In [27]: 
```python
1  !pip install graphviz
```

Requirement already satisfied: graphviz in /Users/ghitabouzida/opt/anacon
da3/lib/python3.9/site-packages (0.20.1)

In [30]: 
```python
 1  import pandas as pd
 2  # Assuming the data is stored in a CSV file called 'titanic_data.csv'
 3  data = pd.read_csv('titanic_dataset.csv')
 4  # Let's first check for missing data
 5  print(data.isnull().sum())
 6  # Drop unnecessary columns (such as PassengerId, Name, Ticket, and Cabi
 7  data = data.drop(columns=['PassengerId', 'Name', 'Ticket', 'Cabin'])
 8  # Handle missing values in Age and Embarked columns (you may choose to
 9  data['Age'].fillna(data['Age'].mean(), inplace=True)
10  data['Embarked'].fillna(data['Embarked'].mode()[0], inplace=True)
11  # Convert categorical features into numerical representations using one
12  data = pd.get_dummies(data, columns=['Sex', 'Embarked'], drop_first=Tru
13  # Separate the target variable (Survived) from the features
14  X = data.drop(columns=['Survived'])
15  y = data['Survived']
```

```
PassengerId      0
Survived         0
Pclass           0
Name             0
Sex              0
Age            177
SibSp            0
Parch            0
Ticket           0
Fare             0
Cabin          687
Embarked         2
dtype: int64
```

