

# Predictions on Iris dataset

## Data Scientist : Simon Pierre Charbel

In this code we will apply 5 machine learning models to the iris dataset.

```
In [44]: 1 # Import necessary Libraries and packages
2 import pandas as pd # For data manipulation and analysis
3 import numpy as np # For numerical operations
4
5 from sklearn.datasets import load_iris # To Load the Iris dataset
6
7 from sklearn.model_selection import train_test_split # For splitting t
8
9 from sklearn.svm import SVC # Import the Support Vector Machine (SVM) c
10 from sklearn.linear_model import LogisticRegression
11 from sklearn.tree import DecisionTreeClassifier
12 from sklearn.neighbors import KNeighborsClassifier
13 from sklearn.naive_bayes import GaussianNB
14
15 from sklearn.metrics import accuracy_score # For calculating the accur
16 from sklearn.metrics import confusion_matrix # For creating a confusion
17 from sklearn.metrics import classification_report # For generating a cl
18
19 import pickle # Save the module
20
21 # For data visualization
22 import seaborn as sns
23 import matplotlib.pyplot as plt
24
25 # Suppress all warnings
26 import warnings
27 warnings.filterwarnings("ignore")
28
29 file_path = r'C:\Users\simon\AppData\Roaming\JetBrains\PyCharmCE2022.2\
```

### 1. Loading and pre-processing the data

```
In [45]: 1 # Load the Iris dataset
2 iris = load_iris()
3
4 # Create a pandas DataFrame
5 iris_df = pd.DataFrame(data=iris.data, columns=iris.feature_names)
6
7 # Add the target column to the DataFrame
8 iris_df['target'] = iris.target
9 # Add the target names as a new column
10 iris_df['species'] = iris.target_names[iris.target]
11
12 # Drop the Duplicates
13 iris_df.drop_duplicates(inplace=True)
14
15 # Print shape and the first few rows of the dataframe
16 print(f"DataFrame Shape: {iris_df.shape[0]} rows and {iris_df.shape[1]}")
17 iris_df.head()
```

DataFrame Shape: 149 rows and 6 columns (shape).

```
Out[45]:
```

	sepal length (cm)	sepal width (cm)	petal length (cm)	petal width (cm)	target	species
0	5.1	3.5	1.4	0.2	0	setosa
1	4.9	3.0	1.4	0.2	0	setosa
2	4.7	3.2	1.3	0.2	0	setosa
3	4.6	3.1	1.5	0.2	0	setosa
4	5.0	3.6	1.4	0.2	0	setosa

## 2. Splitting the data into train (75%) and test (25%)

```
In [46]: 1 # Separating the independent variables (features) from the dependent variable
2 X = iris_df.iloc[:, 0:4] # X contains the feature columns (sepal length, sepal width, petal length, petal width)
3 y = iris_df.iloc[:, -2] # y contains the target column (species)
4
5 # Splitting the dataset into training and testing sets
6 # Shuffle the dataset randomly
7 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
8
```

## 3. Prediction with different models

```
In [47]: 1 # Mapping dictionary to convert class labels
2 class_labels = {0: 'Iris Setosa', 1: 'Iris Versicolor', 2: 'Iris Virginica'}
```

### 3.1 SVM

```
In [48]: 1 # Create a Support Vector Machine (SVM) classifier
          2 svm = SVC()
          3
          4 # Train the SVM classifier on the training data
          5 svm.fit(X_train, y_train)
```

Out[48]: SVC()

