

```

# Import necessary libraries for data analysis and visualization
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import urllib.request

# Download dataset
url = "https://sp8138-heart-attack-dataset.s3.us-east-2.amazonaws.com/heart_attack_data_1000.csv"
local_file = "heart_attack_data.csv"
urllib.request.urlretrieve(url, local_file)

('heart_attack_data.csv', <http.client.HTTPMessage at 0x7f5128540150>)

# Load the dataset into a pandas DataFrame
df = pd.read_csv(local_file)

# Display basic information about the dataset
print("==== DATASET INFO ====")
print(df.info())
print("\n==== FIRST 5 ROWS ====")
display(df.head())

==== DATASET INFO ====
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 999 entries, 0 to 998
Data columns (total 32 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   Age                                   999 non-null    int64
1   Gender                               999 non-null    object
2   Cholesterol                           999 non-null    int64
3   BloodPressure                         999 non-null    int64
4   HeartRate                             999 non-null    int64
5   BMI                                   999 non-null    float64
6   Smoker                               999 non-null    int64
7   Diabetes                             999 non-null    int64
8   Hypertension                         999 non-null    int64
9   FamilyHistory                       999 non-null    int64
10  PhysicalActivity                     999 non-null    int64
11  AlcoholConsumption                  999 non-null    int64
12  Diet                                 999 non-null    object
13  StressLevel                         999 non-null    int64
14  Ethnicity                           999 non-null    object
15  Income                               999 non-null    int64
16  EducationLevel                      999 non-null    object
17  Medication                          999 non-null    object
18  ChestPainType                       999 non-null    object
19  ECGResults                          999 non-null    object

```

```

20  MaxHeartRate          999 non-null    int64
21  ST_Depression         999 non-null    float64
22  ExerciseInducedAngina 999 non-null    object
23  Slope                  999 non-null    object
24  NumberOfMajorVessels  999 non-null    int64
25  Thalassemia           999 non-null    object
26  PreviousHeartAttack   999 non-null    int64
27  StrokeHistory         999 non-null    int64
28  Residence             999 non-null    object
29  EmploymentStatus      999 non-null    object
30  MaritalStatus         999 non-null    object
31  Outcome               999 non-null    object
dtypes: float64(2), int64(16), object(14)
memory usage: 249.9+ KB
None

```

```
===== FIRST 5 ROWS =====
```

```

{"type": "dataframe"}

# Check for missing values in the dataset
print("\n===== MISSING VALUES =====")
print(df.isnull().sum())

```

```
===== MISSING VALUES =====
```

```

Age          0
Gender       0
Cholesterol  0
BloodPressure 0
HeartRate    0
BMI          0
Smoker       0
Diabetes     0
Hypertension 0
FamilyHistory 0
PhysicalActivity 0
AlcoholConsumption 0
Diet        0
StressLevel  0
Ethnicity    0
Income       0
EducationLevel 0
Medication   0
ChestPainType 0
ECGResults   0
MaxHeartRate 0
ST_Depression 0
ExerciseInducedAngina 0
Slope        0

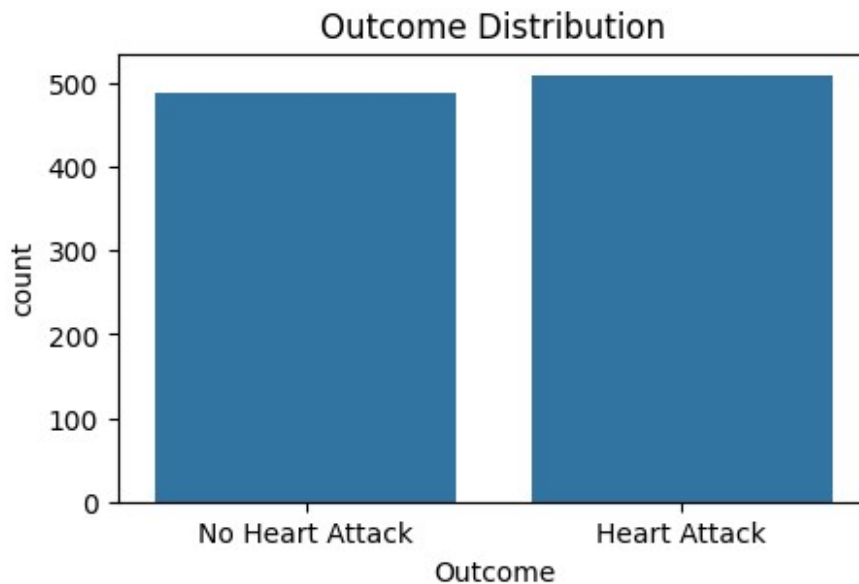
```

```

NumberOfMajorVessels    0
Thalassemia             0
PreviousHeartAttack      0
StrokeHistory           0
Residence               0
EmploymentStatus        0
MaritalStatus           0
Outcome                 0
dtype: int64

# Visualize the distribution of the outcome variable
plt.figure(figsize=(5,3))
sns.countplot(data=df, x="Outcome")
plt.title("Outcome Distribution")
plt.show()
print(df['Outcome'].value_counts())

```



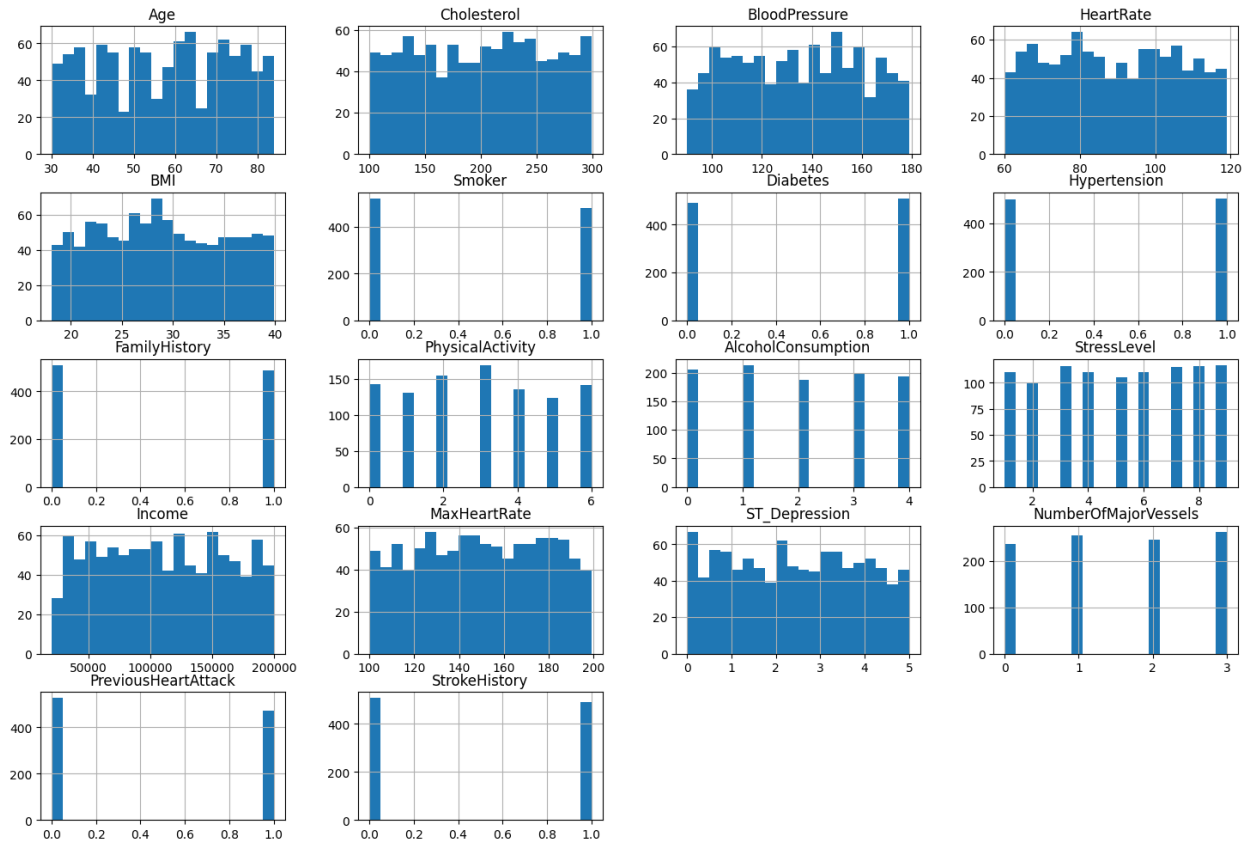
```

Outcome
Heart Attack    510
No Heart Attack 489
Name: count, dtype: int64

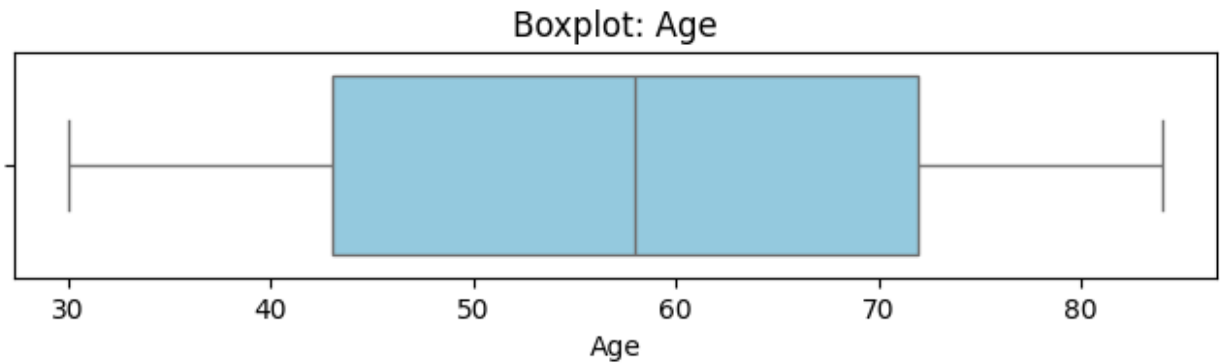
# Plot histograms for all numeric features
num_cols = df.select_dtypes(include=[np.number]).columns
df[num_cols].hist(figsize=(18, 12), bins=20)
plt.suptitle("Numeric Feature Distributions")
plt.show()

