DataSet: Pima Indians Diabetes Dataset

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import seaborn as sns
```

```
1 from sklearn.model_selection import train_test_split
2 from sklearn.preprocessing import LabelEncoder, StandardScaler
3 from sklearn.metrics import classification_report, confusion_matr
4 from sklearn.metrics import ConfusionMatrixDisplay
5 from sklearn.model_selection import learning_curve
```

```
1 from sklearn.linear_model import LogisticRegression
2 from sklearn.neighbors import KNeighborsClassifier
3 from sklearn.tree import DecisionTreeClassifier
4 from sklearn.ensemble import RandomForestClassifier
5 from sklearn.svm import SVC
```

```
1 df = pd.read_csv("/content/drive/MyDrive/Dataset/Pima Indians Dia
```

```
1 display(df.head())
```

|   | Pregnancies | Glucose | BloodPressure | SkinThickness | Insulin | BMI | Dia |
|---|---|---|---|---|---|---|---|
| 0 | 6 | 148 | 72 | 35 | 0 | 33.6 | |
| 1 | 1 | 85 | 66 | 29 | 0 | 26.6 | |
| 2 | 8 | 183 | 64 | 0 | 0 | 23.3 | |
| 3 | 1 | 89 | 66 | 23 | 94 | 28.1 | |
| 4 | 0 | 137 | 40 | 35 | 168 | 43.1 | |

```
1 print(df.shape)
```
```
(768, 9)
```

```
1 print(df.columns.tolist())
```
```
['Pregnancies', 'Glucose', 'BloodPressure', 'SkinThickness', 'Insulin'
```

```
1 df.info()
```
```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 768 entries, 0 to 767
Data columns (total 9 columns):
```

```
 #   Column                    Non-Null Count  Dtype
---  ------                    --------------  -----
 0   Pregnancies               768 non-null    int64
 1   Glucose                   768 non-null    int64
 2   BloodPressure             768 non-null    int64
 3   SkinThickness             768 non-null    int64
 4   Insulin                   768 non-null    int64
 5   BMI                       768 non-null    float64
 6   DiabetesPedigreeFunction  768 non-null    float64
 7   Age                       768 non-null    int64
 8   Outcome                   768 non-null    int64
dtypes: float64(2), int64(7)
memory usage: 54.1 KB
```

```
1 display(df.describe())
```

|       | Pregnancies | Glucose    | BloodPressure | SkinThickness | Insulin    |
|-------|-------------|------------|---------------|---------------|------------|
| count | 768.000000  | 768.000000 | 768.000000    | 768.000000    | 768.000000 |
| mean  | 3.845052    | 120.894531 | 69.105469     | 20.536458     | 79.799479  |
| std   | 3.369578    | 31.972618  | 19.355807     | 15.952218     | 115.244002 |
| min   | 0.000000    | 0.000000   | 0.000000      | 0.000000      | 0.000000   |
| 25%   | 1.000000    | 99.000000  | 62.000000     | 0.000000      | 0.000000   |
| 50%   | 3.000000    | 117.000000 | 72.000000     | 23.000000     | 30.500000  |
| 75%   | 6.000000    | 140.250000 | 80.000000     | 32.000000     | 127.250000 |
| max   | 17.000000   | 199.000000 | 122.000000    | 99.000000     | 846.000000 |

```
1 print((df == 0).sum())
```

```
Pregnancies               111
Glucose                     5
BloodPressure              35
SkinThickness             227
Insulin                   374
BMI                        11
DiabetesPedigreeFunction    0
Age                         0
Outcome                   500
dtype: int64
```

```
1 zero_to_nan_cols = ["Glucose", "BloodPressure", "SkinThickness", "
2
```

```
1 df[zero_to_nan_cols] = df[zero_to_nan_cols].replace(0, np.nan)
```

```
1 print(df.isnull().sum())
```

```
Pregnancies                 0
Glucose                     5
```

```
BloodPressure                      35
SkinThickness                     227
Insulin                           374
BMI                                11
DiabetesPedigreeFunction            0
Age                                 0
Outcome                             0
dtype: int64
```

```
1 df.fillna(df.median(), inplace=True)
```

```
1 print(df.isnull().sum())
```

```
Pregnancies                 0
Glucose                     0
BloodPressure               0
SkinThickness               0
Insulin                     0
BMI                         0
DiabetesPedigreeFunction    0
Age                         0
Outcome                     0
dtype: int64
```

```
1 before_rows = df.shape[0]
2 df.drop_duplicates(inplace=True)
3 after_rows = df.shape[0]
```

```
1 print(f"\nRows Before Removing Duplicates: {before_rows}")
2 print(f"Rows After Removing Duplicates: {after_rows}")
```

