



Assessment Report
on
“Diabetes Prediction”
submitted as partial fulfilment for the award of
BACHELOR OF TECHNOLOGY DEGREE
SESSION 2024-25

in
CSE (AI)

By:

Group Members:-

[Samarth Shukla] 202401100300212
[Sagar Shrivastava] 202401100300209
[Satwik Shaw] 202401100300217
[Prateek Kumar] 202401100300178
[Satyam Tiwari] 202401100300219

Section: [C]
Under the supervision of
[Mayank Lakhtoia]

KIET Group of Institutions, Ghaziabad
May, 2025

1. Introduction

As the prevalence of diabetes continues to rise, especially in developing nations, early prediction through data-driven healthcare is becoming increasingly vital. This project addresses the problem of predicting the likelihood of diabetes using supervised machine learning techniques. Utilizing a dataset of diagnostic attributes such as glucose level, BMI, age, and insulin levels, the model aims to support clinical decision-making through predictive analysis.

2. Problem Statement

To predict whether a patient is likely to be diagnosed with diabetes using available medical and personal health data. This classification task aims to assist healthcare providers in early diagnosis and preventive care.

3. Objectives

- Preprocess the dataset for training a machine learning model.
- Train a Logistic Regression model to classify diabetic and non-diabetic cases.
- Evaluate model performance using standard classification metrics.
- Visualize the confusion matrix using a heatmap for interpretability.

4. Methodology

Data Collection: The user uploads a CSV file containing the dataset (Pima Indians Diabetes Dataset).

Data Preprocessing:

- Handling missing values using mean imputation for columns with zero entries.
- Feature scaling using StandardScaler.

Model Building:

- Splitting the dataset into training and testing sets.
- Training a Logistic Regression classifier.

Model Evaluation:

- Evaluating accuracy, precision, recall, and F1-score.
- Generating a confusion matrix and visualizing it with a heatmap.

5. Data Preprocessing

The dataset is cleaned and prepared as follows:

- Missing numerical values in fields such as Glucose and BloodPressure are replaced with column means.
- Data is scaled using StandardScaler to normalize feature values.
- The dataset is split into 80% training and 20% testing.

6. Model Implementation

Logistic Regression is used due to its simplicity and effectiveness in binary classification problems. The model is trained on the processed dataset and used to predict diabetes status on the test set.

7. Evaluation Metrics

The following metrics are used to evaluate the model:

- **Accuracy:** Measures overall correctness.
- **Precision:** Indicates the proportion of predicted positives that are actual positives.
- **Recall:** Shows the proportion of actual positives that were correctly identified.
- **F1 Score:** Harmonic mean of precision and recall.
- **Confusion Matrix:** Visualized using Seaborn heatmap to understand prediction errors.

8. Results and Analysis

- The model provided reasonable performance on the test set.
- Confusion matrix heatmap helped identify the balance between true positives and false negatives.
- Precision and recall indicated how well the model detected diabetic patients versus false negatives.

9. Conclusion

The logistic regression model successfully classified diabetic and non-diabetic cases with satisfactory performance metrics. The project demonstrates the potential of using machine learning for automating healthcare risk assessments. However, improvements can be made by exploring more advanced models and handling data imbalance.

10. References

- scikit-learn documentation
- pandas documentation
- Seaborn visualization library
- Research articles on diabetes prediction

```

import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.model_selection import train_test_split
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import accuracy_score, confusion_matrix, ConfusionMatrixDisplay, classification_report
import time

# Load dataset
df = pd.read_csv('diabetes_new.csv')

# === EDA: Distribution Plots for each feature ===
numeric_columns = df.select_dtypes(include=['float64', 'int64']).columns.tolist()
numeric_columns.remove('Outcome') # Exclude target column

for col in numeric_columns:
    plt.figure(figsize=(8, 4))
    sns.histplot(df[col], bins=20, kde=True, color='skyblue', edgecolor='black')
    plt.title(f'{col} Distribution')
    plt.xlabel(col)
    plt.ylabel('Frequency')
    plt.grid(True)
    plt.tight_layout()
    plt.show()

# === EDA: Correlation Heatmap ===
plt.figure(figsize=(10, 8))
corr = df.corr()
sns.heatmap(corr, annot=True, fmt=".2f", cmap='coolwarm', square=True)
plt.title("Feature Correlation Heatmap")
plt.tight_layout()
plt.show()

