

✓ Task

predict insurance premiums based on various factors.

Here is all the data you need: "sample_submission.csv" "test.csv" "train.csv"

✓ Data loading

Subtask:

Load the provided CSV files into pandas DataFrames.

Reasoning: Load the three CSV files into pandas DataFrames.

```
1 import pandas as pd
2
3 df_train = pd.read_csv('train.csv')
4 df_test = pd.read_csv('test.csv')
5 df_submission = pd.read_csv('sample_submission.csv')
6
7 display(df_train.head())
8 display(df_test.head())
9 display(df_submission.head())
```



th re	Location	...	Previous Claims	Vehicle Age	Credi Scor
61	Urban	...	2.0	17.0	372
31	Rural	...	1.0	12.0	694
49	Suburban	...	1.0	14.0	Ne
44	Rural	...	1.0	0.0	367
94	Rural	...	0.0	8.0	598

Health Score	Location	Policy Type	Previous Claims
7.657981	Rural	Basic	NaN
3.381379	Suburban	Premium	NaN
4.354527	Urban	Comprehensive	NaN
5.136225	Suburban	Comprehensive	1.0
1.844155	Suburban	Premium	NaN

✓ Data exploration

Subtask:

Explore the training dataset (`df_train`) to understand its characteristics.

Reasoning: Explore the training dataset (`df_train`) by examining its shape, data types, descriptive statistics, missing values, and data distributions of numerical and categorical features. This will provide initial insights into the data's characteristics, potential outliers, and relationships between variables.

```
1 import matplotlib.pyplot as plt
2 import seaborn as sns
3
4 # Data Shape and Info
5 print("Shape of df_train:", df_train.shape)
6 print("\nInfo of df_train:")
7 display(df_train.info())
8
9 # Descriptive Statistics
10 print("\nDescriptive Statistics of Numerical Features:")
11 display(df_train.describe())
12
13 # Missing Values
14 print("\nMissing Values:")
15 missing_values = df_train.isnull().sum()
16 missing_percentage = (missing_values / len(df_train)) * 100
17 missing_df = pd.DataFrame({'Missing Values': missing_values, 'Percentage':
```

```
missing_percentage})
18 display(missing_df)
19
20 # Visualization of missing data pattern (heatmap)
21 plt.figure(figsize=(12, 6))
22 sns.heatmap(df_train.isnull(), cbar=False, yticklabels=False, cmap='viridis')
23 plt.title('Missing Values Heatmap')
24 plt.show()
25
26 # Data Distribution
27 numerical_features = df_train.select_dtypes(include=['number']).columns
28 categorical_features = df_train.select_dtypes(exclude=['number']).columns
29
30 # Histograms for numerical features
31 for col in numerical_features:
32     plt.figure(figsize=(8, 4))
33     sns.histplot(df_train[col], kde=True)
34     plt.title(f'Distribution of {col}')
35     plt.show()
36
37 # Value counts for categorical features
38 for col in categorical_features:
39     print(f"\nValue counts for {col}:")
40     display(df_train[col].value_counts())
41
42 # Relationships between numerical features and target variable
43 for col in numerical_features:
44     if col != 'Premium Amount' and col != 'id':
45         plt.figure(figsize=(8, 4))
46         sns.scatterplot(x=col, y='Premium Amount', data=df_train)
47         plt.title(f'Relationship between {col} and Premium Amount')
```

```
48     plt.show()
49
50 # Relationships between categorical features and target variable
51 for col in categorical_features:
52     if col != 'Policy Start Date': # Exclude date column for boxplot
53         plt.figure(figsize=(10, 6))
54         sns.boxplot(x=col, y='Premium Amount', data=df_train)
55         plt.title(f'Relationship between {col} and Premium Amount')
56         plt.xticks(rotation=45, ha='right') # Rotate x-axis labels for better
        readability
57     plt.show()
58
59 # Outlier detection (using boxplots)
60 for col in numerical_features:
61     if col != 'Premium Amount' and col != 'id':
62         plt.figure(figsize=(8, 4))
63         sns.boxplot(x=col, data=df_train)
64         plt.title(f'Boxplot of {col}')
65         plt.show()
```

➡ Shape of df_train: (1200000, 21)

Info of df_train:

<class 'pandas.core.frame.DataFrame'>

RangeIndex: 1200000 entries, 0 to 1199999

Data columns (total 21 columns):

#	Column	Non-Null Count	Dtype
0	id	1200000 non-null	int64
1	Age	1181295 non-null	float64
2	Gender	1200000 non-null	object
3	Annual Income	1155051 non-null	float64
4	Marital Status	1181471 non-null	object
5	Number of Dependents	1090328 non-null	float64
6	Education Level	1200000 non-null	object
7	Occupation	841925 non-null	object
8	Health Score	1125924 non-null	float64
9	Location	1200000 non-null	object
10	Policy Type	1200000 non-null	object
11	Previous Claims	835971 non-null	float64
12	Vehicle Age	1199994 non-null	float64
13	Credit Score	1062118 non-null	float64
14	Insurance Duration	1199999 non-null	float64
15	Policy Start Date	1200000 non-null	object
16	Customer Feedback	1122176 non-null	object
17	Smoking Status	1200000 non-null	object
18	Exercise Frequency	1200000 non-null	object
19	Property Type	1200000 non-null	object
20	Premium Amount	1200000 non-null	float64

dtypes: float64(9), int64(1), object(11)

memory usage: 192.3+ MB



None

Descriptive Statistics of Numerical Features:

	id	Age	Annual Income	Number of Dependents	Health Score	Previous Claims	Vehicle Age	Credit Score	Insurance Duration
count	1.200000e+06	1.181295e+06	1.155051e+06	1.090328e+06	1.125924e+06	835971.000000	1.199994e+06	1.062118e+06	1.1999
mean	5.999995e+05	4.114556e+01	3.274522e+04	2.009934e+00	2.561391e+01	1.002689	9.569889e+00	5.929244e+02	5.0182
std	3.464103e+05	1.353995e+01	3.217951e+04	1.417338e+00	1.220346e+01	0.982840	5.776189e+00	1.499819e+02	2.5943
min	0.000000e+00	1.800000e+01	1.000000e+00	0.000000e+00	2.012237e+00	0.000000	0.000000e+00	3.000000e+02	1.0000

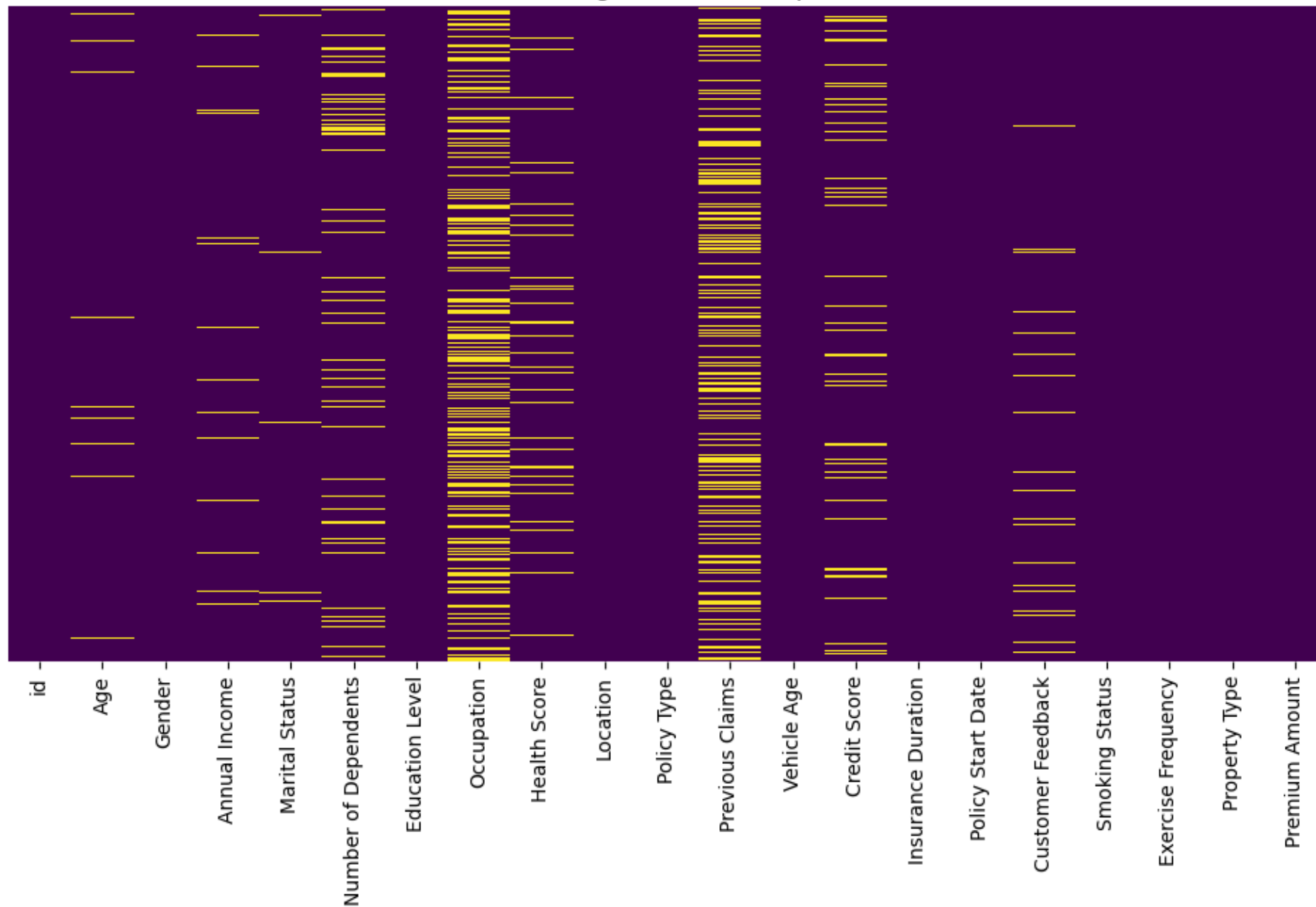
25%	2.999998e+05	3.000000e+01	8.001000e+03	1.000000e+00	1.591896e+01	0.000000	5.000000e+00	4.680000e+02	3.0000
50%	5.999995e+05	4.100000e+01	2.391100e+04	2.000000e+00	2.457865e+01	1.000000	1.000000e+01	5.950000e+02	5.0000
75%	8.999992e+05	5.300000e+01	4.463400e+04	3.000000e+00	3.452721e+01	2.000000	1.500000e+01	7.210000e+02	7.0000
max	1.199999e+06	6.400000e+01	1.499970e+05	4.000000e+00	5.897591e+01	9.000000	1.900000e+01	8.490000e+02	9.0000

Missing Values:

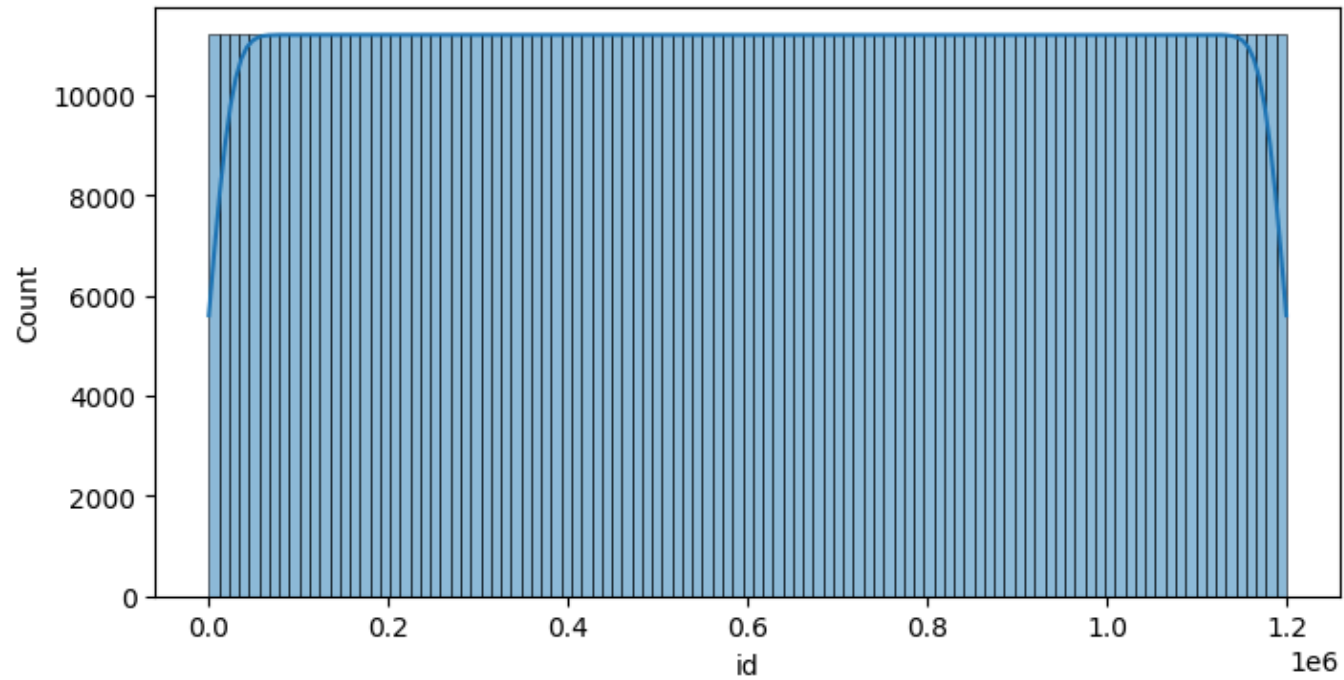
	Missing Values	Percentage	
id	0	0.000000	
Age	18705	1.558750	
Gender	0	0.000000	
Annual Income	44949	3.745750	
Marital Status	18529	1.544083	
Number of Dependents	109672	9.139333	
Education Level	0	0.000000	
Occupation	358075	29.839583	
Health Score	74076	6.173000	
Location	0	0.000000	
Policy Type	0	0.000000	
Previous Claims	364029	30.335750	
Vehicle Age	6	0.000500	
Credit Score	137882	11.490167	
Insurance Duration	1	0.000083	
Policy Start Date	0	0.000000	
Customer Feedback	77824	6.485333	
Smoking Status	0	0.000000	
Exercise Frequency	0	0.000000	

Property Type	0	0.000000
Premium Amount	0	0.000000

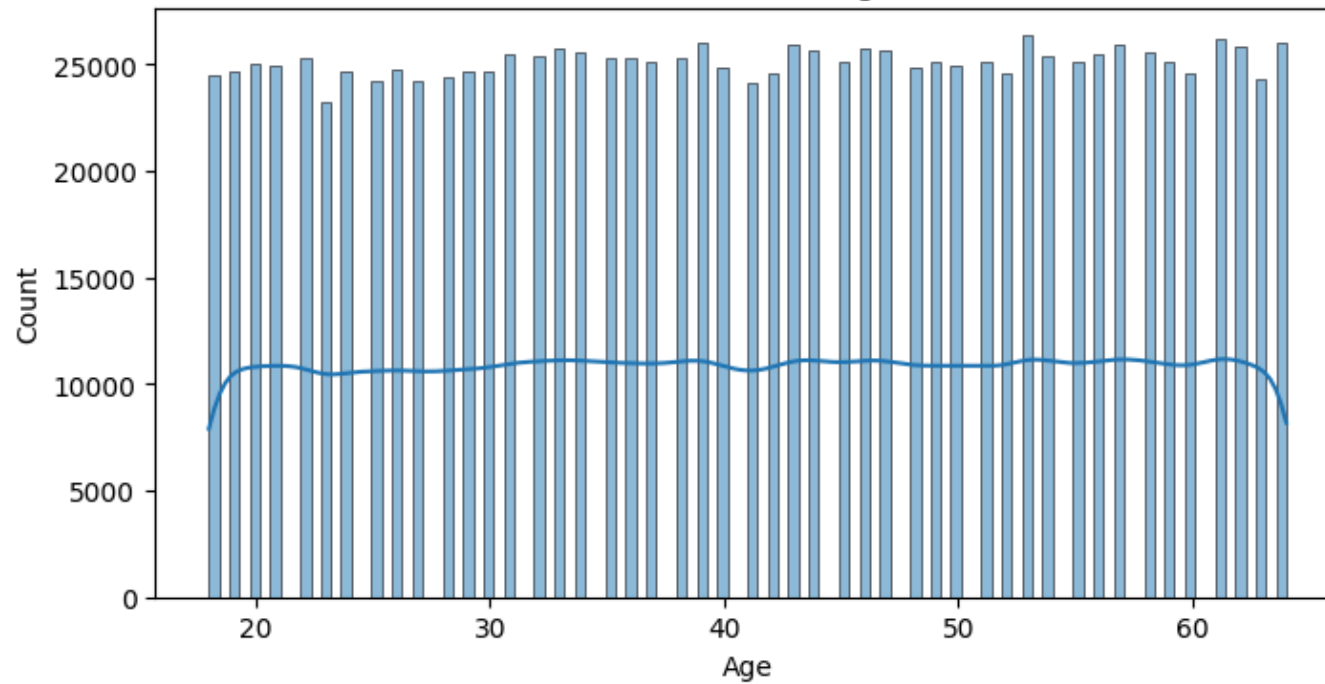
Missing Values Heatmap



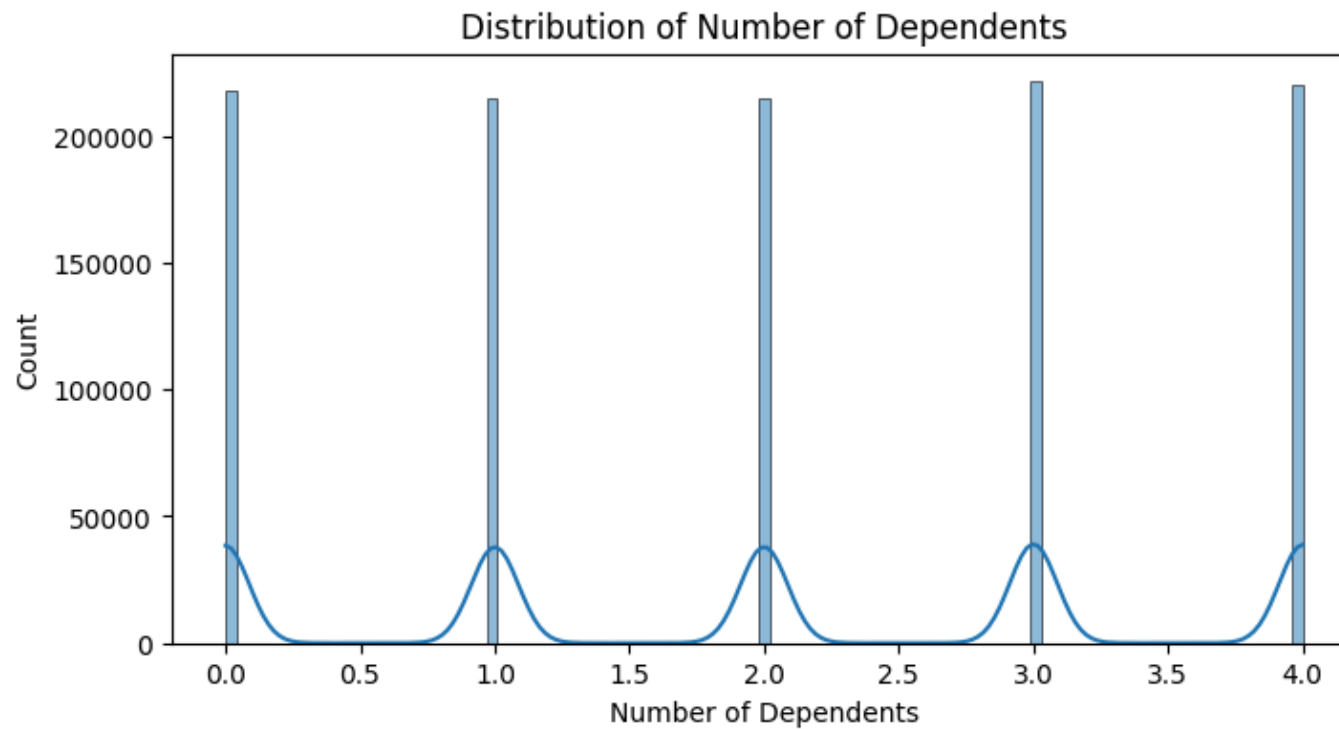
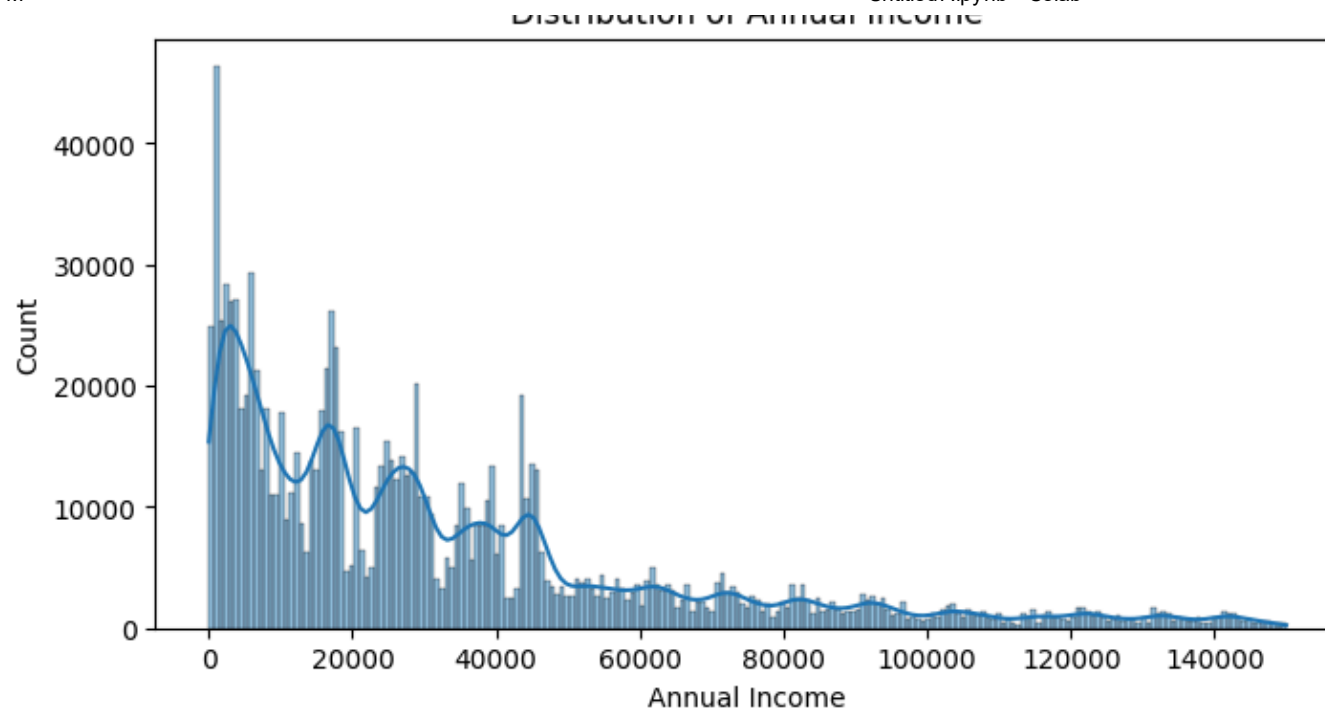
Distribution of id



