

Medical Insurance Cost Analysis

Importing Libraries

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 %matplotlib inline
5 import seaborn as sns
6 from scipy import stats
7 import requests
8 from sklearn.linear_model import LinearRegression
9 from sklearn.preprocessing import StandardScaler, PolynomialFeatures
10 from sklearn.pipeline import Pipeline
11 from sklearn.metrics import mean_squared_error, r2_score
12 import warnings
13 warnings.filterwarnings("ignore", category = UserWarning)
14 warnings.filterwarnings("ignore", category=FutureWarning, message=".*use_inf_as_na.*")
15 from ipywidgets import interact,interactive,fixed,interact_manual
16 from sklearn.model_selection import train_test_split
17 from sklearn.model_selection import cross_val_score
18 from sklearn.model_selection import cross_val_predict
19 from sklearn.linear_model import Ridge
20 from sklearn.model_selection import GridSearchCV
21 from tqdm import tqdm
22 import requests
```

Importing Datasets

```
1 data= pd.read_csv("insurance.csv")
2 df = data
```

```
1 df.head(10)
```

	Age	Diabetes	BloodPressureProblems	AnyTransplants	AnyChronicDiseases	Height	Weight	KnownAllergies	HistoryOfCancerInFamily	NumberOfMajorSurgeries	Premium
0	45	0	0	0	0	155	57	0	0	0	
1	60	1	0	0	0	180	73	0	0	0	
2	36	1	1	0	0	158	59	0	0	1	
3	52	1	1	0	1	183	93	0	0	2	
4	38	0	0	0	1	166	88	0	0	1	
5	30	0	0	0	0	160	69	1	0	1	
6	33	0	0	0	0	150	54	0	0	0	
7	23	0	0	0	0	181	79	1	0	0	
8	48	1	0	0	0	169	74	1	0	0	
9	38	0	0	0	0	182	93	0	0	0	

Data wrangling

```
1 missing_data = df.isnull()
2 missing_data
```

	Age	Diabetes	BloodPressureProblems	AnyTransplants	AnyChronicDiseases	Height	Weight	KnownAllergies	HistoryOfCancerInFamily	NumberOfMajorSurgeries	Premium
0	False	False	False	False	False	False	False	False	False	False	
1	False	False	False	False	False	False	False	False	False	False	
2	False	False	False	False	False	False	False	False	False	False	
3	False	False	False	False	False	False	False	False	False	False	
4	False	False	False	False	False	False	False	False	False	False	
...	
981	False	False	False	False	False	False	False	False	False	False	
982	False	False	False	False	False	False	False	False	False	False	
983	False	False	False	False	False	False	False	False	False	False	
984	False	False	False	False	False	False	False	False	False	False	
985	False	False	False	False	False	False	False	False	False	False	

986 rows × 11 columns

```
1 for column in missing_data.columns.values.tolist():
2     print(column)
3     print(missing_data[column].value_counts())
4     print("")
```

```
Age
Age
False    986
Name: count, dtype: int64

Diabetes
Diabetes
False    986
Name: count, dtype: int64

BloodPressureProblems
BloodPressureProblems
False    986
Name: count, dtype: int64

AnyTransplants
AnyTransplants
False    986
Name: count, dtype: int64

AnyChronicDiseases
AnyChronicDiseases
False    986
Name: count, dtype: int64

Height
Height
False    986
Name: count, dtype: int64

Weight
Weight
False    986
Name: count, dtype: int64

KnownAllergies
KnownAllergies
False    986
Name: count, dtype: int64

HistoryOfCancerInFamily
HistoryOfCancerInFamily
False    986
Name: count, dtype: int64

NumberOfMajorSurgeries
NumberOfMajorSurgeries
False    986
Name: count, dtype: int64

PremiumPrice
PremiumPrice
False    986
Name: count, dtype: int64
```

```
1 df.dtypes
```

	0
Age	int64
Diabetes	int64
BloodPressureProblems	int64
AnyTransplants	int64
AnyChronicDiseases	int64
Height	int64
Weight	int64
KnownAllergies	int64
HistoryOfCancerInFamily	int64
NumberOfMajorSurgeries	int64
PremiumPrice	int64
dtype:	object

```
1 # Map 'Diabetes' column: yes -> 1, no -> 2
2 df['Diabetes'] = df['Diabetes'].map({'yes': 1, 'no': 2})
3 df[['PremiumPrice']] = np.round(df[['PremiumPrice']],2)
```

```
1 df.head(10)
```

	Age	Diabetes	BloodPressureProblems	AnyTransplants	AnyChronicDiseases	Height	Weight	KnownAllergies	HistoryOfCancerInFamily	NumberOfMajorSurgeries	Premium
0	45	NaN	0	0	0	155	57	0	0	0	
1	60	NaN	0	0	0	180	73	0	0	0	
2	36	NaN	1	0	0	158	59	0	0	1	
3	52	NaN	1	0	1	183	93	0	0	2	
4	38	NaN	0	0	1	166	88	0	0	1	
5	30	NaN	0	0	0	160	69	1	0	1	
6	33	NaN	0	0	0	150	54	0	0	0	
7	23	NaN	0	0	0	181	79	1	0	0	
8	48	NaN	0	0	0	169	74	1	0	0	
9	38	NaN	0	0	0	182	93	0	0	0	

Exploratory Data Analysis (EDA)

Distribution analysis