In [31]: 
```python
1  from sklearn.model_selection import train_test_split
2  # Split the data into training and testing sets (80% for training, 20%
3  X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2
```

```
In [48]:    1  import pandas as pd
            2  from sklearn.model_selection import train_test_split
            3  from sklearn.preprocessing import StandardScaler
            4  from sklearn.linear_model import LogisticRegression
            5  from sklearn.pipeline import make_pipeline
            6  from sklearn.metrics import roc_curve, roc_auc_score
            7  import matplotlib.pyplot as plt
            8  from sklearn.neighbors import KNeighborsClassifier
            9  from sklearn.metrics import accuracy_score, classification_report, conf
           10  import numpy as np
           11  from sklearn.model_selection import cross_val_score
           12  import matplotlib.pyplot as plt
           13  from sklearn.tree import DecisionTreeClassifier
           14  from sklearn.tree import plot_tree
           15
           16
           17  # Assuming the data is stored in a CSV file called 'titanic_data.csv'
           18  data = pd.read_csv('titanic_dataset.csv')
           19
           20  # Drop unnecessary columns (such as PassengerId, Name, Ticket, and Cabi
           21  data = data.drop(columns=['PassengerId', 'Name', 'Ticket', 'Cabin'])
           22
           23  # Handle missing values in Age and Embarked columns (you may choose to
           24  data['Age'].fillna(data['Age'].mean(), inplace=True)
           25  data['Embarked'].fillna(data['Embarked'].mode()[0], inplace=True)
           26
           27  # Convert categorical features into numerical representations using one
           28  data = pd.get_dummies(data, columns=['Sex', 'Embarked'], drop_first=Tru
           29
