

```
1 from google.colab import drive
2 drive.mount('/content/drive')
```

Drive already mounted at /content/drive; to attempt to forcibly remount, call dr

```
1 import pandas as pd
2 import numpy as np
3
```

```
1 df= pd.read_csv("/content/drive/MyDrive/Lab_performance/Breast Cancer Wisc
```

```
1 #handle missing value
2
3 missing_counts = df.isna().sum()
4 print(missing_counts[missing_counts > 0])
5
```

```
Unnamed: 32    569
dtype: int64
```

```
1 #Duplicate rows remove
2
3 duplicate_count = df.duplicated().sum()
4 print("Number of duplicate rows found:", duplicate_count)
5 df_cleaned = df.drop_duplicates(keep='first')
6 df_cleaned.reset_index(drop=True, inplace=True)
7
8 print("New shape after removing duplicates:", df_cleaned.shape)
9
10
```

```
Number of duplicate rows found: 0
New shape after removing duplicates: (569, 33)
```

```
1 # Correct data type
2
3 df['diagnosis'] = df.get('diagnosis').astype('category')
4 print(df['diagnosis'].value_counts())
5
```

```
diagnosis
B    357
M    212
Name: count, dtype: int64
```

```

1 # Identify categorical & numeric columns
2
3 cat_cols = df.select_dtypes(include=['object', 'category']).columns
4 num_cols = df.select_dtypes(include='number').columns
5
6 num_cols = num_cols.drop('id') if 'id' in num_cols else num_cols
7
8 print("Categorical columns:", list(cat_cols))
9 print("Total numeric columns:", len(num_cols))
10

```

Categorical columns: ['diagnosis']  
Total numeric columns: 31

```

1 # Label Encoding for binary categorical features
2
3 from sklearn.preprocessing import LabelEncoder
4
5 if 'diagnosis' in df and df['diagnosis'].nunique() == 2:
6     df['diagnosis_le'] = LabelEncoder().fit_transform(df['diagnosis'])
7     print(df['diagnosis'].unique())
8

```

['M', 'B']  
Categories (2, object): ['B', 'M']

```

1 #One-Hot Encoding for multi-class categorical features
2
3 from sklearn.preprocessing import OneHotEncoder
4 import pandas as pd
5
6 multi_class_cols = ['area_category', 'texture_group']
7 multi_class_cols = [col for col in multi_class_cols if col in df.columns]
8
9 if multi_class_cols:
10     print("Applying One-Hot Encoding on:", multi_class_cols)
11
12     ohe = OneHotEncoder(sparse=False, drop='first')
13     encoded = ohe.fit_transform(df[multi_class_cols])
14
15     encoded_df = pd.DataFrame(
16         encoded,
17         columns=ohe.get_feature_names_out(multi_class_cols),
18         index=df.index
19     )
20
21     df = df.drop(columns=multi_class_cols)
22     df = pd.concat([df, encoded_df], axis=1)
23
24 else:
25     print("⚠ No multi-class categorical columns found for One-Hot Encoding

```

```

26
27 df.info()
28

```

⚠ No multi-class categorical columns found for One-Hot Encoding.

<class 'pandas.core.frame.DataFrame'>

RangeIndex: 569 entries, 0 to 568

Data columns (total 34 columns):

#	Column	Non-Null Count	Dtype
0	id	569 non-null	int64
1	diagnosis	569 non-null	category
2	radius_mean	569 non-null	float64
3	texture_mean	569 non-null	float64
4	perimeter_mean	569 non-null	float64
5	area_mean	569 non-null	float64
6	smoothness_mean	569 non-null	float64
7	compactness_mean	569 non-null	float64
8	concavity_mean	569 non-null	float64
9	concave points_mean	569 non-null	float64
10	symmetry_mean	569 non-null	float64
11	fractal_dimension_mean	569 non-null	float64
12	radius_se	569 non-null	float64
13	texture_se	569 non-null	float64
14	perimeter_se	569 non-null	float64
15	area_se	569 non-null	float64
16	smoothness_se	569 non-null	float64
17	compactness_se	569 non-null	float64
18	concavity_se	569 non-null	float64
19	concave points_se	569 non-null	float64
20	symmetry_se	569 non-null	float64
21	fractal_dimension_se	569 non-null	float64
22	radius_worst	569 non-null	float64
23	texture_worst	569 non-null	float64
24	perimeter_worst	569 non-null	float64
25	area_worst	569 non-null	float64
26	smoothness_worst	569 non-null	float64
27	compactness_worst	569 non-null	float64
28	concavity_worst	569 non-null	float64
29	concave points_worst	569 non-null	float64
30	symmetry_worst	569 non-null	float64
31	fractal_dimension_worst	569 non-null	float64
32	Unnamed: 32	0 non-null	float64
33	diagnosis_le	569 non-null	int64