```
1 df.head(), df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 986 entries, 0 to 985
Data columns (total 11 columns):
#   Column                Non-Null Count  Dtype
---  -
0   Age                   986 non-null   int64
1   Diabetes              0 non-null     float64
2   BloodPressureProblems 986 non-null   int64
3   AnyTransplants        986 non-null   int64
4   AnyChronicDiseases    986 non-null   int64
5   Height                986 non-null   int64
6   Weight                986 non-null   int64
7   KnownAllergies        986 non-null   int64
8   HistoryOfCancerInFamily 986 non-null   int64
9   NumberOfMajorSurgeries 986 non-null   int64
10  PremiumPrice          986 non-null   int64
dtypes: float64(1), int64(10)
memory usage: 84.9 KB
(   Age  Diabetes  BloodPressureProblems  AnyTransplants  AnyChronicDiseases  \
0    45         NaN                     0                0                0
1    60         NaN                     0                0                0
2    36         NaN                     1                0                0
3    52         NaN                     1                0                1
4    38         NaN                     0                0                1

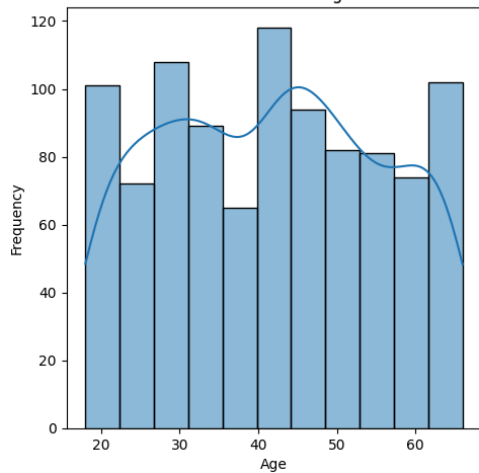
      Height  Weight  KnownAllergies  HistoryOfCancerInFamily  \
0     155     57                0                0
1     180     73                0                0
2     158     59                0                0
3     183     93                0                0
4     166     88                0                0

