# Diabetes Prediction Using Machine Learning

Analyzing diagnostic measurements for early diabetes detection through advanced machine learning techniques.

Group- 5 :

1.Reddivari Abigna - 20026960

2.Pravallika palavayi – 20027660

3.Durga Prasanth guttula - 20027396

4.Ranjith Koduri - 20034703

# Project Overview

## Problem

Early diabetes identification enhances patient outcomes. Delayed diagnosis leads to more severe complications.

## Methodology

Dataset of 768 samples prepared. Models trained include Random Forest, SVM, and Logistic Regression.

## Results

Random Forest achieved 91% accuracy and 95.99% ROC-AUC. Other models showed competitive performance.

## Future Work

Further enhancement through additional models and hyperparameter tuning.

# INTRODUCTION

## Predictive Modeling

Build ML models to determine diabetes status using diagnostic measurements from patient data.

## Data Handling

Overhaul missing data and enhance features for improved prediction accuracy.

## Model Optimization

Select and improve appropriate machine learning algorithms for diabetes classification.

## Clinical Application

Create a tool for early intervention, improving health outcomes and reducing costs.

# Sample Dataset Visualization



## 768
### Total Samples
Complete patient records

## 9
### Features
Diagnostic measurements

## 2
### Classes
Diabetic (1) and non-diabetic (0)

# Data Set Diagnostic

Our dataset consists of 768 patient records with 8 diagnostic measurements and diabetes outcome classification.

| Feature | Description | Clinical Relevance |
|---|---|---|
| Pregnancies | Number of times pregnant | Gestational diabetes risk factor |
| Glucose | Plasma glucose concentration | Key diagnostic indicator |
| Blood Pressure | Diastolic (mm Hg) | Cardiovascular complication indicator |
| Skin Thickness | Triceps skinfold (mm) | Body fat distribution measure |
| Insulin | 2-hour serum insulin (mu U/ml) | Insulin resistance marker |
| BMI | Body Mass Index (kg/m²) | Obesity correlation |
| Diabetes Pedigree | Hereditary influence score | Genetic risk assessment |
| Age | Age in years | Age-related risk factor |

# Methodology

**Identify Issues**

Detect zero values in critical measurements

**Clean Data**

Replace zeros with median values

**Validate**

Ensure data quality for modeling

**Standardize**

Apply Standard Scaler to numerical features

# Model Selection and Training:

1. Develop and assess models :

- Logistic Regression

- Random Forest

- SVM, or Support Vector Machine

 2. Evaluation metrics:

- Accuracy

- Precision

- Recall

- F1-score

- ROC-AUC

# Exploratory Data Analysis

## Correlation Matrix

Identifies relationships between features. Glucose shows strongest correlation with diabetes outcome.

Age and pregnancies also show significant positive correlations with diabetes diagnosis.

## Distribution Plots

Visualize feature distributions across diabetic and non-diabetic groups.

Reveals clear separation in glucose levels, BMI, and age distributions between the two classes.

# Co-relation heatmap


Correlation Matrix

- •Heatmap Observations:

- •Outcome (target variable):

- •substantial positive association between **BMI** and **glucose**.

A moderate relationship exists between **age** and **pregnancy.**

# Feature Distribution Plots

| Feature | Distribution Pattern | Clinical Significance |
|---|---|---|
| Glucose | Strongest separation between groups | Most predictive feature for diabetes diagnosis |
| BMI | Clear rightward shift in diabetic group | Higher values strongly associated with diabetes |
| Insulin | Bimodal distribution in diabetic patients | Indicates insulin resistance patterns |
| Age | Gradual rightward shift in diabetic group | Risk increases progressively with age |
| Blood Pressure | Moderate separation between groups | Secondary indicator of comorbidity risk |

Diabetic patients consistently display elevated values across all features. Glucose levels provide the most distinct separation, followed by BMI and insulin measurements.