dtypes: category(1), float64(31), int64(2)

memory usage: 147.5 KB

```

1 # Feature Scaling (StandardScaler)
2
3 from sklearn.preprocessing import StandardScaler
4
5 remove_cols = [ "diagnosis_le", "Unnamed: 32"]
6
7 numeric_cols = [col for col in df.columns if col not in remove_cols]

```

```

8
9 numeric_cols = df[numeric_cols].select_dtypes(include='number').columns.tolist
10
11 print("Numeric Columns to Scale:", numeric_cols)
12
13 scaler = StandardScaler()
14 df[numeric_cols] = scaler.fit_transform(df[numeric_cols])
15
16 print("Scaling Completed!")
17

```

Numeric Columns to Scale: ['id', 'radius\_mean', 'texture\_mean', 'perimeter\_mean']  
Scaling Completed!

```

1 #Train-Test Split
2
3 import pandas as pd
4 from sklearn.model_selection import train_test_split
5 from sklearn.preprocessing import StandardScaler, LabelEncoder
6
7 df= pd.read_csv("/content/drive/MyDrive/Lab_performance/Breast Cancer Wiscc
8 df = df.loc[:, ~df.columns.str.contains("Unnamed")]
9
10 le = LabelEncoder()
11 df['diagnosis_le'] = le.fit_transform(df['diagnosis'])
12 print("Label Classes:", le.classes_)      # ['B' 'M']
13
14
15 numeric_cols = [c for c in df.select_dtypes(include='number').columns
16                 if c not in ['id', 'diagnosis_le']]
17 print("Numeric Columns to Scale:", numeric_cols)
18
19 scaler = StandardScaler()
20 df[numeric_cols] = scaler.fit_transform(df[numeric_cols])
21 print("Scaling Completed!\n")
22
23 X = df.drop(columns=['id', 'diagnosis', 'diagnosis_le'])
24 y = df['diagnosis_le']
25
26 X_train, X_test, y_train, y_test = train_test_split(
27     X, y, test_size=0.2, stratify=y, random_state=42
28 )
29
30 print("Train shape:", X_train.shape, y_train.shape)
31 print("Test shape:", X_test.shape, y_test.shape)
32
33 print("\nTrain class distribution:\n",
34       y_train.value_counts(normalize=True))
35 print("\nTest class distribution:\n",

```

```

36     y_test.value_counts(normalize=True))
37

```

Label Classes: ['B' 'M']  
 Numeric Columns to Scale: ['radius\_mean', 'texture\_mean', 'perimeter\_mean', 'area']  
 Scaling Completed!

Train shape: (455, 30) (455,)  
 Test shape: (114, 30) (114,)

Train class distribution:  
 diagnosis\_le  
 0 0.626374  
 1 0.373626  
 Name: proportion, dtype: float64