```

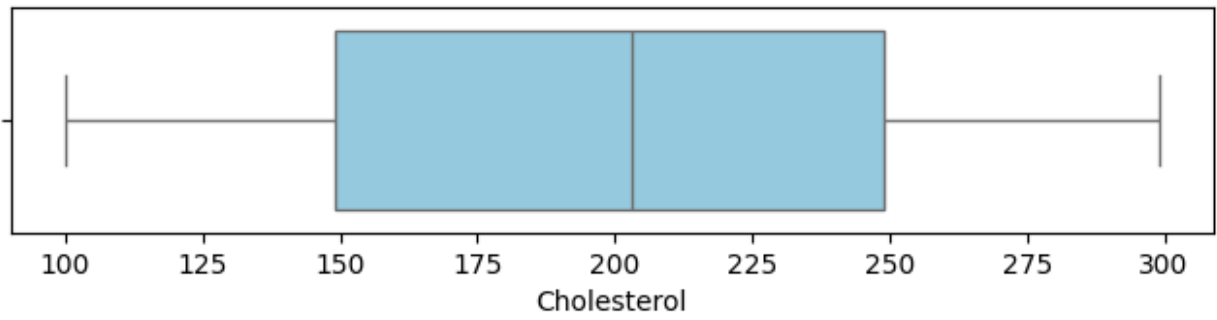
Numeric Feature Distributions



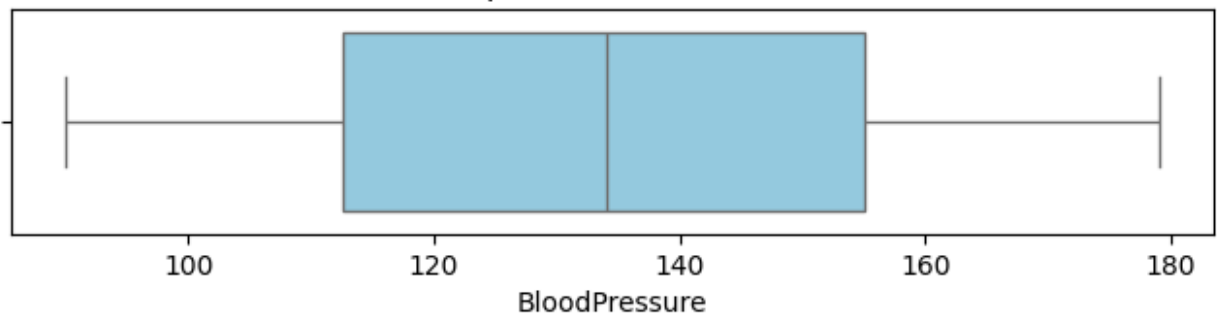
```
# Create boxplots for each numeric feature to identify outliers
for col in num_cols:
    plt.figure(figsize=(8,1.5))
    sns.boxplot(x=df[col], color='skyblue')
    plt.title(f"Boxplot: {col}")
    plt.show()
```



