```
1  Name: Mustafa Aydoğan
2  ID: 191101002
3  Course: BIL570 / BIL470
```

# Exploratory Data Analysis (EDA)

In [1]:

```python
 1  import pandas as pd
 2  from sklearn.utils import shuffle
 3  import seaborn as sns
 4  import matplotlib.pyplot as plt
 5  import numpy as np
 6
 7  from sklearn.metrics import confusion_matrix, classification_report, roc_curve, roc_
 8  from sklearn.metrics import accuracy_score, precision_score
 9
10  from dt import DecisionTreeClassifier
```

In [2]:

```python
 1  # Veri setini yükleme
 2  data = pd.read_csv('iris.csv')
 3
 4  # Veri setinin başlığını kontrol etme
 5  print(data.head())
```

```
   Id  SepalLengthCm  SepalWidthCm  PetalLengthCm  PetalWidthCm      Speci
es
0  1            5.1           3.5            1.4           0.2  Iris-seto
sa
1  2            4.9           3.0            1.4           0.2  Iris-seto
sa
2  3            4.7           3.2            1.3           0.2  Iris-seto
sa
3  4            4.6           3.1            1.5           0.2  Iris-seto
sa
4  5            5.0           3.6            1.4           0.2  Iris-seto
sa
```

In [3]:

```python
 1  # Tür adlarını tam sayılara dönüştürme
 2  data['Species'] = data['Species'].map({'Iris-setosa': 0, 'Iris-versicolor': 1, 'Iris
 3
 4  # Sonuçları kontrol etme
 5  print(data.head())
```

```
   Id  SepalLengthCm  SepalWidthCm  PetalLengthCm  PetalWidthCm  Species
0  1            5.1           3.5            1.4           0.2        0
1  2            4.9           3.0            1.4           0.2        0
2  3            4.7           3.2            1.3           0.2        0
3  4            4.6           3.1            1.5           0.2        0
4  5            5.0           3.6            1.4           0.2        0
```

```
1  # Veri setinin özeti
2  print(data.describe())
```

|       | Id | SepalLengthCm | SepalWidthCm | PetalLengthCm | PetalWidthCm \ |
|-------|-----|-----|-----|-----|-----|
| count | 150.000000 | 150.000000 | 150.000000 | 150.000000 | 150.000000 |
| mean | 75.500000 | 5.843333 | 3.054000 | 3.758667 | 1.198667 |
| std | 43.445368 | 0.828066 | 0.433594 | 1.764420 | 0.763161 |
| min | 1.000000 | 4.300000 | 2.000000 | 1.000000 | 0.100000 |
| 25% | 38.250000 | 5.100000 | 2.800000 | 1.600000 | 0.300000 |
| 50% | 75.500000 | 5.800000 | 3.000000 | 4.350000 | 1.300000 |
| 75% | 112.750000 | 6.400000 | 3.300000 | 5.100000 | 1.800000 |
| max | 150.000000 | 7.900000 | 4.400000 | 6.900000 | 2.500000 |

