

# D208 Performance Assessment Task 1

This is the code for my d208 performance assessment task 1

## A1. Research Question

For this performance assessment, my research question is "What quantitative variables most significantly contribute to Tenure?". This is an important question to ask because understanding what factors contribute to tenure can allow the business to retain customers for a longer period of time. Focusing on retaining customers is an important practice because often times it is easier to retain customers rather than to try and attract new ones.

## A2. State Objectives and Goals for Analysis

The goal of this analysis is to gain greater insight into what factors directly correlate with Tenure. In this analysis, I will be using linear regression modeling. Using linear regression, we can identify what multiple quantitative variables correlate with a single quantitative variable. The objective is to finish this analysis with a list of quantitative variables that significantly correlate with churn.

## B1. Assumptions

There are four assumptions of a multiple linear regression model that we should consider. These are:

- There is a linear relationship between the dependent variable and the independent variable. This is an issue because the assumption of linearity is violated if it is not true, which violates the fundamental assumptions of the model and brings its accuracy into question.
- The independent variables are not too highly correlated with each other. This is Multicollinearity. Multicollinearity occurs when two or more independent variables are highly correlated with each other. The reason that this causes problems is that it can be difficult to determine which of these correlated variables is the one that is affecting the dependent variable.
- Observations are selected independently and randomly from the population. This is a similar issue to multicollinearity where various independent variables are working together to create an effect on the dependent, and one's effect may be overstated.
- Residuals should be normally distributed with a mean of zero. This is called Homoscedasticity. It is an issue because it can lead to biased and inefficient parameter estimates.

## B2. Programming Language and Benefits

The programming language that I used for this analysis is Python. Two reasons why I am using this language are:

- I am familiar with the language. I am not as comfortable with using R and understand how to code in Python. This will make the project completion more timely and efficient.
- Access to python libraries that can do multiple linear regression. There are a widespread list of libraries that I can use to finish my analysis for this project. This flexibility assists in timely

The libraries that I will be using in this analysis are as follows:

- Pandas: This library is essential to import the CSV and apply analysis to the data.
- numpy: We use numpy to use arrays and set up the dataframe to be used for statistical analysis
- scipy.stats: We use scipy for many of the statistical models. For instance using zscores in order to detect outliers
- matplotlib: We use matplotlib for visualization such as histograms
- statsmodels.api: We use statsmodel to run our multiple linear regression model

## Import Libraries

```
In [2]: # import the libraries
import pandas as pd
from pandas import DataFrame
import numpy as np
import scipy.stats as stats
import matplotlib.pyplot as plt
import statsmodels.api as sm
from statsmodels.formula.api import ols
from statsmodels.compat import lzip
```

## B3. Justification of Using Regression

Multiple linear regression is an appropriate technique to use for this analysis because the research question focuses on quantitative variables. My independent and dependent variables are all quantitative so this an ideal method to utilize. If I were to ask a question about categorical variables, then a more ideal model to use would be the logistic regression model.

## C1. Data Cleaning

For my data cleaning, I am going to start by focusing on null data, outliers, and duplicates. We can start by importing our data from a CSV. I am also going to drop all the categorical

variables since we are only focusing on quantitative variables for this analysis. This means dropping Customer\_id, Interaction, UID, City, State, County, Zip, Lat, Lng, Area, TimeZone, Job, Marital, Gender, Churn, Techie, Contract, Port\_modem, Tablet, InternetService, Phone, Multiple, OnlineSecurity, OnlineBackup, DeviceProtection, TechSupport, StreamingTV, StreamingMovies, PaperlessBilling, and PaymentMethod.

```
In [3]: df = pd.read_csv('churn_clean.csv')
```

```
In [4]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 10000 entries, 0 to 9999
Data columns (total 50 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   CaseOrder                            10000 non-null  int64
1   Customer_id                          10000 non-null  object
2   Interaction                           10000 non-null  object
3   UID                                  10000 non-null  object
4   City                                 10000 non-null  object
5   State                                10000 non-null  object
6   County                               10000 non-null  object
7   Zip                                  10000 non-null  int64
8   Lat                                  10000 non-null  float64
9   Lng                                  10000 non-null  float64
10  Population                           10000 non-null  int64
11  Area                                 10000 non-null  object
12  TimeZone                             10000 non-null  object
13  Job                                  10000 non-null  object
14  Children                             10000 non-null  int64
15  Age                                  10000 non-null  int64
16  Income                               10000 non-null  float64
17  Marital                              10000 non-null  object
18  Gender                               10000 non-null  object
19  Churn                                10000 non-null  object
20  Outage_sec_perweek                   10000 non-null  float64
21  Email                                10000 non-null  int64
22  Contacts                             10000 non-null  int64
23  Yearly_equip_failure                 10000 non-null  int64
24  Techie                               10000 non-null  object
25  Contract                             10000 non-null  object
26  Port_modem                           10000 non-null  object
27  Tablet                               10000 non-null  object
28  InternetService                     10000 non-null  object
29  Phone                                10000 non-null  object
30  Multiple                             10000 non-null  object
31  OnlineSecurity                       10000 non-null  object
32  OnlineBackup                         10000 non-null  object
33  DeviceProtection                     10000 non-null  object
34  TechSupport                          10000 non-null  object
35  StreamingTV                          10000 non-null  object
36  StreamingMovies                      10000 non-null  object
37  PaperlessBilling                     10000 non-null  object
38  PaymentMethod                        10000 non-null  object
39  Tenure                               10000 non-null  float64
40  MonthlyCharge                        10000 non-null  float64
41  Bandwidth_GB_Year                   10000 non-null  float64
42  Item1                                10000 non-null  int64
43  Item2                                10000 non-null  int64
44  Item3                                10000 non-null  int64
45  Item4                                10000 non-null  int64
46  Item5                                10000 non-null  int64
47  Item6                                10000 non-null  int64
48  Item7                                10000 non-null  int64
49  Item8                                10000 non-null  int64
dtypes: float64(7), int64(16), object(27)
memory usage: 3.8+ MB
```

```
In [5]: dfq = df.drop(['Customer_id', 'Interaction', 'UID', 'City', 'State', 'County', 'Zip'])
dfq.info()
```

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 10000 entries, 0 to 9999
Data columns (total 20 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   CaseOrder                             10000 non-null  int64
1   Population                             10000 non-null  int64
2   Children                               10000 non-null  int64
3   Age                                    10000 non-null  int64
4   Income                                 10000 non-null  float64
5   Outage_sec_perweek                     10000 non-null  float64
6   Email                                  10000 non-null  int64
7   Contacts                               10000 non-null  int64
8   Yearly equip_failure                   10000 non-null  int64
9   Tenure                                 10000 non-null  float64
10  MonthlyCharge                          10000 non-null  float64
11  Bandwidth_GB_Year                      10000 non-null  float64
12  Item1                                  10000 non-null  int64
13  Item2                                  10000 non-null  int64
14  Item3                                  10000 non-null  int64
15  Item4                                  10000 non-null  int64
16  Item5                                  10000 non-null  int64
17  Item6                                  10000 non-null  int64
18  Item7                                  10000 non-null  int64
19  Item8                                  10000 non-null  int64
dtypes: float64(5), int64(15)
memory usage: 1.5 MB

```

## Treating Nulls

We begin by checking the dataframe for nulls. We can use `.isnull().sum()` to look through the variables and see if there is any missing data. Using this function we can see that there are no nulls present in the data. Another thing I would like to check is population. This is because for a value like this, it cannot be 0 since it should count the customer. Using the `nsmallest()` function, we can see that zeroes do exist within the data. I would like to drop those zeroes and replace it with the median as the distribution is skewed right. We determine the distribution by creating a histogram of population. After dropping all the zero values from population and replacing them with median, we can see our minimum is no longer zero.

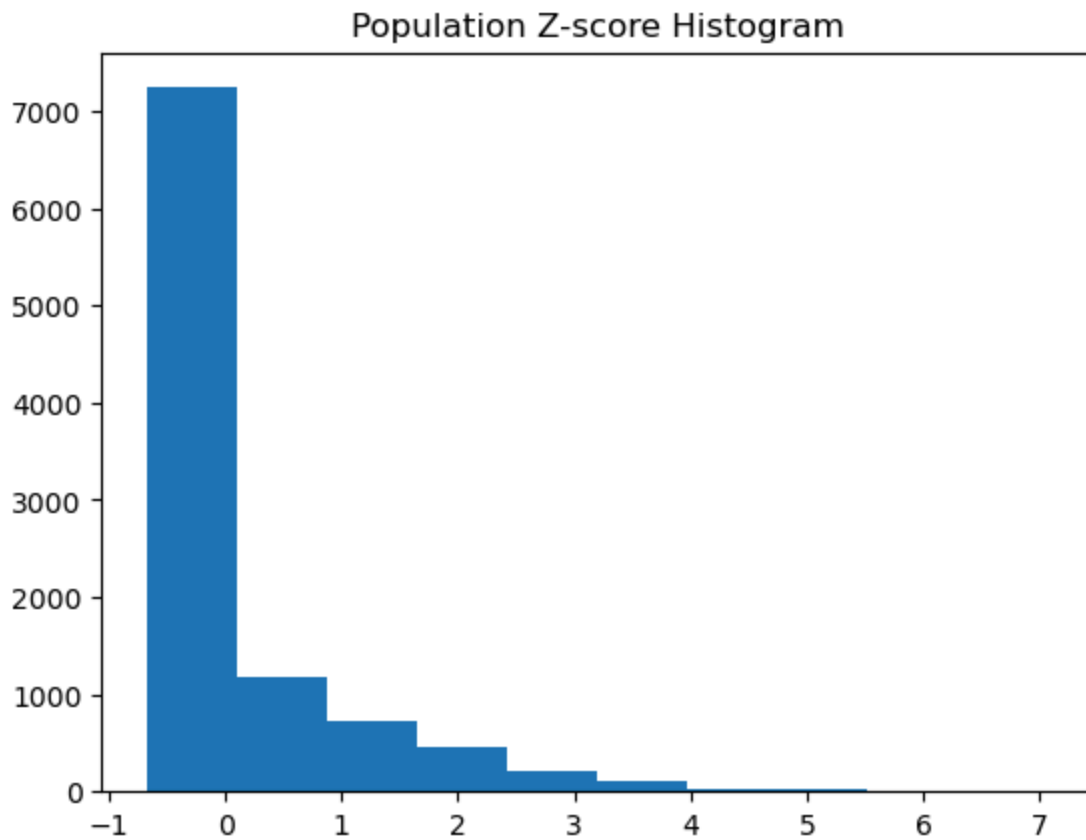
```
In [6]: dfq.isnull().sum()
```

```
Out[6]: CaseOrder      0
Population    0
Children      0
Age           0
Income        0
Outage_sec_perweek  0
Email         0
Contacts      0
Yearly_equip_failure  0
Tenure        0
MonthlyCharge 0
Bandwidth_GB_Year 0
Item1         0
Item2         0
Item3         0
Item4         0
Item5         0
Item6         0
Item7         0
Item8         0
dtype: int64
```

```
In [7]: # Check the poulation for zeroes
dfq.Population.nsmallest(n=10)
```

```
Out[7]: 13      0
422     0
428     0
434     0
446     0
682     0
694     0
719     0
814     0
839     0
Name: Population, dtype: int64
```

```
In [8]: # create hist for population
dfq['zscore'] = stats.zscore(dfq['Population'])
plt.hist(dfq['zscore'])
plt.title('Population Z-score Histogram')
plt.show()
```



```
In [9]: # drop all zeroes
dfq['Population'] = np.where(dfq['Population'] == 0, np.nan, dfq['Population'])
# fill with median as it is skewed right
dfq['Population'] = dfq['Population'].fillna(dfq['Population'].median())
```

```
In [10]: # Check the poulation for zeroes
dfq.Population.nsmallest(n=10)
```

```
Out[10]: 4453    2.0
261      4.0
3475     4.0
6018     4.0
2613     5.0
2092     6.0
2192     6.0
5054     6.0
5149     6.0
6048     6.0
Name: Population, dtype: float64
```

```
In [11]: #drop zscore
dfq = dfq.drop(['zscore'],axis=1)
```

## Finding & Treating Duplicates

Next, we will check to see if there are any duplicates in the data. We can do this by using `.duplicated().value_counts()` which will output a true or false depending on whether or not duplicates exist within the dataframe. We can see from the output of false 10,000 times that there are no duplicates within the data.