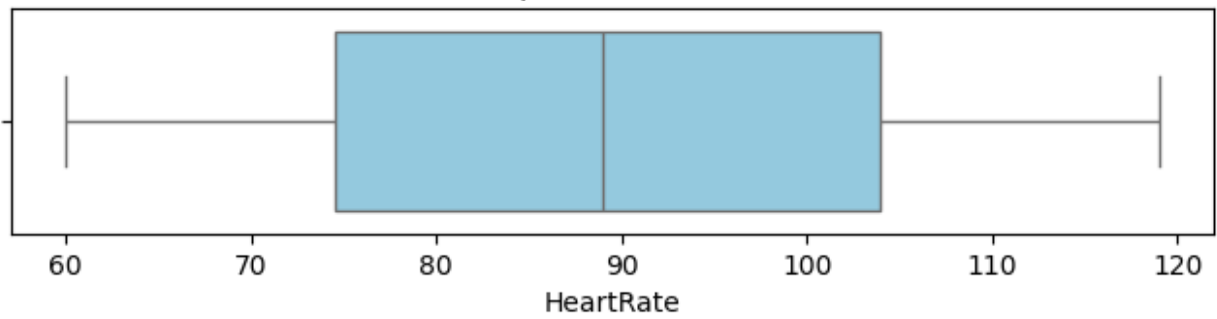
Boxplot: Cholesterol



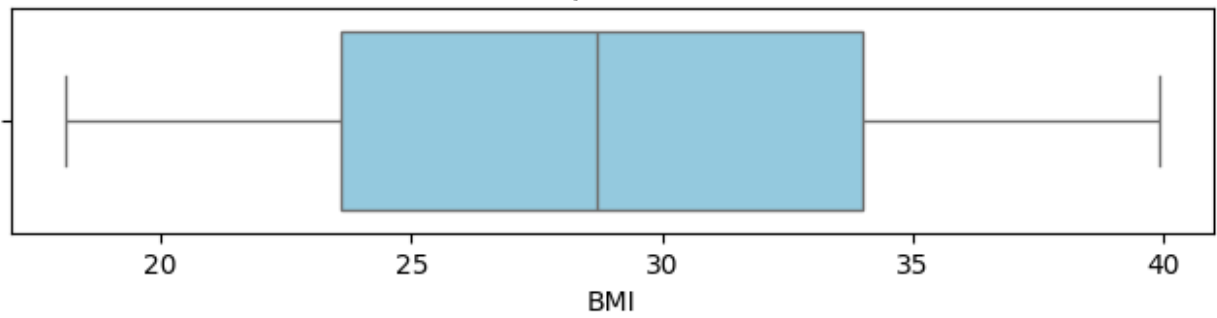
Boxplot: BloodPressure



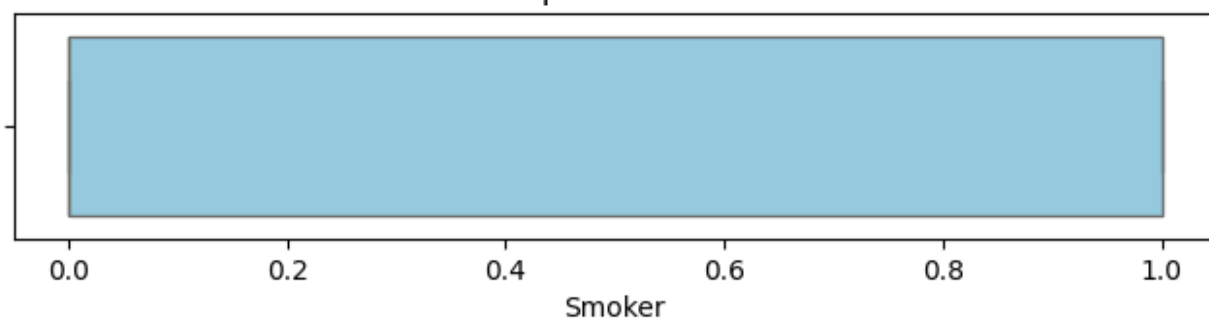
Boxplot: HeartRate



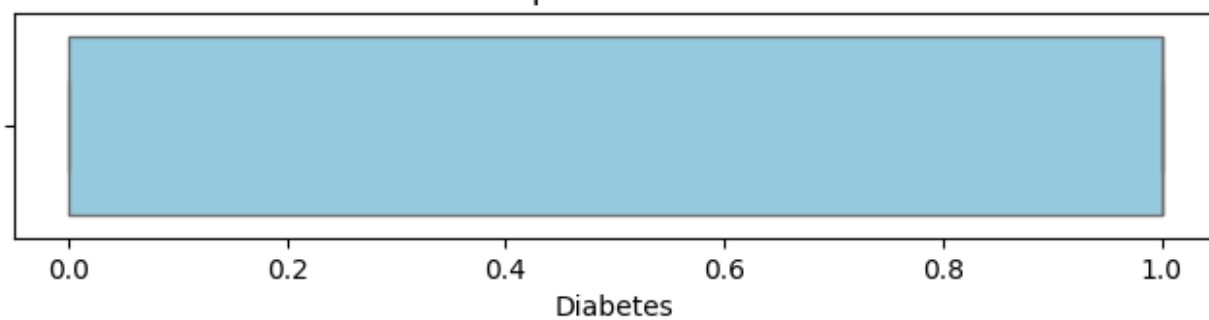
Boxplot: BMI



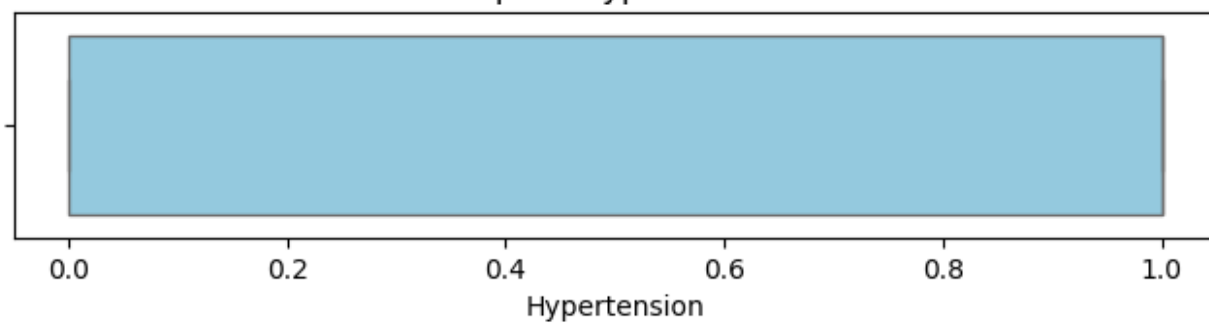
Boxplot: Smoker



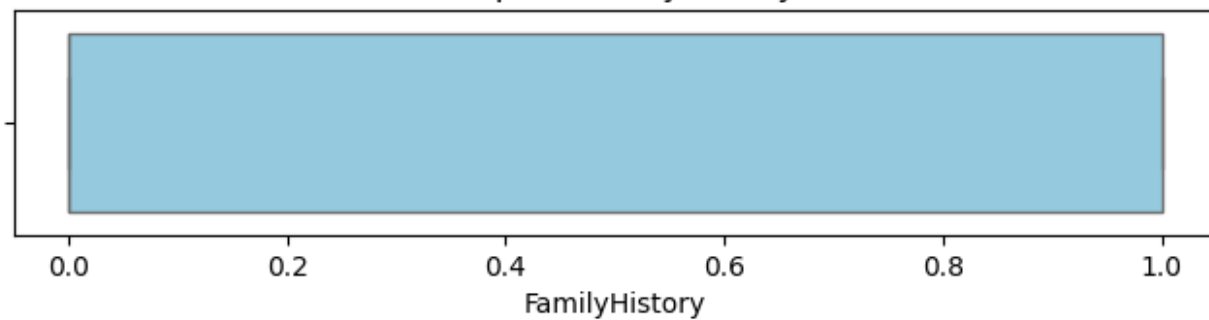
Boxplot: Diabetes



