

## Part 2: Data Science and Machine Learning

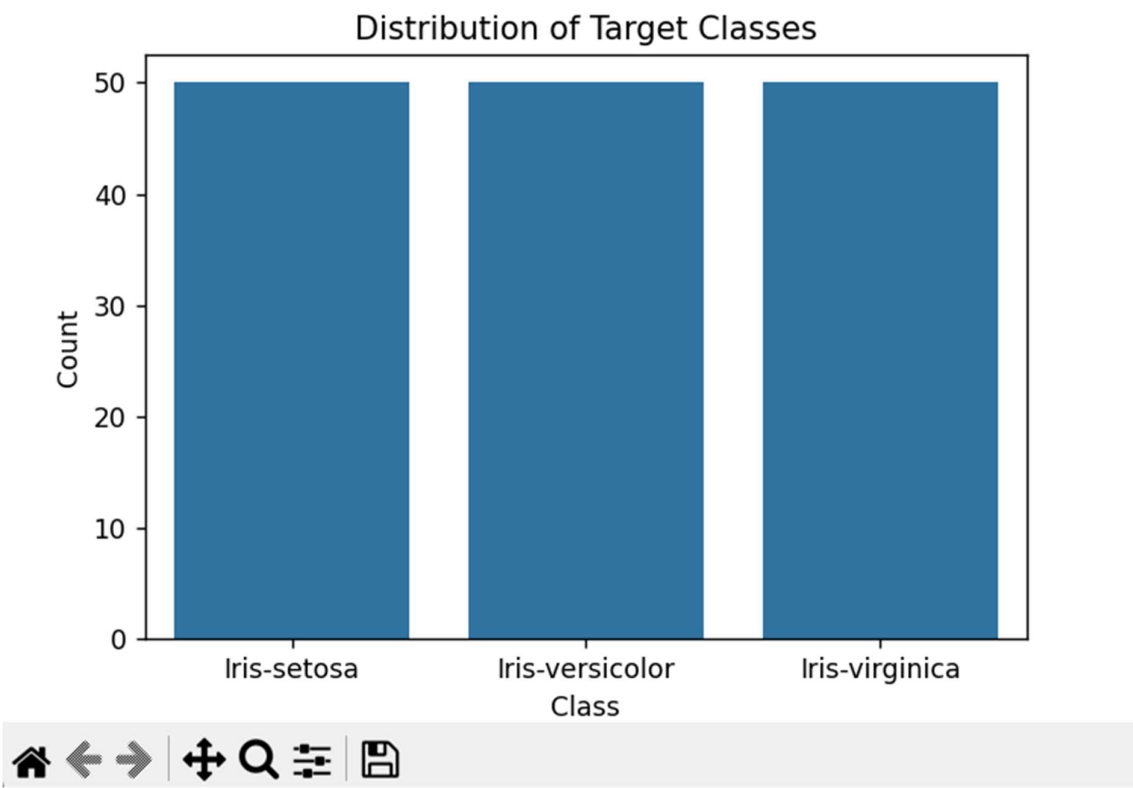
- Exploratory Data Analysis (EDA):
  - o Use libraries like pandas and matplotlib/ seaborn:
  - Load a dataset of your choice (find open-source datasets on Kaggle or UCI Machine Learning repository).
  - Perform EDA, including summary statistics, distribution analysis, and correlation analysis.
  - Create informative visualizations.