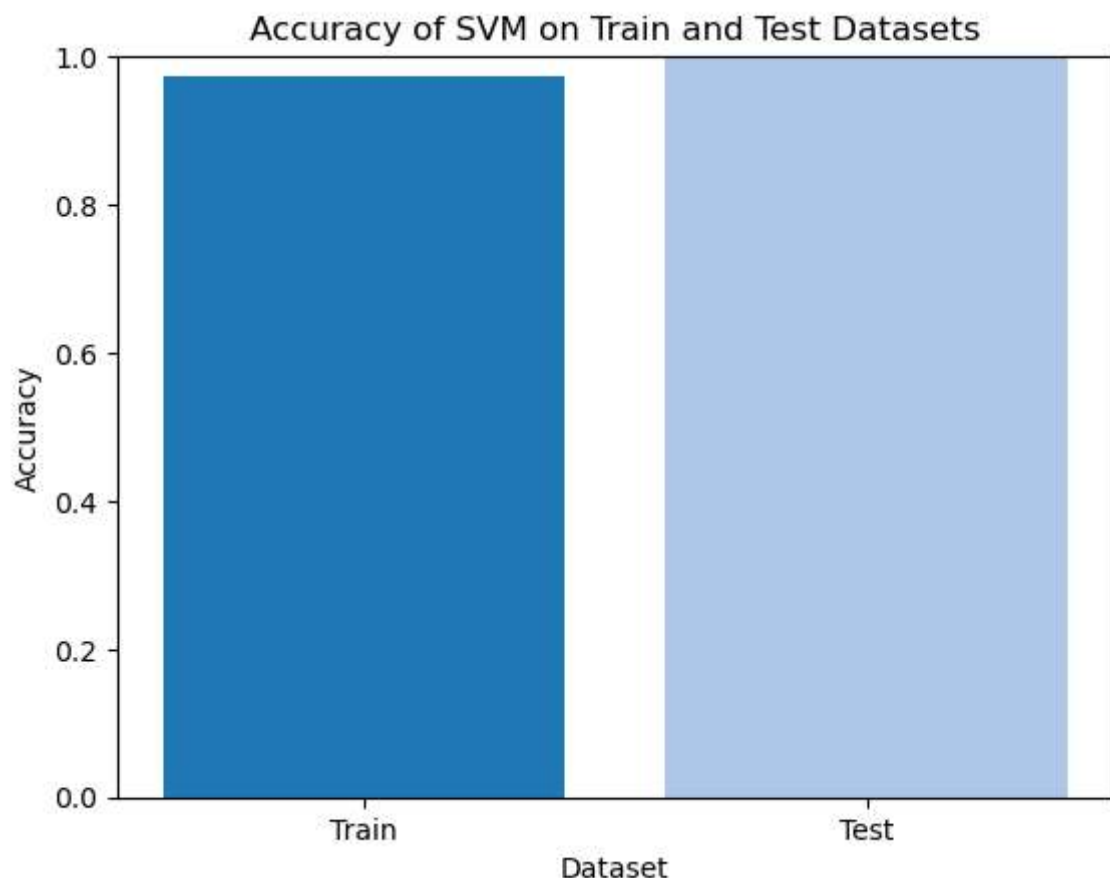
**In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.**

**On GitHub, the HTML representation is unable to render, please try loading this page with nbviewer.org.**

```
In [49]: 1 # Make predictions on both the training and test data
2 train_pred_svm = svm.predict(X_train)
3 test_pred_svm = svm.predict(X_test)
4
5 # Evaluate the SVM classifier
6 # Calculate accuracy for both datasets
7 train_accuracy_svm = accuracy_score(y_train, train_pred_svm)
8 print(f'Train Accuracy : {round(train_accuracy_svm, 2)*100}%')
9 test_accuracy_svm = accuracy_score(y_test, test_pred_svm)
10 print(f'Test Accuracy : {round(test_accuracy_svm, 2)*100}%')
11
12
13 # Create a bar chart to visualize the accuracy of the train and test da
14 labels = ['Train', 'Test']
15 accuracy_values_svm = [train_accuracy_svm, test_accuracy_svm]
16
17 blue_palette = ['#1f77b4', '#aec7e8']
18 plt.bar(labels, accuracy_values_svm, color=blue_palette)
19 plt.ylim(0, 1) # Set the y-axis limits to range from 0 to 1 (100%)
20 plt.xlabel('Dataset')
21 plt.ylabel('Accuracy')
22 plt.title('Accuracy of SVM on Train and Test Datasets')
23 plt.show()
```