```
Rows Before Removing Duplicates: 768
Rows After Removing Duplicates: 768
```

```
1 print(df["Outcome"].unique())
```

```
[1 0]
```

```
1 le = LabelEncoder()
2 df["Outcome"] = le.fit_transform(df["Outcome"])
```

```
1
2 print("\nUnique values in Outcome after encoding:")
3 print(df["Outcome"].unique())
```

```
Unique values in Outcome after encoding:
[1 0]
```

```
1 X = df.drop("Outcome", axis=1)
2 y = df["Outcome"]
```

```
1 print("\nFeature Matrix (X) Shape:", X.shape)
2 print("Target Vector (y) Shape:", y.shape)
```

```
Feature Matrix (X) Shape: (768, 8)
Target Vector (y) Shape: (768,)
```

```
1 scaler = StandardScaler()
2 X_scaled = scaler.fit_transform(X)
3
4 print("\nFeature Scaling Completed.")
```
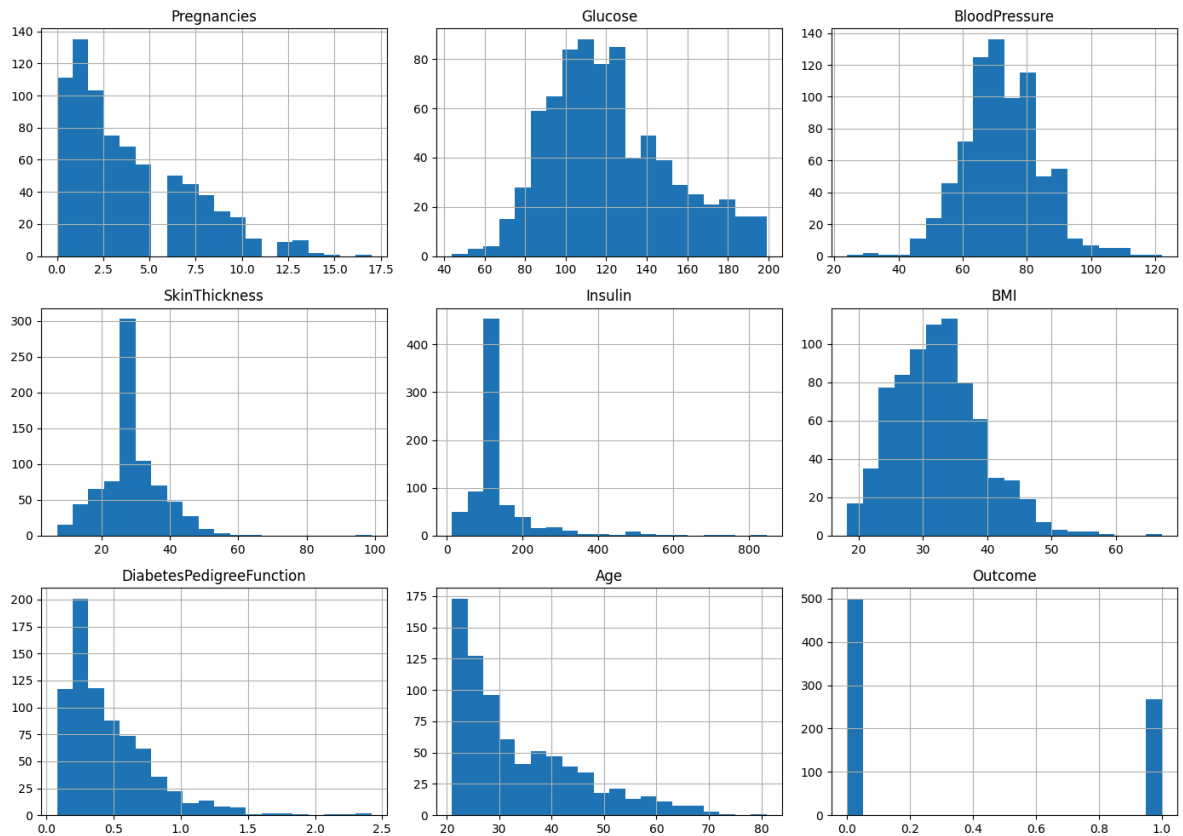
```
Feature Scaling Completed.
```

```
1 X_train, X_test, y_train, y_test = train_test_split(
2     X_scaled, y, test_size=0.20, random_state=42
3 )
```

```
1 print("\nTraining Set Shape:", X_train.shape)
2 print("Testing Set Shape:", X_test.shape)
```