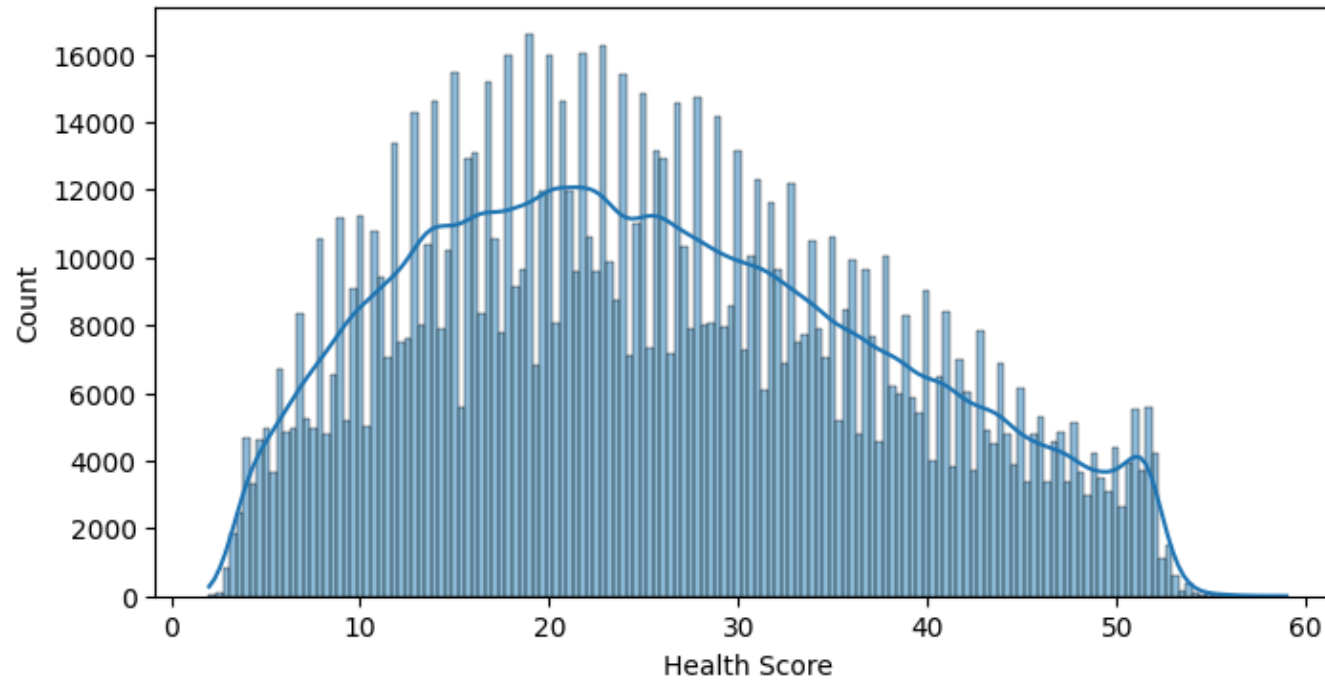
Distribution of Age



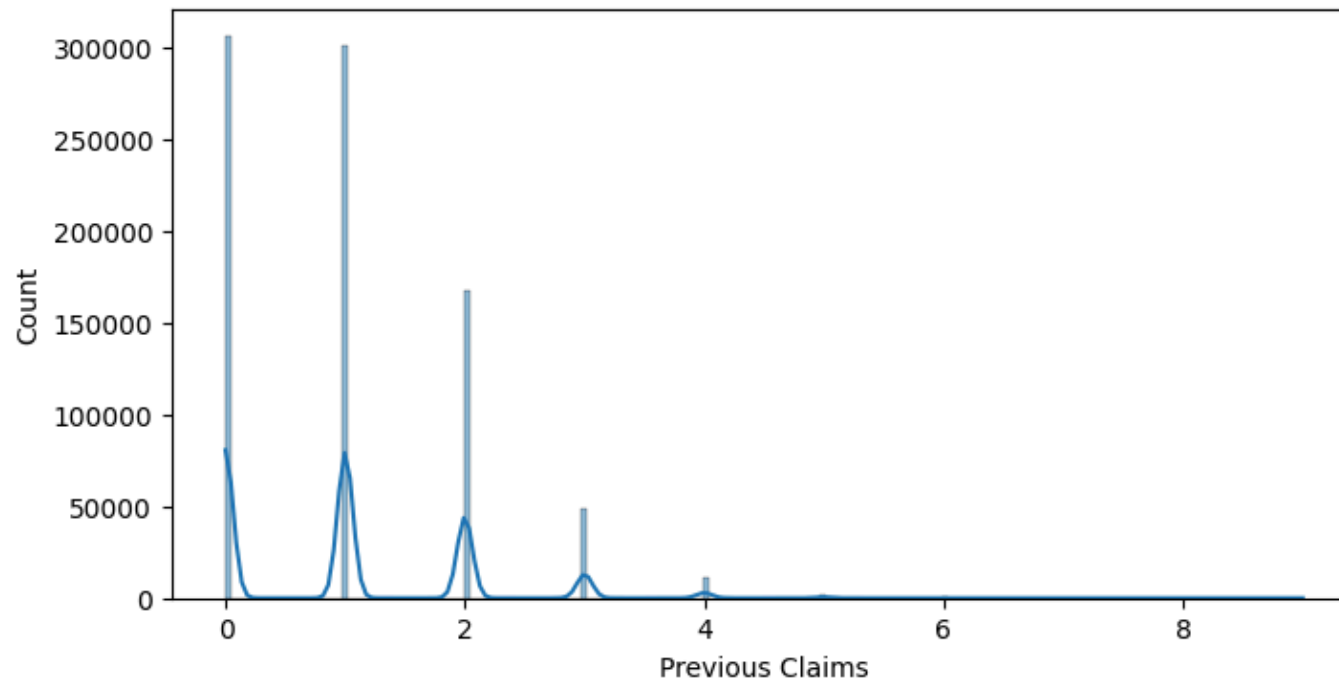
Distribution of Annual Income



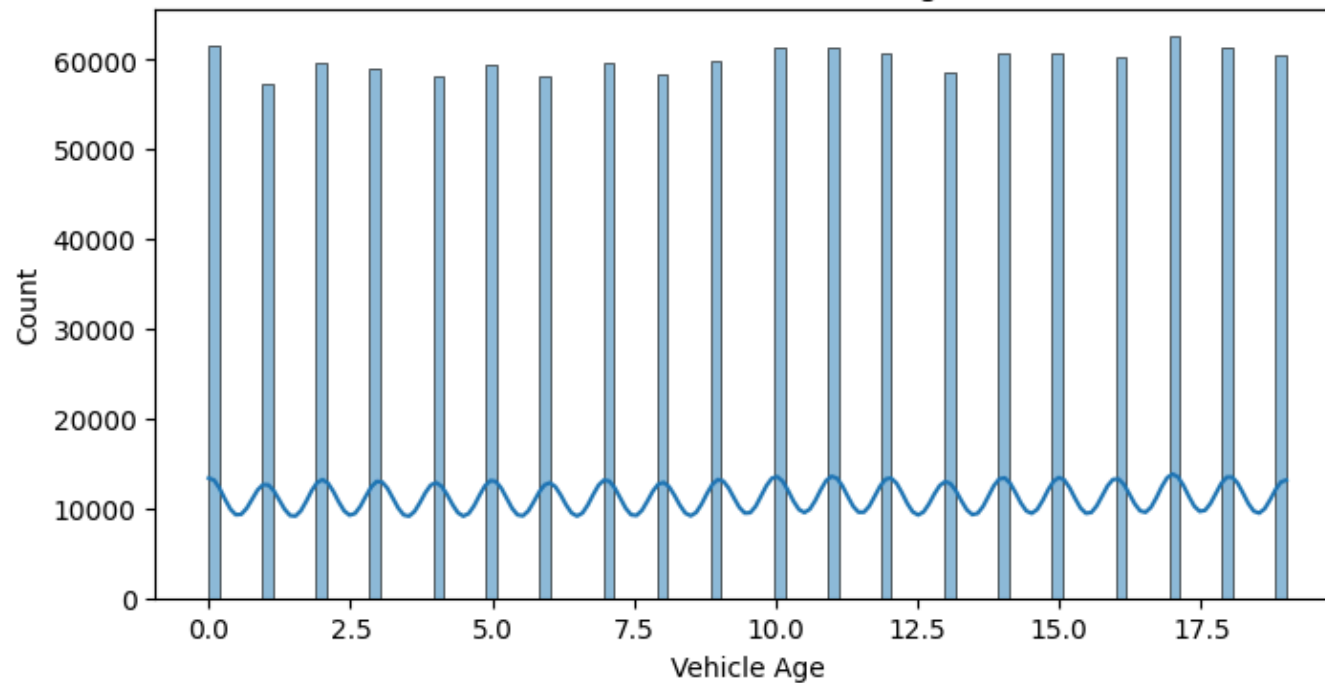
Distribution of Health Score



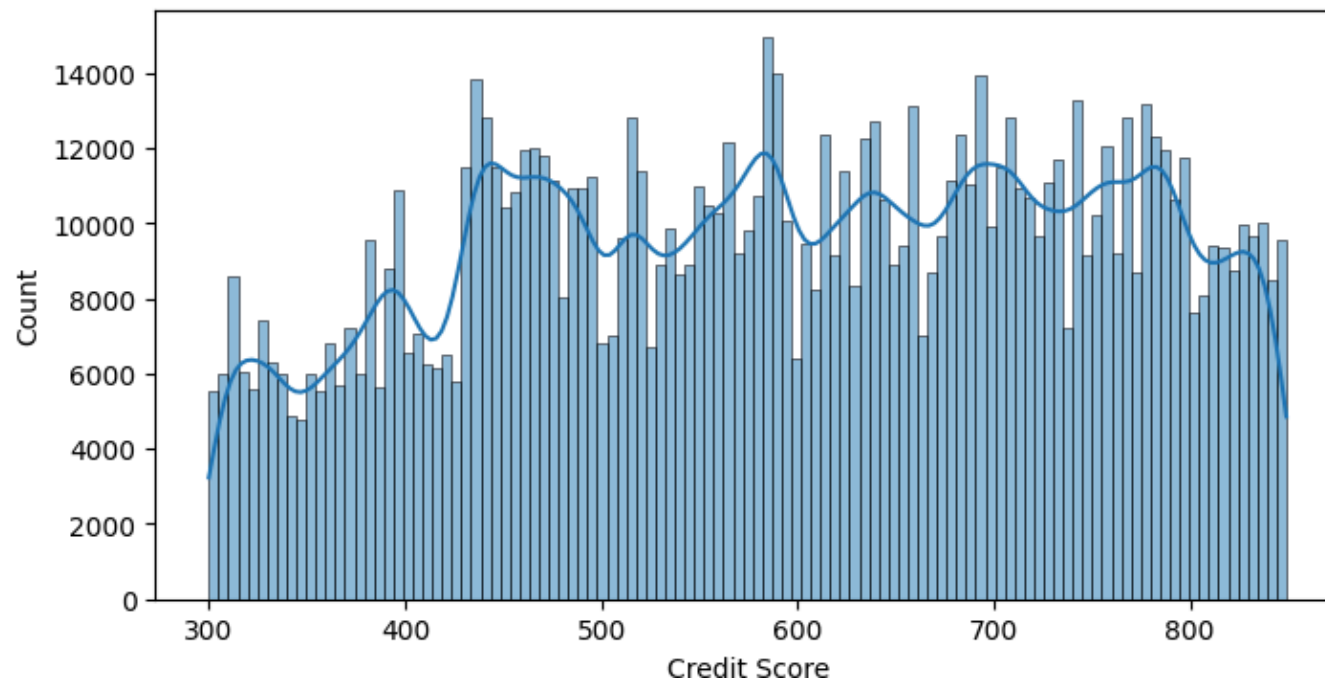
Distribution of Previous Claims



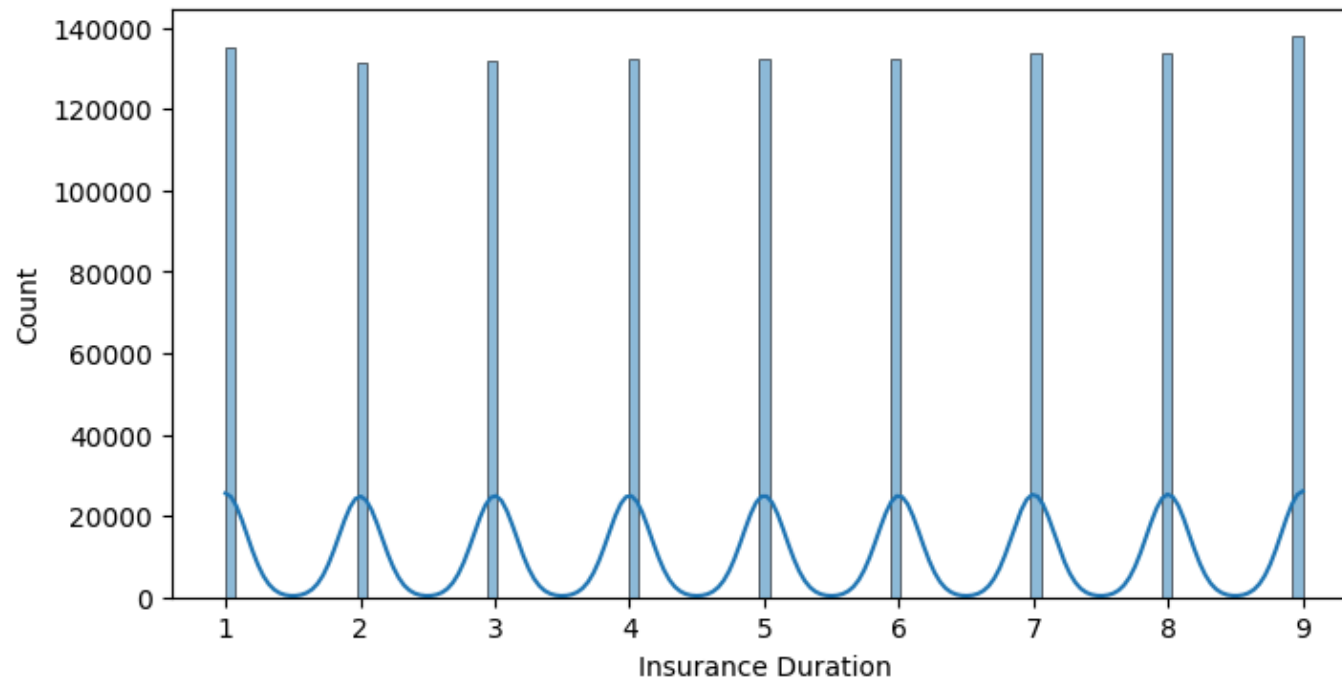
Distribution of Vehicle Age



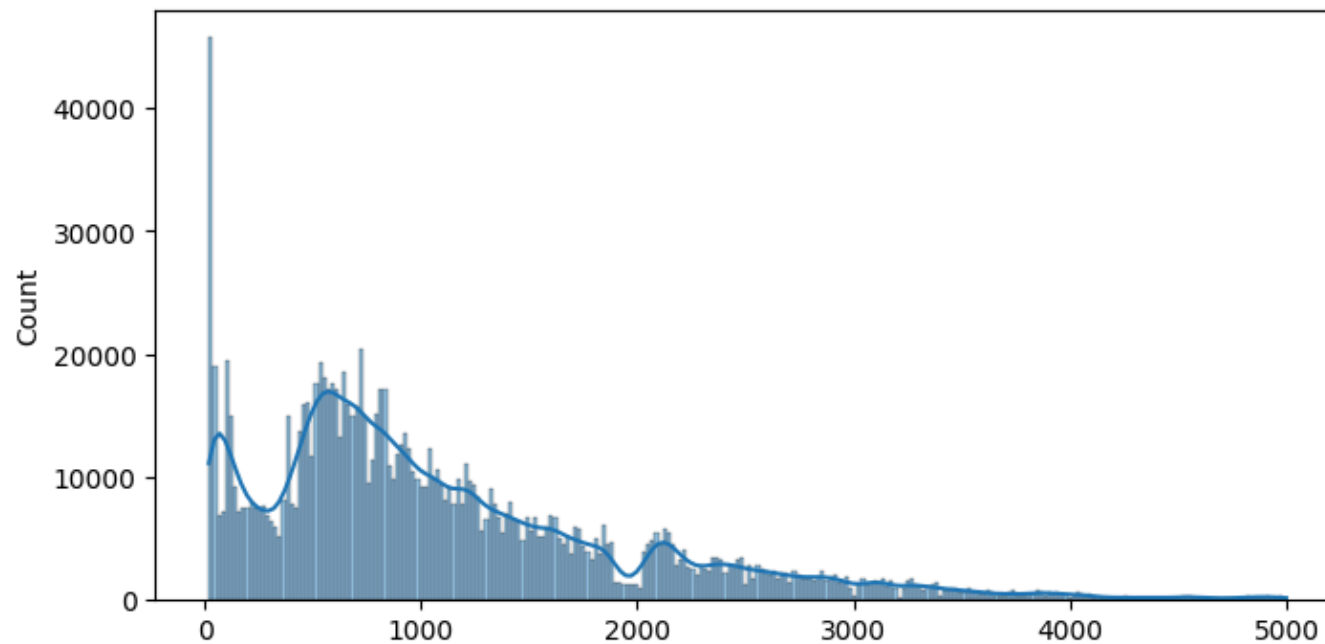
Distribution of Credit Score



Distribution of Insurance Duration



Distribution of Premium Amount



Premium Amount

Value counts for Gender:

count	
Gender	
Male	602571
Female	597429

dtype: int64

Value counts for Marital Status:

count	
Marital Status	
Single	395391
Married	394316
Divorced	391764

dtype: int64

Value counts for Education Level:

count	
Education Level	
Master's	303818
PhD	303507
Bachelor's	303234
High School	289441

dtype: int64

Value counts for Occupation:

count

Occupation

Employed	282750
Self-Employed	282645
Unemployed	276530

dtype: int64

Value counts for Location:

count

Location

Suburban	401542
Rural	400947
Urban	397511

dtype: int64

Value counts for Policy Type:

count

Policy Type

Premium	401846
Comprehensive	399600
Basic	398554

dtype: int64

Value counts for Policy Start Date:

count

Policy Start Date

2020-02-08 15:21:39.134960	142
2023-08-13 15:21:39.155231	137

```
2022-02-02 15:21:39.134960    137
2022-08-30 15:21:39.134960    134
2023-11-02 15:21:39.134960    118
...
2021-06-07 15:21:39.104139      1
2024-07-19 15:21:39.233998      1
2019-12-14 15:21:39.110557      1
2020-07-23 15:21:39.217387      1
2020-10-19 15:21:39.118178      1
```

167381 rows × 1 columns

dtype: int64

Value counts for Customer Feedback:

	count
Customer Feedback	
Average	377905
Poor	375518
Good	368753

dtype: int64

Value counts for Smoking Status:

	count
Smoking Status	
Yes	601873
No	598127

dtype: int64

Value counts for Exercise Frequency:

	count
Exercise Frequency	
Weekly	306179
Monthly	299830
Rarely	299420
Daily	294571

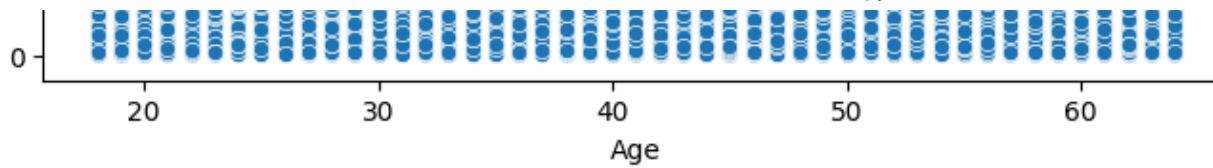
dtype: int64

Value counts for Property Type:

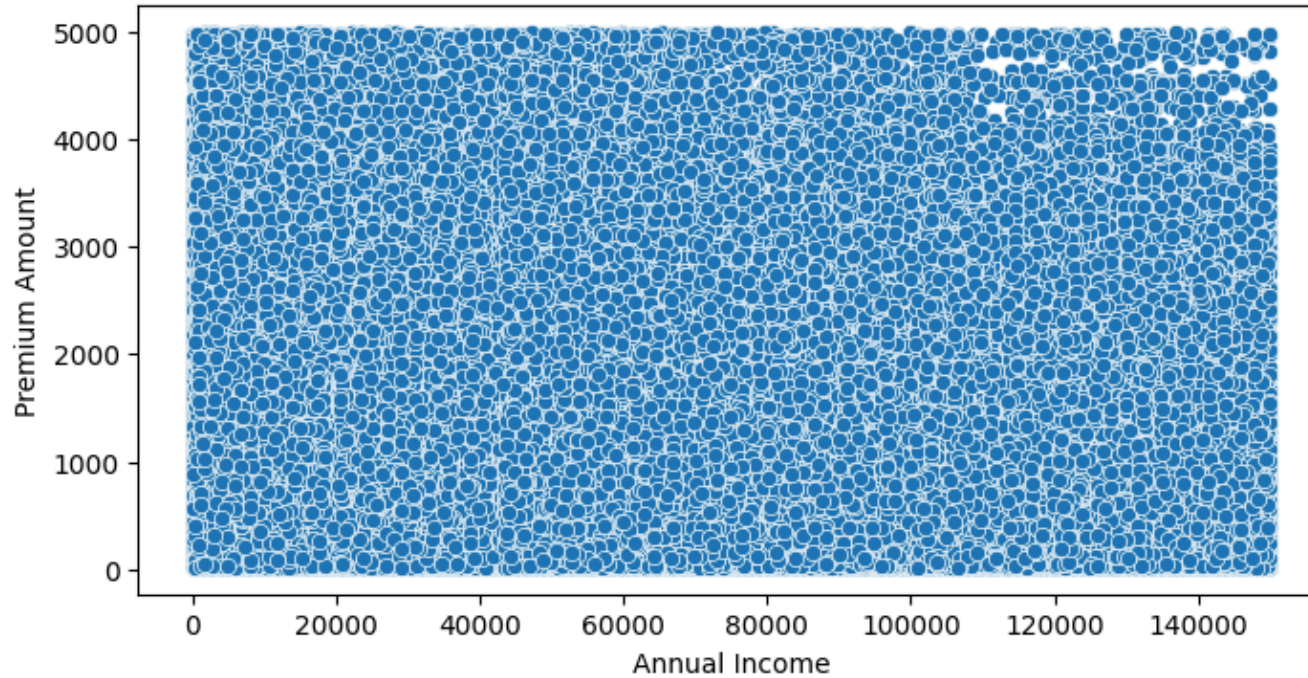
	count
Property Type	
House	400349
Apartment	399978
Condo	399673

