

```
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
from sklearn.linear_model import LogisticRegression
from sklearn.tree import DecisionTreeClassifier
from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score, classification_report, confusion_matrix, roc
import seaborn as sns
import matplotlib.pyplot as plt
```

```
df = pd.read_csv("/content/diabetes.csv")
print(df.head())
```

```

      Pregnancies  Glucose  BloodPressure  SkinThickness  Insulin   BMI   \
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

      DiabetesPedigreeFunction  Age  Outcome
0              0.627         50         1
1              0.351         31         0
2              0.672         32         1
3              0.167         21         0
4              2.288         33         1
```

```
print(df.info())
print("Missing values:\n", df.isnull().sum())
```

```
<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
None
Missing values:
Pregnancies      0
Glucose           0
BloodPressure     0
SkinThickness     0
Insulin           0
BMI               0
DiabetesPedigreeFunction 0
Age               0
Outcome           0
dtype: int64
```

```
print(df.describe())
sns.countplot(x="Outcome", data=df, palette="Set2")
plt.title("Count of Diabetes Outcome")
plt.show()
```

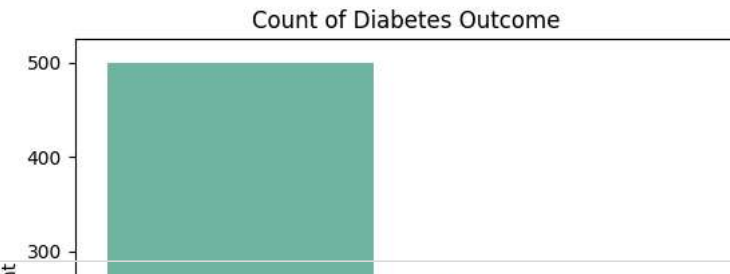
	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin \
count	768.000000	768.000000	768.000000	768.000000	768.000000
mean	3.845052	120.894531	69.105469	20.536458	79.799479
std	3.369578	31.972618	19.355807	15.952218	115.244002
min	0.000000	0.000000	0.000000	0.000000	0.000000
25%	1.000000	99.000000	62.000000	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

	BMI	DiabetesPedigreeFunction	Age	Outcome
count	768.000000	768.000000	768.000000	768.000000
mean	31.992578	0.471876	33.240885	0.348958
std	7.884160	0.331329	11.760232	0.476951
min	0.000000	0.078000	21.000000	0.000000
25%	27.300000	0.243750	24.000000	0.000000
50%	32.000000	0.372500	29.000000	0.000000
75%	36.600000	0.626250	41.000000	1.000000
max	67.100000	2.420000	81.000000	1.000000