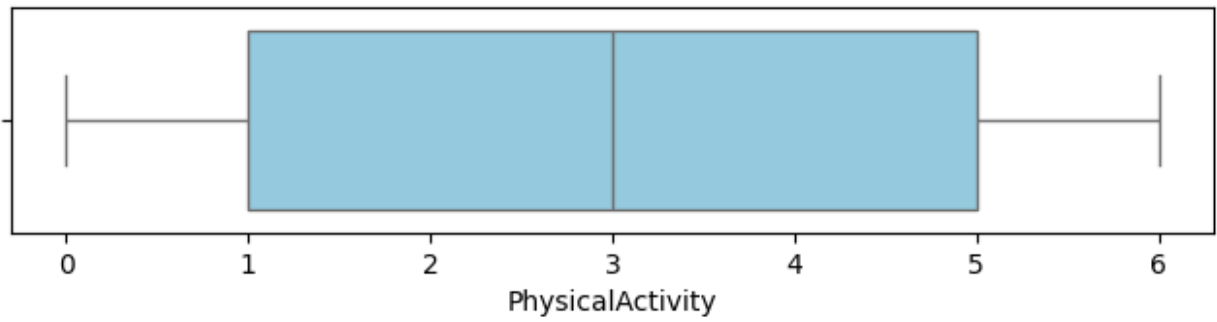
Boxplot: Hypertension



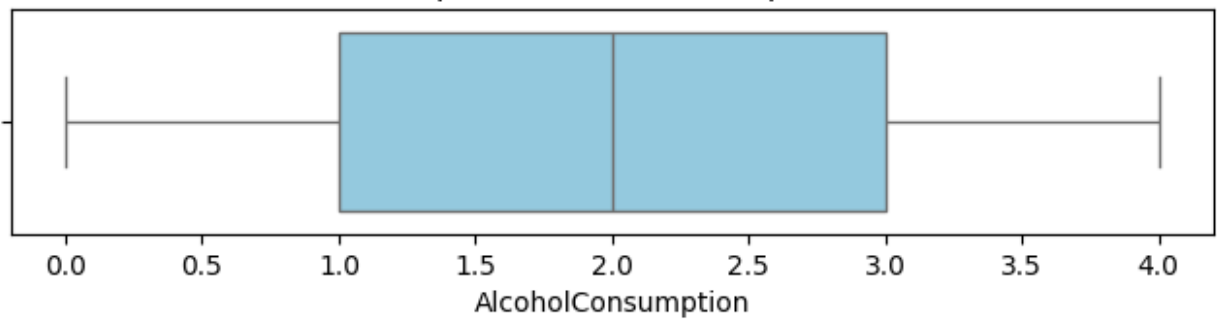
Boxplot: FamilyHistory



Boxplot: PhysicalActivity



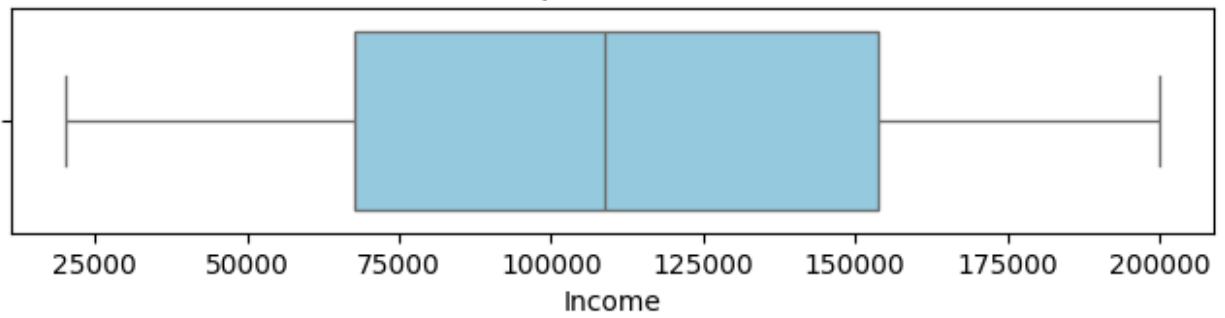
Boxplot: AlcoholConsumption

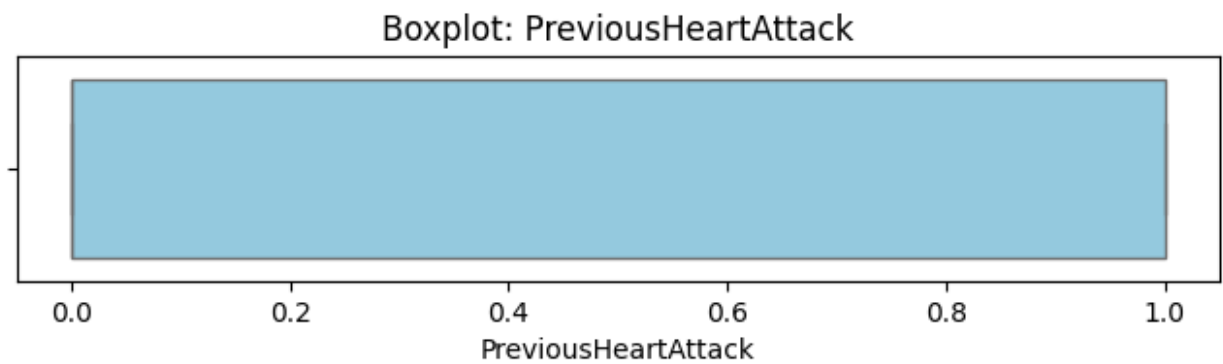
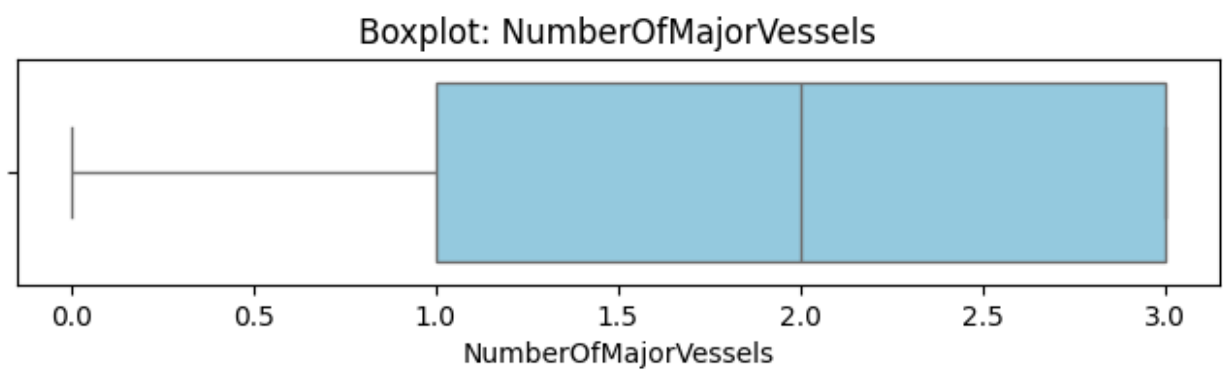
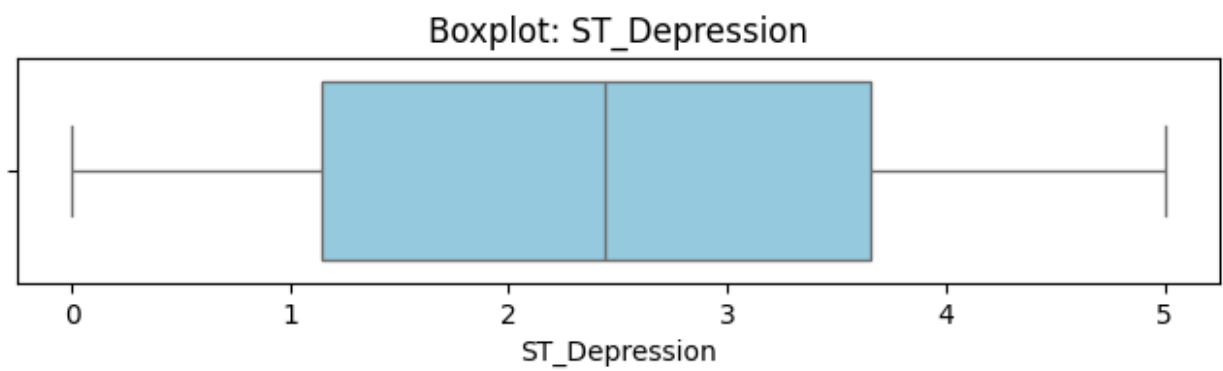
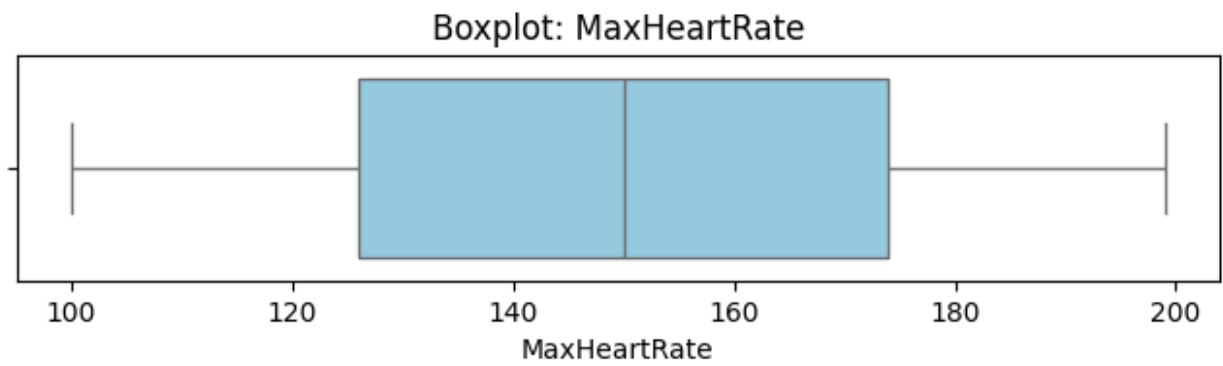