      NumberOfMajorSurgeries  PremiumPrice
0                0          25000
1                0          29000
2                1          23000
3                2          28000
4                1          23000 ,
None)
```

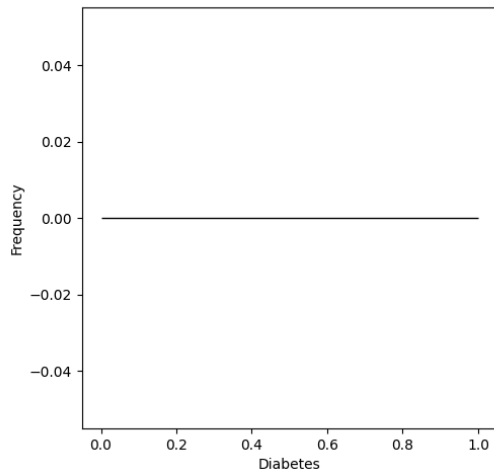
```
1 import matplotlib.pyplot as plt
2 import seaborn as sns
3
4 # Plotting distribution for all numerical variables
5 fig, axes = plt.subplots(nrows=4, ncols=3, figsize=(15, 20))
6 axes = axes.flatten()
7
8 for i, column in enumerate(df.columns):
9     sns.histplot(df[column], kde=True, ax=axes[i])
10    axes[i].set_title(f'Distribution of {column}')
11    axes[i].set_xlabel(column)
12    axes[i].set_ylabel('Frequency')
13
14 # Remove empty subplots
15 for j in range(len(df.columns), len(axes)):
16     fig.delaxes(axes[j])
17
18 plt.tight_layout()
19 plt.show()
```



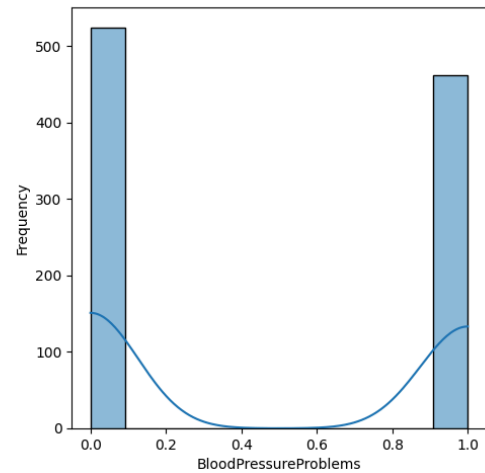
Distribution of Age



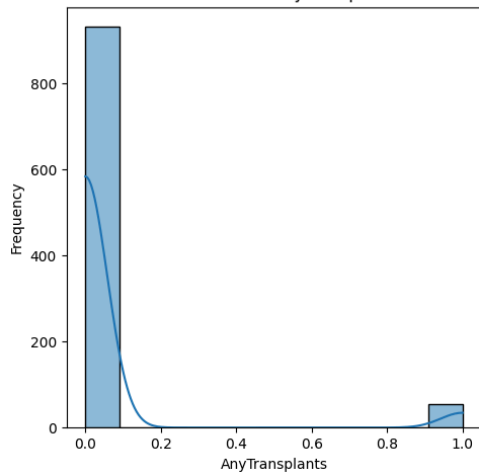
