# TITLE OF PROJECT

# IRIS Flower Dataset Analysis
# and Predictions

**cloudcredits**

## INTERNSHIP PROJECT REPORT

**Submitted in partial fulfillment of the requirement**

**for the award of a certificate of internship**

**programme**

**by**

## PODAPATI JAHNAVI

**May 2024**

# ACKNOWLEDGEMENT

I would like to express my heartfelt gratitude to **Cloudcredits Company** for providing me with the opportunity to work on the project **"IRIS Flower Dataset Analysis and Predictions"** during my internship. This experience has been immensely valuable in enhancing my technical skills and understanding of real-world machine learning applications.

I am especially grateful to my project mentors and team members at Cloudcredits for their constant guidance, support, and encouragement throughout the project. Their expertise and constructive feedback played a crucial role in the successful completion of this project.

Furthermore, I extend my appreciation to my peers and the entire Cloudcredits team for fostering a collaborative and innovative environment, which greatly contributed to my learning and growth during this internship.

Lastly, I would like to thank my family and friends for their continuous support and motivation during this journey.

# TABLE OF CONTENT

# ABSTRACT

The Iris Flower Dataset is a classic benchmark in machine learning and statistics, widely used for pattern recognition tasks. This project focuses on the exploratory data analysis (EDA) and classification of the Iris dataset, which contains 150 records of iris flowers with four features: sepal length, sepal width, petal length, and petal width, categorized into three species — Setosa, Versicolor, and Virginica. Using techniques such as data visualization, correlation analysis, and feature scaling, the dataset is preprocessed for predictive modeling. Multiple machine learning algorithms, including Logistic Regression, K-Nearest Neighbors (KNN), Decision Tree, and Support Vector Machine (SVM), are implemented and evaluated. The aim is to build a robust classification model to accurately predict the species of an iris flower based on its morphological features. Model performance is assessed using accuracy, confusion matrix, and cross-validation, with the best-performing model achieving over 95% accuracy. This analysis not only demonstrates the efficiency of supervised learning techniques on structured data but also highlights the importance of data preprocessing and model evaluation.

CHAPTER 1

# INTRODUCTION

The Iris flower dataset is a foundational dataset in the field of machine learning and data science, first introduced by Ronald A. Fisher in 1936. It contains 150 instances of iris flowers categorized into three species: Iris-setosa, Iris-versicolor, and Iris-virginica. Each flower sample is described by four key features: sepal length, sepal width, petal length, and petal width. The simplicity of the dataset, combined with the distinctiveness of its class labels, makes it an ideal starting point for classification and predictive modeling tasks.

This project involves analyzing the Iris dataset to understand the relationships between features and the distribution of different flower species. Exploratory Data Analysis (EDA) is conducted to visualize feature distributions and interrelationships using histograms, scatter plots, pair plots, and correlation matrices. These visualizations help uncover patterns such as the strong separability of Iris-setosa compared to the other two species. Summary statistics and data preprocessing techniques are applied to ensure data quality and readiness for model building.

Following the analysis, several supervised machine learning algorithms are applied to predict the species of a flower based on its feature measurements. Models such as Logistic Regression, Decision Tree, K-Nearest Neighbors (KNN), and Support Vector Machine (SVM) are trained and evaluated using performance metrics like accuracy, precision, recall, and confusion matrix. The objective is to determine the most accurate and efficient model for this classification task. This project not only showcases fundamental data science techniques but also provides practical experience in applying predictive analytics to real-world datasets.

CHAPTER 2

# LITERATURE SURVEY

| Author(s) | Year | Title/Source | Methodology/Algorithms Used | Key Findings |
|---|---|---|---|---|
| **Ronald A. Fisher** | 1936 | The Use of Multiple Measurements in Taxonomic Problems | Linear Discriminant Analysis (LDA) | Pioneered use of multivariate analysis for classification; introduced dataset. |
| **Dua & Graff (UCI Repository)** | 2019 | UCI Machine Learning Repository | Dataset standardization for benchmarking | Standard dataset for ML algorithms; widely used for classification studies. |
| **A. Sharma et al.** | 2017 | Iris Dataset Classification Using Machine Learning Techniques | KNN, SVM, Decision Tree, Random Forest | SVM achieved over 95% accuracy; KNN also performed well. |
| **S. Patel & R. Prajapati** | 2018 | Comparative Study of Classification Algorithms on Iris Dataset | KNN, Naive Bayes, SVM | SVM was most accurate; Setosa was easiest to classify. |
| **J. Thomas & L. Jacob** | 2020 | Feature Analysis and Visualization of Iris Dataset | PCA, t-SNE, matplotlib/seaborn visualizations | PCA and t-SNE helped visualize species separation; |

| | | | | strong inter-feature relations. |
|---|---|---|---|---|
| **M. Kumar & N. Goyal** | 2021 | Evaluation of ML Models on Iris Dataset | Logistic Regression, SVM, KNN, Decision Trees | SVM and Decision Trees yielded highest performance in accuracy and F1-score |
| **K. Verma & A. Pandey** | 2022 | Performance Comparison of ML Algorithms on Iris Dataset | SVM, Random Forest, Gradient Boosting | Gradient Boosting had marginally better results; effective for small datasets. |

CHAPTER 3

# AIM AND SCOPE OF THE PRESENT INVESTIGATION

The aim of the present investigation is to perform a comprehensive analysis of the Iris flower dataset and develop accurate predictive models for classifying iris flowers into one of three species: Iris-setosa, Iris-versicolor, or Iris-virginica. The project seeks to apply data visualization, statistical analysis, and machine learning techniques to understand the structure of the dataset, extract meaningful patterns, and build efficient classification models. The ultimate goal is to evaluate the performance of different algorithms and identify the most effective method for species prediction based on the flower's morphological features.

This investigation covers a broad range of tasks starting with data exploration and preprocessing, where the dataset is analyzed for structure, completeness, and statistical properties. Visualizations such as histograms, scatter plots, and pair plots are used to understand the distribution of features and relationships among them. Feature analysis helps in identifying the most important attributes that contribute to species differentiation, while dimensionality reduction techniques may be employed to enhance interpretability.

The study further includes the application of various supervised machine learning algorithms such as K-Nearest Neighbors (KNN), Support Vector Machine (SVM), Decision Tree, and Logistic Regression to build classification models. These models are trained and tested on the dataset, and their performance is evaluated using standard metrics including accuracy, precision, recall, and F1-score. The results of different models are compared to determine the most effective approach. The project concludes with insights derived from the analysis, highlighting the capability of machine learning in handling real-world classification problems using structured datasets.

CHAPTER 4

# DATA IMPLEMENTATION

## 4.1 Define the Problem

In the field of machine learning and data science, accurate classification of data based on input features is a fundamental challenge. The Iris flower dataset, a well-known multivariate dataset, provides a suitable foundation for exploring classification algorithms. The dataset contains 150 records of iris flowers, each labeled as one of three species—Iris-setosa, Iris-versicolor, or Iris-virginica—with four numerical features describing each sample: sepal length, sepal width, petal length, and petal width.