dtype: int64



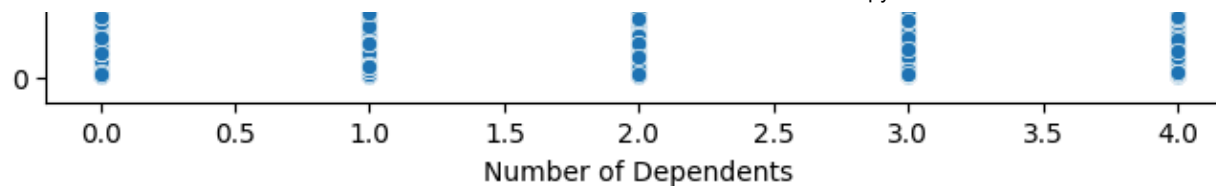


Relationship between Annual Income and Premium Amount

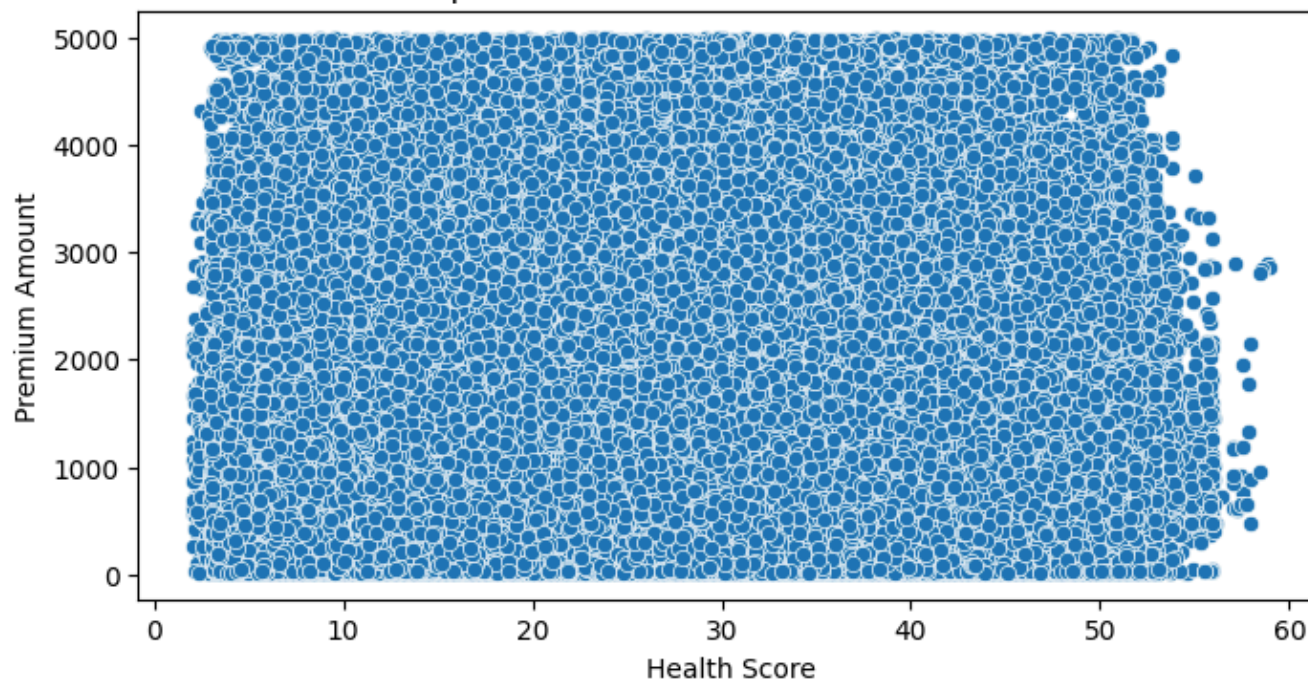


Relationship between Number of Dependents and Premium Amount

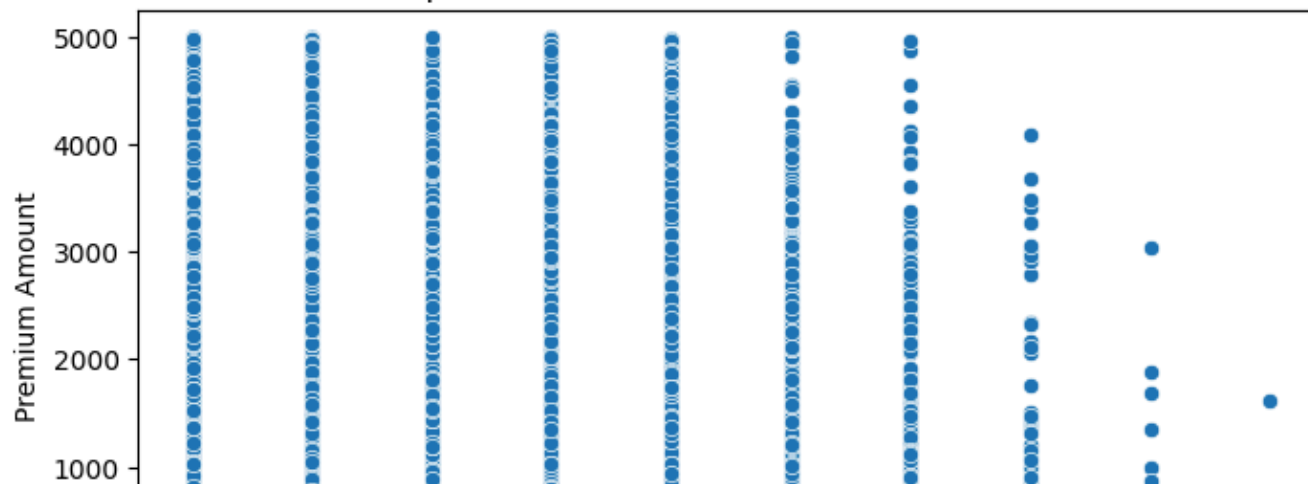




Relationship between Health Score and Premium Amount

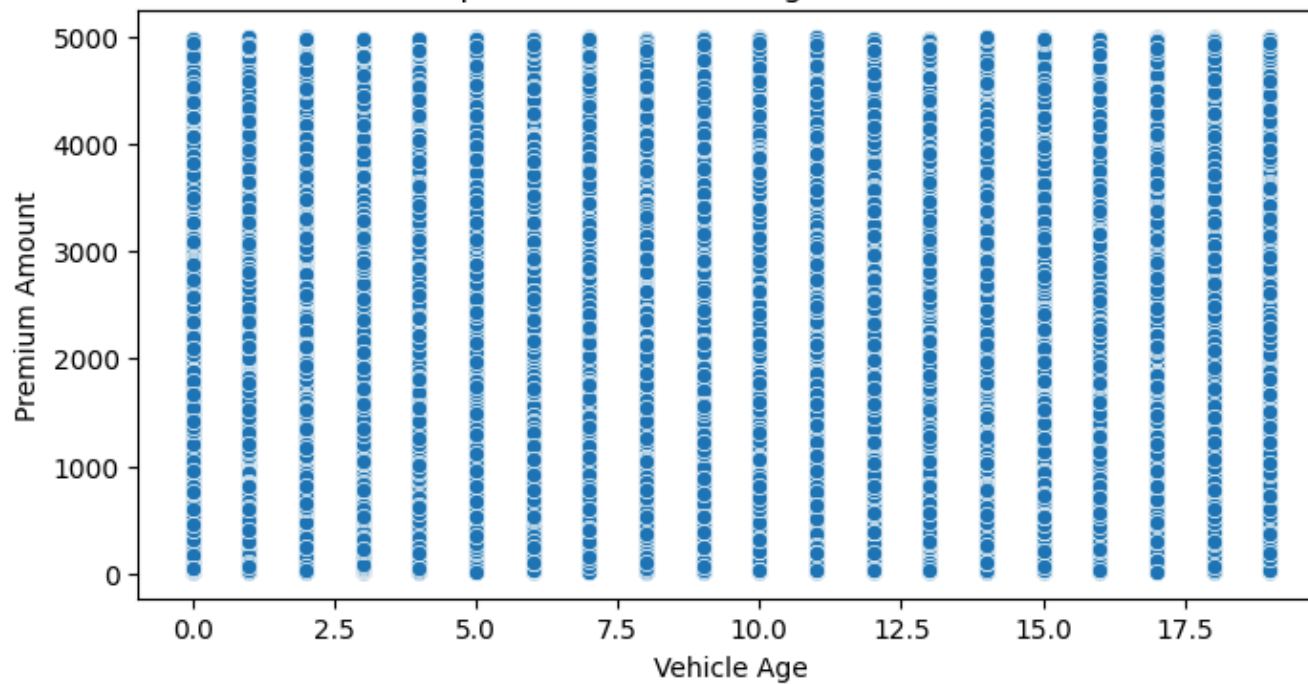


Relationship between Previous Claims and Premium Amount

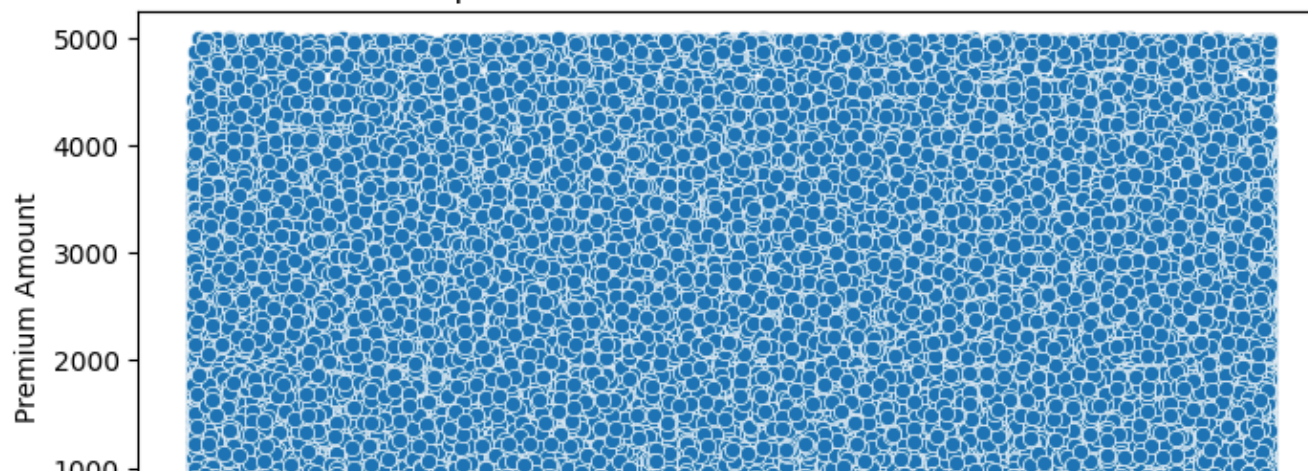


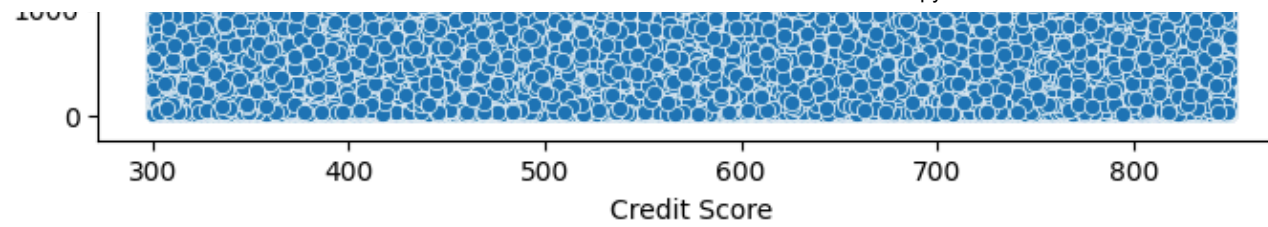


Relationship between Vehicle Age and Premium Amount

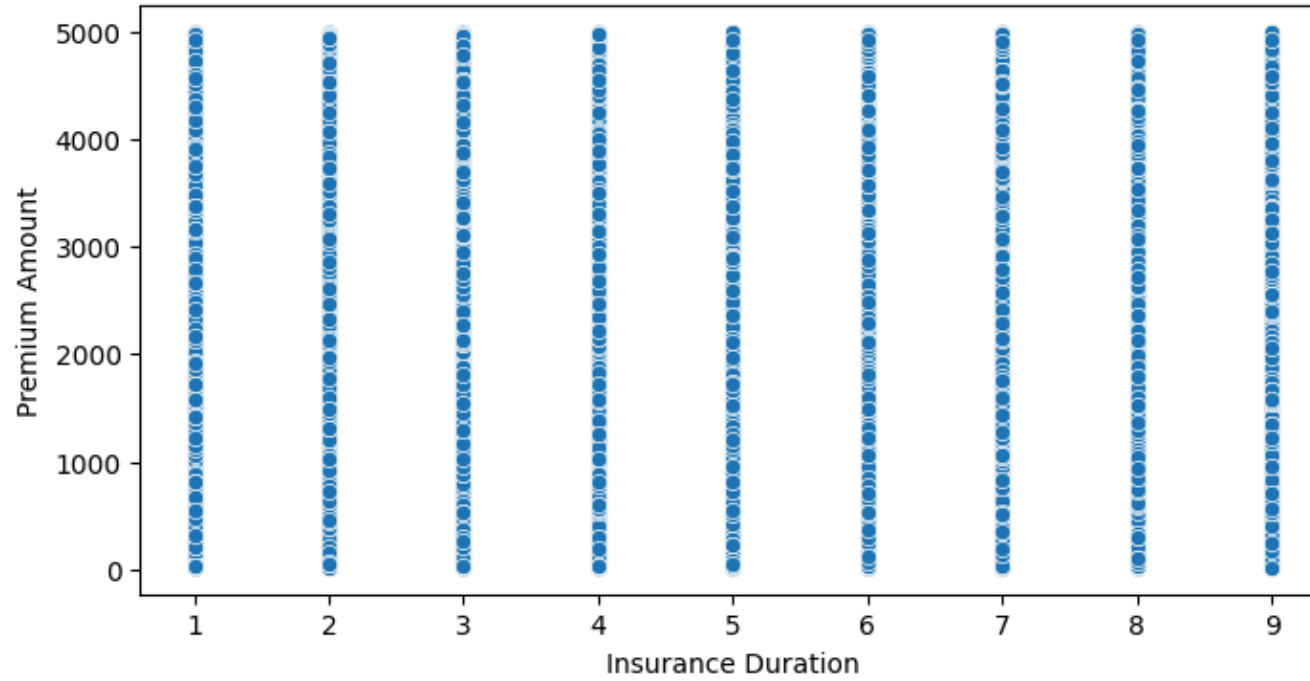


Relationship between Credit Score and Premium Amount



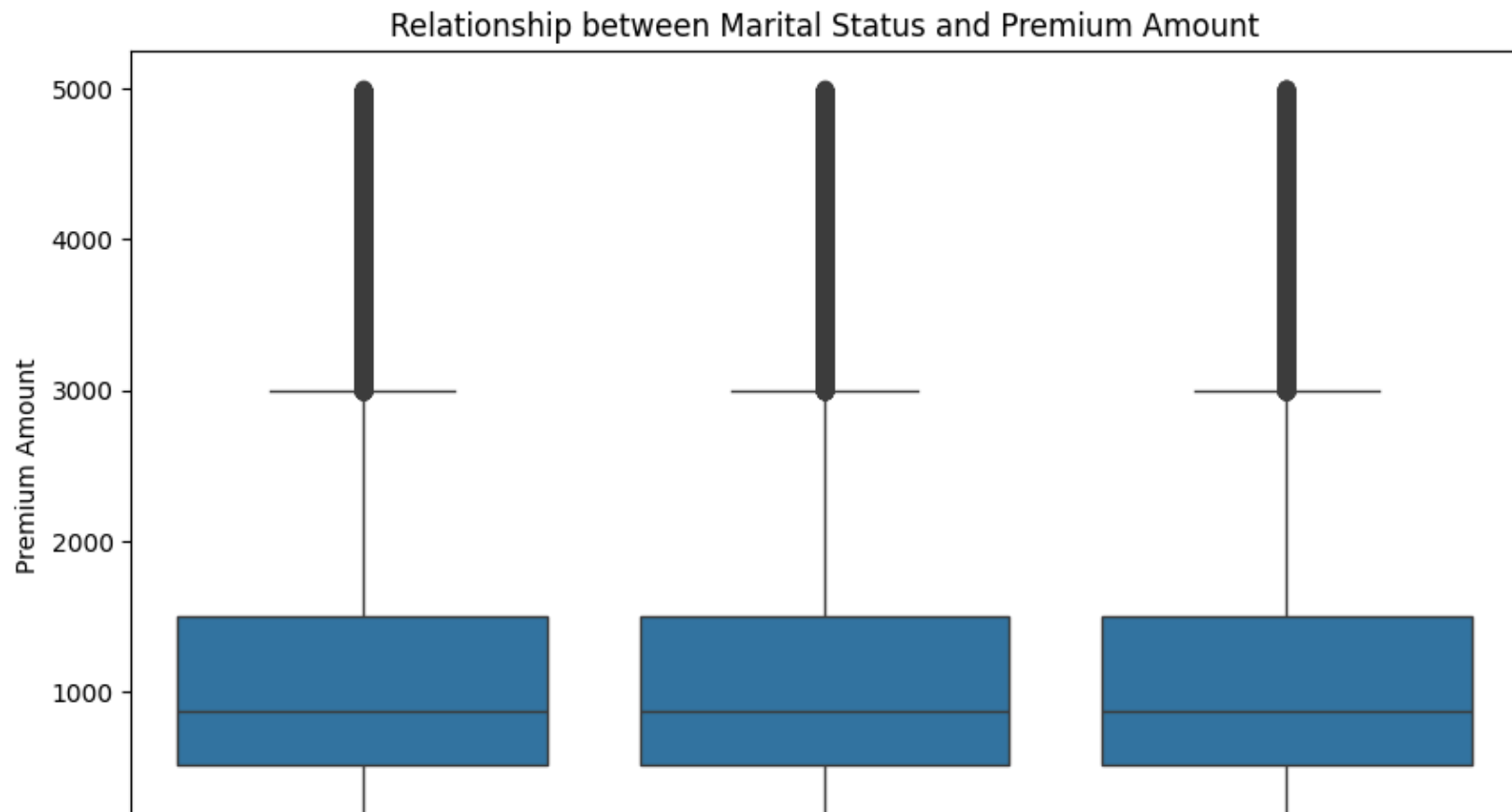
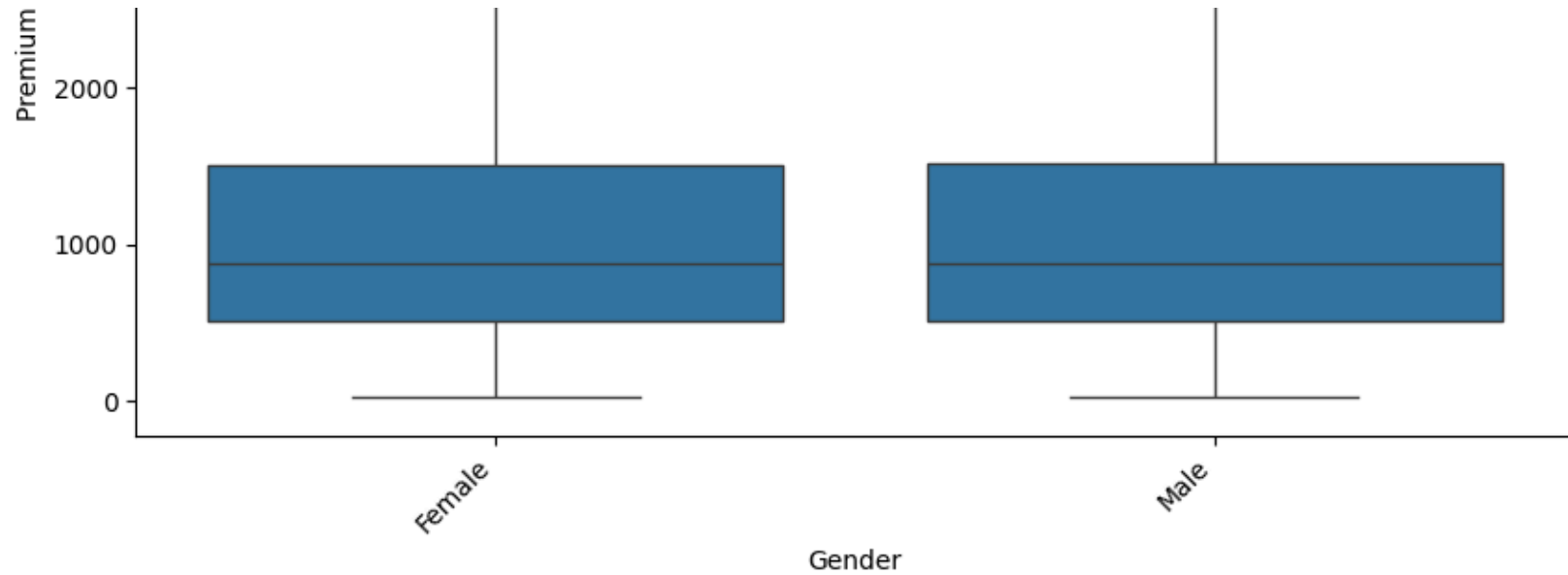


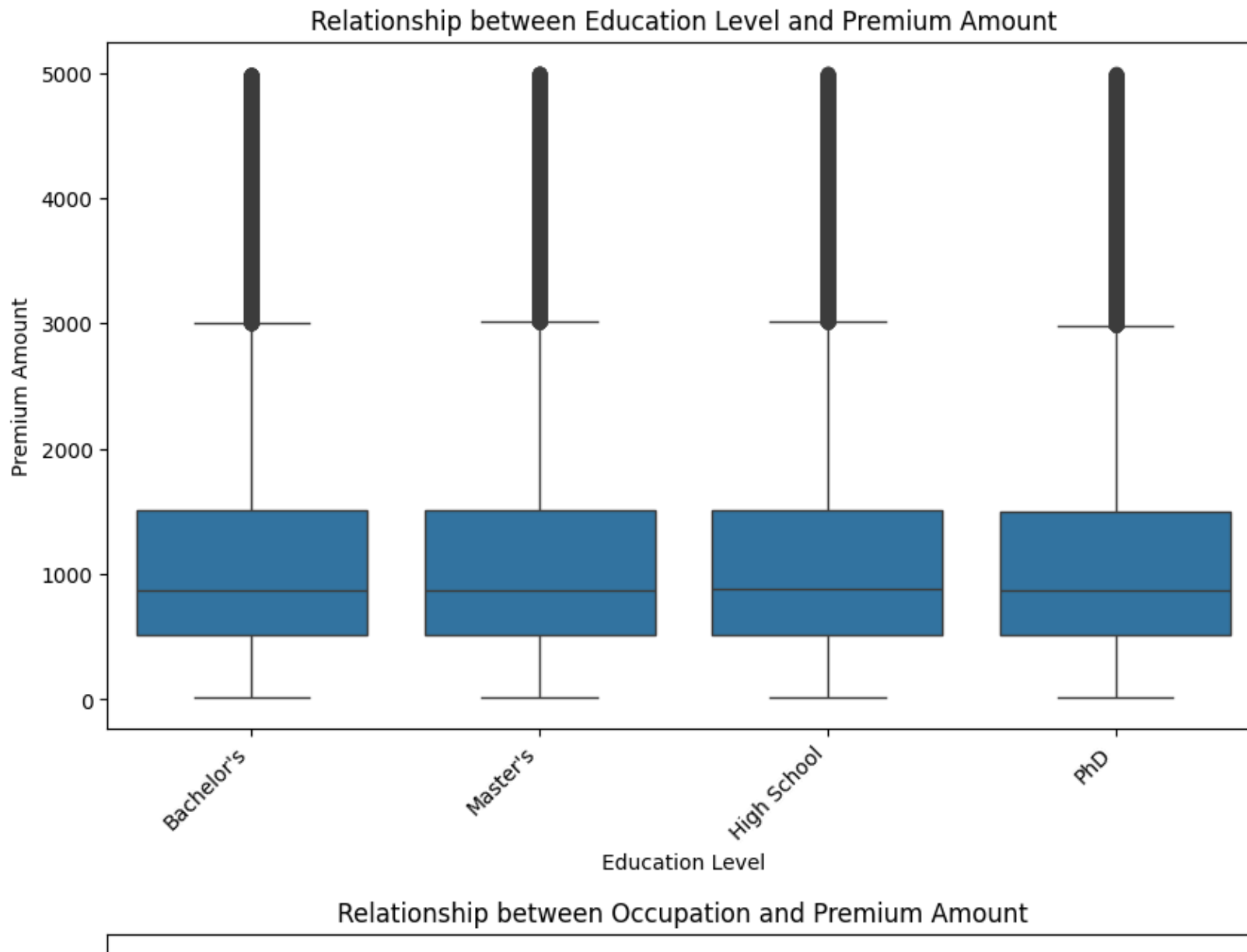
Relationship between Insurance Duration and Premium Amount

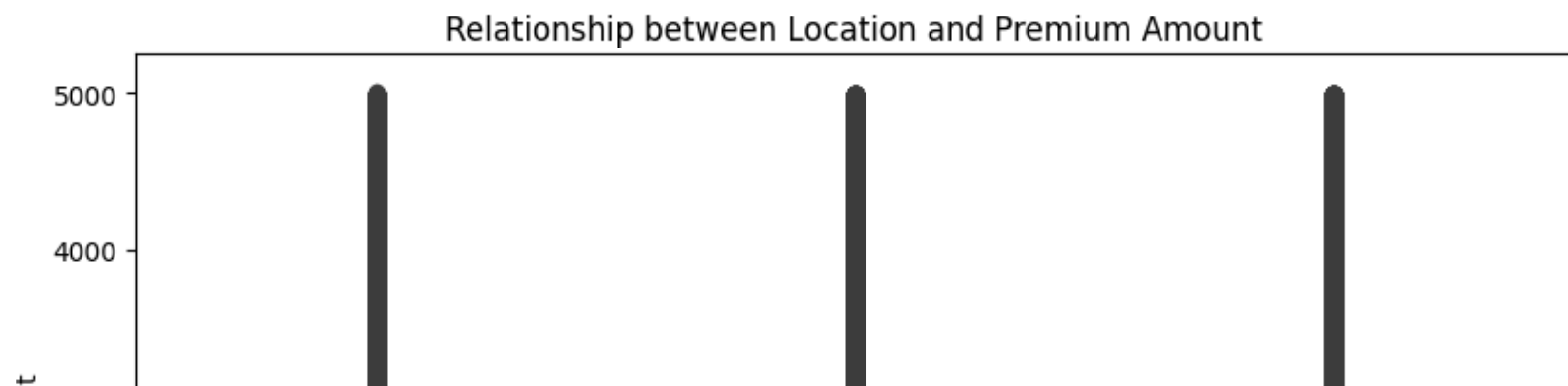
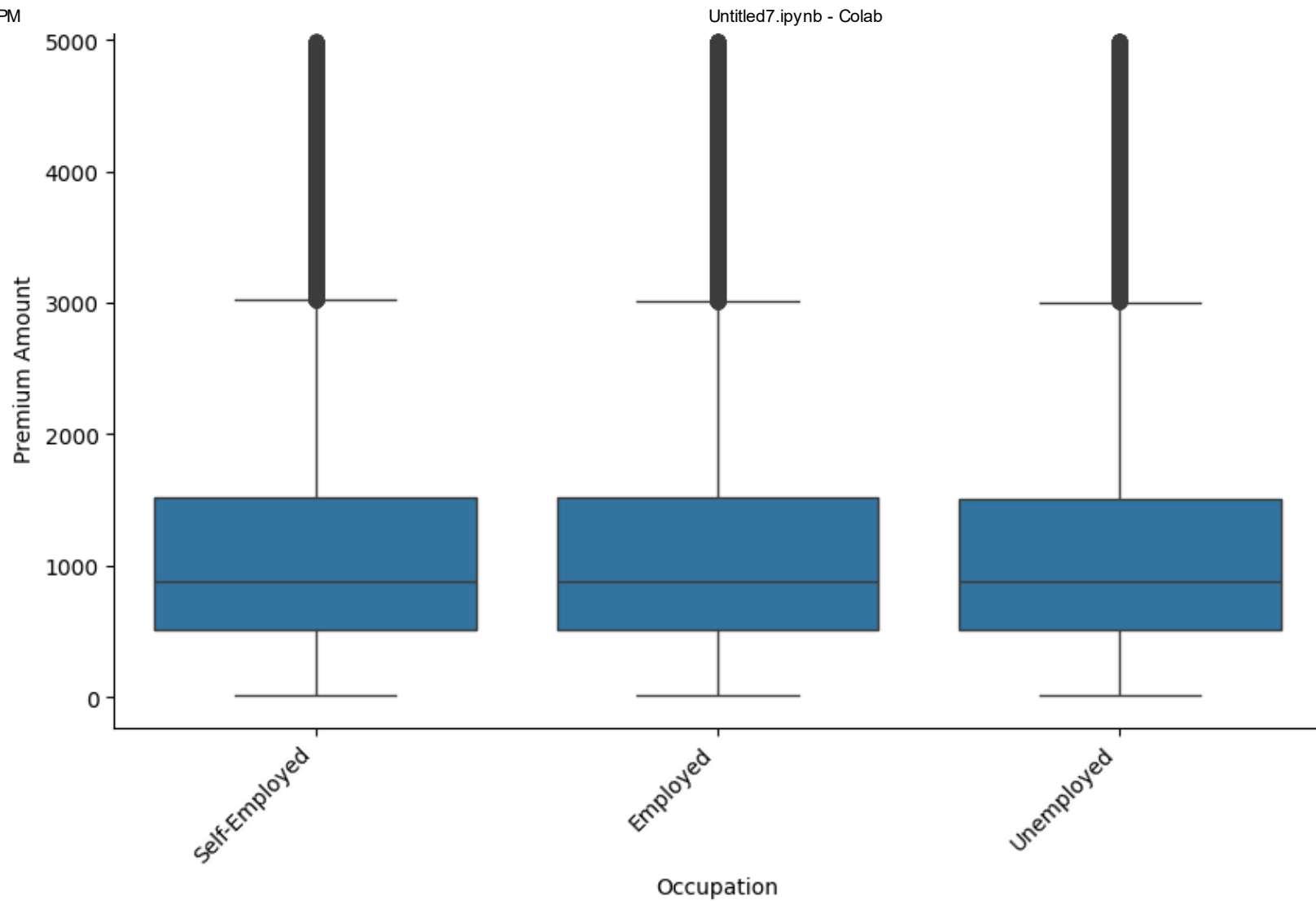


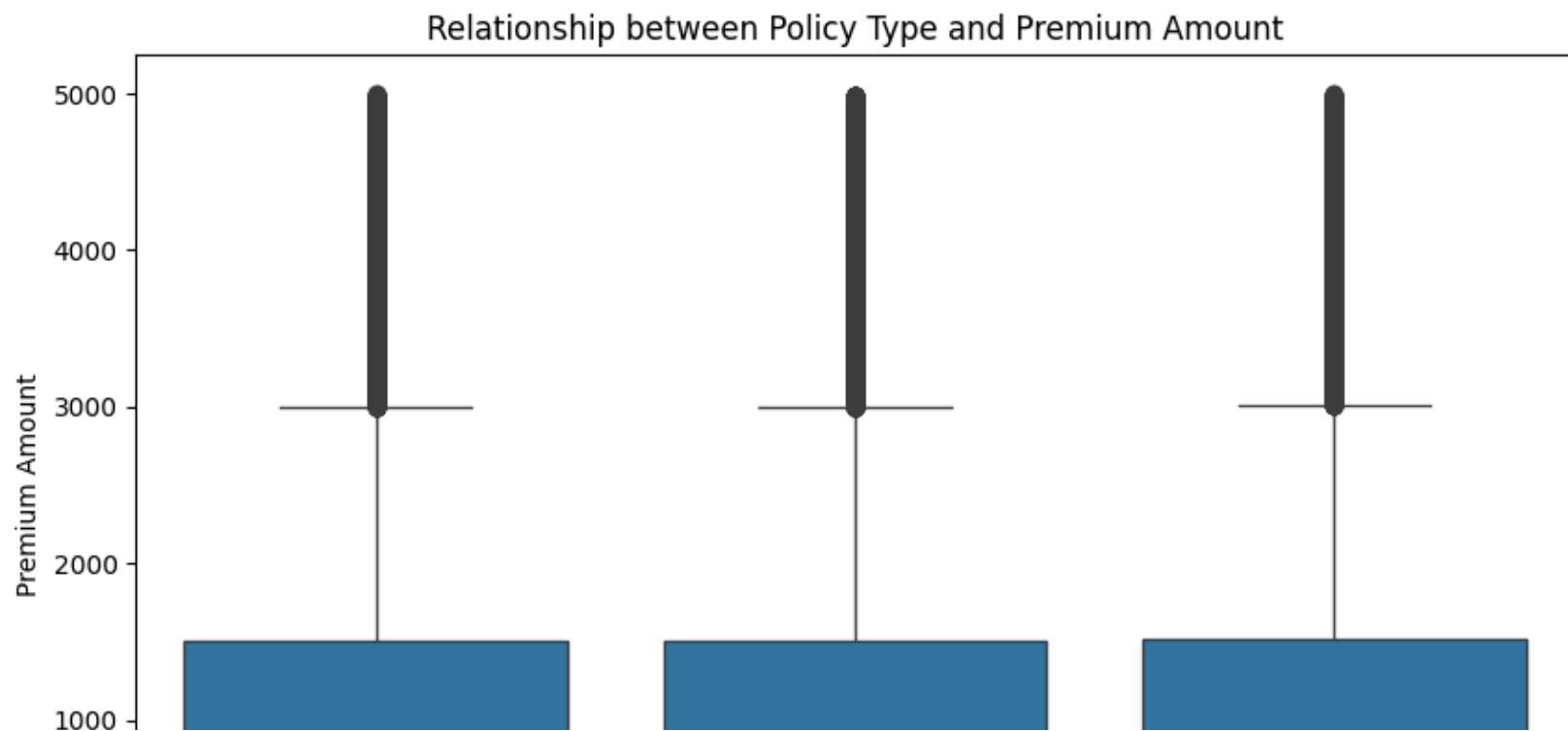
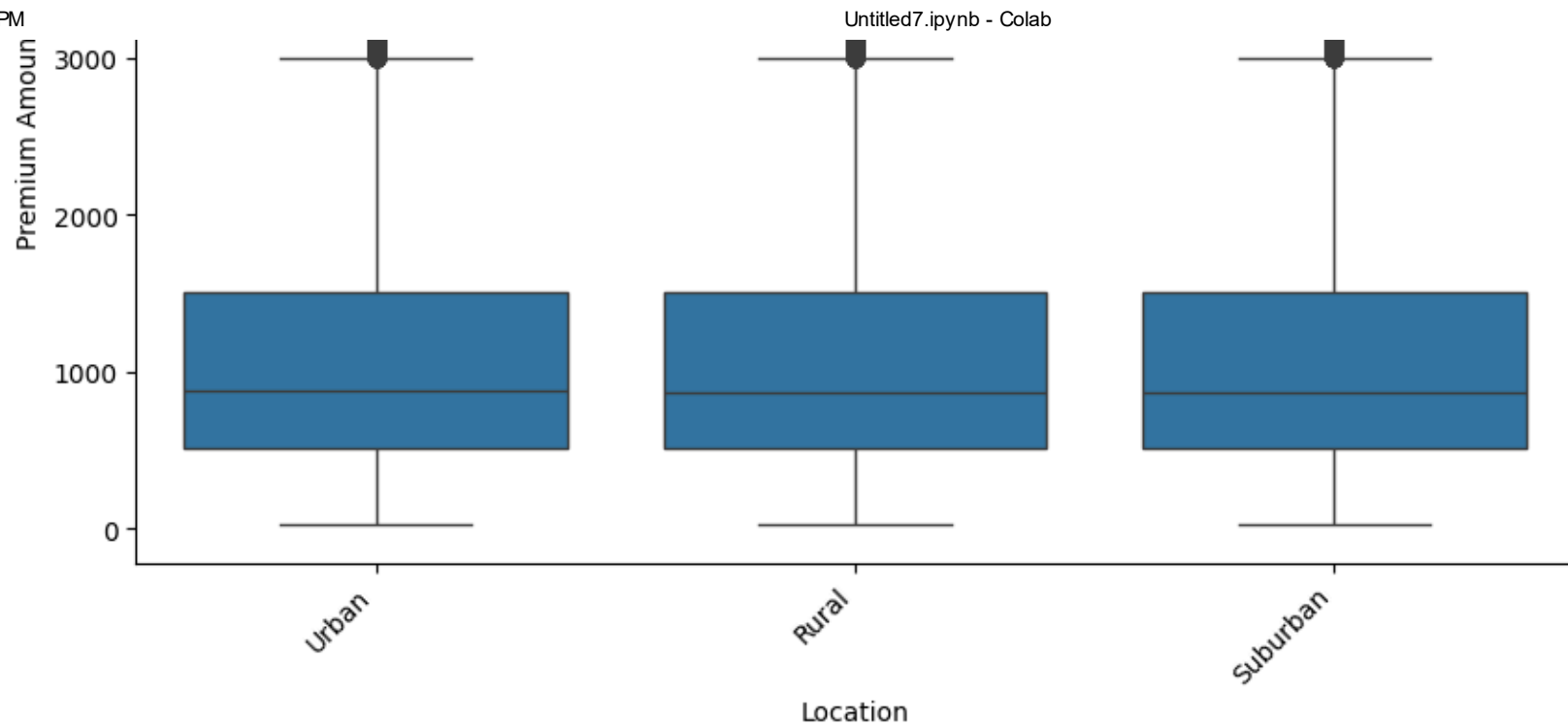
Relationship between Gender and Premium Amount