Distribution of Diabetes



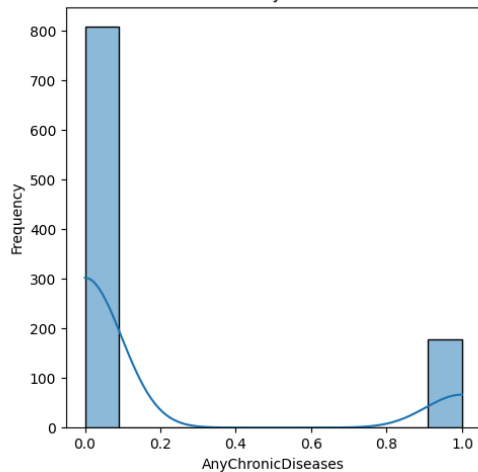
Distribution of BloodPressureProblems



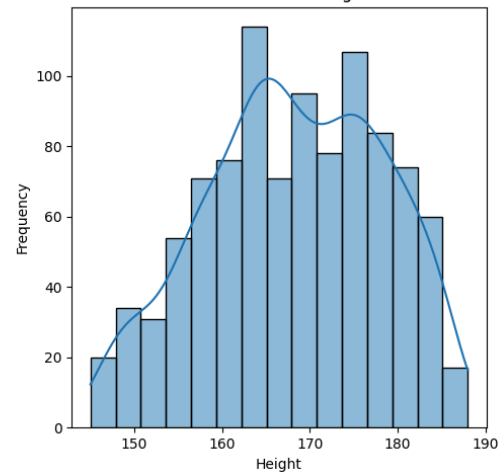
Distribution of AnyTransplants



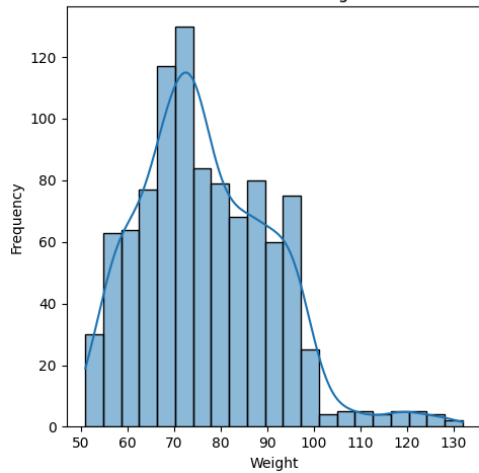
Distribution of AnyChronicDiseases



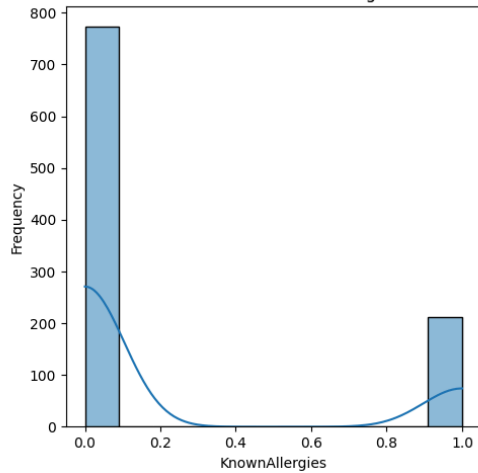
Distribution of Height



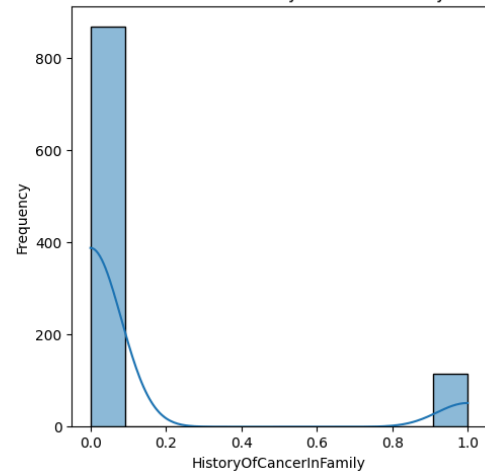
Distribution of Weight



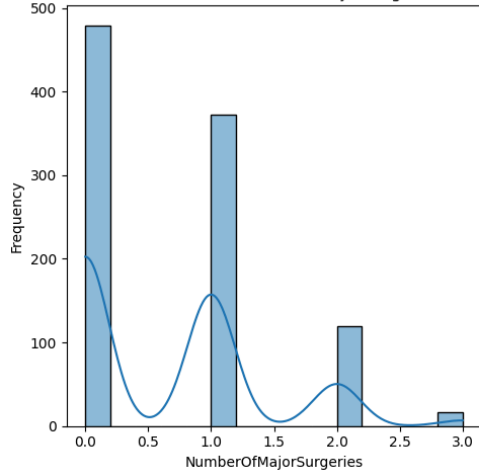
Distribution of KnownAllergies



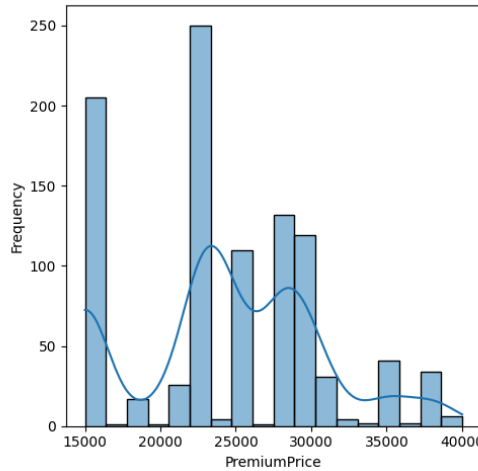
Distribution of HistoryOfCancerInFamily



Distribution of NumberOfMajorSurgeries



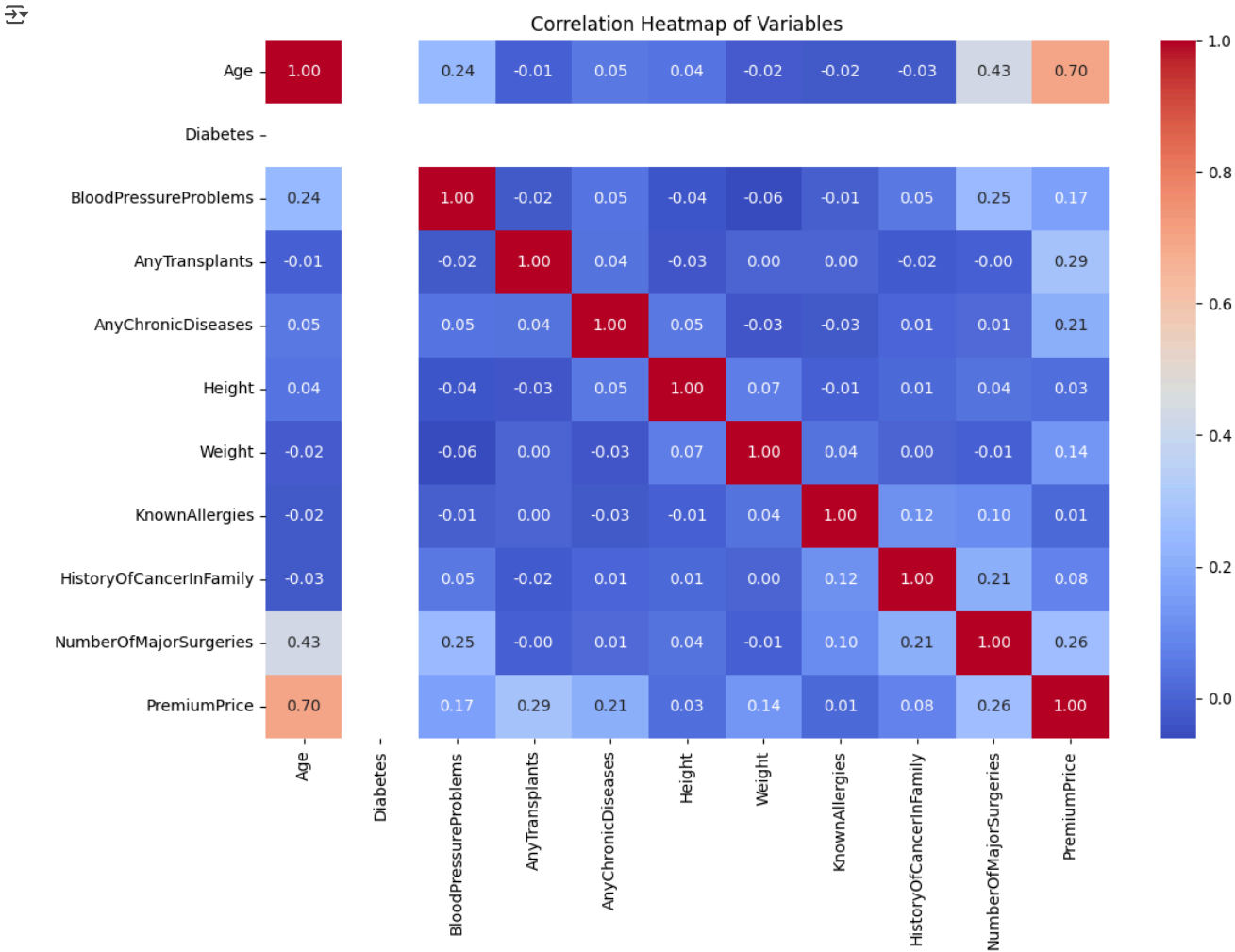
Distribution of PremiumPrice



The histograms above show the distributions of all variables in the dataset. Some variables like PremiumPrice and Age have broader distributions, while others, such as Diabetes and BloodPressureProblems, are binary.

Correlation Analysis