```

```

# === Model Training ===
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)

start_train = time.time()
model = RandomForestClassifier(random_state=42)
model.fit(X_train, y_train)
end_train = time.time()

start_pred = time.time()
y_pred = model.predict(X_test)
end_pred = time.time()

# === Evaluation ===
accuracy = accuracy_score(y_test, y_pred)
cm = confusion_matrix(y_test, y_pred)
report = classification_report(y_test, y_pred)

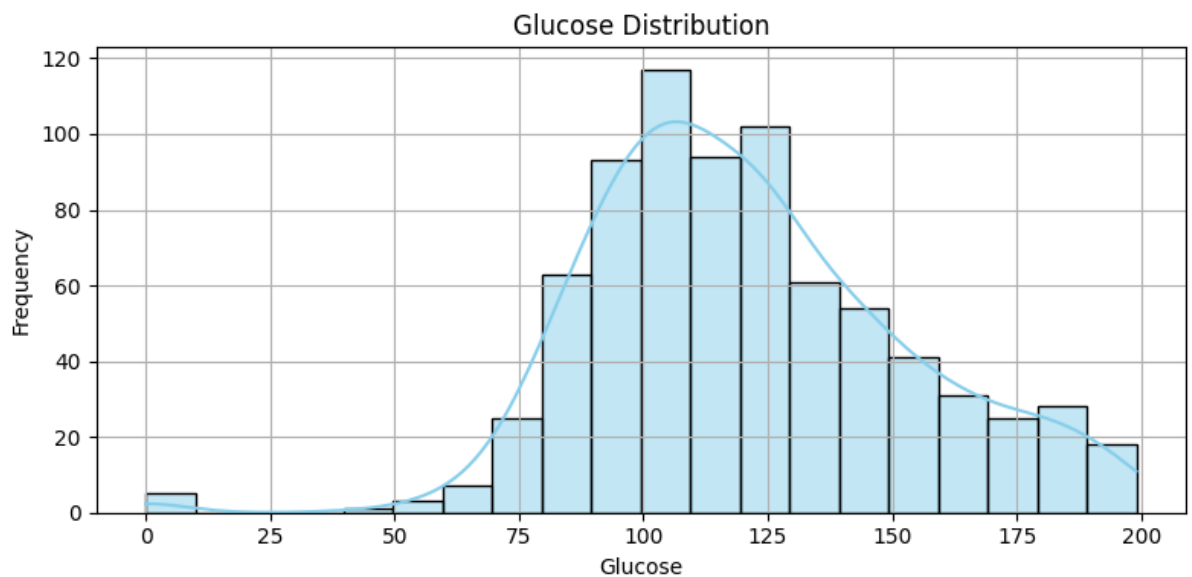
