# PREDICTION MODEL 1

February 13, 2025

```python
[1]:  # Import necessary libraries
      import pandas as pd
      import numpy as np
      import matplotlib.pyplot as plt
      import seaborn as sns
      from sklearn.model_selection import train_test_split
      from sklearn.preprocessing import StandardScaler
      from sklearn.linear_model import LogisticRegression, Lasso
      from sklearn.tree import DecisionTreeClassifier, plot_tree
      from sklearn.metrics import (
          classification_report, confusion_matrix, roc_auc_score, accuracy_score,
          precision_score, recall_score, f1_score, roc_curve
      )
      from sklearn.decomposition import PCA
      from sklearn.ensemble import VotingClassifier
      from imblearn.over_sampling import SMOTE
      import tensorflow as tf
      from tensorflow.keras.models import Sequential
      from tensorflow.keras.layers import Dense

      # 1. Load the dataset
      print("Loading dataset...")
      df = pd.read_csv('diabetes.csv')
      print("\nDataset successfully loaded!")

      # 2. Data Exploration
      print("\nDataset Overview:")
      print(df.head())
      print("\nMissing Values:")
      print(df.isnull().sum())
      print("\nStatistical Summary:")
      print(df.describe())

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

```python
plt.show()

# Histograms for feature distributions
df.hist(figsize=(12, 10), bins=20, edgecolor='black')
plt.suptitle('Feature Distributions')
plt.show()

# Boxplots for key features
plt.figure(figsize=(10, 6))
sns.boxplot(data=df[['Glucose', 'BMI', 'Insulin']])
plt.title('Boxplot of Key Features')
plt.show()

# 3. Data Preprocessing
# Handling missing values
print("\nHandling missing values...")
df.fillna(df.mean(), inplace=True)
print("Missing values handled.\n")

# Splitting features and target
X = df.drop('Outcome', axis=1)
y = df['Outcome']

# Normalizing data
print("Normalizing data...")
scaler = StandardScaler()
X_scaled = scaler.fit_transform(X)

# Addressing class imbalance using SMOTE
print("Addressing class imbalance using SMOTE...")
smote = SMOTE(random_state=42)
X_resampled, y_resampled = smote.fit_resample(X_scaled, y)
print("Class imbalance addressed.\n")

# Train-test split
X_train, X_test, y_train, y_test = train_test_split(
    X_resampled, y_resampled, test_size=0.3, random_state=42
)

# 4. Dimensionality Reduction with PCA
print("Applying PCA...")
pca = PCA(n_components=2)
X_train_pca = pca.fit_transform(X_train)
X_test_pca = pca.transform(X_test)

# Visualizing PCA
plt.figure(figsize=(8, 6))
```

```python
plt.scatter(X_train_pca[:, 0], X_train_pca[:, 1], c=y_train, cmap='coolwarm',␣
 ↪edgecolors='k', alpha=0.7)
plt.title('PCA Visualization of Training Data')
plt.xlabel('Principal Component 1')
plt.ylabel('Principal Component 2')
plt.colorbar(label='Outcome')
plt.show()
print(f'Explained Variance Ratio by PCA Components: {pca.
 ↪explained_variance_ratio_}\n')

# 5. Logistic Regression Model
print("Training Logistic Regression model...")
log_reg = LogisticRegression()
log_reg.fit(X_train, y_train)
y_pred_log_reg = log_reg.predict(X_test)

print("\nLogistic Regression Metrics:")
print(f"Accuracy: {accuracy_score(y_test, y_pred_log_reg)}")
print(f"Precision: {precision_score(y_test, y_pred_log_reg)}")
print(f"Recall: {recall_score(y_test, y_pred_log_reg)}")
print(f"F1 Score: {f1_score(y_test, y_pred_log_reg)}")
print(f"ROC-AUC: {roc_auc_score(y_test, log_reg.predict_proba(X_test)[:, 1])}")
print(confusion_matrix(y_test, y_pred_log_reg))