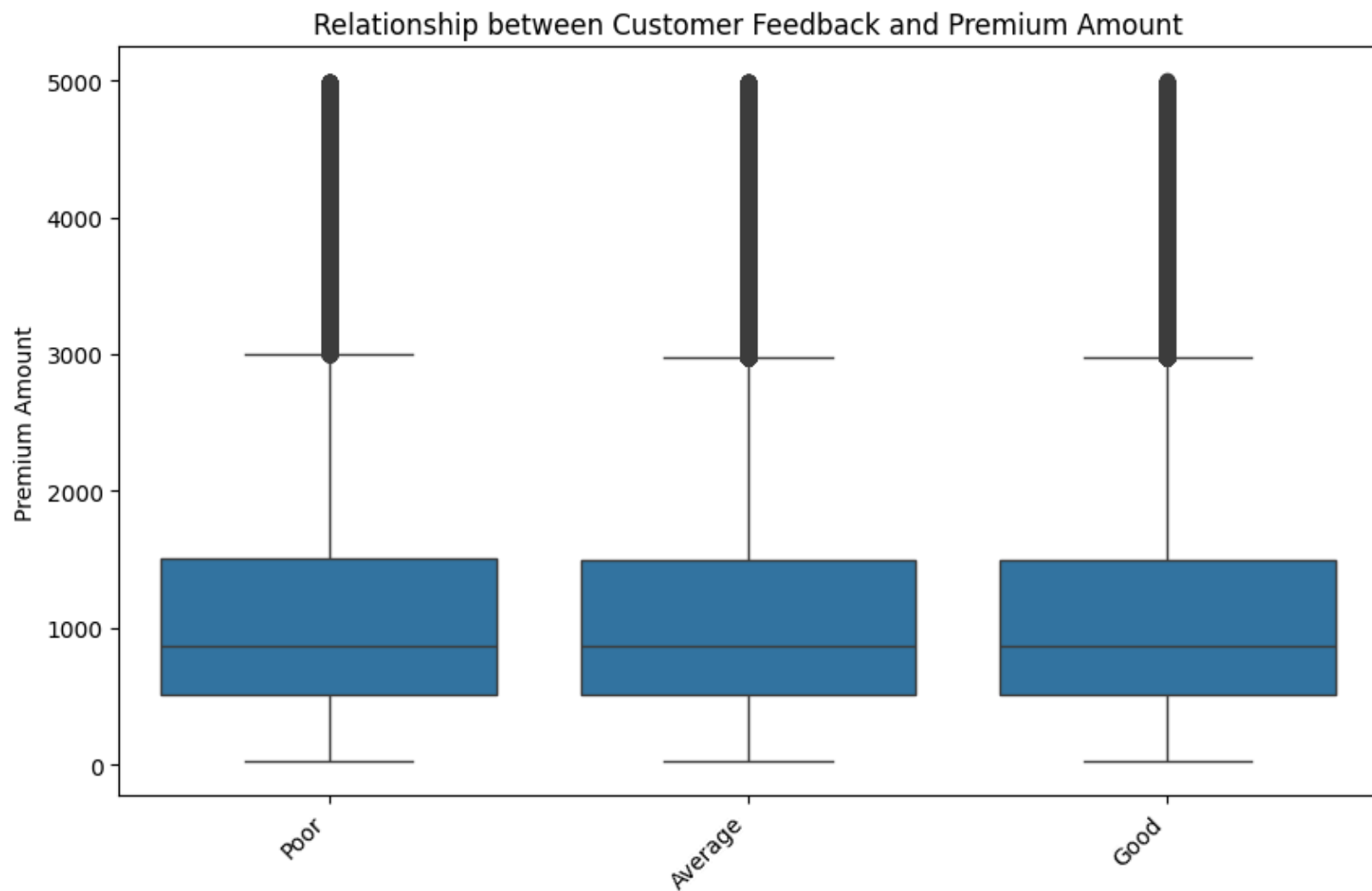
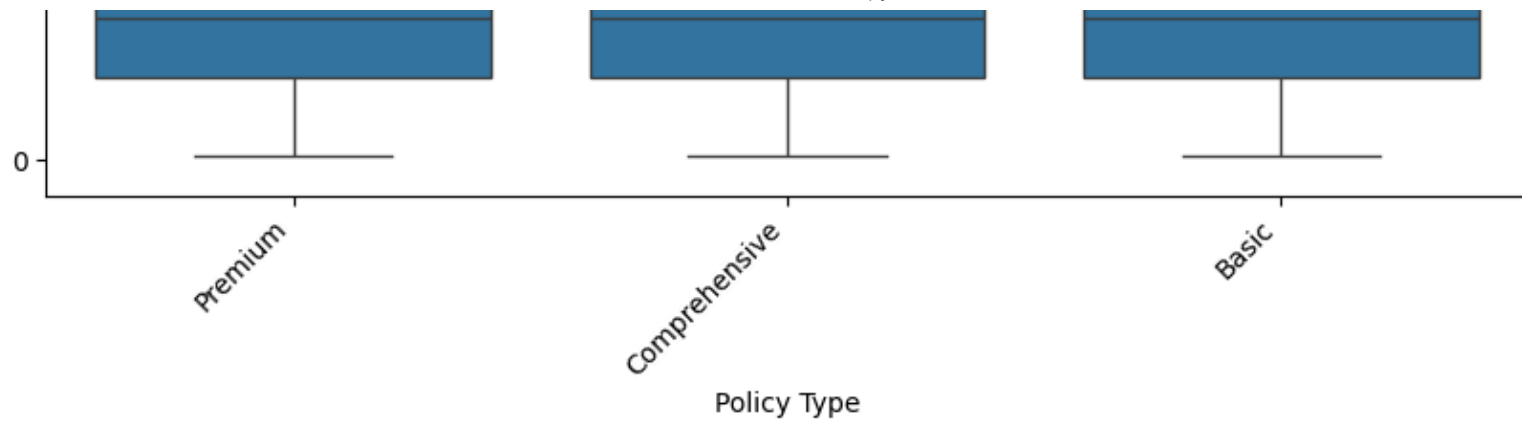




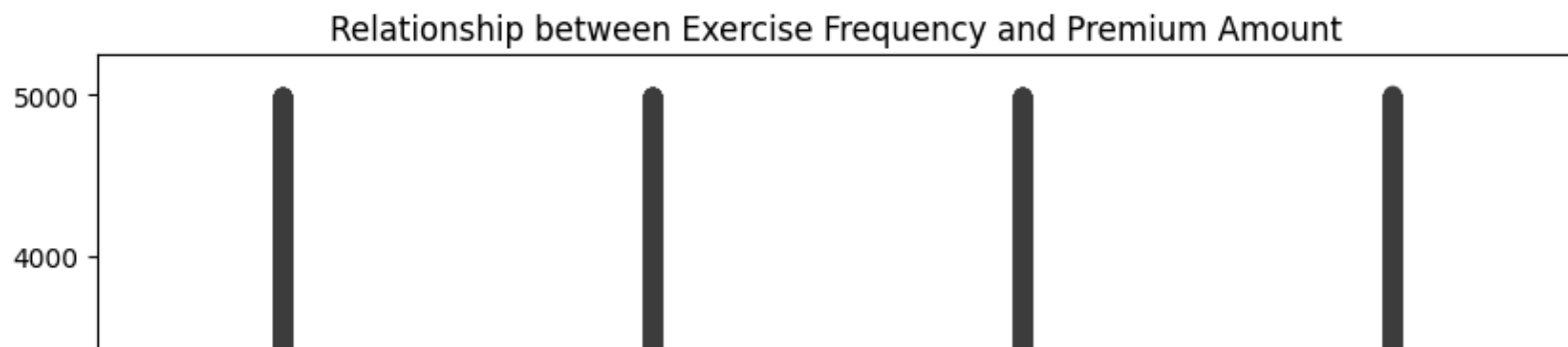
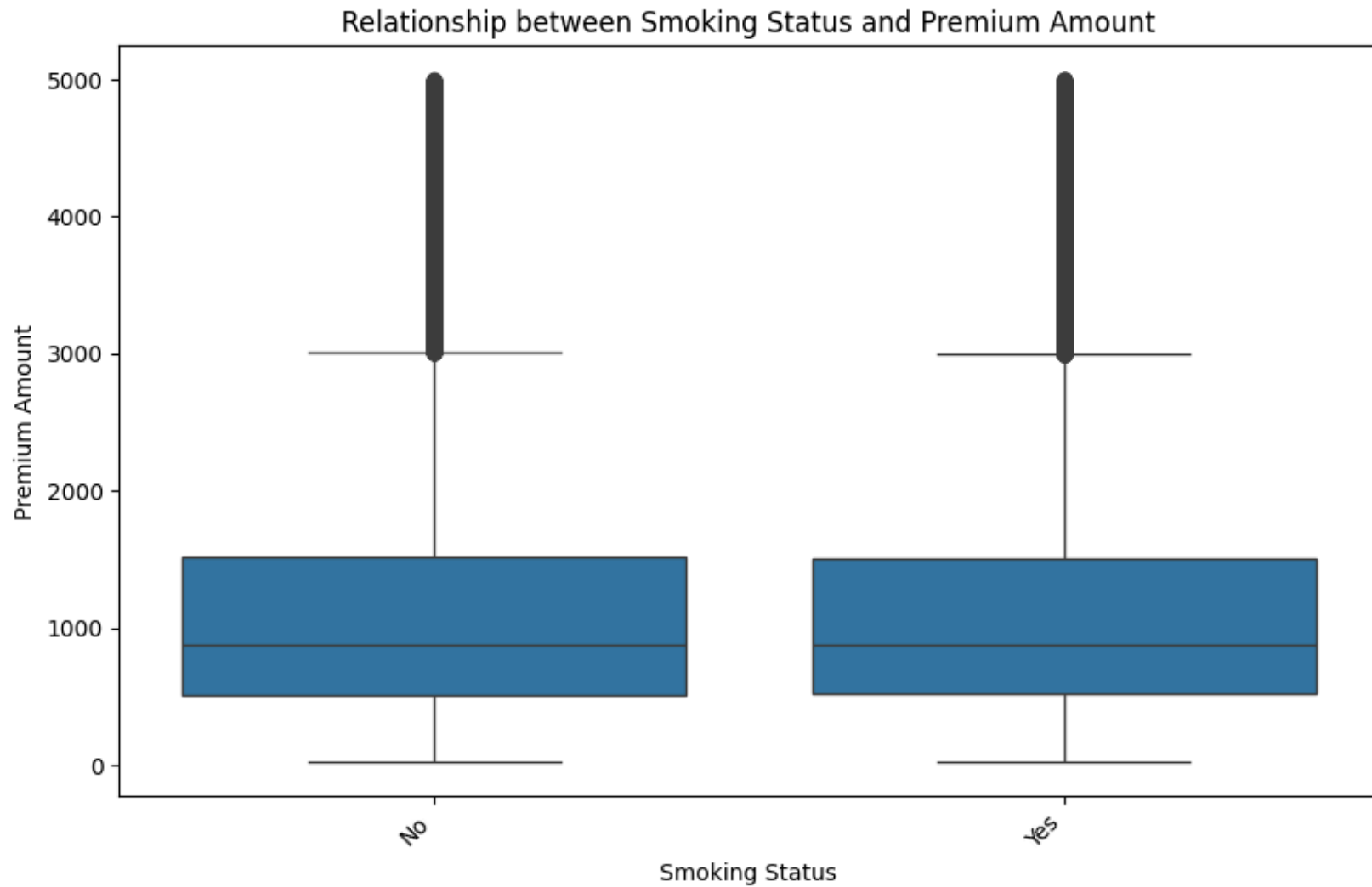


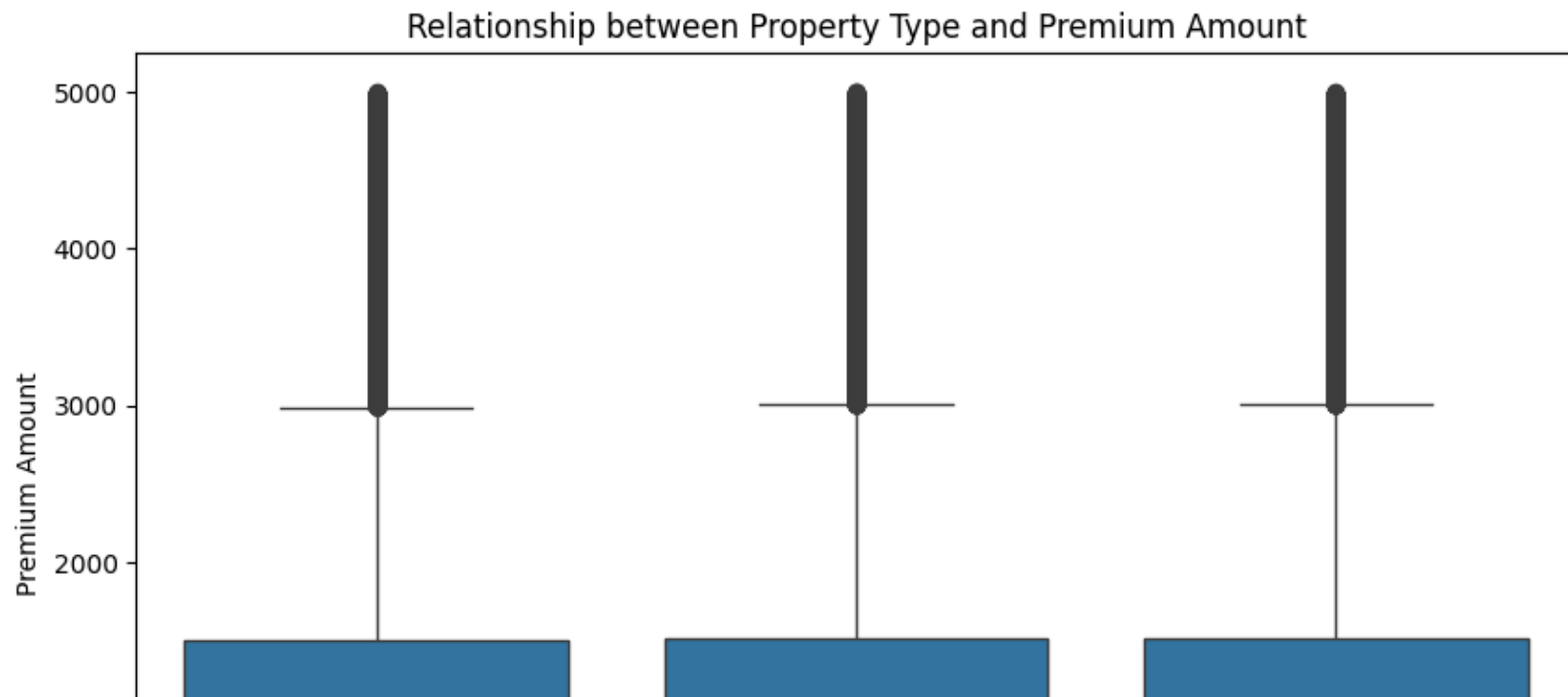
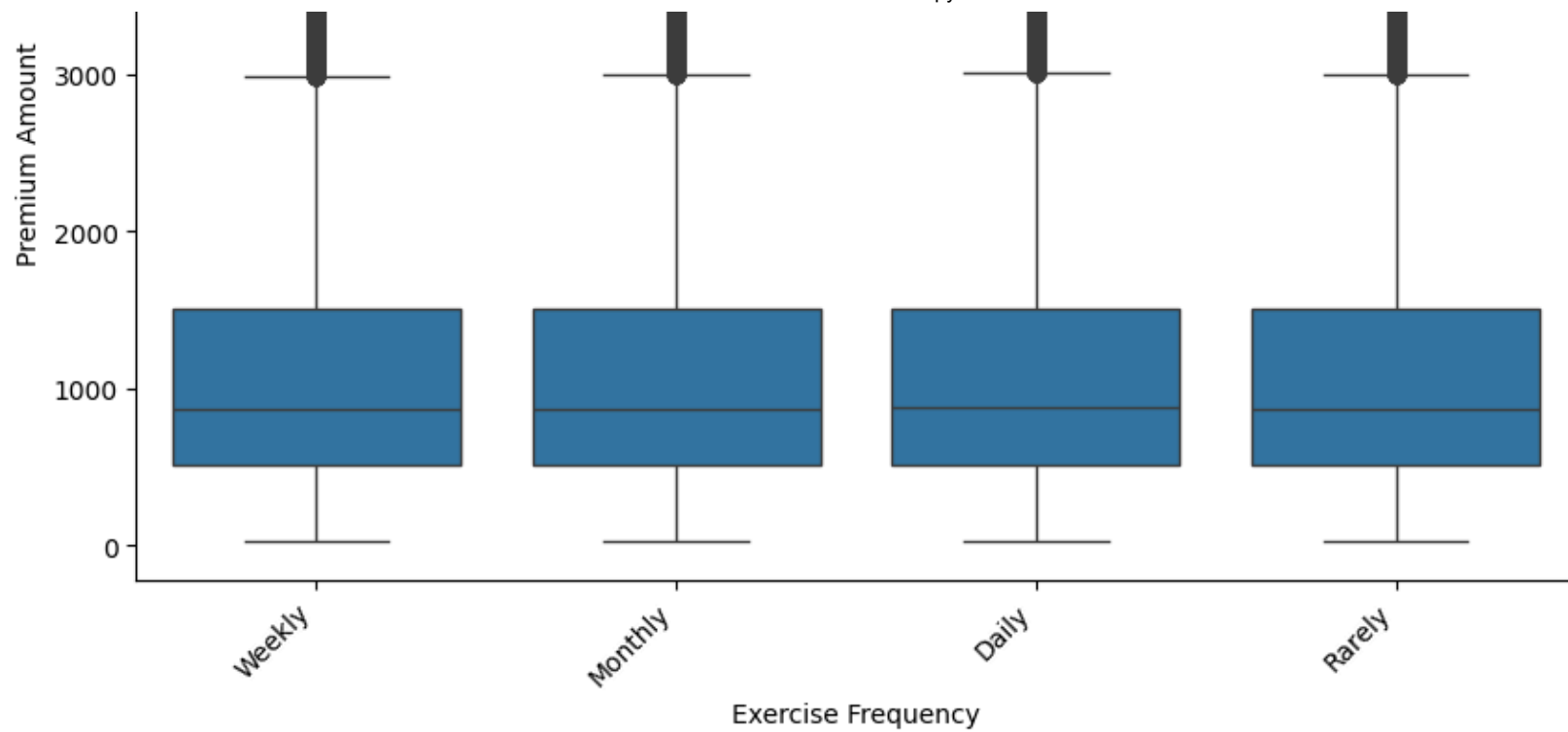