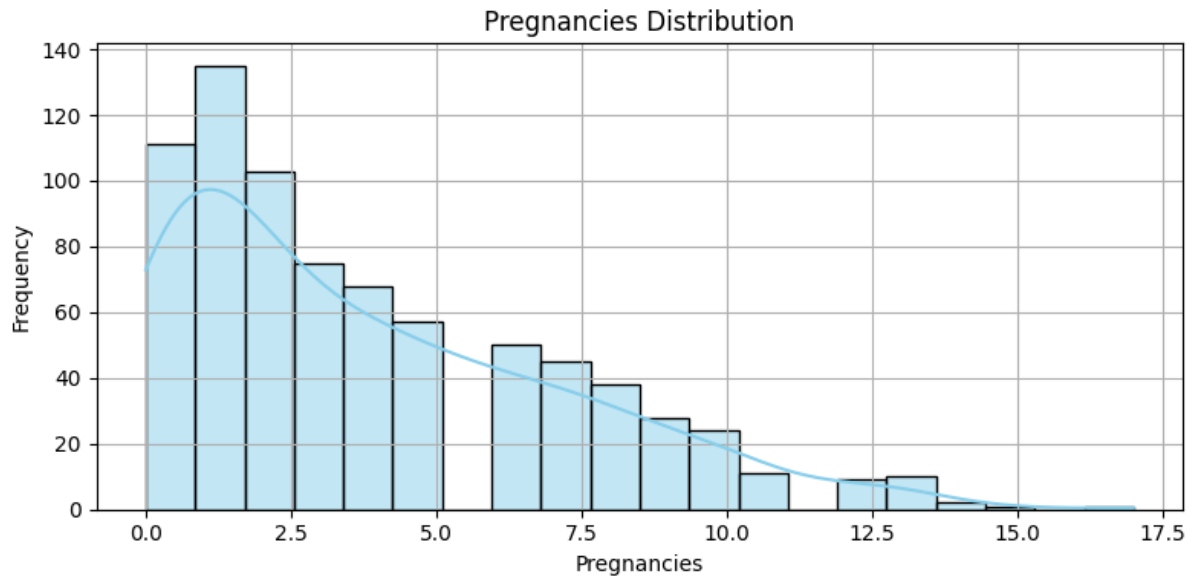
print(f"✅ Model Accuracy: {accuracy:.2f}")
print(f"🕒 Training time: {end_train - start_train:.4f} seconds")
print(f"🕒 Prediction time: {end_pred - start_pred:.4f} seconds")
print(f"\n📄 Classification Report:")
print(report)