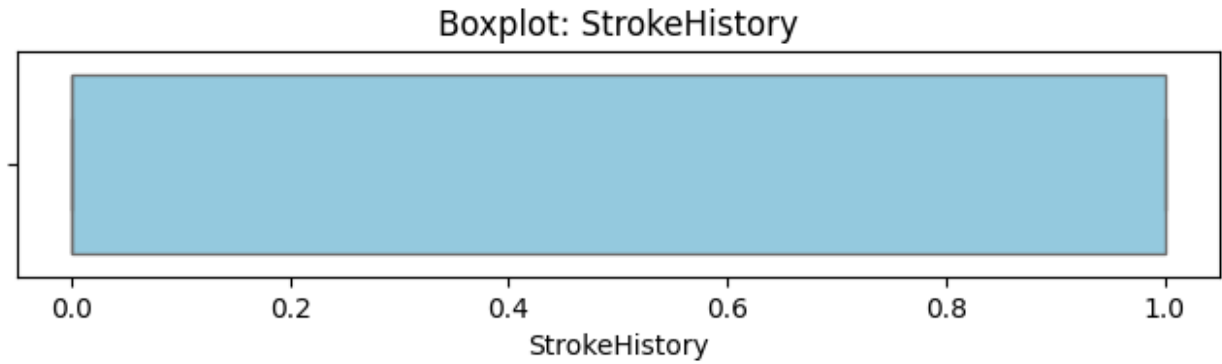
Boxplot: StressLevel



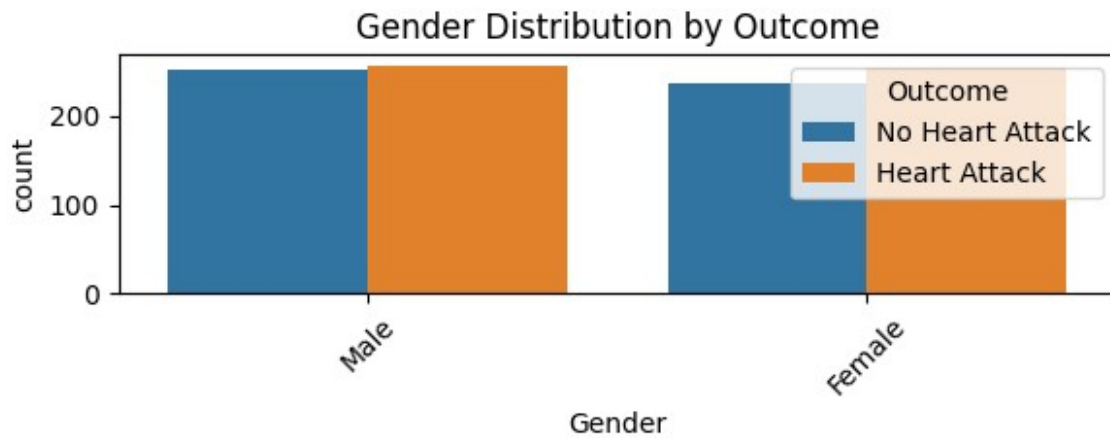
Boxplot: Income

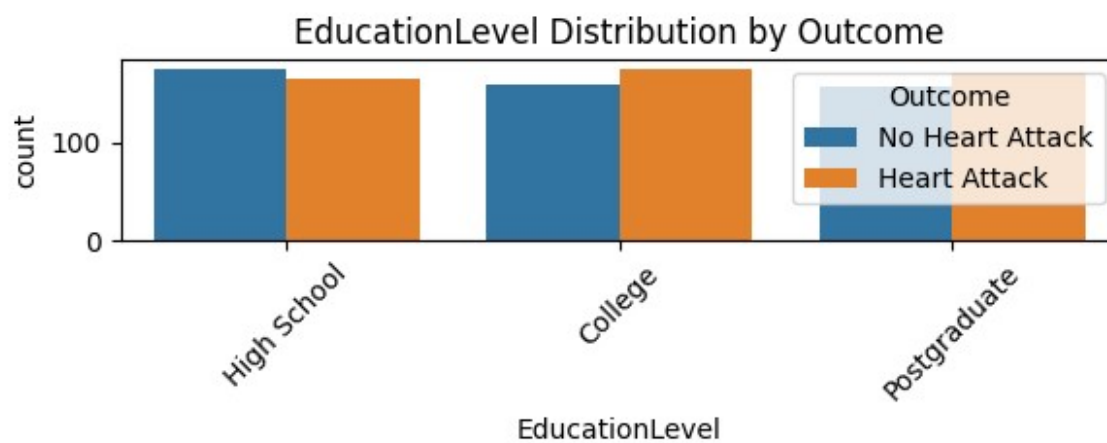
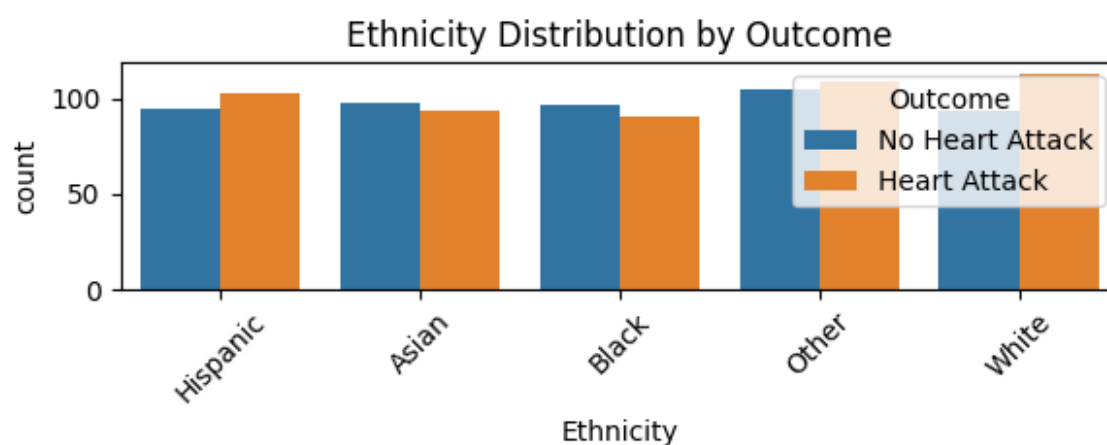
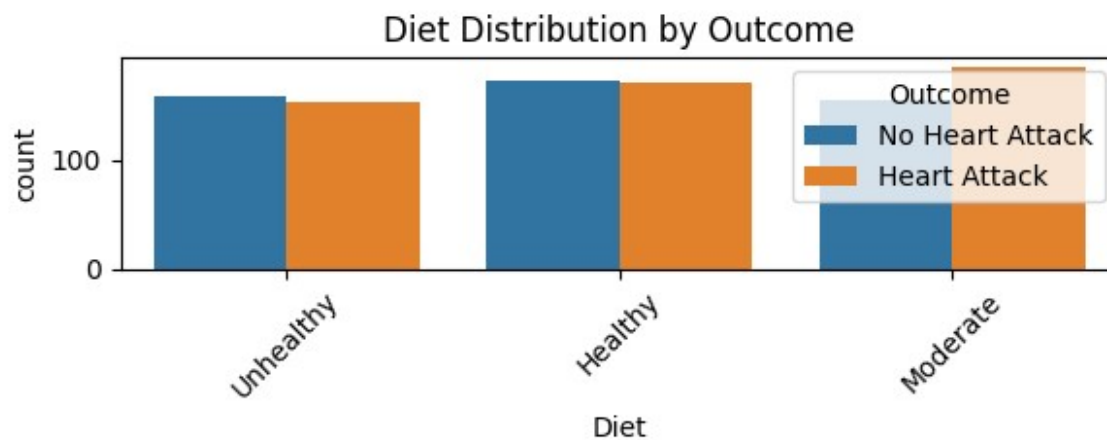