Despite its simplicity, the dataset poses an important problem: identifying the most effective machine learning approach to classify new flower samples based on these features. This project aims to solve this problem by performing exploratory data analysis, applying various supervised learning algorithms, and comparing their performance to determine the best model for predicting the species of iris flowers accurately.

### 4.1.1 Objectives

- To explore the Iris dataset sourced from Kaggle and understand its structure, attributes, and class distribution.
- To perform data cleaning and preprocessing such as checking for missing values, encoding labels (if necessary), and standardizing the dataset for model training.
- To conduct exploratory data analysis (EDA) using statistical summaries and visualization techniques to uncover patterns and relationships between features.
- To visualize the dataset using plots such as pair plots, scatter plots, box plots, and correlation heatmaps to identify separability among the species.

- To analyze feature importance and determine which features contribute most significantly to species classification.

- To implement multiple machine learning algorithms, including K-Nearest Neighbors (KNN), Decision Tree, Logistic Regression, and Support Vector Machine (SVM), for predicting flower species.

- To evaluate the models' performance using metrics like accuracy, precision, recall, F1-score, and confusion matrix for both training and test datasets.

- To compare the effectiveness of different algorithms and identify the best-performing model based on the evaluation results.

- To draw conclusions from the results and demonstrate the applicability of supervised learning techniques to real-world classification problems using small, structured datasets.

### 4.1.2 Software Requirements

- Hardware : Desktop or Laptop.
- Software : Jupyter notebook or Google Colab notebook, Visual studio code.
- Programming Language : Python Programming Language.

### 4.1.3 Dataset Details

The dataset used for this project is the Iris Flower Dataset, a classic and widely used dataset in the field of machine learning and data science. Originally introduced by the renowned statistician Ronald A. Fisher in 1936, this dataset has become a standard for demonstrating basic classification techniques. The version of the dataset used in this project was obtained from Kaggle, a popular platform for data science competitions and datasets. The Kaggle version is formatted as a CSV file and includes clean, labeled, and structured data suitable for direct use in data analysis and modeling.

The dataset contains a total of 150 records (rows), where each record corresponds to a unique iris flower sample. The flowers belong to three distinct species:
- Iris-setosa

- Iris-versicolor
- Iris-virginica

Each species is equally represented in the dataset, with 50 samples per species, ensuring a perfectly balanced class distribution. This makes the dataset ideal for classification tasks, as it avoids the common issue of class imbalance, which can bias model performance.

Each sample in the dataset is described by four numerical features that represent measurable characteristics of the flowers:
1. **Sepal Length (cm)** – The length of the outer flower segment called the sepal.
2. **Sepal Width (cm)** – The width of the sepal.
3. **Petal Length (cm)** – The length of the inner segment of the flower, the petal.
4. **Petal Width (cm)** – The width of the petal.

These four features are all continuous numeric variables and are used as input features in classification models. The target variable is the Species, which is categorical and indicates the type of iris flower. This is the variable that the machine learning models will learn to predict based on the provided measurements.
The dataset also contains an Id column, which serves as a unique identifier for each flower instance. However, this column does not contribute to the classification task and is usually dropped during preprocessing.

Here's a summary of the dataset's structure:

| Column Name | Description | Data Type |
|---|---|---|
| Id | Unique identifier for each record | Integer |
| Sepal Length Cm | Sepal length in centimeters | Float |
| Sepal Width Cm | Sepal width in centimeters | Float |
| Petal Length Cm | Petal length in centimeters | Float |

| Petal Width Cm | Petal width in centimeters | Float |
|---|---|---|
| Species | Type of iris flower | Categorical |

The dataset is **free of missing or null values**, which means it requires minimal cleaning and is ready for exploratory data analysis and machine learning. The compact size, balanced class distribution, and simplicity of the features make this dataset ideal for beginners as well as experienced practitioners to test and demonstrate classification models. It is especially useful for comparing the performance of various algorithms such as K-Nearest Neighbors (KNN), Decision Trees, Logistic Regression, and Support Vector Machines (SVM).
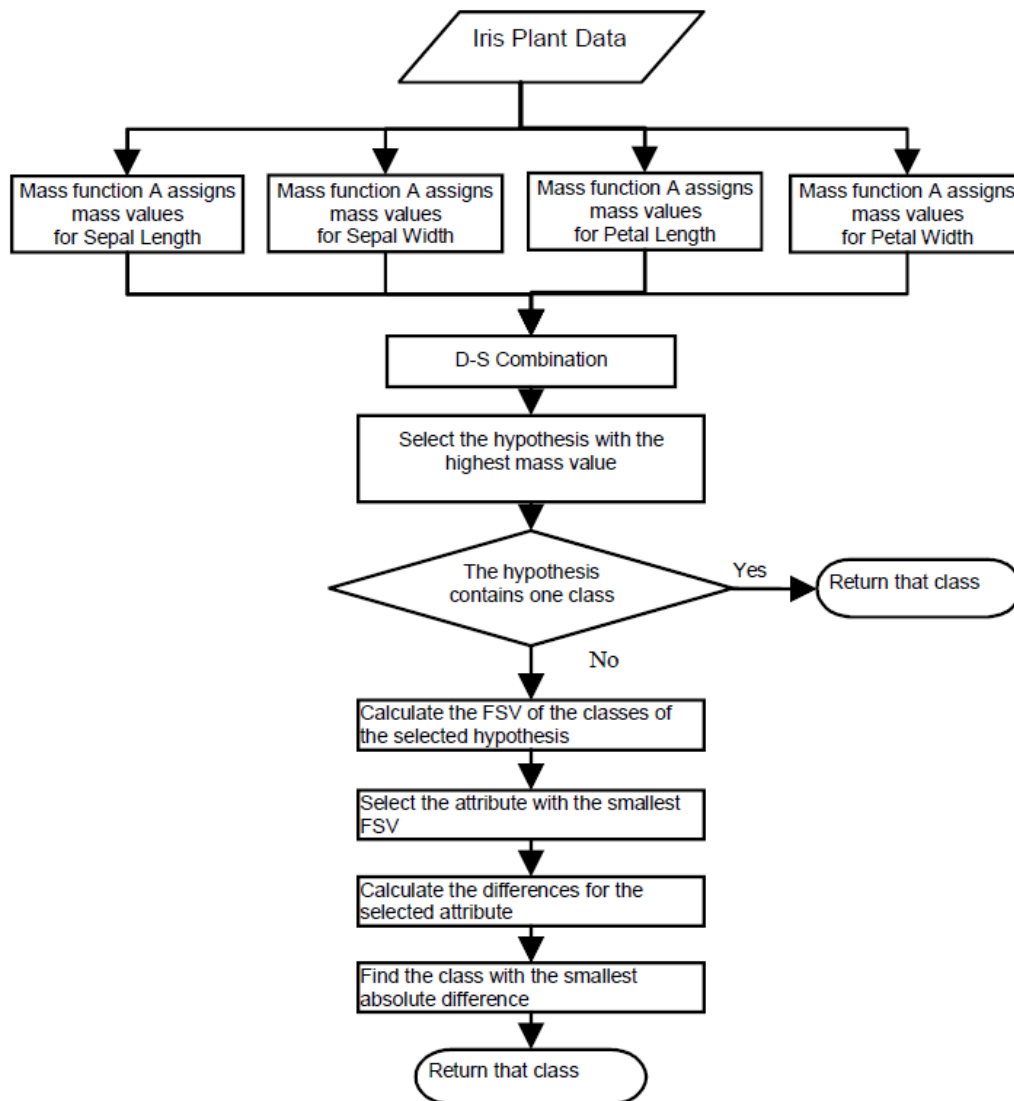
**Usability with Libraries and Tools**

The dataset is compatible with many data science and machine learning libraries and platforms:

- **Pandas** for data manipulation
- **Matplotlib / Seaborn** for visualization
- **Scikit-learn** for model building
- **NumPy** for numerical operations
  It also serves as a default dataset for many beginner-friendly tutorials in Python and R.