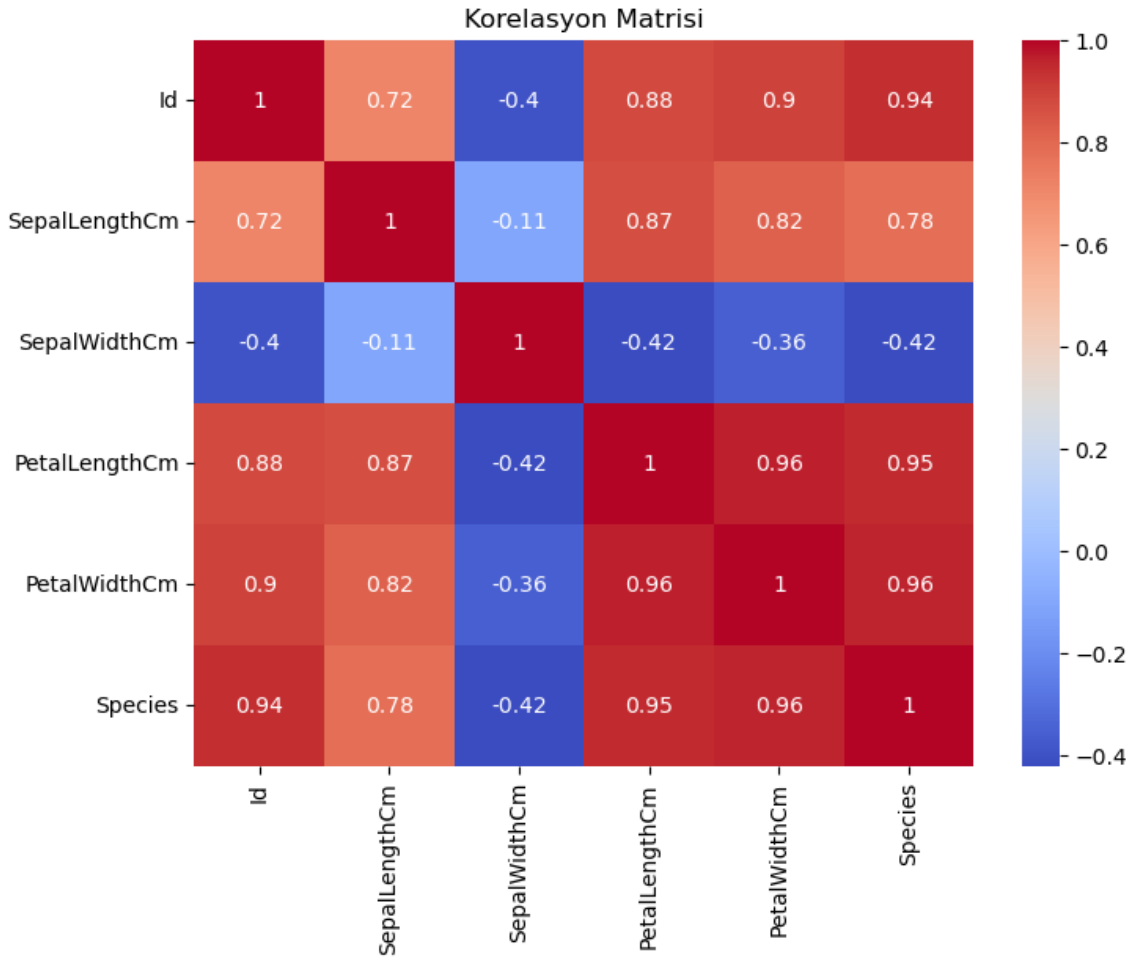|       | Species |
|-------|---------|
| count | 150.000000 |
| mean | 1.000000 |
| std | 0.819232 |
| min | 0.000000 |
| 25% | 0.000000 |
| 50% | 1.000000 |
| 75% | 2.000000 |
| max | 2.000000 |

```
1  # Veri setinin karıştırılması
2  data = shuffle(data)
3  print(data.head())
```

