support for each class.

- **Confusion Matrix**: Shows the counts of true positives, true negatives, false positives, and false negatives.

Code:

```
from sklearn.tree import DecisionTreeClassifier

model = DecisionTreeClassifier()

model.fit(x_train,y_train)

print("Accuracy:",model.score(x_test,y_test)*100)
```

**K-Nearest Neighbors:**

- **Classification Report**: Provides precision, recall, F1-score, and support for each class.

- **Confusion Matrix**: Shows the counts of true positives, true negatives, false positives, and false negatives.

Code:

```
from sklearn.neighbors import KNeighborsClassifier

model = KNeighborsClassifier()

model.fit(x_train,y_train)

print("Accuracy:",model.score(x_test,y_test)*100)
```

**2.7 Performance Evaluation:**

To evaluate the performance of each model based on accuracy, we calculate the accuracy score for the Logistic Regression, Decision Tree, and K-Nearest Neighbors models on the test dataset. Accuracy is the ratio of correctly predicted instances to the total instances in the dataset.

Steps for Performance Evaluation

1. Import Required Libraries

2. Calculate Accuracy for Each Model

3. Compare the Accuracies

Based on the accuracy scores, all three models (Logistic Regression, Decision Tree, and K-Nearest Neighbors) perform exceptionally well on the Iris dataset. Given the simplicity of the Iris dataset and the perfect accuracy achieved, it is essential to test these models on more complex datasets to evaluate their true performance and generalizability.

**3. Patterns Identified**

3.1 Pattern Identification in the Iris Dataset

In the exploratory data analysis (EDA) phase, we identified several patterns and relationships in the Iris dataset. These patterns help us understand the characteristics of the data and provide insights that are valuable for building predictive models. Here are the key patterns identified:

3.1.1  Distributions of Individual Features

**Sepal Length:**

Distribution: The sepal length is approximately normally distributed with a slight right skew.

Typical Values: Most sepal lengths are in the range of 4.5 to 7.5 cm.

**Sepal Width**

Distribution: The sepal width has a somewhat bimodal distribution, with peaks around 3.0 and 3.5 cm.

Typical Values: Sepal width values primarily range from 2.0 to 4.0 cm.

**Petal Length**

Distribution: Petal length is right-skewed, indicating more flowers with shorter petal lengths.

Typical Values: Most petal lengths are between 1.0 and 6.5 cm.

**Petal Width**

Distribution: Similar to petal length, petal width is also right-skewed.

Typical Values: Petal width values mostly range from 0.1 to 2.5 cm.

**3.2 Pairwise Relationships Between Features**

- Sepal Length vs. Sepal Width

Pattern: There is no clear separation between species based on sepal length and sepal width alone.

Species Distribution: Setosa species tend to have higher sepal width compared to Versicolor and Virginica.

- Petal Length vs. Petal Width

Pattern: There is a clear separation between the three species. Setosa has the smallest petal dimensions, Versicolor has intermediate values, and Virginica has the largest.

Species Distribution This feature pair provides a strong discriminatory power for distinguishing between species.

- Sepal Length vs. Petal Length

Pattern: Sepal length and petal length show a distinct separation between species, particularly for Setosa.

Species Distribution: Setosa is clearly separated from Versicolor and Virginica, with the latter two overlapping more.

- Sepal Width vs. Petal Width

Pattern: There is some overlap between Versicolor and Virginica, but Setosa is distinctly separated.

Species Distribution: Setosa has a significantly different distribution compared to the other two species.

**3.3 Species Characteristics**

- Setosa

Sepal Dimensions: Generally higher sepal width compared to the other species.

Petal Dimensions: Smallest petal length and width, providing clear separation from Versicolor and Virginica.

- Versicolor

Sepal Dimensions: Intermediate values for both sepal length and width.

Petal Dimensions: Intermediate petal length and width, overlapping with Virginica but distinguishable from Setosa.

- Virginica

Sepal Dimensions: Longer sepal length compared to Setosa and Versicolor.

Petal Dimensions: Largest petal length and width, providing clear separation from Setosa and some overlap with Versicolor.

**4. Conclusion**

These patterns highlight the distinguishing characteristics of each species in the Iris dataset. Understanding these patterns is crucial for building effective classification models, as they provide insights into which features are most informative for distinguishing between species. The clear separation observed in petal dimensions, particularly petal length and width, suggests that these features are highly discriminative and should be emphasized in model training and feature selection processes.