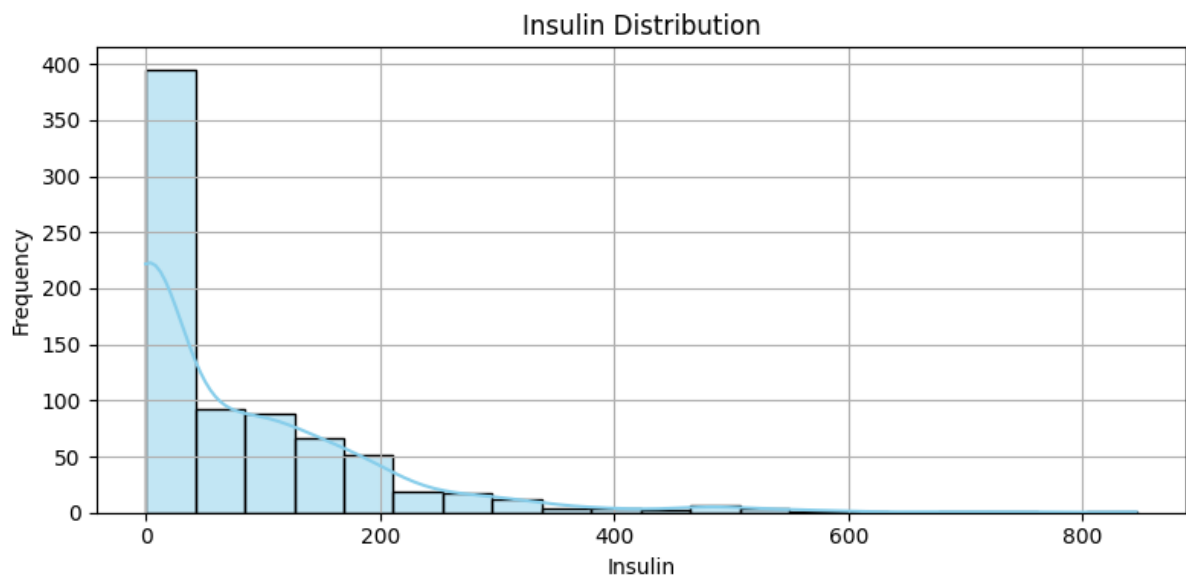
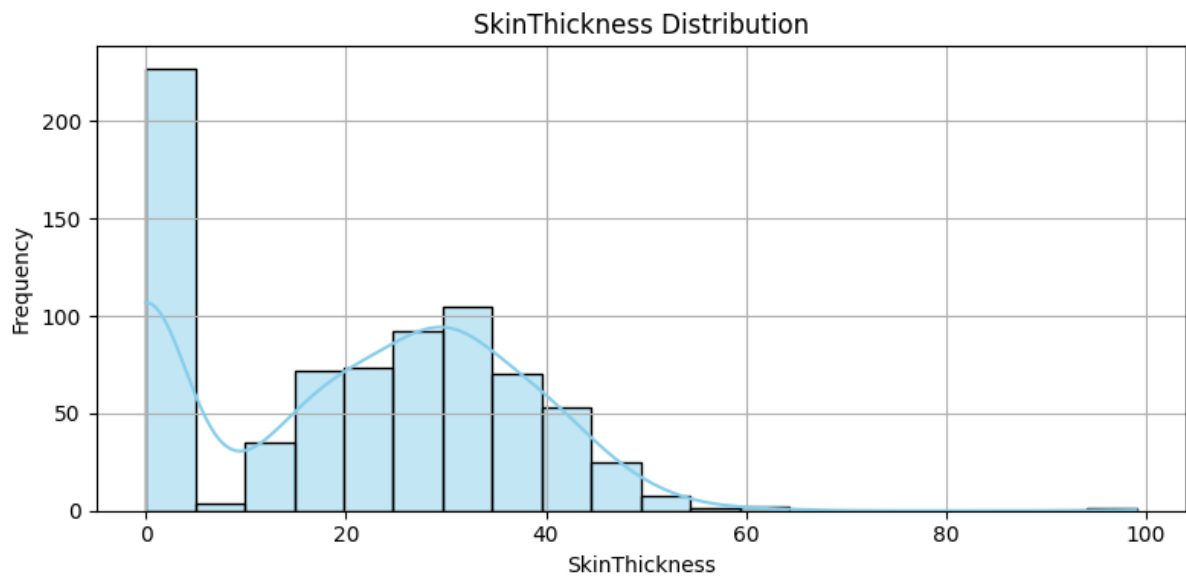
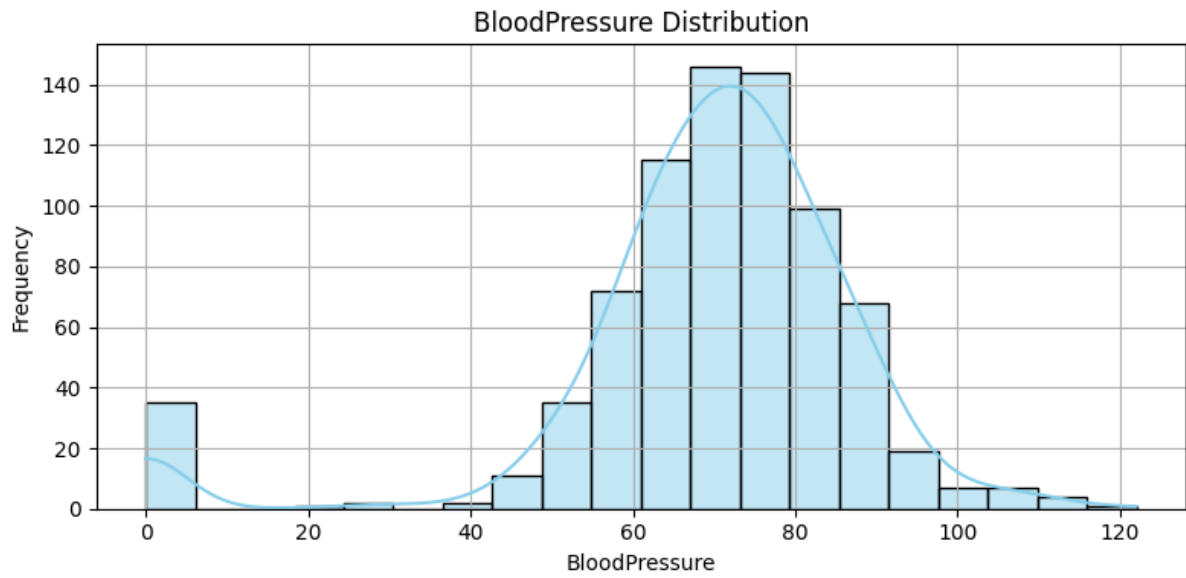
# === Confusion Matrix Plot ===
disp = ConfusionMatrixDisplay(confusion_matrix=cm, display_labels=model.classes_)
disp.plot(cmap='Blues', values_format='d')