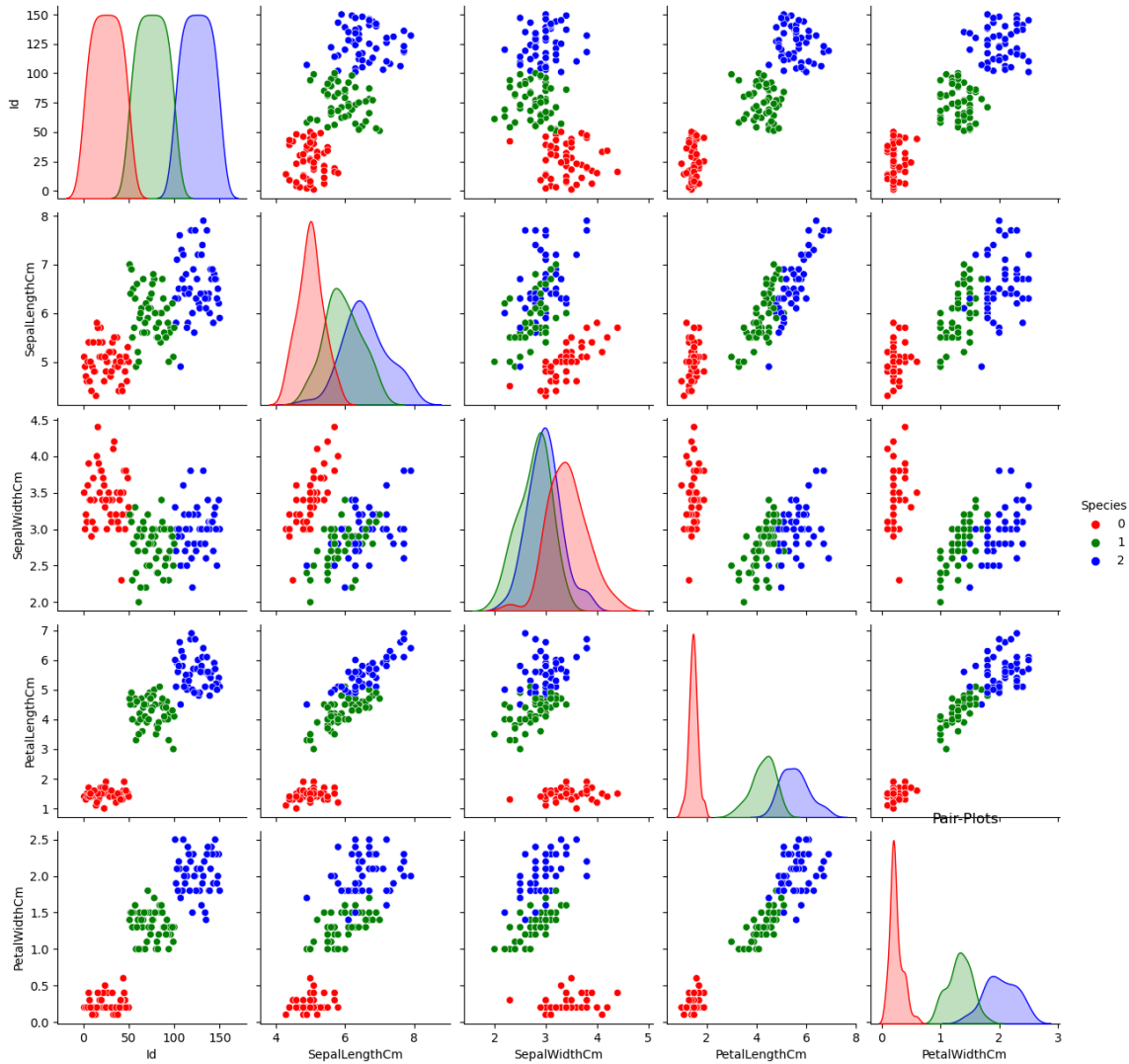|     | Id | SepalLengthCm | SepalWidthCm | PetalLengthCm | PetalWidthCm | Species |
|-----|-----|-----|-----|-----|-----|-----|
| 107 | 108 | 7.3 | 2.9 | 6.3 | 1.8 | 2 |
| 101 | 102 | 5.8 | 2.7 | 5.1 | 1.9 | 2 |
| 59 | 60 | 5.2 | 2.7 | 3.9 | 1.4 | 1 |
| 126 | 127 | 6.2 | 2.8 | 4.8 | 1.8 | 2 |
| 44 | 45 | 5.1 | 3.8 | 1.9 | 0.4 | 0 |

```python
# Korelasyon matrisi
correlation_matrix = data.corr()

# Korelasyon matrisini görselleştirme
plt.figure(figsize=(8, 6))
sns.heatmap(correlation_matrix, annot=True, cmap='coolwarm')
plt.title('Korelasyon Matrisi')
plt.show()
```



Korelasyon Matrisi

```
1  # Öznitelikler için pair-plots
2  custom_palette = ['red', 'green', 'blue']
3  sns.pairplot(data, hue='Species', palette=custom_palette)
4  plt.title('Pair-Plots')
5  plt.show()
```



## Veriyi eğitim ve test setlerine ayırma

In [8]:

```
1  X = data.iloc[:, :-1].values
2  y = data.iloc[:, -1].values
3
4  train_size = int(0.8 * len(X))
5  X_train, X_test = X[:train_size], X[train_size:]
6  y_train, y_test = y[:train_size], y[train_size:]
```

### DecisionTreeClassifier'ı eğitme

In [9]:

```python
classifier = DecisionTreeClassifier(max_depth=5)
classifier.fit(X_train.tolist(), y_train.tolist())
```