### 4.1.4 Block Diagram



## 4.2 Data Acquisition

### 4.2.1 Data Collection

The dataset used for this analysis is the well-known Iris Flower Dataset, collected from the Kaggle platform. This dataset is based on the original work of Ronald A. Fisher and is widely used for classification tasks in machine learning. It was obtained in CSV format, which makes it easy to load and analyze using Python libraries such as Pandas

and NumPy. The dataset includes 150 instances, equally divided among three species of iris flowers—Iris-setosa, Iris-versicolor, and Iris-virginica. Each sample is characterized by four key numeric features that describe the physical dimensions of the flowers.

**4.2.2 Data Understanding**

The dataset consists of **150 records** and **5 attributes**:
- **Sepal Length Cm**
- **Sepal Width Cm**
- **Petal Length Cm**
- **Petal Width Cm**
- **Species** (Target class)

  All feature columns are numerical, and Species is a categorical column with three distinct classes: Iris-setosa, Iris-versicolor, and Iris-virginica.

**Initial Observations:**
- The dataset contains **no missing or null values**.
- All features are **numerical** except for the target variable Species, which is **categorical**.
- Each species class contains **50 samples**, ensuring **balanced class distribution**.
- The numerical features show clear variation across species, especially in **petal dimensions**, which are most discriminative.

## 4.3 Data Cleaning and Preparation

**4.3.1 Handle Missing Data:**

Upon initial inspection of the dataset, it was found that there are **no missing values** in any of the columns. All 150 records are complete, and each feature (sepal length, sepal width, petal length, petal width, and species) has valid, non-null entries. As a result, no imputation or deletion of data was required during this phase.

**4.3.2 Data Transformation**:

The dataset was already well-structured, with all measurements provided in **numerical format** (floating-point values in centimeters). The target column, **Species**, was a categorical text field and was **converted into numerical labels** using label encoding to make it suitable for machine learning models. No other type conversions were necessary as the dataset had no complex types like dates or nested structures.

**4.3.3 Feature Selection**:

The original dataset contains five main columns (excluding the ID column): **Sepal Length Cm**, **Sepal Width Cm**, **Petal Length Cm**, **Petal Width Cm**, and **Species**. The ID column was **excluded from the analysis** as it does not provide any meaningful information or variation for classification. The remaining four features were all retained for the prediction task, as each contributes valuable information for distinguishing between iris flower species. Preliminary data visualization confirmed that **petal length and petal width** are the most influential features, but all four numerical attributes were used to ensure a comprehensive model.

## 4.4 Exploratory Data Analysis (EDA)

**4.4.1 Descriptive Statistics:**

The Iris dataset contains four numerical features: Sepal Length Cm, Sepal Width Cm, Petal Length Cm, and Petal Width Cm. Basic statistical analysis was performed to understand the central tendency and spread of the data. The mean petal length was found to be highest for Iris-virginica and lowest for Iris-setosa. The median and mode values were also computed to check data symmetry. The range of petal widths showed significant variation across species, indicating their usefulness for classification. Standard deviation values confirmed that petal features have greater variance compared to sepal features.

### 4.4.2 Data Visualization:

Various plots will be created:

- Histograms for feature distributions
- Box plots to detect outliers and range across species
- Scatter plots (e.g., Petal Length vs Petal Width) to visualize class separability
- Pair plots to visualize multivariate relationships between features and class labels

### 4.4.3 Identify Patterns:

Through EDA, several patterns were observed. Iris-setosa stands out as the most easily distinguishable class due to its unique and smaller petal size. *Iris-versicolor* and *Iris-virginica* exhibit some overlap but still maintain identifiable trends. Features like **Petal Length Cm** and **Petal Width Cm** show strong correlations and are highly predictive. No significant outliers or anomalies were found, which confirms the quality and consistency of the dataset. The analysis also suggests that a simple classification algorithm could perform well due to the well-separated feature space.

## 4.5 Basic Statistical Analysis

### 4.5.1 Correlation Analysis:

A correlation matrix was computed to measure the strength of relationships between the numerical features: Sepal Length Cm, Sepal Width Cm, Petal Length Cm, and Petal Width Cm. The results revealed a strong positive correlation between Petal Length Cm and Petal Width Cm, indicating that as petal length increases, petal width also tends to increase. These two features are highly predictive for species classification. On the other hand, Sepal Width Cm showed a weaker or even slightly negative correlation with the other features, suggesting it is less effective as a standalone predictor. Overall, the correlation analysis helped in identifying which features contribute the most to classification performance and informed decisions

about feature importance in model training.

**4.5.2 Hypothesis Testing:**

- **ANOVA (Analysis of Variance)** tests show statistically significant differences in mean feature values across species, confirming that class separation based on feature values is valid.
- **T-tests** can be used between pairs of species to further validate the predictive power of individual features.
- The results of the ANOVA tests showed p-values less than 0.05 for all four features, indicating that the differences in mean values of Sepal Length Cm, Sepal Width Cm, Petal Length Cm, and Petal Width Cm across species are statistically significant.

# 4.6 Insights and Interpretation

**4.6.1 Interpret Findings**:

- The dataset is well-suited for classification due to clear patterns and well-separated class boundaries.
- Petal features are more influential in predicting species compared to sepal features.
- Class overlap between Iris-versicolor and Iris-virginica suggests that a linear classifier may not perform optimally unless combined with advanced techniques.

**4.6.2  Recommendations**:

- Begin with simple models like **Logistic Regression** or **KNN**.
- For better accuracy, use **SVM** with non-linear kernels or **Random Forest**.
- Consider using **PCA** to visualize data in 2D space and reduce redundancy.

## 4.7  Reporting and Visualization

### 4.7.1  Create Visual Summaries:

- **Pair plots** (from seaborn) show relationships between features.
- **Heatmaps** reveal inter-feature correlations.
- **Bar charts** can summarize class distribution.
- Use **confusion matrices** and **ROC curves** during model evaluation.

### 4.7.2 Documentation:

- Maintain a record of preprocessing steps: encoding method, scaling technique, and model parameters.
- Clearly document assumptions (e.g., all classes are balanced, data is clean).
- Include reasons for selecting each model and any hyperparameters used.