plt.title(f'Confusion Matrix (Accuracy: {accuracy:.2f})')
plt.grid(False)
plt.tight_layout()
plt.show()

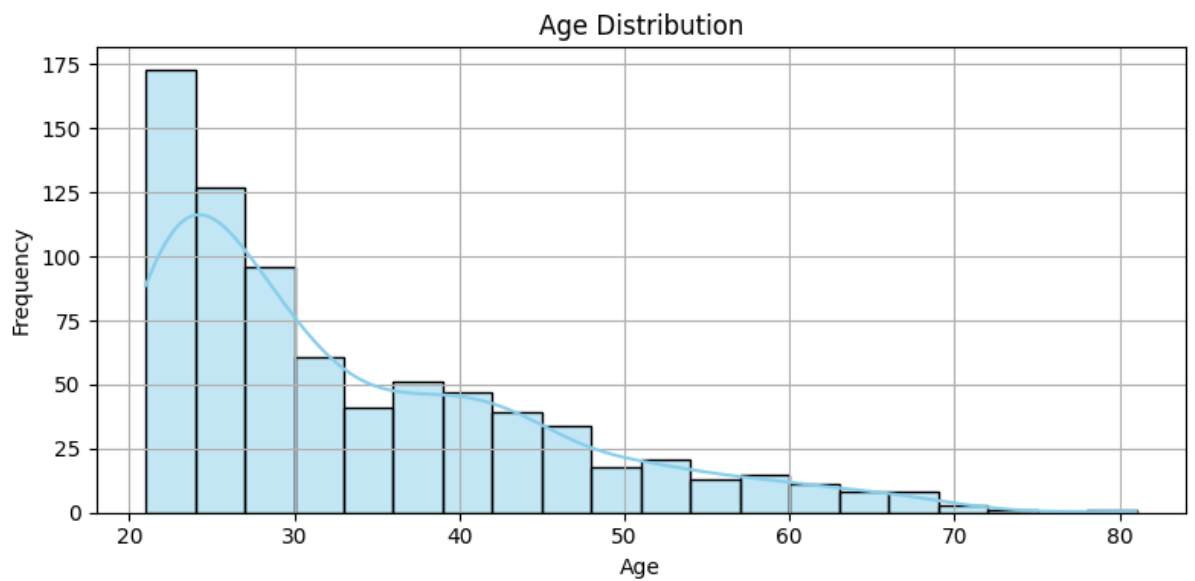
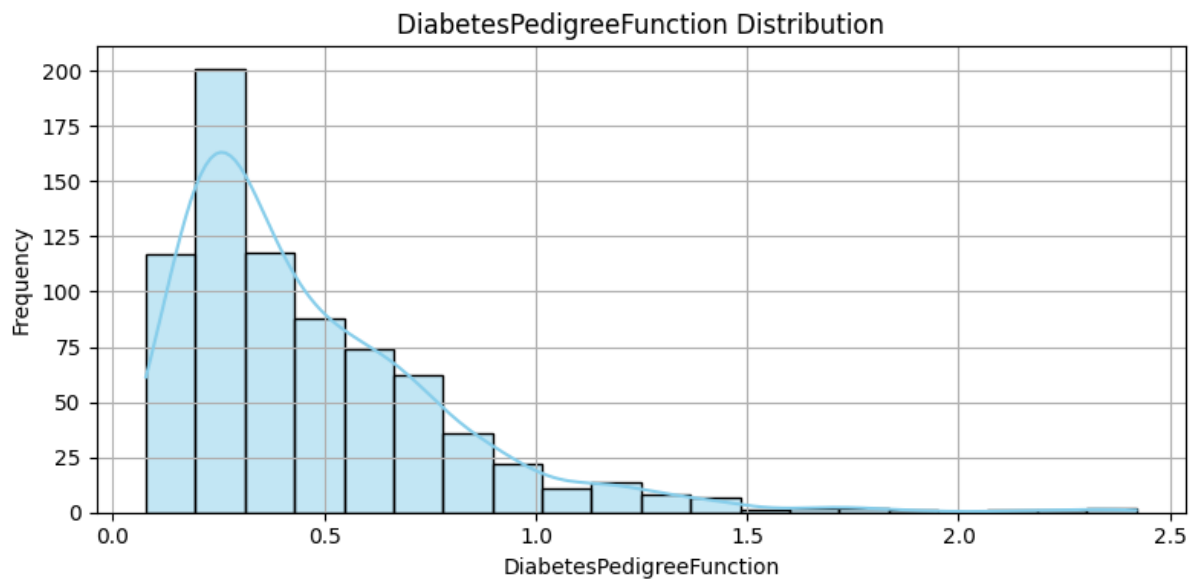
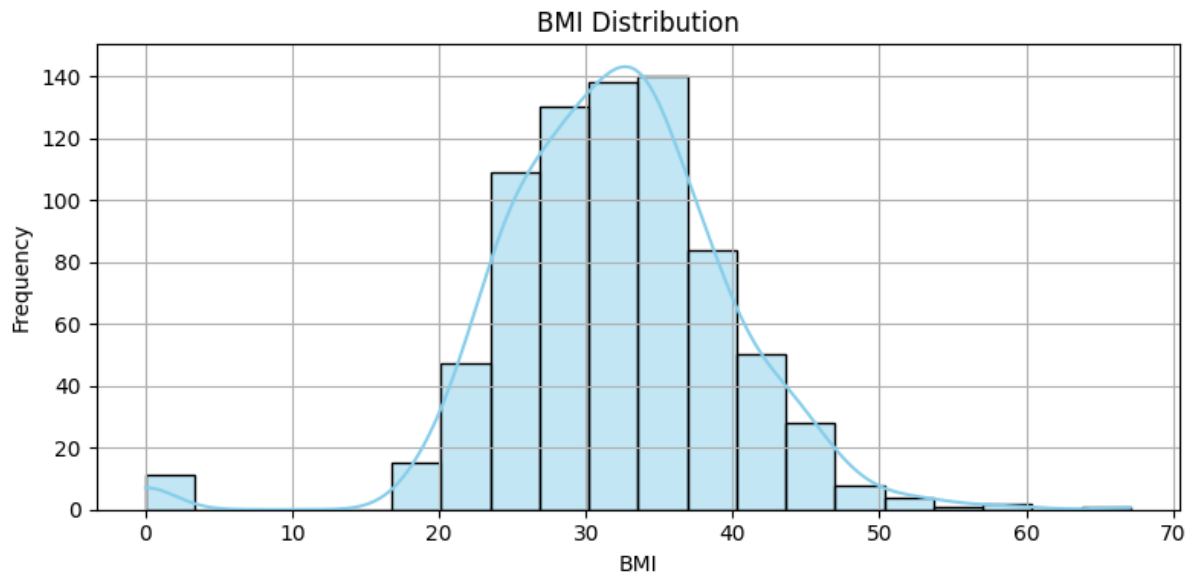
```