           30  # Separate the target variable (Survived) from the features
           31  X = data.drop(columns=['Survived'])
           32  y = data['Survived']
           33
           34  # Split the data into training and testing sets (80% for training, 20%
           35  X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2
           36
           37  # Create a pipeline with feature scaling and logistic regression
           38  pipe = make_pipeline(StandardScaler(), LogisticRegression(random_state=
           39
           40  # Fit the pipeline on the training data (feature scaling is applied)
           41  pipe.fit(X_train, y_train)
           42
           43  # Make predictions on the test set
           44  y_pred = pipe.predict(X_test)
           45
           46  # Evaluate the model
           47  from sklearn.metrics import accuracy_score, classification_report, conf
           48
           49  # Calculate the accuracy of the model
           50  accuracy = accuracy_score(y_test, y_pred)
           51  print("Accuracy:", accuracy)
           52
           53  # Get a classification report
           54  print(classification_report(y_test, y_pred))
           55
           56  # Get the confusion matrix
           57  conf_matrix = confusion_matrix(y_test, y_pred)
```

```
58  print("Confusion Matrix:")
59  print(conf_matrix)
60
```
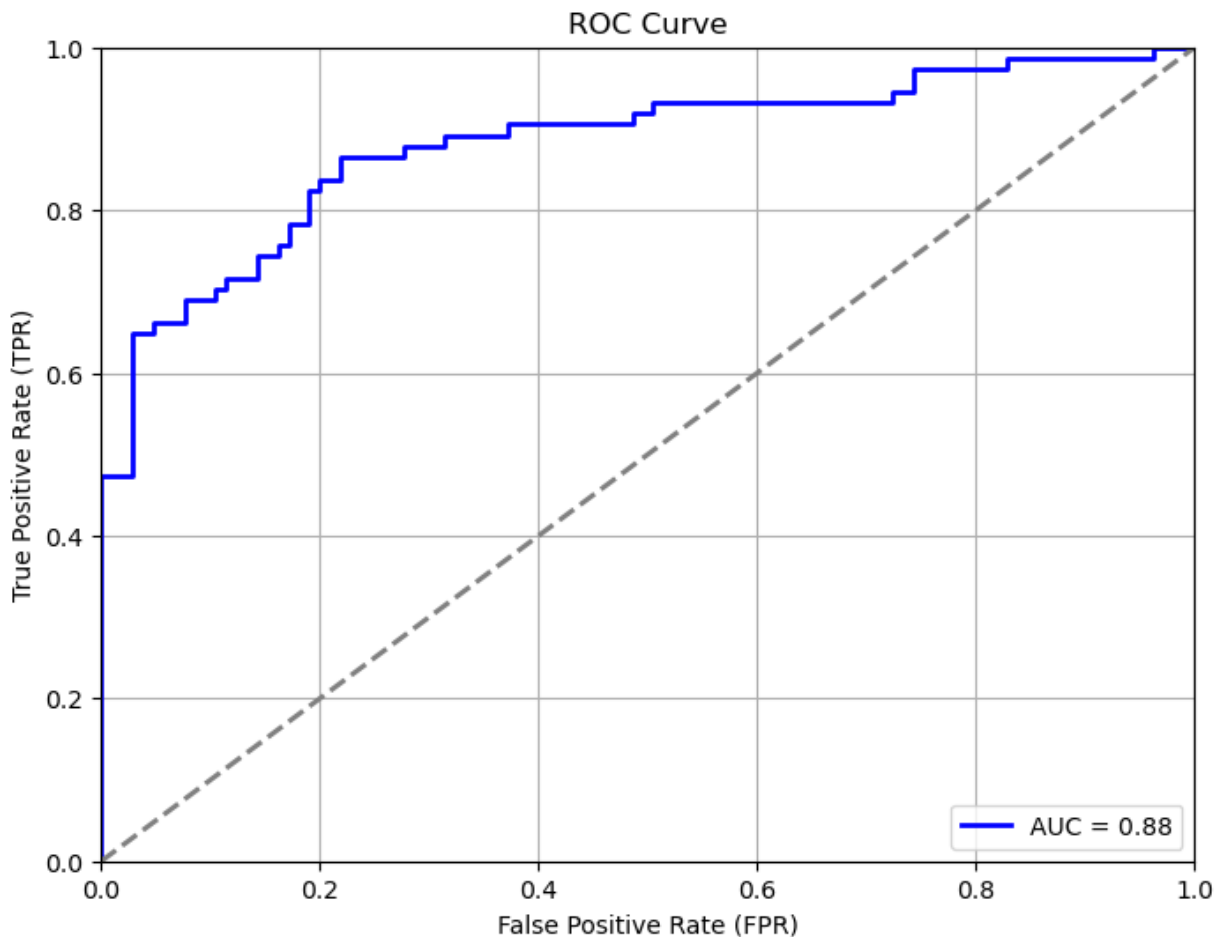
Accuracy: 0.8100558659217877

|              | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| 0            | 0.83      | 0.86   | 0.84     | 105     |
| 1            | 0.79      | 0.74   | 0.76     | 74      |
|              |           |        |          |         |
| accuracy     |           |        | 0.81     | 179     |
| macro avg    | 0.81      | 0.80   | 0.80     | 179     |
| weighted avg | 0.81      | 0.81   | 0.81     | 179     |

```
Confusion Matrix:
[[90 15]
 [19 55]]
```