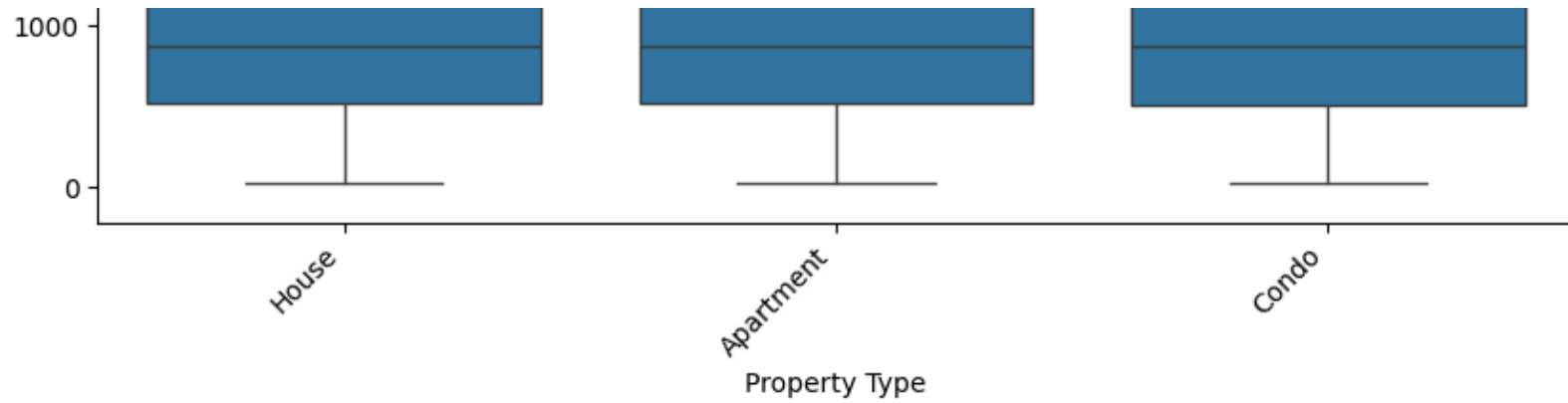




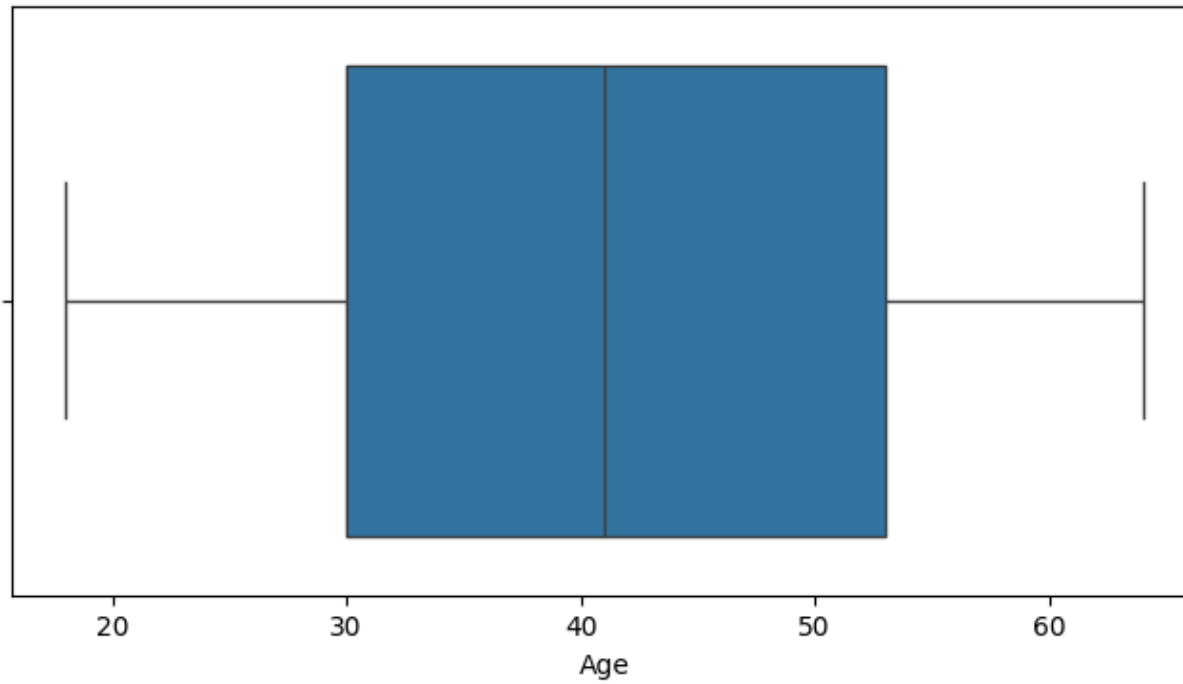
Customer Feedback





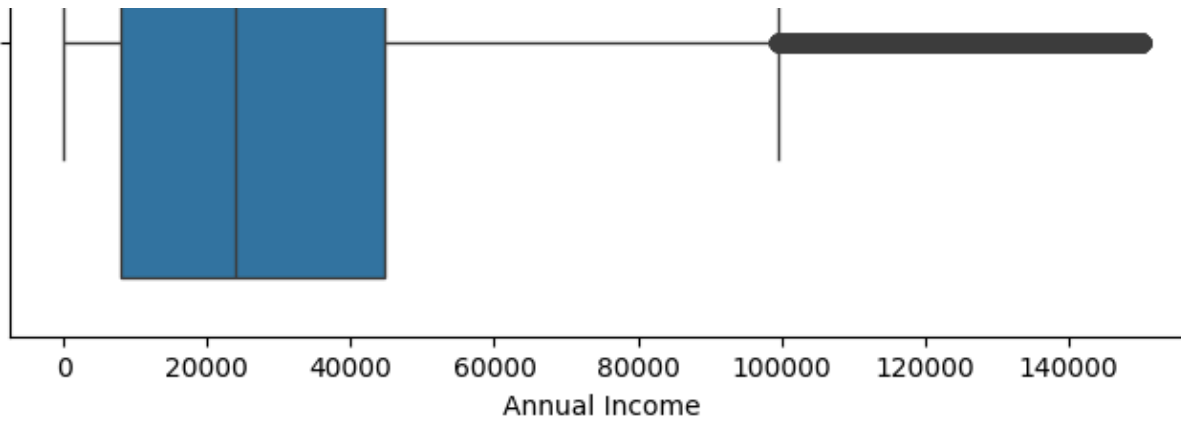


Boxplot of Age

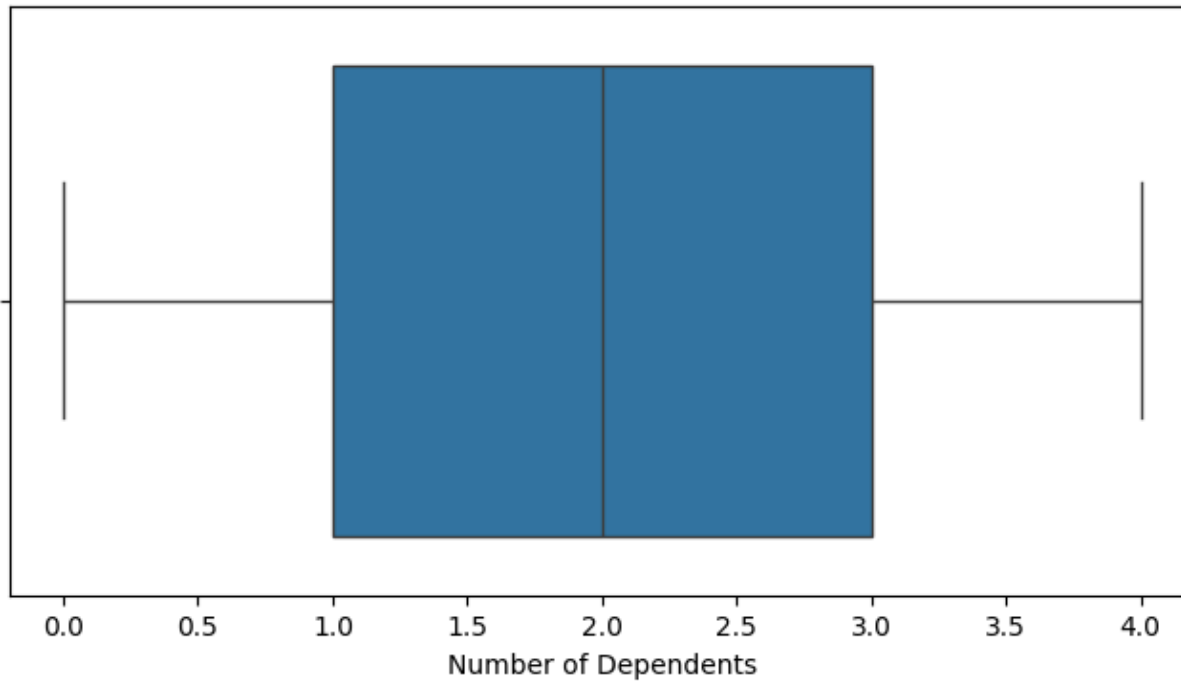


Boxplot of Annual Income

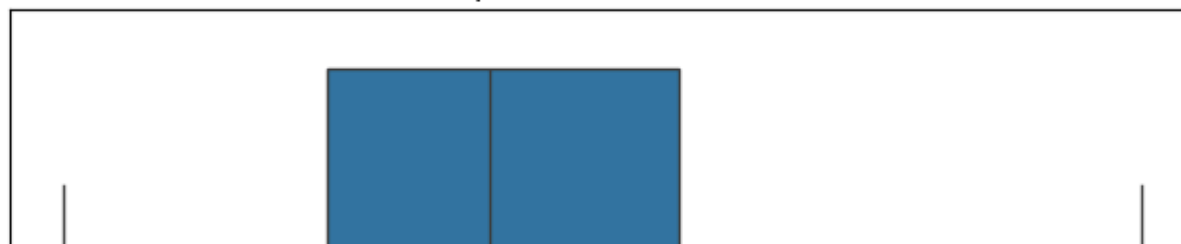


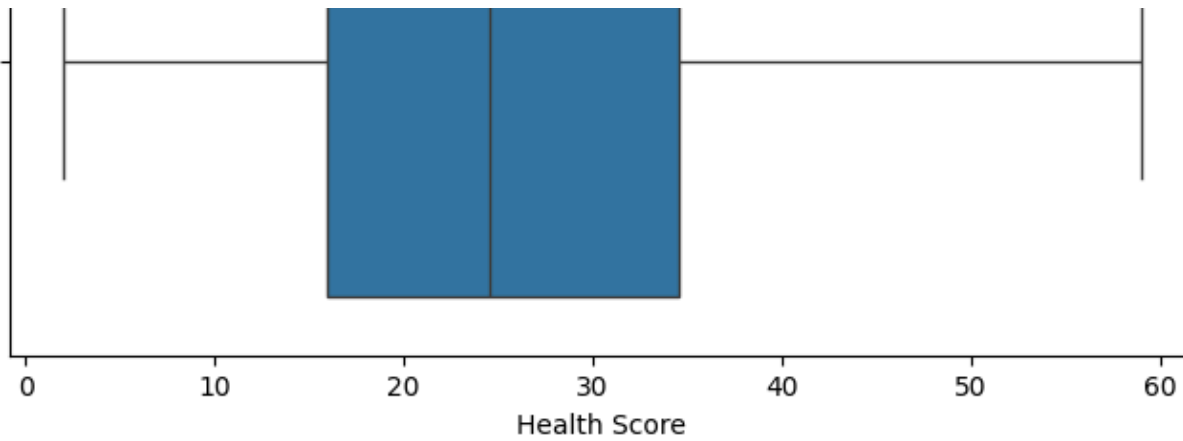


Boxplot of Number of Dependents

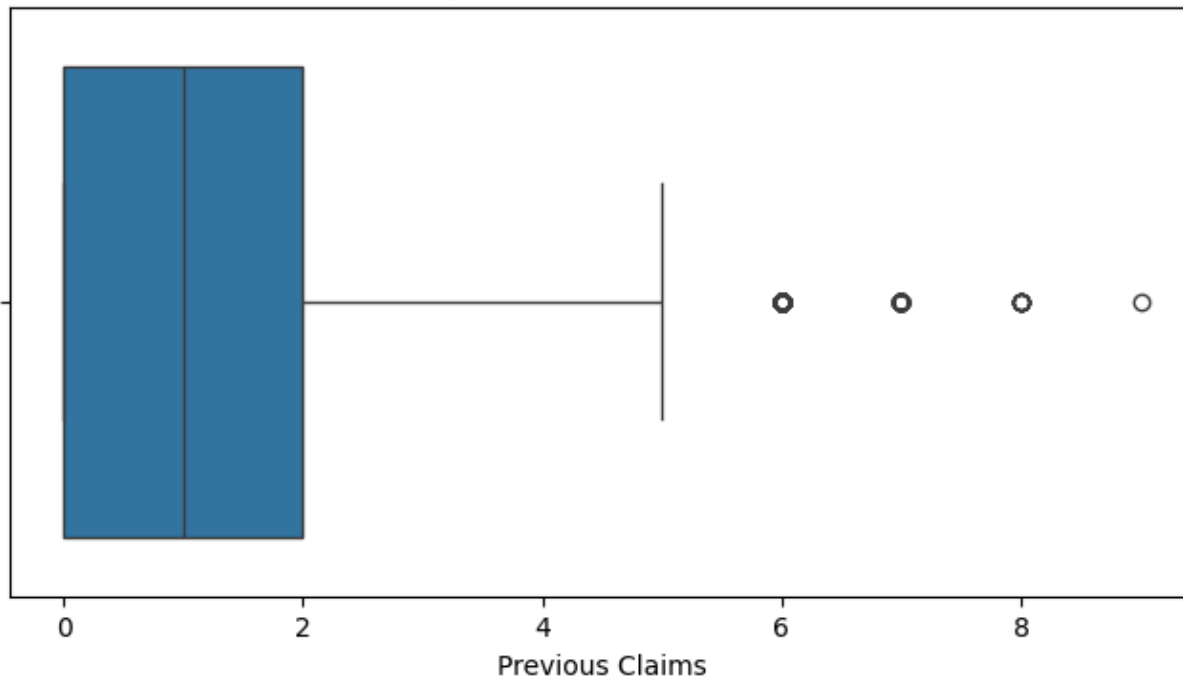


Boxplot of Health Score

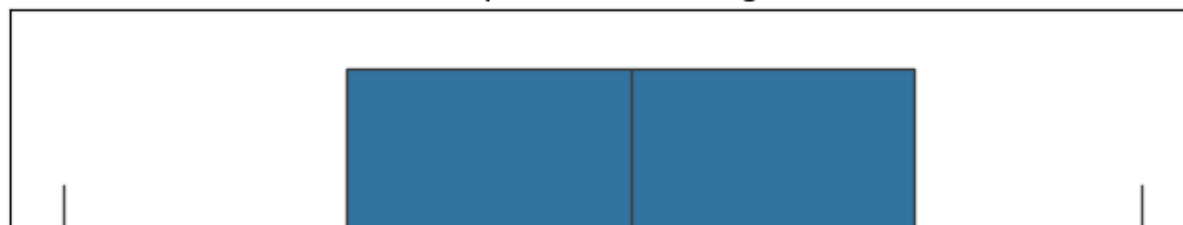


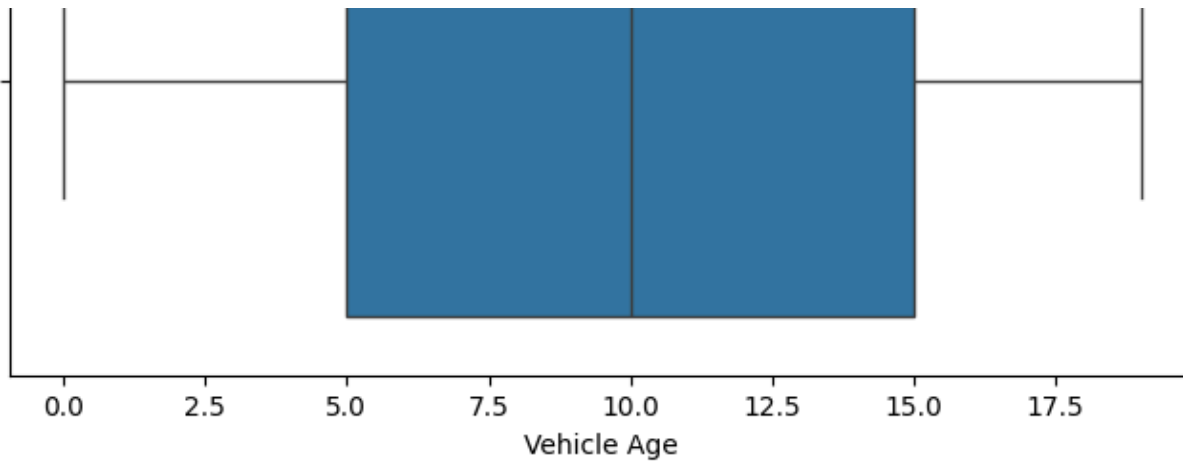


Boxplot of Previous Claims

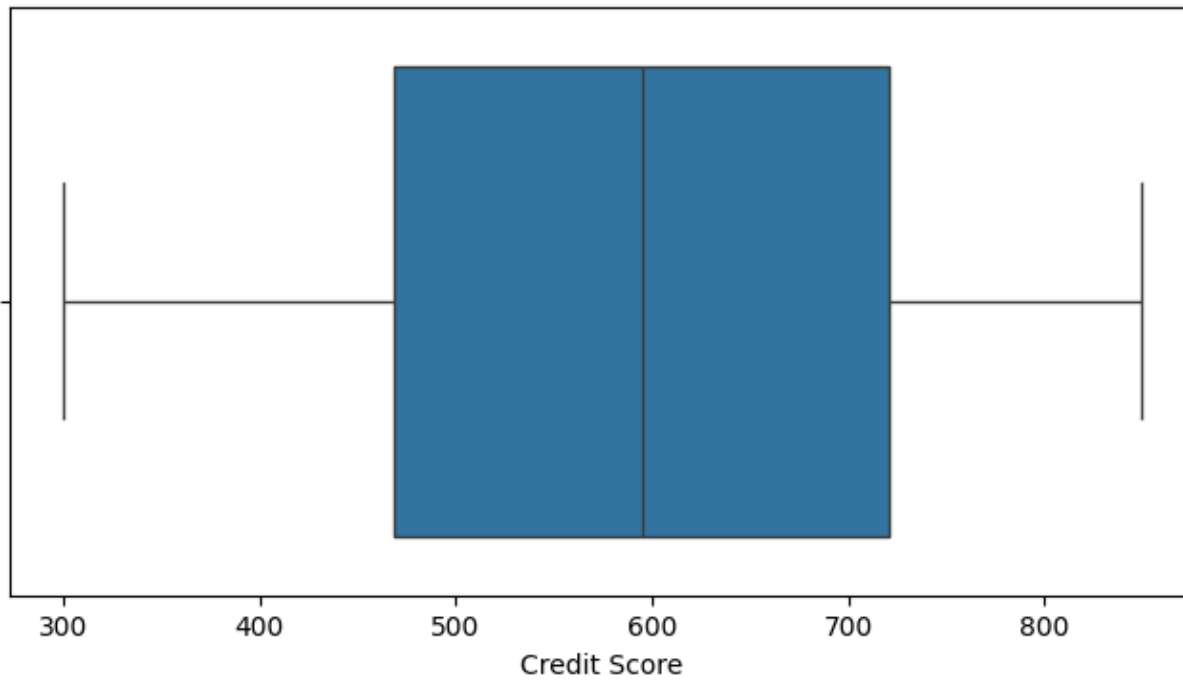


Boxplot of Vehicle Age

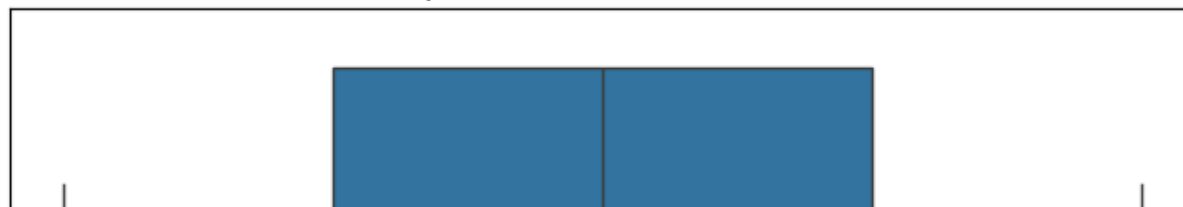


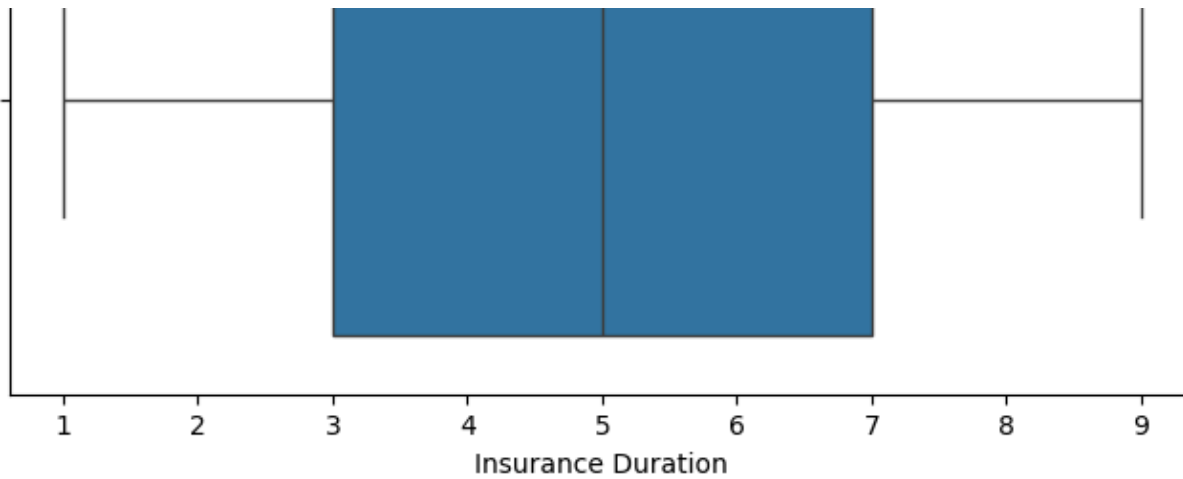


Boxplot of Credit Score



Boxplot of Insurance Duration





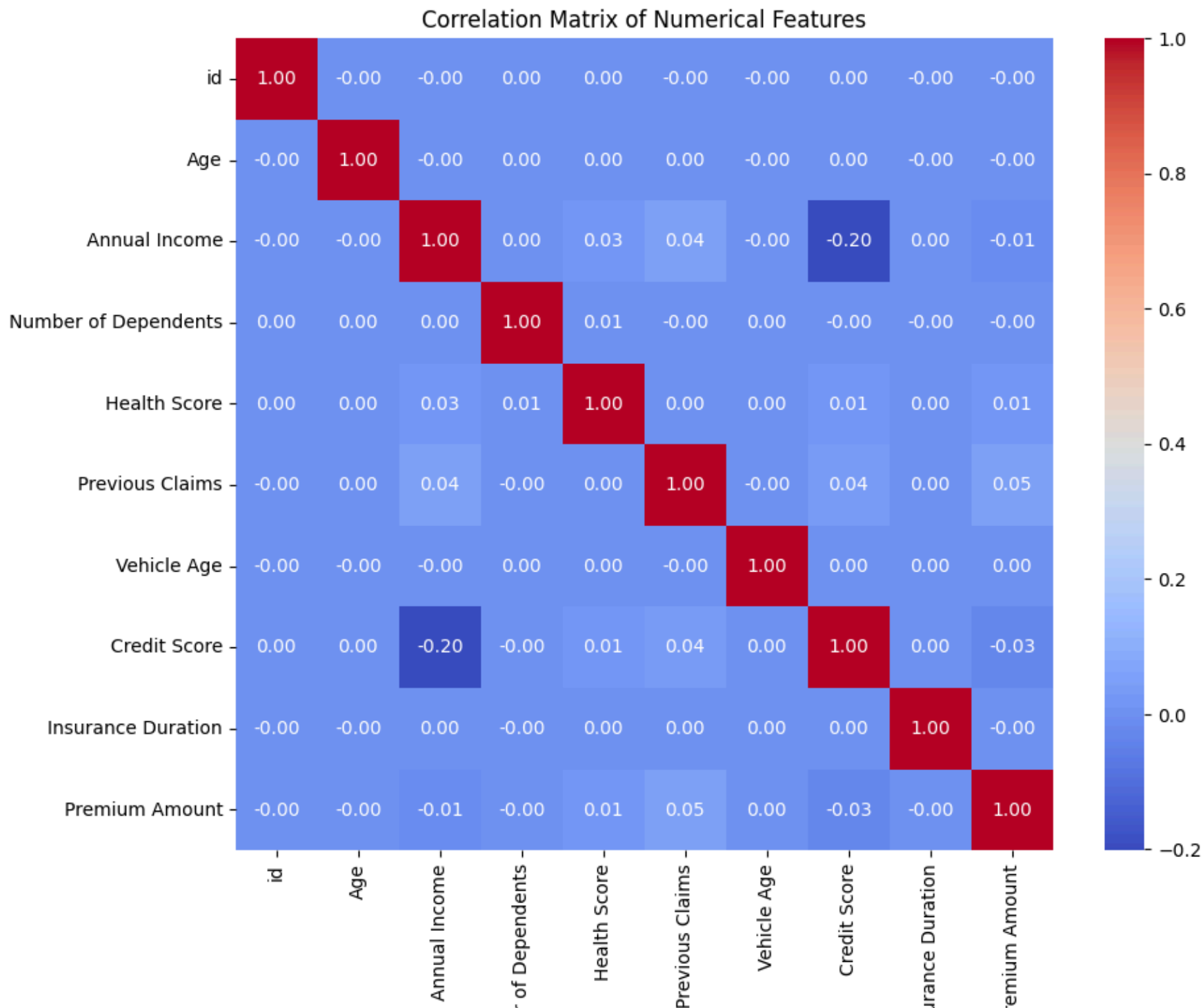
Next steps:

[Generate code with missing_df](#)[View recommended plots](#)[New interactive sheet](#)

Reasoning: The previous code block generated a large number of plots, many of which are difficult to interpret or irrelevant. I will refine the code to focus on more informative visualizations and statistical summaries related to the target variable "Premium Amount" and key features, addressing missing values and outliers effectively.

```
1 import matplotlib.pyplot as plt
2 import seaborn as sns
3
4 # Focus on relationships with 'Premium Amount'
5 numerical_features = df_train.select_dtypes(include=['number']).columns
6 categorical_features = df_train.select_dtypes(exclude=['number']).columns
7
8 # Correlation matrix (heatmap) for numerical features
9 plt.figure(figsize=(10, 8))
10 correlation_matrix = df_train[numerical_features].corr()
11 sns.heatmap(correlation_matrix, annot=True, cmap='coolwarm', fmt=".2f")
12 plt.title('Correlation Matrix of Numerical Features')
13 plt.show()
14
15
16 # Pairplot for selected numerical features and 'Premium Amount'
17 selected_numerical_features = ['Age', 'Annual Income', 'Health Score', 'Credit Score',
18 sns.pairplot(df_train[selected_numerical_features], diag_kind='kde')
19 plt.suptitle('Pairplot of Selected Numerical Features', y=1.02)
20 plt.show()
21
22 # Box plots for categorical features vs. 'Premium Amount'
23 for col in categorical_features:
24     if col not in ['Policy Start Date', 'id']:
25         plt.figure(figsize=(10, 6))
26         sns.boxplot(x=col, y='Premium Amount', data=df_train)
```

```
27     plt.title(f'{col} vs. Premium Amount')
28     plt.xticks(rotation=45, ha='right')
29     plt.tight_layout()
30     plt.show()
31
32 # Missing Value Analysis for crucial features
33 missing_cols = ['Age', 'Annual Income', 'Health Score', 'Credit Score', 'Previous Claim']
34 for col in missing_cols:
35     plt.figure(figsize=(8, 4))
36     sns.histplot(df_train[col], kde=True)
37     plt.title(f"Distribution of {col} with Missing Values")
38     plt.show()
39
40     plt.figure(figsize=(8, 4))
41     sns.boxplot(y=col, data=df_train)
42     plt.title(f"Box Plot for {col} with Missing Values")
43     plt.show()
```

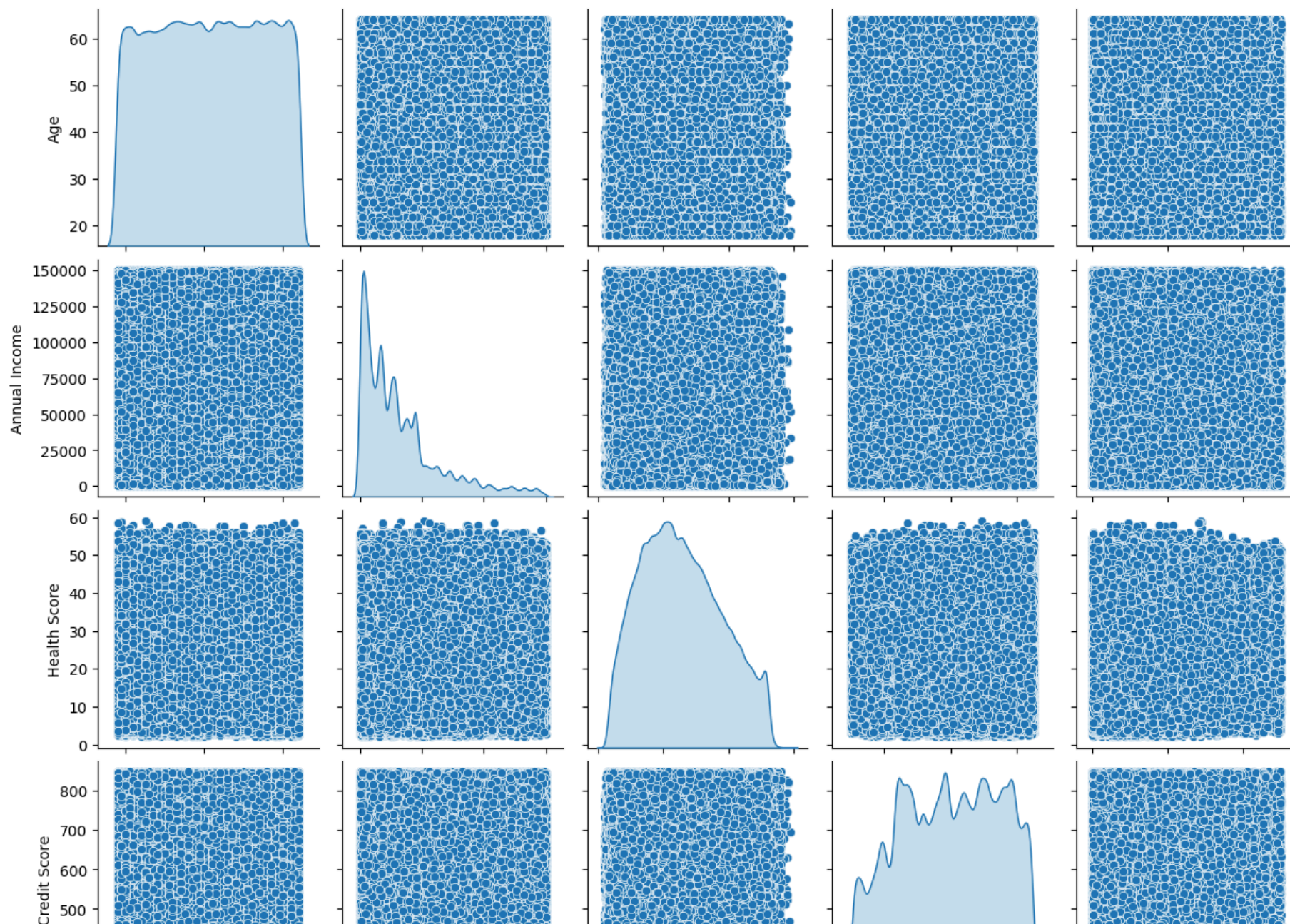


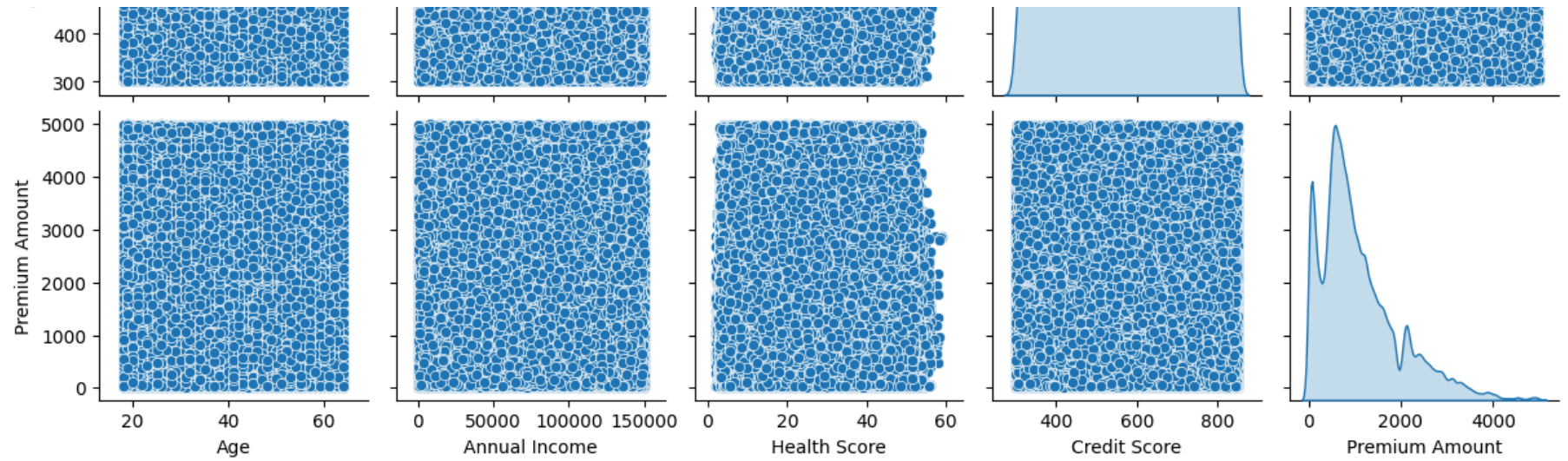
Number

Insu

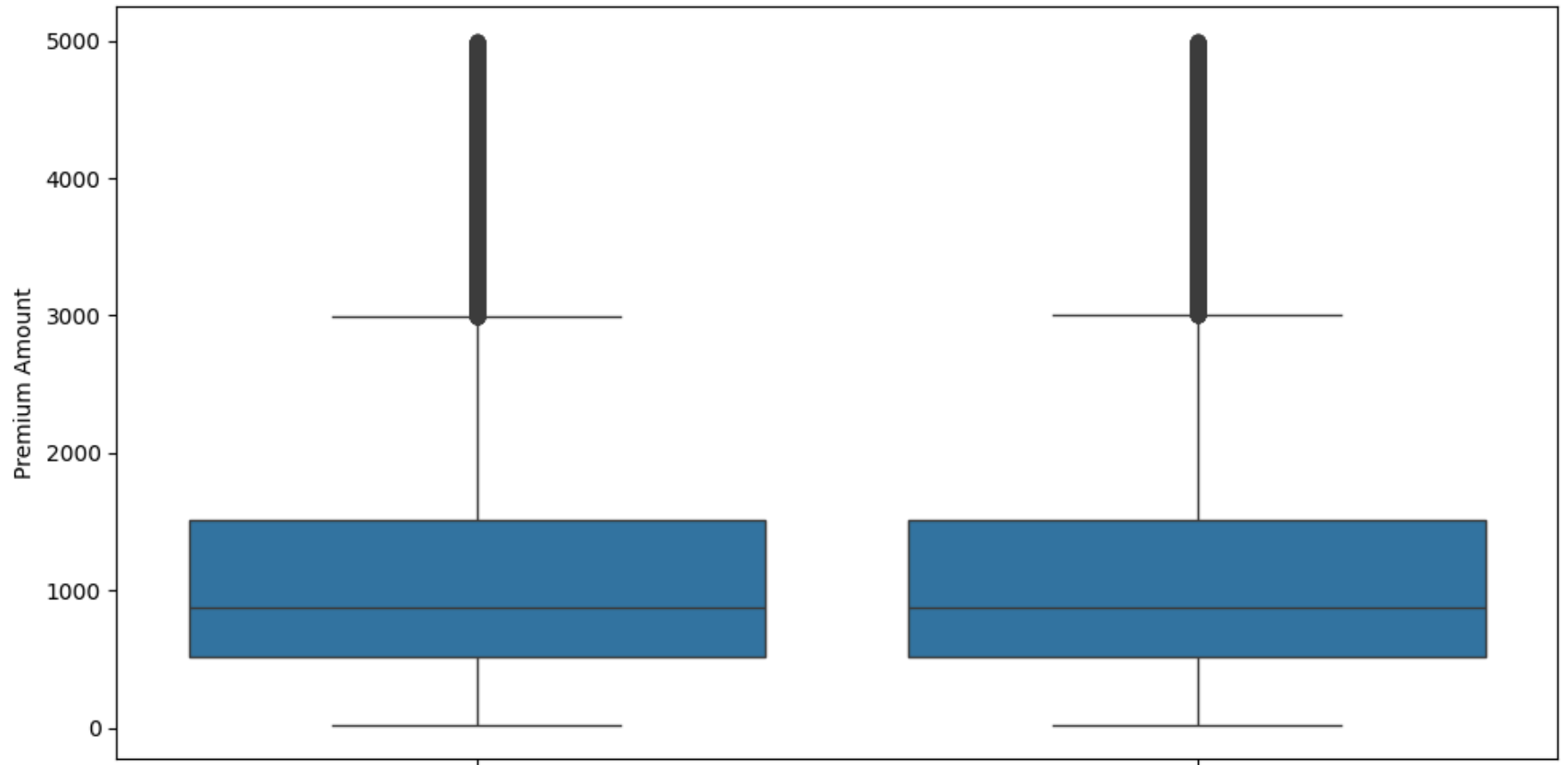
PI

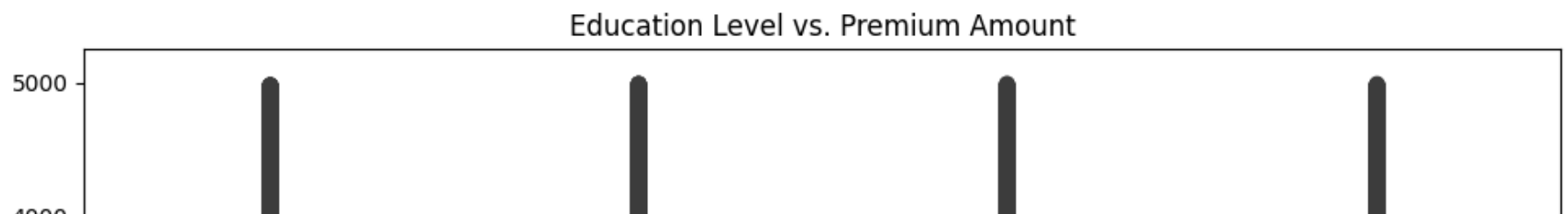
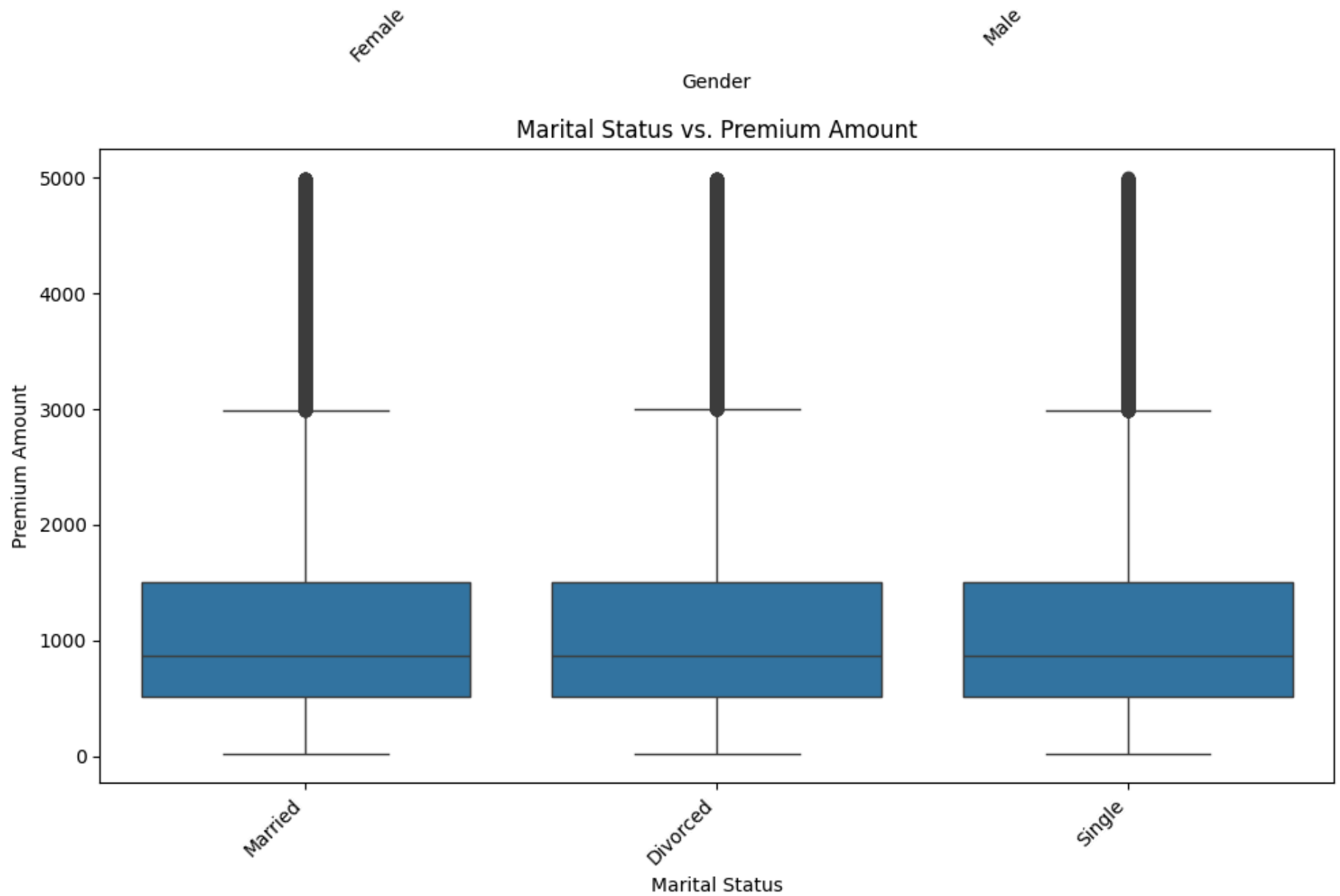
Pairplot of Selected Numerical Features

