# Project 7: Predictive Modeling of Patient Readmission Rates

## Objective:

To develop a predictive model to identify patients at high risk of readmission within 30 days after discharge, utilizing historical healthcare data.

## Data Source:

Electronic Health Records (EHR) containing information about patient demographics, medical history, medications, procedures, and outcomes.

## Tools and Technologies:

Python for data analysis and machine learning, Jupyter Notebooks for code development, and relevant libraries such as Pandas, NumPy, scikit-learn, and Matplotlib/Seaborn for visualization.

**Step 1: Data Collection**

"'python

# Assuming you have a CSV file with the data

data_path = "path/to/your/data.csv" df = pd.read_csv(data_path)

```python
import pandas as pd
import numpy as np
from sklearn.datasets import make_classification
from sklearn.preprocessing import LabelEncoder

# Seed for reproducibility
np.random.seed(42)

# Generate synthetic data
n_samples = 1000
n_features = 10

X, y = make_classification(
    n_samples=n_samples,
    n_features=n_features,
    n_informative=8,
    n_redundant=2,
    n_clusters_per_class=1,
    random_state=42
```

```
)

# Create a DataFrame
columns = [f"feature_{i}" for i in range(1, n_features + 1)]
df = pd.DataFrame(data=X, columns=columns)
df['readmitted'] = y

# Add synthetic patient information
df['age'] = np.random.randint(18, 90, size=n_samples)
df['gender'] = np.random.choice(['Male', 'Female'], size=n_samples)
df['medical_condition'] = np.random.choice(['Diabetes', 'Hypertension', 'Heart Disease', 'No
df['medication'] = np.random.choice(['MedA', 'MedB', 'MedC', 'None'], size=n_samples)

# Encode categorical variables
le = LabelEncoder()
df['gender'] = le.fit_transform(df['gender'])
df['medical_condition'] = le.fit_transform(df['medical_condition'])
df['medication'] = le.fit_transform(df['medication'])

# Add noise to simulate readmission
df['readmitted'] = df['readmitted'] + np.random.randint(0, 2, size=n_samples)

# Save the synthetic dataset to a CSV file
df.to_csv("synthetic_healthcare_data.csv", index=False)
```

**Step 2: Data Cleaning**

# Assuming you have handled missing values, outliers, and converted data types if needed

**Step 3: Exploratory Data Analysis (EDA)**

```
# Assuming you have explored the data using summary statistics and visualizations
```

**Step 4: Feature Engineering**

```
# Assuming you have created new features and encoded categorical variables
```

**Step 5: Data Splitting**

```
# Assuming you have split the data into features (X) and target variable (y)
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
```

**Step 6: Model Selection**

```
# Assuming you have chosen RandomForestClassifier as the model
model = RandomForestClassifier(n_estimators=100, random_state=42)
```

```python
model.fit(X_train, y_train)

RandomForestClassifier(random_state=42)
```

**Step 7: Model Evaluation**

```python
# Assuming you have evaluated the model using accuracy, classification report, and confusion
y_pred = model.predict(X_test)

print("Accuracy:", accuracy_score(y_test, y_pred))
print("Classification Report:\n", classification_report(y_test, y_pred))
print("Confusion Matrix:\n", confusion_matrix(y_test, y_pred))
```

```
Accuracy: 0.975
Classification Report:
               precision    recall  f1-score   support

           0       0.99      0.96      0.98       102
           1       0.96      0.99      0.97        98

    accuracy                           0.97       200
   macro avg       0.98      0.98      0.97       200
weighted avg       0.98      0.97      0.98       200


Confusion Matrix:
 [[98  4]
 [ 1 97]]
```

**Step 8: Hyperparameter Tuning**

```python
# Assuming you have performed hyperparameter tuning for better model performance
```