In [45]:

```python
# Get the predicted probabilities for the positive class (survived) fro
y_pred_prob = pipe.predict_proba(X_test)[:, 1]

# Calculate the ROC curve values: FPR, TPR, and the corresponding thres
fpr, tpr, thresholds = roc_curve(y_test, y_pred_prob)

# Calculate the AUC score
auc_score = roc_auc_score(y_test, y_pred_prob)

# Plot the ROC curve
plt.figure(figsize=(8, 6))
plt.plot(fpr, tpr, color='blue', lw=2, label=f'AUC = {auc_score:.2f}')
plt.plot([0, 1], [0, 1], color='gray', linestyle='--', lw=2)
plt.xlim([0.0, 1.0])
plt.ylim([0.0, 1.0])
plt.xlabel('False Positive Rate (FPR)')
plt.ylabel('True Positive Rate (TPR)')
plt.title('ROC Curve')
plt.legend(loc='lower right')
plt.grid(True)
plt.show()
```

In [47]:
```python
import warnings
warnings.simplefilter(action='ignore', category=FutureWarning)

# Your KNN code here...
# Create a KNN classifier
knn_model = KNeighborsClassifier()

# Fit the KNN model to the training data
knn_model.fit(X_train, y_train)

# Make predictions on the test set
y_pred_knn = knn_model.predict(X_test)

# Evaluate the KNN model
accuracy_knn = accuracy_score(y_test, y_pred_knn)
print("KNN Accuracy:", accuracy_knn)

# Get a classification report for KNN
print("Classification Report (KNN):")
print(classification_report(y_test, y_pred_knn))

# Get the confusion matrix for KNN
conf_matrix_knn = confusion_matrix(y_test, y_pred_knn)
print("Confusion Matrix (KNN):")
print(conf_matrix_knn)
```

```
KNN Accuracy: 0.6927374301675978
Classification Report (KNN):
              precision    recall  f1-score   support

           0       0.71      0.80      0.75       105
           1       0.66      0.54      0.59        74

    accuracy                           0.69       179
   macro avg       0.68      0.67      0.67       179
weighted avg       0.69      0.69      0.69       179

Confusion Matrix (KNN):
[[84 21]
 [34 40]]
```
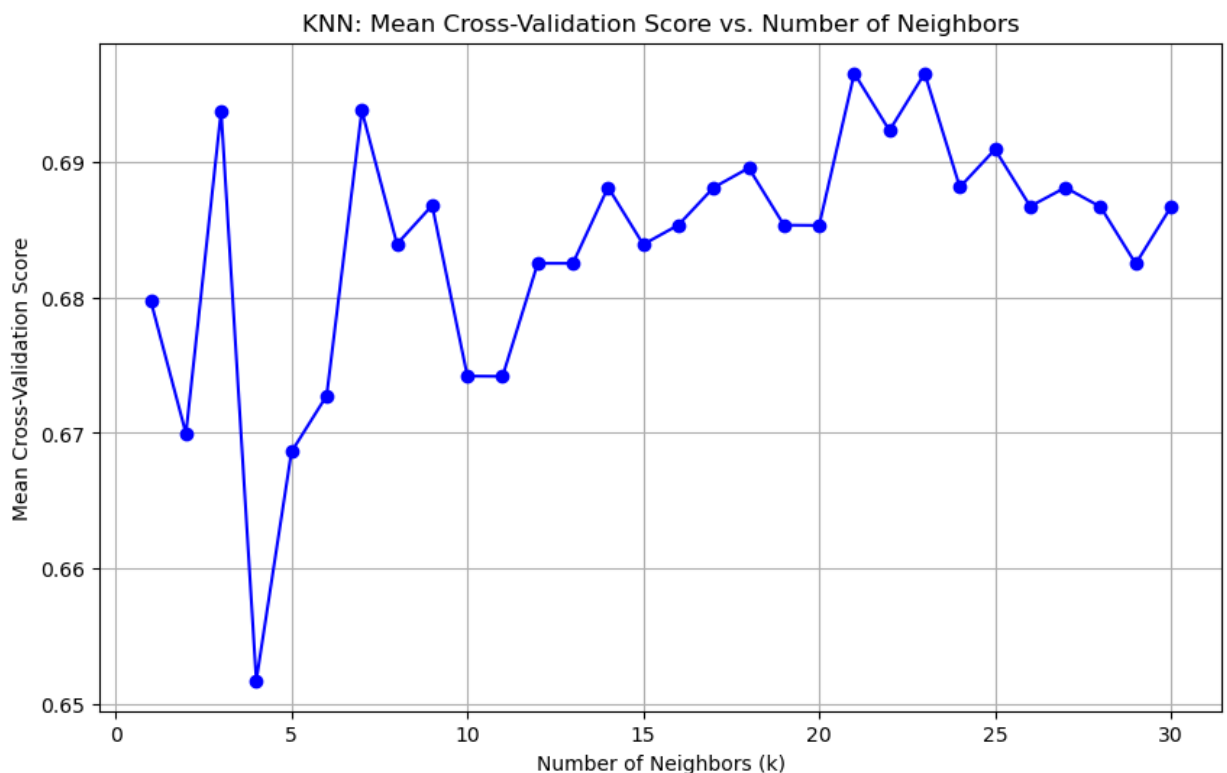
In [49]:
```python
# Define a range of k values to try (e.g., from 1 to 30)
k_values = np.arange(1, 31)

# Initialize lists to store the mean cross-validation scores for each k
cv_scores = []

# Perform k-fold cross-validation for each k value
for k in k_values:
    knn_model = KNeighborsClassifier(n_neighbors=k)
    scores = cross_val_score(knn_model, X_train, y_train, cv=5)
    cv_scores.append(scores.mean())

# Find the optimal k with the highest mean cross-validation score
optimal_k = k_values[np.argmax(cv_scores)]

# Plot the mean cross-validation scores for each k value
import matplotlib.pyplot as plt

plt.figure(figsize=(10, 6))
plt.plot(k_values, cv_scores, marker='o', linestyle='-', color='b')
plt.xlabel('Number of Neighbors (k)')
plt.ylabel('Mean Cross-Validation Score')
plt.title('KNN: Mean Cross-Validation Score vs. Number of Neighbors')
plt.grid(True)
plt.show()

print("Optimal Number of Neighbors (k):", optimal_k)
```