Test class distribution:  
 diagnosis\_le  
 0 0.631579  
 1 0.368421  
 Name: proportion, dtype: float64

```

1 #TRAIN FIVE ML MODELS
2
3 import pandas as pd
4 from sklearn.model_selection import train_test_split
5 from sklearn.preprocessing import StandardScaler
6 from sklearn.metrics import accuracy_score
7 from sklearn.linear_model import LogisticRegression
8 from sklearn.svm import SVC
9 from sklearn.neighbors import KNeighborsClassifier
10 from sklearn.ensemble import RandomForestClassifier
11 from sklearn.tree import DecisionTreeClassifier
12 df= pd.read_csv("/content/drive/MyDrive/Lab_performance/Breast Cancer Wiscc
13 df = df.loc[:, ~df.columns.str.contains('^Unnamed')]
14 df = df.fillna(df.median(numeric_only=True))
15 X = df.drop("diagnosis", axis=1)
16 y = df["diagnosis"]
17 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, r
18 scaler = StandardScaler()
19 X_train = scaler.fit_transform(X_train)
20 X_test = scaler.transform(X_test)
21 models = {
22 "Logistic Regression": LogisticRegression(max_iter=500),
23 "SVM": SVC(),
24 "KNN": KNeighborsClassifier(),
25 "Random Forest": RandomForestClassifier(),
26 "Decision Tree": DecisionTreeClassifier()
27 }
28 for name, model in models.items():
29     model.fit(X_train, y_train)
30     pred = model.predict(X_test)

```

```

31     acc = accuracy_score(y_test, pred)
32     print(f"{name} Accuracy: {acc:.4f}")

```

```

Logistic Regression Accuracy: 0.9737
SVM Accuracy: 0.9825
KNN Accuracy: 0.9474
Random Forest Accuracy: 0.9649
Decision Tree Accuracy: 0.9123

```

```

1 #Confusion Matrix + Classification Report (Precision, Recall, F1-score)
2
3 from sklearn.metrics import confusion_matrix, classification_report, accuracy_score
4 import seaborn as sns
5 import matplotlib.pyplot as plt
6
7 def evaluate_model(model, X_test, y_test, model_name):
8     y_pred = model.predict(X_test)
9
10    print(f"\n=====")
11    print(f"Model: {model_name}")
12    print("=====")
13
14    # Accuracy
15    acc = accuracy_score(y_test, y_pred)
16    print("Accuracy:", acc)
17
18    # Classification Report
19    print("\nClassification Report:")
20    print(classification_report(y_test, y_pred))
21
22    # Confusion Matrix
23    cm = confusion_matrix(y_test, y_pred)
24    plt.figure(figsize=(5,4))
25    sns.heatmap(cm, annot=True, cmap="Blues", fmt="d")
26    plt.title(f"{model_name} - Confusion Matrix")
27    plt.xlabel("Predicted")
28    plt.ylabel("Actual")
29    plt.show()
30

```

```

1 for model_name, model in models.items():
2     evaluate_model(model, X_test, y_test, model_name)
3

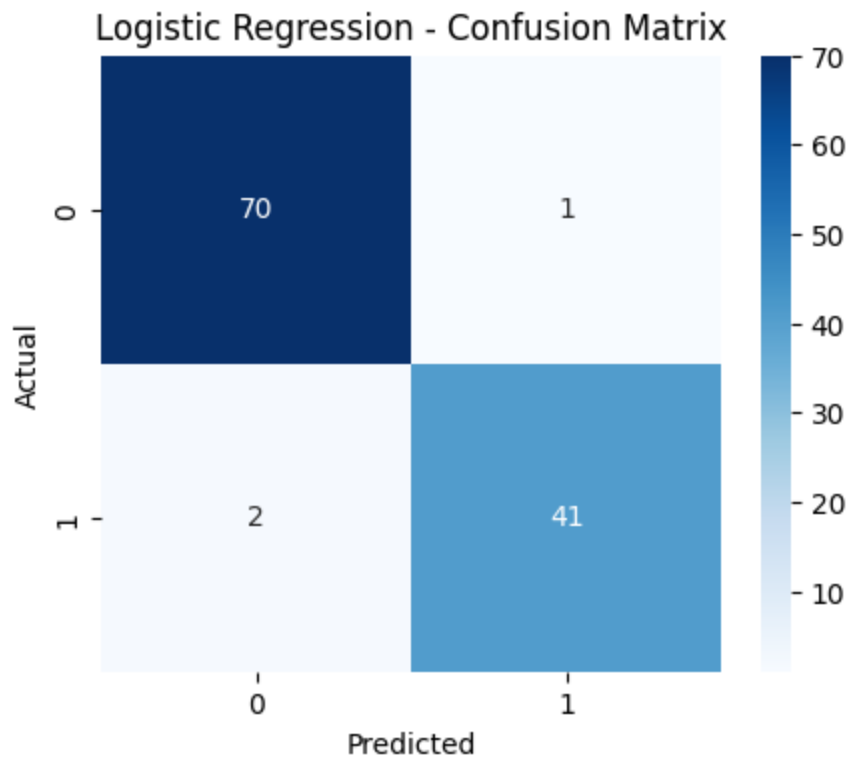
```



```
=====
Model: Logistic Regression
=====
Accuracy: 0.9736842105263158
```