Train Accuracy : 97.0%

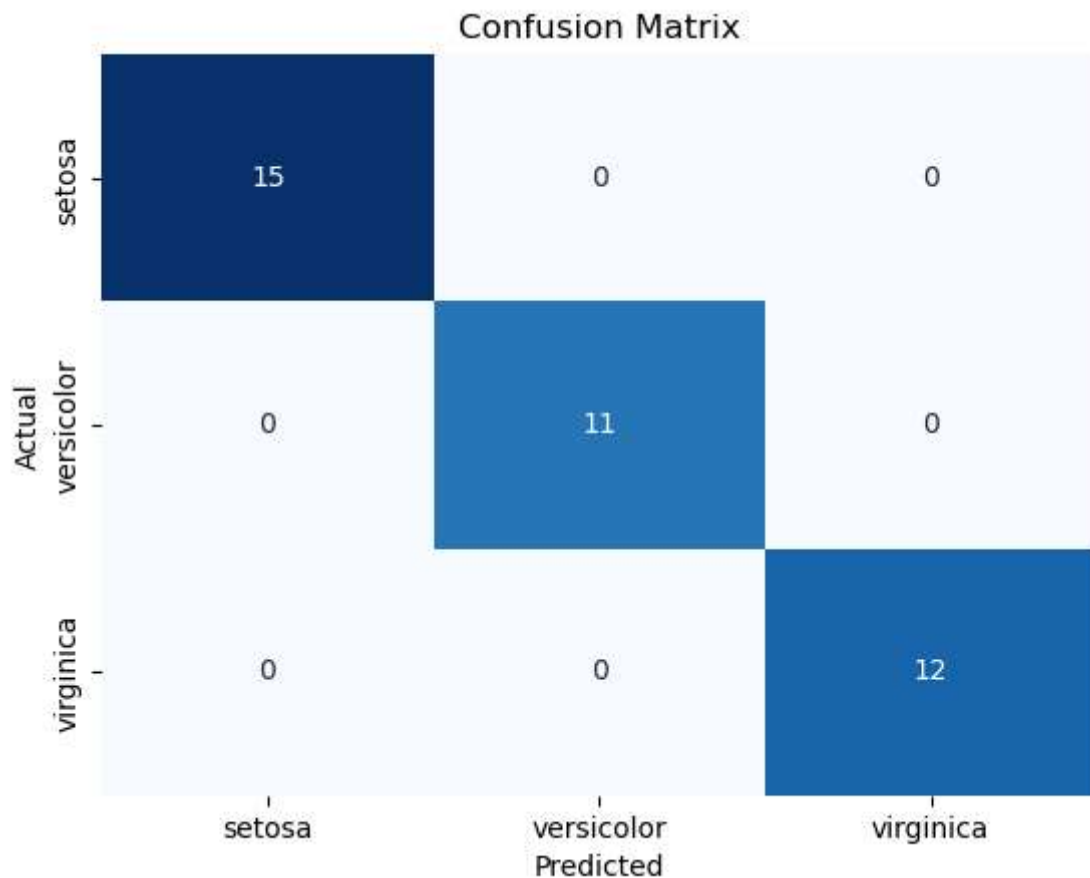
Test Accuracy : 100.0%



```

In [50]: 1 # Visualize the confusion matrix as a heatmap
2 conf_matrix = confusion_matrix(y_test, test_pred_svm)
3 sns.heatmap(conf_matrix, annot=True, fmt='d',
4             cmap='Blues', cbar=False, xticklabels=iris.target_names,
5             yticklabels=iris.target_names)
6 plt.xlabel('Predicted')
7 plt.ylabel('Actual')
8 plt.title('Confusion Matrix')
9 plt.figure(figsize=(4, 4))
10 plt.show()
11
12 # Generate a classification report
13 class_report = classification_report(y_test, test_pred_svm)
14 print("\nClassification Report:")
15 print(class_report)

```



<Figure size 400x400 with 0 Axes>

## Classification Report:

	precision	recall	f1-score	support
0	1.00	1.00	1.00	15
1	1.00	1.00	1.00	11
2	1.00	1.00	1.00	12
accuracy			1.00	38
macro avg	1.00	1.00	1.00	38
weighted avg	1.00	1.00	1.00	38

```
In [51]: 1 # Prediction
2 # Create a new data point with custom input values
3 sepal_length = 6
4 sepal_width = 4
5 petal_length = 4
6 petal_width = 4
7
8 new_data_point = np.array([[sepal_length, sepal_width, petal_length, pe
9
10 # Use the model to make predictions
11 prediction = svm.predict(new_data_point)
12 predicted_class_label = svm.predict(new_data_point)[0] # Get the predi
13 predicted_species = class_labels[predicted_class_label] # Map to speci
14
15 print(f'Predicted Species: {predicted_species} ({predicted_class_label})')
```

Predicted Species: Iris Virginica (2)

```
In [52]: 1 # Save the mode
2 with open(file_path+'svm_model.pkl', 'wb') as model_file:
3     pickle.dump(svm, model_file)
4     print('Model Saved!')
5
6
7 # Load the SVM model from a file
8 #with open(file_path+'svm_model.pkl', 'rb') as model_file:
9 #    loaded_svm = pickle.load(model_file)
10
```

Model Saved!