```
In [12]: dfq.duplicated().value_counts()
```

```
Out[12]: False      10000  
dtype: int64
```

## Finding & Treating Outliers

We can start by checking the histograms of all of our quantitative variables. After looking through it, the distributions are as follows:

- 'CaseOrder' - normal
- 'Population' - skewed right
- 'Children' - skewed right
- 'Age' - uniform
- 'Income' - skewed right
- 'Outage\_sec\_perweek' - normal
- 'Email' - normal
- 'Contacts' - skewed right
- 'Yearly\_equip\_failure' - skewed right
- 'Tenure' - bimodal
- 'MonthlyCharge' - normal
- 'Bandwidth\_GB\_Year' - bimodal
- 'Item1' - normal
- 'Item2' - normal
- 'Item3' - normal
- 'Item4' - normal
- 'Item5' - normal
- 'Item6' - normal
- 'Item7' - normal
- 'Item8' - normal

This is useful information to note for later. We can also identify from our histograms if the data passes 3 standard deviations. I will use that as a cutoff for what we identify as outliers. Using this benchmark, the following variables contain outliers:

- Population, Children, Income, Outage\_sec\_perweek, Email, Contacts, Yearly\_equip\_failure, Item1, Item2, Item3, Item4, Item5, Item6, Item7, and Item8