Classification Report:

	precision	recall	f1-score	support
B	0.97	0.99	0.98	71
M	0.98	0.95	0.96	43
accuracy			0.97	114
macro avg	0.97	0.97	0.97	114
weighted avg	0.97	0.97	0.97	114



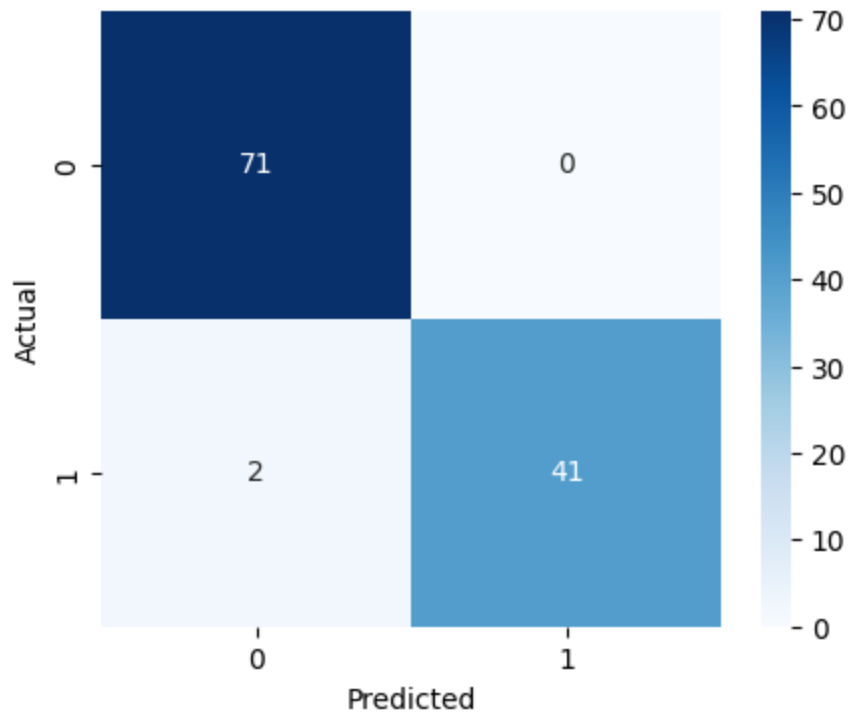
```
=====
Model: SVM
=====
Accuracy: 0.9824561403508771
```

Classification Report:

	precision	recall	f1-score	support
B	0.97	1.00	0.99	71
M	1.00	0.95	0.98	43
accuracy			0.98	114
macro avg	0.99	0.98	0.98	114
weighted avg	0.98	0.98	0.98	114