### 3.2 Logistic Regression

```
In [53]: 1 # Create a Logistic Regression classifier
          2 logistic_reg = LogisticRegression(max_iter=1000)
          3
          4 # Train the LR classifier on the training data
          5 logistic_reg.fit(X_train, y_train)
```

Out[53]: LogisticRegression(max\_iter=1000)

**In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.**

**On GitHub, the HTML representation is unable to render, please try loading this page with nbviewer.org.**

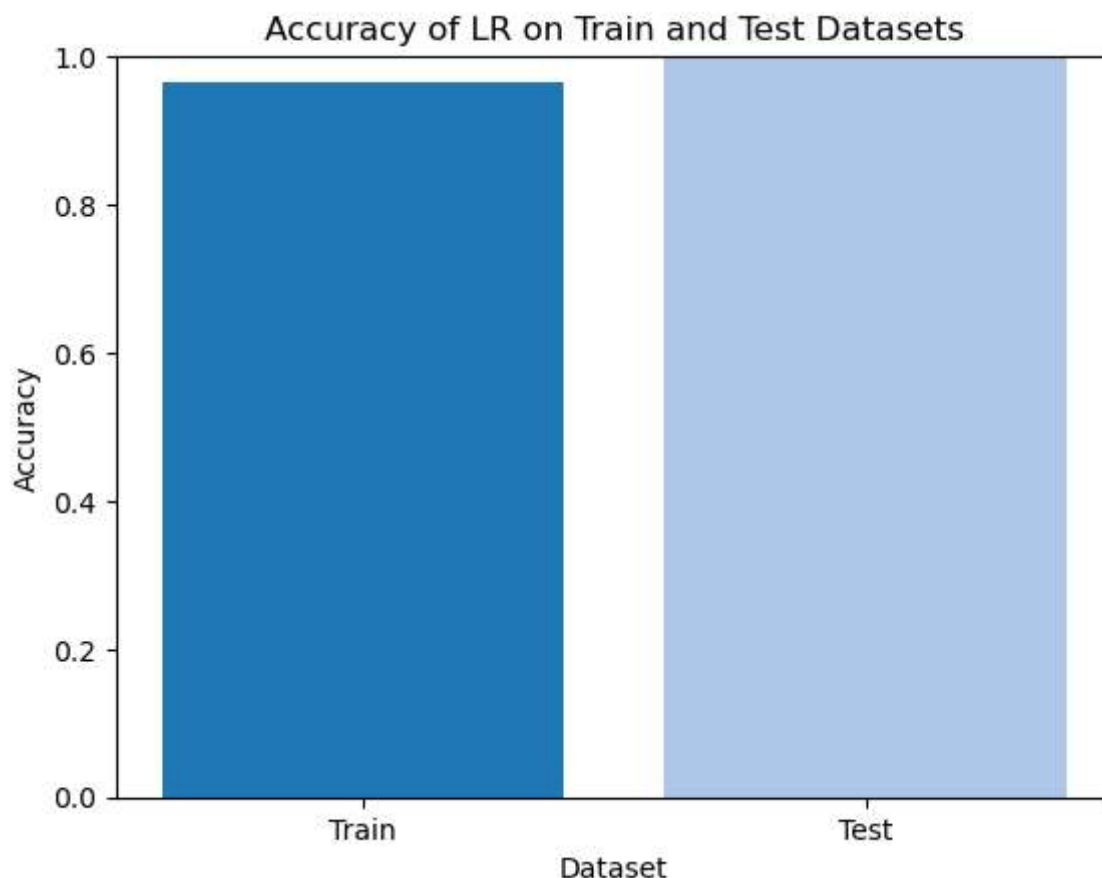
```

In [54]: 1 # Make predictions on both the training and test data
2 train_pred_logistic_reg = logistic_reg.predict(X_train)
3 test_pred_logistic_reg = logistic_reg.predict(X_test)
4
5 # Evaluate the LR classifier
6 # Calculate accuracy for both datasets
7 train_accuracy_logistic_reg = accuracy_score(y_train, train_pred_logist
8 print(f'Train Accuracy : {round(train_accuracy_logistic_reg, 2)*100}%')
9 test_accuracy_logistic_reg = accuracy_score(y_test, test_pred_logistic_
10 print(f'Test Accuracy : {round(test_accuracy_logistic_reg, 2)*100}%')
11
12
13 # Create a bar chart to visualize the accuracy of the train and test da
14 labels = ['Train', 'Test']
15 accuracy_values_logistic_reg = [train_accuracy_logistic_reg, test_accu
16
17 blue_palette = ['#1f77b4', '#aec7e8']
18 plt.bar(labels, accuracy_values_logistic_reg, color=blue_palette)
19 plt.ylim(0, 1) # Set the y-axis limits to range from 0 to 1 (100%)
20 plt.xlabel('Dataset')
21 plt.ylabel('Accuracy')
22 plt.title('Accuracy of LR on Train and Test Datasets')
23 plt.show()

