

Data Science - Project 1 Documentation

Jiten Nilawar

January 20, 2024

Introduction

This document provides documentation for the data science task performed on the Iris dataset. The task includes Exploratory Data Analysis (EDA) and the implementation of a basic machine learning model for flower species prediction.

Exploratory Data Analysis (EDA)

Data Loading

The Iris dataset was loaded from the provided CSV file using the pandas library.

Summary Statistics and Visualization

Summary statistics were generated to understand the central tendency and dispersion of the data. Visualizations, such as histograms, box plots, and pair plots, were created to explore the distribution of each feature and relationships between features.

```
1 # Import necessary libraries
2 from sklearn.model_selection import train_test_split
3 from sklearn.tree import DecisionTreeClassifier
4 from sklearn.metrics import accuracy_score, precision_score, recall_score,
   classification_report, confusion_matrix
5
6 # Load the Iris dataset
7 url = "https://archive.ics.uci.edu/ml/machine-learning-databases/iris/iris.data"
8 column_names = ["sepal_length", "sepal_width", "petal_length", "petal_width", "
   class"]
9 iris_df = pd.read_csv(url, header=None, names=column_names)
10
11 # Split the dataset into features (X) and target variable (y)
12 X = iris_df.drop("class", axis=1)
13 y = iris_df["class"]
14
15 # Split the dataset into training and testing sets (80% training, 20% testing)
16 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
   random_state=42)
17
18 # Initialize the Decision Tree classifier
19 clf = DecisionTreeClassifier(random_state=42)
20
21 # Train the model on the training set
22 clf.fit(X_train, y_train)
23
24 # Make predictions on the testing set
25 y_pred = clf.predict(X_test)
26
27 # Evaluate the model
28 accuracy = accuracy_score(y_test, y_pred)
29 precision = precision_score(y_test, y_pred, average="weighted")
```

```

30 recall = recall_score(y_test, y_pred, average="weighted")
31
32 # Display evaluation metrics
33 print("Accuracy:", accuracy)
34 print("Precision:", precision)
35 print("Recall:", recall)
36
37 # Display classification report and confusion matrix
38 print("\nClassification Report:")
39 print(classification_report(y_test, y_pred))
40
41 print("\nConfusion Matrix:")
42 print(confusion_matrix(y_test, y_pred))

```

Listing 1: EDA Code

Data Science Task

Algorithm Choice

For the predictive modeling task, a Decision Tree classifier was chosen due to its simplicity and effectiveness for classification tasks, especially when dealing with small to medium-sized datasets.

Feature Selection

All four features (sepal length, sepal width, petal length, and petal width) were used for training the model, as they are essential for distinguishing between different iris species.

Training and Evaluation

The dataset was split into training and testing sets. The Decision Tree classifier was trained on the training set and evaluated on the testing set. Evaluation metrics included accuracy, precision, recall, classification report, and confusion matrix.

```

1 # Import necessary libraries
2 import seaborn as sns
3 import matplotlib.pyplot as plt
4
5 # Load the Iris dataset
6 url = "https://archive.ics.uci.edu/ml/machine-learning-databases/iris/iris.data"
7 column_names = ["sepal_length", "sepal_width", "petal_length", "petal_width", "
8                 class"]
9 iris_df = pd.read_csv(url, header=None, names=column_names)
10
11 # Set the style for seaborn plots
12 sns.set(style="whitegrid")
13
14 # Distribution of each feature using histograms
15 plt.figure(figsize=(12, 8))
16 for i, column in enumerate(iris_df.columns[:-1]):
17     plt.subplot(2, 2, i+1)
18     sns.histplot(iris_df[column], kde=True, bins=20, color='skyblue')
19     plt.title(f"{column} Distribution")
20
21 plt.tight_layout()
22 plt.show()
23
24 # Box plots for each feature
25 plt.figure(figsize=(12, 8))
26 for i, column in enumerate(iris_df.columns[:-1]):
27     plt.subplot(2, 2, i+1)
28     sns.boxplot(x="class", y=column, data=iris_df)

```

```

28     plt.title(f"{column} Distribution by Class")
29
30 plt.tight_layout()
31 plt.show()
32
33 # Pairplot to visualize relationships between features
34 sns.pairplot(iris_df, hue="class", markers=["o", "s", "D"])
35 plt.show()

```

Listing 2: Data Science Task Code

Challenges Faced

During the EDA, ensuring proper visualization and interpretation of pair plots for relationships between features can be challenging, especially in larger datasets. In the data science task, tuning hyperparameters of the model for better performance might be required in more complex datasets.

Conclusion

In conclusion, this data science project effectively explored the Iris dataset through Exploratory Data Analysis (EDA) and implemented a Decision Tree classifier for flower species prediction. The chosen features and algorithm demonstrated satisfactory performance, as indicated by accuracy, precision, and recall metrics. The project provides valuable insights into dataset characteristics and lays the groundwork for further applications in machine learning and data science.