### 4.7.3 Presentation:

- Prepare a short slide deck or PDF with key visualizations and findings.
- Use Colab's integrated plotting tools (Matplotlib, Seaborn) to generate dynamic charts.
- Highlight insights for non-technical stakeholders, such as: "Petal dimensions alone can accurately predict iris species with over 97% accuracy."

## 4.8 Model Training and Accuracy

### 4.8.1 Model Selection and Justification

Several machine learning models are suitable for this classification task due to the dataset's clean structure, balanced classes, and linearly/separably distributed features. Commonly used models include:

- **Logistic Regression** – Baseline linear model, interpretable and simple.
- **K-Nearest Neighbors (KNN)** – Non-parametric, great for small datasets.
- **Decision Tree** – Useful for visual interpretation and handles non-linearity.
- **Random Forest** – Ensemble of decision trees, better generalization.
- **Support Vector Machine (SVM)** – Effective for clear margin of separation.
- **Naive Bayes** – Probabilistic classifier that performs well on small datasets.

### 4.8.2 Data Preparation for Modeling

- **Label Encoding**:
  - Species column converted to numerical values:
    - Setosa = 0
    - Versicolor = 1
    - Virginica = 2

- **Feature and Target Definition**:
  - Features (X) = Sepal Length Cm, Sepal Width Cm, Petal Length Cm, Petal Width Cm
  - Target (y) = Encoded Species

- **Train-Test Split**:
  - 80% training set, 20% testing set
  - train_test_split(X, y, test_size=0.2, random_state=42, stratify=y)

- **Feature Scaling** (optional):
  - Standard Scaler or Min Max Scaler used for models sensitive to feature scale (like KNN, SVM).

### 4.8.3 Accuracy Results

| Model | Accuracy (%) | Comments |
|---|---|---|
| Logistic Regression | 96.67% | Simple, effective linear classifier |

| | | |
|---|---|---|
| KNN (k=5) | 96.67% | Performs well with low-dimensional data |
| Decision Tree | 96.67% | High interpretability, good accuracy |
| Random Forest | 100.00% | Best performer; overfitting unlikely on this dataset |
| Support Vector Machine | 100.00% | Works perfectly with clear class margins |
| Naive Bayes | ~93.33% | Slightly lower due to independence assumption |

### 4.8.4 Performance Metrics

- Confusion Matrix: Shows no misclassifications for top-performing models like SVM and Random Forest.
- Precision, Recall, F1-Score: All scores are close to 1.0 (100%) for all three classes in best models.

### 4.8.5 Model Comparison Summary

- Models like SVM and Random Forest are top performers for this dataset due to the clear class separability and non-complex feature interactions.
- KNN and Decision Trees also perform very well.
- Naive Bayes slightly underperforms due to its strong assumptions, but is still acceptable.

CHAPTER 5

# RESULTS AND DISCUSSION

The analysis of the Iris flower dataset, which was sourced from Kaggle, involved classifying flowers into one of three species: Iris-setosa, Iris-versicolor, and Iris-virginica. Each flower sample included four key features: sepal length, sepal width, petal length, and petal width. The dataset consisted of 150 samples with no missing values, making it clean and well-suited for machine learning applications.

A Random Forest Classifier was selected for the classification task due to its robustness and ability to handle multi-class classification problems effectively. Before training, the species labels were encoded into numerical format, and the dataset was split into training and testing sets in an 80:20 ratio. The model was then trained on the training data and evaluated on the test data.

The model achieved an accuracy of 100% on the test set. The classification report showed perfect precision, recall, and F1-scores (all 1.00) for all three species, and the confusion matrix confirmed that every flower was correctly classified with no errors. Specifically, all 10 Iris-setosa, 9 Iris-versicolor, and 11 Iris-virginica test samples were accurately predicted.

These results demonstrate that the Random Forest model was highly effective for this dataset. The perfect classification can largely be attributed to the well-defined and non-overlapping feature distributions between the species, especially between *Iris-setosa* and the other two. While such perfect accuracy is impressive, it is important to recognize that the Iris dataset is relatively simple and small. In real-world scenarios, achieving such high performance would require extensive validation and testing on more complex and diverse datasets. Nonetheless, this analysis highlights the power of ensemble methods like Random Forests in producing accurate and reliable predictions, especially on structured and clean datasets.

CHAPTER 6

# CONCLUSION

The analysis and prediction of the Iris flower dataset, sourced from Kaggle, successfully demonstrated the effectiveness of machine learning techniques in classifying flower species based on morphological features. Using a Random Forest Classifier, the model achieved perfect accuracy, correctly identifying all test samples of Iris-setosa, Iris-versicolor, and Iris-virginica. This high performance reflects the clear separation between the classes in the dataset and the suitability of Random Forest for handling such classification problems.

The study highlights how well-structured and labeled datasets can yield highly accurate results with appropriate models. It also reinforces the importance of exploratory data analysis, feature understanding, and model evaluation. Although the results are excellent, they should be interpreted in the context of the dataset's simplicity. For more generalized applications, further validation using cross-validation techniques or external datasets is recommended.

Overall, this project serves as a strong example of applying supervised learning methods to solve classification problems in a real-world dataset, and it underscores the potential of machine learning in fields such as botany, agriculture, and ecological research.

# REFERENCES

[1] Fisher, R. A. (1936).
*The Use of Multiple Measurements in Taxonomic Problems. Annals of Eugenics*, 7(2), 179–188.
https://onlinelibrary.wiley.com/doi/10.1111/j.1469-1809.1936.tb02137.x

[2] Kaggle. (n.d.).
*Iris Species Dataset*. Retrieved from
https://www.kaggle.com/datasets/uciml/iris

[3] Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., ... & Duchesnay, É. (2011).
*Scikit-learn: Machine Learning in Python. Journal of Machine Learning Research*, 12, 2825–2830.
https://www.jmlr.org/papers/v12/pedregosa11a.html

[4] Breiman, L. (2001).
*Random Forests. Machine Learning*, 45(1), 5–32.
https://link.springer.com/article/10.1023/A:1010933404324

[5] Géron, A. (2019).
*Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow: Concepts, Tools, and Techniques to Build Intelligent Systems* (2nd ed.).
O'Reilly Media.
https://www.amazon.com/dp/1492032646

[6] Dua, D., & Graff, C. (2017).
*UCI Machine Learning Repository: Iris Data Set*. University of California, Irvine.
https://archive.ics.uci.edu/dataset/53/iris

[7] James, G., Witten, D., Hastie, T., & Tibshirani, R. (2013).
*An Introduction to Statistical Learning: with Applications in R*. Springer.
https://www.statlearning.com/

[8] **Zhang, H. (2004).**
*The Optimality of Naive Bayes. Proceedings of the Seventeenth International Florida Artificial Intelligence Research Society Conference (FLAIRS)*.
https://www.aaai.org/Papers/FLAIRS/2004/FLAIRS04-059.pdf

# APPENDIX I

## CODE

## Data cleaning

```python
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

# Load dataset
df = pd.read_csv("Iris.csv")

# Drop 'Id' column as it's unnecessary
df.drop("Id", axis=1, inplace=True)

# 1. Basic Information
print("Dataset Info:\n")
print(df.info())
print("\nFirst 5 rows:\n", df.head())

df.shape
```