```

Train Accuracy : 96.0%

Test Accuracy : 100.0%

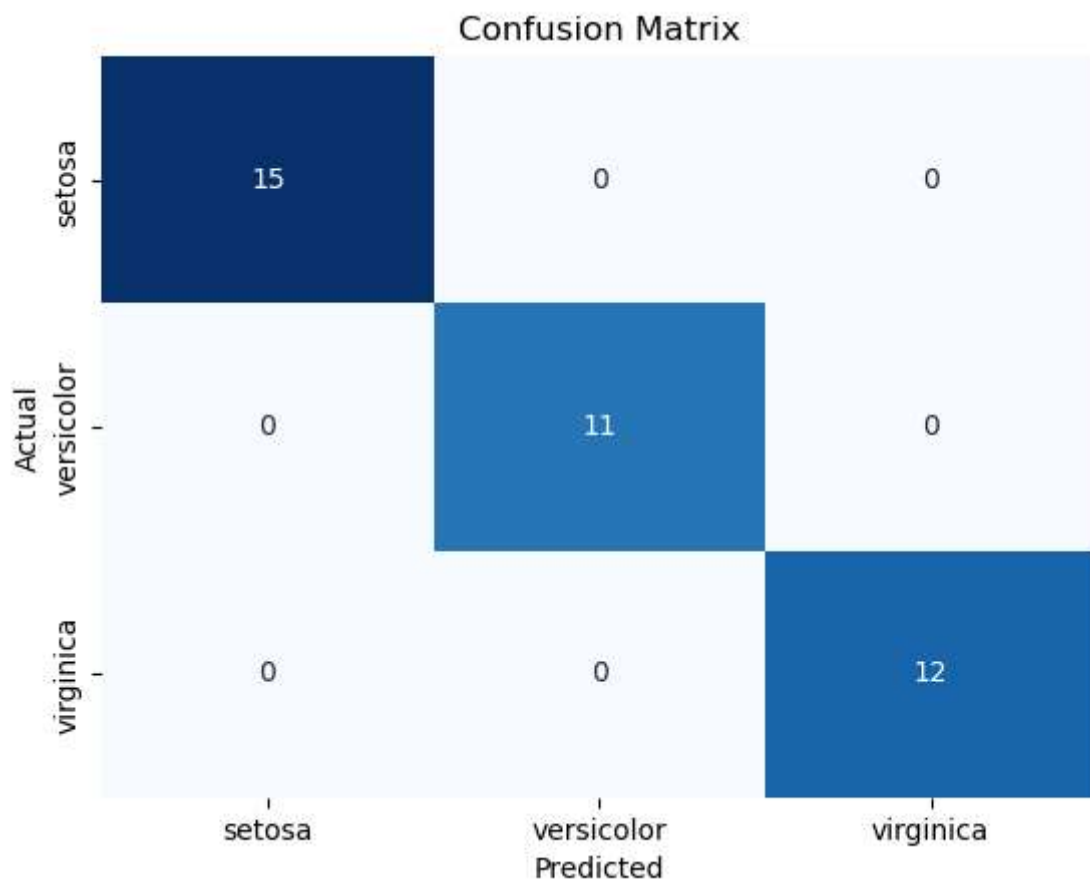




```

In [55]: 1 # Visualize the confusion matrix as a heatmap
2 conf_matrix = confusion_matrix(y_test, test_pred_logistic_reg)
3 sns.heatmap(conf_matrix, annot=True, fmt='d',
4             cmap='Blues', cbar=False, xticklabels=iris.target_names,
5             yticklabels=iris.target_names)
6 plt.xlabel('Predicted')
7 plt.ylabel('Actual')
8 plt.title('Confusion Matrix')
9 plt.figure(figsize=(4, 4))
10 plt.show()
11
12 # Generate a classification report
13 class_report = classification_report(y_test, test_pred_logistic_reg)
14 print("\nClassification Report:")
15 print(class_report)