Now that I know what variables are the ones that need to be solved, I can run a for loop to drop the outliers which are values equivalent to a z-score greater or less than 3 and -3 three respectively. We also need to know the distribution to understand what we need to impute these variables with. For population we impute with median since it is skewed right. For children, we use median since it's skewed right. For income we use median. For outage\_sec\_perweek we use mean since it is distributed normally. For Email we use mean. For Contacts we use median. For Yearly\_equip\_failure we use median. For item1 through item8 we use mean.

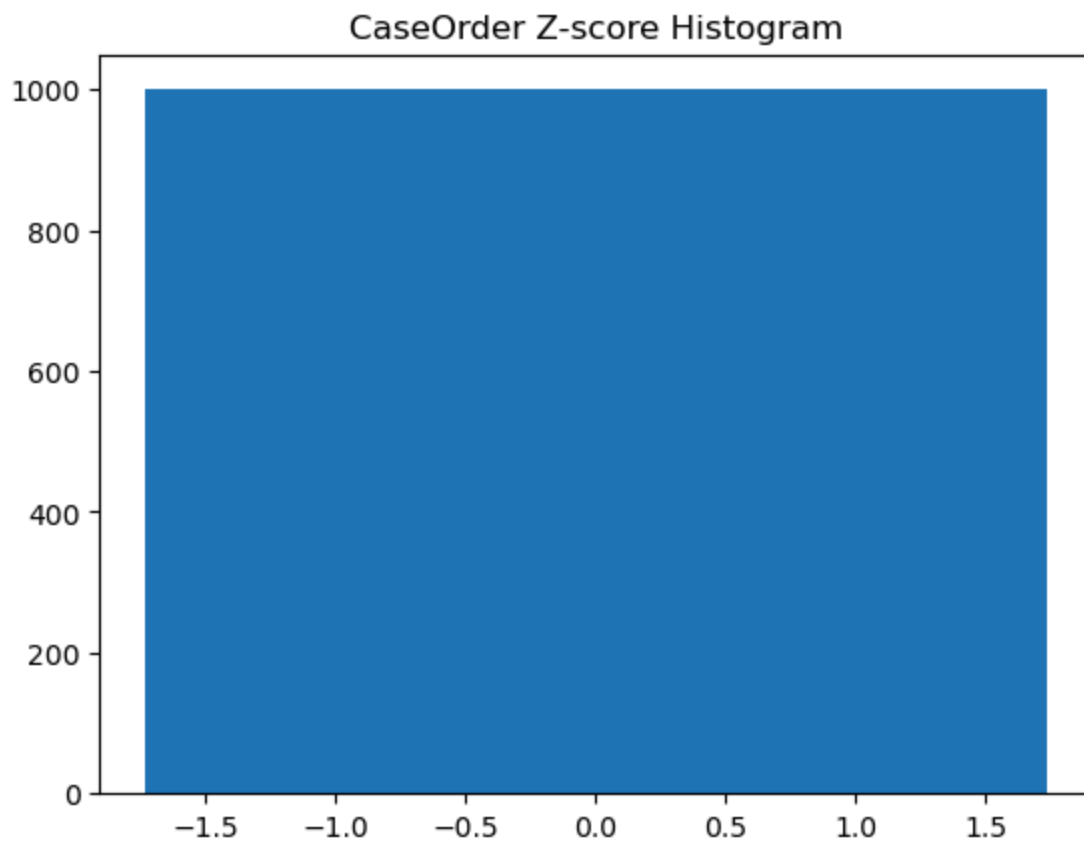


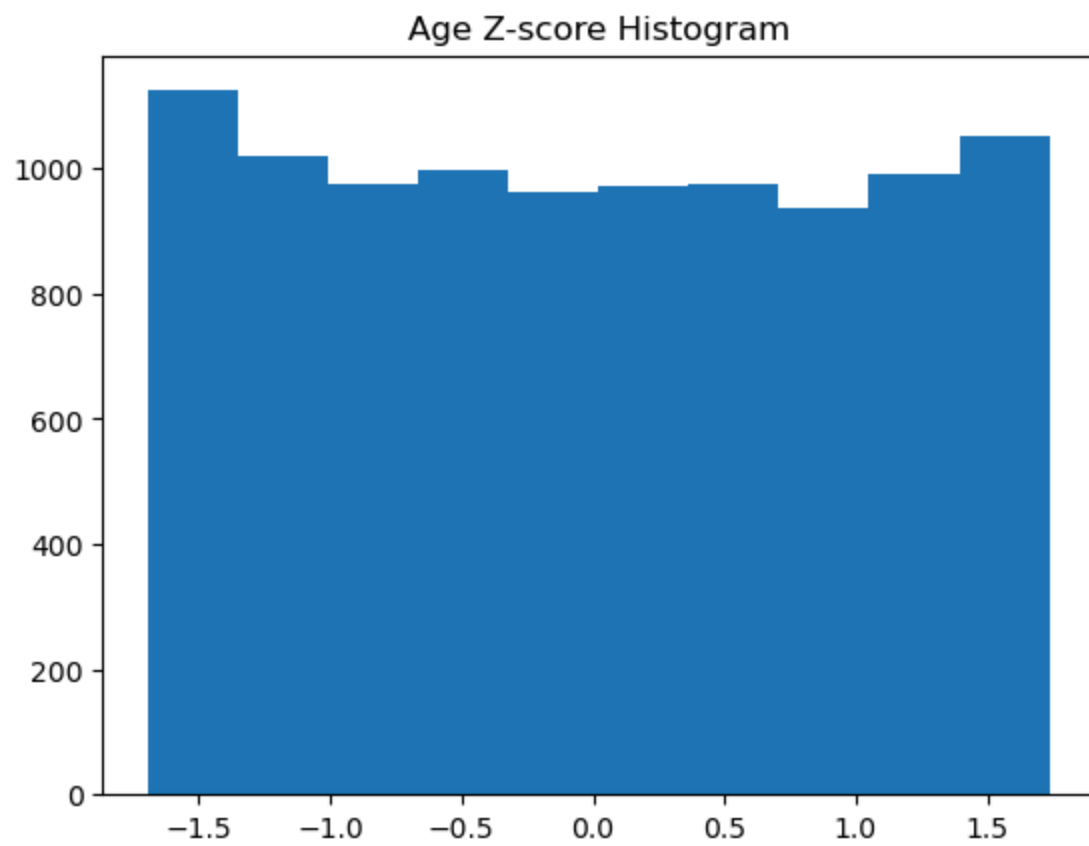
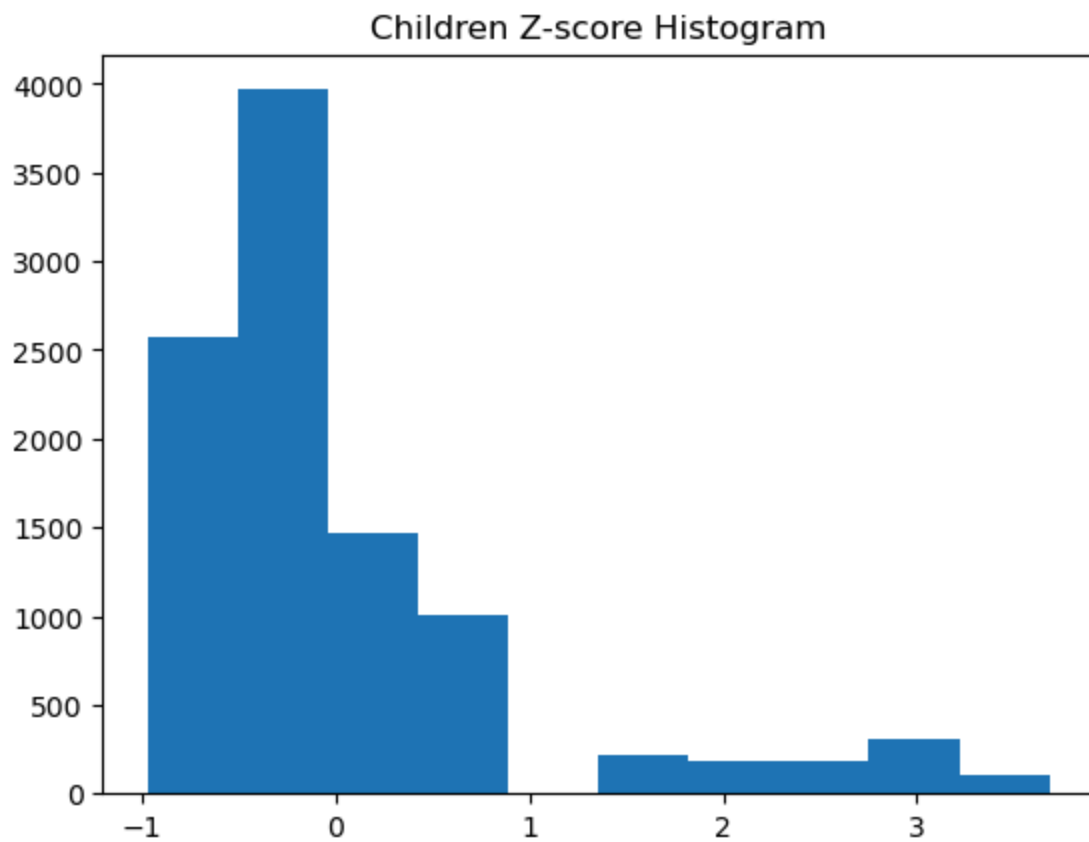
After the for loop runs for both median and for mean, we are able to see that the histograms are fixed and the outliers have been treated.

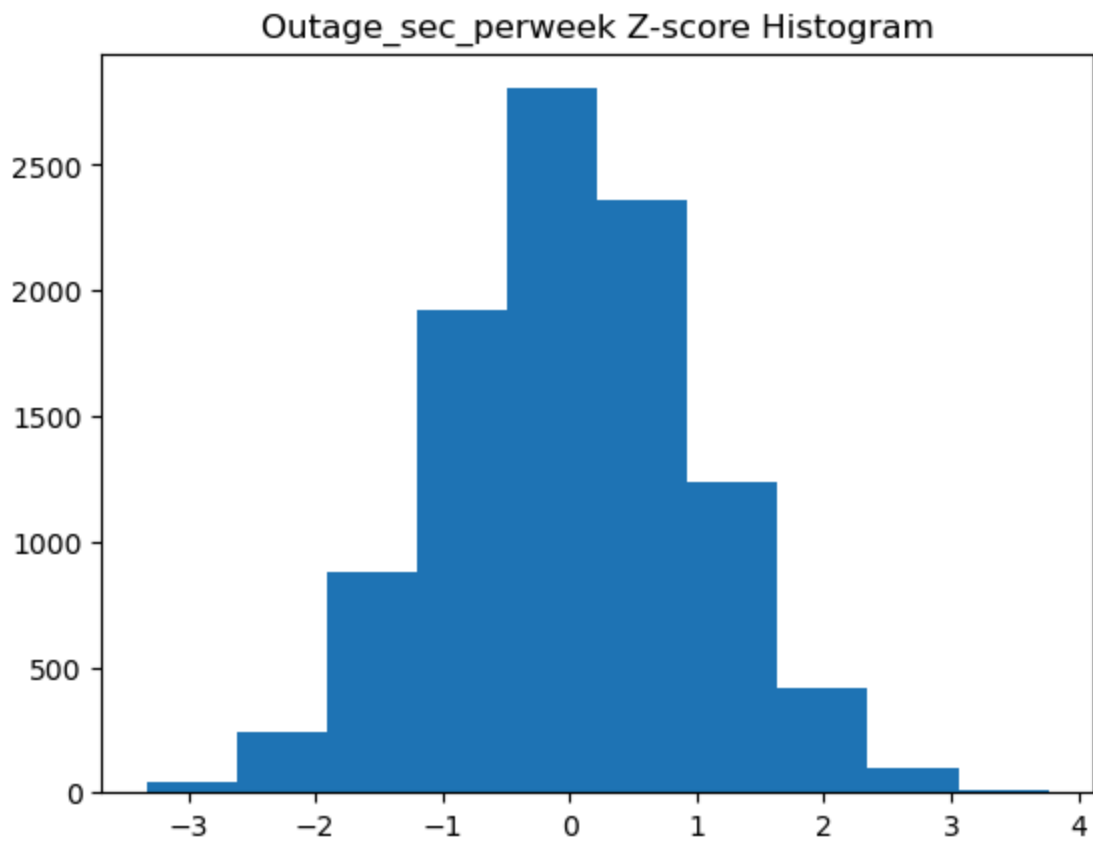
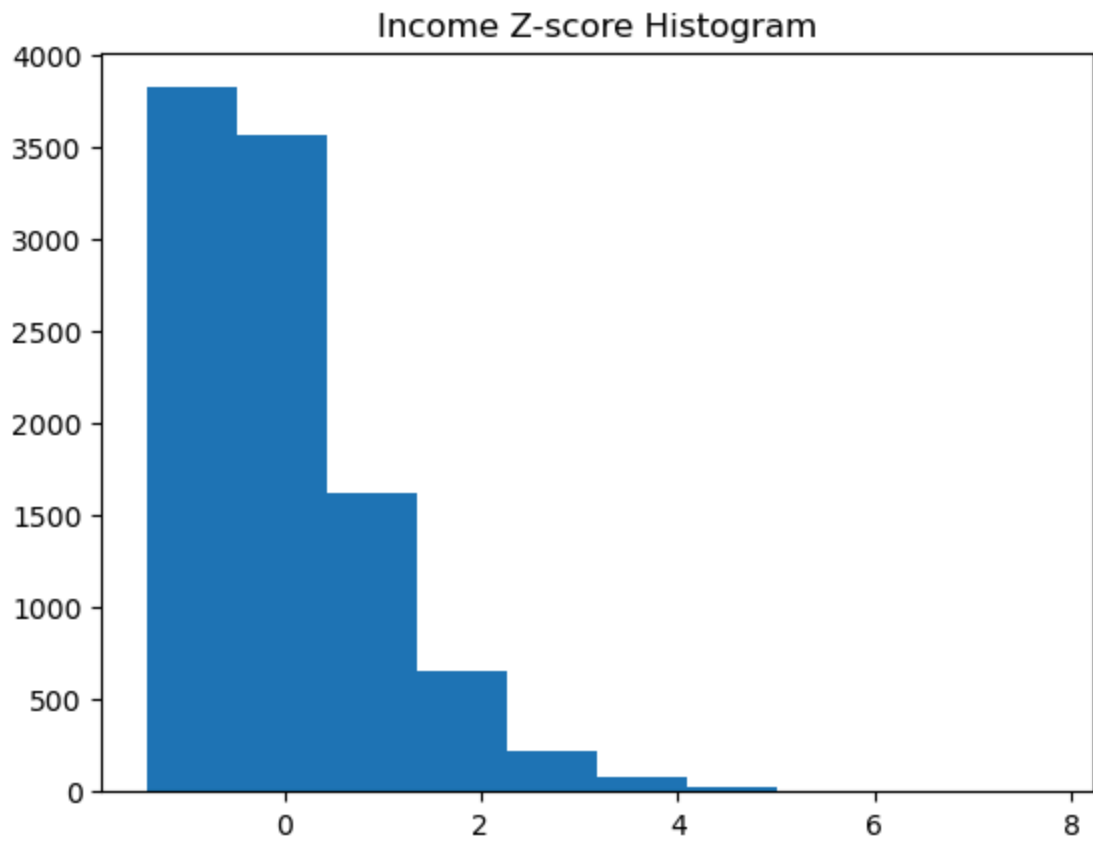
```
In [13]: # create a list of columns
dfq_c = dfq.columns.tolist()
dfq_c
```

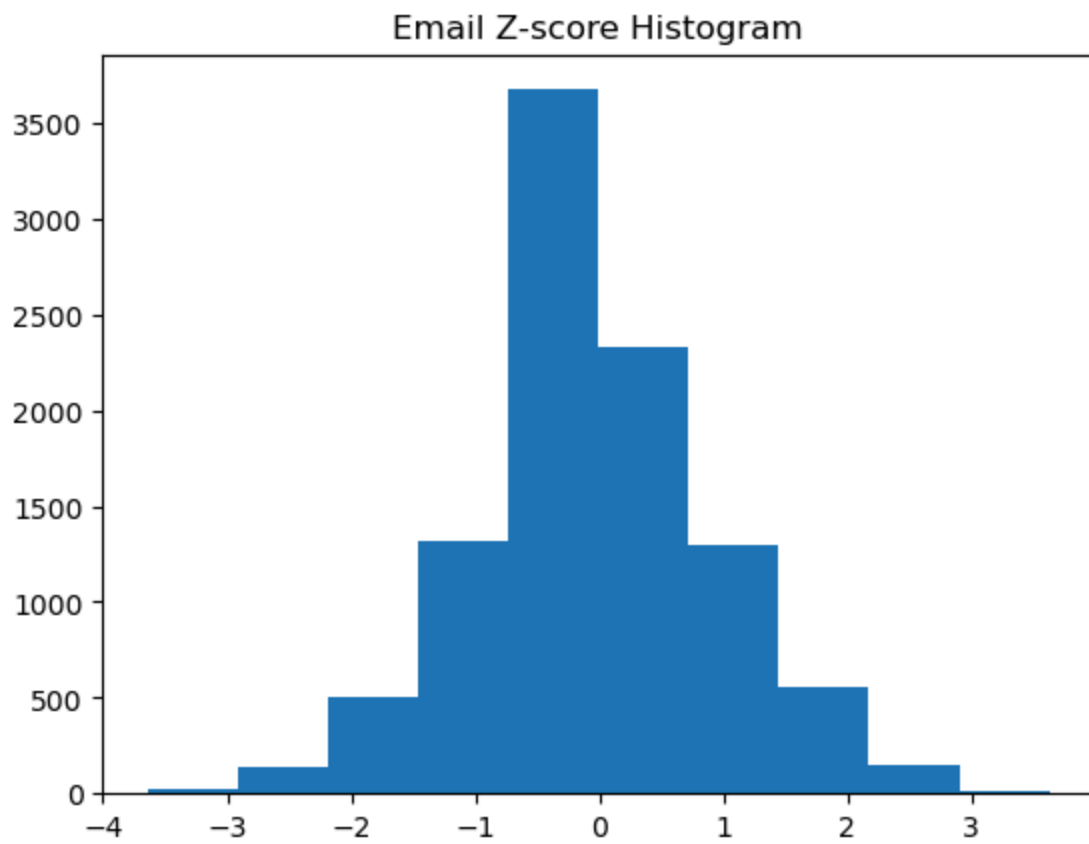
```
Out[13]: ['CaseOrder',