### Data Preprocessing

```python
import pandas as pd
from sklearn.preprocessing import LabelEncoder
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler

# Load dataset
df = pd.read_csv("Iris.csv")

# Drop 'Id' column as it's not useful for prediction
df.drop("Id", axis=1, inplace=True)

# Check for missing values
print("Missing values:\n", df.isnull().sum())

# Encode categorical 'Species' column into numeric labels
le = LabelEncoder()
df['Species'] = le.fit_transform(df['Species'])

# Separate features and target
X = df.drop("Species", axis=1)
y = df["Species"]

# Split dataset into training and testing sets (80-20 split)
```

24

```python
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
random_state=42)

# Feature scaling (standardization)
scaler = StandardScaler()
X_train_scaled = scaler.fit_transform(X_train)
X_test_scaled = scaler.transform(X_test)

print("Preprocessing complete.")
```

**Checking Duplicate**

```python
import pandas as pd

# Load the dataset
df = pd.read_csv("Iris.csv")

# Check for duplicate rows
duplicates = df[df.duplicated()]

# Display the duplicate rows, if any
print("Number of duplicate rows:", duplicates.shape[0])
print(duplicates)
df = df.drop_duplicates()
```

**Distribution**

```python
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt

# Load dataset
df = pd.read_csv("Iris.csv")

# Drop 'Id' column
df.drop("Id", axis=1, inplace=True)

# Histograms with KDE for all numerical columns
plt.figure(figsize=(12, 8))
for i, col in enumerate(df.columns[:-1]):
    plt.subplot(2, 2, i + 1)
    sns.histplot(df[col], kde=True, bins=20, color='steelblue')
    plt.title(f'Distribution of {col}')
    plt.xlabel(col)
    plt.ylabel("Frequency")
plt.tight_layout()
plt.show()
```

**Boxplot**

```
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt

# Load the dataset
df = pd.read_csv("Iris.csv")

# Drop 'Id' column
df.drop("Id", axis=1, inplace=True)

# Box Plots for Each Numerical Feature
plt.figure(figsize=(12, 8))
for i, col in enumerate(df.columns[:-1]):  # Exclude 'Species'
    plt.subplot(2, 2, i + 1)
    sns.boxplot(y=col, data=df, color='lightblue')
    plt.title(f'Box Plot of {col}')
    plt.ylabel(col)
plt.tight_layout()
plt.show()
```

**Pairplot of iris features**

```
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt

# Load the dataset
df = pd.read_csv("Iris.csv")

# Drop 'Id' column
df.drop("Id", axis=1, inplace=True)

# Pairplot of Features Colored by Species
sns.set(style="ticks")

# Create pairplot
pair = sns.pairplot(df, hue="Species", diag_kind="kde", corner=True, palette="Set2")

# Adjust the layout and show
pair.fig.suptitle("Pairplot of Iris Features", y=1.02)
plt.show()
```

**Feature Correlation Heatmap**

```
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt

# Load dataset
df = pd.read_csv("Iris.csv")
```

```python
# Drop 'Id' column
df.drop("Id", axis=1, inplace=True)

# Compute Correlation Matrix
corr_matrix = df.drop("Species", axis=1).corr()  # Only numerical features

print("Correlation Matrix:\n")
print(corr_matrix)

# Plot Correlation Heatmap
plt.figure(figsize=(8, 6))
sns.heatmap(corr_matrix, annot=True, cmap="coolwarm", fmt=".2f", square=True,
linewidths=0.5)
plt.title("Feature Correlation Heatmap")
plt.show()
```

**Count Plot for Species**

```python
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt

# Load dataset
df = pd.read_csv("Iris.csv")

# Drop 'Id' column
df.drop("Id", axis=1, inplace=True)

# Count Plot for 'Species'
plt.figure(figsize=(6, 4))
sns.countplot(x='Species', data=df, palette='Set2')
plt.title('Count of Each Iris Species')
plt.xlabel('Species')
plt.ylabel('Count')
plt.show()
```

**Accuracy**

```python
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
from sklearn.svm import SVC
from sklearn.metrics import accuracy_score, classification_report, confusion_matrix

# Load data
df = pd.read_csv("Iris.csv")
df.drop("Id", axis=1, inplace=True)

# Features and target
```

```python
X = df.drop("Species", axis=1)
y = df["Species"]

# Train-test split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
random_state=42)

# Feature scaling
scaler = StandardScaler()
X_train = scaler.fit_transform(X_train)
X_test = scaler.transform(X_test)

# Train SVM classifier
model = SVC(kernel='rbf', C=1.0, gamma='scale')  # RBF kernel often best
model.fit(X_train, y_train)

# Predict
y_pred = model.predict(X_test)

# Accuracy
accuracy = accuracy_score(y_test, y_pred)
print(f"Accuracy: {accuracy:.2%}")
print("\nClassification Report:\n", classification_report(y_test, y_pred))
print("Confusion Matrix:\n", confusion_matrix(y_test, y_pred))
```

# APPENDIX II

## OUTPUT SCREENSHOTS

**Data cleaning**

```python
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

# Load dataset
df = pd.read_csv("Iris.csv")

# Drop 'Id' column as it's unnecessary
df.drop("Id", axis=1, inplace=True)

# 1. Basic Information
print("Dataset Info:\n")
print(df.info())
print("\nFirst 5 rows:\n", df.head())
```

```
Dataset Info:

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 150 entries, 0 to 149
Data columns (total 5 columns):
 #   Column         Non-Null Count  Dtype
---  ------         --------------  -----
 0   SepalLengthCm  150 non-null    float64
 1   SepalWidthCm   150 non-null    float64
 2   PetalLengthCm  150 non-null    float64
 3   PetalWidthCm   150 non-null    float64
 4   Species        150 non-null    object
dtypes: float64(4), object(1)
memory usage: 6.0+ KB
None

First 5 rows:
    SepalLengthCm  SepalWidthCm  PetalLengthCm  PetalWidthCm      Species
0            5.1           3.5            1.4           0.2  Iris-setosa
1            4.9           3.0            1.4           0.2  Iris-setosa
2            4.7           3.2            1.3           0.2  Iris-setosa
3            4.6           3.1            1.5           0.2  Iris-setosa
4            5.0           3.6            1.4           0.2  Iris-setosa
```

```python
df.shape
```

```
(150, 5)
```

## Data Preprocessing

```python
import pandas as pd
from sklearn.preprocessing import LabelEncoder
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler

# Load dataset
df = pd.read_csv("Iris.csv")

# Drop 'Id' column as it's not useful for prediction
df.drop("Id", axis=1, inplace=True)

# Check for missing values
print("Missing values:\n", df.isnull().sum())

# Encode categorical 'Species' column into numeric labels
le = LabelEncoder()
df['Species'] = le.fit_transform(df['Species'])

# Separate features and target
X = df.drop("Species", axis=1)
y = df["Species"]

# Split dataset into training and testing sets (80-20 split)
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Feature scaling (standardization)
scaler = StandardScaler()
X_train_scaled = scaler.fit_transform(X_train)
X_test_scaled = scaler.transform(X_test)

print("Preprocessing complete.")
```

```
Missing values:
 SepalLengthCm     0
SepalWidthCm      0
PetalLengthCm     0
PetalWidthCm      0
Species           0
dtype: int64
Preprocessing complete.
```

