

Step-1: Project Setup & Dataset Download

```
!pip install pandas numpy scikit-learn matplotlib seaborn

Requirement already satisfied: pandas in /usr/local/lib/python3.12/dist-packages (2.2.2)
Requirement already satisfied: numpy in /usr/local/lib/python3.12/dist-packages (2.0.2)
Requirement already satisfied: scikit-learn in /usr/local/lib/python3.12/dist-packages (1.6.1)
Requirement already satisfied: matplotlib in /usr/local/lib/python3.12/dist-packages (3.10.0)
Requirement already satisfied: seaborn in /usr/local/lib/python3.12/dist-packages (0.13.2)
Requirement already satisfied: python-dateutil>=2.8.2 in /usr/local/lib/python3.12/dist-packages (from pandas) (2.9.0.post0)
Requirement already satisfied: pytz>=2020.1 in /usr/local/lib/python3.12/dist-packages (from pandas) (2025.2)
Requirement already satisfied: tzdata>=2022.7 in /usr/local/lib/python3.12/dist-packages (from pandas) (2025.2)
Requirement already satisfied: scipy>=1.6.0 in /usr/local/lib/python3.12/dist-packages (from scikit-learn) (1.16.3)
Requirement already satisfied: joblib>=1.2.0 in /usr/local/lib/python3.12/dist-packages (from scikit-learn) (1.5.2)
Requirement already satisfied: threadpoolctl>=3.1.0 in /usr/local/lib/python3.12/dist-packages (from scikit-learn) (3.6.0)
Requirement already satisfied: contourpy>=1.0.1 in /usr/local/lib/python3.12/dist-packages (from matplotlib) (1.3.3)
Requirement already satisfied: cycler>=0.10 in /usr/local/lib/python3.12/dist-packages (from matplotlib) (0.12.1)
Requirement already satisfied: fonttools>=4.22.0 in /usr/local/lib/python3.12/dist-packages (from matplotlib) (4.60.1)
Requirement already satisfied: kiwisolver>=1.3.1 in /usr/local/lib/python3.12/dist-packages (from matplotlib) (1.4.9)
Requirement already satisfied: packaging>=20.0 in /usr/local/lib/python3.12/dist-packages (from matplotlib) (25.0)
Requirement already satisfied: pillow>=8 in /usr/local/lib/python3.12/dist-packages (from matplotlib) (11.3.0)
Requirement already satisfied: pyparsing>=2.3.1 in /usr/local/lib/python3.12/dist-packages (from matplotlib) (3.2.5)
Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.12/dist-packages (from python-dateutil>=2.8.2->pandas) (1.17.0)
```

Import Libraries

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
```

Upload dataset

```
from google.colab import files
uploaded=files.upload()
```

Choose files No file chosen

Upload widget is only available when the cell has been executed in the current browser session. Please rerun this cell to enable.

Load Dataset

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

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age	Outcome
0	6	148	72	35	0	33.6	0.627	50	1
1	1	85	66	29	0	26.6	0.351	31	0
2	8	183	64	0	0	23.3	0.672	32	1
3	1	89	66	23	94	28.1	0.167	21	0
4	0	137	40	35	168	43.1	2.288	33	1

Step-2: Exploratory Data Analysis(EDA)

2.1- Basic Info about Data

```
#Check shape (rows and columns)
print("Shape of dataset:", df.shape)

#Display column names
print("\nColumn names")
print(df.columns)

#Summary info
print("\nInfo:")
df.info()

#Basic statistics
```

```
print("\nSummary Statistics:")
df.describe()
```

Shape of dataset: (768, 9)

Column names

```
Index(['Pregnancies', 'Glucose', 'BloodPressure', 'SkinThickness', 'Insulin',
       'BMI', 'DiabetesPedigreeFunction', 'Age', 'Outcome'],
      dtype='object')
```

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
```

Summary Statistics:

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age	
count	768.000000	768.000000	768.000000	768.000000	768.000000	768.000000	768.000000	768.000000	768
mean	3.845052	120.894531	69.105469	20.536458	79.799479	31.992578	0.471876	33.240885	
std	3.369578	31.972618	19.355807	15.952218	115.244002	7.884160	0.331329	11.760232	
min	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.078000	21.000000	
25%	1.000000	99.000000	62.000000	0.000000	0.000000	27.300000	0.243750	24.000000	
50%	3.000000	117.000000	72.000000	23.000000	30.500000	32.000000	0.372500	29.000000	
75%	6.000000	140.250000	80.000000	32.000000	127.250000	36.600000	0.626250	41.000000	
max	17.000000	199.000000	122.000000	99.000000	846.000000	67.100000	2.420000	81.000000	

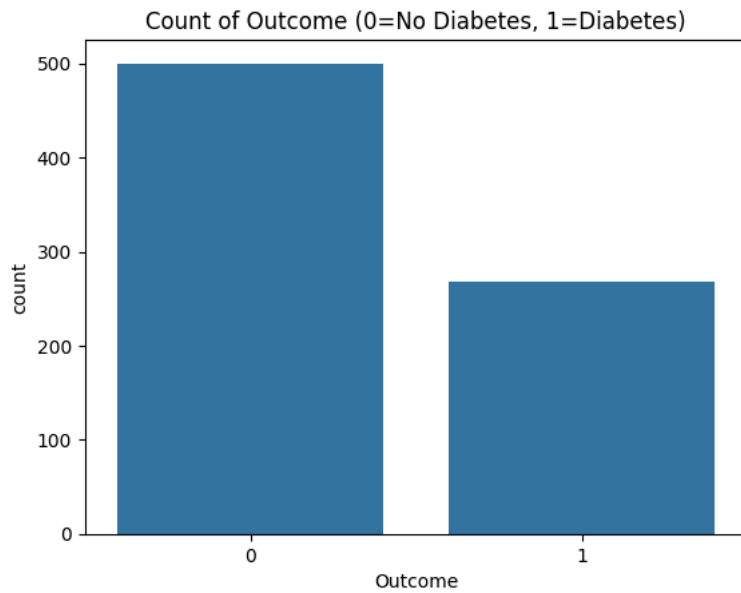
2.2-Check Missing Value

```
#Check for null values
print(df.isnull().sum())
```

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

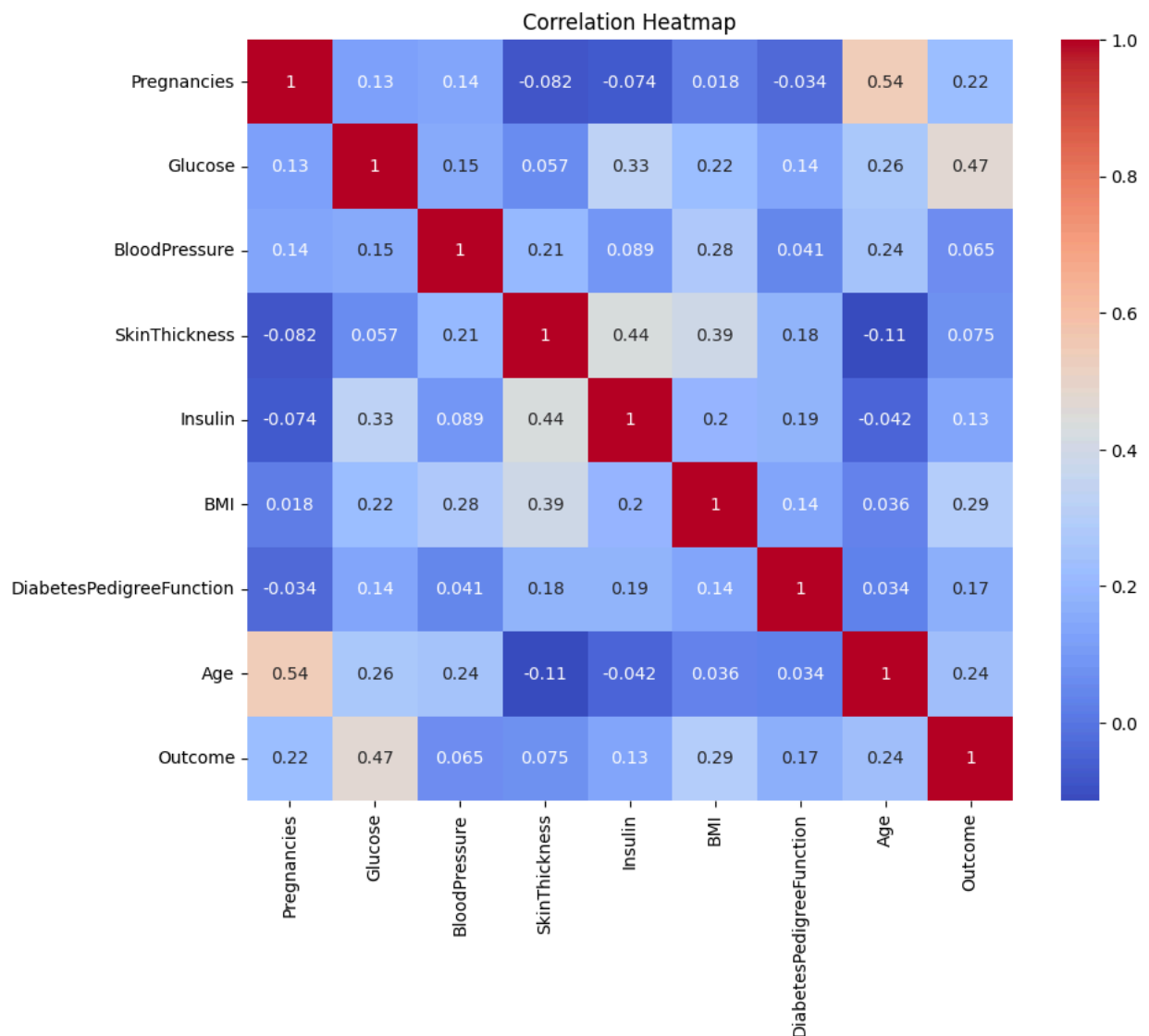
2.3-Understand the Target column

```
#Count of Diabetic vs non-Diabetic
sns.countplot(x='Outcome',data=df)
plt.title("Count of Outcome (0=No Diabetes, 1=Diabetes)")
plt.show()
```