```
Training Set Shape: (614, 8)
Testing Set Shape: (154, 8)
```

```
1 print("\nGenerating Histograms...")
2 df.hist(figsize=(14, 10), bins=20)
3 plt.tight_layout()
4 plt.show()
```
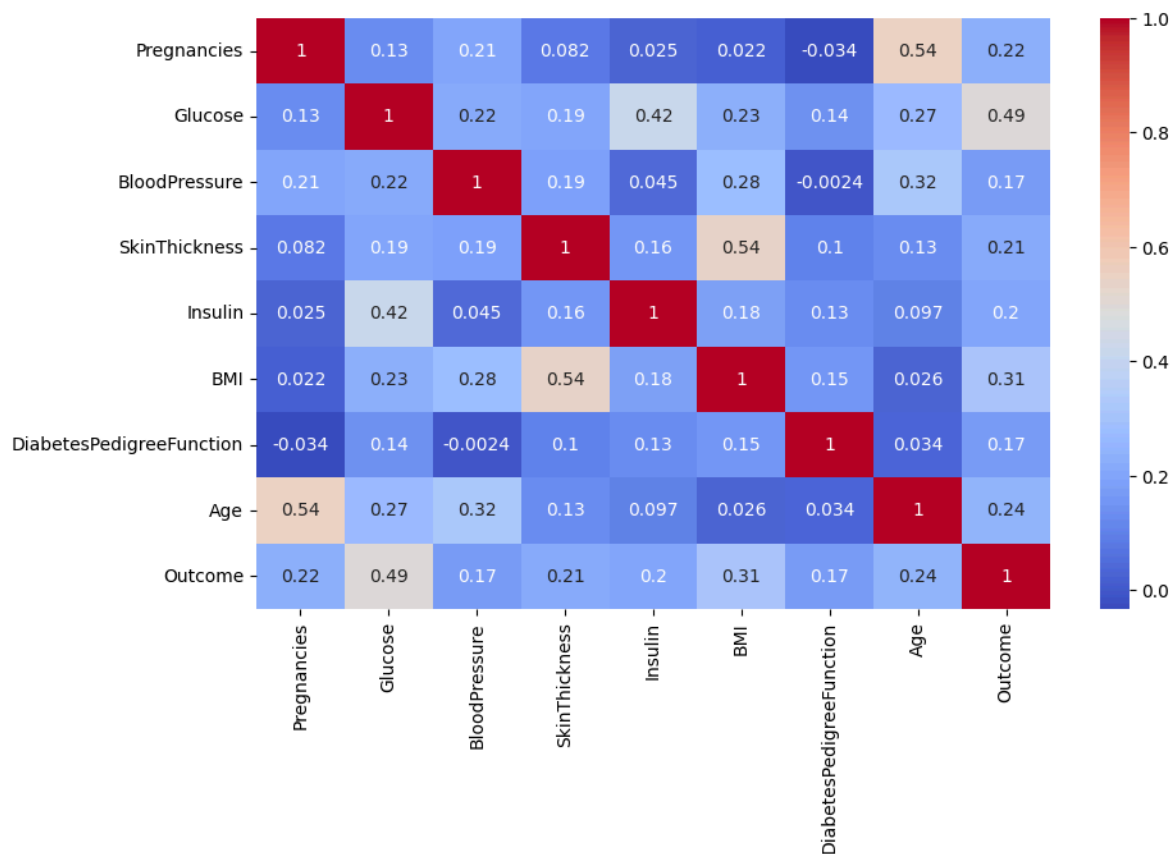
Generating Histograms...