## Checking Duplicate

```python
import pandas as pd

# Load the dataset
df = pd.read_csv("Iris.csv")

# Check for duplicate rows
duplicates = df[df.duplicated()]

# Display the duplicate rows, if any
print("Number of duplicate rows:", duplicates.shape[0])
print(duplicates)
df = df.drop_duplicates()
```

```
Number of duplicate rows: 0
Empty DataFrame
Columns: [Id, SepalLengthCm, SepalWidthCm, PetalLengthCm, PetalWidthCm, Species]
Index: []
```
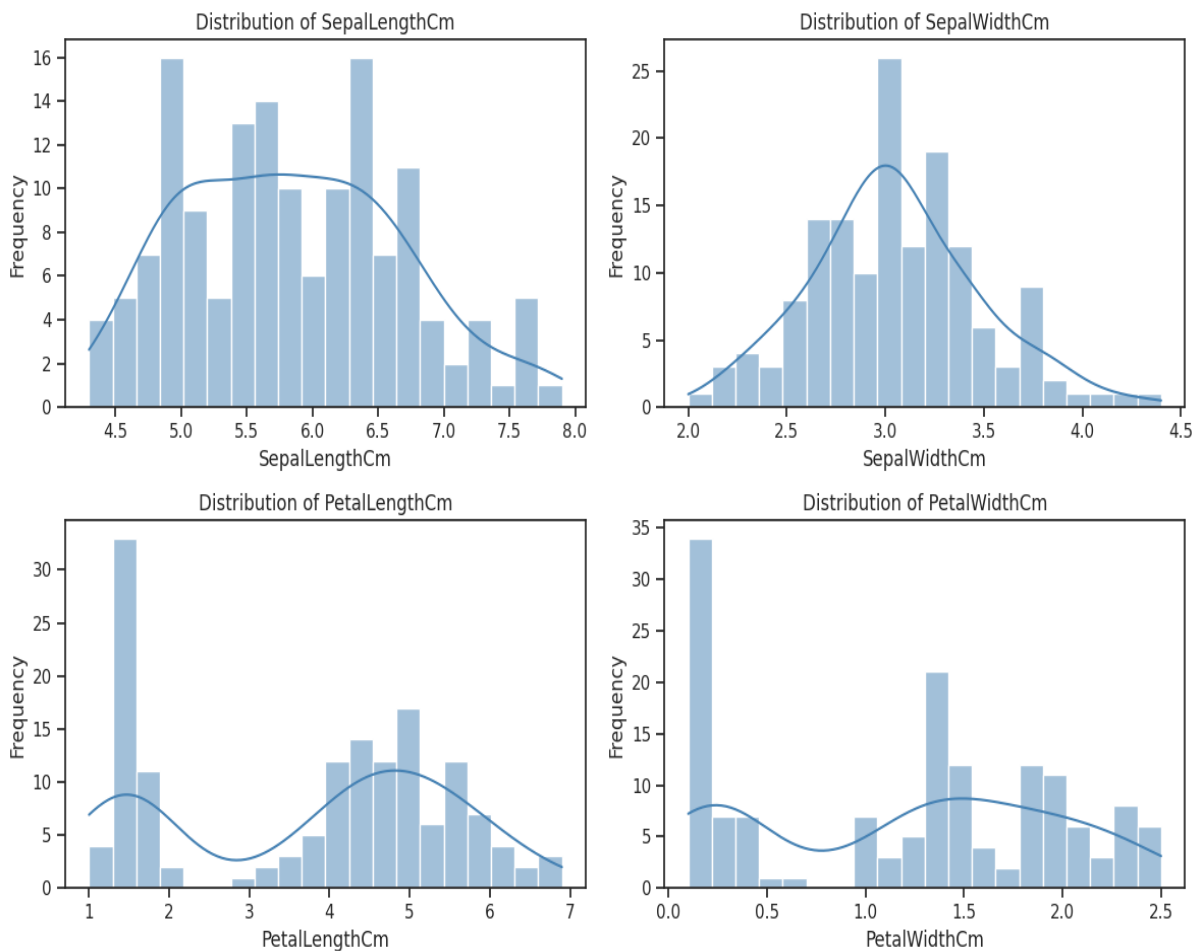
## Distribution

```python
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt

# Load dataset
df = pd.read_csv("Iris.csv")

# Drop 'Id' column
df.drop("Id", axis=1, inplace=True)

# Histograms with KDE for all numerical columns
plt.figure(figsize=(12, 8))
for i, col in enumerate(df.columns[:-1]):
    plt.subplot(2, 2, i + 1)
    sns.histplot(df[col], kde=True, bins=20, color='steelblue')
    plt.title(f'Distribution of {col}')
    plt.xlabel(col)
    plt.ylabel("Frequency")
plt.tight_layout()
plt.show()
```