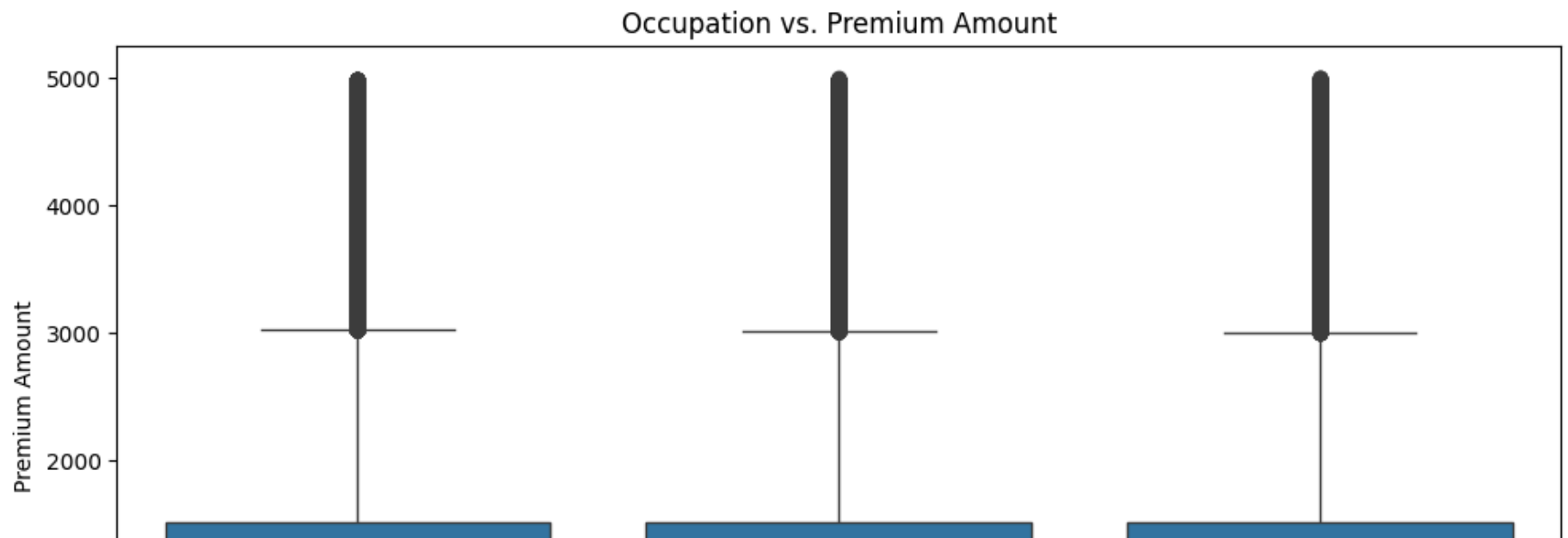
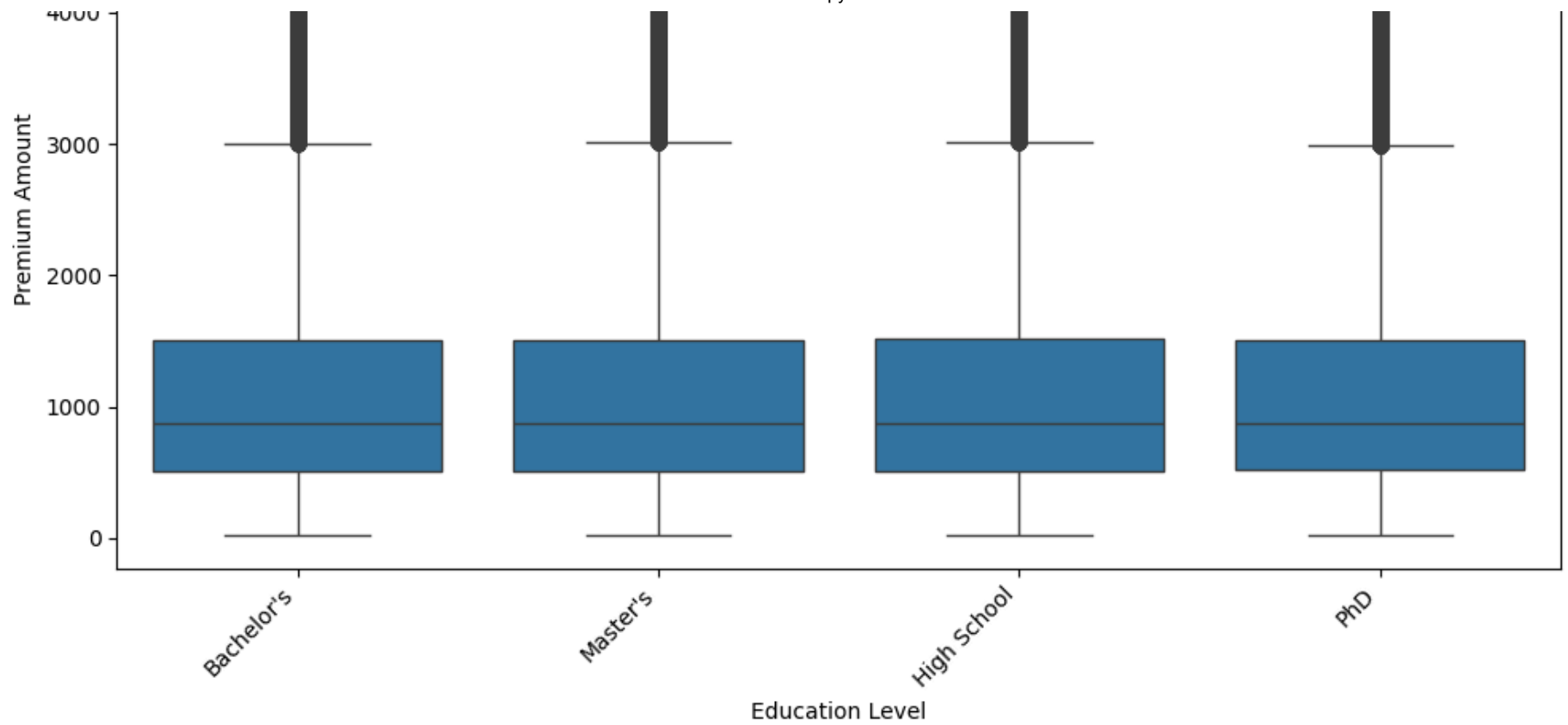


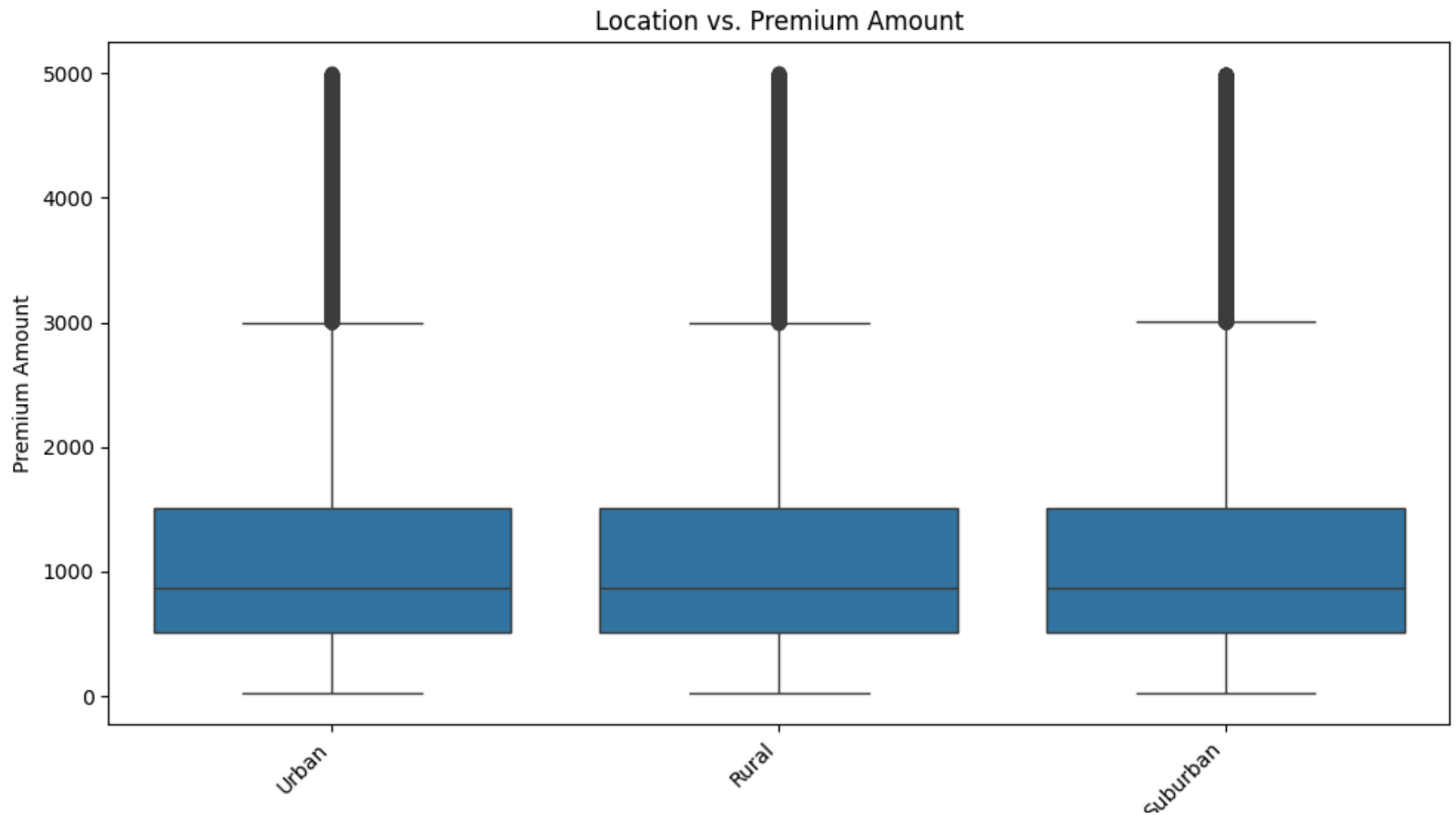
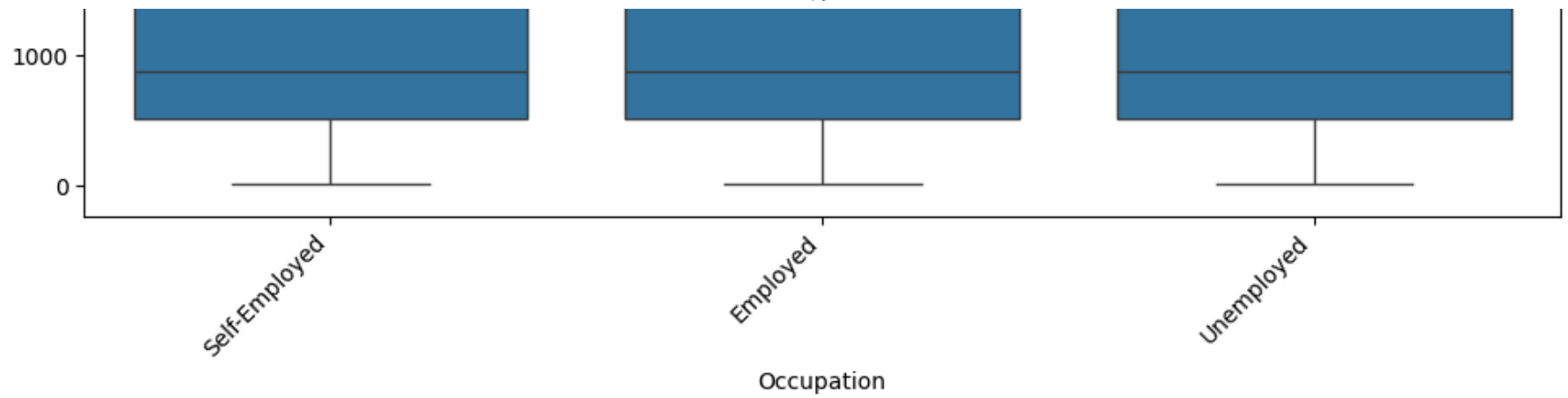


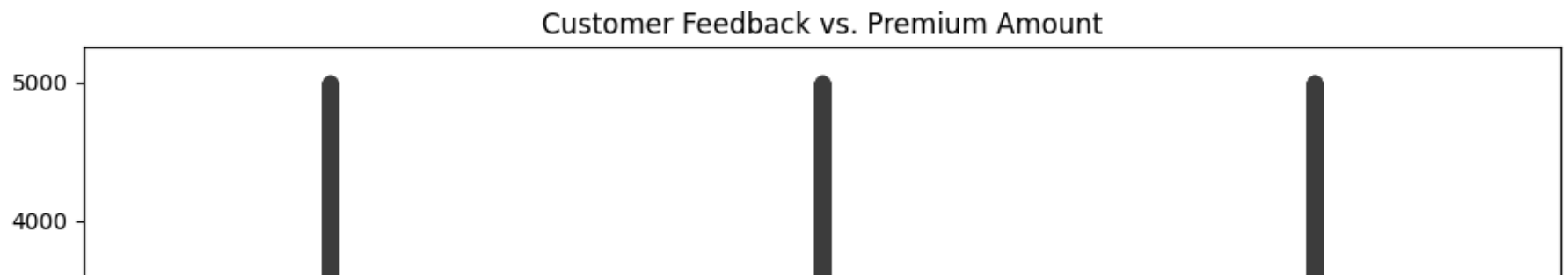
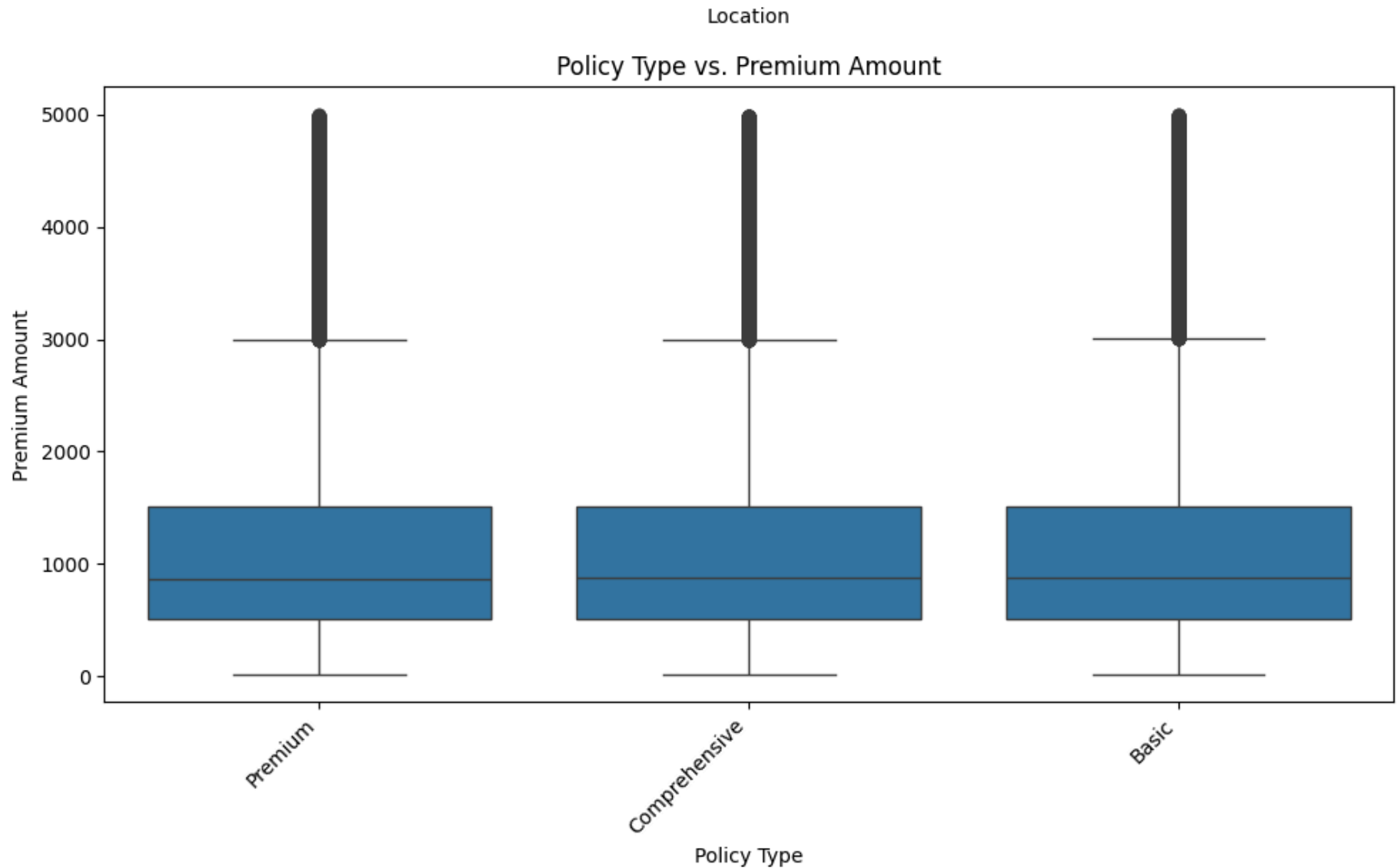
Gender vs. Premium Amount

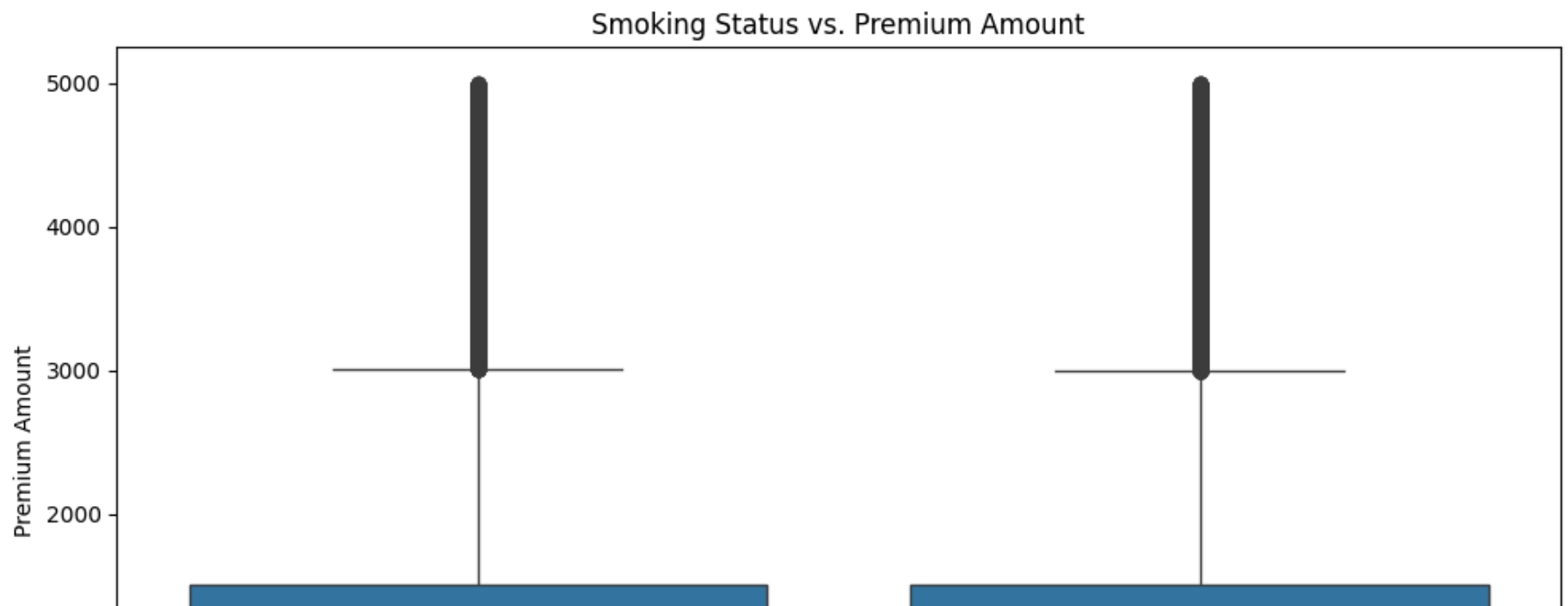
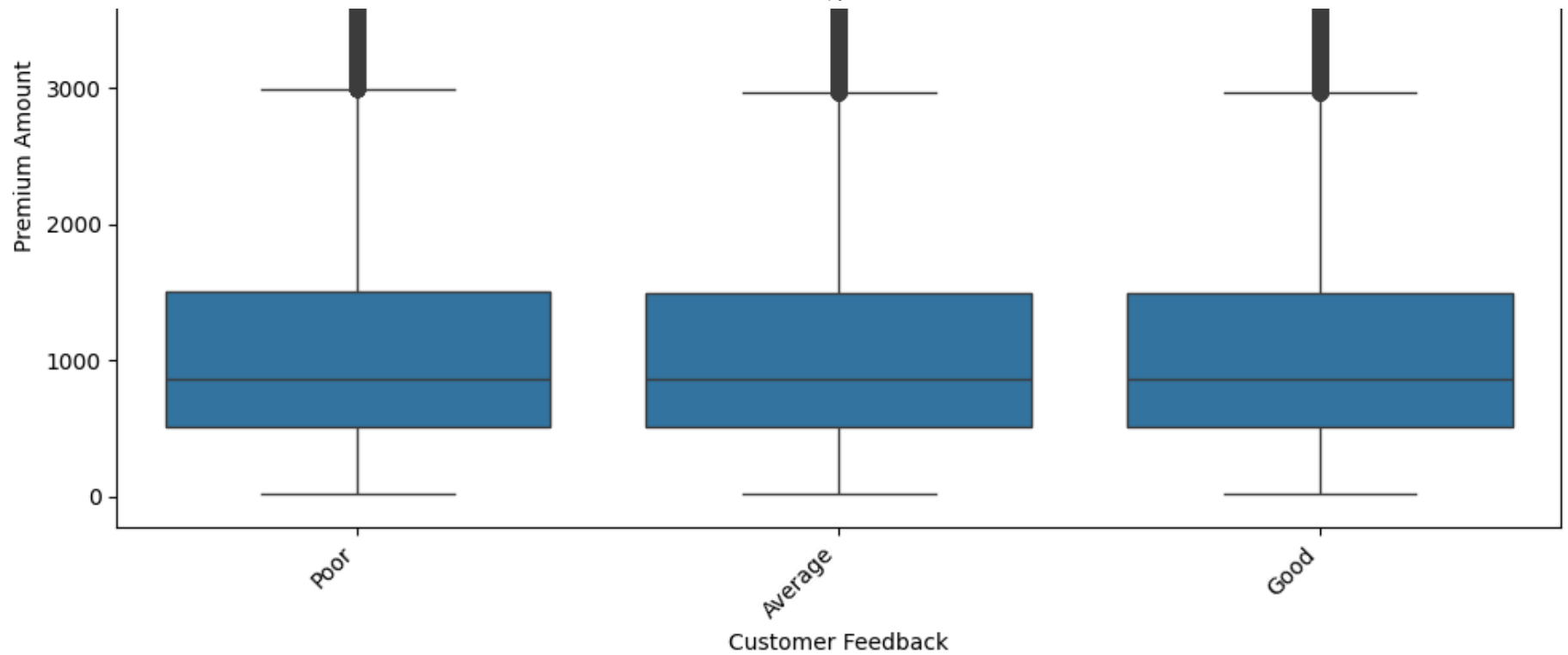


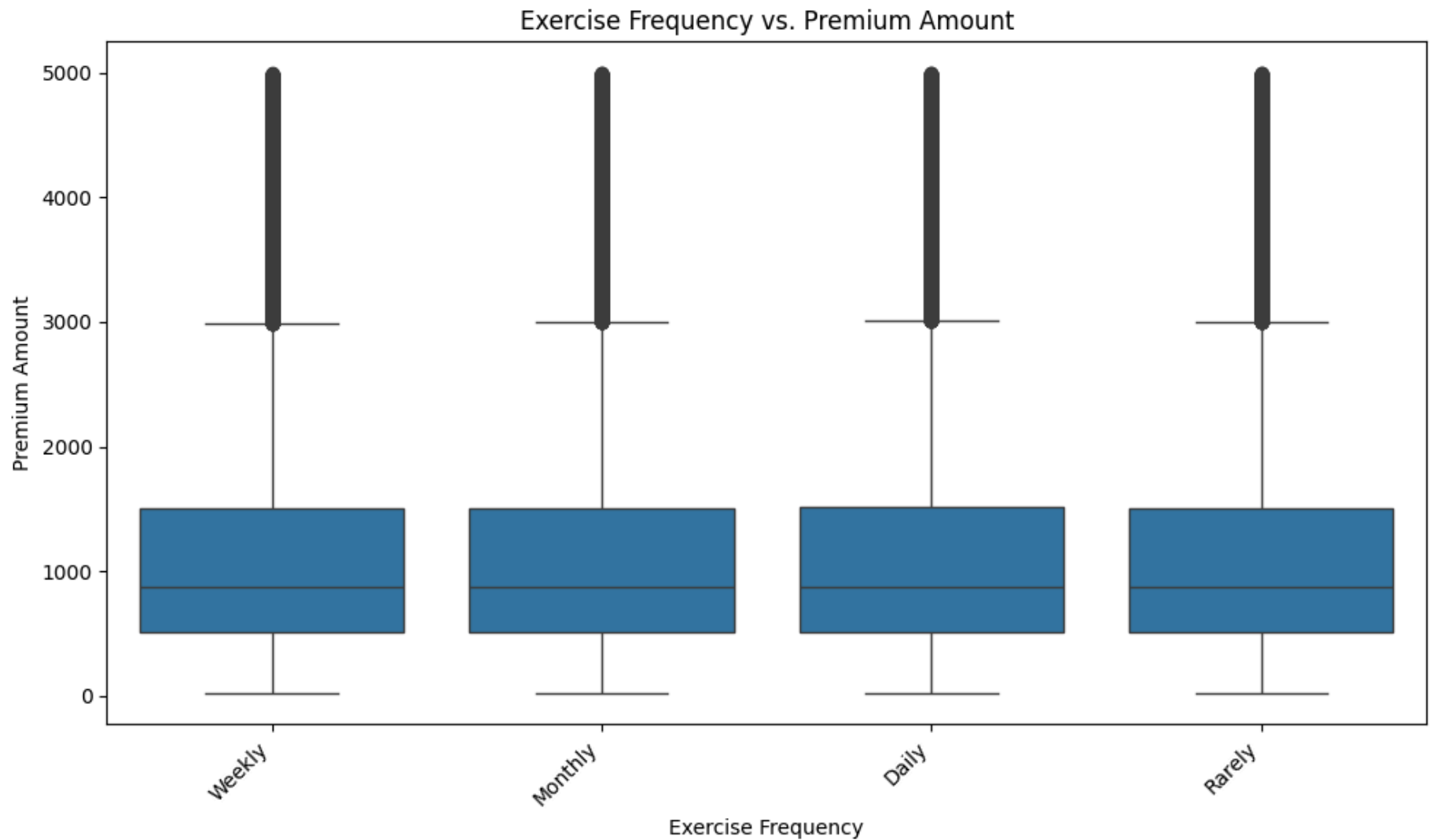
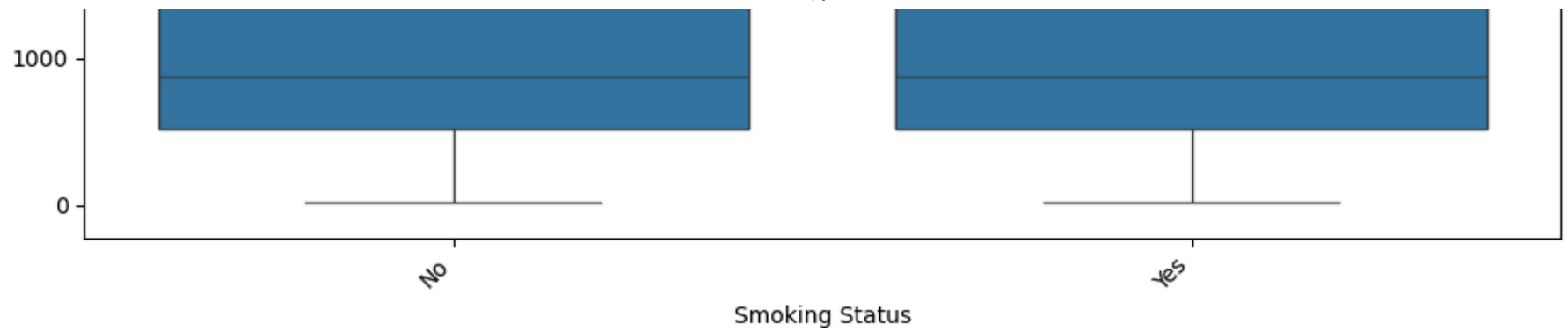