```
1 print("\nCorrelation Heatmap:")
2 plt.figure(figsize=(10, 6))
3 sns.heatmap(df.corr(), annot=True, cmap="coolwarm")
4 plt.show()
```
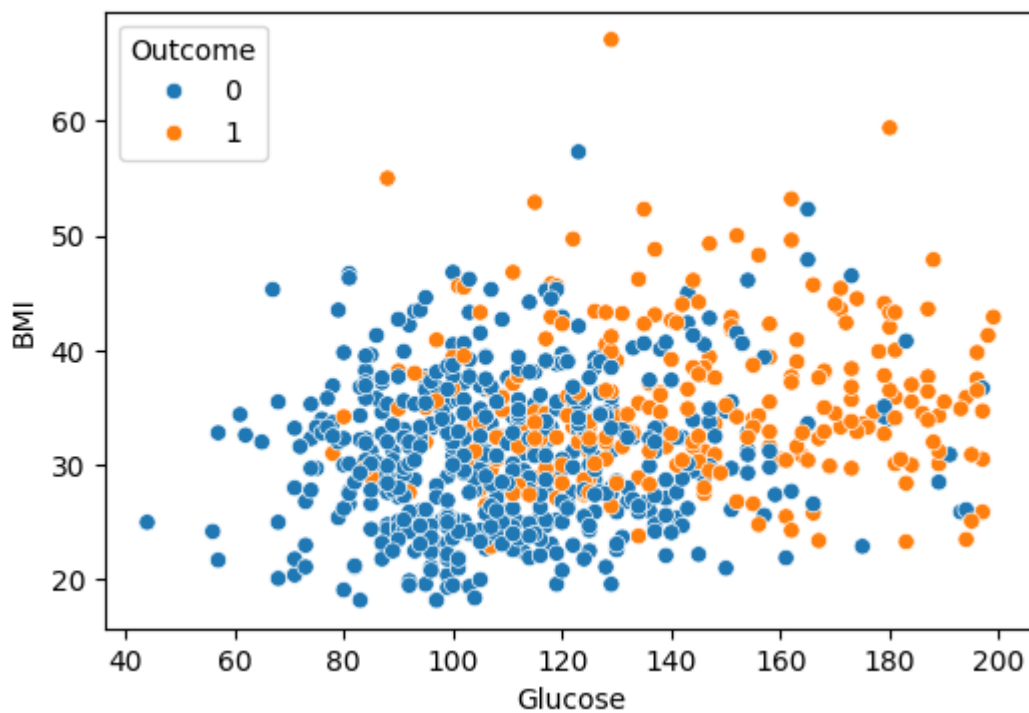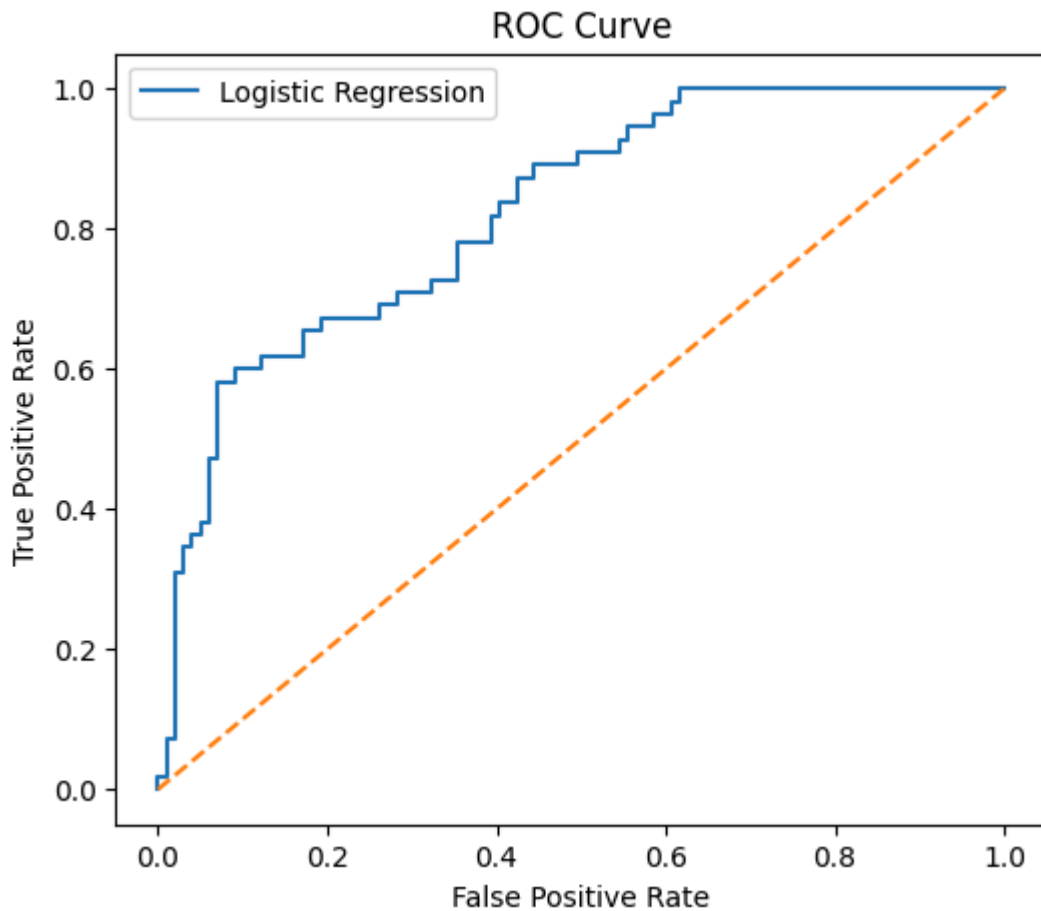
Correlation Heatmap:



```
1 print("\nScatter Plot: Glucose vs BMI")
2 plt.figure(figsize=(6, 4))
3 sns.scatterplot(x=df["Glucose"], y=df["BMI"], hue=df["Outcome"])
4 plt.show()
```

Scatter Plot: Glucose vs BMI



```
1 plt.figure(figsize=(6, 5))
2 plt.plot(fpr, tpr, label="Logistic Regression")
3 plt.plot([0, 1], [0, 1], linestyle="--")
4 plt.xlabel("False Positive Rate")
5 plt.ylabel("True Positive Rate")
6 plt.title("ROC Curve")
7 plt.legend()
8 plt.show()
```

ROC Curve

```
1 print("\nGenerating Learning Curve for SVM...")
2
3 train_sizes, train_scores, test_scores = learning_curve(
4     SVC(), X_scaled, y, cv=5, scoring='accuracy'
5 )
```

Generating Learning Curve for SVM...