# 6. Decision Tree Model
print("\nTraining Decision Tree model...")
dt = DecisionTreeClassifier(max_depth=5, min_samples_split=5)
dt.fit(X_train, y_train)
y_pred_dt = dt.predict(X_test)

print("\nDecision Tree Metrics:")
print(f"Accuracy: {accuracy_score(y_test, y_pred_dt)}")
print(f"Precision: {precision_score(y_test, y_pred_dt)}")
print(f"Recall: {recall_score(y_test, y_pred_dt)}")
print(f"F1 Score: {f1_score(y_test, y_pred_dt)}")
print(f"ROC-AUC: {roc_auc_score(y_test, dt.predict_proba(X_test)[:, 1])}")
print(confusion_matrix(y_test, y_pred_dt))

# Visualizing the Decision Tree
plt.figure(figsize=(12, 8))
plot_tree(dt, filled=True, feature_names=df.columns[:-1],␣
 ↪class_names=['Non-diabetic', 'Diabetic'], rounded=True)
plt.title('Decision Tree Visualization')
plt.show()

# 7. Neural Network Model
print("\nTraining Neural Network model...")
```

```python
model = Sequential([
    Dense(16, input_dim=X_train.shape[1], activation='relu'),
    Dense(8, activation='relu'),
    Dense(1, activation='sigmoid')
])
model.compile(optimizer='adam', loss='binary_crossentropy',
 ↪metrics=['accuracy'])
history = model.fit(X_train, y_train, epochs=50, batch_size=10,
 ↪validation_data=(X_test, y_test), verbose=0)

# Neural Network Evaluation
nn_accuracy = model.evaluate(X_test, y_test, verbose=0)[1]
print(f"\nNeural Network Accuracy: {nn_accuracy}\n")

# 8. Ensemble Model: Voting Classifier
print("\nTraining Voting Classifier...")
voting_clf = VotingClassifier(estimators=[('log_reg', log_reg), ('dt', dt)],
 ↪voting='soft')
voting_clf.fit(X_train, y_train)
y_pred_voting = voting_clf.predict(X_test)

print("\nVoting Classifier Metrics:")
print(f"Accuracy: {accuracy_score(y_test, y_pred_voting)}")
print(f"ROC-AUC: {roc_auc_score(y_test, voting_clf.predict_proba(X_test)[:,
 ↪1])}")
print(classification_report(y_test, y_pred_voting))

# 9. Feature Importance Using Lasso Regularization
print("\nApplying Lasso Regularization...")
lasso = Lasso(alpha=0.1)
lasso.fit(X_train, y_train)
lasso_coefficients = pd.Series(lasso.coef_, index=df.columns[:-1])

# Plot Lasso coefficients
plt.figure(figsize=(10, 6))
lasso_coefficients.plot(kind='bar', color='skyblue')
plt.title('Lasso Regularization Feature Importance')
plt.ylabel('Coefficient Value')
plt.show()

# 10. Insights and Conclusions
print("\nKey Insights:")
print("1. Logistic Regression and Decision Tree models highlight Glucose, BMI,
 ↪and Insulin as key predictors.")
print("2. PCA visualization shows good separability between diabetic and
 ↪non-diabetic classes.")
```