```
1 # Generate correlation matrix
2 correlation_matrix = df.corr()
3
4 # Plot the heatmap
5 plt.figure(figsize=(12, 8))
6 sns.heatmap(correlation_matrix, annot=True, fmt=".2f", cmap="coolwarm", cbar=True)
7 plt.title("Correlation Heatmap of Variables")
8 plt.show()
```



The heatmap displays the correlations between all variables. Notable points include:

- PremiumPrice appears moderately correlated with NumberOfMajorSurgeries and Weight.
- Other variables show low or no correlation with PremiumPrice.

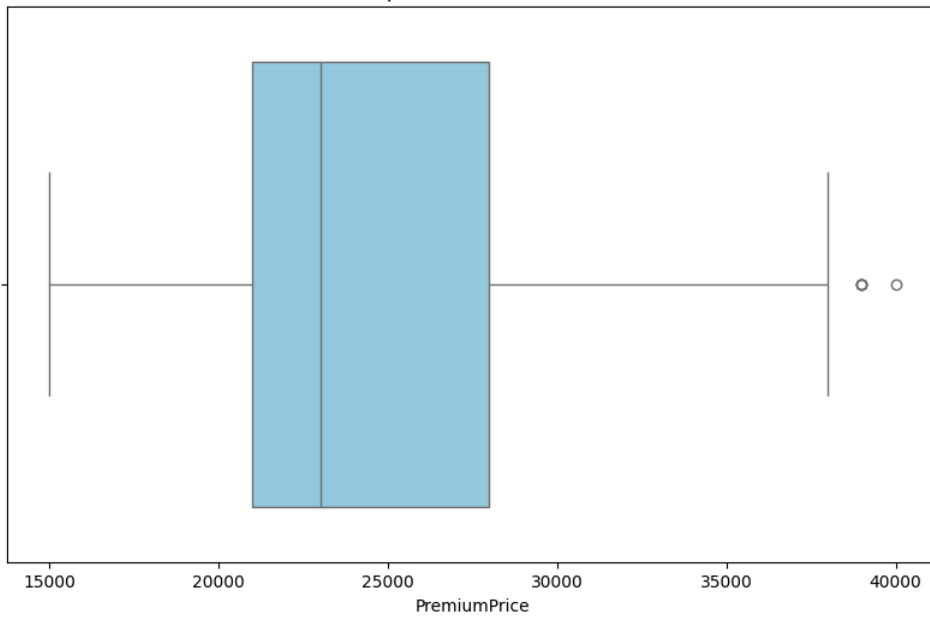
Outlier Detection

