**Exploratory Data Analysis (EDA)**

**Univariate Analysis**

```python
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import numpy as np

# Load the dataset (replace with your file path)
data = pd.read_csv('/content/heart disease prediction.csv')

# Display basic information about the dataset
print(data.head())
print(data.info())
print(data.describe())

# Univariate Analysis
# Distribution of the target variable
plt.figure(figsize=(6, 4))
sns.countplot(x='HeartDisease', data=data)
plt.title('Distribution of Heart Disease')
plt.show()

# Distribution of numerical features
data.hist(bins=30, figsize=(15, 10))
plt.show()
```

```
        Age Sex ChestPainType  RestingBP  Cholesterol  FastingBS RestingECG  MaxHR  \
0        40   M           ATA        140          289          0     Normal    172
1        49   F           NAP        160          180          0     Normal    156
2        37   M           ATA        130          283          0         ST     98
3        48   F           ASY        138          214          0     Normal    108
4        54   M           NAP        150          195          0     Normal    122

   ExerciseAngina  Oldpeak ST_Slope  HeartDisease
0               N      0.0       Up             0
1               N      1.0     Flat             1
2               N      0.0       Up             0
3               Y      1.5     Flat             1
4               N      0.0       Up             0
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 918 entries, 0 to 917
Data columns (total 12 columns):
 #   Column         Non-Null Count  Dtype
---  ------         --------------  -----
 0   Age            918 non-null    int64
 1   Sex            918 non-null    object
 2   ChestPainType  918 non-null    object
 3   RestingBP      918 non-null    int64
 4   Cholesterol    918 non-null    int64
 5   FastingBS      918 non-null    int64
 6   RestingECG     918 non-null    object
 7   MaxHR          918 non-null    int64
 8   ExerciseAngina 918 non-null    object
 9   Oldpeak        918 non-null    float64
 10  ST_Slope       918 non-null    object
 11  HeartDisease   918 non-null    int64
dtypes: float64(1), int64(6), object(5)
memory usage: 86.2+ KB
None
              Age    RestingBP  Cholesterol    FastingBS        MaxHR  \
count  918.000000   918.000000   918.000000   918.000000   918.000000
mean    53.510893   132.396514   198.799564     0.233115   136.809368
std      9.432617    18.514154   109.384145     0.423046    25.460334
min     28.000000     0.000000     0.000000     0.000000    60.000000
25%     47.000000   120.000000   173.250000     0.000000   120.000000
50%     54.000000   130.000000   223.000000     0.000000   138.000000
75%     60.000000   140.000000   267.000000     0.000000   156.000000
max     77.000000   200.000000   603.000000     1.000000   202.000000

          Oldpeak  HeartDisease
count  918.000000    918.000000
mean     0.887364      0.553377
std      1.066570      0.497414
min     -2.600000      0.000000