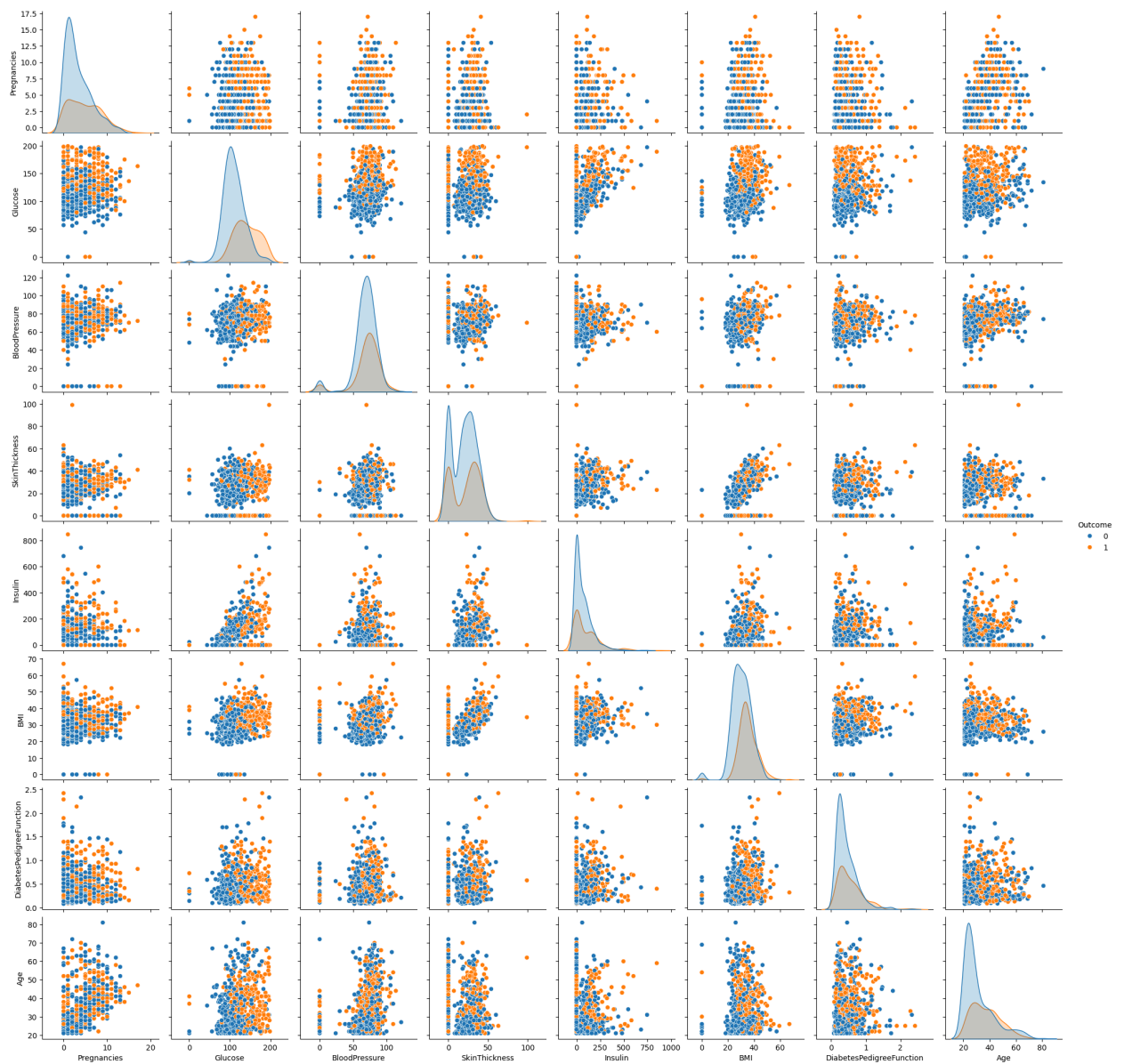
2.4-Check Correlations

```
plt.figure(figsize=(10,8))
sns.heatmap(df.corr(), annot=True, cmap="coolwarm")
plt.title("Correlation Heatmap")
plt.show()
```