```
1 train_mean = train_scores.mean(axis=1)
2 test_mean = test_scores.mean(axis=1)
```
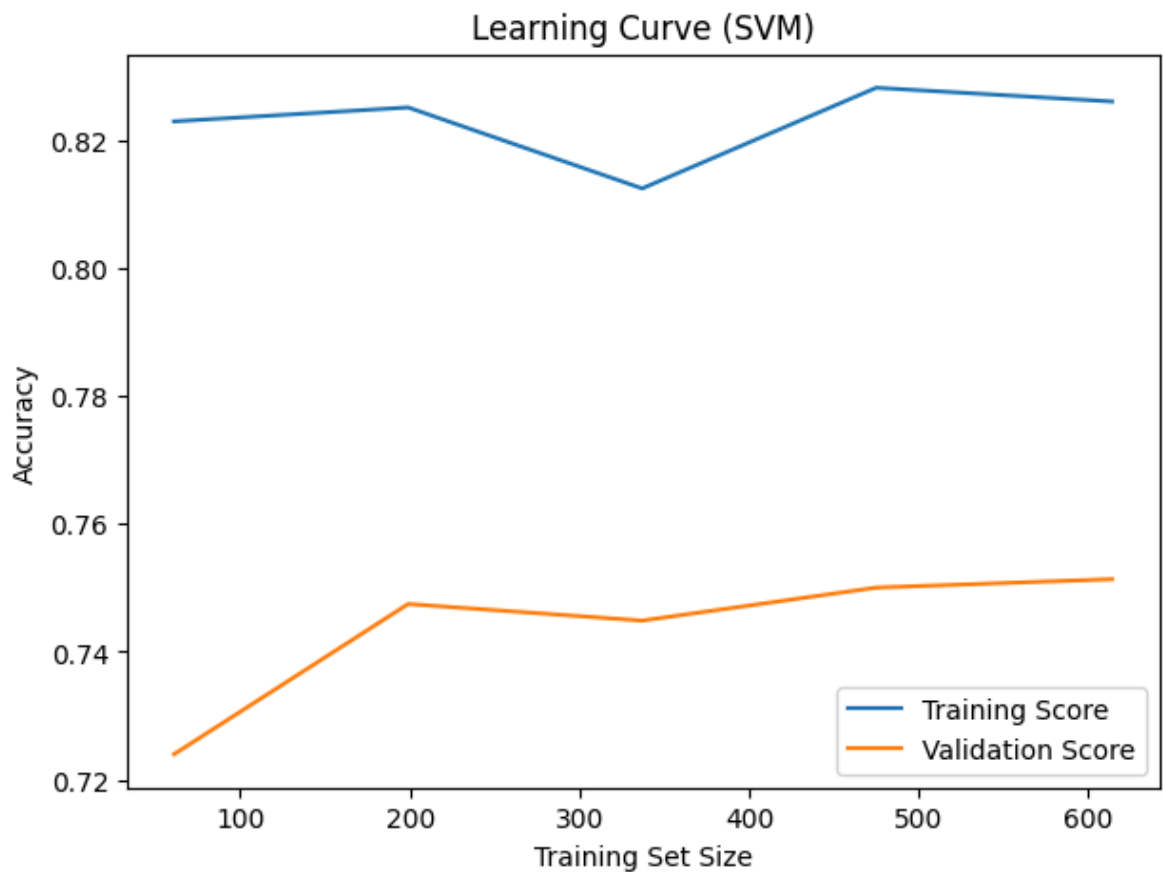
```
1 plt.figure(figsize=(7, 5))
2 plt.plot(train_sizes, train_mean, label="Training Score")
3 plt.plot(train_sizes, test_mean, label="Validation Score")
4 plt.xlabel("Training Set Size")
5 plt.ylabel("Accuracy")
6 plt.title("Learning Curve (SVM)")
7 plt.legend()
8 plt.show()
```
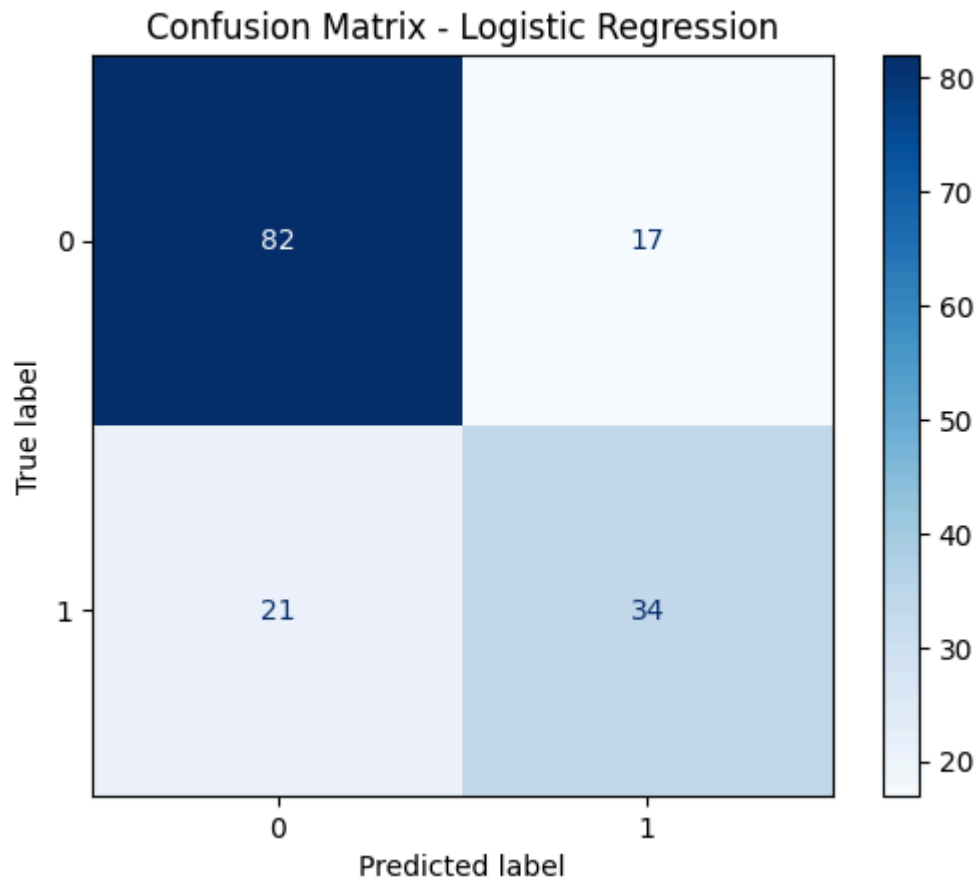
## Learning Curve (SVM)