### SVM - Confusion Matrix



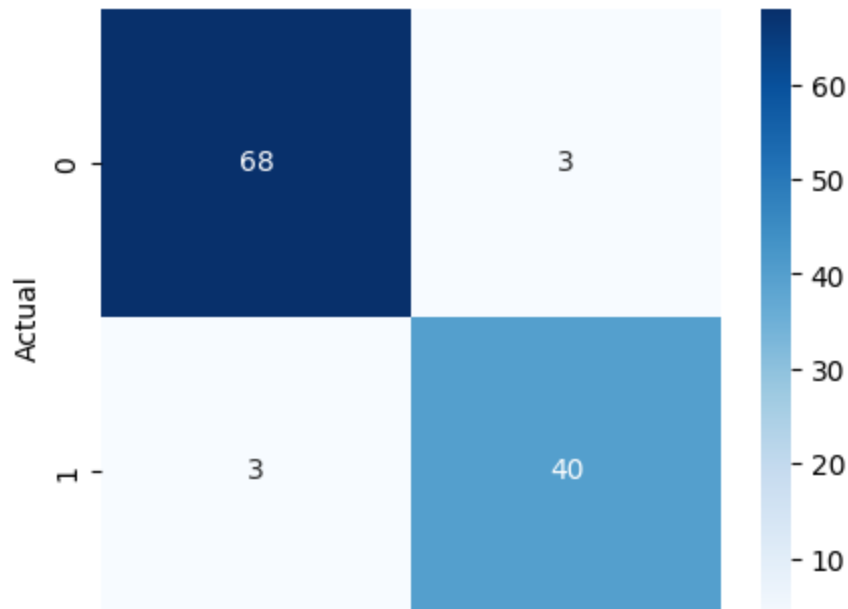


```
=====  
Model: KNN  
=====  
Accuracy: 0.9473684210526315
```

Classification Report:

	precision	recall	f1-score	support
B	0.96	0.96	0.96	71
M	0.93	0.93	0.93	43
accuracy			0.95	114
macro avg	0.94	0.94	0.94	114
weighted avg	0.95	0.95	0.95	114

KNN - Confusion Matrix



0

1

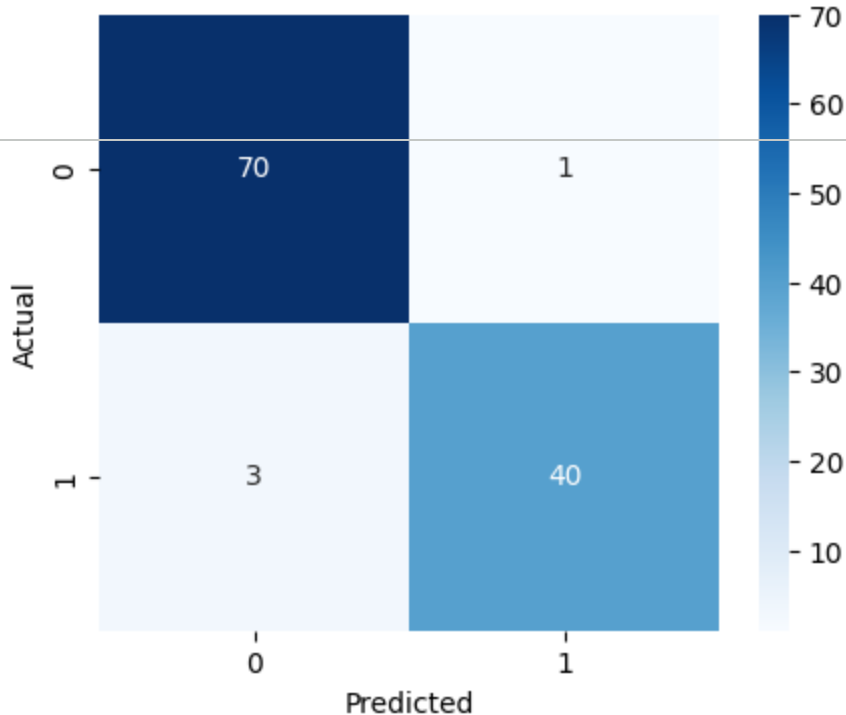
Predicted

```
=====
Model: Random Forest
=====
Accuracy: 0.9649122807017544
```

Classification Report:

	precision	recall	f1-score	support
B	0.96	0.99	0.97	71
M	0.98	0.93	0.95	43
accuracy			0.96	114
macro avg	0.97	0.96	0.96	114
weighted avg	0.97	0.96	0.96	114

Random Forest - Confusion Matrix



```
=====
Model: Decision Tree
=====
Accuracy: 0.9122807017543859
```