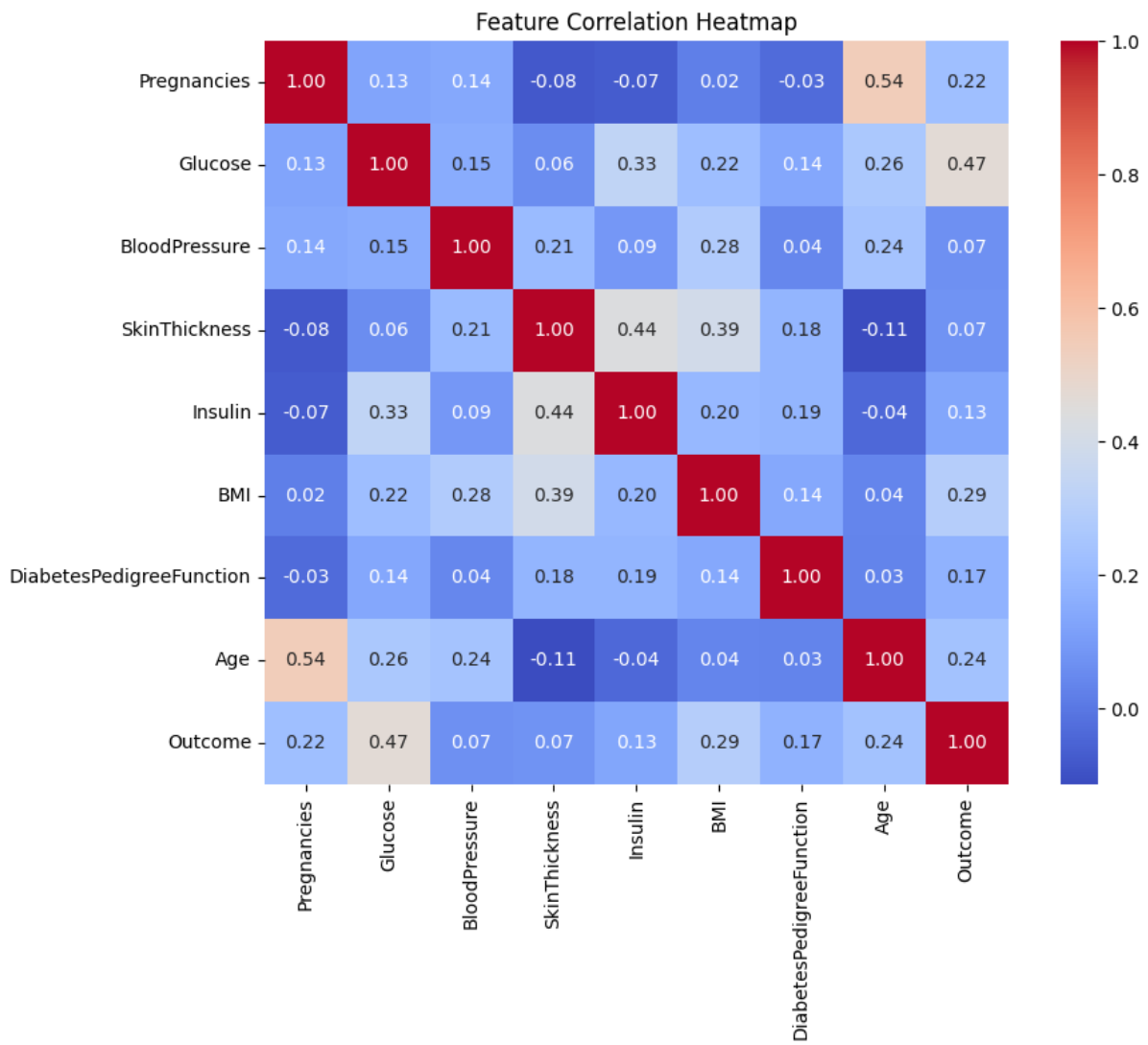
```
# === Feature Importance - Horizontal Bar Chart ===
importances = model.feature_importances_
feat_imp_df = pd.DataFrame({'Feature': X.columns, 'Importance': importances})
feat_imp_df = feat_imp_df.sort_values(by='Importance', ascending=True) # For horizontal chart

plt.figure(figsize=(10, 6))
plt.barh(feat_imp_df['Feature'], feat_imp_df['Importance'], color='blue') # Horizontal bars
plt.xlabel('Importance')
plt.ylabel('Feature')
plt.title('Feature Importance - Random Forest')
plt.tight_layout()
plt.show()
```









✔ Model Accuracy: 0.72

🕒 Training time: 0.2410 seconds

🕒 Prediction time: 0.0106 seconds

📄 Classification Report:

	precision	recall	f1-score	support
0	0.79	0.78	0.78	99
1	0.61	0.62	0.61	55

accuracy			0.72	154
macro avg	0.70	0.70	0.70	154
weighted avg	0.72	0.72	0.72	154