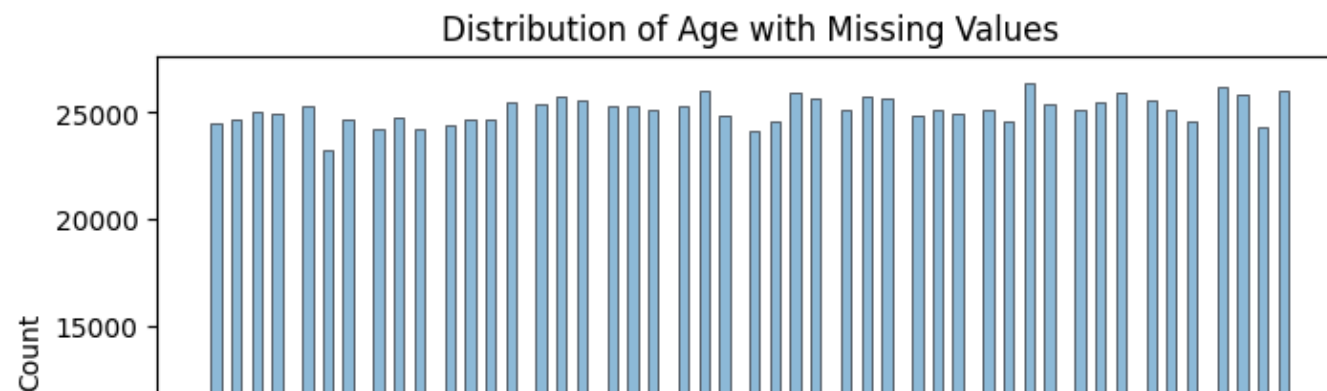
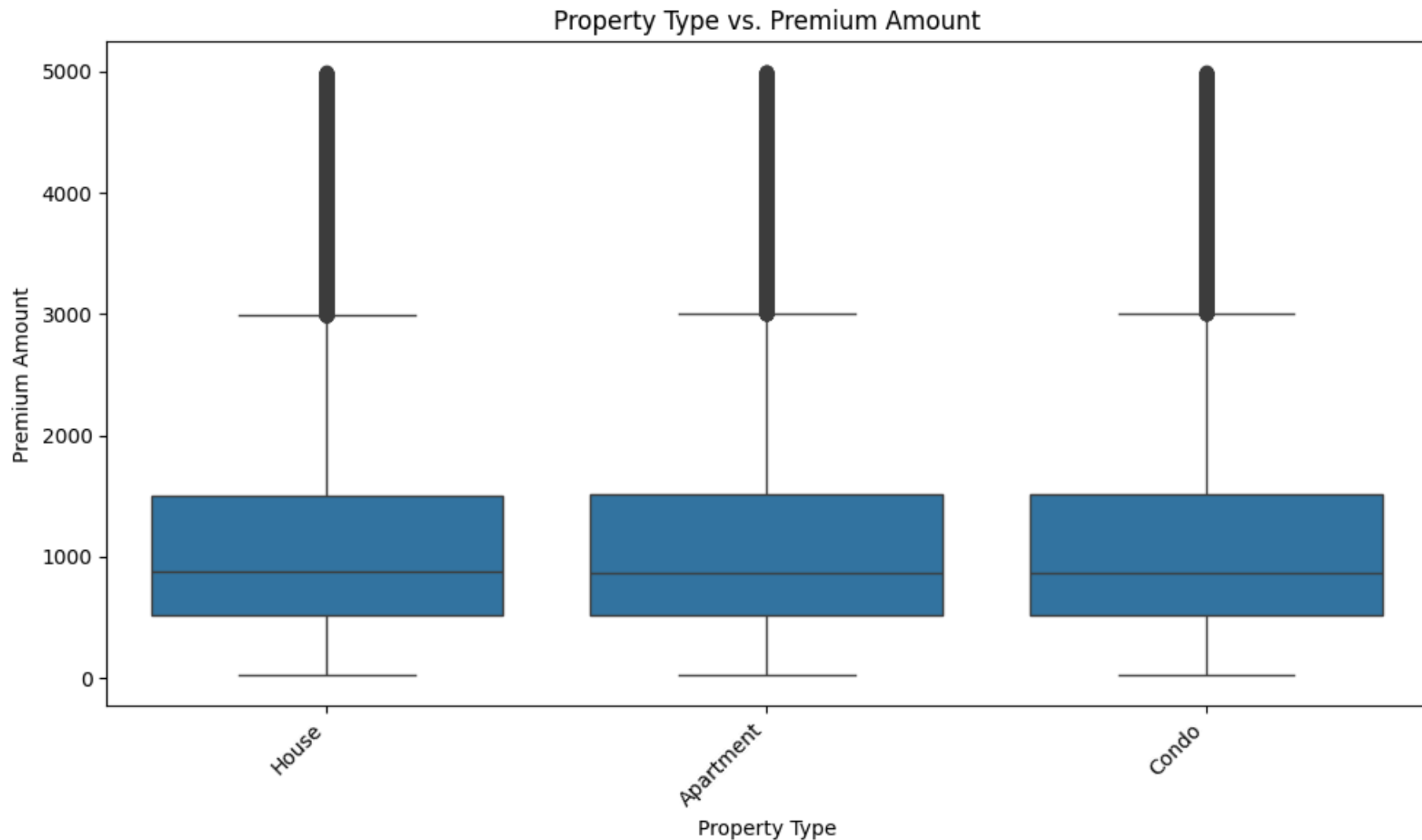