```
1 def detect_outliers_iqr(data, column):
2     Q1 = data[column].quantile(0.25)
3     Q3 = data[column].quantile(0.75)
4     IQR = Q3 - Q1
5     lower_bound = Q1 - 1.5 * IQR
6     upper_bound = Q3 + 1.5 * IQR
7     outliers = data[(data[column] < lower_bound) | (data[column] > upper_bound)]
8     return outliers, lower_bound, upper_bound
9
10 # Detecting outliers for PremiumPrice, Weight, and NumberOfMajorSurgeries
11 outliers_premium, lb_premium, ub_premium = detect_outliers_iqr(df, "PremiumPrice")
12 outliers_weight, lb_weight, ub_weight = detect_outliers_iqr(df, "Weight")
13 outliers_surgeries, lb_surgeries, ub_surgeries = detect_outliers_iqr(df, "NumberOfMajorSurgeries")
14
15 # Outputting the results
16 len(outliers_premium), lb_premium, ub_premium, \
17 len(outliers_weight), lb_weight, ub_weight, \
18 len(outliers_surgeries), lb_surgeries, ub_surgeries
19
20 # Boxplot for PremiumPrice
21 plt.figure(figsize=(10, 6))
22 sns.boxplot(data=df, x="PremiumPrice", color="skyblue")
23 plt.title("Boxplot of PremiumPrice")
```

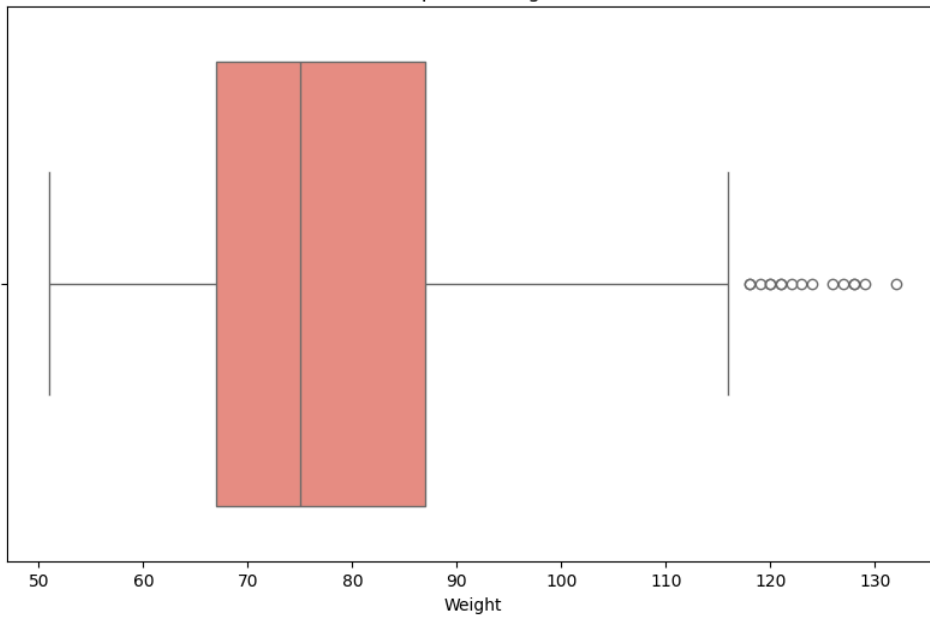
```
24 plt.xlabel("PremiumPrice")
25 plt.show()
26
27 # Boxplot for Weight
28 plt.figure(figsize=(10, 6))
29 sns.boxplot(data=df, x="Weight", color="salmon")
30 plt.title("Boxplot of Weight")
31 plt.xlabel("Weight")
32 plt.show()
33
34 # Boxplot for NumberOfMajorSurgeries
35 plt.figure(figsize=(10, 6))
36 sns.boxplot(data=df, x="NumberOfMajorSurgeries", color="lightgreen")
37 plt.title("Boxplot of NumberOfMajorSurgeries")
38 plt.xlabel("Number of Major Surgeries")
39 plt.show()
```



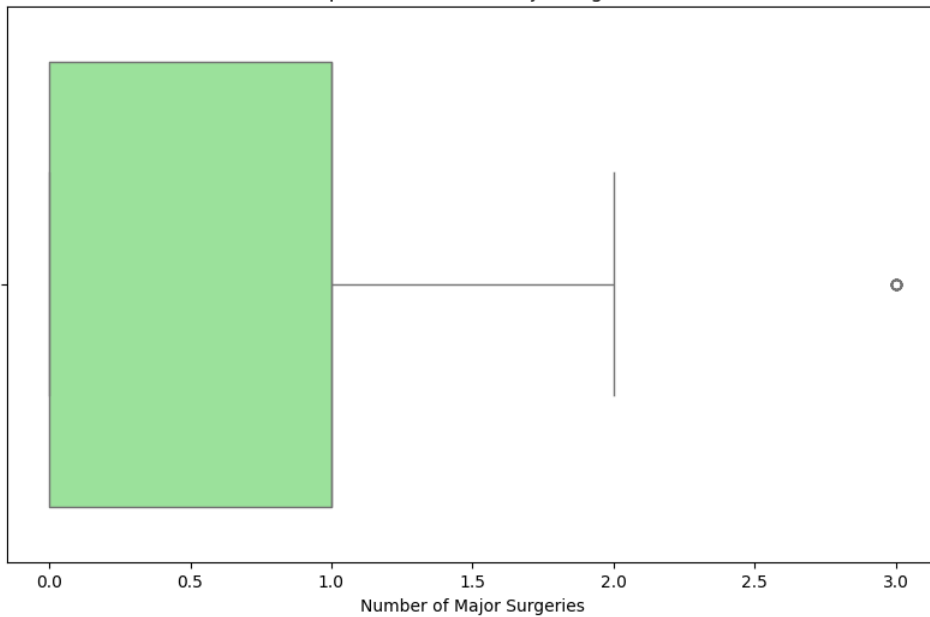
Boxplot of PremiumPrice



Boxplot of Weight



Boxplot of NumberOfMajorSurgeries

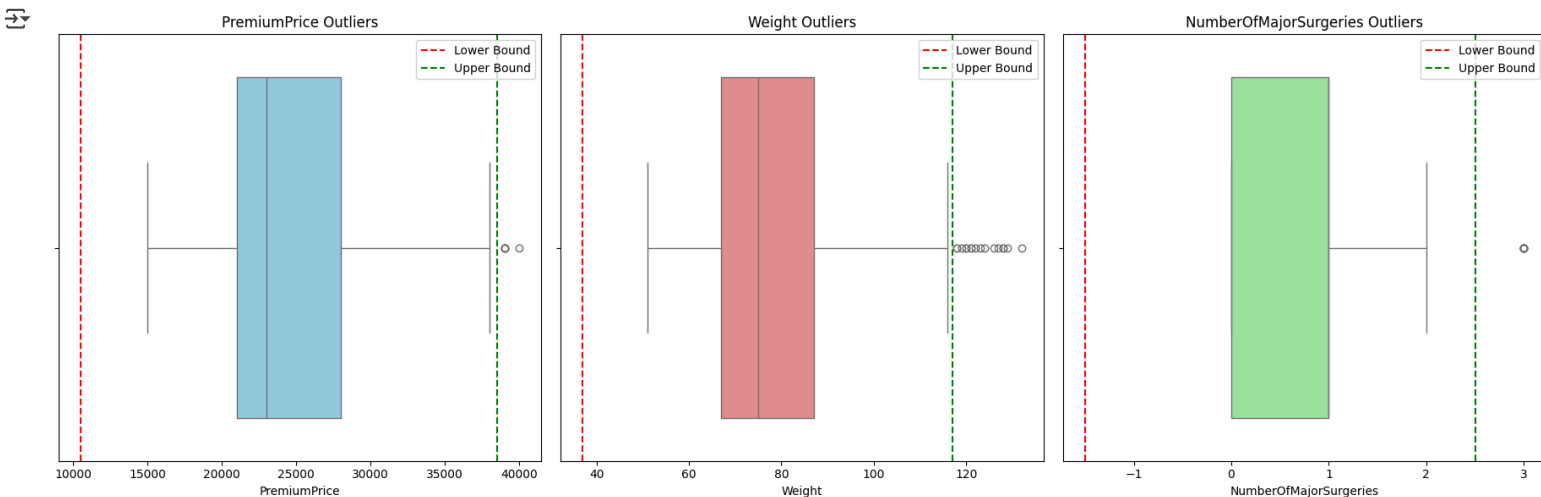