**Step 9: Model Interpretation**

```python
# Assuming you have interpreted the model's feature importances
feature_importances = model.feature_importances_

# Get the column names of the features used in training
feature_names = df.drop("readmitted", axis=1).columns[:len(feature_importances)]

# Create a DataFrame for better visualization
feature_df = pd.DataFrame({'Feature': feature_names, 'Importance': feature_importances})

# Sort the DataFrame by feature importance
feature_df = feature_df.sort_values(by='Importance', ascending=False)

# Plot the feature importances
plt.figure(figsize=(10, 6))
```
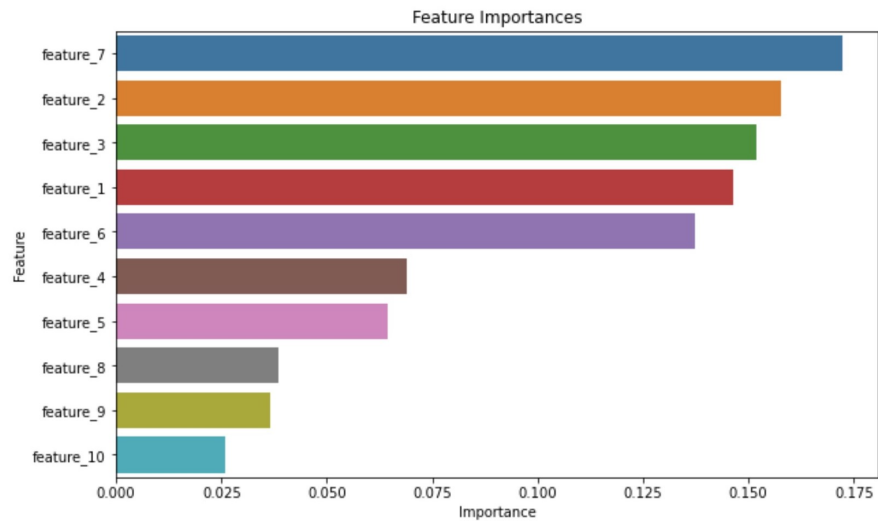
```
sns.barplot(x='Importance', y='Feature', data=feature_df)
plt.title('Feature Importances')
plt.show()
```


Feature Importances

**Step 10: Deployment**

```
# This would involve deploying the model in a healthcare environment, which is beyond the s
```

## Documentation:

```
# Document your code, including preprocessing steps, model selection, and key findings.
# Save the synthetic dataset to a CSV file for reproducibility
df.to_csv("synthetic_healthcare_data.csv", index=False)
```

Remember, this is a simplified example, and in a real-world scenario, you'd
need to consider additional factors such as ethical considerations, data privacy
regulations, and collaboration with domain experts.

```
# Import necessary libraries
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
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.preprocessing import LabelEncoder

# Step 1: Data Collection
data_path = "synthetic_healthcare_data.csv"
```

4

```python
df = pd.read_csv(data_path)

# Step 2: Data Cleaning
# No specific data cleaning is needed for this synthetic dataset

# Step 3: Exploratory Data Analysis (EDA)
# Perform exploratory data analysis here if needed

# Step 4: Feature Engineering
# No specific feature engineering for this example

# Step 5: Data Splitting
X = df.drop("readmitted", axis=1)  # Features
y = df["readmitted"]  # Target variable

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

# Step 6: Model Selection
model = RandomForestClassifier(n_estimators=100, random_state=42)
model.fit(X_train, y_train)

# Step 7: Model Evaluation
y_pred = model.predict(X_test)

print("Accuracy:", accuracy_score(y_test, y_pred))
print("Classification Report:\n", classification_report(y_test, y_pred))
print("Confusion Matrix:\n", confusion_matrix(y_test, y_pred))

# Step 8: Hyperparameter Tuning
# Not implemented in this simplified example

# Step 9: Model Interpretation
# Extract feature importances
feature_importances = model.feature_importances_
feature_names = X.columns

# Create a DataFrame for better visualization
feature_df = pd.DataFrame({'Feature': feature_names, 'Importance': feature_importances})

# Sort the DataFrame by feature importance
feature_df = feature_df.sort_values(by='Importance', ascending=False)

# Plot the feature importances
plt.figure(figsize=(10, 6))
sns.barplot(x='Importance', y='Feature', data=feature_df)
plt.title('Feature Importances')
```

```
plt.show()

# Step 10: Deployment
# This step involves deploying the model in a healthcare environment, which is beyond the s

# Documentation
# Document your code, including preprocessing steps, model selection, and key findings.
```