2.5-Pairplot for Relationship

```
sns.pairplot(df, hue="Outcome")
plt.show()
```



✓ Step-3:Data preprocessing & Splitting

3.1-Split Features & Target

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

3.2-Split into Train & Test Sets

```
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split( X, y, test_size=0.2, random_state=42)
print("Training set size:", X_train.shape)
print("Testing set size:", X_test.shape)
```

```
Training set size: (614, 8)
Testing set size: (154, 8)
```

3.3-Feature Scaling

```
from sklearn.preprocessing import StandardScaler

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

▼ Step-4:Model Training

4.1-Import and Train the Model

```
from sklearn.linear_model import LogisticRegression

#Create model
model=LogisticRegression()

#Train model
model.fit(X_train, y_train)
```

```
▼ LogisticRegression ⓘ ?
LogisticRegression()
```

4.2-Make predictions

```
#Predict on test data
y_pred=model.predict(X_test)

#Show firts few predictions
print("Prediction", y_pred[:10])
```

```
Prediction [0 0 0 0 0 0 0 1 1 1]
```

4.3-Evaluate Accuracy

```
from sklearn.metrics import accuracy_score, confusion_matrix, classification_report

#Calculate accuracy
acc=accuracy_score(y_test, y_pred)
print("Accuracy:", acc)

#Confusion Matrix
print("\nConfusion Matrix:\n", confusion_matrix(y_test, y_pred))

#Classification Report
print("\nClassification Report:\n", classification_report(y_test, y_pred))
```