```
1 import matplotlib.pyplot as plt
2 import seaborn as sns
3 # Visualizing outliers for selected variables
4 fig, axes = plt.subplots(nrows=1, ncols=3, figsize=(18, 6))
5
6 # PremiumPrice
7 # Ensure 'df' is defined by running the preceding cells that load and prepare the data.
8 sns.boxplot(data=df, x="PremiumPrice", ax=axes[0], color="skyblue")
9 axes[0].axvline(lb_premium, color="red", linestyle="--", label="Lower Bound")
10 axes[0].axvline(ub_premium, color="green", linestyle="--", label="Upper Bound")
```

```

1 axes[0].set_title("PremiumPrice Outliers")
2 axes[0].legend()
3
4 # Weight
5 sns.boxplot(data=df, x="Weight", ax=axes[1], color="lightcoral")
6 axes[1].axvline(lb_weight, color="red", linestyle="--", label="Lower Bound")
7 axes[1].axvline(ub_weight, color="green", linestyle="--", label="Upper Bound")
8 axes[1].set_title("Weight Outliers")
9 axes[1].legend()
10
11 # NumberOfMajorSurgeries
12 sns.boxplot(data=df, x="NumberOfMajorSurgeries", ax=axes[2], color="lightgreen")
13 axes[2].axvline(lb_surgeries, color="red", linestyle="--", label="Lower Bound")
14 axes[2].axvline(ub_surgeries, color="green", linestyle="--", label="Upper Bound")
15 axes[2].set_title("NumberOfMajorSurgeries Outliers")
16 axes[2].legend()
17
18 plt.tight_layout()
19 plt.show()

```



The boxplots above illustrate the detected outliers:

- PremiumPrice: Outliers lie significantly beyond the upper bound (38,500), indicating high premiums.
- Weight: Outliers include both underweight and overweight individuals beyond the range [37, 117].
- NumberOfMajorSurgeries: Outliers are above 2 surgeries, as expected.

Handling strategies might include:

- PremiumPrice: Winsorizing or capping at a threshold.
- Weight: Normalizing or segmenting into categories.
- NumberOfMajorSurgeries: Considering separate models for groups with many surgeries.

Hypothesis Testing

```

1 from scipy.stats import ttest_ind, chi2_contingency, f_oneway
2
3 # Example: T-test for smokers vs non-smokers (assuming data has such columns)
4 group1 = df[df['Diabetes'] == 0]['PremiumPrice']
5 group2 = df[df['Diabetes'] == 1]['PremiumPrice']
6 t_stat, p_value = ttest_ind(group1, group2)
7 print(f"T-test: t-statistic = {t_stat}, p-value = {p_value}")
8
9 # Example: ANOVA for NumberOfMajorSurgeries and PremiumPrice
10 anova_stat, anova_p = f_oneway(
11     df[df['NumberOfMajorSurgeries'] == 0]['PremiumPrice'],
12     df[df['NumberOfMajorSurgeries'] == 1]['PremiumPrice'],
13     df[df['NumberOfMajorSurgeries'] >= 2]['PremiumPrice']
14 )
15 print(f"ANOVA: F-statistic = {anova_stat}, p-value = {anova_p}")
16
17 # Example: Chi-square test for HistoryOfCancerInFamily and AnyChronicDiseases
18 contingency_table = pd.crosstab(
19     df['HistoryOfCancerInFamily'],
20     df['AnyChronicDiseases']
21 )
22 chi2_stat, chi2_p, _, _ = chi2_contingency(contingency_table)
23 print(f"Chi-square test: Chi2-statistic = {chi2_stat}, p-value = {chi2_p}")

```

```

T-test: t-statistic = nan, p-value = nan
ANOVA: F-statistic = 39.2415277169692, p-value = 4.0158930660447635e-17
Chi-square test: Chi2-statistic = 0.02062393388215223, p-value = 0.8858081638149811
/usr/local/lib/python3.11/dist-packages/scipy/_lib/deprecation.py:234: SmallSampleWarning: One or more sample arguments is too small; all returned values will be N
return f(*args, **kwargs)

```

```

1 from scipy.stats import f_oneway, ttest_ind
2

```



```

3 # ANOVA for PremiumPrice across NumberOfMajorSurgeries (0, 1, 2)
4 groups_surgeries = [df[df["NumberOfMajorSurgeries"] == group]["PremiumPrice"]
5                       for group in df["NumberOfMajorSurgeries"].unique()]
6 anova_surgeries = f_oneway(*groups_surgeries)
7
8 # T-tests for PremiumPrice across AnyChronicDiseases (0 vs. 1)
9 group_chronic_0 = df[df["AnyChronicDiseases"] == 0]["PremiumPrice"]
10 group_chronic_1 = df[df["AnyChronicDiseases"] == 1]["PremiumPrice"]
11 ttest_chronic = ttest_ind(group_chronic_0, group_chronic_1, equal_var=False)
12
13 # T-tests for PremiumPrice across HistoryOfCancerInFamily (0 vs. 1)
14 group_cancer_0 = df[df["HistoryOfCancerInFamily"] == 0]["PremiumPrice"]
15 group_cancer_1 = df[df["HistoryOfCancerInFamily"] == 1]["PremiumPrice"]
16 ttest_cancer = ttest_ind(group_cancer_0, group_cancer_1, equal_var=False)
17
18 # Output results
19 anova_surgeries.pvalue, ttest_chronic.pvalue, ttest_cancer.pvalue

```

```

(np.float64(2.8711631377228097e-16),
 np.float64(1.7279736467737666e-13),
 np.float64(0.01982822652964323))

```

The results of the hypothesis tests are as follows:

- ANOVA for PremiumPrice across NumberOfMajorSurgeries:

p-value: 2.87e-16

Interpretation: The mean PremiumPrice differs significantly across groups (0, 1, 2 surgeries).

- T-test for PremiumPrice across AnyChronicDiseases (0 vs. 1):

p-value: 1.73e-13

Interpretation: There is a statistically significant difference in PremiumPrice between individuals with and without chronic diseases.

- T-test for PremiumPrice across HistoryOfCancerInFamily (0 vs. 1):

p-value: 0.0198

Interpretation: There is a statistically significant difference in PremiumPrice between those with and without a family history of cancer, though the effect may be less pronounced.

ML Modeling

Data Preprocessing