          'Population',
          'Children',
          'Age',
          'Income',
          'Outage_sec_perweek',
          'Email',
          'Contacts',
          'Yearly_equip_failure',
          'Tenure',
          'MonthlyCharge',
          'Bandwidth_GB_Year',
          'Item1',
          'Item2',
          'Item3',
          'Item4',
          'Item5',
          'Item6',
          'Item7',
          'Item8']
```

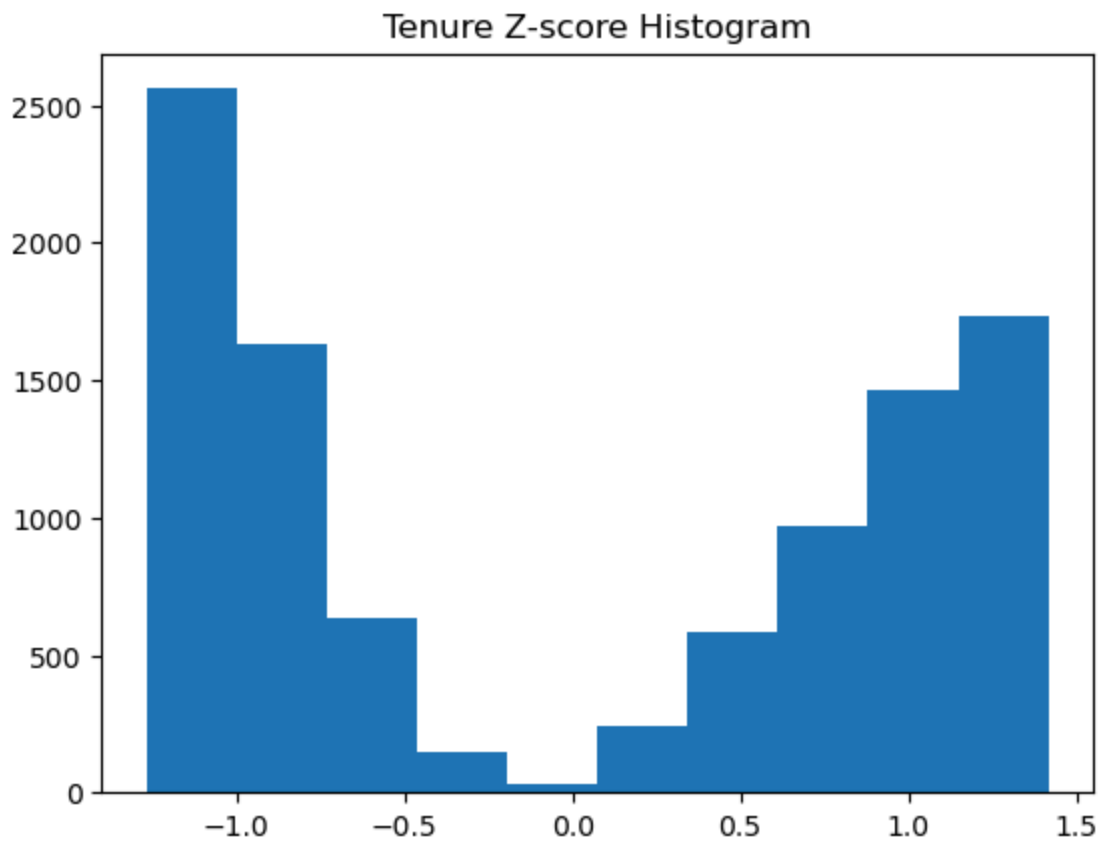
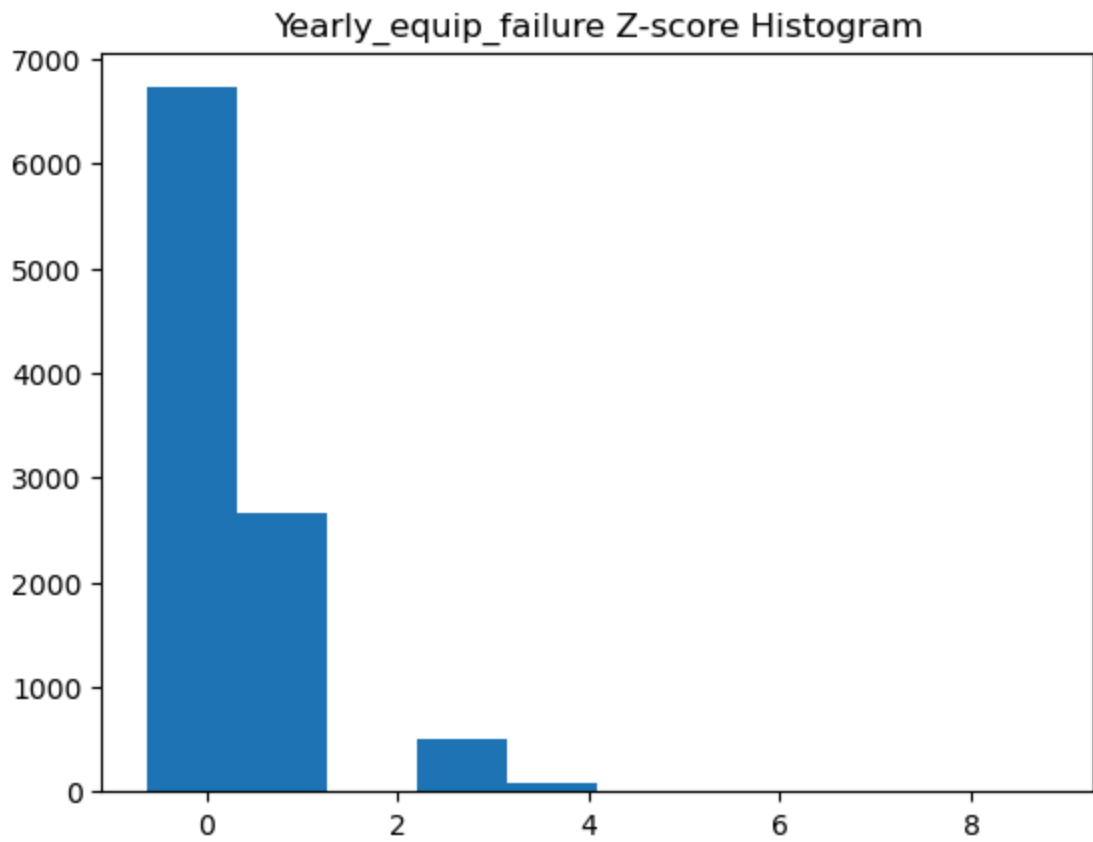
```
In [14]: for column in dfq_c:
          dfq['zscore'] = stats.zscore(dfq[column])
          plt.hist(dfq['zscore'])
          plt.title(column + ' Z-score Histogram')
          plt.show()
```

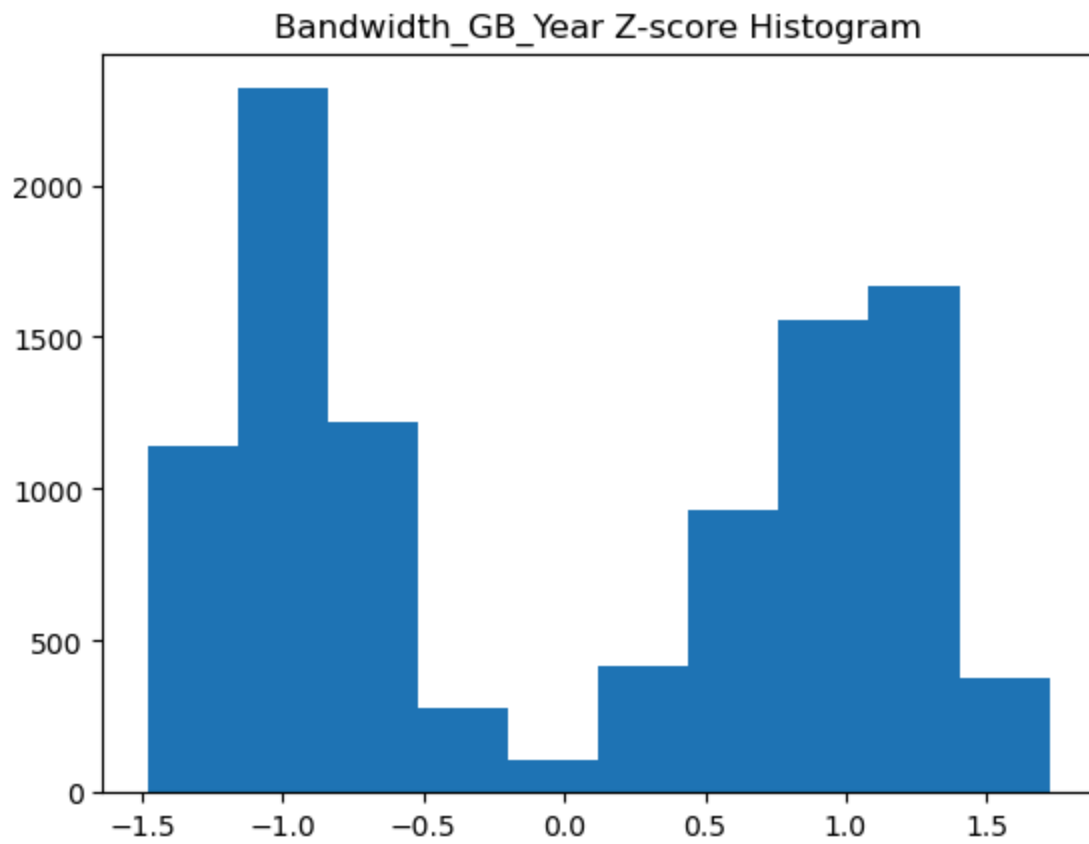
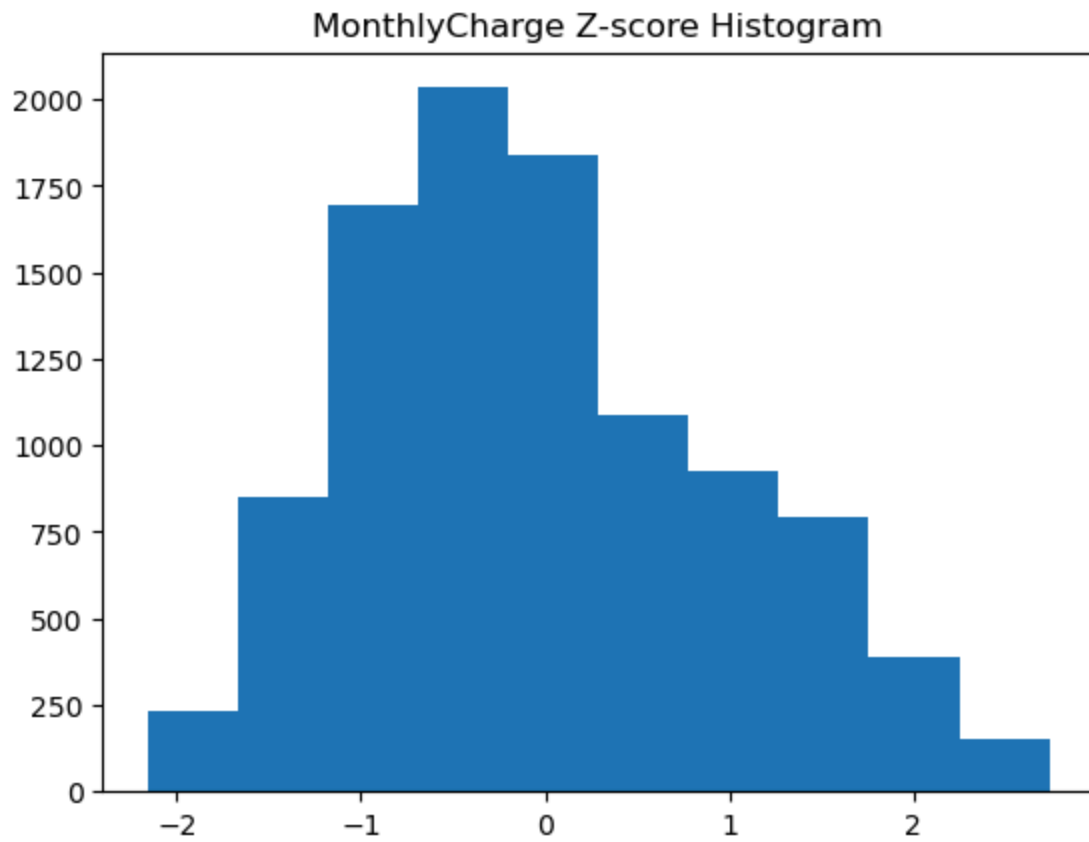


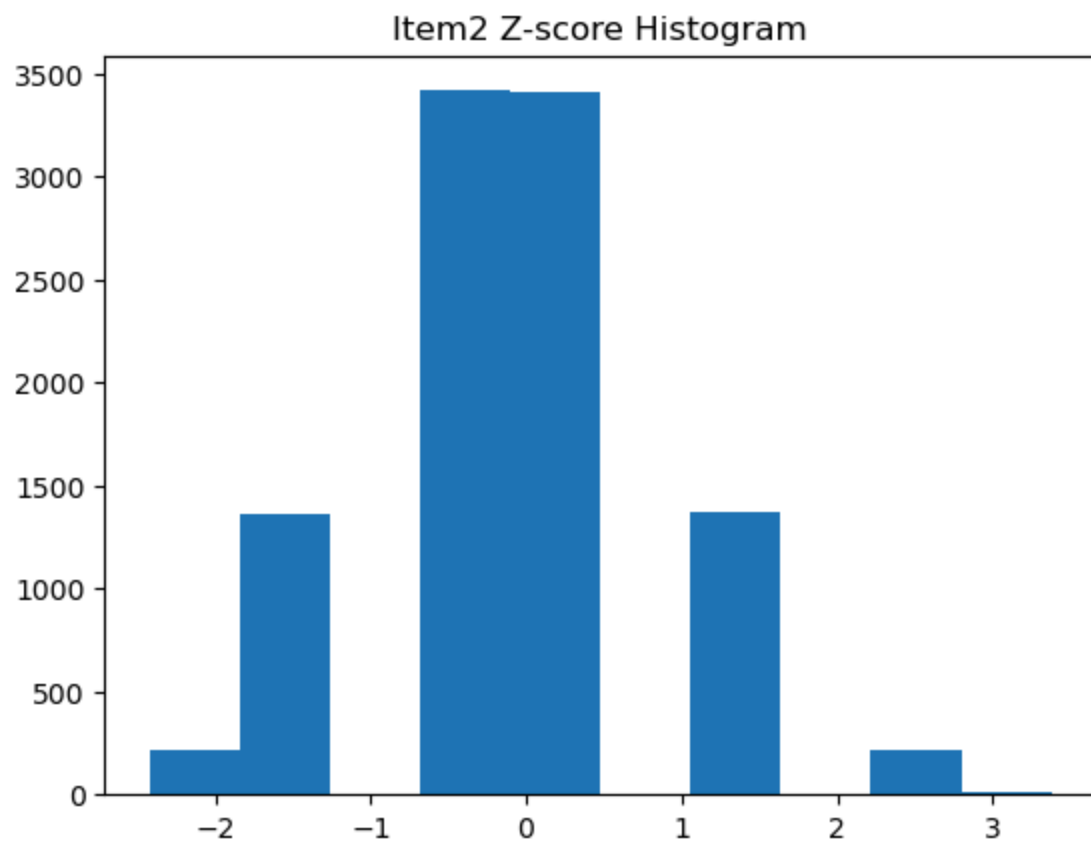
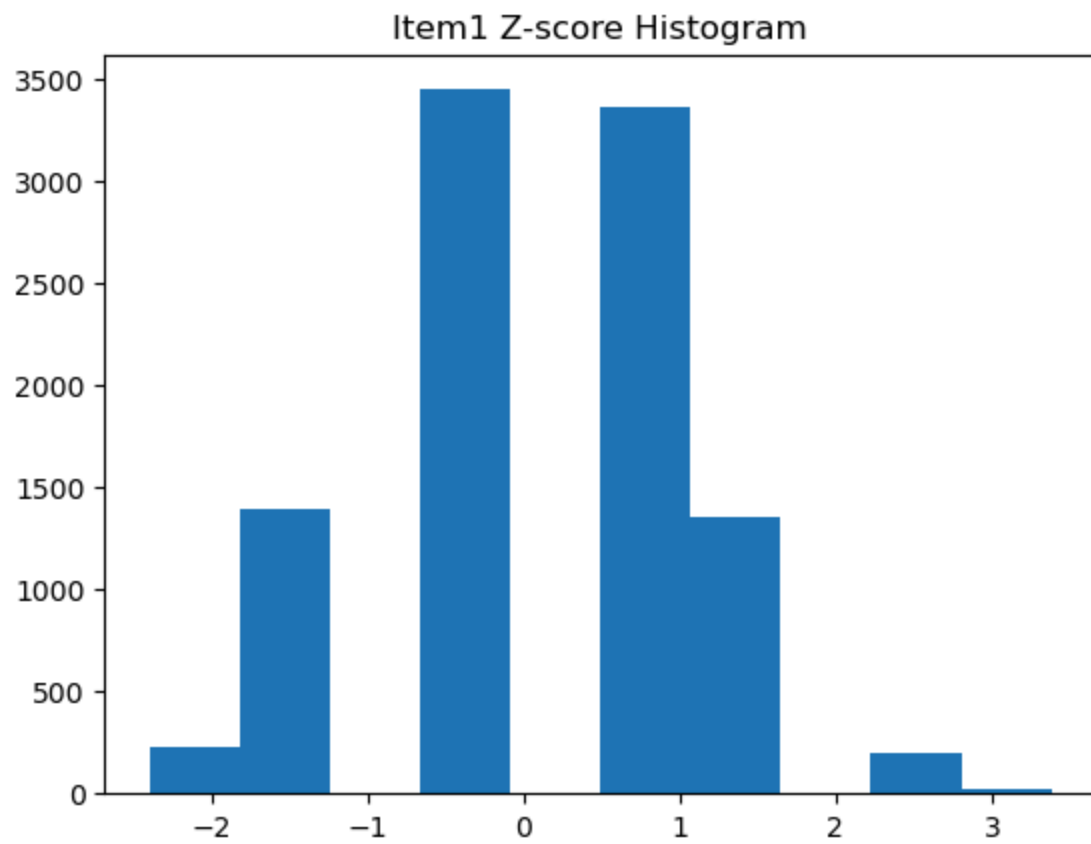




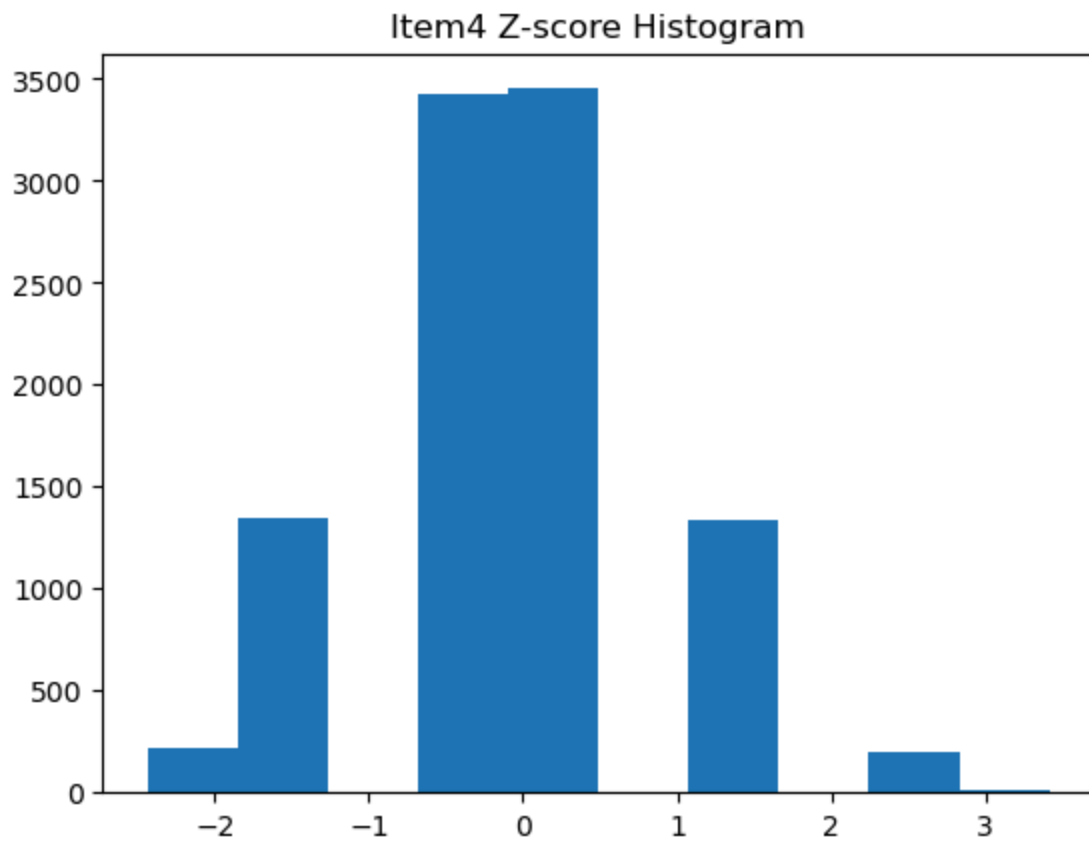
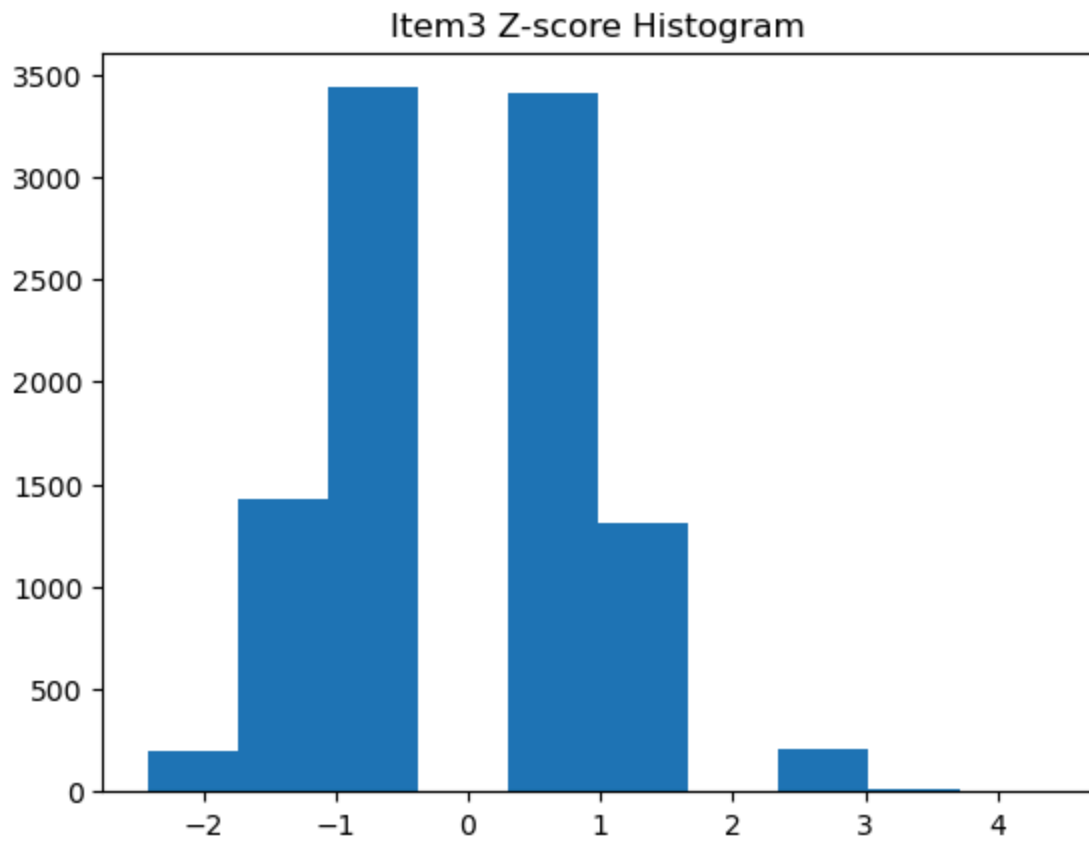


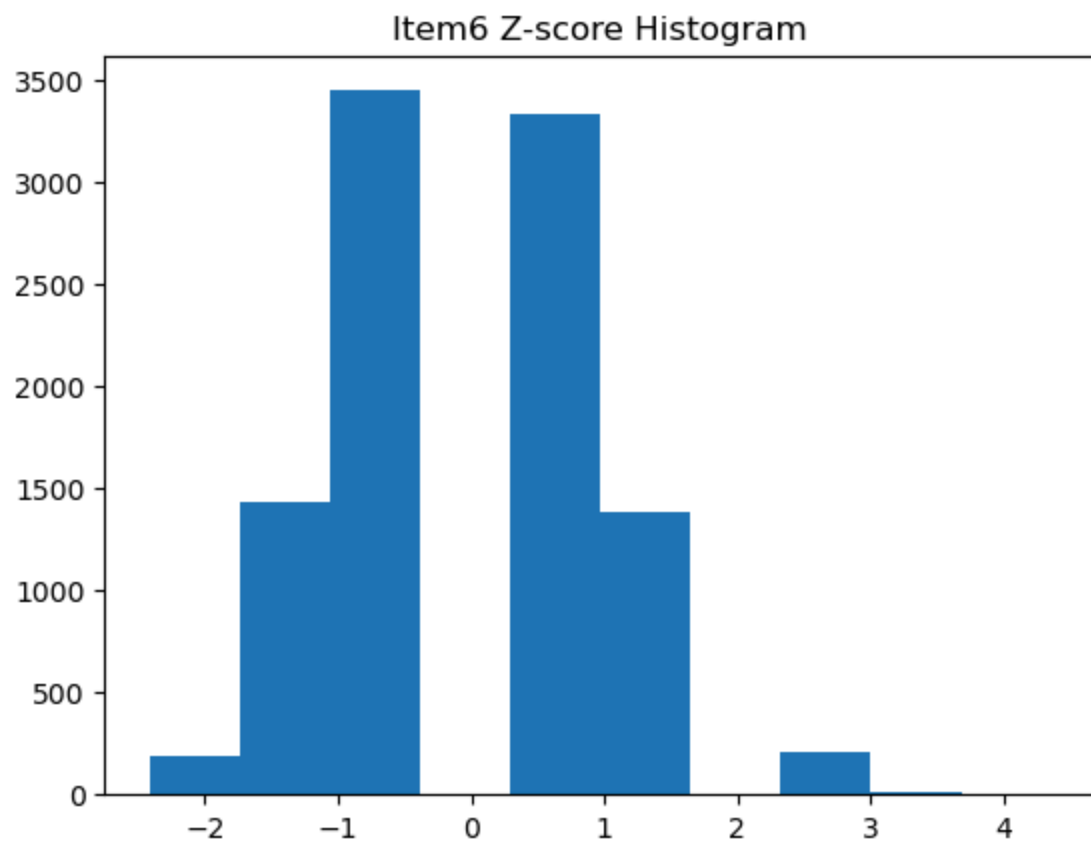
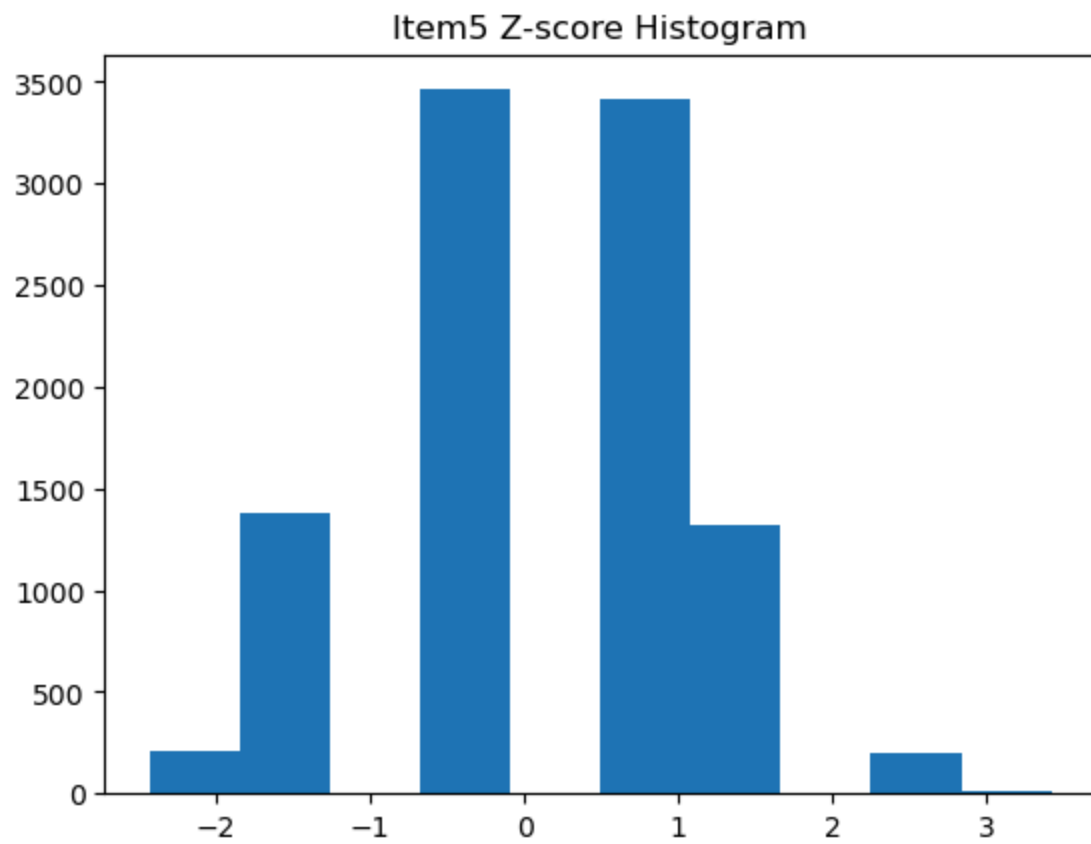


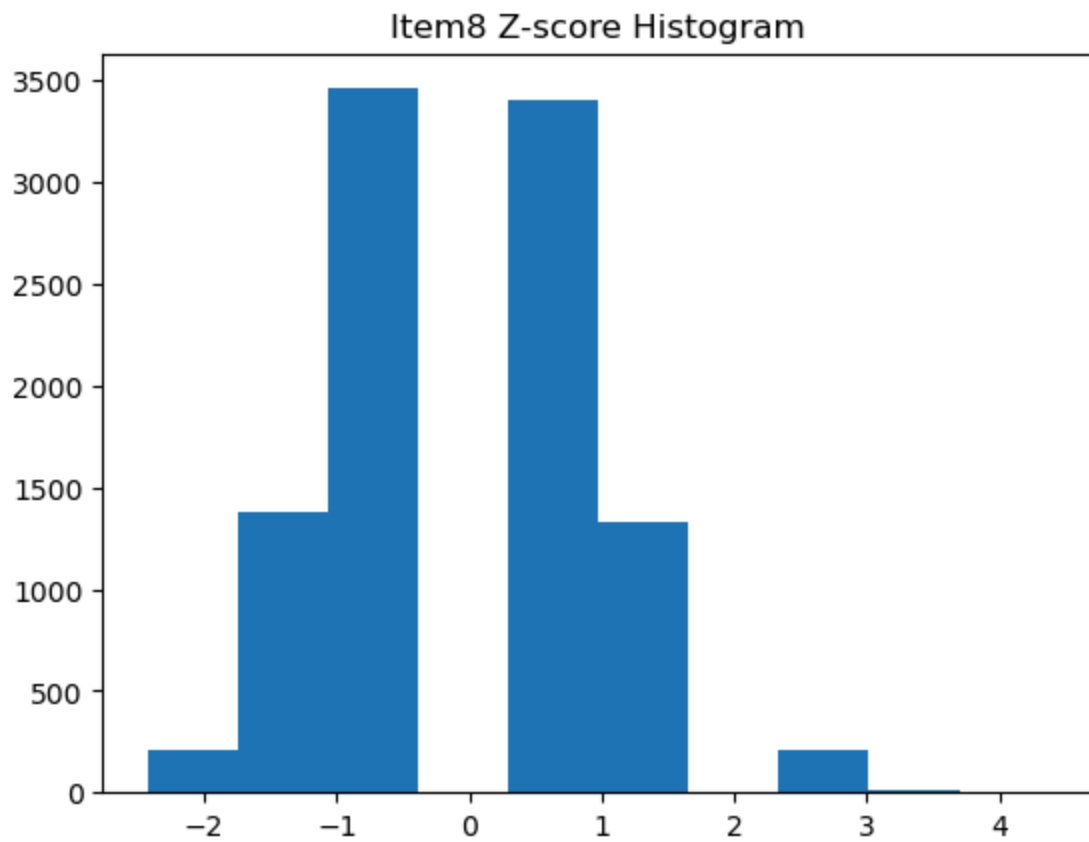
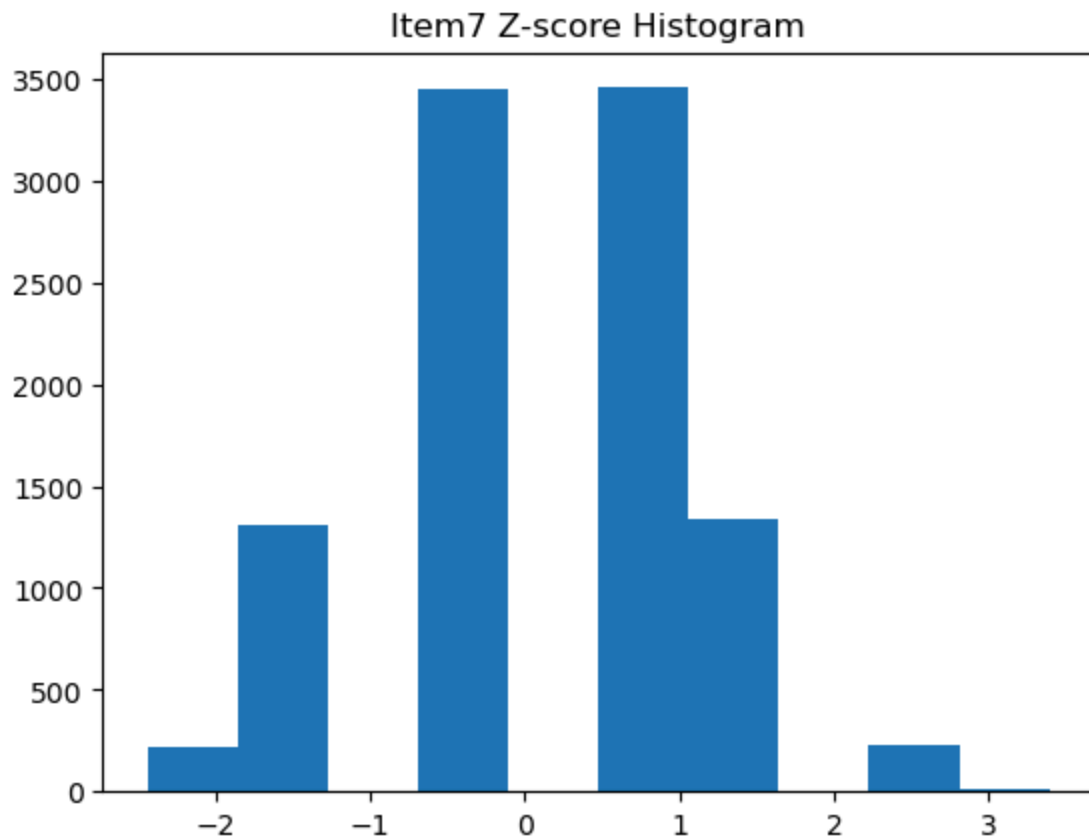