```
print("3. Ensemble models provide improved accuracy and robustness.")
print("4. Lasso Regularization identifies the most important features for⏎
 ↪diabetes prediction.")
print("\nAnalysis complete!")
```

Loading dataset…

Dataset successfully loaded!

Dataset Overview:
```
   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
```

Missing Values:
```
Pregnancies                 0
Glucose                     0
BloodPressure               0
SkinThickness               0
Insulin                     0
BMI                         0
DiabetesPedigreeFunction    0
Age                         0
Outcome                     0
dtype: int64
```
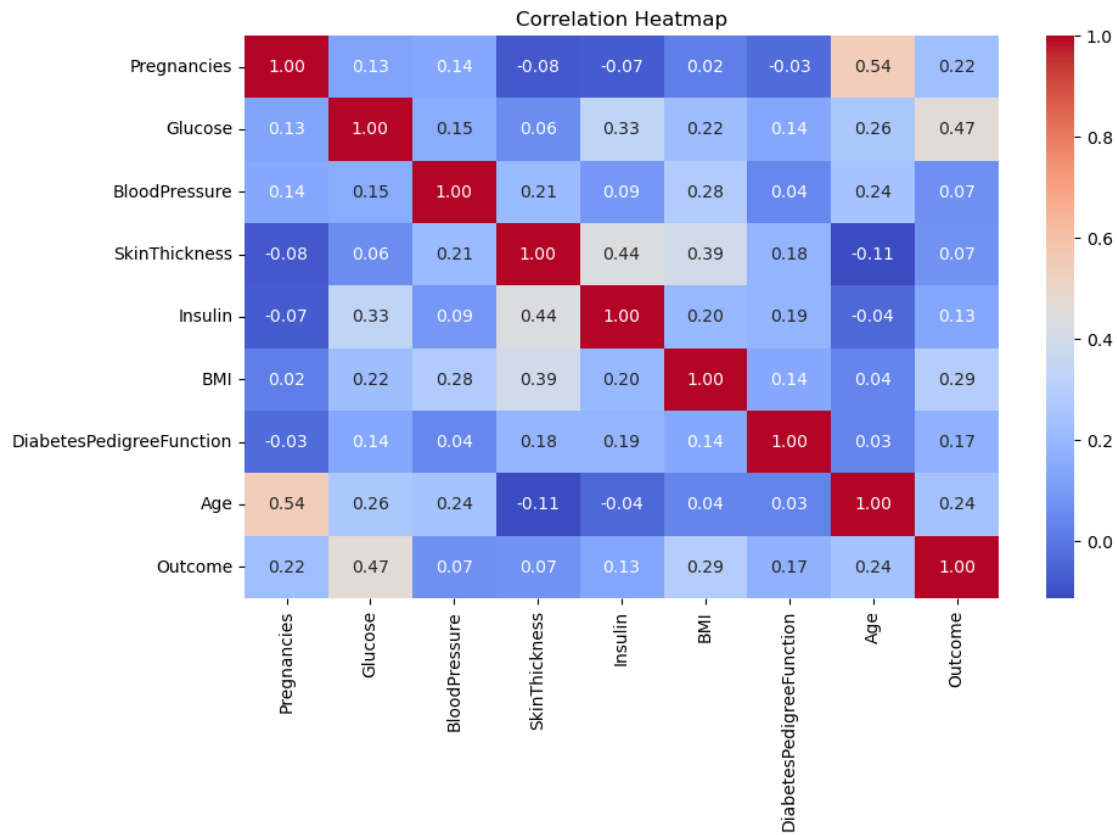
Statistical Summary:
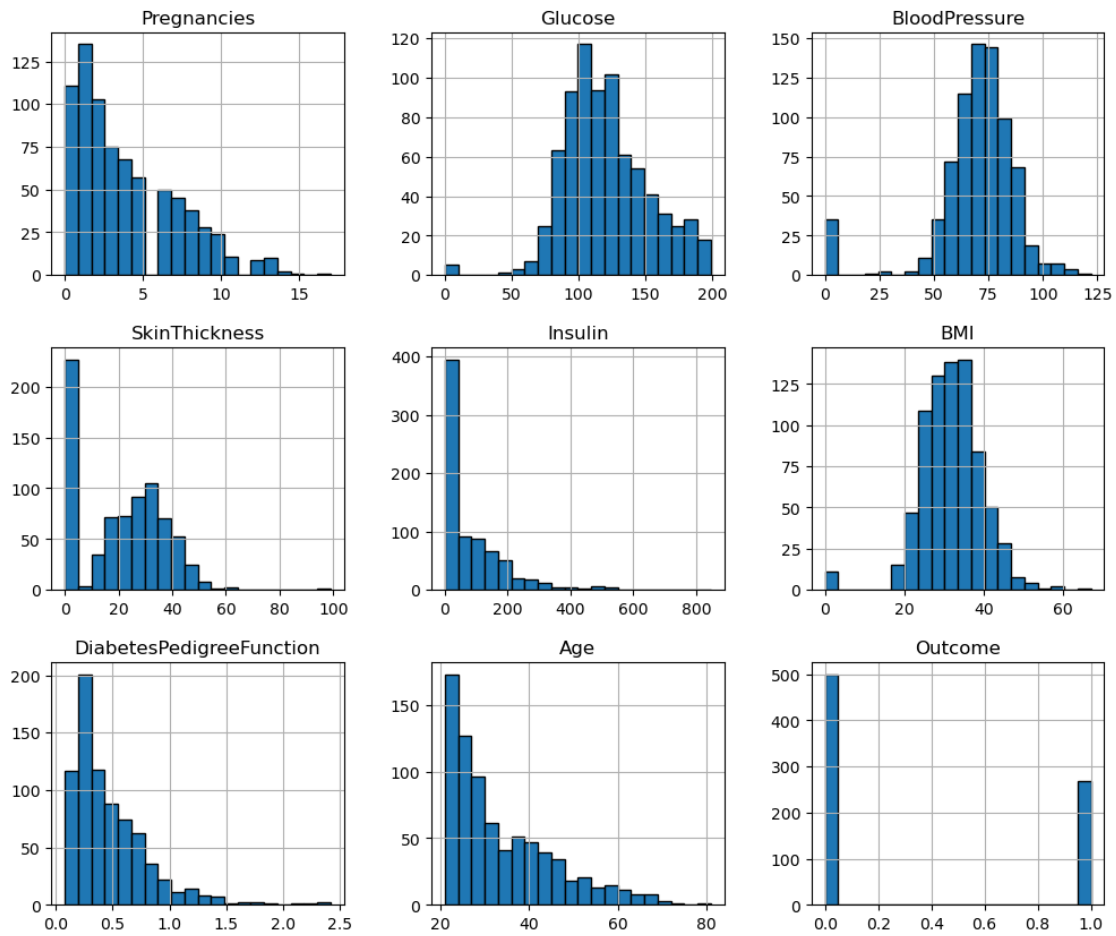```
       Pregnancies     Glucose  BloodPressure  SkinThickness      Insulin  \
count   768.000000  768.000000     768.000000     768.000000   768.000000
mean      3.845052  120.894531      69.105469      20.536458    79.799479
std       3.369578   31.972618      19.355807      15.952218   115.244002