25%      0.000000      0.000000
50%      0.600000      1.000000
75%      1.500000      1.000000
max      6.200000      1.000000
```
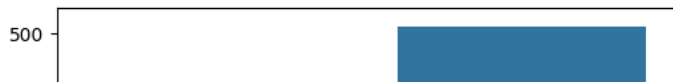
### Distribution of Heart Disease

**Bivariate Analysis**

```
# Correlation heatmap
plt.figure(figsize=(10, 8))
sns.heatmap(data.select_dtypes(include=np.number).corr(), annot=True, cmap='coolwarm', vmin=-1, vmax=1) # Select only numerical columns f
plt.title('Correlation Heatmap')
plt.show()


# Boxplot for numerical features by target
plt.figure(figsize=(12, 6))
sns.boxplot(x='HeartDisease', y='Age', data=data)
plt.title('Age vs Heart Disease')
plt.show()
```
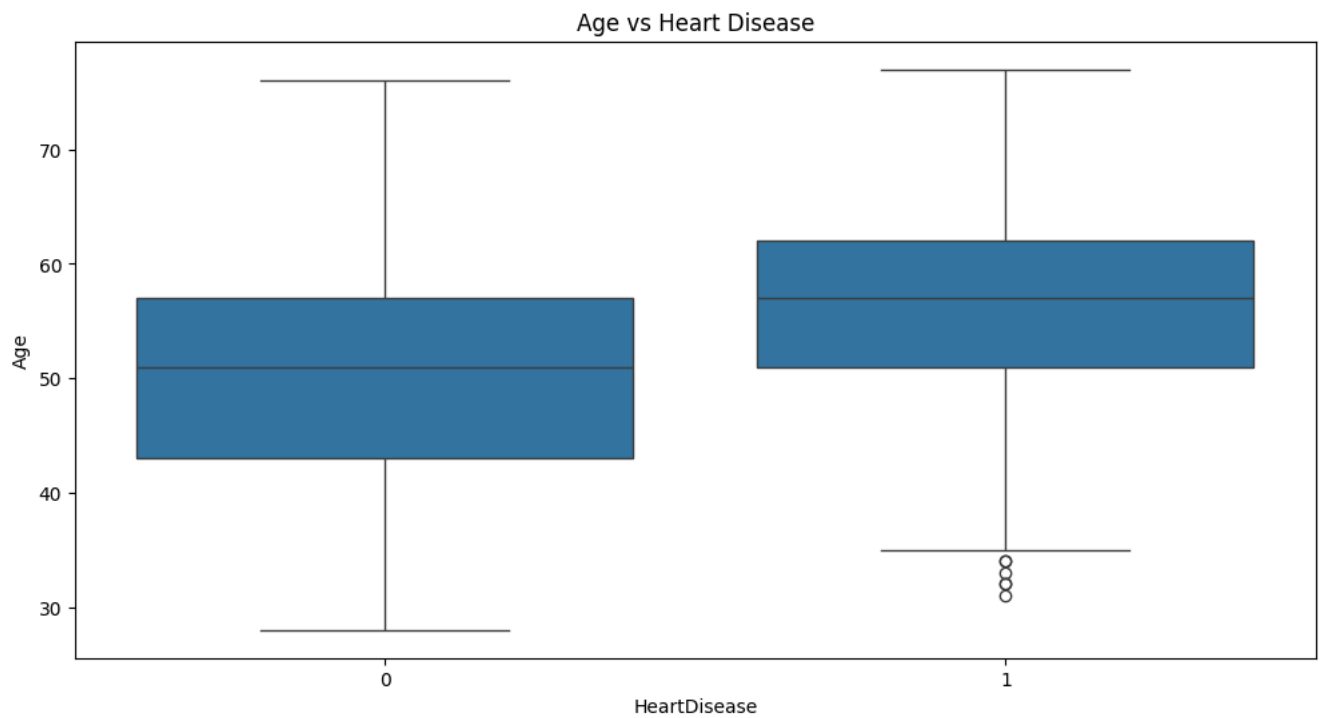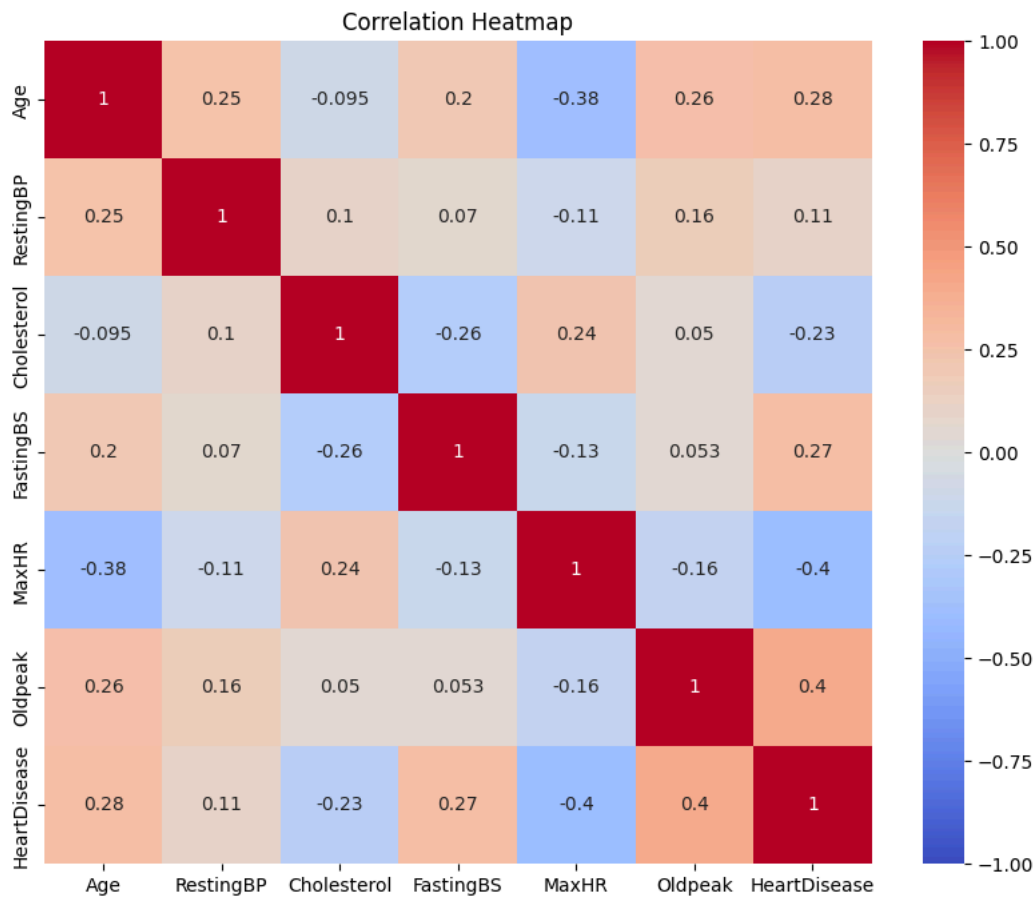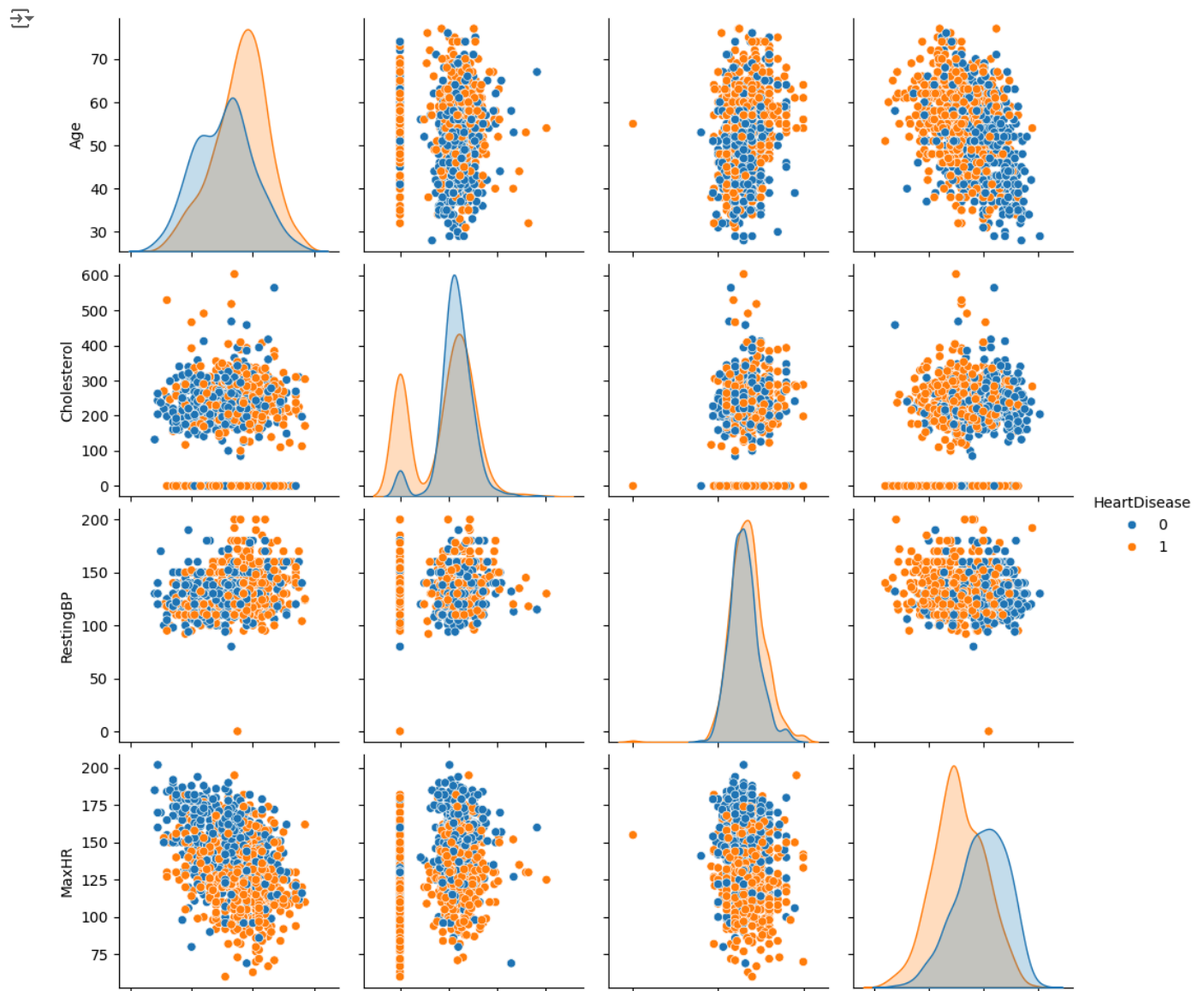
## Correlation Heatmap



## Age vs Heart Disease



**Multivariate Analysis**

```python
# Pairplot for selected features
sns.pairplot(data, hue='HeartDisease', vars=['Age', 'Cholesterol', 'RestingBP', 'MaxHR'])
plt.show()
```

## Data Preprocessing

```python
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler, OneHotEncoder
from sklearn.compose import ColumnTransformer
from sklearn.pipeline import Pipeline
from sklearn.impute import SimpleImputer
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import classification_report, accuracy_score
import joblib

# Preprocessing
# Handle missing values and encode categorical variables
features = data.drop(columns=['HeartDisease'])
target = data['HeartDisease']