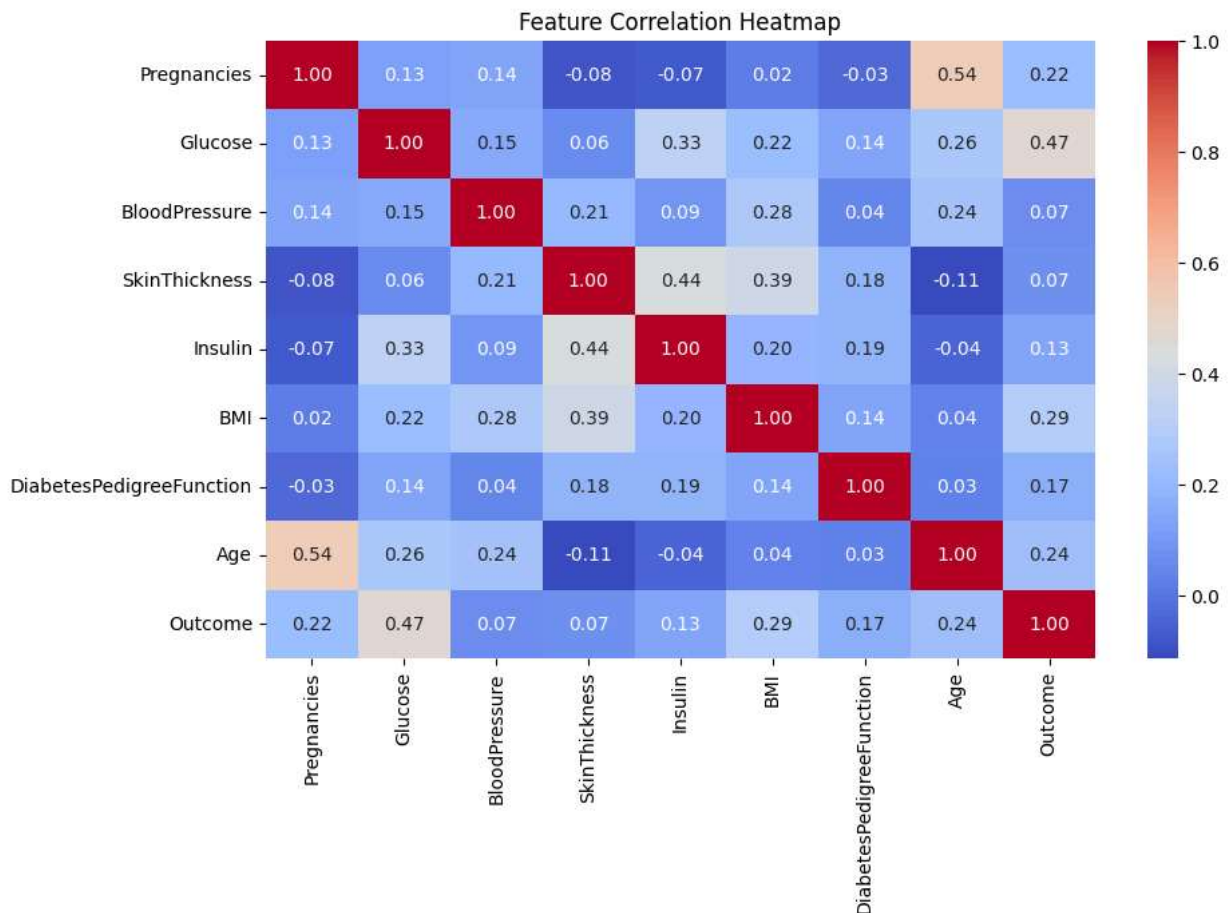
/tmp/ipython-input-3212167674.py:2: FutureWarning:

Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.0. Assign the `x` variable to `hue` and set `l`

```
sns.countplot(x="Outcome", data=df, palette="Set2")
```



```
plt.figure(figsize=(10,6))
sns.heatmap(df.corr(), annot=True, cmap="coolwarm", fmt=".2f")
plt.title("Feature Correlation Heatmap")
plt.show()
```



```
X = df.drop("Outcome", axis=1)
y = df["Outcome"]
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42, stratify=y)
```

```
scaler = StandardScaler()
X_train_scaled = scaler.fit_transform(X_train)
X_test_scaled = scaler.transform(X_test)
```

```
log_model = LogisticRegression(max_iter=1000).fit(X_train_scaled, y_train)
tree_model = DecisionTreeClassifier(max_depth=5, random_state=42).fit(X_train, y_train)
y_pred_log = log_model.predict(X_test_scaled)
y_pred_tree = tree_model.predict(X_test)
```

```
def evaluate_model(y_true, y_pred, name):
    print(f"\n{name} Evaluation")
    print("Accuracy :", accuracy_score(y_true, y_pred))
    print("Precision:", precision_score(y_true, y_pred))
    print("Recall   :", recall_score(y_true, y_pred))
    print("F1-score :", f1_score(y_true, y_pred))
    print("\nClassification Report:\n", classification_report(y_true, y_pred))

evaluate_model(y_test, y_pred_log, "Logistic Regression")
evaluate_model(y_test, y_pred_tree, "Decision Tree")
```

```
Logistic Regression Evaluation
Accuracy : 0.7142857142857143
Precision: 0.6086956521739131
Recall   : 0.5185185185185185
F1-score : 0.56
```