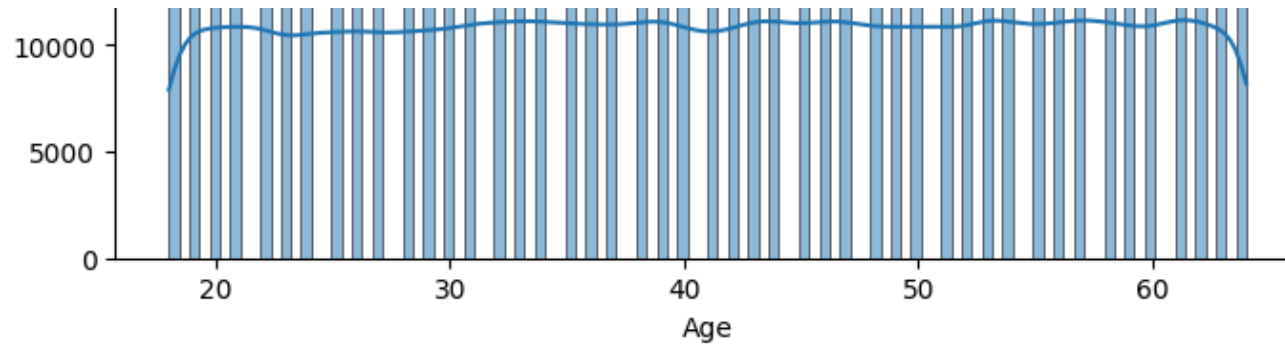




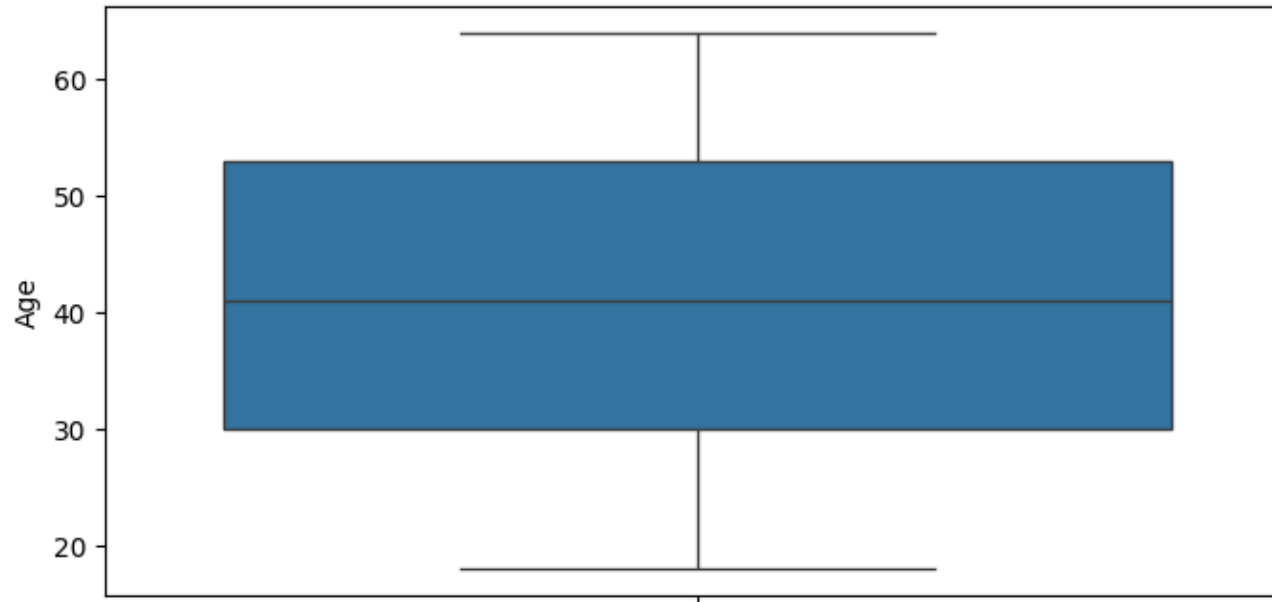




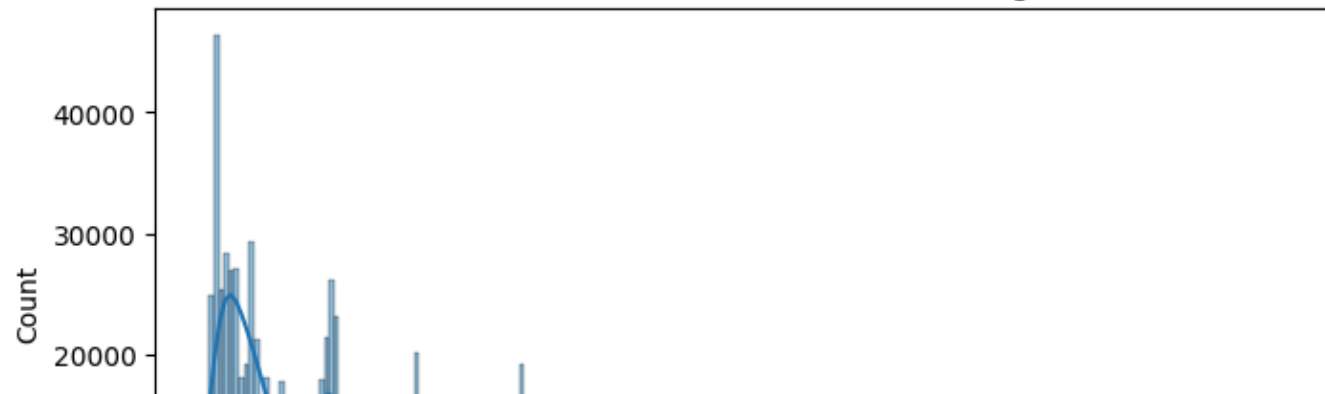


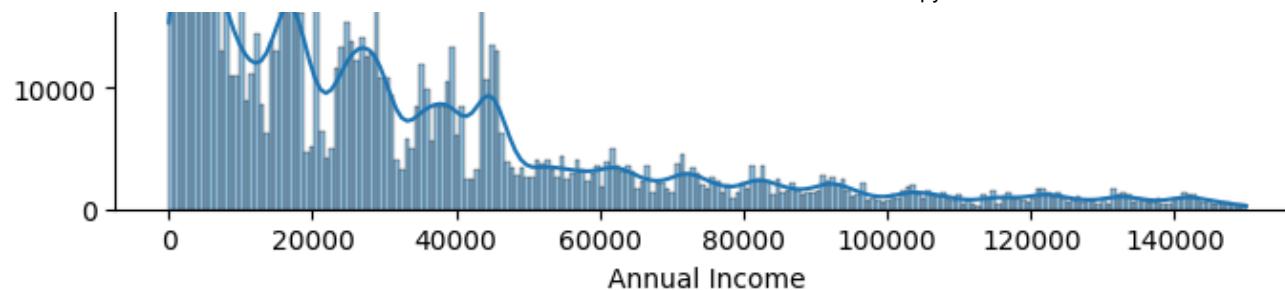


Box Plot for Age with Missing Values

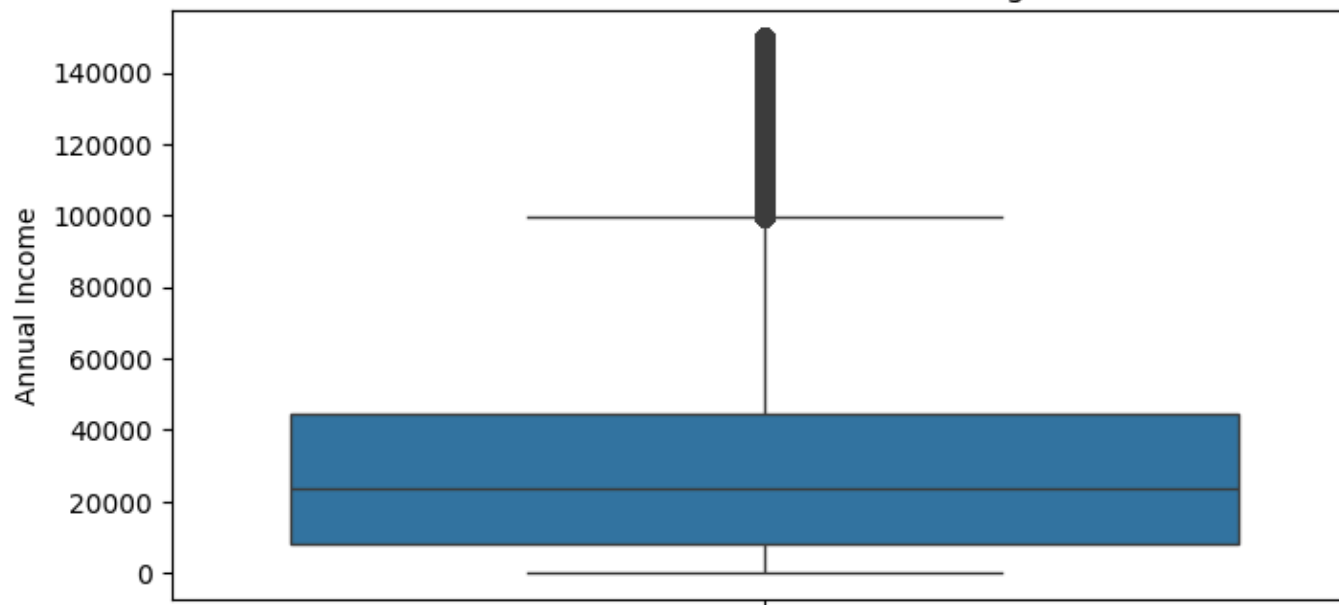


Distribution of Annual Income with Missing Values

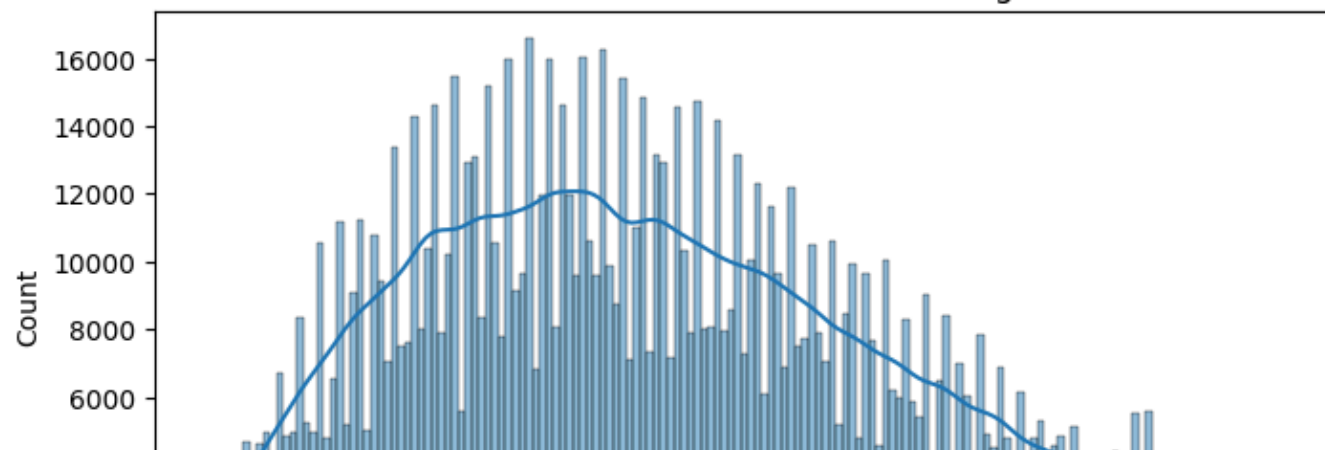


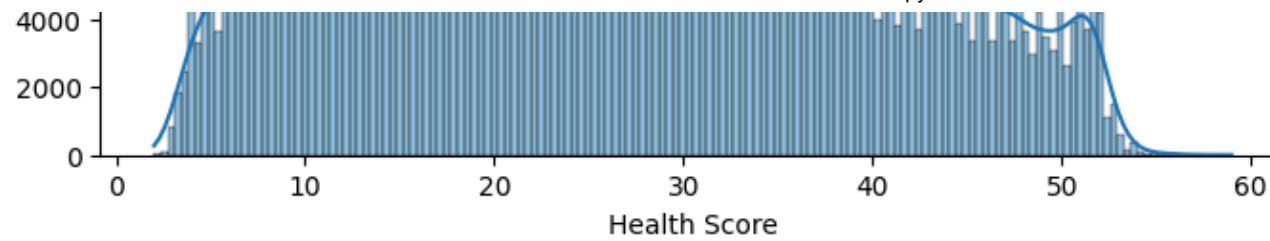


Box Plot for Annual Income with Missing Values

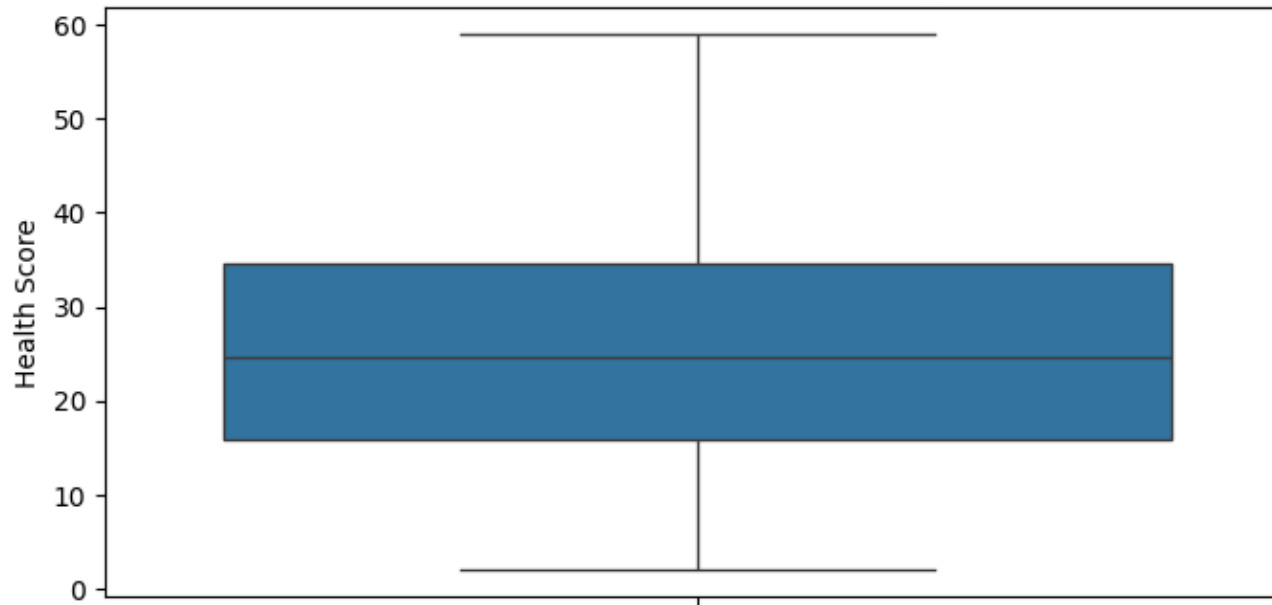


Distribution of Health Score with Missing Values

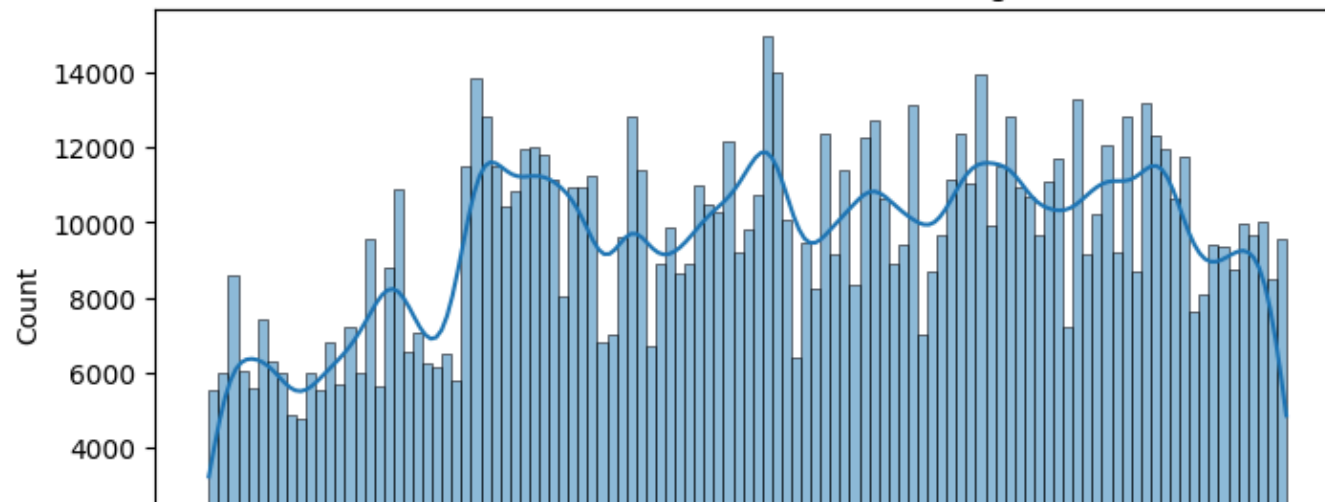


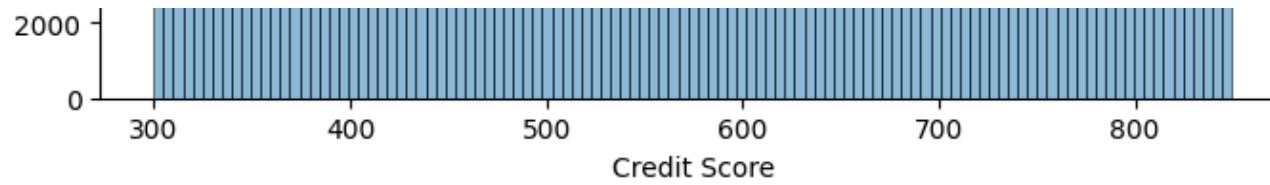


Box Plot for Health Score with Missing Values

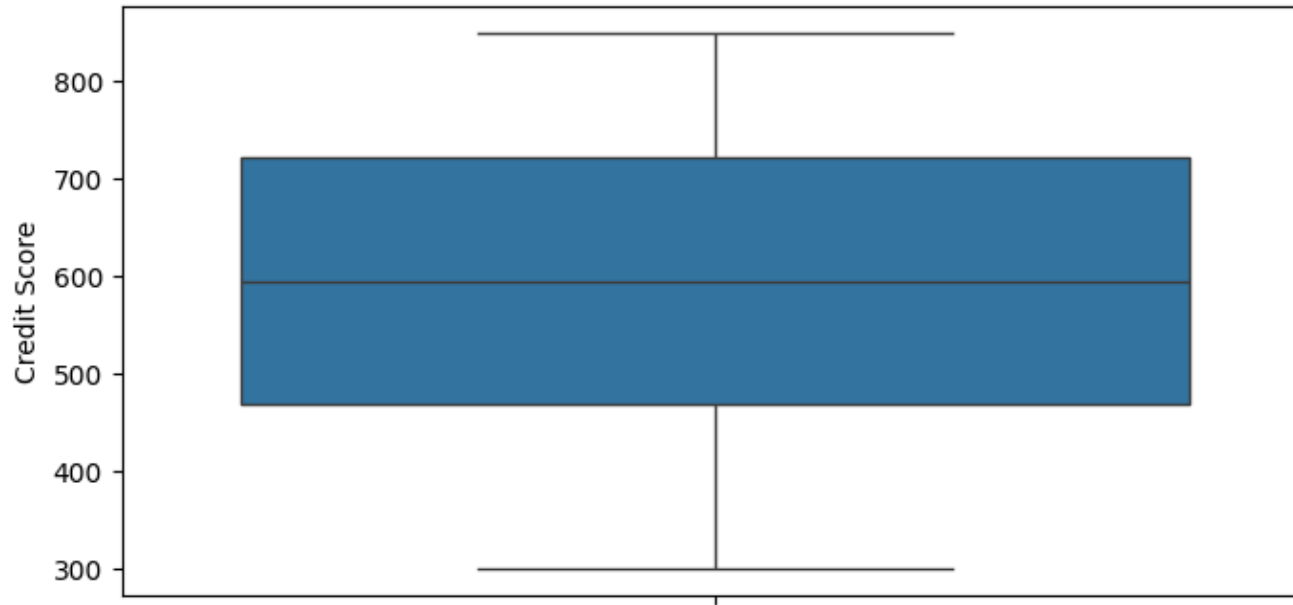


Distribution of Credit Score with Missing Values

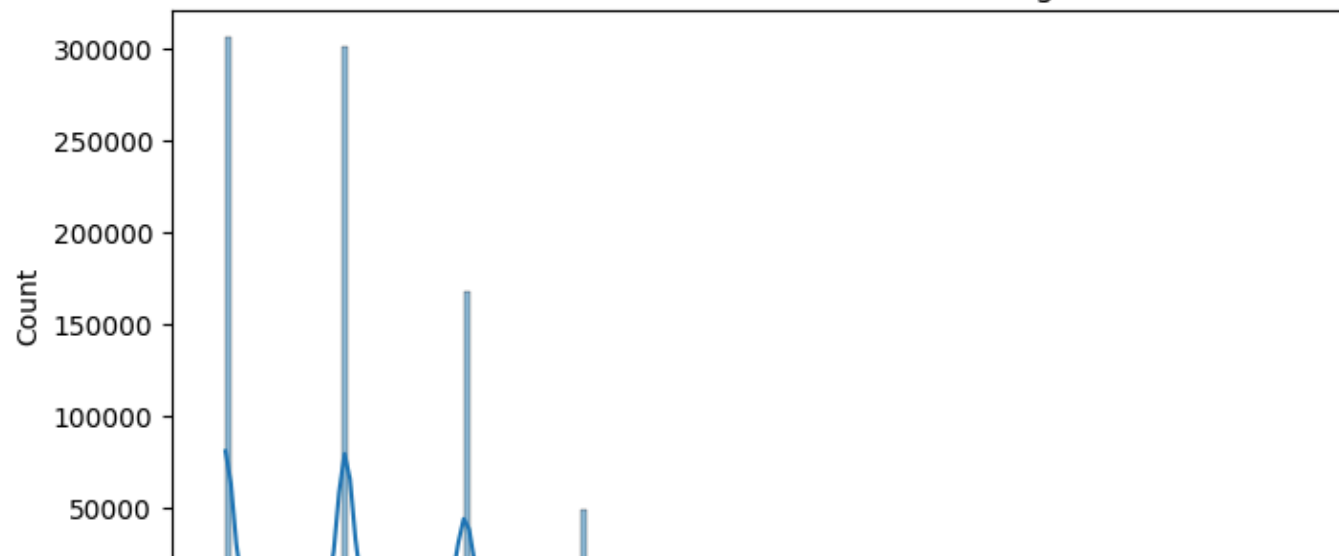


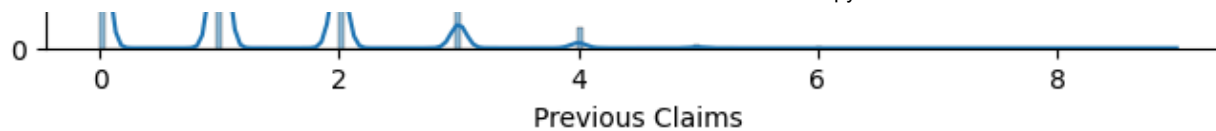


Box Plot for Credit Score with Missing Values

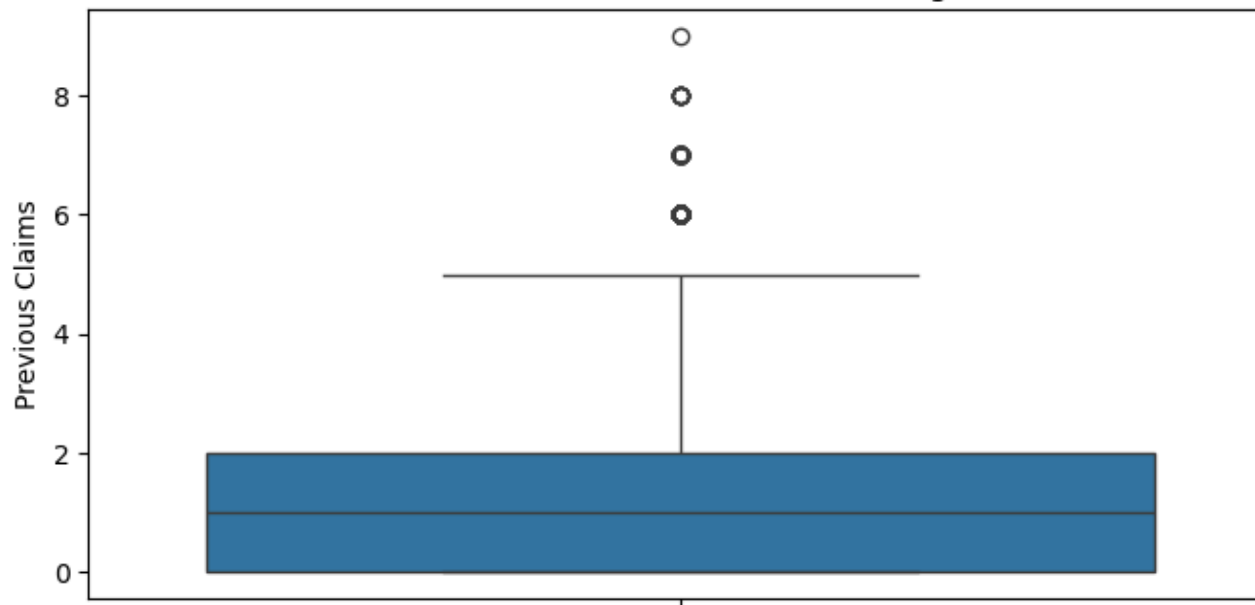


Distribution of Previous Claims with Missing Values

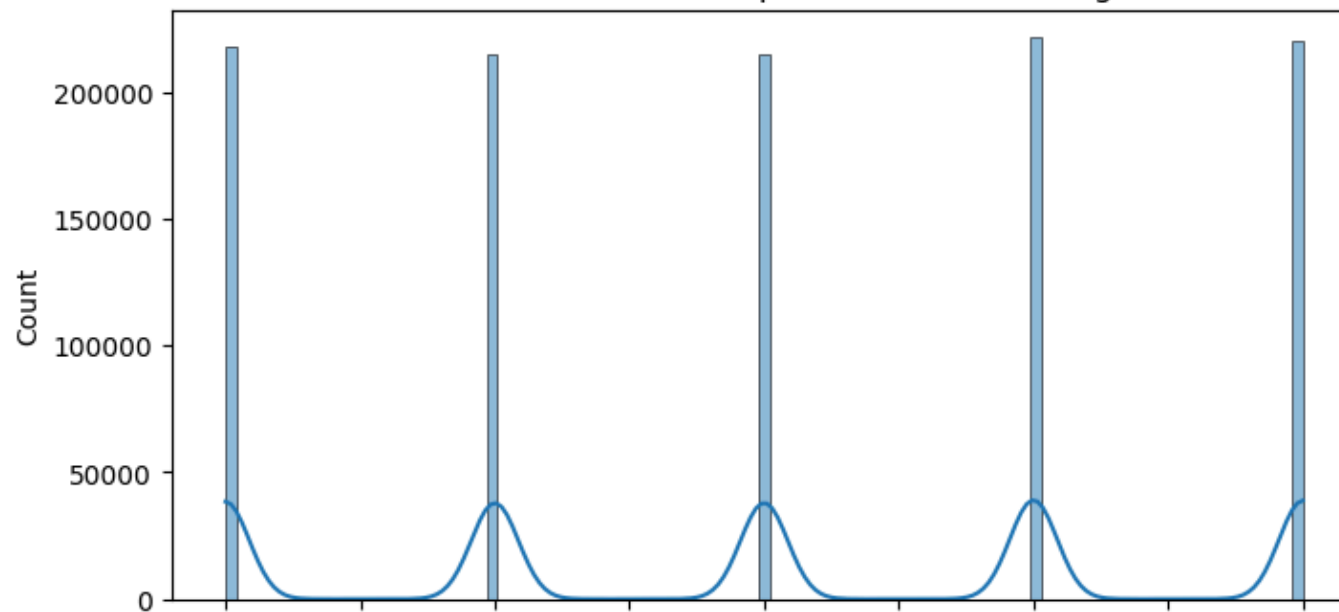




Box Plot for Previous Claims with Missing Values



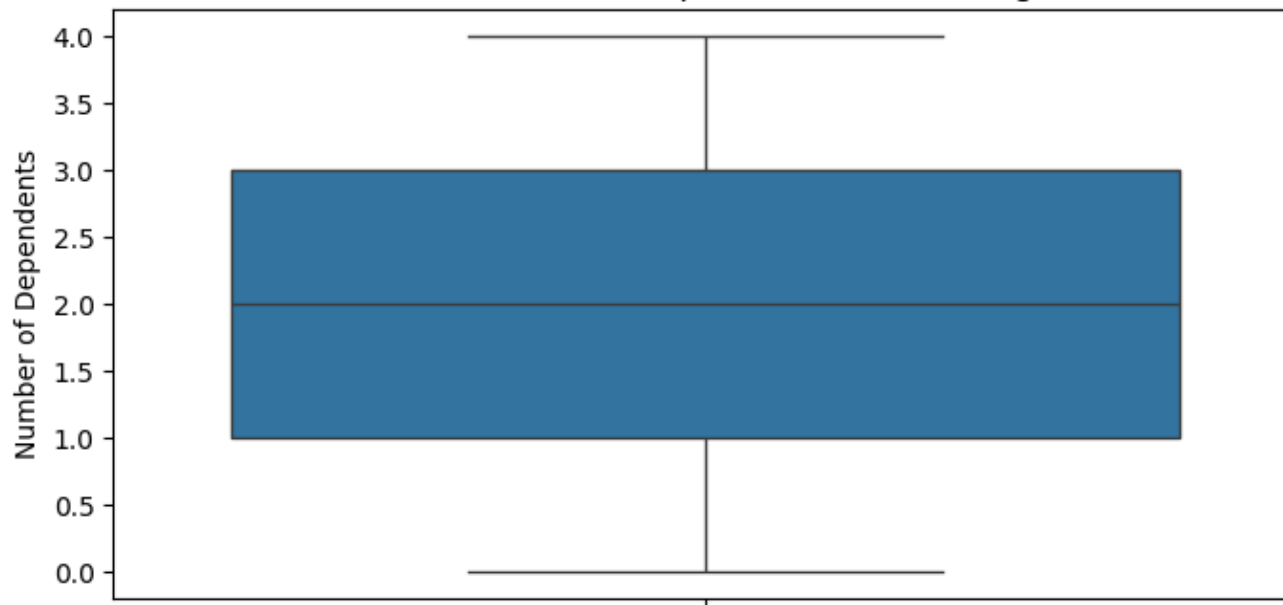
Distribution of Number of Dependents with Missing Values



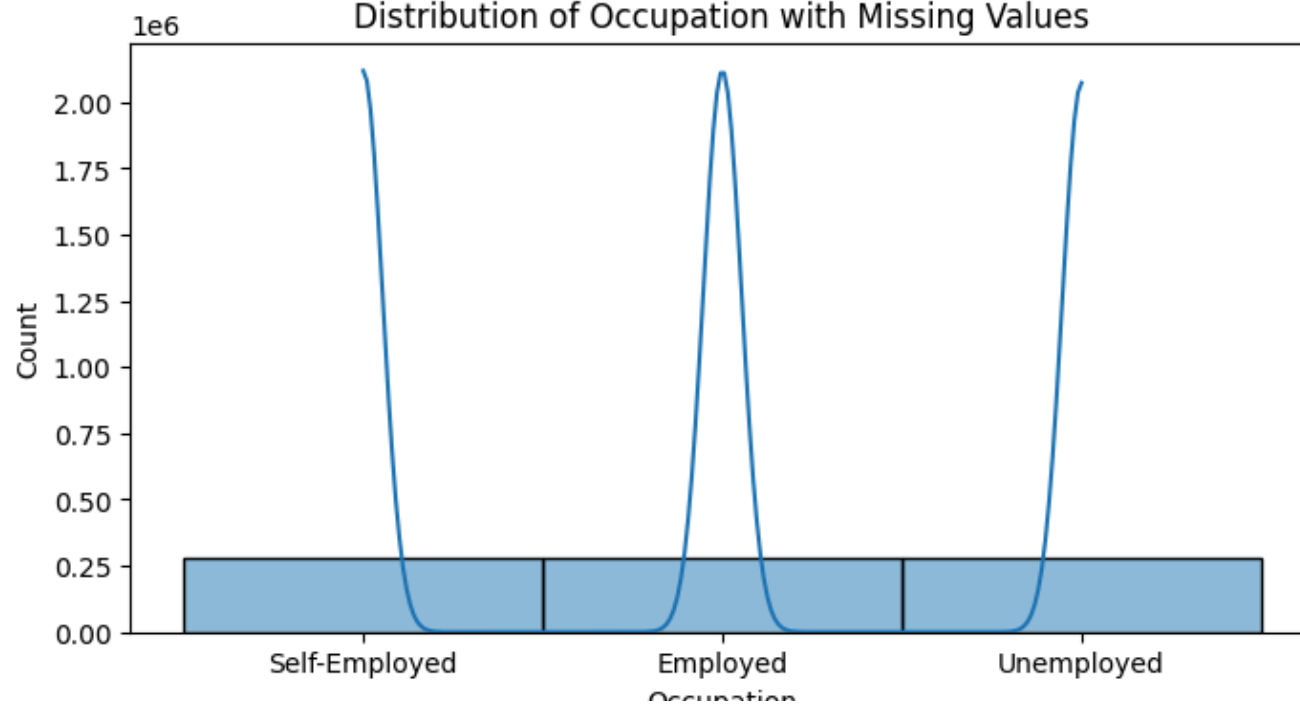
0.0 0.5 1.0 1.5 2.0 2.5 3.0 3.5 4.0

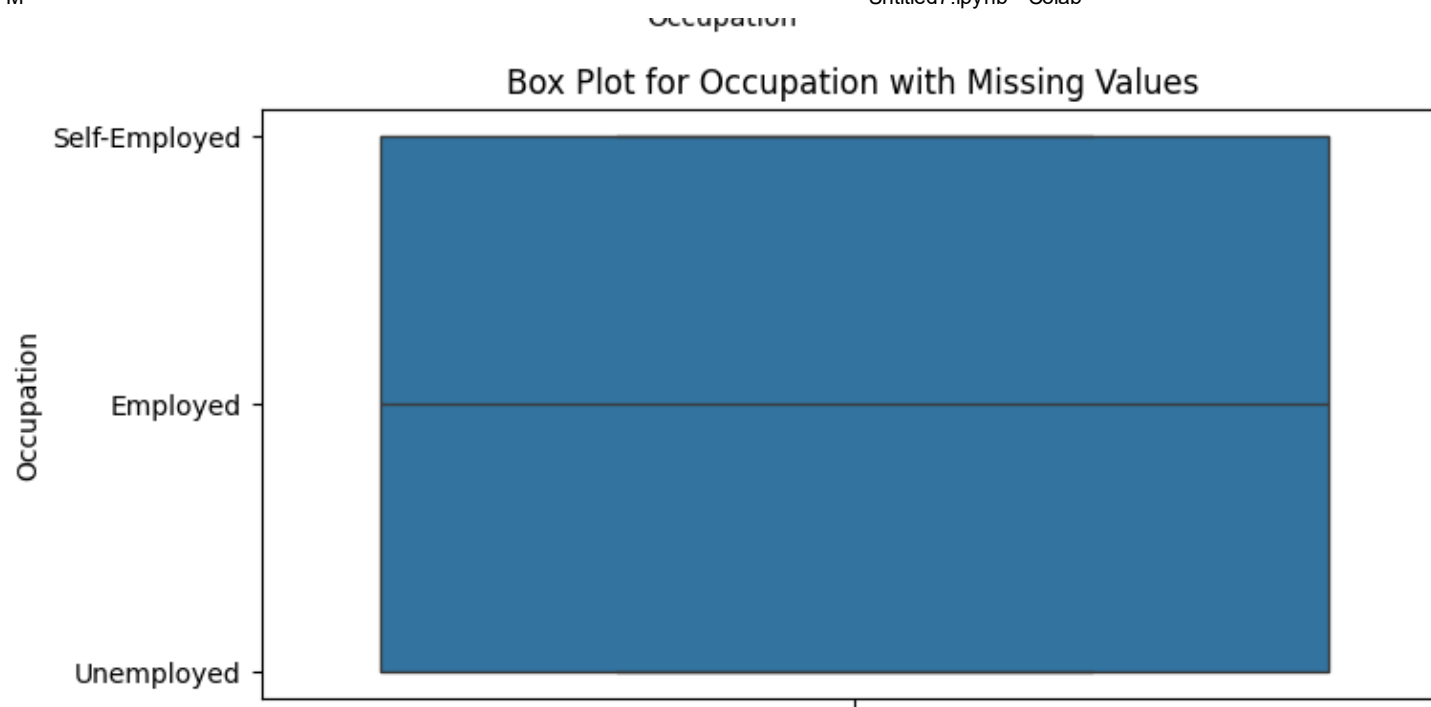
Number of Dependents

Box Plot for Number of Dependents with Missing Values



Distribution of Occupation with Missing Values





✓ Data cleaning


Subtask:

Clean the training and testing datasets (df_train and df_test)

Reasoning: Impute missing values and handle outliers in both df_train and df_test datasets. Use the mean for numerical features and mode for categorical features for imputation. Cap outliers at the 95th percentile. Apply the same strategy to both datasets.

```
1 import pandas as pd
2 import numpy as np
3
4 # Identify numerical and categorical columns
5 numerical_cols = df_train.select_dtypes(include=np.number).columns
6 categorical_cols = df_train.select_dtypes(exclude=np.number).columns
7
8 # Impute missing numerical values with the mean
9 for col in numerical_cols:
10     if df_train[col].isnull().any():
11         mean_val = df_train[col].mean()
12         df_train[col].fillna(mean_val, inplace=True)
13         df_test[col].fillna(mean_val, inplace=True)
14
15 # Impute missing categorical values with the mode
16 for col in categorical_cols:
17     if df_train[col].isnull().any():
18         mode_val = df_train[col].mode()[0]
19         df_train[col].fillna(mode_val, inplace=True)
20         if col in df_test.columns: # Check if the column exists in df_test
21             df_test[col].fillna(mode_val, inplace=True)
```

```
22
23 # Outlier handling (capping at the 95th percentile)
24 for col in numerical_cols:
25     if col != 'id' and col != 'Premium Amount': # Exclude ID and target
26         upper_limit = df_train[col].quantile(0.95)
27         df_train[col] = np.where(df_train[col] > upper_limit, upper_limit, df_train[col])
28         df_test[col] = np.where(df_test[col] > upper_limit, upper_limit, df_test[col])
29
30 display(df_train.head())
31 display(df_test.head())
```

 <ipython-input-4-c81664d17b69>:12: FutureWarning: A value is trying to be set on a copy of a DataFrame or Series through
The behavior will change in pandas 3.0. This inplace method will never work because the intermediate object on which we

For example, when doing 'df[col].method(value, inplace=True)', try using 'df.method({col: value}, inplace=True)' or df[c

```
df_train[col].fillna(mean_val, inplace=True)
```

<ipython-input-4-c81664d17b69>:13: FutureWarning: A value is trying to be set on a copy of a DataFrame or Series through
The behavior will change in pandas 3.0. This inplace method will never work because the intermediate object on which we

For example, when doing 'df[col].method(value, inplace=True)', try using 'df.method({col: value}, inplace=True)' or df[c

```
df_test[col].fillna(mean_val, inplace=True)
```

<ipython-input-4-c81664d17b69>:19: FutureWarning: A value is trying to be set on a copy of a DataFrame or Series through
The behavior will change in pandas 3.0. This inplace method will never work because the intermediate object on which we

For example, when doing 'df[col].method(value, inplace=True)', try using 'df.method({col: value}, inplace=True)' or df[c

```
df_train[col].fillna(mode_val, inplace=True)
```

<ipython-input-4-c81664d17b69>:21: FutureWarning: A value is trying to be set on a copy of a DataFrame or Series through
The behavior will change in pandas 3.0. This inplace method will never work because the intermediate object on which we

For example, when doing 'df[col].method(value, inplace=True)', try using 'df.method({col: value}, inplace=True)' or df[c

```
df_test[col].fillna(mode_val, inplace=True)
```

	id	Age	Gender	Annual Income	Marital Status	Number of Dependents	Education Level	Occupation	Health Score	Location	...	Previous Claims	Vehicle Age	Cr s
0	0	19.0	Female	10049.0	Married	1.0	Bachelor's	Self-Employed	22.598761	Urban	...	2.0	17.0	372.0
1	1	39.0	Female	31678.0	Divorced	3.0	Master's	Employed	15.569731	Rural	...	1.0	12.0	694.0
2	2	23.0	Male	25602.0	Divorced	3.0	High School	Self-Employed	47.177549	Suburban	...	1.0	14.0	592.0
3	3	21.0	Male	103528.0	Married	2.0	Bachelor's	Employed	10.938144	Rural	...	1.0	0.0	367.0
4	4	21.0	Male	39651.0	Single	1.0	Bachelor's	Self-Employed	20.376094	Rural	...	0.0	8.0	598.0

5 rows × 21 columns

	id	Age	Gender	Annual Income	Marital Status	Number of Dependents	Education Level	Occupation	Health Score	Location	Policy Type	Previous Claims
0	1200000	28.0	Female	2310.0	Single	4.0	Bachelor's	Self-Employed	7.657981	Rural	Basic	1.002689
1	1200001	31.0	Female	103528.0	Married	2.0	Master's	Self-Employed	13.381379	Suburban	Premium	1.002689
2	1200002	47.0	Female	17092.0	Divorced	0.0	PhD	Unemployed	24.354527	Urban	Comprehensive	1.002689
3	1200003	28.0	Female	30424.0	Divorced	3.0	PhD	Self-Employed	5.136225	Suburban	Comprehensive	1.000000
4	1200004	24.0	Male	10863.0	Divorced	2.0	High School	Unemployed	11.844155	Suburban	Premium	1.002689

Reasoning: Fix the chained assignment warnings from the previous code block and re-run the data cleaning steps.

```
1 import pandas as pd
2 import numpy as np
3
4 # Identify numerical and categorical columns
5 numerical_cols = df_train.select_dtypes(include=np.number).columns
6 categorical_cols = df_train.select_dtypes(exclude=np.number).columns
7
8 # Impute missing numerical values with the mean
9 for col in numerical_cols:
10     if df_train[col].isnull().any():
11         mean_val = df_train[col].mean()
12         df_train[col] = df_train[col].fillna(mean_val)
13         df_test[col] = df_test[col].fillna(mean_val)
14
15 # Impute missing categorical values with the mode
16 for col in categorical_cols:
17     if df_train[col].isnull().any():
18         mode_val = df_train[col].mode()[0]
19         df_train[col] = df_train[col].fillna(mode_val)
20         if col in df_test.columns: # Check if the column exists in df_test
21             df_test[col] = df_test[col].fillna(mode_val)
22
23 # Outlier handling (capping at the 95th percentile)
24 for col in numerical_cols:
25     if col != 'id' and col != 'Premium Amount': # Exclude ID and target
26         upper_limit = df_train[col].quantile(0.95)
27         df_train[col] = np.where(df_train[col] > upper_limit, upper_limit, df_train[col])
28         df_test[col] = np.where(df_test[col] > upper_limit, upper_limit, df_test[col])
```

29

```
30 display(df_train.head())
```

```
31 display(df_test.head())
```



	id	Age	Gender	Annual Income	Marital Status	Number of Dependents	Education Level	Occupation	Health Score	Location	...	Previous Claims	Vehicle Age	Cr
0	0	19.0	Female	10049.0	Married	1.0	Bachelor's	Self-Employed	22.598761	Urban	...	2.0	17.0	372.0
1	1	39.0	Female	31678.0	Divorced	3.0	Master's	Employed	15.569731	Rural	...	1.0	12.0	694.0
2	2	23.0	Male	25602.0	Divorced	3.0	High School	Self-Employed	47.177549	Suburban	...	1.0	14.0	592.9
3	3	21.0	Male	103528.0	Married	2.0	Bachelor's	Employed	10.938144	Rural	...	1.0	0.0	367.0
4	4	21.0	Male	39651.0	Single	1.0	Bachelor's	Self-Employed	20.376094	Rural	...	0.0	8.0	598.0

5 rows × 21 columns

	id	Age	Gender	Annual Income	Marital Status	Number of Dependents	Education Level	Occupation	Health Score	Location	Policy Type	Previous Claims
0	1200000	28.0	Female	2310.0	Single	4.0	Bachelor's	Self-Employed	7.657981	Rural	Basic	1.002689
1	1200001	31.0	Female	103528.0	Married	2.0	Master's	Self-Employed	13.381379	Suburban	Premium	1.002689
2	1200002	47.0	Female	17092.0	Divorced	0.0	PhD	Unemployed	24.354527	Urban	Comprehensive	1.002689
3	1200003	28.0	Female	30424.0	Divorced	3.0	PhD	Self-Employed	5.136225	Suburban	Comprehensive	1.000000
4	1200004	24.0	Male	10863.0	Divorced	2.0	High School	Unemployed	11.844155	Suburban	Premium	1.002689

✓ Feature engineering

Subtask:

Engineer new features from existing ones in both the training (`df_train`) and testing (`df_test`) datasets to potentially improve model performance.

Reasoning: Apply one-hot encoding to categorical features, create interaction and polynomial features, and then scale numerical features using standardization. This addresses all parts of the subtask instructions in one go.

```
1 import pandas as pd
2 from sklearn.preprocessing import StandardScaler
3
4 # 1. Categorical Encoding
5 categorical_cols = ['Gender', 'Marital Status', 'Education Level', 'Occupation', 'Loca
6 df_train = pd.get_dummies(df_train, columns=categorical_cols, drop_first=True)
7 df_test = pd.get_dummies(df_test, columns=categorical_cols, drop_first=True)
8
9 # Align columns in df_test with df_train
10 missing_cols = set(df_train.columns) - set(df_test.columns)
11 for c in missing_cols:
12     df_test[c] = 0
13 df_test = df_test[df_train.columns] # Ensure same order
14 df_test = df_test.drop(columns=['Premium Amount'])
15
16
17 # 2. Interaction Features (Example: Age * Annual Income)
18 df_train['Age_AnnualIncome'] = df_train['Age'] * df_train['Annual Income']
19 df_test['Age_AnnualIncome'] = df_test['Age'] * df_test['Annual Income']
20
21 # 3. Polynomial Features (Example: Age^2)
22 df_train['Age_squared'] = df_train['Age'] ** 2
23 df_test['Age_squared'] = df_test['Age'] ** 2
24
```

```
25
26 # 4. Feature Scaling
27 numerical_cols = ['Age', 'Annual Income', 'Health Score', 'Number of Dependents', 'Pre
28 scaler = StandardScaler()
29 df_train[numerical_cols] = scaler.fit_transform(df_train[numerical_cols])
30 df_test[numerical_cols] = scaler.transform(df_test[numerical_cols])
31
32 display(df_train.head())
33 display(df_test.head())
```



▼ Data splitting

Subtask:

Split the training data (df_train) into training and validation sets.

Reasoning: Split the training data into training and validation sets, stratifying by the target variable 'Premium Amount'.

```
1 from sklearn.model_selection import train_test_split
2
3 # Separate features (X) and target variable (y)
4 X = df_train.drop(columns=['Premium Amount', 'Policy Start Date', 'id'])
5 y = df_train['Premium Amount']
6
7 # Split data into training and validation sets
8 X_train, X_val, y_train, y_val = train_test_split(X, y, test_size=0.2, random_state=42)
```

```
-----
ValueError                                Traceback (most recent call last)
<ipython-input-7-fe444ea34ee4> in <cell line: 0>()
      6
      7 # Split data into training and validation sets
----> 8 X_train, X_val, y_train, y_val = train_test_split(X, y, test_size=0.2, random_state=42, stratify=y)
```

3 frames

```
/usr/local/lib/python3.11/dist-packages/sklearn/model_selection/_split.py in _iter_indices(self, X, y, groups)
    2316         class_counts = np.bincount(y_indices)
    2317         if np.min(class_counts) < 2:
-> 2318             raise ValueError(
    2319                 "The least populated class in y has only 1"
    2320                 " member, which is too few. The minimum"
```

ValueError: The least populated class in y has only 1 member, which is too few. The minimum number of groups for any class cannot be less than 2.

Next steps: [Explain error](#)

Reasoning: The error message indicates that stratification is failing because at least one class in the target variable 'Premium Amount' has only one member. Since stratification is used to maintain the class distribution in the training and validation sets, and this is not possible with a single member class, try splitting without stratification.

```
1 from sklearn.model_selection import train_test_split
2
3 # Separate features (X) and target variable (y)
4 X = df_train.drop(columns=['Premium Amount', 'Policy Start Date', 'id'])
5 y = df_train['Premium Amount']
6
7 # Split data into training and validation sets without stratification
```