# Project Report: Predicting Patient Diagnosis Using Random Forest Classifier

**Akhil Sharma (102203028)**
**Amol Singh (102203354)**

# Problem Statement

The healthcare industry generates vast amounts of data, which, if analyzed effectively, can uncover patterns to aid in early diagnosis of diseases. This project addresses the problem of predicting patient diagnoses—categorized into **Control**, **Benign**, and **PDAC (Pancreatic Ductal Adenocarcinoma)**—using clinical and demographic data. Early detection is crucial for diseases like PDAC, as it significantly improves prognosis and treatment outcomes.

The dataset for this project comprises patient health records, including both numerical and categorical features. Challenges include handling missing data, converting categorical data into a machine-readable format, and building a robust predictive model.

# Proposed Solution

To tackle this problem, we employ a **Random Forest Classifier**, an ensemble machine learning method. Random forests are particularly well-suited for healthcare datasets because they:

- Handle mixed data types (numerical and categorical).
- Are resistant to overfitting, by using ensembles of decision trees.
- Provide feature importance metrics to identify the most influential factors in predictions.

Our approach includes:

1. **Data Preprocessing**: Handling missing values, encoding categorical variables, and splitting the data into training and testing subsets.
2. **Model Training**: Using a Random Forest Classifier to train on the processed data.
3. **Evaluation**: Assessing the model's performance using accuracy, precision, recall, F1-score, and a confusion matrix.
4. **Insights**: Visualizing feature importance to understand the impact of different variables on predictions.

# Methodology and Approach

## 1. Data Preprocessing

The dataset required significant preprocessing:

- Missing values in the `diagnosis` column were filled with `1` (assuming the patient is healthy), while `stage` was replaced with "Unknown."
- Numerical columns had missing values replaced with their mean.
- Categorical features (e.g., `Diabetes Status`, `Smoking History`) were encoded using one-hot encoding to create machine-readable binary variables.
- Irrelevant columns, such as `sample_id`, `patient_cohort`, and `sample_origin`, were excluded from the feature set.

## 2. Feature and Target Selection

The `diagnosis` column was selected as the target variable (`y`), with all other relevant columns forming the feature matrix (`X`). Boolean columns were converted to integer type for compatibility with the Random Forest model.

## 3. Model Implementation

A Random Forest Classifier was trained on an 80-20 train-test split. The `stratify` parameter ensured class distribution consistency between training and testing sets.