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
```
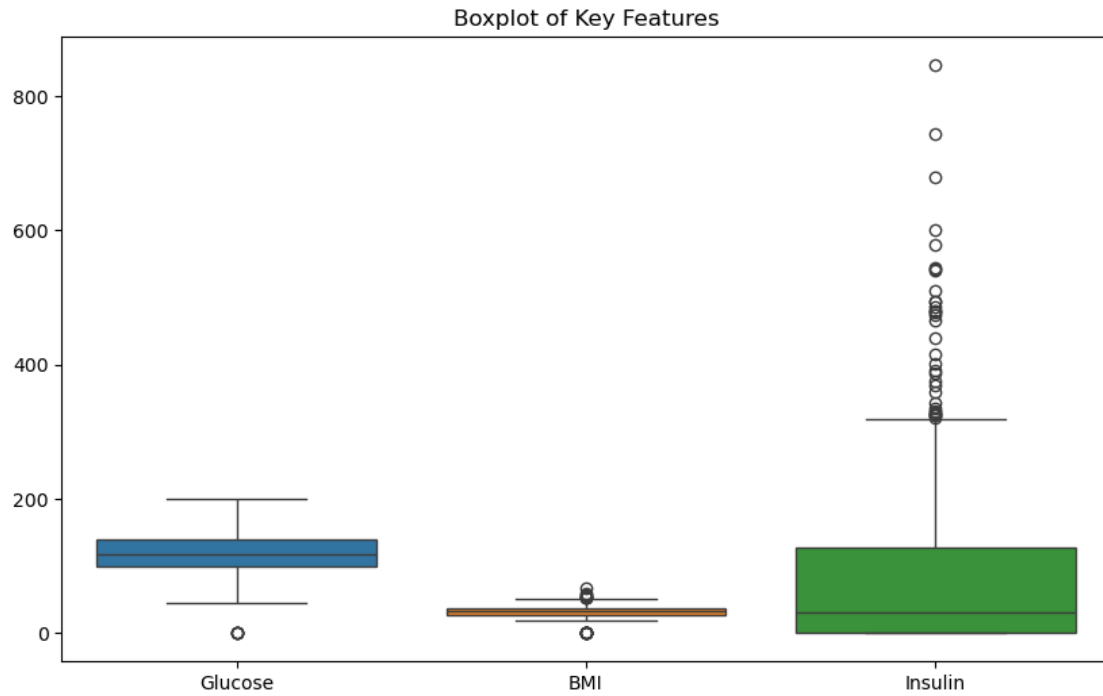
```
count   768.000000              768.000000  768.000000  768.000000
mean     31.992578                0.471876   33.240885    0.348958
std       7.884160                0.331329   11.760232    0.476951
min       0.000000                0.078000   21.000000    0.000000
25%      27.300000                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
```
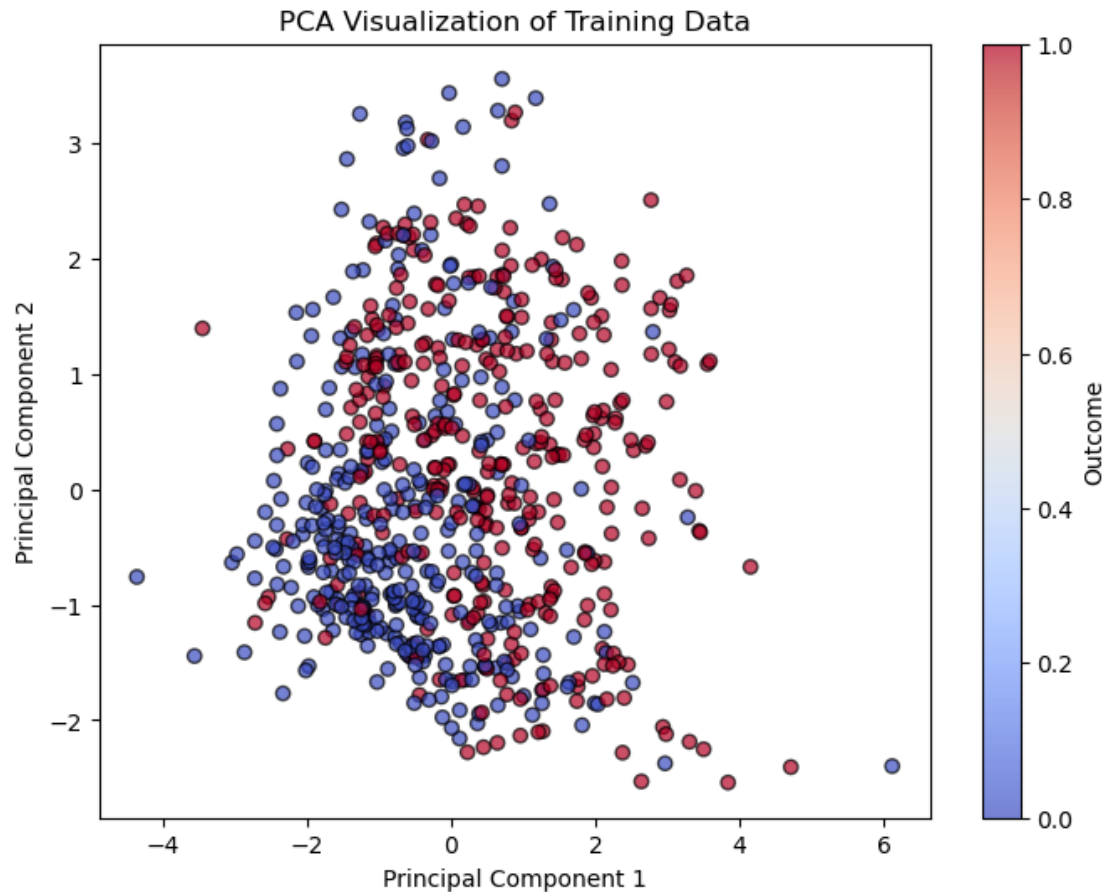