```
In [15]: # Run a for loop for all the identified variables that need to be fixed that w.  
dfq_z_median = ['Population', 'Children', 'Income', 'Contacts', 'Yearly equip_'  
dfq_z_median  
for column in dfq_z_median:  
    # create nulls for outliers in population
```

```

dfq['zscore'] = stats.zscore(dfq[column])
dfq[column] = np.where(dfq['zscore'] > 2, np.nan, dfq[column])
dfq[column] = np.where(dfq['zscore'] < -2, np.nan, dfq[column])
# use fillna function to impute outliers with median
dfq[column] = dfq[column].fillna(dfq[column].median())

```

```

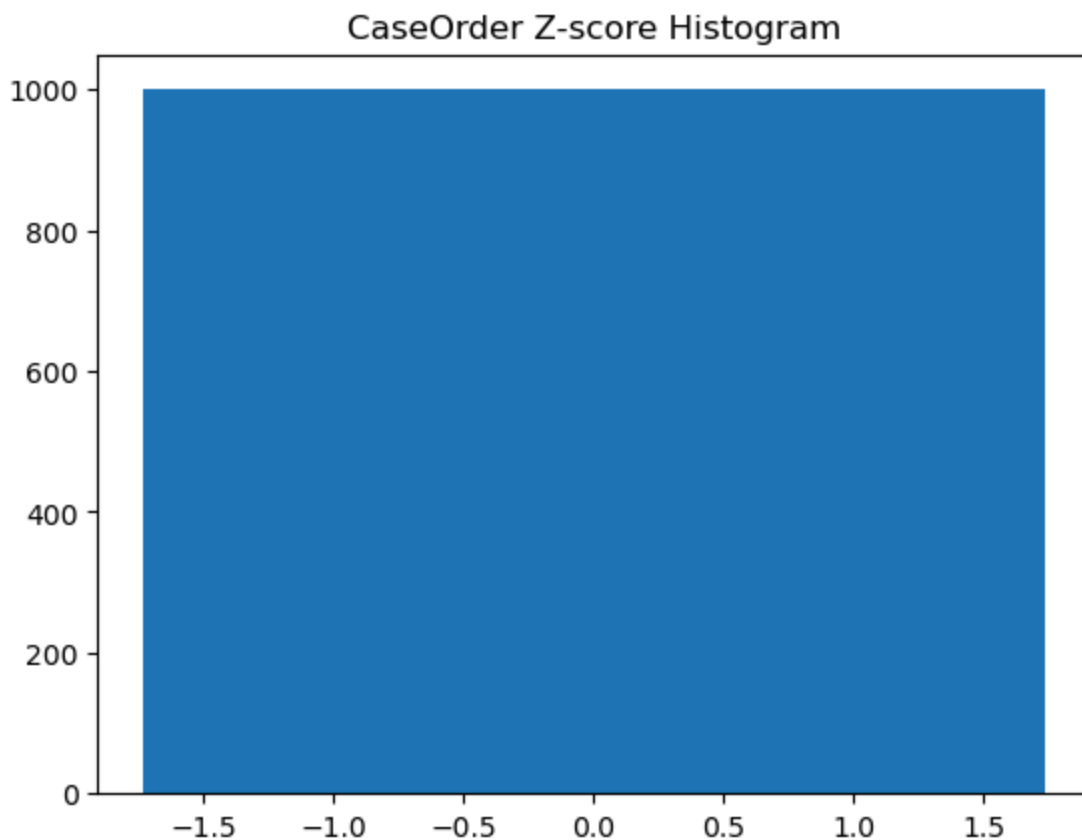
In [16]: # Run a for loop for all the identified variables that need to be fixed that w.
dfq_z_mean = ['Outage_sec_perweek', 'Email', 'Item1', 'Item2', 'Item3', 'Item4']
dfq_z_mean
for column in dfq_z_mean:
    # create nulls for outliers in population
    dfq['zscore'] = stats.zscore(dfq[column])
    dfq[column] = np.where(dfq['zscore'] > 2, np.nan, dfq[column])
    dfq[column] = np.where(dfq['zscore'] < -2, np.nan, dfq[column])
    # use fillna function to impute outliers with mean
    dfq[column] = dfq[column].fillna(dfq[column].mean())

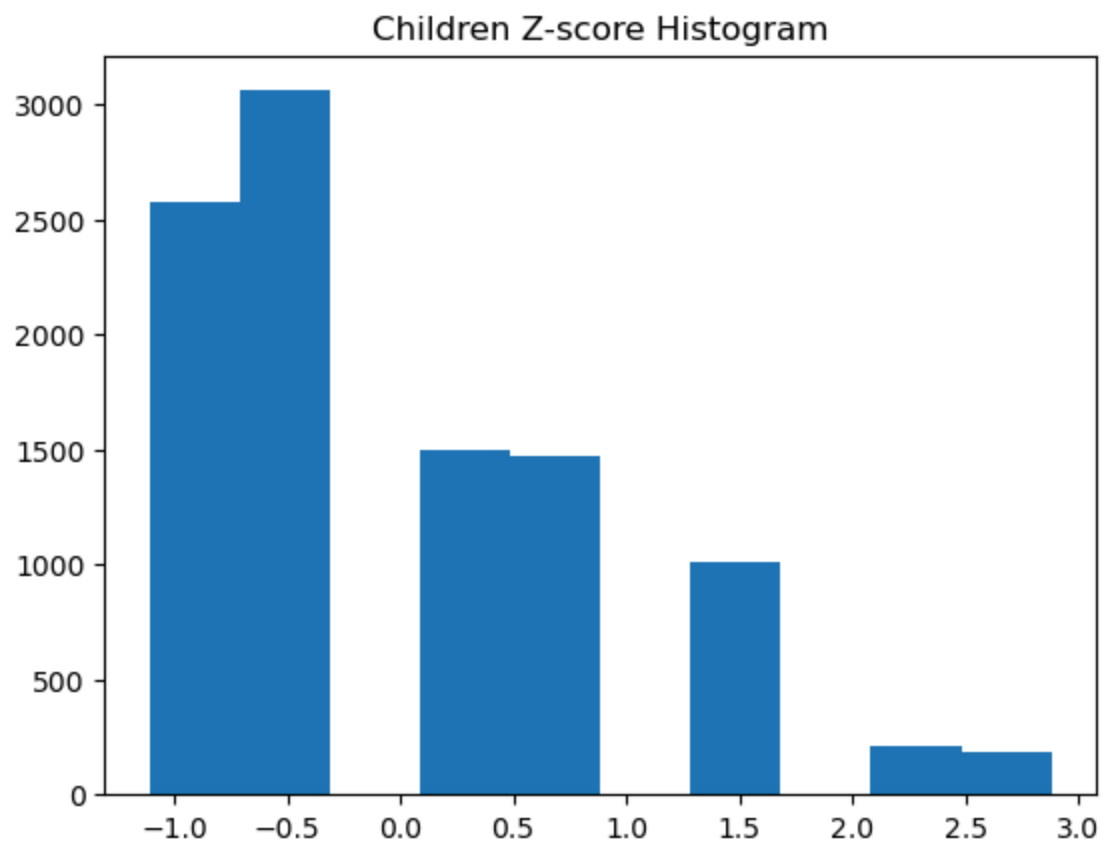
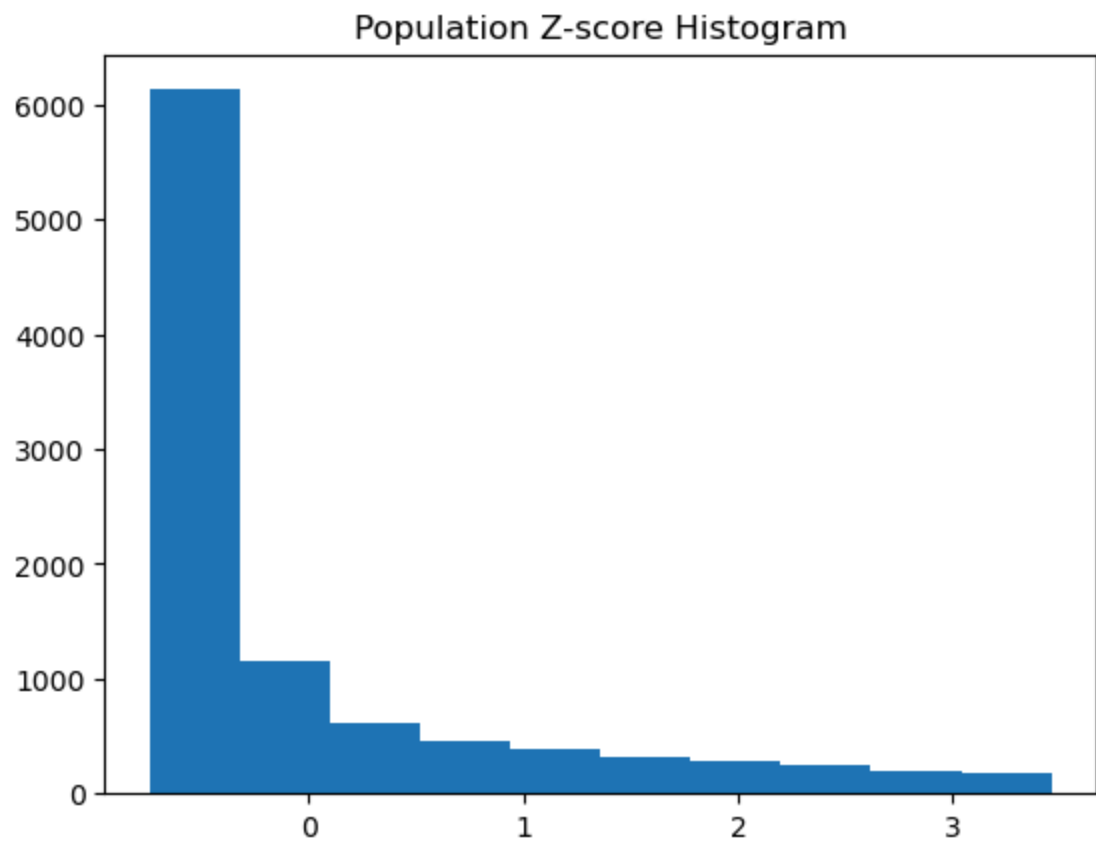
```

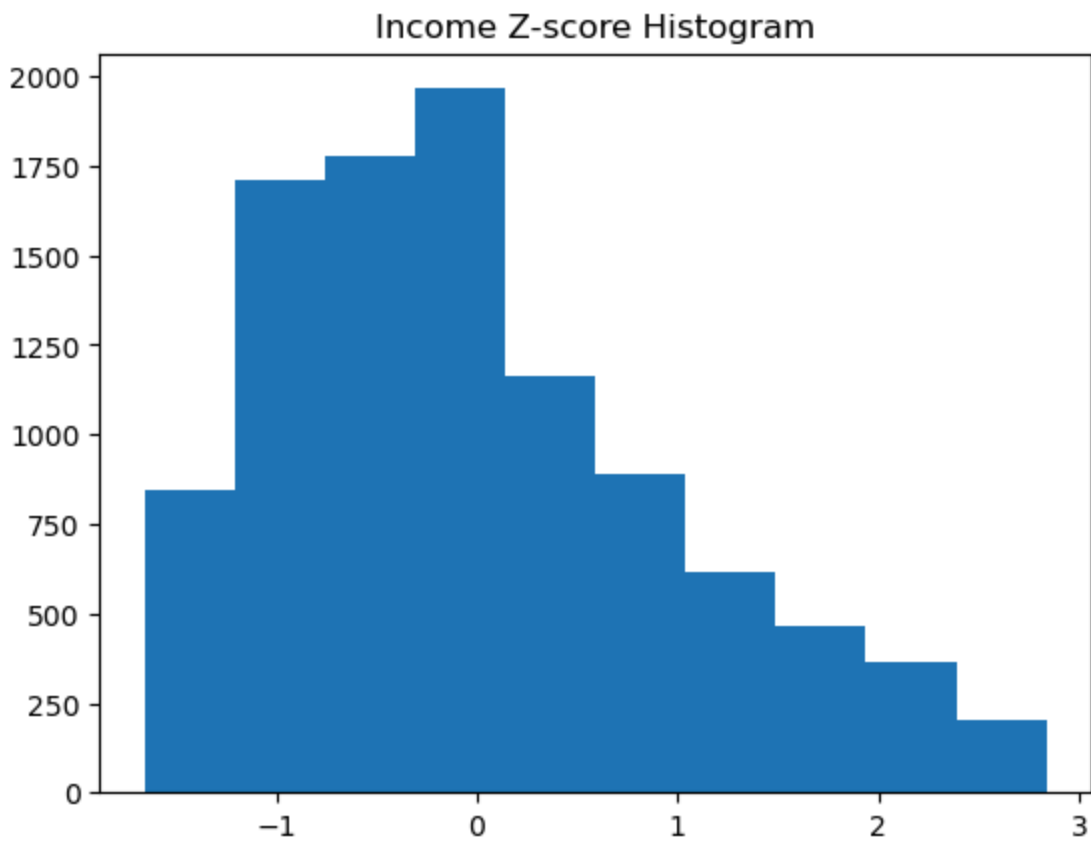
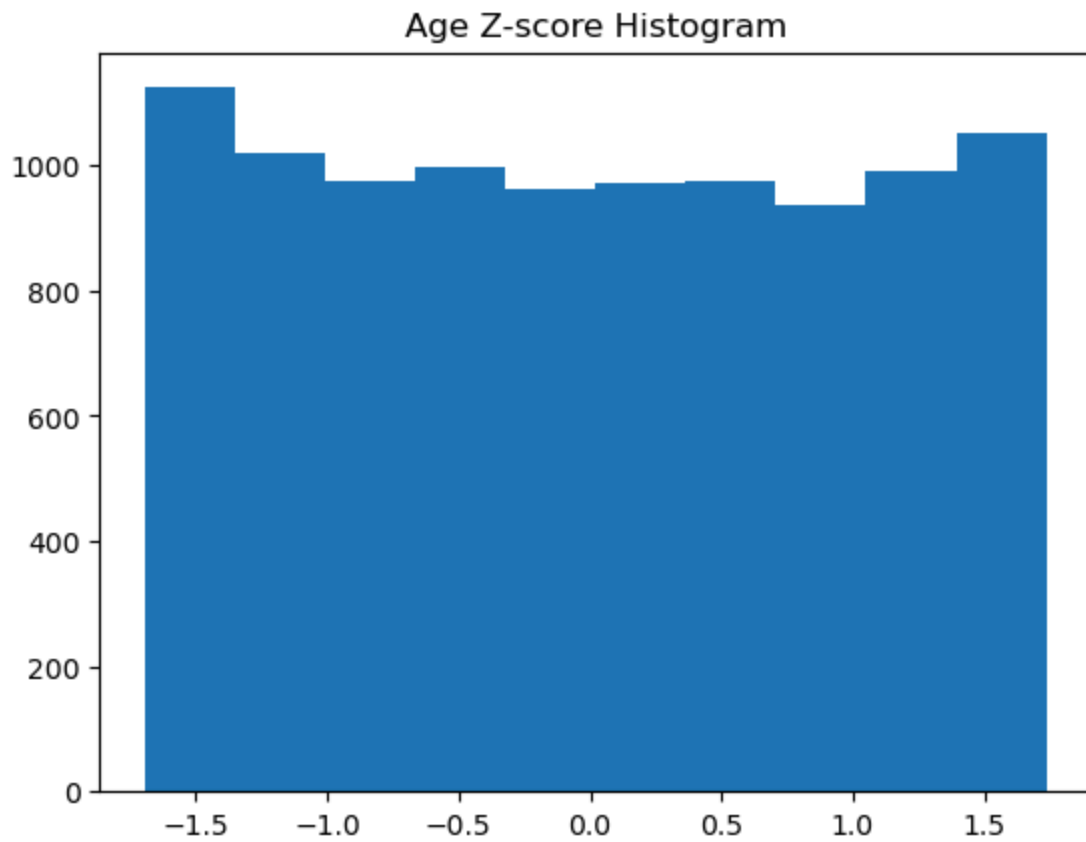
```

In [17]: # check new histograms
for column in dfq_c:
    dfq['zscore'] = stats.zscore(dfq[column])
    plt.hist(dfq['zscore'])
    plt.title(column + ' Z-score Histogram')
    plt.show()