## 4. Model Evaluation

```
Accuracy: 0.81
              precision    recall  f1-score   support

           1       0.76      0.59      0.67        37
           2       0.69      0.83      0.76        41
           3       1.00      1.00      1.00        40

    accuracy                           0.81       118
   macro avg       0.82      0.81      0.81       118
weighted avg       0.82      0.81      0.81       118
```
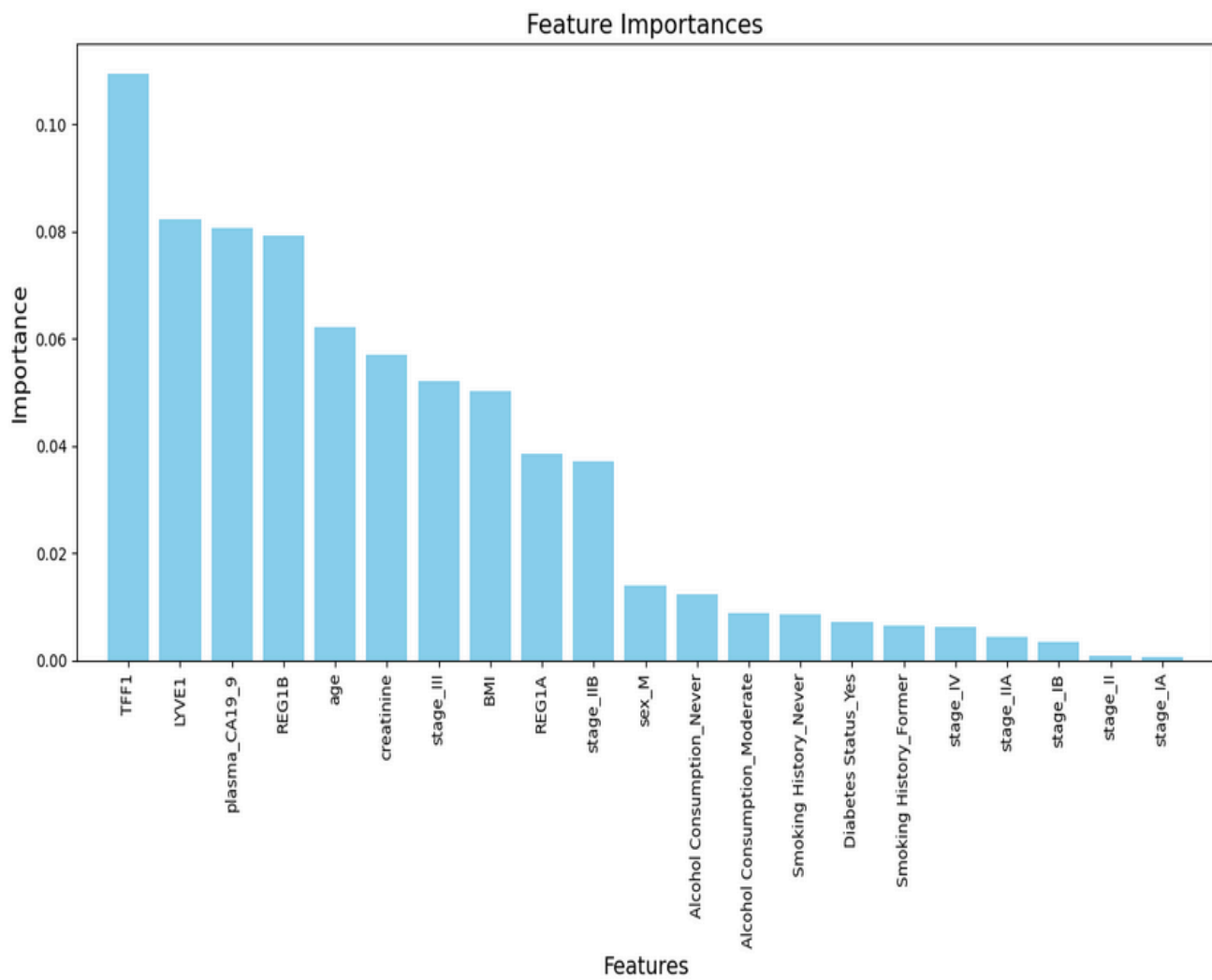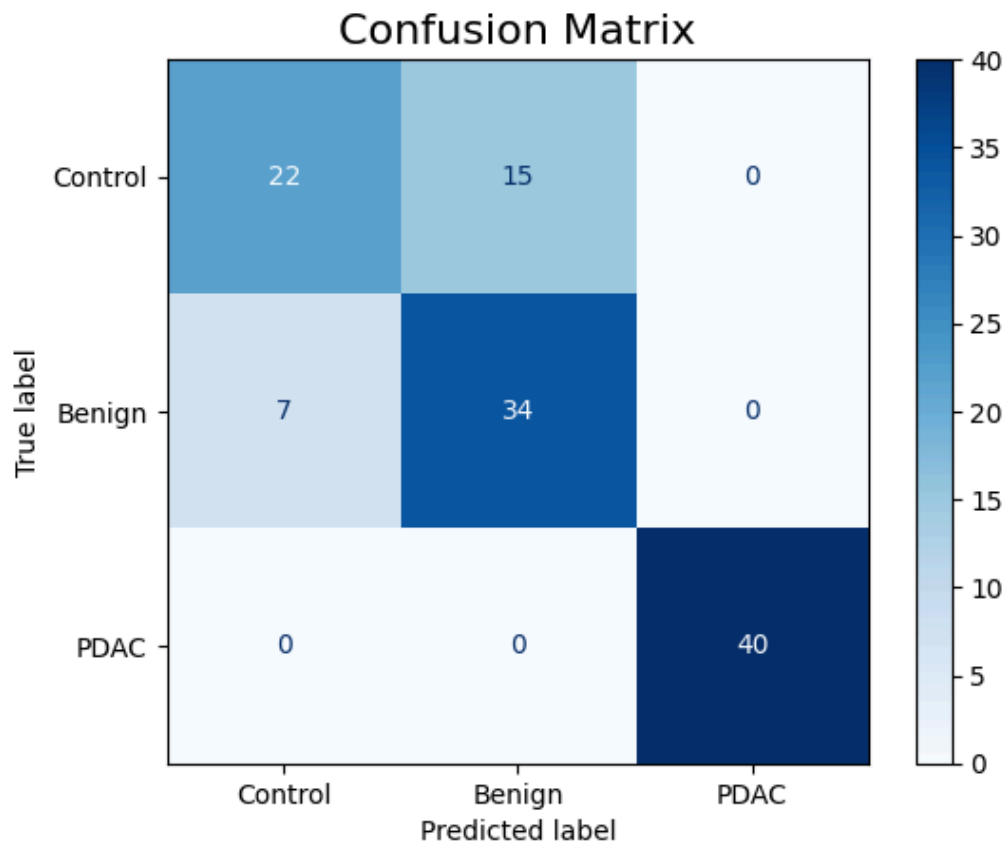
# Results and Observations