Correlation Heatmap

Feature Distributions

Boxplot of Key Features

Handling missing values…
Missing values handled.

Normalizing data…
Addressing class imbalance using SMOTE…
Class imbalance addressed.

Applying PCA…

## PCA Visualization of Training Data



Explained Variance Ratio by PCA Components: [0.26889982 0.21347626]

Training Logistic Regression model…

Logistic Regression Metrics:
Accuracy: 0.7533333333333333
Precision: 0.7655172413793103
Recall: 0.7350993377483444
F1 Score: 0.75
ROC-AUC: 0.8327036757189209
[[115  34]
 [ 40 111]]

Training Decision Tree model…

Decision Tree Metrics:
Accuracy: 0.7366666666666667
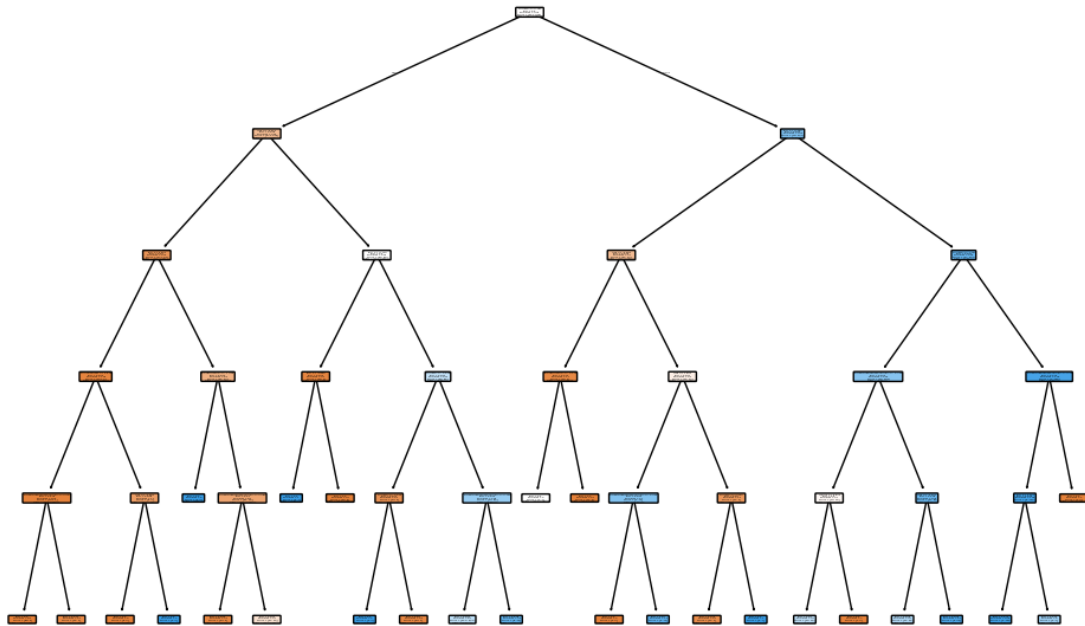Precision: 0.711764705882353
Recall: 0.8013245033112583

```
F1 Score: 0.7538940809968847
ROC-AUC: 0.804080181341393
[[100  49]
 [ 30 121]]
```

Decision Tree Visualization



Training Neural Network model…

```
/opt/anaconda3/lib/python3.12/site-packages/keras/src/layers/core/dense.py:87:
UserWarning: Do not pass an `input_shape`/`input_dim` argument to a layer. When
using Sequential models, prefer using an `Input(shape)` object as the first
layer in the model instead.
  super().__init__(activity_regularizer=activity_regularizer, **kwargs)
```