In [50]: ▶

```python
# density graph
fig,ax = plt.subplots(4,2, figsize=(20,20))
sns.distplot(df.Pregnancies, bins=20, ax=ax[0,0], color="red")
sns.distplot(df.Glucose, bins=20, ax=ax[0,1], color="red")
sns.distplot(df.BloodPressure, bins=20, ax=ax[1,0], color="red")
sns.distplot(df.SkinThickness, bins=20, ax=ax[1,1], color="red")
sns.distplot(df.Insulin, bins=20, ax=ax[2,0], color="red")
sns.distplot(df.BMI, bins=20, ax=ax[2,1], color="red")
sns.distplot(df.DiabetesPedigreeFunction, bins=20, ax=ax[3,0], color="red")
sns.distplot(df.Age, bins=20, ax=ax[3,1], color="red")
```
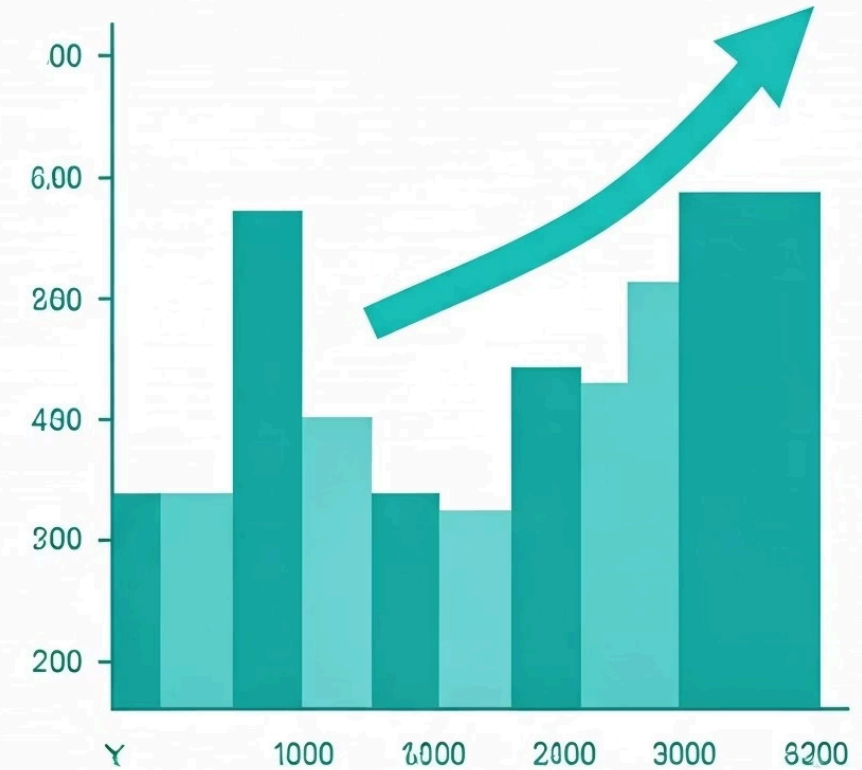
Out[50]: <Axes: xlabel='Age', ylabel='Density'>

```
In [81]:  # pair plot
          p = sns.pairplot(df, hue="Outcome")
```

# OUTCOME AND ACCURACY

# Model1:LOgistic Regression

```
In [123]:    # Machine Learning Algorithms
             # Logistic Regreesion

             log_reg = LogisticRegression()
             log_reg.fit(X_train, y_train)

Out[123]:    ▾ LogisticRegression
             LogisticRegression()
```

```
In [128]:    # Machine Learning Algorithms
             # Logistic Regression
             # Machine Learning Algorithms
             # Logistic Regression


             y_pred = log_reg.predict(X_test)
             accuracy_score(y_train, log_reg.predict(X_train))

             #Out[103]:  0.8470394736842195

             log_reg_acc = accuracy_score(y_test, log_reg.predict(X_test))
             confusion_matrix(y_test, y_pred)

             print(classification_report(y_test, y_pred))

             # Corrected ROC-AUC score calculation
             print("ROC-AUC Score:", roc_auc_score(y_test, log_reg.predict_proba(X_test)[:, 1]))
```

```
                   precision    recall  f1-score   support

               0       0.94      0.90      0.92        98
               1       0.83      0.89      0.86        54

        accuracy                           0.89       152
       macro avg       0.88      0.89      0.89       152
    weighted avg       0.90      0.89      0.90       152

ROC-AUC Score: 0.9504913076341648
```