```python
1  log_model = LogisticRegression(max_iter=300)
2  log_model.fit(X_train, y_train)
3
4  y_pred_log = log_model.predict(X_test)
5
6  print("Confusion Matrix – Logistic Regression:")
7  cm_log = confusion_matrix(y_test, y_pred_log)
8  print(cm_log)
9
10 disp = ConfusionMatrixDisplay(confusion_matrix=cm_log)
11 disp.plot(cmap="Blues")
12 plt.title("Confusion Matrix – Logistic Regression")
13 plt.show()
14
15 print("\nClassification Report – Logistic Regression:")
16 print(classification_report(y_test, y_pred_log))
17
```

```
Confusion Matrix — Logistic Regression:
[[82 17]
 [21 34]]
```

## Confusion Matrix - Logistic Regression



```
Classification Report — Logistic Regression:
              precision    recall  f1-score   support

           0       0.80      0.83      0.81        99
           1       0.67      0.62      0.64        55

    accuracy                           0.75       154
   macro avg       0.73      0.72      0.73       154
weighted avg       0.75      0.75      0.75       154
```
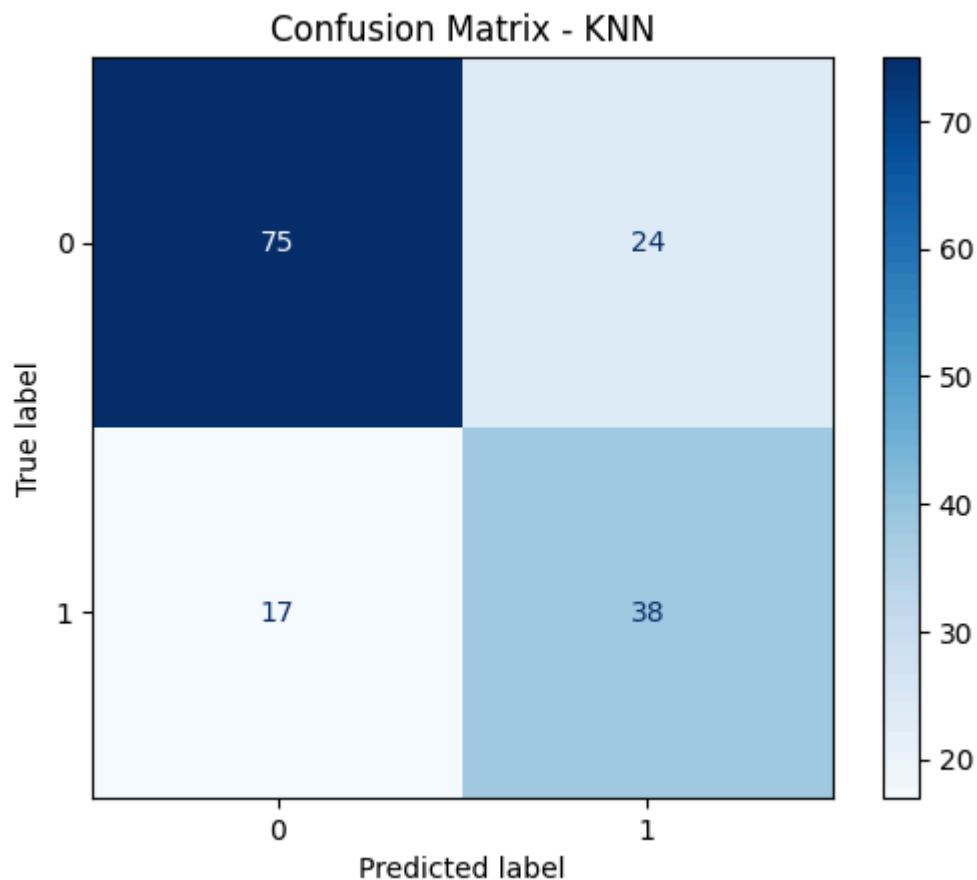
```
 1 knn_model = KNeighborsClassifier()
 2 knn_model.fit(X_train, y_train)
 3
 4 y_pred_knn = knn_model.predict(X_test)
 5
 6 print("Confusion Matrix — KNN:")
 7 cm_knn = confusion_matrix(y_test, y_pred_knn)
 8 print(cm_knn)
 9
10 disp = ConfusionMatrixDisplay(confusion_matrix=cm_knn)
11 disp.plot(cmap="Blues")
12 plt.title("Confusion Matrix — KNN")
13 plt.show()
14
15 print("\nClassification Report — KNN:")
```