```
Accuracy: 0.7532467532467533
```

```
Confusion Matrix:
[[79 20]
 [18 37]]
```

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

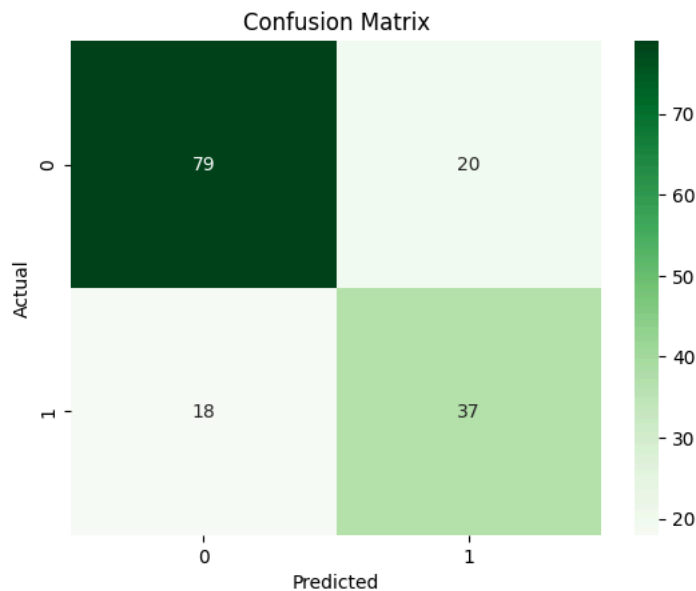
     0       0.81      0.80      0.81        99
     1       0.65      0.67      0.66        55

 accuracy          0.75
 macro avg          0.73
weighted avg          0.76
```

4.4-Visualize the Confusion Matrix

```
sns.heatmap(confusion_matrix(y_test, y_pred), annot=True, fmt='d', cmap='Greens')
plt.title("Confusion Matrix")
plt.xlabel("Predicted")
```

```
plt.ylabel("Actual")
plt.show()
```



▼ Step-5: Save the Model & Make New Predictions

5.1-Save the Model

```
import joblib

#Save model to file
joblib.dump(model, "diabetes_model.pkl")

print("Model saved successfully")
```

Model saved successfully

5.2-Load the Model Again

```
#Load the saved model
loaded_model = joblib.load("diabetes_model.pkl")
```

5.3-Predict for New Input

```
#Example patient data: [Pregnancies, Glucose, BloodPressure, SkinThickness, Insulin, BMI, DiabetesPedigreeFunction, Age]
sample = np.array([[2,120,70,20,79,25.0,0.5,32]])

#Scale it (important)
sample_scaled = scaler.transform(sample)

#Predict
prediction = loaded_model.predict(sample_scaled)
print("Prediction", "Diabetic" if prediction[0] == 1 else "Non-Diabetic")
```

Prediction Non-Diabetic
 /usr/local/lib/python3.12/dist-packages/sklearn/utils/validation.py:2739: UserWarning: X does not have valid feature names, warnings.warn()

▼ Diabetes_model.pkl

```
from google.colab import files
files.download("diabetes_model.pkl")
```

▼ requirements.txt

```
!pip freeze > requirements.txt
```

```
!cat requirements.txt
```

```
absl-py==1.4.0
absolutify-imports==0.3.1
accelerate==1.11.0
aiofiles==24.1.0
aiohappyeyeballs==2.6.1
aiohttp==3.13.1
aiosignal==1.4.0
alabaster==1.0.0
albucore==0.0.24
albumintations==2.0.8
ale-py==0.11.2
alembic==1.17.0
altair==5.5.0
annotated-doc==0.0.3
annotated-types==0.7.0
antlr4-python3-runtime==4.9.3
anyio==4.11.0
anywidget==0.9.18
argon2-cffi==25.1.0
argon2-cffi-bindings==25.1.0
array_record==0.8.2
arrow==1.4.0
arviz==0.22.0
astropy==7.1.1
astropy-iers-data==0.2025.10.27.0.39.10
astunparse==1.6.3
atpublic==5.1
attrs==25.4.0
audioread==3.1.0
Authlib==1.6.5
autograd==1.8.0
babel==2.17.0
backcall==0.2.0
beartype==0.22.4
beautifulsoup4==4.13.5
betterproto==2.0.0b6
bigframes==2.27.0
bigquery-magics==0.10.3
bleach==6.3.0
blinker==1.9.0
blis==1.3.0
blobfile==3.1.0
blosc2==3.11.0
bokeh==3.7.3
Bottleneck==1.4.2
bqplot==0.12.45
branca==0.8.2
Brotli==1.1.0
build==1.3.0
CacheControl==0.14.3
cachetools==5.5.2
catalogue==2.0.10
certifi==2025.10.5
cffi==2.0.0
chardet==5.2.0
charset-normalizer==3.4.4
chex==0.1.90
clashol==0.11.1
```

```
from google.colab import files
files.download("requirements.txt")
```