# Model2:SVM

ROC-AUC Score: 0.9504913076341648

```
In [143]:  # SVM

svc = SVC(probability=True)
parameter = {
    "gamma": [0.0001, 0.001, 0.01, 0.1],
    'C': [0.01, 0.05, 0.5, 1, 10, 15, 20]  # Removed duplicate 0.01
}
grid_search = GridSearchCV(svc, parameter)
grid_search.fit(X_train, y_train)
grid_search.best_params_


{'C': 10, 'gamma': 0.01}


grid_search.best_score_



svc = SVC(C=10, gamma=0.01, probability=True)
svc.fit(X_train, y_train)
y_pred = svc.predict(X_test)
print("Training Accuracy:", accuracy_score(y_train, svc.predict(X_train)))
print("Test Accuracy:", accuracy_score(y_test, y_pred))
print("\nConfusion Matrix:")
print(confusion_matrix(y_test, y_pred))
print("\nClassification Report:")
print(classification_report(y_test, y_pred))
print("\nROC-AUC Score:", roc_auc_score(y_test, svc.predict_proba(X_test)[:, 1]))
```

```
Training Accuracy: 0.875
Test Accuracy: 0.9078947368421053

Confusion Matrix:
[[90  8]
 [ 6 48]]

Classification Report:
              precision    recall  f1-score   support

           0       0.94      0.92      0.93        98
           1       0.86      0.89      0.87        54

    accuracy                           0.91       152
   macro avg       0.90      0.90      0.90       152
weighted avg       0.91      0.91      0.91       152


ROC-AUC Score: 0.9516250944822373
```