### Test seti üzerinde sınıflandırma

In [10]:

```python
predictions = classifier.predict(X_test.tolist())
```
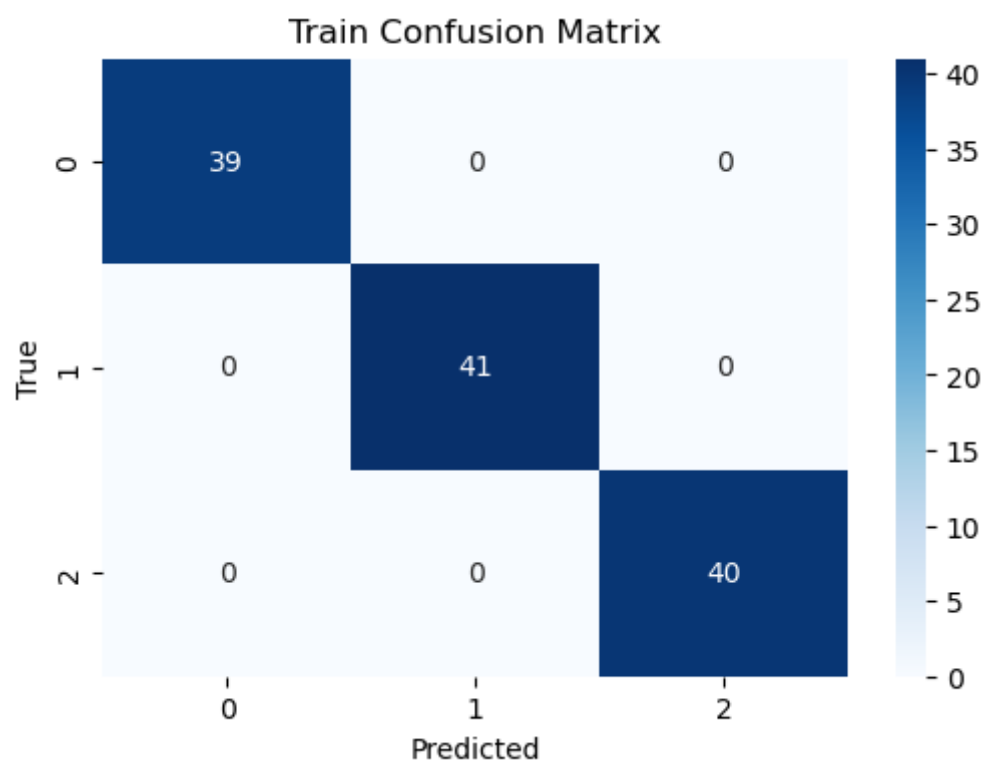
# Sonuçları değerlendirme

In [11]:

```python
accuracy = sum(predictions == y_test) / len(y_test)
print("Accuracy:", accuracy)
```