```

1 from sklearn.model_selection import train_test_split
2 from sklearn.preprocessing import StandardScaler, OneHotEncoder
3 from sklearn.compose import ColumnTransformer
4 from sklearn.pipeline import Pipeline
5 # Step 1: Handle Missing Values
6 # Check for missing values
7 print("Missing Values:\n", data.isnull().sum())
8
9 # If missing values exist, handle them
10 if data.isnull().sum().any():
11     # Fill numerical features with median
12     for col in data.select_dtypes(include=np.number).columns:
13         data[col].fillna(data[col].median(), inplace=True)
14     # Fill categorical features with mode
15     for col in data.select_dtypes(include='object').columns:
16         data[col].fillna(data[col].mode()[0], inplace=True)
17
18 # Step 2: Feature Engineering
19 # Example: Calculate BMI (Body Mass Index) if height and weight columns exist
20 if 'height' in data.columns and 'weight' in data.columns:
21     data['BMI'] = data['weight'] / (data['height'] / 100) ** 2
22
23 # Split the dataset BEFORE identifying numerical/categorical features for the preprocessor
24 X = data.drop('PremiumPrice', axis=1)
25 y = data['PremiumPrice']
26
27 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
28
29 # Step 3: Scaling and Encoding
30 # Identify numerical and categorical columns AFTER splitting
31 numerical_features = X_train.select_dtypes(include=np.number).columns.tolist()
32 categorical_features = X_train.select_dtypes(include='object').columns.tolist()
33
34
35 # Preprocessing Pipeline
36 preprocessor = ColumnTransformer(
37     transformers=[
38         ('num', StandardScaler(), numerical_features),
39         ('cat', OneHotEncoder(handle_unknown='ignore'), categorical_features) # Added handle_unknown for robustness
40     ]
41 )
42
43
44 # Apply the preprocessing pipeline
45 X_train = preprocessor.fit_transform(X_train)
46 X_test = preprocessor.transform(X_test)
47
48 print("Preprocessing Complete. Ready for Model Selection.")

```

```
➦ Missing Values:
  Age      0
  Diabetes 0
  BloodPressureProblems 0
  AnyTransplants 0
  AnyChronicDiseases 0
  Height    0
  Weight    0
  KnownAllergies 0
  HistoryOfCancerInFamily 0
  NumberOfMajorSurgeries 0
  PremiumPrice 0
  dtype: int64
Preprocessing Complete. Ready for Model Selection.
```

Model Selection

```
1 from sklearn.linear_model import LinearRegression
2 from sklearn.tree import DecisionTreeRegressor
3 from sklearn.ensemble import RandomForestRegressor, GradientBoostingRegressor
4 from sklearn.neural_network import MLPRegressor
5
6 # Step 1: Linear Regression
7 linear_model = LinearRegression()
8 linear_model.fit(X_train, y_train)
9 linear_preds = linear_model.predict(X_test)
10
11 # Step 2: Decision Tree
12 tree_model = DecisionTreeRegressor(random_state=42)
13 tree_model.fit(X_train, y_train)
14 tree_preds = tree_model.predict(X_test)
15
16 # Step 3: Random Forest
17 forest_model = RandomForestRegressor(random_state=42)
18 forest_model.fit(X_train, y_train)
19 forest_preds = forest_model.predict(X_test)
20
21 # Step 4: Gradient Boosting
22 gb_model = GradientBoostingRegressor(random_state=42)
23 gb_model.fit(X_train, y_train)
24 gb_preds = gb_model.predict(X_test)
25
26 # Step 5: Neural Network
27 nn_model = MLPRegressor(random_state=42, max_iter=500)
28 nn_model.fit(X_train, y_train)
29 nn_preds = nn_model.predict(X_test)
30
31 print("Model Training Complete. Ready for Evaluation.")
```

```
➦ Model Training Complete. Ready for Evaluation.
```

Model Evaluation and Validation

```
1 from sklearn.metrics import mean_squared_error, mean_absolute_error, r2_score
2
3 # Function to evaluate models
4 def evaluate_model(name, y_test, predictions):
5     mse = mean_squared_error(y_test, predictions)
6     mae = mean_absolute_error(y_test, predictions)
7     r2 = r2_score(y_test, predictions)
8     print(f"{name} Evaluation:")
9     print(f" - MSE: {mse}")
10    print(f" - MAE: {mae}")
11    print(f" - R2: {r2}\n")
12
13 # Evaluate all models
14 evaluate_model("Linear Regression", y_test, linear_preds)
15 evaluate_model("Decision Tree", y_test, tree_preds)
16 evaluate_model("Random Forest", y_test, forest_preds)
17 evaluate_model("Gradient Boosting", y_test, gb_preds)
18 evaluate_model("Neural Network", y_test, nn_preds)
```

```
➦ Linear Regression Evaluation:
- MSE: 12221661.705858208
- MAE: 2586.2253840681074
- R2: 0.713394427027874
```

```
Decision Tree Evaluation:
- MSE: 14893939.393939395
- MAE: 1095.959595959596
- R2: 0.6507278521900168
```

```
Random Forest Evaluation:
- MSE: 5310708.080808081
- MAE: 1034.7474747474748
- R2: 0.8754605904647057
```

```
Gradient Boosting Evaluation:
- MSE: 6166949.130591766
- MAE: 1501.7274850951546
- R2: 0.8553812049783704
```

```
Neural Network Evaluation:
- MSE: 577411808.7458334
- MAE: 23238.412330525465
- R2: -12.540666258758183
```