```



<Figure size 400x400 with 0 Axes>

## Classification Report:

	precision	recall	f1-score	support
0	1.00	1.00	1.00	15
1	1.00	1.00	1.00	11
2	1.00	1.00	1.00	12
accuracy			1.00	38
macro avg	1.00	1.00	1.00	38
weighted avg	1.00	1.00	1.00	38

```
In [56]: 1 # Prediction
2 # Create a new data point with custom input values
3 sepal_length = 1
4 sepal_width = 1
5 petal_length = 1
6 petal_width = 1
7
8 new_data_point = np.array([[sepal_length, sepal_width, petal_length, pe
9
10 # Use the model to make predictions
11 prediction = logistic_reg.predict(new_data_point)
12 predicted_class_label = logistic_reg.predict(new_data_point)[0] # Get
13 predicted_species = class_labels[predicted_class_label] # Map to speci
14
15 print(f'Predicted Species: {predicted_species} ({predicted_class_label})')
```

Predicted Species: Iris Setosa (0)

```
In [57]: 1 # Save the mode
2 with open(file_path+'\\logistic_reg_model.pkl', 'wb') as model_file:
3     pickle.dump(logistic_reg, model_file)
4     print('Model Saved!')
5
6
7 # Load the LR model from a file
8 #with open(file_path+'\\logistic_reg_model.pkl', 'rb') as model_file:
9 #    loaded_svm = pickle.load(model_file)
```

Model Saved!

### 3.3 Decision Tree

```

In [58]: 1 # Create Decision Tree classifier object
2 dt = DecisionTreeClassifier()
3 dt.fit(X_train, y_train)
4
5 # Make predictions on both the training and test data
6 train_pred_dt = dt.predict(X_train)
7 test_pred_dt = dt.predict(X_test)
8
9 # Evaluate the DT classifier
10 # Calculate accuracy for both datasets
11 train_accuracy_dt = accuracy_score(y_train, train_pred_dt)
12 print(f'Train Accuracy : {round(train_accuracy_dt, 2)*100}%')
13 test_accuracy_dt = accuracy_score(y_test, test_pred_dt)
14 print(f'Test Accuracy : {round(test_accuracy_dt, 2)*100}%')
15
16 # Generate a classification report
17 class_report = classification_report(y_test, test_pred_dt)
18 print("\nClassification Report:")
19 print(class_report)
20
21 # Save the model
22 with open(file_path+'dt_model.pkl', 'wb') as model_file:
23     pickle.dump(dt, model_file)
24 print('Model Saved!')

```

Train Accuracy : 100.0%

Test Accuracy : 100.0%

Classification Report:

	precision	recall	f1-score	support
0	1.00	1.00	1.00	15
1	1.00	1.00	1.00	11
2	1.00	1.00	1.00	12
accuracy			1.00	38
macro avg	1.00	1.00	1.00	38
weighted avg	1.00	1.00	1.00	38

Model Saved!