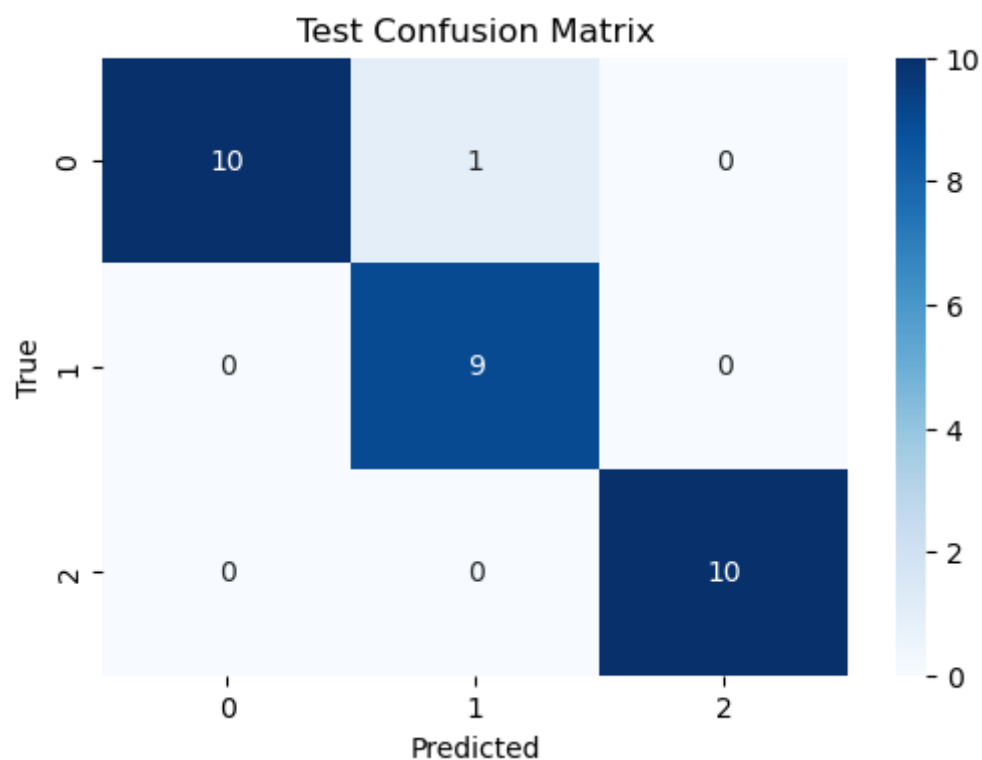Accuracy: 0.9666666666666667

# Confusion Matrix

In [12]:

```python
# Training set evaluation
train_predictions = classifier.predict(X_train.tolist())
train_confusion_matrix = confusion_matrix(y_train.tolist(), train_predictions)
plt.figure(figsize=(6, 4))
plt.title('Train Confusion Matrix')
sns.heatmap(train_confusion_matrix, annot=True, cmap='Blues', fmt='g')
plt.xlabel('Predicted')
plt.ylabel('True')
plt.show()
```

In [13]:

```python
# Test set evaluation
test_predictions = classifier.predict(X_test.tolist())
test_confusion_matrix = confusion_matrix(y_test.tolist(), test_predictions)
plt.figure(figsize=(6, 4))
plt.title('Test Confusion Matrix')
sns.heatmap(test_confusion_matrix, annot=True, cmap='Blues', fmt='g')
plt.xlabel('Predicted')
plt.ylabel('True')
plt.show()
```