Neural Network Accuracy: 0.7633333206176758
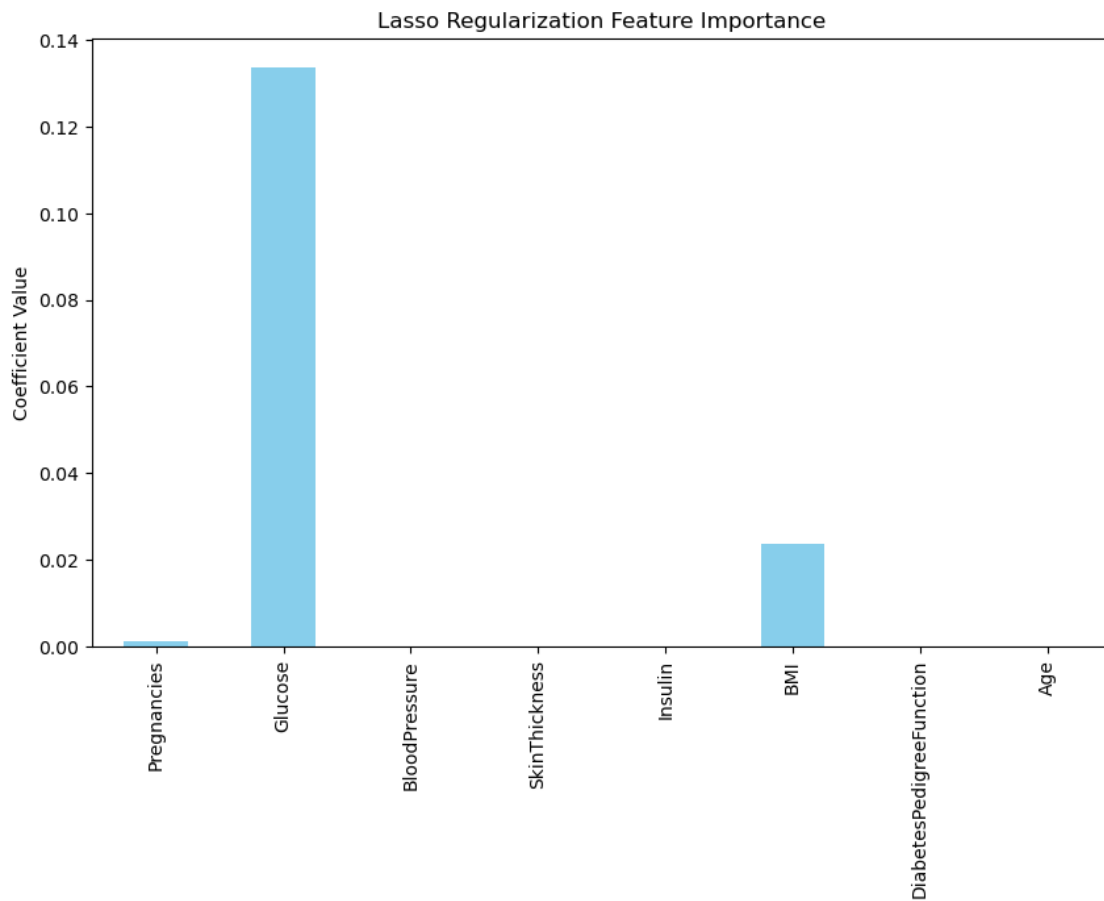
Training Voting Classifier…

Voting Classifier Metrics:
Accuracy: 0.7466666666666667
ROC-AUC: 0.8447931019156407

```
              precision    recall  f1-score   support
```

|              | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| 0            | 0.75      | 0.73   | 0.74     | 149     |
| 1            | 0.74      | 0.76   | 0.75     | 151     |
| accuracy     |           |        | 0.75     | 300     |
| macro avg    | 0.75      | 0.75   | 0.75     | 300     |
| weighted avg | 0.75      | 0.75   | 0.75     | 300     |

Applying Lasso Regularization…



Key Insights:
1. Logistic Regression and Decision Tree models highlight Glucose, BMI, and Insulin as key predictors.
2. PCA visualization shows good separability between diabetic and non-diabetic classes.
3. Ensemble models provide improved accuracy and robustness.
4. Lasso Regularization identifies the most important features for diabetes

prediction.

Analysis complete!