# Identify categorical and numerical columns
categorical_features = features.select_dtypes(include=['object']).columns
numerical_features = features.select_dtypes(exclude=['object']).columns

# Define preprocessing for numerical and categorical data
numerical_transformer = Pipeline(steps=[
    ('imputer', SimpleImputer(strategy='median')),
    ('scaler', StandardScaler())])

categorical_transformer = Pipeline(steps=[
    ('imputer', SimpleImputer(strategy='most_frequent')),
    ('onehot', OneHotEncoder(handle_unknown='ignore'))])

preprocessor = ColumnTransformer(
    transformers=[
        ('num', numerical_transformer, numerical_features),
        ('cat', categorical_transformer, categorical_features)])

# Split data
```

```
    X_train, X_test, y_train, y_test = train_test_split(features, target, test_size=0.2, random_state=42)

    # Define model
    model = Pipeline(steps=[('preprocessor', preprocessor),
                            ('classifier', RandomForestClassifier(random_state=42))])

    # Train model
    model.fit(X_train, y_train)

    # Predict and evaluate
    y_pred = model.predict(X_test)
    print("Train Accuracy:", accuracy_score(y_train, model.predict(X_train)))
    print("Test Accuracy:", accuracy_score(y_test, y_pred))
    print(classification_report(y_test, y_pred))
```

```
Train Accuracy: 1.0
Test Accuracy: 0.8804347826086957
              precision    recall  f1-score   support

           0       0.85      0.87      0.86        77
           1       0.90      0.89      0.90       107
```