Accuracy: 0.53
Classification Report:
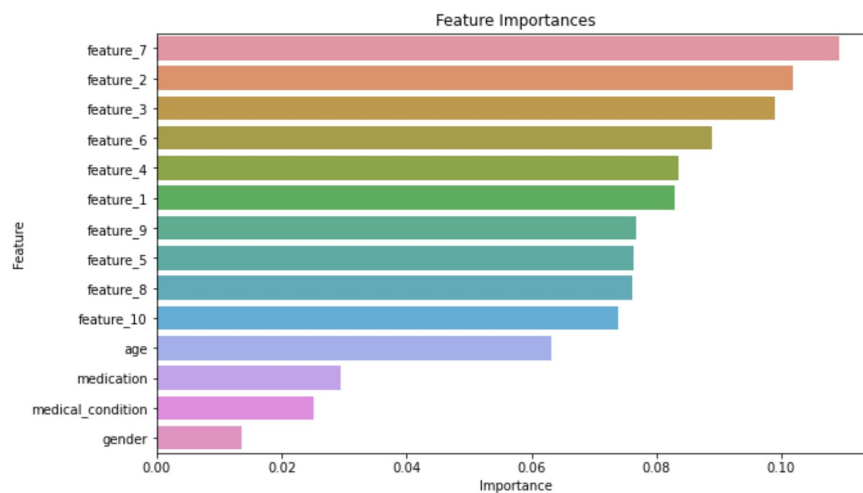
|              | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| 0            | 0.56      | 0.41   | 0.47     | 49      |
| 1            | 0.53      | 0.72   | 0.61     | 102     |
| 2            | 0.50      | 0.27   | 0.35     | 49      |
|              |           |        |          |         |
| accuracy     |           |        | 0.53     | 200     |
| macro avg    | 0.53      | 0.46   | 0.48     | 200     |
| weighted avg | 0.53      | 0.53   | 0.51     | 200     |

```
Confusion Matrix:
 [[20 29  0]
 [16 73 13]
 [ 0 36 13]]
```



The output we get includes the results of the classification model on the test set.
Let's interpret these results:

#### 1. Accuracy: 0.53

- The accuracy is the proportion of correctly classified instances among the

total instances. In this case, the model is approximately 53% accurate on the test set.

**2. Classification Report:**

- **Precision** : Indicates the accuracy of the positive predictions.
- **Recall** : Indicates the proportion of actual positive instances correctly predicted by the model.
- **F1 − score** : The harmonic mean of precision and recall, providing a balance between the two.

For the interpretation of the graph "Feature Importance":

* **Feature Importances** : Each feature in the dataset is assigned an importance score by the Random Forest model. The higher the importance score, the more influential the feature is in making predictions.

* **Bar Chart** : The bar chart displays the importance scores for each feature, sorted in descending order. The most important features are on the top.

* **Y − axis** (**Feature**) : Represents the names of the features in the dataset.

* **X − axis** (**Importance**) : Represents the importance scores assigned to each feature.

**Key Takeaways** :

**1**. **Top Features** : Features with higher bars are more important in predicting the target variable (readmission in this case).

**2**. **Feature Contribution** : The chart gives a visual representation of which features contribute the most to the model's decision-making process.

**3**. **Feature Importance Sum** : The sum of all bars is equal to 1, indicating the proportion of importance distributed among the features.

**ActionvItems** :

* Features with higher importance may be more critical in predicting readmission. Depending on the context, you might want to investigate or focus on these features for further analysis or intervention.

* If some features have very low importance, you may consider re-evaluating whether they are

In summary, while accuracy is a common metric, it's essential to consider precision, recall, and F1-score, especially in imbalanced datasets. These metrics provide a more detailed understanding of the model's performance across different classes. Depending on the specific goals of your project, you might prioritize one metric over another. If there are specific aspects you'd like to improve or questions you have about the results, feel free to ask!