```
16 print(classification_report(y_test, y_pred_knn))
17
```

```
Confusion Matrix – KNN:
[[75 24]
 [17 38]]
```



Confusion Matrix - KNN

```
Classification Report – KNN:
              precision    recall  f1-score   support

           0       0.82      0.76      0.79        99
           1       0.61      0.69      0.65        55

    accuracy                           0.73       154
   macro avg       0.71      0.72      0.72       154
weighted avg       0.74      0.73      0.74       154
```

```
 1 dt_model = DecisionTreeClassifier()
 2 dt_model.fit(X_train, y_train)
 3
 4 y_pred_dt = dt_model.predict(X_test)
 5
 6 print("Confusion Matrix – Decision Tree:")
 7 cm_dt = confusion_matrix(y_test, y_pred_dt)
 8 print(cm_dt)
 9
10 disp = ConfusionMatrixDisplay(confusion_matrix=cm_dt)
11 disp.plot(cmap="Greens")
12 plt.title("Confusion Matrix – Decision Tree")
13 plt.show()
```

```
14
15 print("\nClassification Report – Decision Tree:")
16 print(classification_report(y_test, y_pred_dt))
```
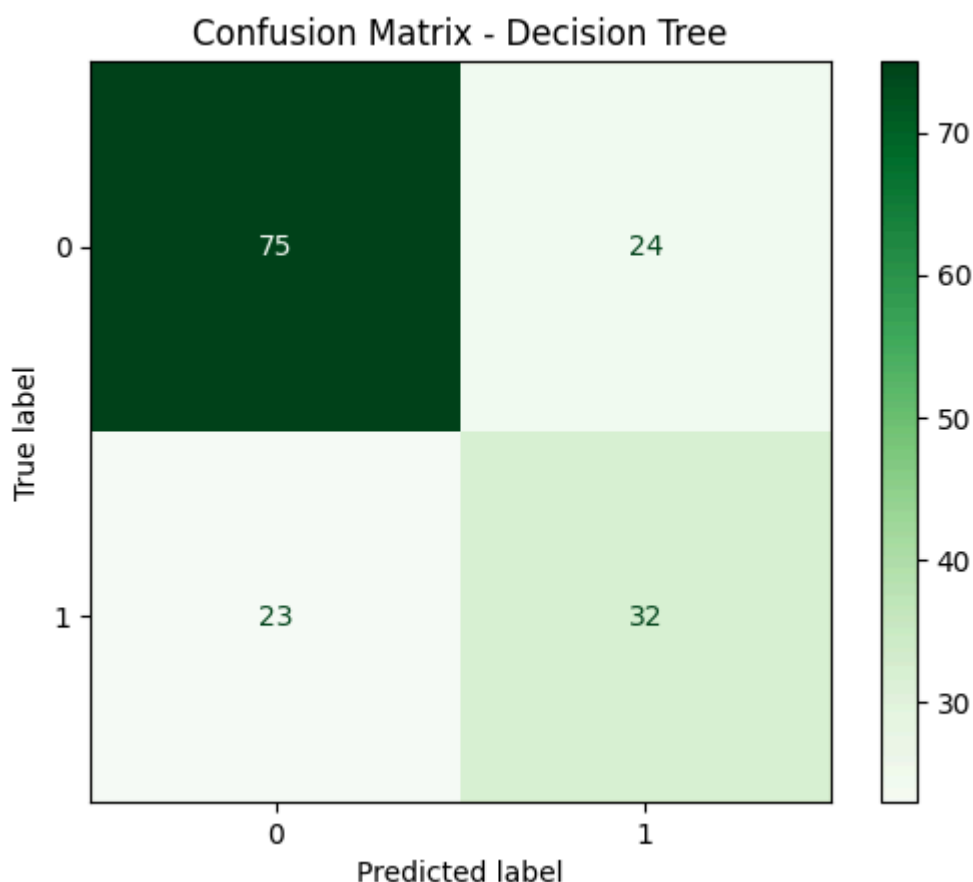
```
Confusion Matrix – Decision Tree:
[[75 24]
 [23 32]]
```



Confusion Matrix - Decision Tree