## Model Performance

The Random Forest Classifier achieved an accuracy of **83%** on the test data. Detailed metrics from the classification report highlight the model's effectiveness across all classes.

## Key Insights

1. The top features influencing the diagnosis were:
   - **TFF1**
   - **LYVE1**


Feature Importances

2. The confusion matrix revealed that **PDAC cases were correctly classified in 100% of instances**.

## Confusion Matrix



## Conclusion

This project demonstrates the effectiveness of machine learning, specifically Random Forests, in predicting patient diagnoses using clinical data. By leveraging the feature importance insights, healthcare providers can focus on the most influential factors for early disease detection.

# Code Implementation and Results

```python
import pandas as pd
import numpy as np
from sklearn.model_selection import train_test_split
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import accuracy_score, classification_report,
confusion_matrix, ConfusionMatrixDisplay
import matplotlib.pyplot as plt

# Load the modified dataset
file_path = 'Debernardi_modified_data.csv'
data = pd.read_csv(file_path)

data['diagnosis'] = data['diagnosis'].fillna(1).astype(int)  # Assume
patient is healthy

data['stage'] = data['stage'].fillna('Unknown')  # Fill NaNs with a
placeholder or consider dropping

numeric_cols = data.select_dtypes(include='number').columns # picking
the numeric data colums
data[numeric_cols] =
data[numeric_cols].fillna(data[numeric_cols].mean())  # Fill NaNs with
mean

# Convert categorical features to numerical using one-hot encoding
categorical_columns = ['Diabetes Status', 'Smoking History', 'Alcohol
Consumption', 'sex', 'stage']
data = pd.get_dummies(data, columns=categorical_columns,
drop_first=True)

# Define features and target variable
X = data.drop(columns=['sample_id', 'patient_cohort', 'sample_origin',
'diagnosis', 'benign_sample_diagnosis'])
y = data['diagnosis']

# Convert boolean columns to int (True -> 1, False -> 0)
bool_cols = X.select_dtypes(include='bool').columns
X[bool_cols] = X[bool_cols].astype(int)

# Split the data into training and testing sets
```

```python
X_train, X_test, y_train, y_test = train_test_split(X, y,
test_size=0.2, random_state=42, stratify=y)

# Train a Random Forest classifier
rf_model = RandomForestClassifier(random_state=42)

# Fit the model
rf_model.fit(X_train, y_train)

# Make predictions
y_pred = rf_model.predict(X_test)

# Calculate accuracy
accuracy = accuracy_score(y_test, y_pred)
print(f'Accuracy: {accuracy:.2f}')

# Print classification report for more detailed metrics
print(classification_report(y_test, y_pred))

# Get feature importances
importances = rf_model.feature_importances_
indices = np.argsort(importances)[::-1]

# Plot the feature importances
plt.figure(figsize=(12, 8))
plt.title('Feature Importances', fontsize=16)
plt.bar(range(X.shape[1]), importances[indices], align='center',
color='skyblue')
plt.xticks(range(X.shape[1]), X.columns[indices], rotation=90,
fontsize=10)
plt.xlim([-1, X.shape[1]])
plt.ylabel('Importance', fontsize=14)
plt.xlabel('Features', fontsize=14)
plt.tight_layout()
plt.show()

# Plot confusion matrix
cm = confusion_matrix(y_test, y_pred)
disp = ConfusionMatrixDisplay(confusion_matrix=cm,
display_labels=['Control', 'Benign', 'PDAC'])
disp.plot(cmap=plt.cm.Blues)
plt.title('Confusion Matrix', fontsize=16)
plt.show()
```