Classification Report:

	precision	recall	f1-score	support
B	0.94	0.92	0.93	71
M	0.87	0.91	0.89	43
accuracy			0.91	114
macro avg	0.90	0.91	0.91	114
weighted avg	0.91	0.91	0.91	114

```

1 #ROC Curve
2
3 from sklearn.metrics import roc_curve, auc
4
5 def plot_roc_curve(model, X_test, y_test, model_name):
6     # Probability predictions (required for ROC)
7     y_prob = model.predict_proba(X_test)[:, 1]
8
9     fpr, tpr, thresholds = roc_curve(y_test, y_prob)
10    roc_auc = auc(fpr, tpr)
11
12    plt.figure(figsize=(6,5))
13    plt.plot(fpr, tpr, label=f"{model_name} (AUC = {roc_auc:.3f})")
14    plt.plot([0,1], [0,1], linestyle="--")
15    plt.xlabel("False Positive Rate")
16    plt.ylabel("True Positive Rate")
17    plt.title(f"ROC Curve - {model_name}")
18    plt.legend()
19    plt.show()

```

```

1 from sklearn.preprocessing import LabelEncoder
2
3 le = LabelEncoder()
4 df['diagnosis'] = le.fit_transform(df['diagnosis'])
5
6 # Check mapping:
7 print(dict(zip(le.classes_, le.transform(le.classes_))))
8

```

```
{'B': np.int64(0), 'M': np.int64(1)}
```

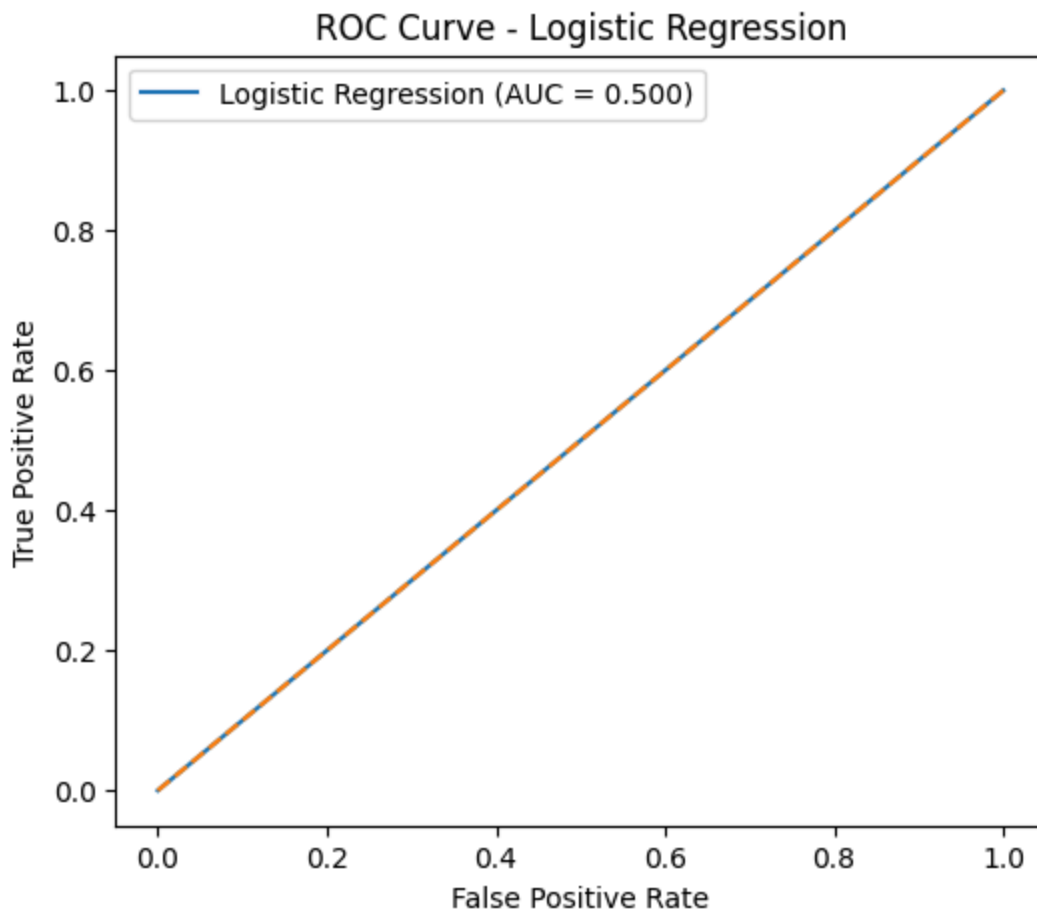
```

1 X = df.drop("diagnosis", axis=1)
2 y = df["diagnosis"]
3
4 from sklearn.model_selection import train_test_split
5 X_train, X_test, y_train, y_test = train_test_split(
6     X, y, test_size=0.2, random_state=42
7 )