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

     0       0.76       0.82       0.79       100
     1       0.61       0.52       0.56        54

   accuracy          0.68
  macro avg       0.68       0.67       0.67       154
 weighted avg       0.71       0.71       0.71       154
```

```
Decision Tree Evaluation
Accuracy : 0.7922077922077922
Precision: 0.7037037037037037
Recall   : 0.7037037037037037
F1-score : 0.7037037037037037
```

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

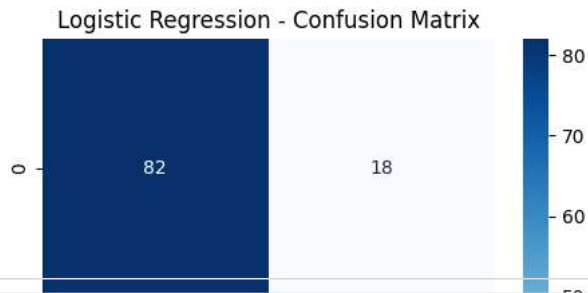
     0       0.84       0.84       0.84       100
     1       0.70       0.70       0.70        54

   accuracy          0.79
  macro avg       0.77       0.77       0.77       154
 weighted avg       0.79       0.79       0.79       154
```

```
plt.figure(figsize=(12,5))

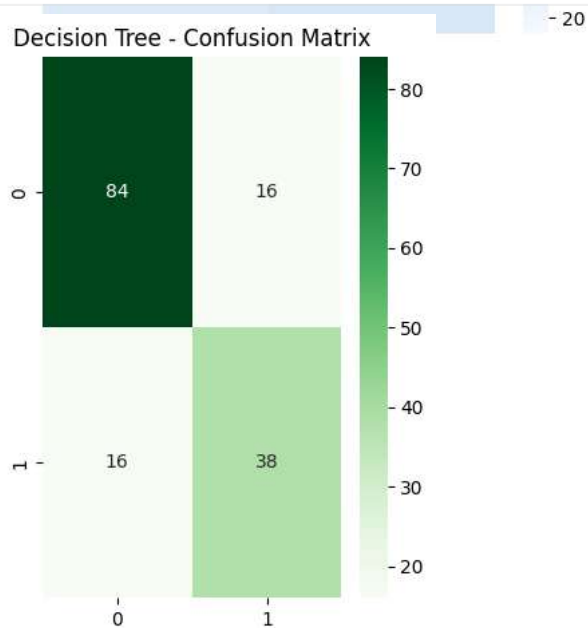
plt.subplot(1,2,1)
sns.heatmap(confusion_matrix(y_test, y_pred_log), annot=True, fmt="d", cmap="Blues")
plt.title("Logistic Regression - Confusion Matrix")
```

```
Text(0.5, 1.0, 'Logistic Regression - Confusion Matrix')
```



```
plt.subplot(1,2,2)
sns.heatmap(confusion_matrix(y_test, y_pred_tree), annot=True, fmt="d", cmap="Greens")
plt.title("Decision Tree - Confusion Matrix")

plt.tight_layout()
plt.show()
```



```
y_prob_log = log_model.predict_proba(X_test_scaled)[:,-1]
fpr, tpr, _ = roc_curve(y_test, y_prob_log)
roc_auc = auc(fpr, tpr)

plt.plot(fpr, tpr, color="blue", label=f"Logistic Regression (AUC = {roc_auc:.2f})")
plt.plot([0,1], [0,1], linestyle="--", color="gray")
plt.xlabel("False Positive Rate")
plt.ylabel("True Positive Rate")
plt.title("ROC Curve")
plt.legend()
plt.show()
```