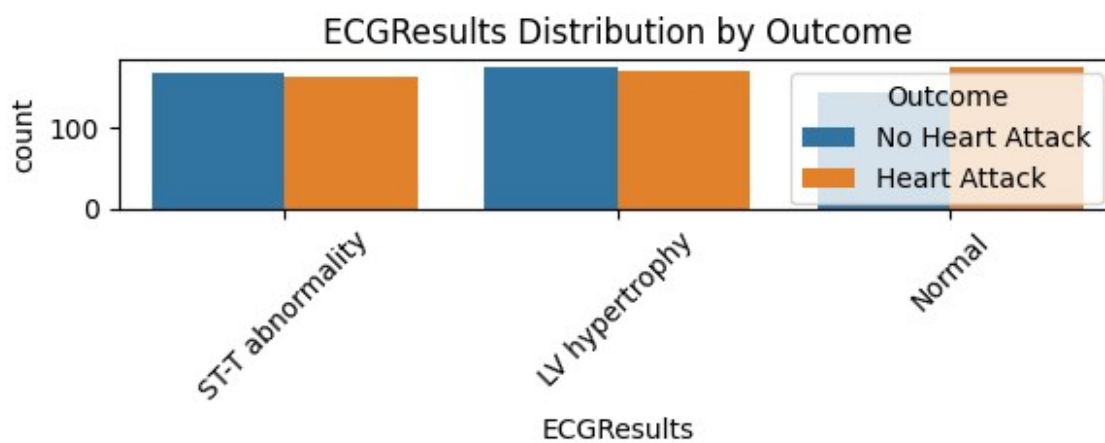
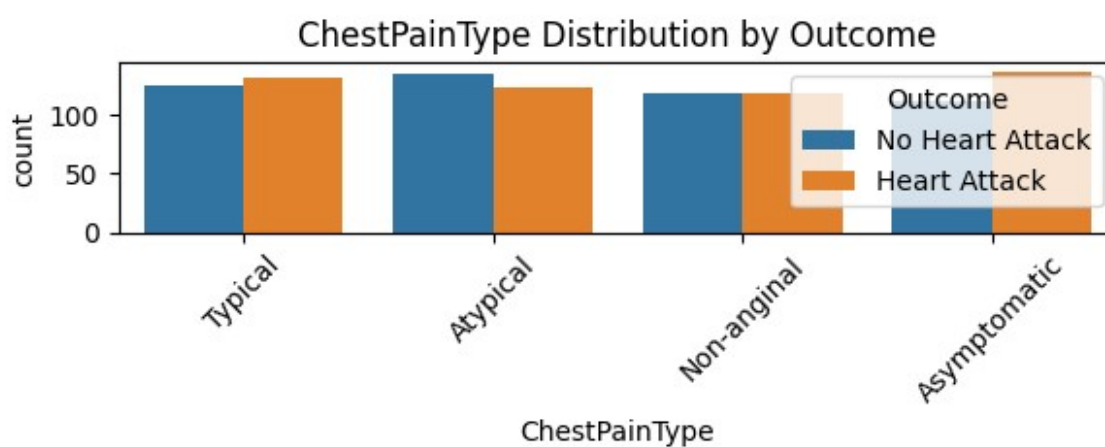
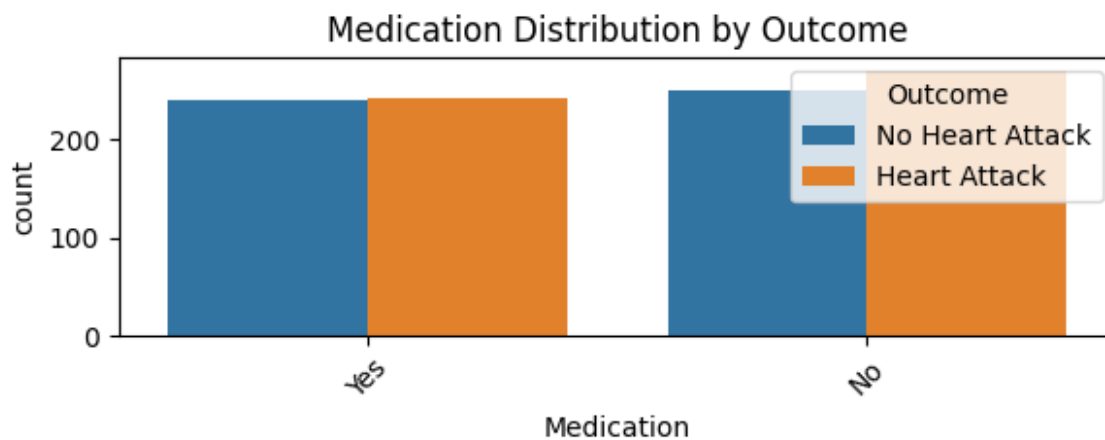


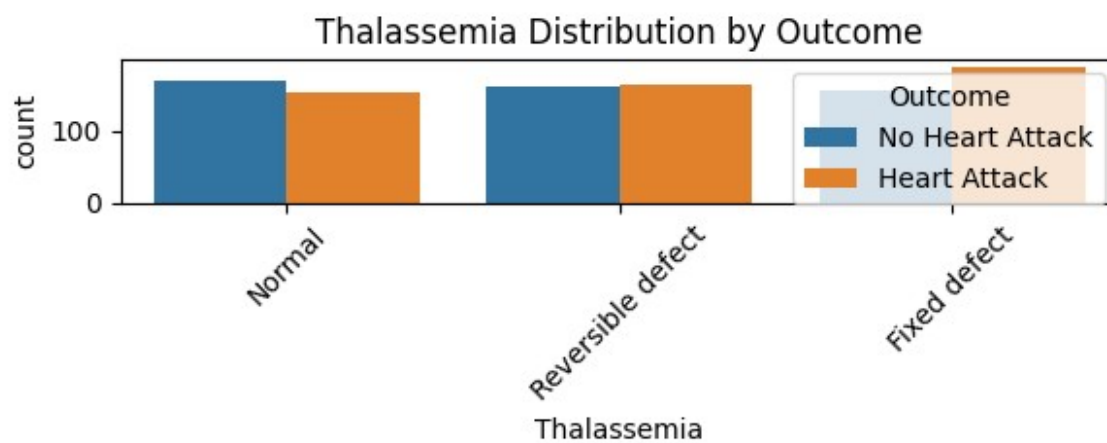
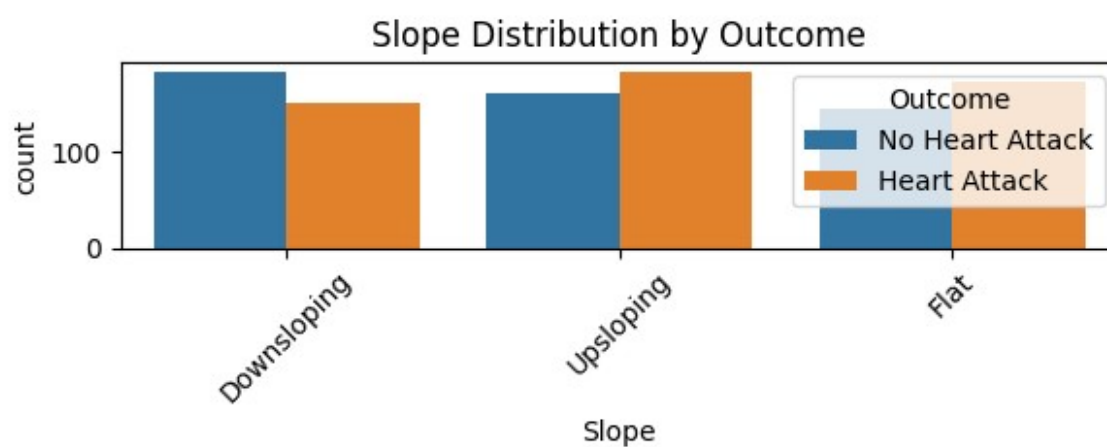
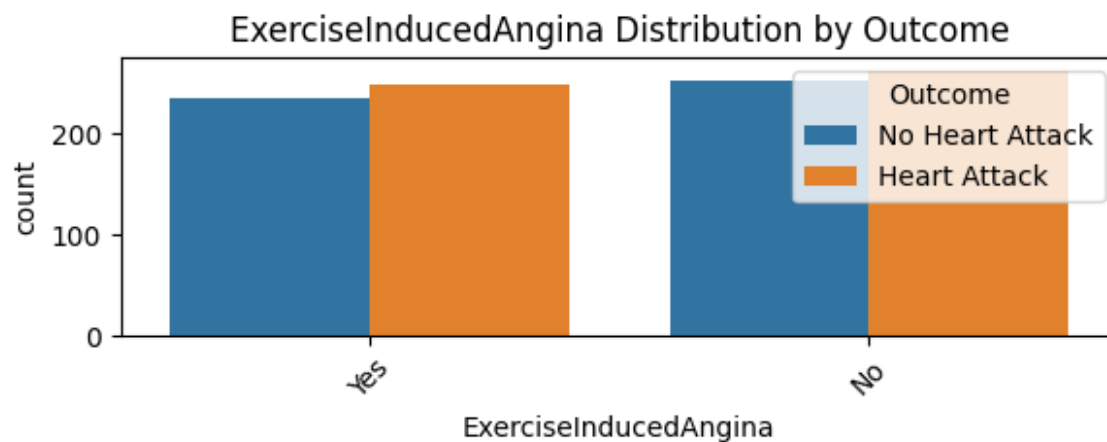