Optimal Number of Neighbors (k): 23

In [50]:
```python
from sklearn.tree import DecisionTreeClassifier

# Create a Decision Tree classifier
decision_tree_model = DecisionTreeClassifier(random_state=42)

# Fit the Decision Tree model to the training data
decision_tree_model.fit(X_train, y_train)

# Make predictions on the test set using the Decision Tree model
y_pred_decision_tree = decision_tree_model.predict(X_test)

# Evaluate the Decision Tree model
accuracy_decision_tree = accuracy_score(y_test, y_pred_decision_tree)
print("Decision Tree Accuracy:", accuracy_decision_tree)

# Get a classification report for Decision Tree
print("Classification Report (Decision Tree):")
print(classification_report(y_test, y_pred_decision_tree))

# Get the confusion matrix for Decision Tree
conf_matrix_decision_tree = confusion_matrix(y_test, y_pred_decision_tr
print("Confusion Matrix (Decision Tree):")
print(conf_matrix_decision_tree)
```

```
Decision Tree Accuracy: 0.7932960893854749
Classification Report (Decision Tree):
              precision    recall  f1-score   support

           0       0.81      0.84      0.83       105
           1       0.76      0.73      0.74        74

    accuracy                           0.79       179
   macro avg       0.79      0.78      0.79       179
weighted avg       0.79      0.79      0.79       179


Confusion Matrix (Decision Tree):
[[88 17]
 [20 54]]
```
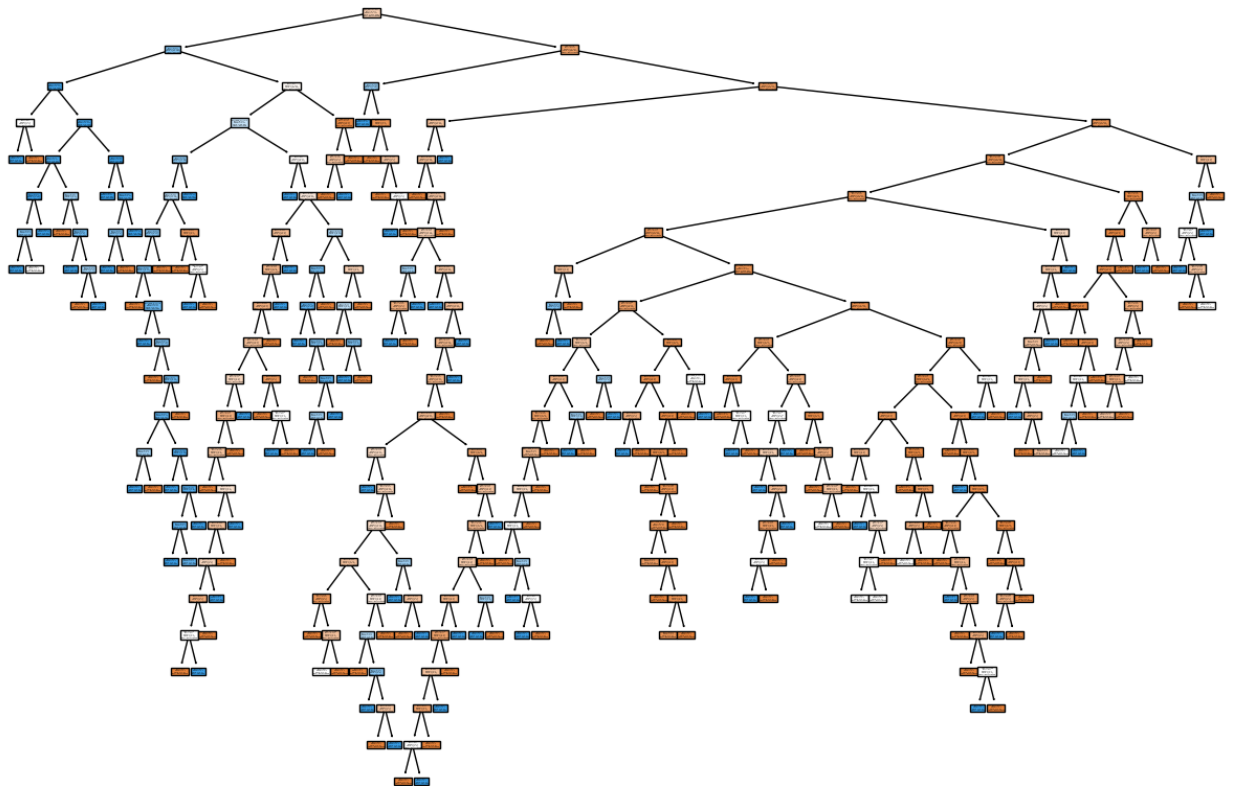
```
In [55]:    1  plt.rcParams['figure.figsize'] = (15, 10)
            2  plot_tree(decision_tree_model, feature_names=X.columns, class_names=['N
            3  plt.show()
```



```
In [56]:    1  # Create a Decision Tree classifier with modified parameters
            2  decision_tree_model_modified = DecisionTreeClassifier(max_depth=5, min_
            3
            4  # Fit the modified Decision Tree model to the training data
            5  decision_tree_model_modified.fit(X_train, y_train)
            6
            7  # Make predictions on the test set using the modified Decision Tree mod
            8  y_pred_decision_tree_modified = decision_tree_model_modified.predict(X_
            9
           10  # Evaluate the modified Decision Tree model
           11  accuracy_decision_tree_modified = accuracy_score(y_test, y_pred_decisio
           12  print("Modified Decision Tree Accuracy:", accuracy_decision_tree_modifi
```

```
Modified Decision Tree Accuracy: 0.7988826815642458
```