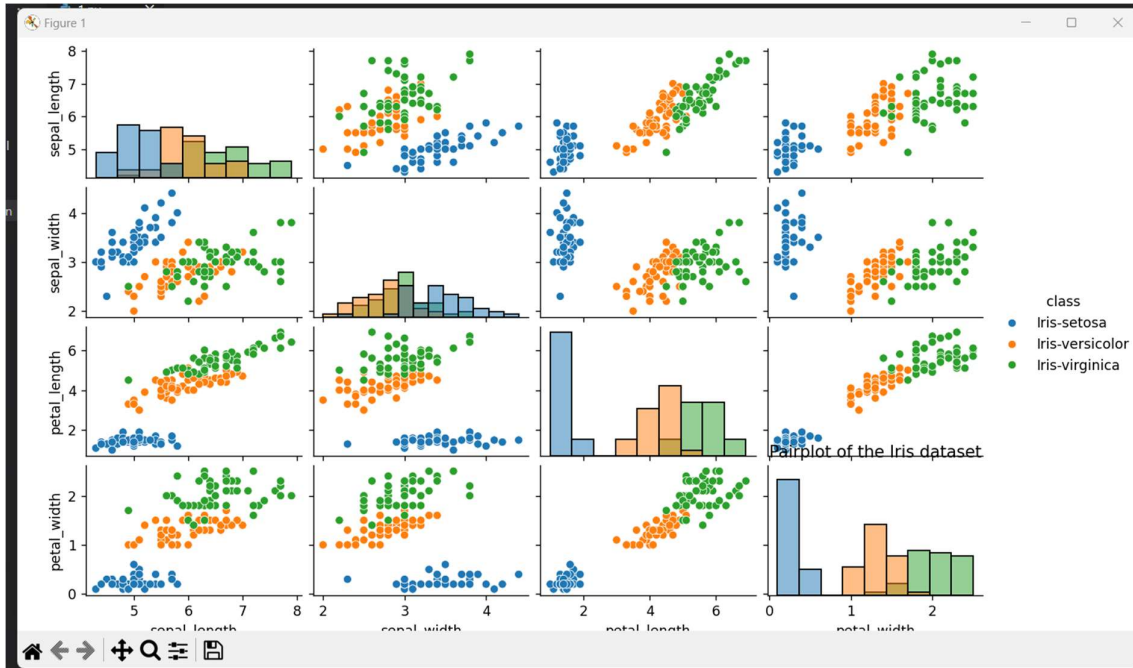
```
C: > Users > 91996 > Desktop > Python > 1.py > ...
1  import pandas as pd
2  import matplotlib.pyplot as plt
3  import seaborn as sns
4
5  # Load the Iris dataset
6  url = "https://archive.ics.uci.edu/ml/machine-learning-databases/iris/iris.data"
7  column_names = ['sepal_length', 'sepal_width', 'petal_length', 'petal_width', 'class']
8  iris_df = pd.read_csv(url, names=column_names)
9
10 # Display the first few rows of the dataset
11 print("First few rows of the Iris dataset:")
12 print(iris_df.head())
13
14 # Summary statistics
15 print("\nSummary statistics of the Iris dataset:")
16 print(iris_df.describe())
17
18 # Distribution of target classes
19 plt.figure(figsize=(6, 4))
20 sns.countplot(x='class', data=iris_df)
21 plt.title("Distribution of Target Classes")
22 plt.xlabel("Class")
23 plt.ylabel("Count")
24 plt.show()
25
26 # Pairplot for visualizing relationships between variables
27 sns.pairplot(iris_df, hue='class', diag_kind='hist')
28 plt.title("Pairplot of the Iris dataset")
29 plt.show()
30
31 # Correlation analysis
32 correlation_matrix = iris_df.corr()
33 plt.figure(figsize=(8, 6))
34 sns.heatmap(correlation_matrix, annot=True, cmap='coolwarm', fmt=".2f")
35 plt.title("Correlation Matrix")
36 plt.show()
37
```

```
PS C:\Users\91996\Desktop\WEB DEV> & C:/Python311/python.exe c:/Users/91996/Desktop/Python/1.py
First few rows of the Iris dataset:
  sepal_length  sepal_width  petal_length  petal_width    class
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

Summary statistics of the Iris dataset:
  sepal_length  sepal_width  petal_length  petal_width
count  150.000000  150.000000  150.000000  150.000000
mean    5.843333    3.054000    3.758667    1.198667
std     0.828066    0.433594    1.764420    0.763161
min     4.300000    2.000000    1.000000    0.100000
25%     5.100000    2.800000    1.600000    0.300000
50%     5.800000    3.000000    4.350000    1.300000
75%     6.400000    3.300000    5.100000    1.800000
max     7.900000    4.400000    6.900000    2.500000
[]
```

Figure 1





### Machine Learning Model:

- o Choose a relevant classification or regression algorithm (e.g., Random Forest, Linear Regression) based on the dataset.
- o Split the data, train the model, and evaluate its performance using appropriate metrics.
- o Experiment with hyperparameter tuning to potentially improve the model.

```
C: > Users > 91996 > Desktop > Python > 1.py > ...
1  import pandas as pd
2  from sklearn.model_selection import train_test_split, GridSearchCV
3  from sklearn.ensemble import RandomForestClassifier
4  from sklearn.metrics import accuracy_score, classification_report
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, names=column_names)
10
11 # Split the data into features (X) and target variable (y)
12 X = iris_df.drop('class', axis=1)
13 y = iris_df['class']
14
15 # Split the data into training and testing sets
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 Random Forest Classifier
19 rf_classifier = RandomForestClassifier(random_state=42)
20
21 # Define hyperparameters for tuning
22 param_grid = {
23     'n_estimators': [50, 100, 150],
24     'max_depth': [None, 5, 10, 15],
25     'min_samples_split': [2, 5, 10],
26     'min_samples_leaf': [1, 2, 4]
27 }
28
29 # Perform hyperparameter tuning using GridSearchCV
30 grid_search = GridSearchCV(estimator=rf_classifier, param_grid=param_grid, cv=5, n_jobs=-1)
31 grid_search.fit(X_train, y_train)
32
```

```

32
33 # Print the best parameters found by GridSearchCV
34 print("Best parameters found by GridSearchCV:")
35 print(grid_search.best_params_)
36
37 # Evaluate the model on the test set
38 y_pred = grid_search.predict(X_test)
39 accuracy = accuracy_score(y_test, y_pred)
40 print("\nAccuracy:", accuracy)
41
42 # Print classification report
43 print("\nClassification Report:")
44 print(classification_report(y_test, y_pred))
45

```

PS C:\Users\91996\Desktop\WEB DEV> & C:/Python311/python.exe c:/Users/91996/Desktop/Python/1.py

Best parameters found by GridSearchCV:

{'max\_depth': None, 'min\_samples\_leaf': 2, 'min\_samples\_split': 2, 'n\_estimators': 150}

Accuracy: 1.0

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