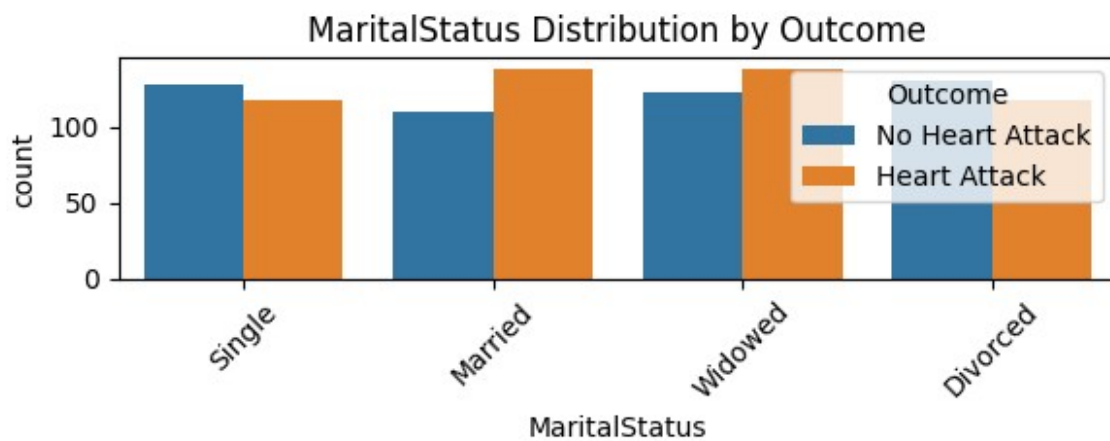
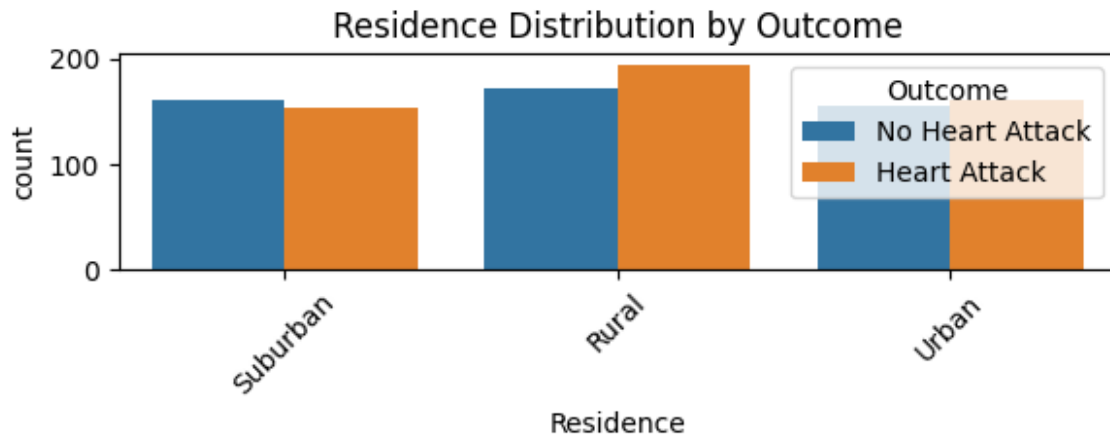
```
# Visualize the distribution of categorical features by outcome
cat_cols = df.select_dtypes(include='object').columns.drop('Outcome')
for col in cat_cols:
    plt.figure(figsize=(6,2.5))
    sns.countplot(data=df, x=col, hue="Outcome")
    plt.title(f"{col} Distribution by Outcome")
    plt.xticks(rotation=45)
    plt.tight_layout()
    plt.show()
```



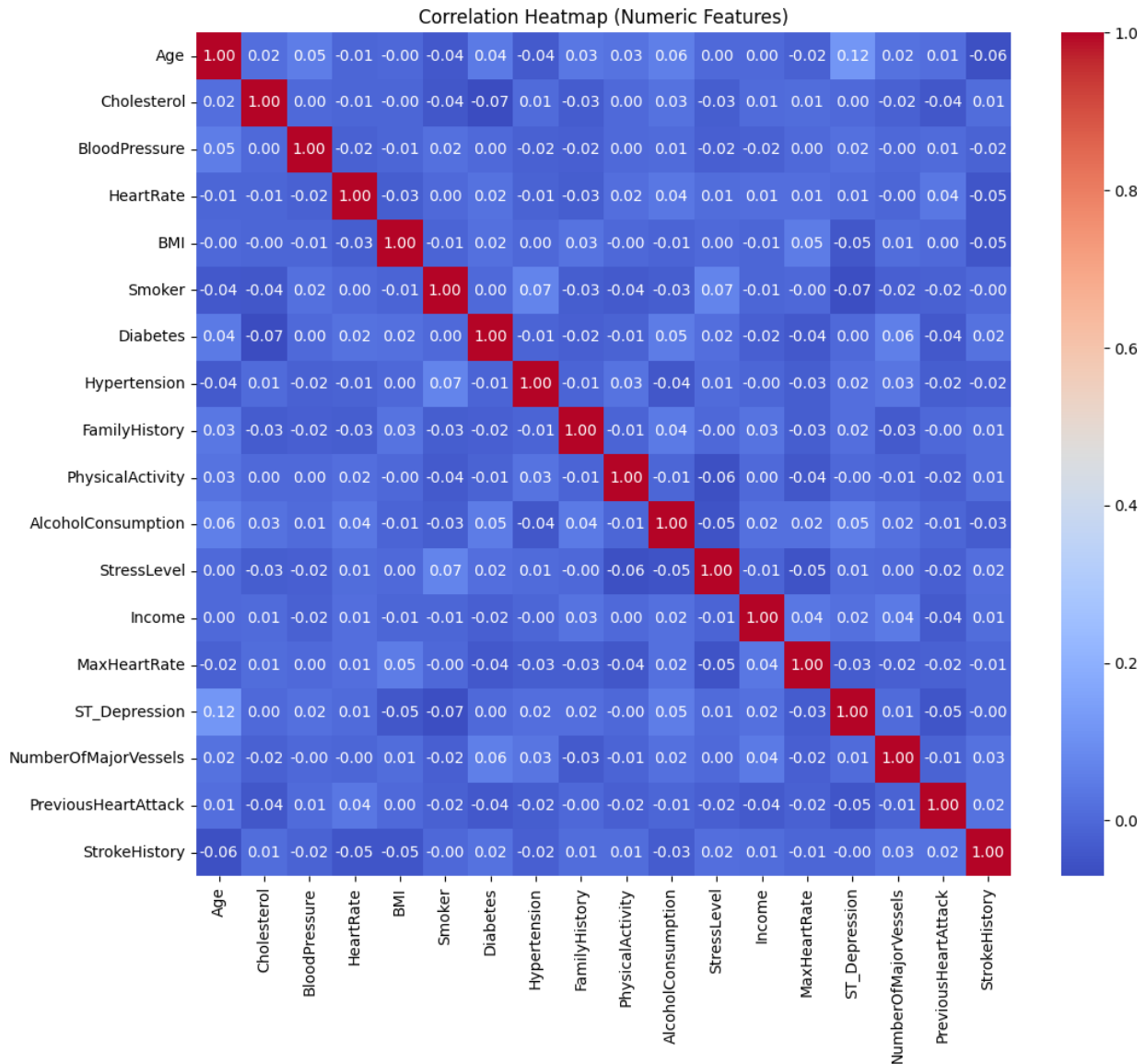








```
# Create a heatmap to visualize correlations between numeric features
plt.figure(figsize=(12,10))
corr = df[num_cols].corr()
sns.heatmap(corr, annot=True, cmap='coolwarm', fmt=".2f")
plt.title("Correlation Heatmap (Numeric Features)")
plt.show()
```



```
# Calculate and display the top features correlated with the outcome
df_corr = df.copy()
df_corr['OutcomeCode'] =
df_corr['Outcome'].astype('category').cat.codes
corr_with_outcome =
df_corr[num_cols].corrwith(df_corr['OutcomeCode']).abs().sort_values(a
scending=False)
print("\n===== TOP 10 NUMERIC FEATURES CORRELATED WITH OUTCOME =====")
print(corr_with_outcome.head(10))

===== TOP 10 NUMERIC FEATURES CORRELATED WITH OUTCOME =====
MaxHeartRate          0.050449
PhysicalActivity       0.042116
Hypertension          0.036982
```

```

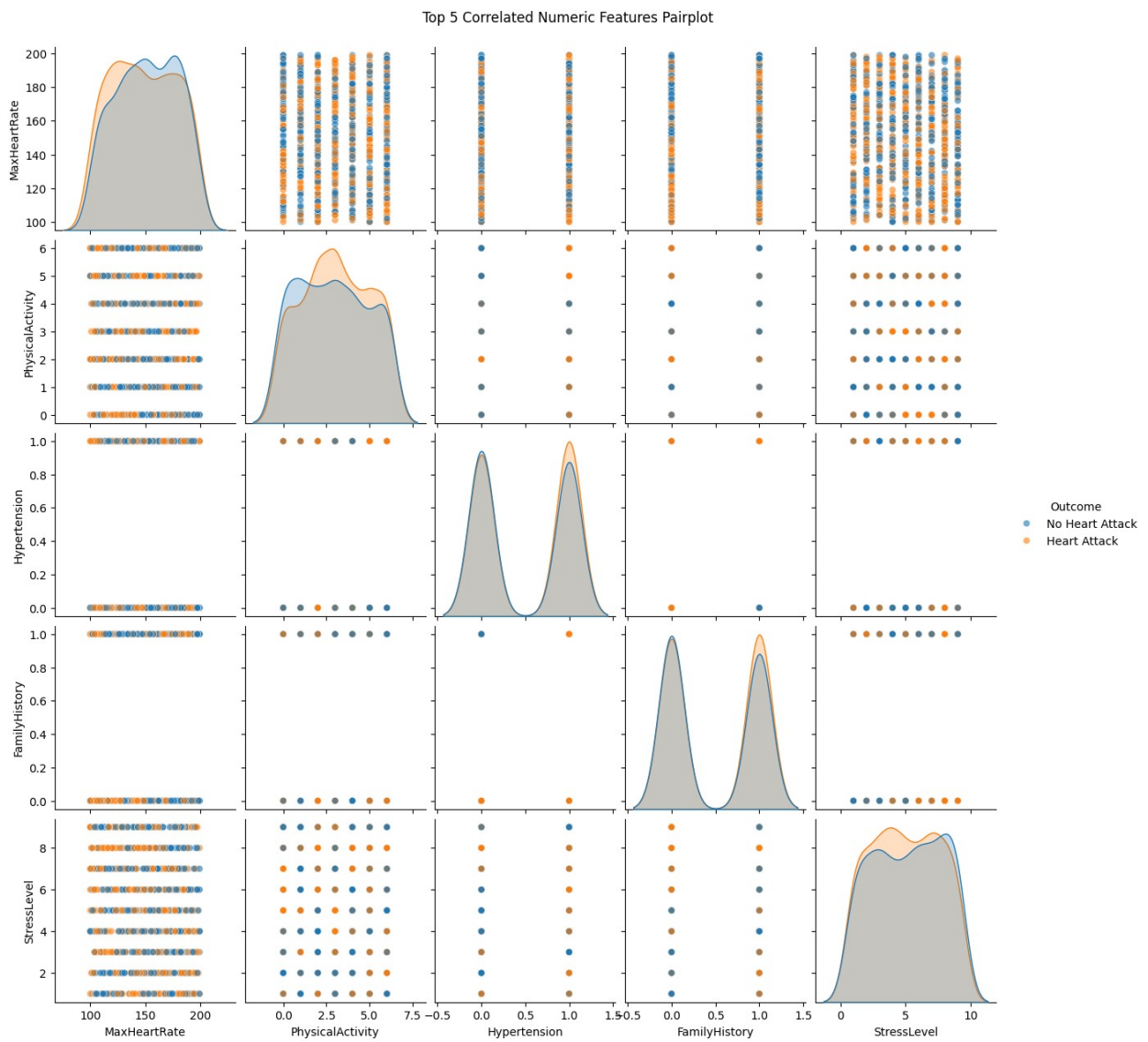
FamilyHistory      0.033490
StressLevel        0.033041
Smoker             0.031883
AlcoholConsumption 0.030888
Cholesterol         0.020881
StrokeHistory      0.020712
ST_Depression      0.018160
dtype: float64