In [57]:
```python
from sklearn.ensemble import RandomForestClassifier

# Create a Random Forest classifier with default parameters (n_estimato
random_forest_model = RandomForestClassifier(random_state=42)

# Fit the Random Forest model to the training data
random_forest_model.fit(X_train, y_train)

# Make predictions on the test set using the Random Forest model
y_pred_random_forest = random_forest_model.predict(X_test)

# Evaluate the Random Forest model
accuracy_random_forest = accuracy_score(y_test, y_pred_random_forest)
print("Random Forest Accuracy (Default n_estimators=100):", accuracy_ra
```
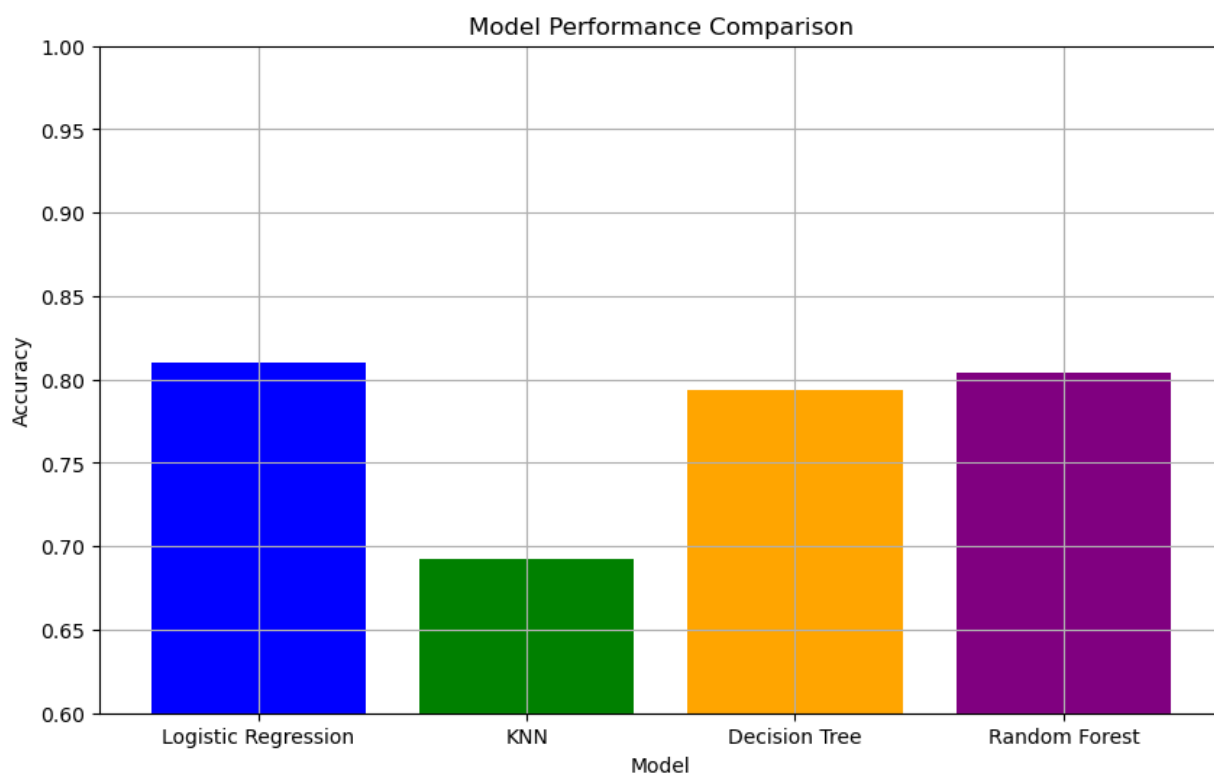
Random Forest Accuracy (Default n_estimators=100): 0.8044692737430168

In [58]:
```python
# Create a Random Forest classifier with modified number of estimators
random_forest_model_modified = RandomForestClassifier(n_estimators=200,

# Fit the modified Random Forest model to the training data
random_forest_model_modified.fit(X_train, y_train)

# Make predictions on the test set using the modified Random Forest mod
y_pred_random_forest_modified = random_forest_model_modified.predict(X_

# Evaluate the modified Random Forest model
accuracy_random_forest_modified = accuracy_score(y_test, y_pred_random_
print("Random Forest Accuracy (Modified n_estimators=200):", accuracy_r
```

Random Forest Accuracy (Modified n_estimators=200): 0.8100558659217877

In [65]:
```python
import matplotlib.pyplot as plt

# List of model names and their respective accuracy scores
model_names = ['Logistic Regression', 'KNN', 'Decision Tree', 'Random F
acc = [accuracy_knn, accuracy_decision_tree, accuracy_random_forest]
```

In [69]:

```python
import matplotlib.pyplot as plt

# List of model names and their respective accuracy scores
model_names = ['Logistic Regression', 'KNN', 'Decision Tree', 'Random F
accuracy_scores = [accuracy, accuracy_knn, accuracy_decision_tree, accu

# Plot the summary graph
plt.figure(figsize=(10, 6))
plt.bar(model_names, accuracy_scores, color=['blue', 'green', 'orange',
plt.ylim(0.6, 1.0)   # Set the y-axis limits to better compare the accur
plt.xlabel('Model')
plt.ylabel('Accuracy')
plt.title('Model Performance Comparison')
plt.grid(True)
plt.show()
```



In [ ]:

```
1
```