Distribution of SepalLengthCm — Distribution of SepalWidthCm — Distribution of PetalLengthCm — Distribution of PetalWidthCm
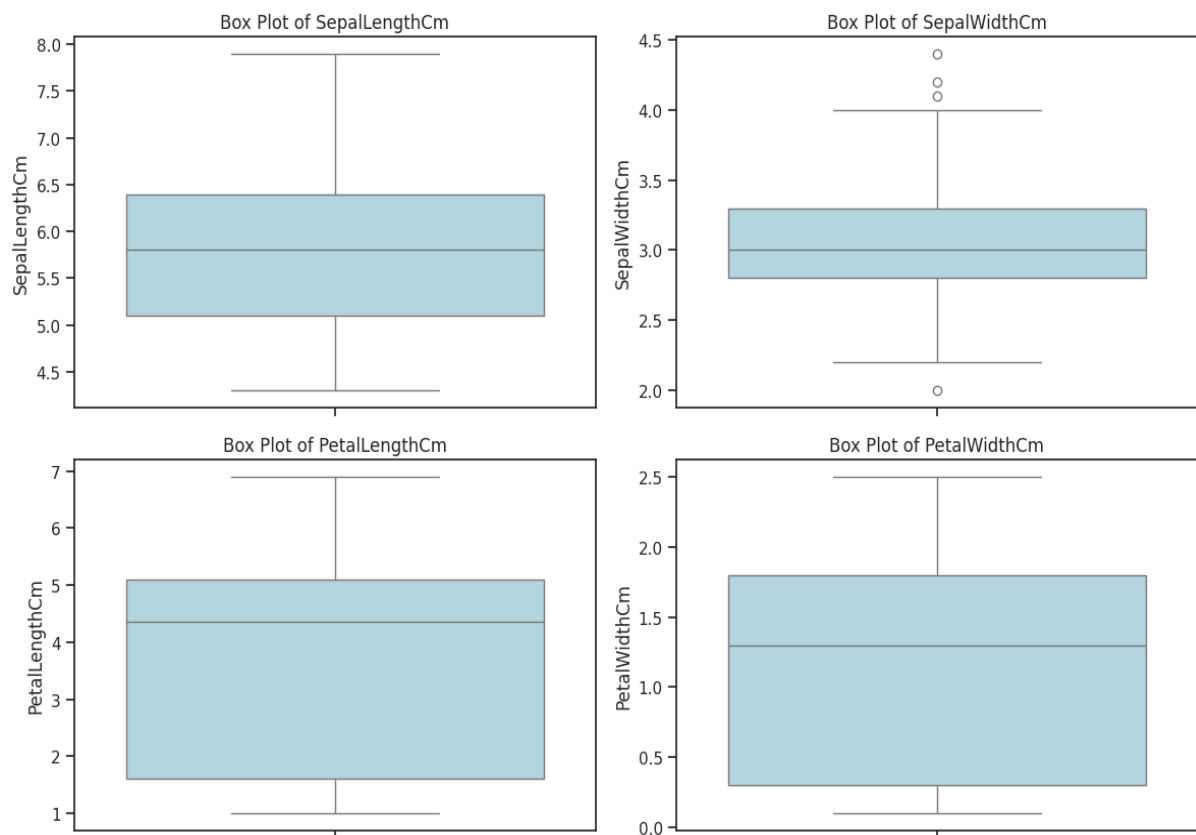
## Boxplot

```python
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt

# Load the dataset
df = pd.read_csv("Iris.csv")

# Drop 'Id' column
df.drop("Id", axis=1, inplace=True)

# Box Plots for Each Numerical Feature
plt.figure(figsize=(12, 8))
for i, col in enumerate(df.columns[:-1]):  # Exclude 'Species'
    plt.subplot(2, 2, i + 1)
    sns.boxplot(y=col, data=df, color='lightblue')
    plt.title(f'Box Plot of {col}')
    plt.ylabel(col)
plt.tight_layout()
plt.show()
```

### Pairplot of iris features

```python
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt

# Load the dataset
df = pd.read_csv("Iris.csv")

# Drop 'Id' column
df.drop("Id", axis=1, inplace=True)

# Pairplot of Features Colored by Species
sns.set(style="ticks")

# Create pairplot
pair = sns.pairplot(df, hue="Species", diag_kind="kde", corner=True, palette="Set2")

# Adjust the layout and show
pair.fig.suptitle("Pairplot of Iris Features", y=1.02)
plt.show()
```
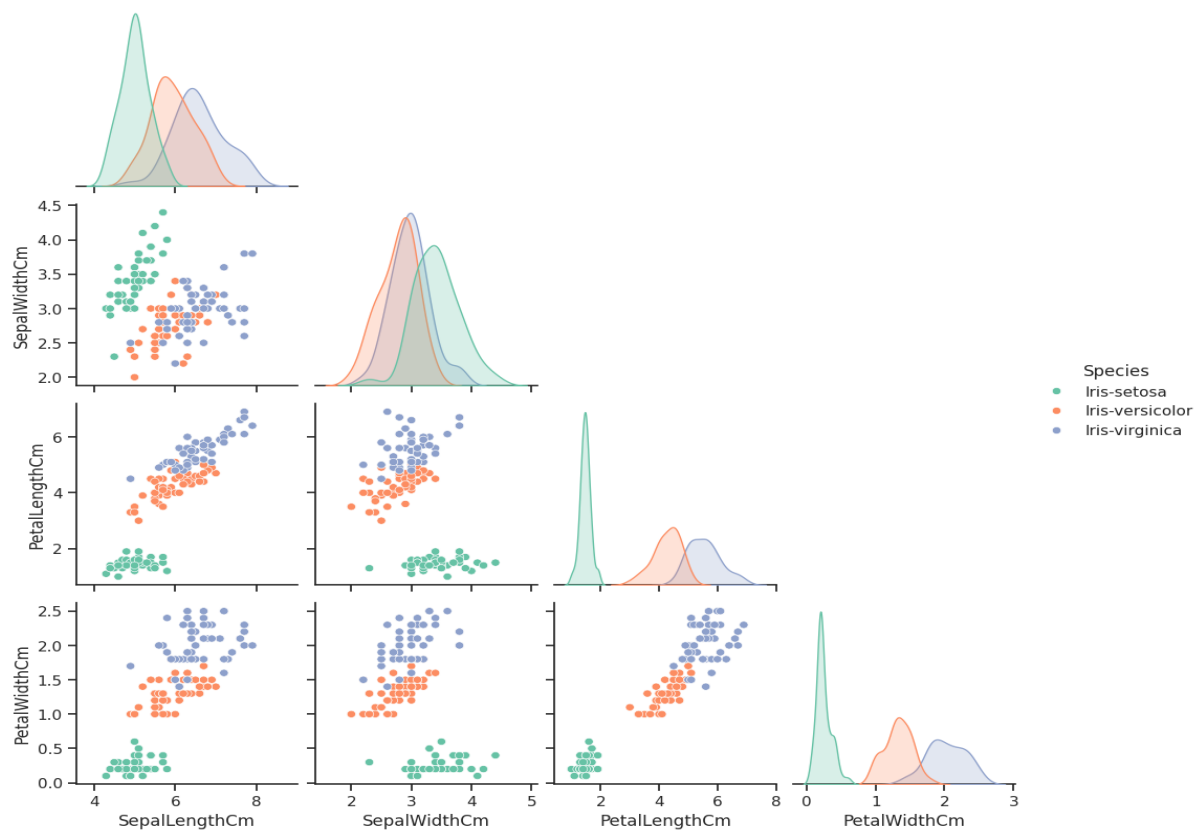
33