### 3.4 K-NN

```

In [59]: 1 # Create Decision Tree classifier object
          2 k = 3
          3 knn = KNeighborsClassifier(n_neighbors=k)
          4 knn.fit(X_train, y_train)
          5
          6 # Make predictions on both the training and test data
          7 train_pred_knn = knn.predict(X_train)
          8 test_pred_knn = knn.predict(X_test)
          9
          10 # Evaluate the DT classifier
          11 # Calculate accuracy for both datasets
          12 train_accuracy_knn = accuracy_score(y_train, train_pred_knn)
          13 print(f'Train Accuracy : {round(train_accuracy_knn, 2)*100}%')
          14 test_accuracy_knn = accuracy_score(y_test, test_pred_knn)
          15 print(f'Test Accuracy : {round(test_accuracy_knn, 2)*100}%')
          16
          17 # Generate a classification report
          18 class_report = classification_report(y_test, test_pred_knn)
          19 print("\nClassification Report:")
          20 print(class_report)
          21
          22 # Save the model
          23 with open(file_path+'knn_model.pkl', 'wb') as model_file:
          24     pickle.dump(knn, model_file)
          25 print('Model Saved!')

```

Train Accuracy : 95.0%

Test Accuracy : 100.0%

Classification Report:

	precision	recall	f1-score	support
0	1.00	1.00	1.00	15
1	1.00	1.00	1.00	11
2	1.00	1.00	1.00	12
accuracy			1.00	38
macro avg	1.00	1.00	1.00	38
weighted avg	1.00	1.00	1.00	38

Model Saved!

### 3.5 Gaussian Naive Bayes

```

In [60]: 1 # Create Decision Tree classifier object
2 gnb = GaussianNB()
3 gnb.fit(X_train, y_train)
4
5 # Make predictions on both the training and test data
6 train_pred_gnb = gnb.predict(X_train)
7 test_pred_gnb = gnb.predict(X_test)
8
9 # Evaluate the DT classifier
10 # Calculate accuracy for both datasets
11 train_accuracy_gnb = accuracy_score(y_train, train_pred_gnb)
12 print(f'Train Accuracy : {round(train_accuracy_gnb, 2)*100}%')
13 test_accuracy_gnb = accuracy_score(y_test, test_pred_gnb)
14 print(f'Test Accuracy : {round(test_accuracy_gnb, 2)*100}%')
15
16 # Generate a classification report
17 class_report = classification_report(y_test, test_pred_gnb)
18 print("\nClassification Report:")
19 print(class_report)
20
21 # Save the model
22 with open(file_path+'gnb_model.pkl', 'wb') as model_file:
23     pickle.dump(gnb, model_file)
24 print('Model Saved!')