```

```
1 plot_roc_curve(models["Logistic Regression"], X_test, y_test, "Logistic Reg
```

/usr/local/lib/python3.12/dist-packages/sklearn/utils/validation.py:2732: UserWarning: warnings.warn(

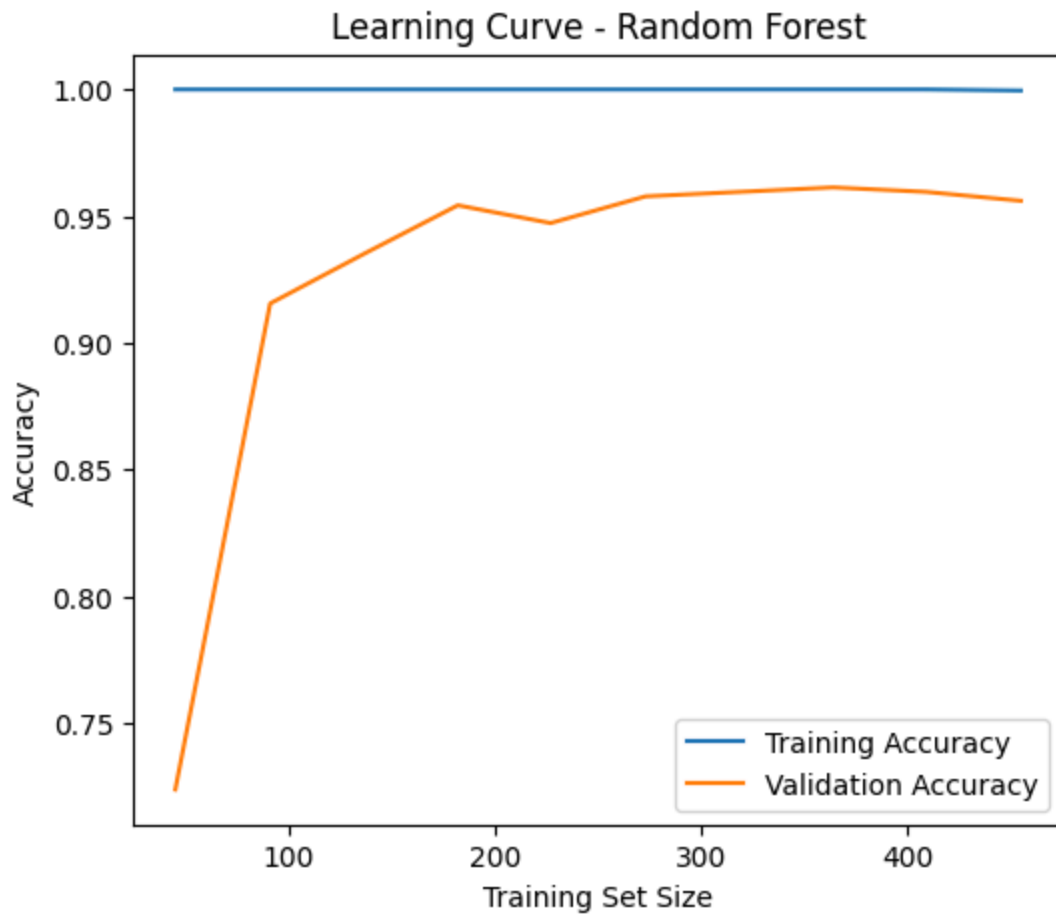


```

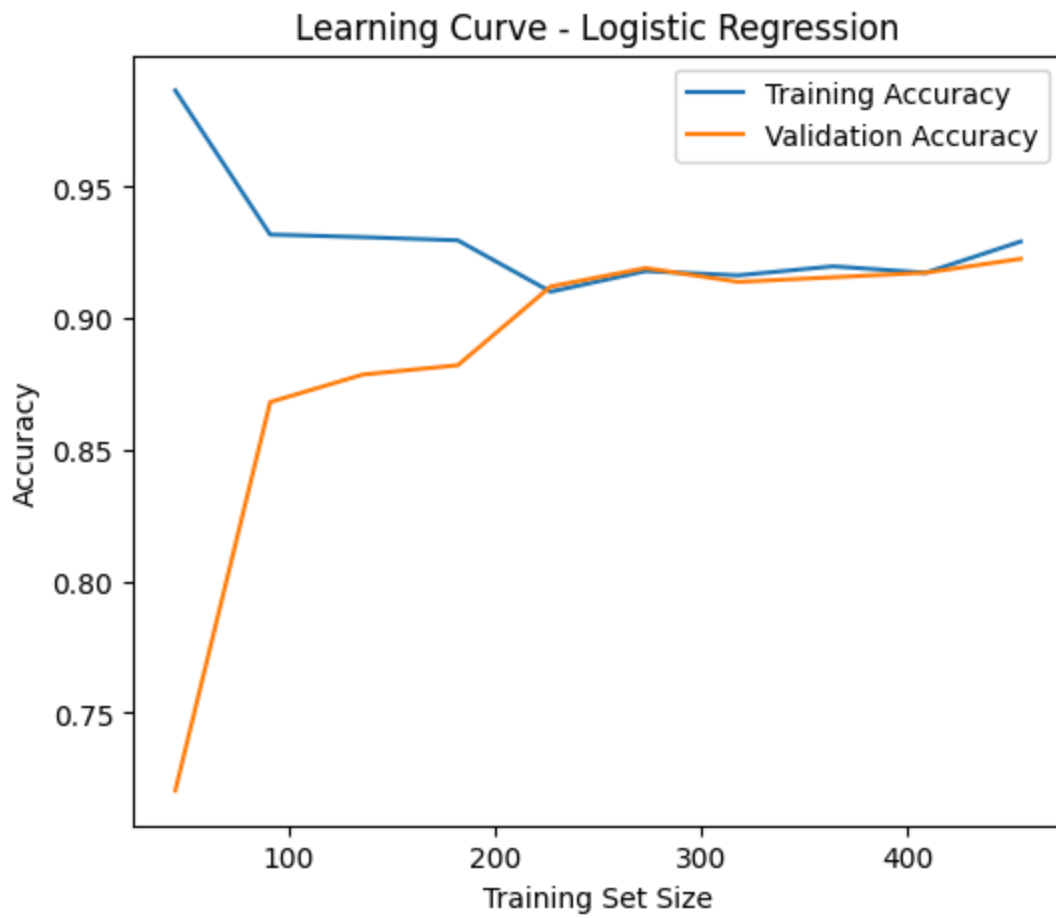
1 #Loss vs Validation Loss Curve
2 from sklearn.model_selection import learning_curve
3 import numpy as np
4
5 def plot_learning_curve(model, X, y, model_name):
6     train_sizes, train_scores, test_scores = learning_curve(
7         model, X, y, cv=5, scoring='accuracy',
8         train_sizes=np.linspace(0.1, 1.0, 10), n_jobs=-1
9     )
10
11     train_mean = train_scores.mean(axis=1)
12     test_mean = test_scores.mean(axis=1)
13
14     plt.figure(figsize=(6,5))
15     plt.plot(train_sizes, train_mean, label="Training Accuracy")
16     plt.plot(train_sizes, test_mean, label="Validation Accuracy")
17     plt.title(f"Learning Curve - {model_name}")
18     plt.xlabel("Training Set Size")
19     plt.ylabel("Accuracy")
20     plt.legend()
21     plt.show()

```

```
1 plot_learning_curve(models["Random Forest"], X, y, "Random Forest")
```



```
1 plot_learning_curve(models["Logistic Regression"], X, y, "Logistic Regressi
```



```
1 plot_learning_curve(models["Decision Tree"], X, y, "Decision Tree")
```

## Learning Curve - Decision Tree

```
1 plot_learning_curve(models["SVM"], X, y, "SVM")
```

## Learning Curve - SVM