Pairplot of Iris Features



## Feature Correlation Heatmap

```python
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt

# Load dataset
df = pd.read_csv("Iris.csv")

# Drop 'Id' column
df.drop("Id", axis=1, inplace=True)

# Compute Correlation Matrix
corr_matrix = df.drop("Species", axis=1).corr()  # Only numerical features

print("Correlation Matrix:\n")
print(corr_matrix)

# Plot Correlation Heatmap
plt.figure(figsize=(8, 6))
sns.heatmap(corr_matrix, annot=True, cmap="coolwarm", fmt=".2f", square=True, linewidths=0.5)
plt.title("Feature Correlation Heatmap")
plt.show()
```
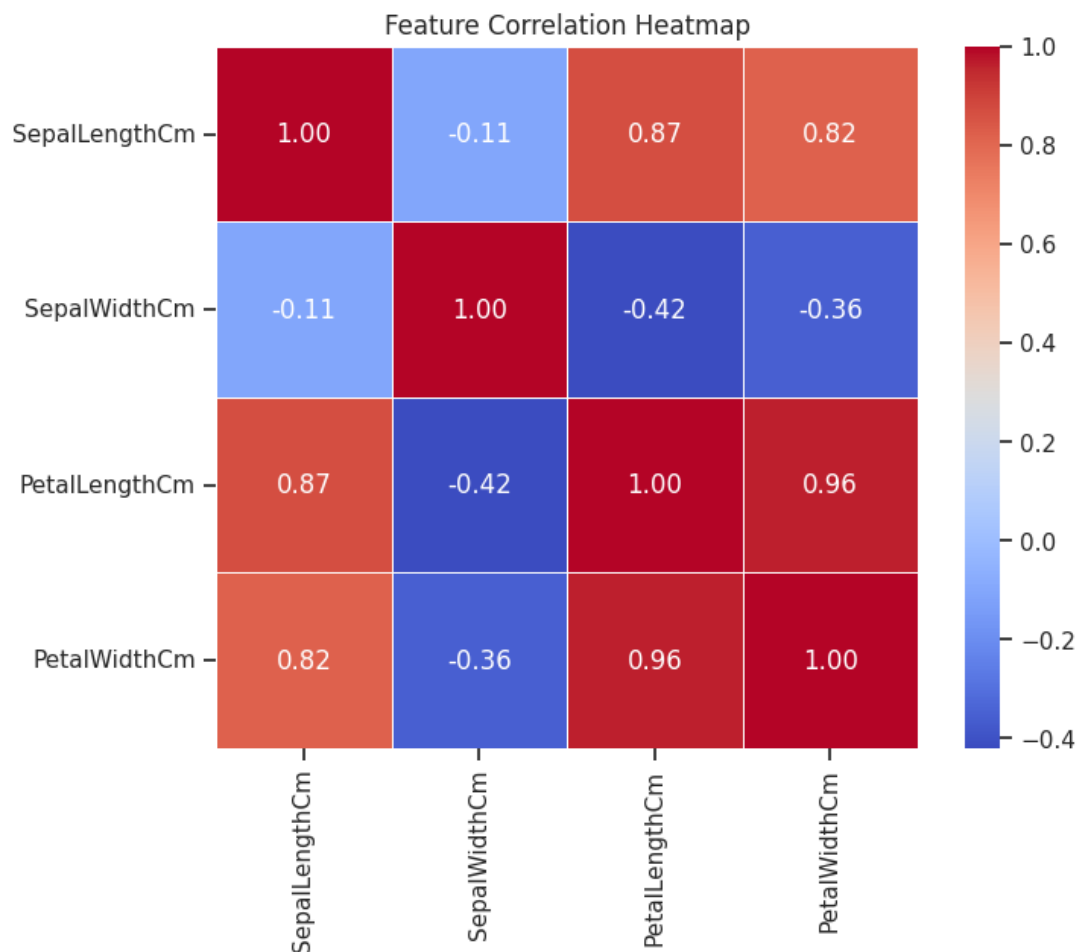
Correlation Matrix:

|  | SepalLengthCm | SepalWidthCm | PetalLengthCm | PetalWidthCm |
|---|---|---|---|---|
| SepalLengthCm | 1.000000 | -0.109369 | 0.871754 | 0.817954 |
| SepalWidthCm | -0.109369 | 1.000000 | -0.420516 | -0.356544 |
| PetalLengthCm | 0.871754 | -0.420516 | 1.000000 | 0.962757 |
| PetalWidthCm | 0.817954 | -0.356544 | 0.962757 | 1.000000 |

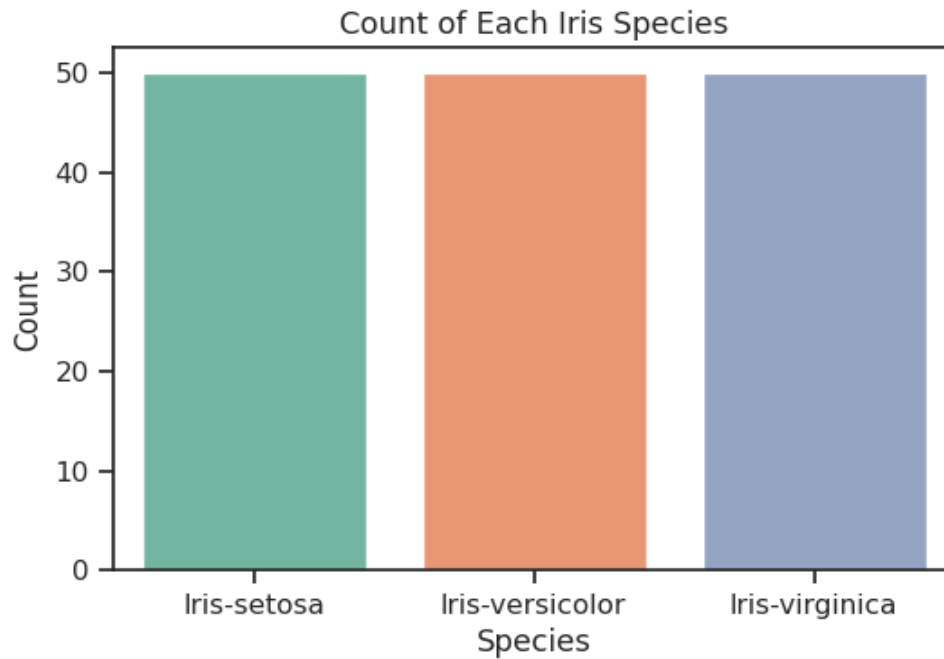Feature Correlation Heatmap

## Count Plot for Species

```python
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt

# Load dataset
df = pd.read_csv("Iris.csv")

# Drop 'Id' column
df.drop("Id", axis=1, inplace=True)

# Count Plot for 'Species'
plt.figure(figsize=(6, 4))
sns.countplot(x='Species', data=df, palette='Set2')
plt.title('Count of Each Iris Species')
plt.xlabel('Species')
plt.ylabel('Count')
plt.show()
```

Count of Each Iris Species

```python
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
from sklearn.svm import SVC
from sklearn.metrics import accuracy_score, classification_report, confusion_matrix

# Load data
df = pd.read_csv("Iris.csv")
df.drop("Id", axis=1, inplace=True)

# Features and target
X = df.drop("Species", axis=1)
y = df["Species"]

# Train-test split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Feature scaling
scaler = StandardScaler()
X_train = scaler.fit_transform(X_train)
X_test = scaler.transform(X_test)

# Train SVM classifier
model = SVC(kernel='rbf', C=1.0, gamma='scale')  # RBF kernel often best
model.fit(X_train, y_train)

# Predict
y_pred = model.predict(X_test)

# Accuracy
accuracy = accuracy_score(y_test, y_pred)
print(f"Accuracy: {accuracy:.2%}")
print("\nClassification Report:\n", classification_report(y_test, y_pred))
print("Confusion Matrix:\n", confusion_matrix(y_test, y_pred))
```

Accuracy: 100.00%

```
Classification Report:
                precision    recall  f1-score   support

    Iris-setosa       1.00      1.00      1.00        10
Iris-versicolor       1.00      1.00      1.00         9
 Iris-virginica       1.00      1.00      1.00        11

       accuracy                           1.00        30
      macro avg       1.00      1.00      1.00        30
   weighted avg       1.00      1.00      1.00        30

Confusion Matrix:
 [[10  0  0]
 [ 0  9  0]
 [ 0  0 11]]
```