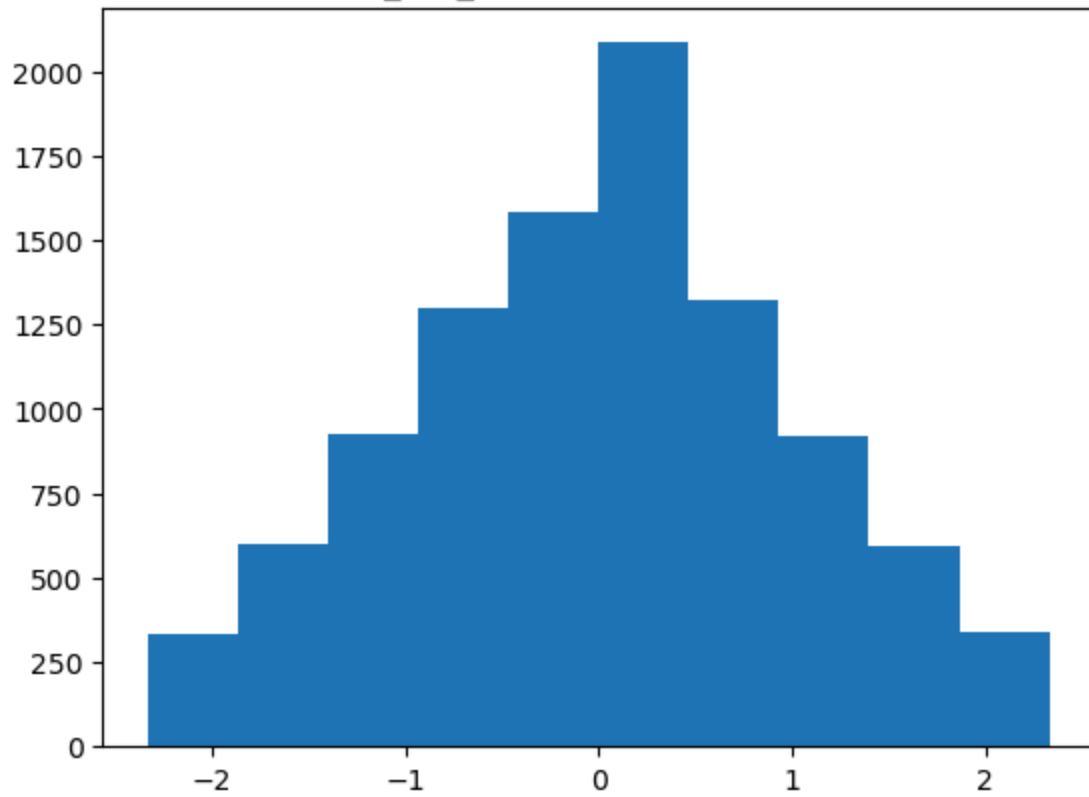
```



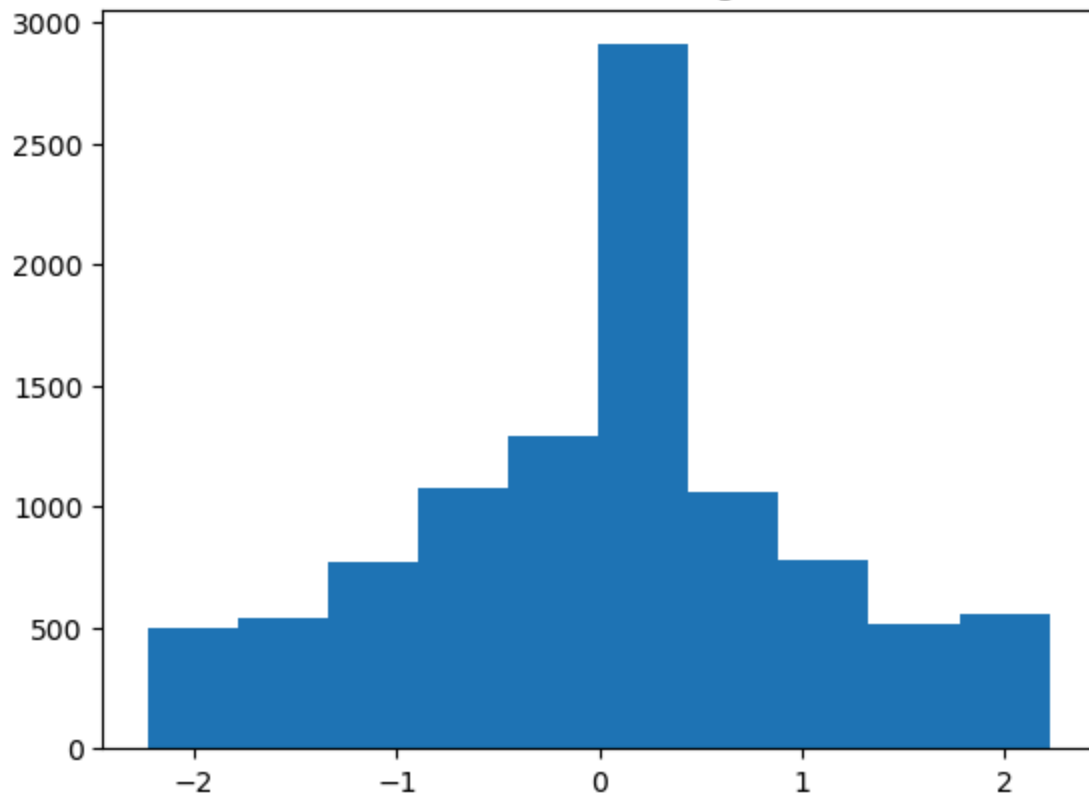


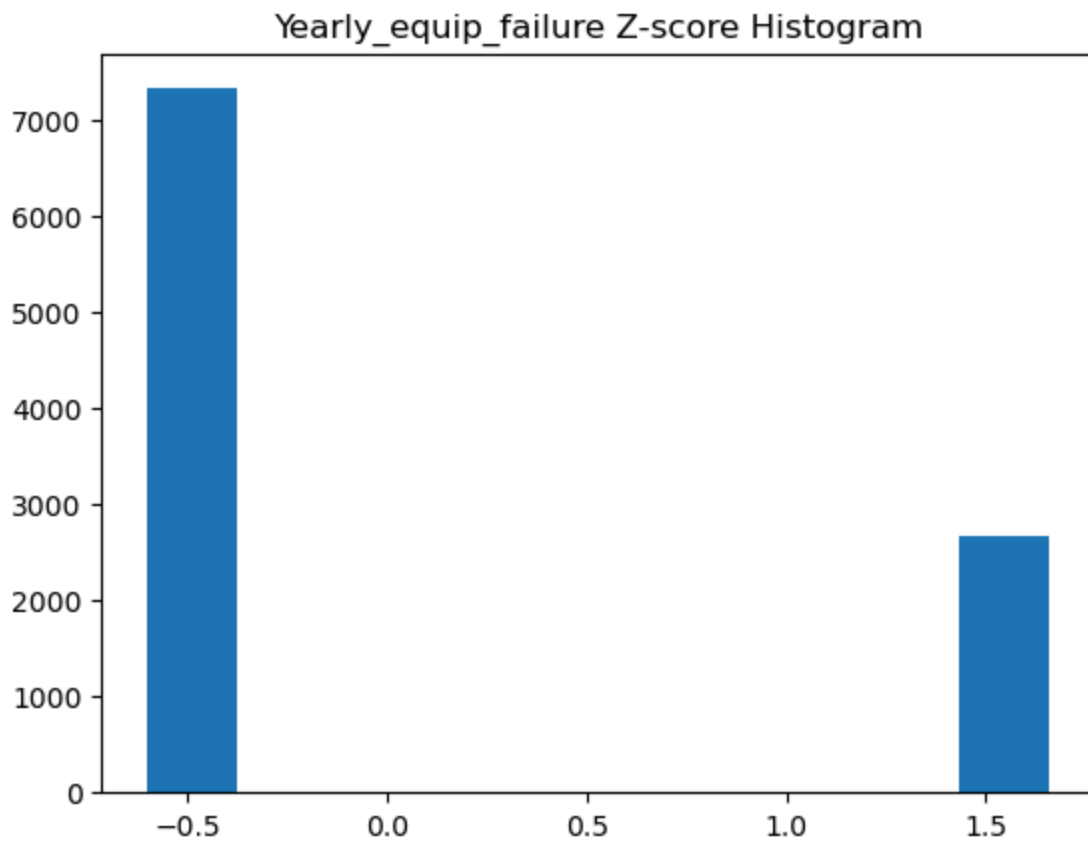
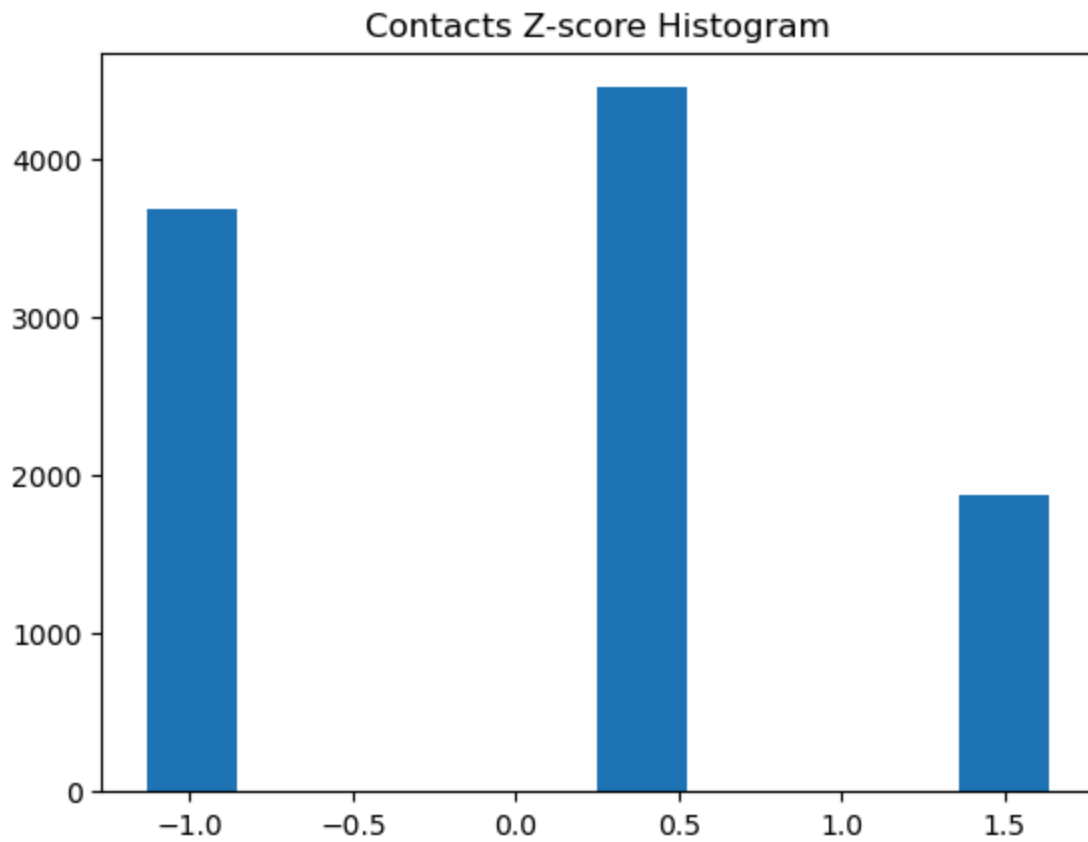


Outage\_sec\_perweek Z-score Histogram

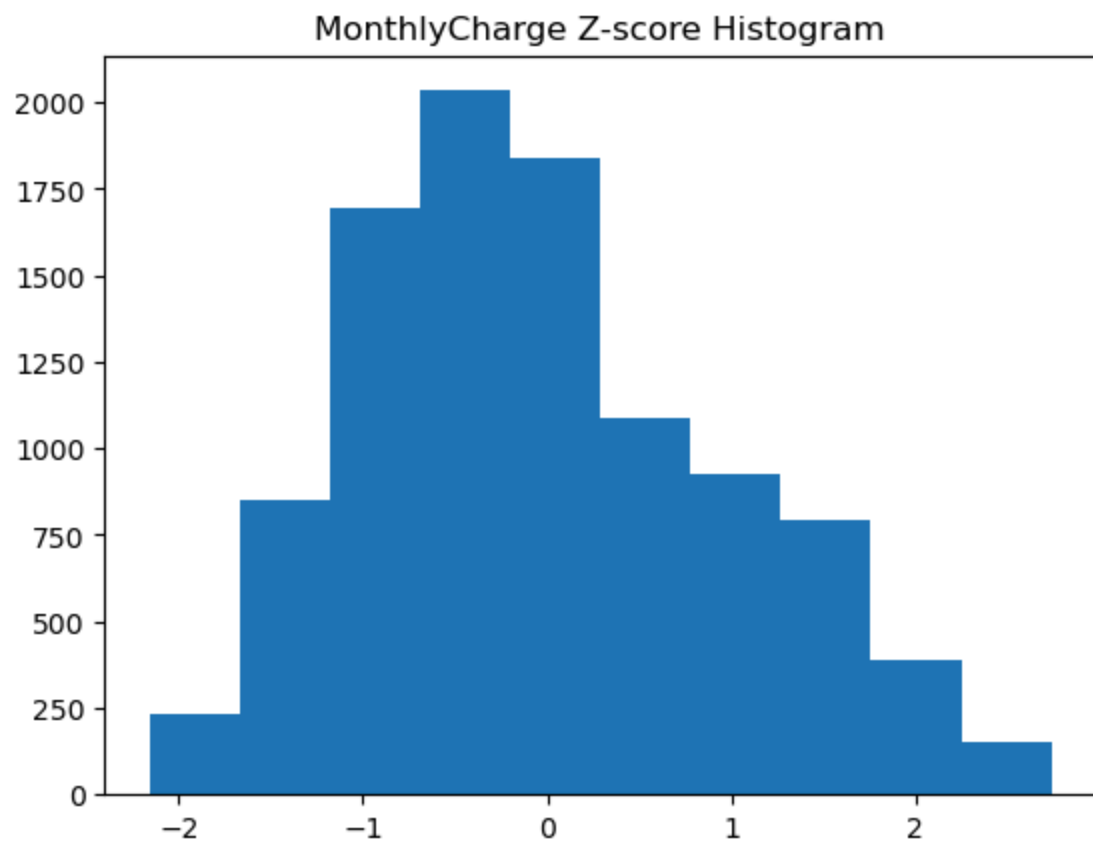
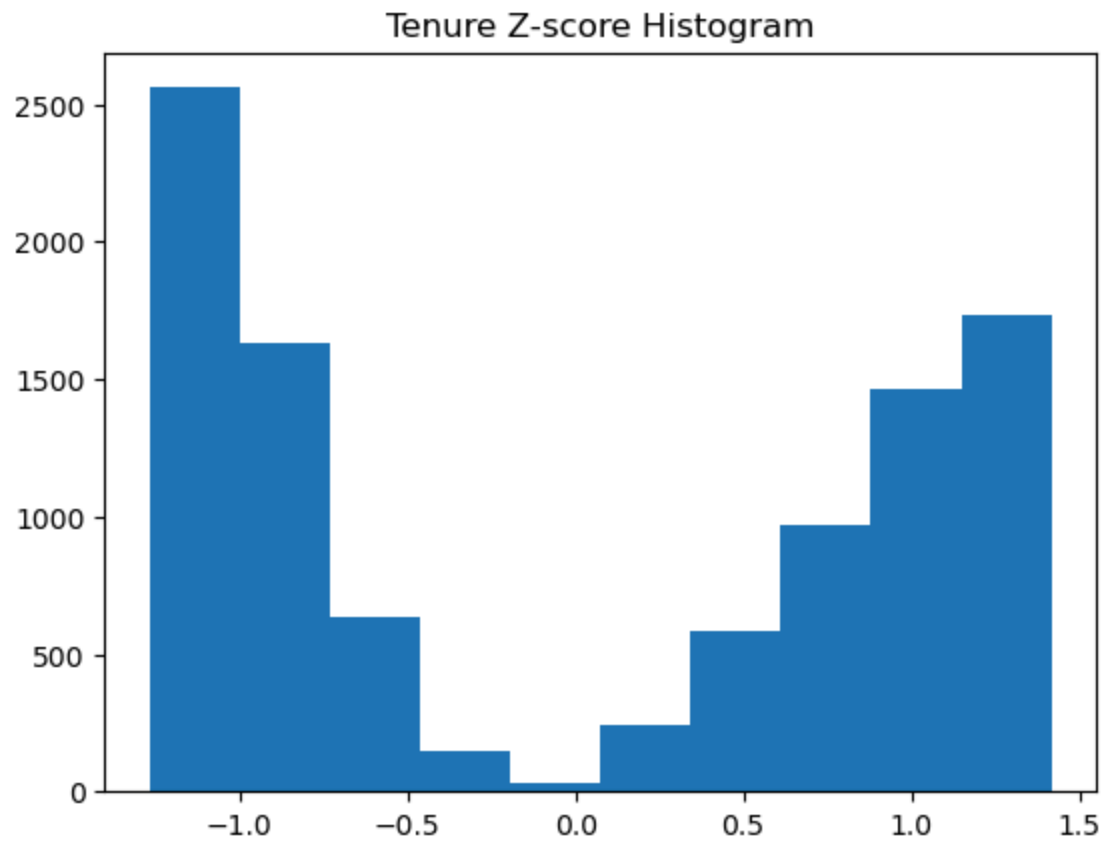