# Model3:Random Forest

```
In [145]:  # random forest
           rand_clf = RandomForestClassifier(criterion='entropy', max_depth=15, max_features=0.75,
                                             min_samples_leaf=2, min_samples_split=3,
                                             n_estimators=130)
           rand_clf.fit(X_train, y_train)

           # Output from model initialization
           RandomForestClassifier
           RandomForestClassifier(criterion='entropy', max_depth=15, max_features=0.75,
           min_samples_leaf=2, min_samples_split=3,
           n_estimators=130)

           # Model evaluation
           y_pred = rand_clf.predict(X_test)
           print(accuracy_score(y_train, rand_clf.predict(X_train)))
           rand_acc = accuracy_score(y_test, rand_clf.predict(X_test))
           print(accuracy_score(y_test, rand_clf.predict(X_test)))
           print(confusion_matrix(y_test, y_pred))
           print(classification_report(y_test, y_pred))
           # Added ROC-AUC score calculation
           print("ROC-AUC Score:", roc_auc_score(y_test, rand_clf.predict_proba(X_test)[:, 1]))
```

```
0.9917763157894737
0.9078947368421053
[[89  9]
 [ 5 49]]
              precision    recall  f1-score   support

           0       0.95      0.91      0.93        98
           1       0.84      0.91      0.88        54

    accuracy                           0.91       152
   macro avg       0.90      0.91      0.90       152
weighted avg       0.91      0.91      0.91       152

ROC-AUC Score: 0.9599395313681028
```
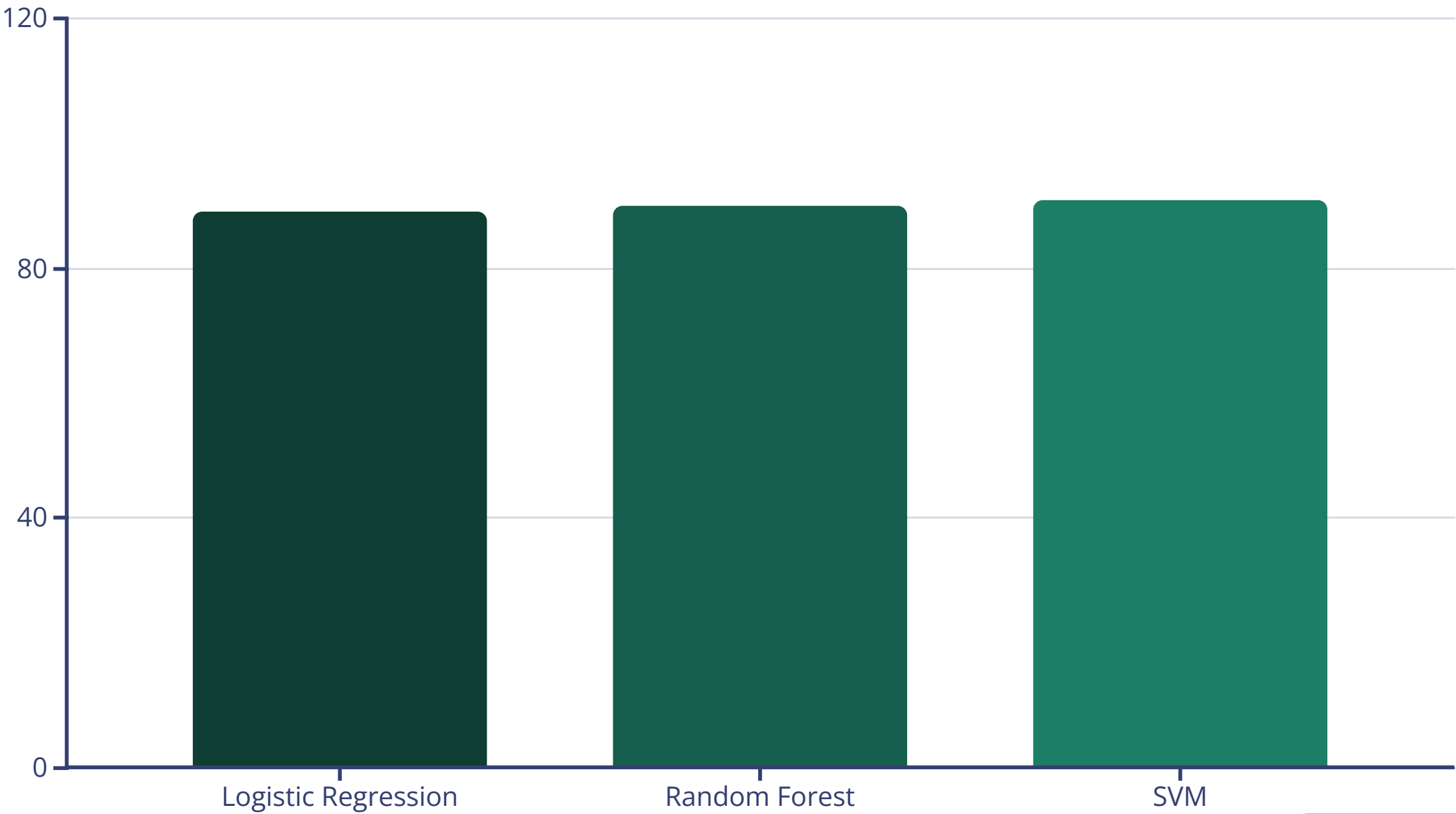
# Analysis

**Logistic Regression:**

Accuracy: 89%

**Random Forest:**

Accuracy: 91%

**SVM:**

Accuracy: 91%

# ROC-AUC

**Logistic Regression:95.04%**

**Random Forest:95.99%**

**SVM:95.16%**

•With the greatest ROC-AUC score of 95.99% and a balanced classification report, the Random Forest model performs the best overall.

# Conclusion

The model's strong ROC-AUC score of 95.99% suggests that it is a good fit for detecting diabetes patients.

With a balanced categorization report, the Random Forest model performed the best in predicting diabetes.

Clinical Relevance:

1. The approach aids in the early intervention identification of those who are at-risk.

2. Better diabetes diagnosis and treatment result in better health and lower medical expenses

Future Improvements:

1. To improve accuracy, adjust the hyperparameters.

2. Examine further ensemble techniques and machine learning models.

# THANK YOU