```

Train Accuracy : 94.0%

Test Accuracy : 100.0%

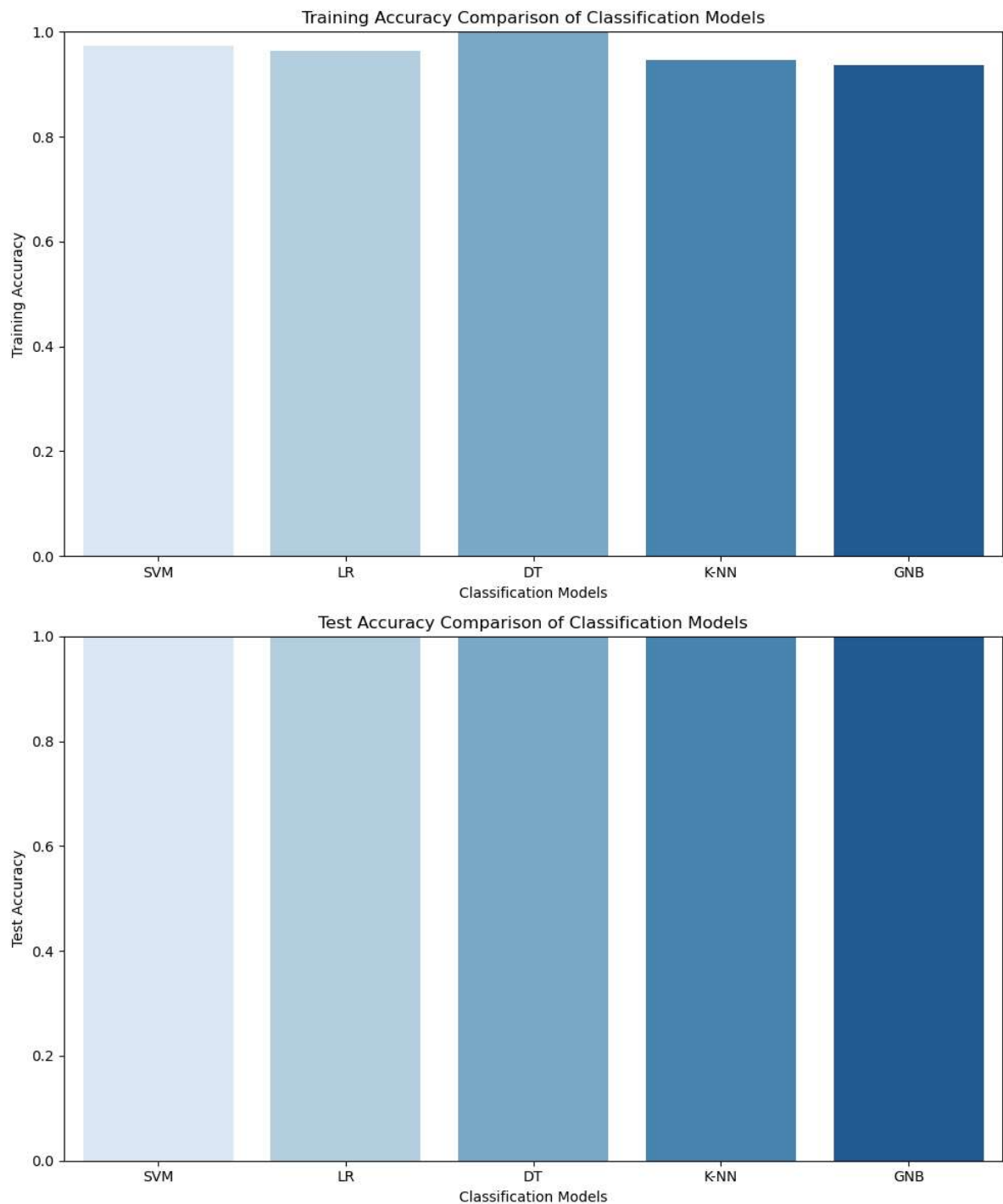
Classification Report:

	precision	recall	f1-score	support
0	1.00	1.00	1.00	15
1	1.00	1.00	1.00	11
2	1.00	1.00	1.00	12
accuracy			1.00	38
macro avg	1.00	1.00	1.00	38
weighted avg	1.00	1.00	1.00	38

Model Saved!

#### 4. Compare the models

```
In [61]: 1 # Model names and their corresponding accuracies
2 models = ['SVM', 'LR', 'DT', 'K-NN', 'GNB']
3 train_accuracies = [train_accuracy_svm, train_accuracy_logistic_reg, tr
4 test_accuracies = [test_accuracy_svm, test_accuracy_logistic_reg, test_
5
6 # Create subplots with two rows and one column
7 fig, axes = plt.subplots(2, 1, figsize=(10, 12))
8
9 # Plot the training accuracies
10 sns.barplot(x=models, y=train_accuracies, ax=axes[0], palette="Blues")
11 axes[0].set_xlabel('Classification Models')
12 axes[0].set_ylabel('Training Accuracy')
13 axes[0].set_title('Training Accuracy Comparison of Classification Model
14 axes[0].set_ylim(0, 1) # Set the y-axis limits to range from 0 to 1 (1
15
16 # Plot the test accuracies
17 sns.barplot(x=models, y=test_accuracies, ax=axes[1], palette="Blues")
18 axes[1].set_xlabel('Classification Models')
19 axes[1].set_ylabel('Test Accuracy')
20 axes[1].set_title('Test Accuracy Comparison of Classification Models')
21 axes[1].set_ylim(0, 1) # Set the y-axis limits to range from 0 to 1 (1
22
23 # Adjust spacing between subplots
24 plt.tight_layout()
25
26 # Show the combined plot
27 plt.show()
28
```



- Conclusion:

In summary, all models achieved 100% accuracy on the testing data, indicating their strong performance in classifying the Iris dataset. However, when evaluating their accuracy on the training data, the Decision Tree (DT) model outperformed the other models, demonstrating its effectiveness in capturing patterns within the dataset.