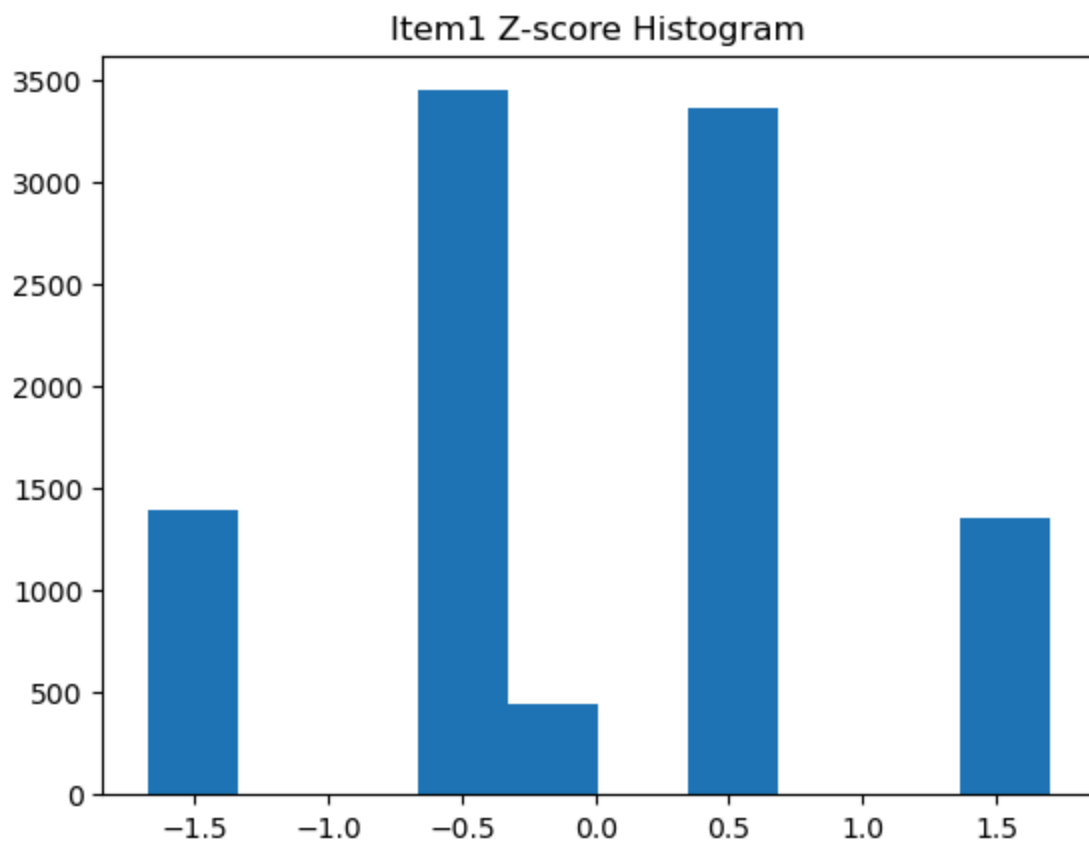
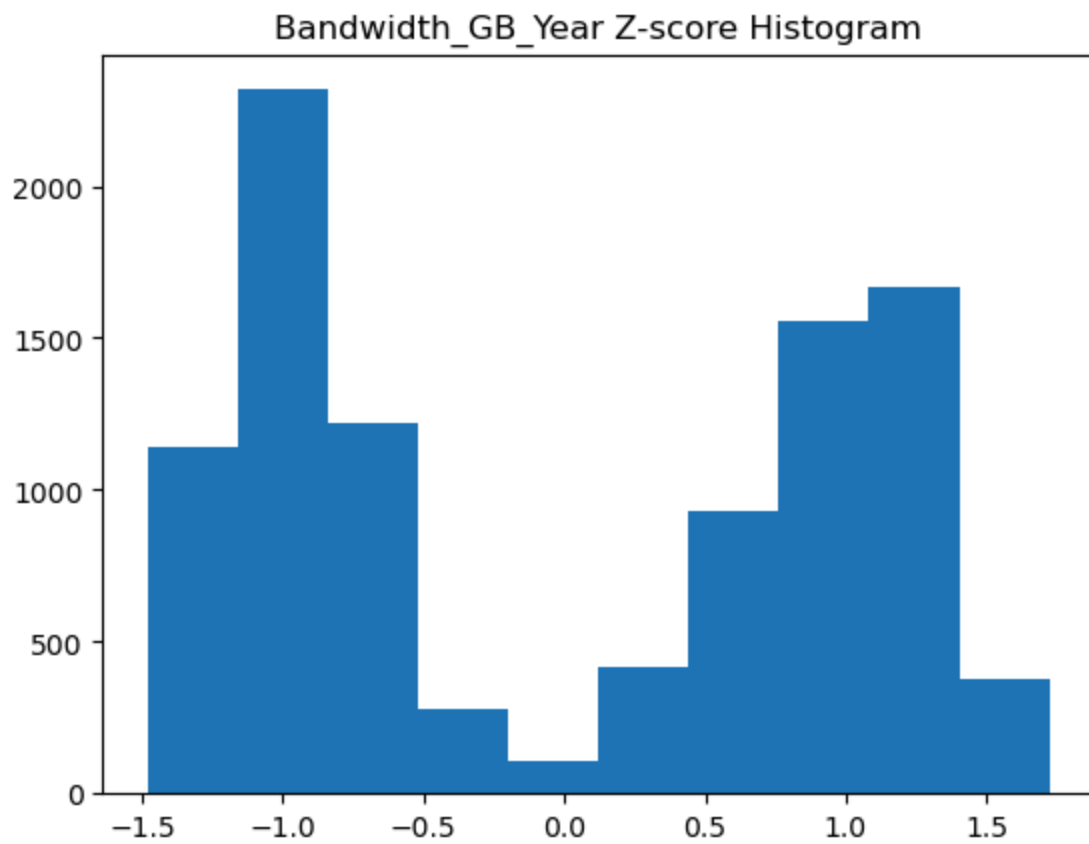
Email Z-score Histogram

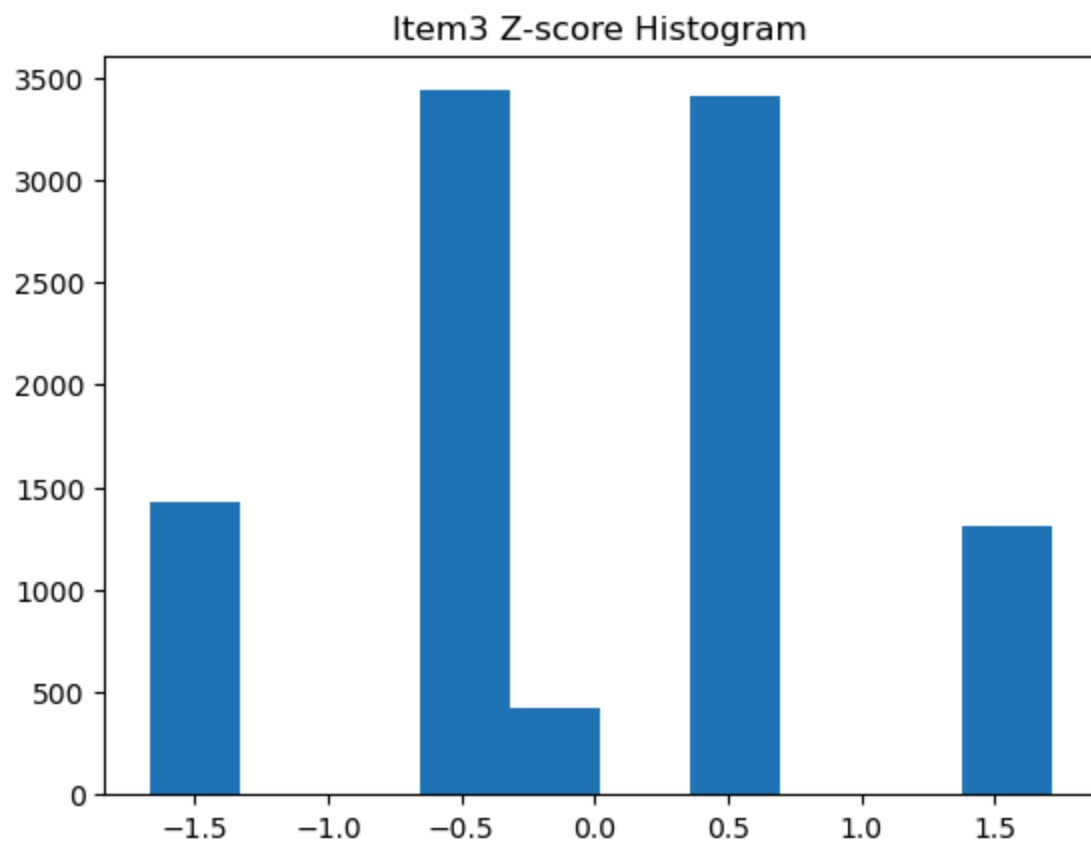
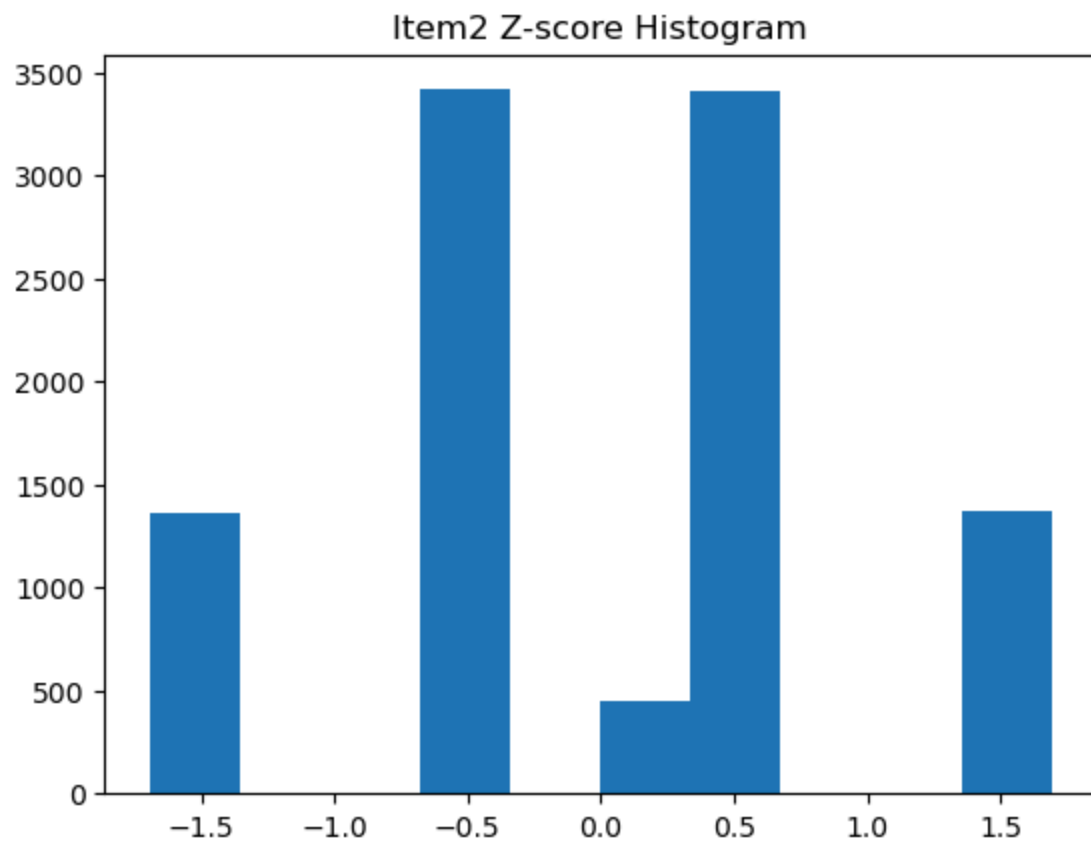


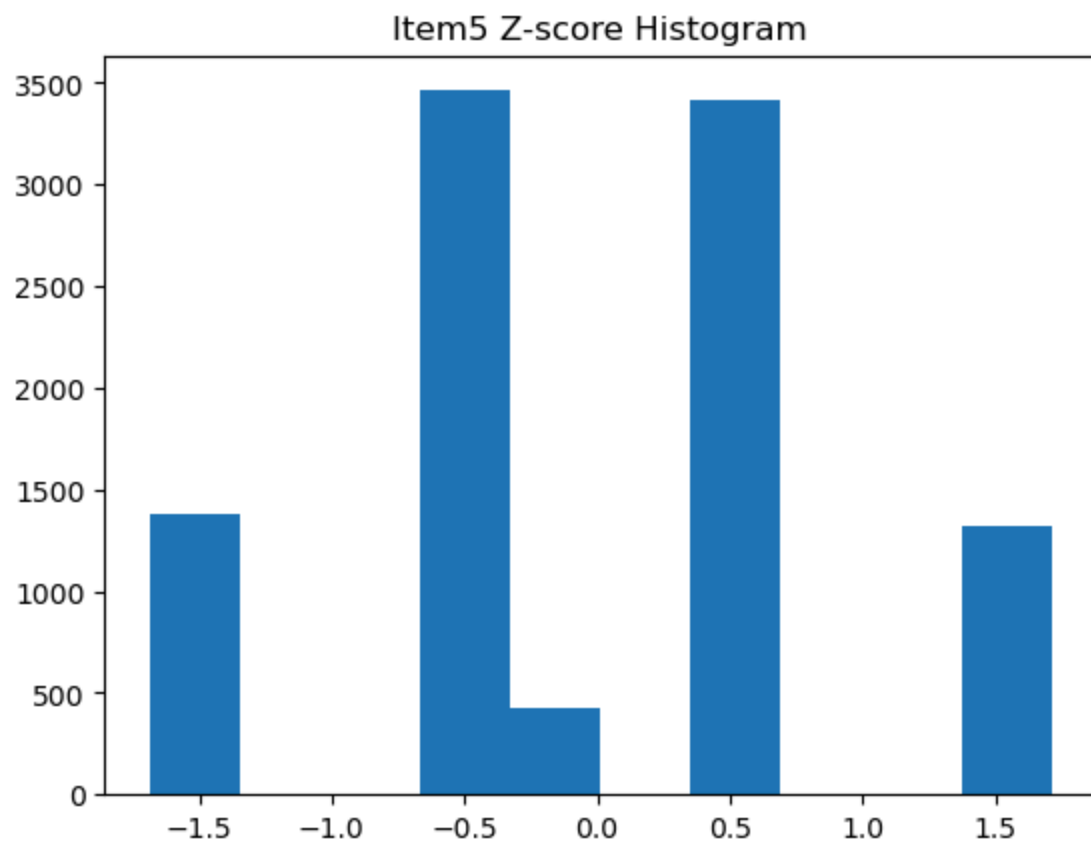
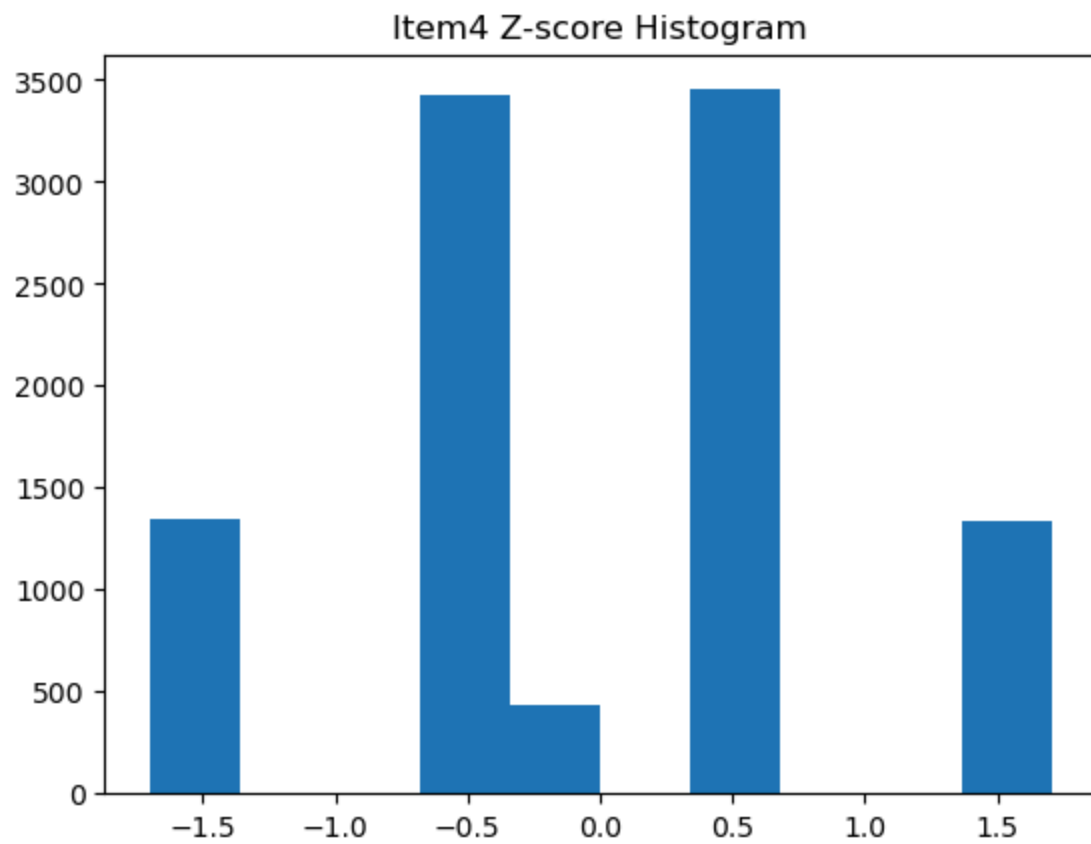