```
Classification Report – Decision Tree:
              precision    recall  f1-score   support

           0       0.77      0.76      0.76        99
           1       0.57      0.58      0.58        55

    accuracy                           0.69       154
   macro avg       0.67      0.67      0.67       154
weighted avg       0.70      0.69      0.70       154
```
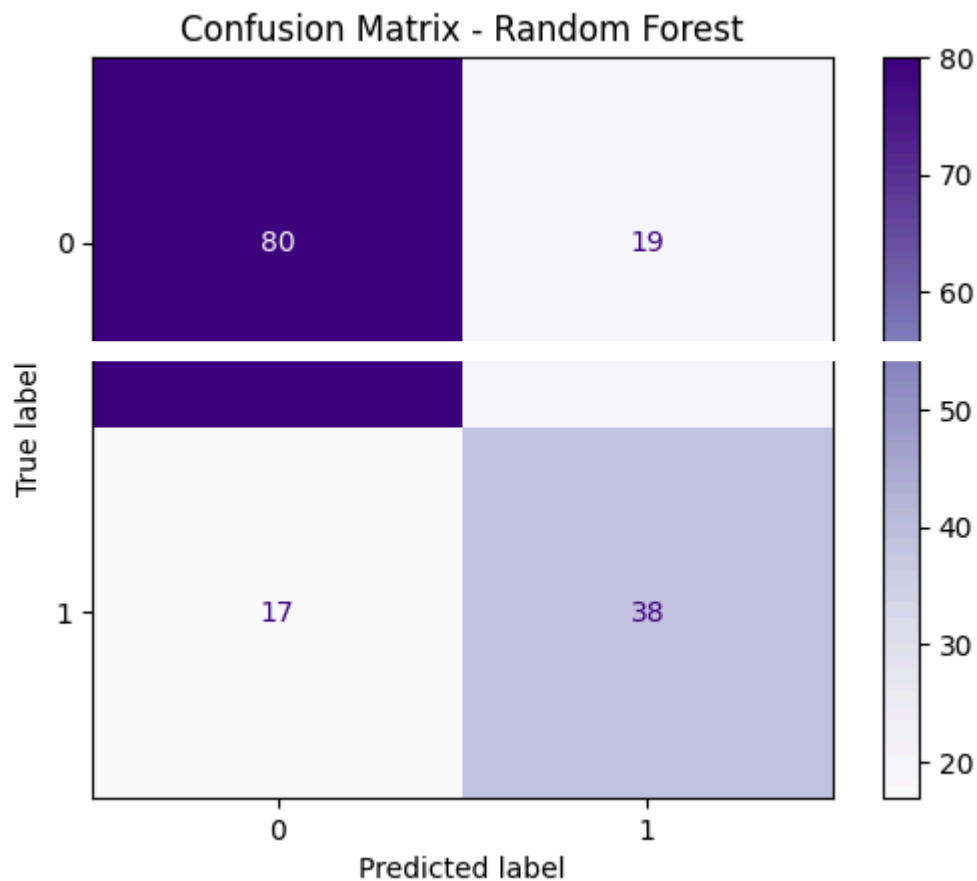
```
1
2 rf_model = RandomForestClassifier()
3 rf_model.fit(X_train, y_train)
4
5 y_pred_rf = rf_model.predict(X_test)
6
7 print("Confusion Matrix – Random Forest:")
8 cm_rf = confusion_matrix(y_test, y_pred_rf)
9 print(cm_rf)
10
11 disp = ConfusionMatrixDisplay(confusion_matrix=cm_rf)
12 disp.plot(cmap="Purples")
```

```
13 plt.title("Confusion Matrix – Random Forest")
14 plt.show()
15
16 print("\nClassification Report – Random Forest:")
```

```
Confusion Matrix – Random Forest:
[[80 19]
 [17 38]]
```

## Confusion Matrix - Random Forest



```
Classification Report – Random Forest:
              precision    recall  f1–score   support

           0       0.82      0.81      0.82        99
           1       0.67      0.69      0.68        55

    accuracy                           0.77       154
```