```

```

# Create a pairplot for the top 5 correlated numeric features
top5 = corr_with_outcome.head(5).index.tolist()
sns.pairplot(df, vars=top5, hue="Outcome", plot_kws={'alpha':0.6})
plt.suptitle("Top 5 Correlated Numeric Features Pairplot", y=1.02)
plt.show()

```



```
# Display the distribution of categorical features by outcome
print("\n===== CATEGORICAL FEATURE DISTRIBUTION BY OUTCOME =====")
for col in cat_cols:
    cross = pd.crosstab(df[col], df['Outcome'], normalize='index')
    print(f"\n{col} (proportion of Outcome):")
    print(cross)
```

===== CATEGORICAL FEATURE DISTRIBUTION BY OUTCOME =====

Gender (proportion of Outcome):

Outcome	Heart Attack	No Heart Attack
Gender		
Female	0.518367	0.481633
Male	0.502947	0.497053

Diet (proportion of Outcome):

Outcome	Heart Attack	No Heart Attack
Diet		
Healthy	0.497110	0.502890
Moderate	0.542522	0.457478
Unhealthy	0.490385	0.509615

Ethnicity (proportion of Outcome):

Outcome	Heart Attack	No Heart Attack
Ethnicity		
Asian	0.489583	0.510417
Black	0.484043	0.515957
Hispanic	0.520202	0.479798
Other	0.509346	0.490654
White	0.545894	0.454106

EducationLevel (proportion of Outcome):

Outcome	Heart Attack	No Heart Attack
EducationLevel		
College	0.525526	0.474474
High School	0.486726	0.513274
Postgraduate	0.519878	0.480122

Medication (proportion of Outcome):

Outcome	Heart Attack	No Heart Attack
Medication		
No	0.519305	0.480695
Yes	0.501040	0.498960

ChestPainType (proportion of Outcome):

Outcome	Heart Attack	No Heart Attack
ChestPainType		
Asymptomatic	0.552419	0.447581
Atypical	0.476744	0.523256

Non-anginal	0.500000	0.500000
Typical	0.513619	0.486381

ECGResults (proportion of Outcome):

Outcome	Heart Attack	No Heart Attack
---------	--------------	-----------------

ECGResults

LV hypertrophy	0.494220	0.505780
Normal	0.548287	0.451713
ST-T abnormality	0.490964	0.509036

ExerciseInducedAngina (proportion of Outcome):

Outcome	Heart Attack	No Heart Attack
---------	--------------	-----------------

ExerciseInducedAngina

No	0.508738	0.491262
Yes	0.512397	0.487603

Slope (proportion of Outcome):

Outcome	Heart Attack	No Heart Attack
---------	--------------	-----------------

Slope

Downsloping	0.453731	0.546269
Flat	0.545455	0.454545
Upsloping	0.533333	0.466667

Thalassemia (proportion of Outcome):

Outcome	Heart Attack	No Heart Attack
---------	--------------	-----------------

Thalassemia

Fixed defect	0.549133	0.450867
Normal	0.473846	0.526154
Reversible defect	0.506098	0.493902

Residence (proportion of Outcome):

Outcome	Heart Attack	No Heart Attack
---------	--------------	-----------------

Residence

Rural	0.529891	0.470109
Suburban	0.488889	0.511111
Urban	0.509494	0.490506

EmploymentStatus (proportion of Outcome):

Outcome	Heart Attack	No Heart Attack
---------	--------------	-----------------

EmploymentStatus

Employed	0.536657	0.463343
Retired	0.489614	0.510386
Unemployed	0.504673	0.495327

MaritalStatus (proportion of Outcome):

Outcome	Heart Attack	No Heart Attack
---------	--------------	-----------------

MaritalStatus

Divorced	0.473684	0.526316
Married	0.556452	0.443548

Single	0.479508	0.520492
Widowed	0.530769	0.469231

```
print("\n==== EDA COMPLETE ====")
```

```
==== EDA COMPLETE ====
```