In [14]:
```python
# Accuracy
train_accuracy = np.mean(train_predictions == y_train) * 100
test_accuracy = np.mean(test_predictions == y_test) * 100
print("\nAccuracy (Train set):", train_accuracy)
print("Accuracy (Test set):", test_accuracy)

# F1-Score
train_f1_score = classification_report(y_train, train_predictions, output_dict=True)
test_f1_score = classification_report(y_test, test_predictions, output_dict=True)['w
print("\nF1-Score (Train set):", train_f1_score)
print("F1-Score (Test set):", test_f1_score)

# Precision
train_precision = classification_report(y_train, train_predictions, output_dict=True
test_precision = classification_report(y_test, test_predictions, output_dict=True)['
print("\nPrecision (Train set):", train_precision)
print("Precision (Test set):", test_precision)

# Recall
train_recall = classification_report(y_train, train_predictions, output_dict=True)['
test_recall = classification_report(y_test, test_predictions, output_dict=True)['wei
print("\nRecall (Train set):", train_recall)
print("Recall (Test set):", test_recall)
```

```
Accuracy (Train set): 100.0
Accuracy (Test set): 96.66666666666667

F1-Score (Train set): 1.0
F1-Score (Test set): 0.966750208855472

Precision (Train set): 1.0
Precision (Test set): 0.9700000000000001

Recall (Train set): 1.0
Recall (Test set): 0.9666666666666667
```
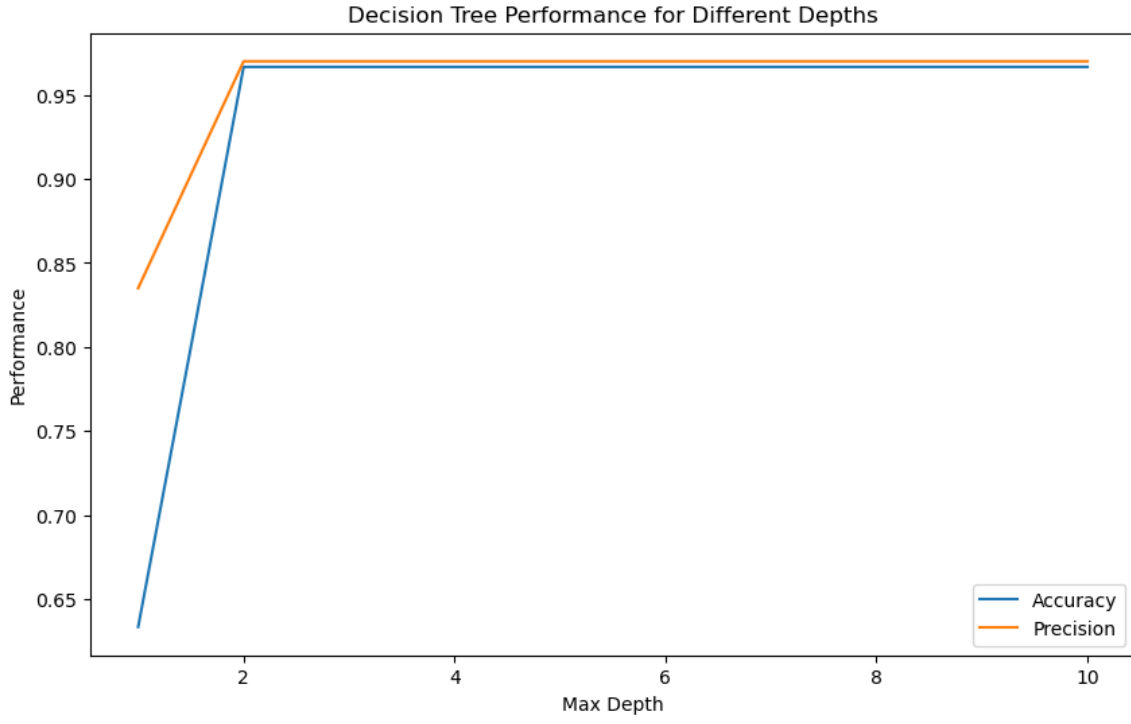
```
In [15]:
 1  L = [1, 2, 3, 4, 5, 6, 7, 8, 9, 10]
 2  accuracies = []
 3  precisions = []
 4
 5  for depth in L:
 6      # DecisionTreeClassifier'ı eğitme
 7      classifier = DecisionTreeClassifier(max_depth=depth)
 8      classifier.fit(X_train.tolist(), y_train.tolist())
 9
10      # Test seti üzerinde sınıflandırma
11      predictions = classifier.predict(X_test.tolist())
12
13      # Accuracy
14      accuracy = accuracy_score(y_test, predictions)
15      accuracies.append(accuracy)
16
17      # Precision
18      precision = precision_score(y_test, predictions, average='weighted', zero_divisi
19      precisions.append(precision)
20
21  plt.figure(figsize=(10, 6))
22  plt.plot(L, accuracies, label='Accuracy')
23  plt.plot(L, precisions, label='Precision')
24  plt.xlabel('Max Depth')
25  plt.ylabel('Performance')
26  plt.title('Decision Tree Performance for Different Depths')
27  plt.legend()
28  plt.show()
29
30  best_accuracy = max(accuracies)
31  best_accuracy_index = accuracies.index(best_accuracy)
32  best_precision = precisions[best_accuracy_index]
33  best_depth = L[best_accuracy_index]
34
35  print("Best Depth:", best_depth)
36  print("Best Accuracy:", best_accuracy)
37  print("Corresponding Precision:", best_precision)
```

Decision Tree Performance for Different Depths

Best Depth: 2
Best Accuracy: 0.9666666666666667
Corresponding Precision: 0.9700000000000001

| 1 | Bu verilere göre, Decision Tree Classifier modeli 2 derinlik seviyesinde en iyi sonuçları verdi. Yani, modelin en doğru tahminleri ve hassas sonuçları bu derinlik seviyesinde elde edildi. |
| 2 | |
| 3 | Modelin daha derin olması ise overfit riskini artırabilir ve genel başarıyı olumsuz etkileyebilir. Dolayısıyla, 2 derinlik seviyesi, modelin dengeli ve iyi bir şekilde genellemesini sağlamıştır. |
| 4 | |
| 5 | Ancak, daha fazla deney yapmak ve farklı veri setlerinde modelin performansını test etmek için daha fazla çalışma gerekmektedir. |