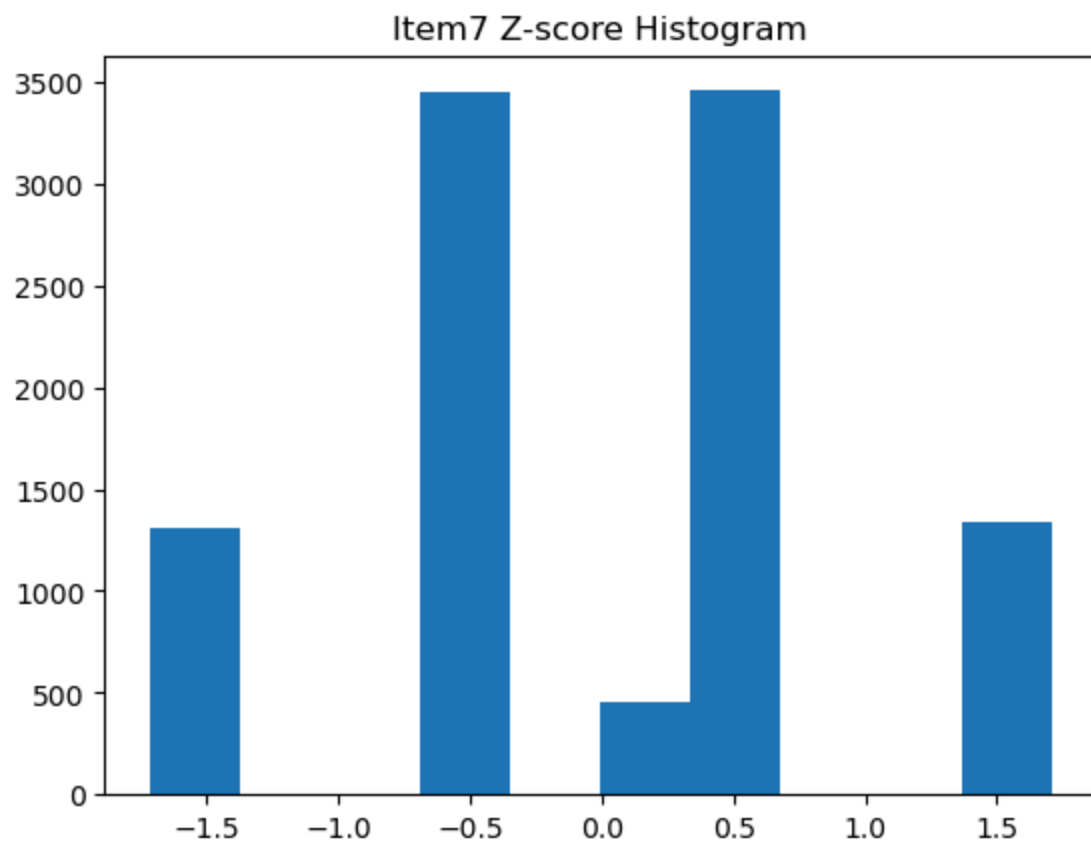
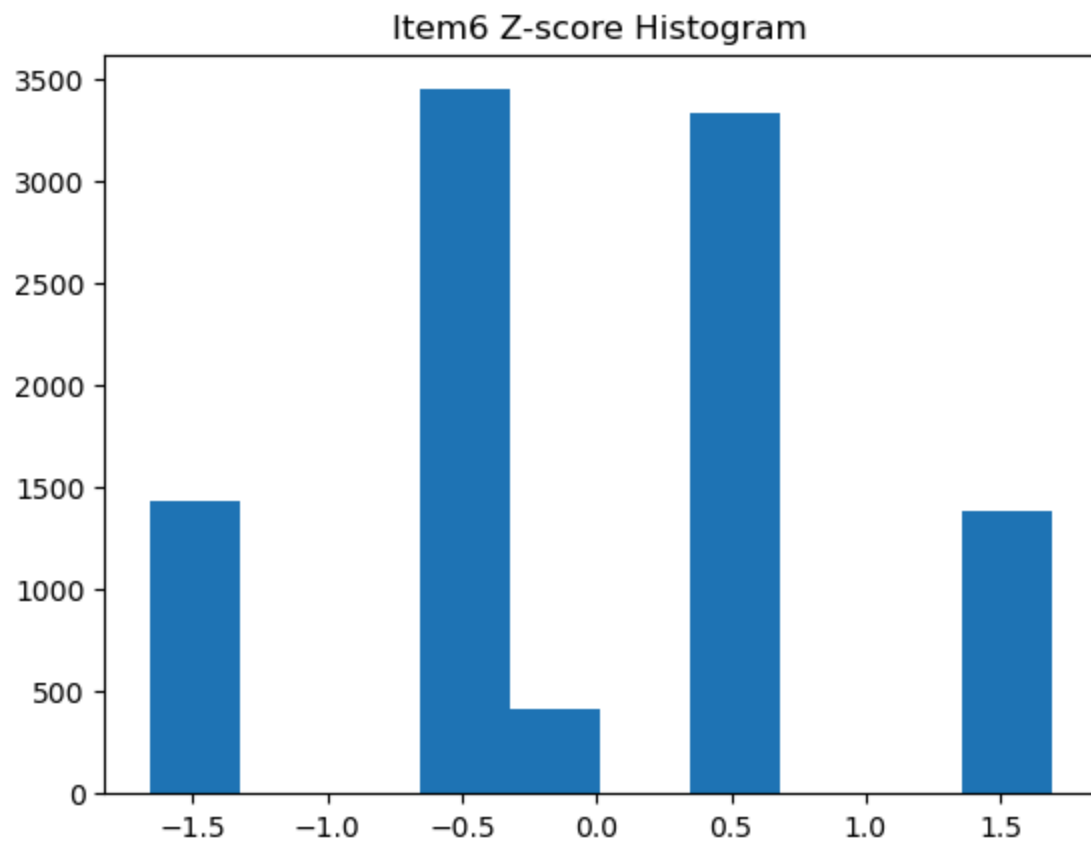


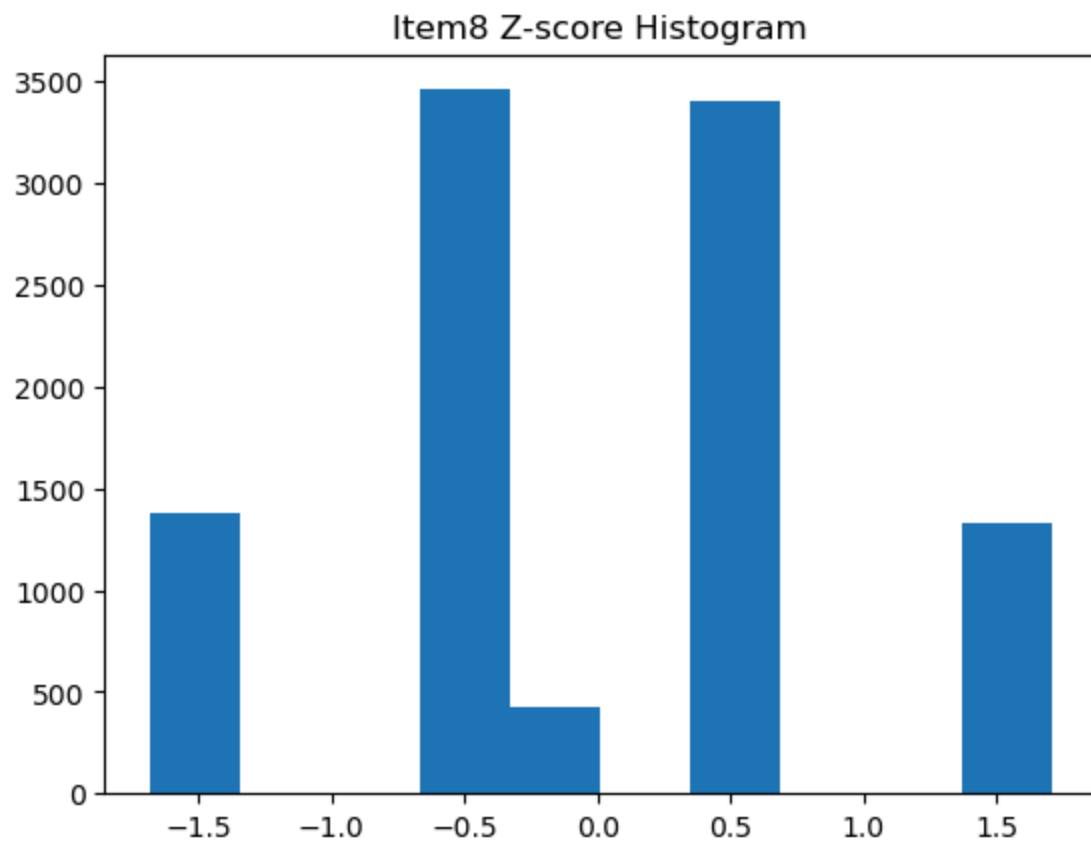




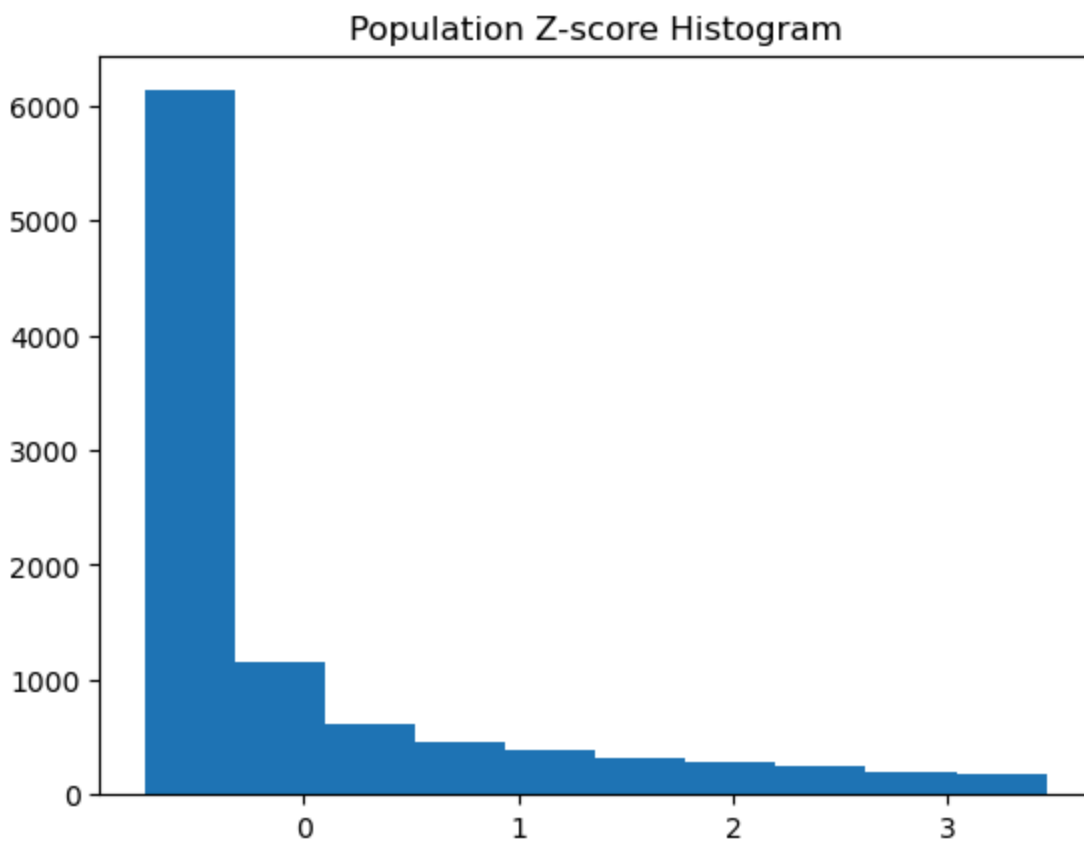
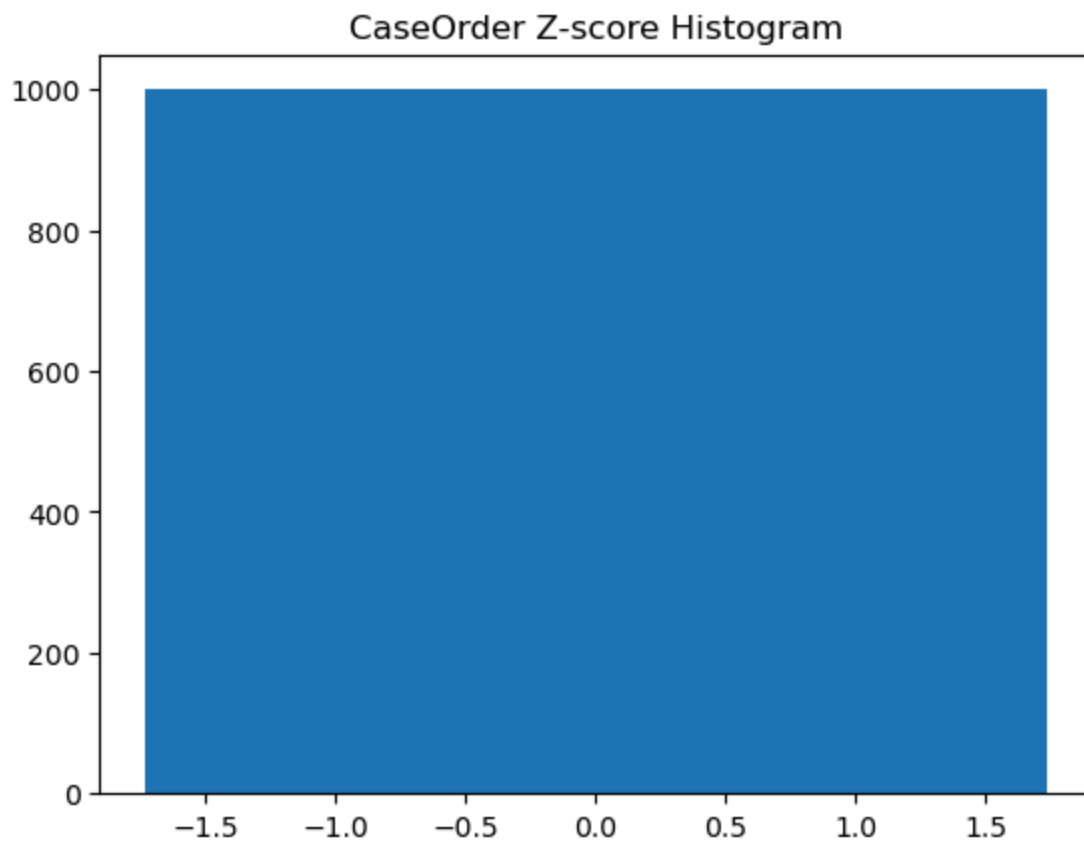


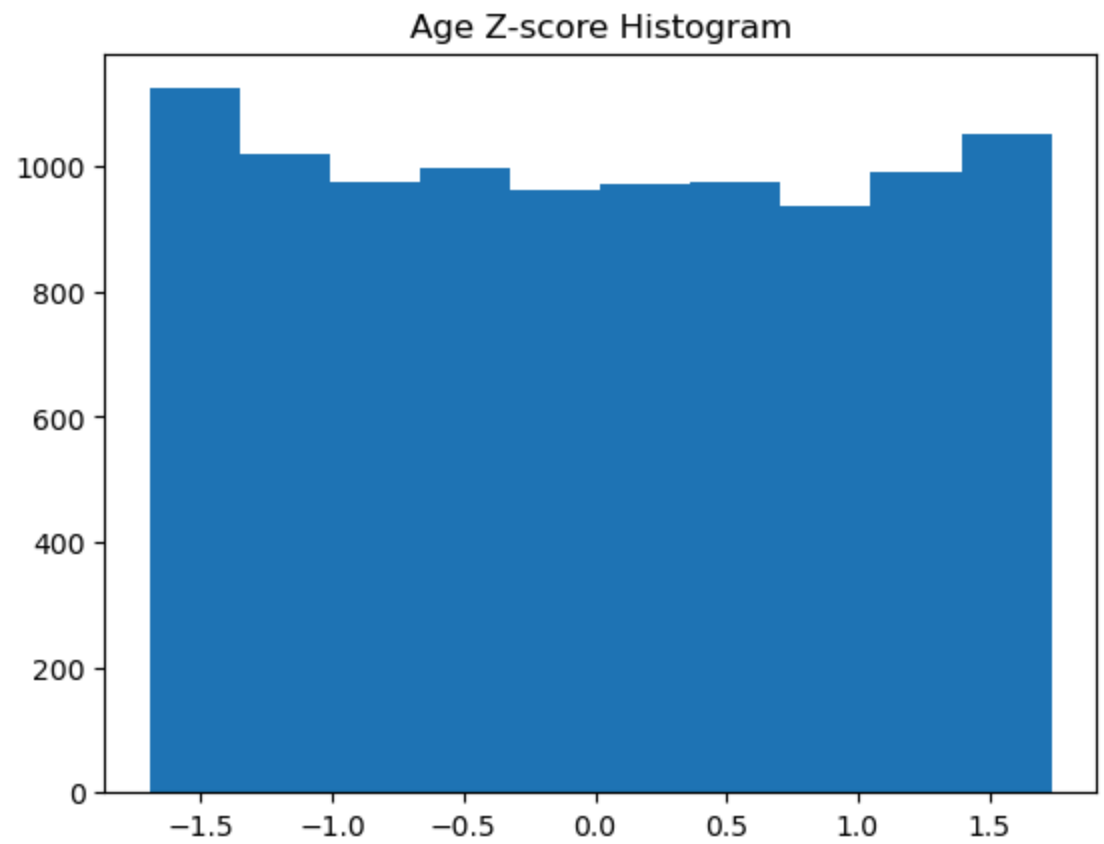
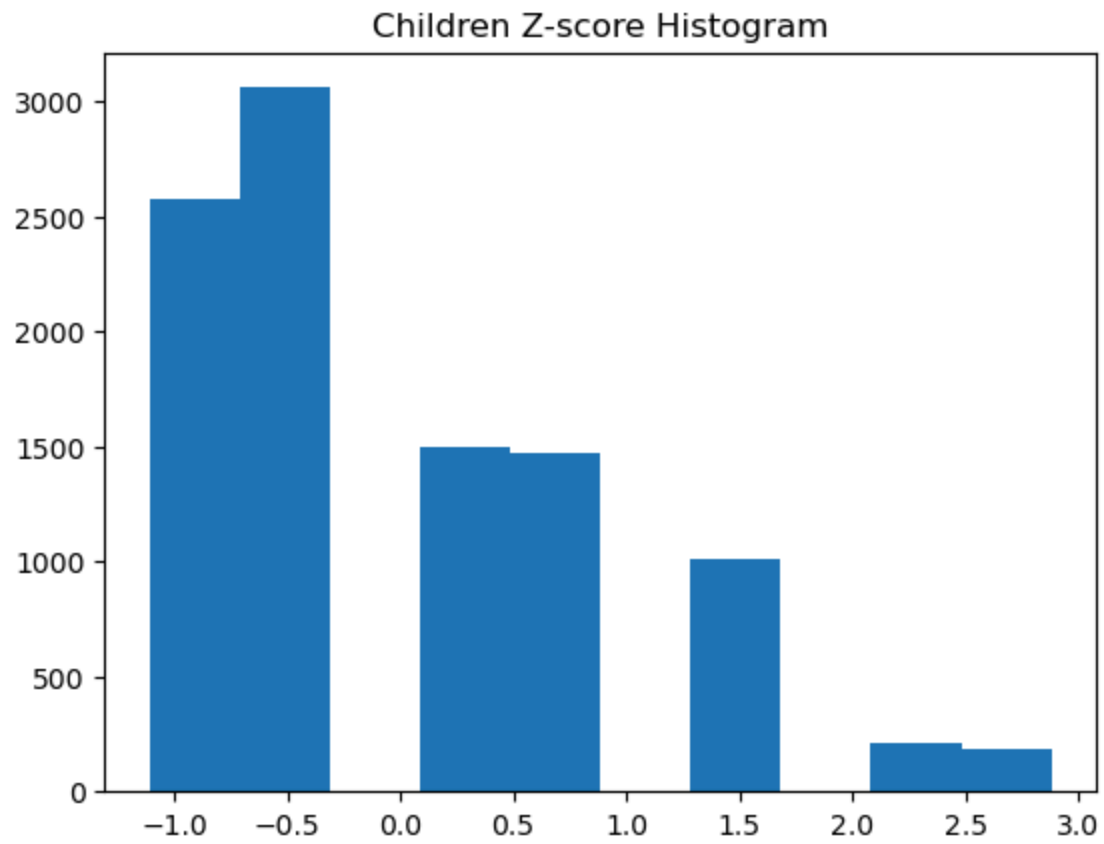




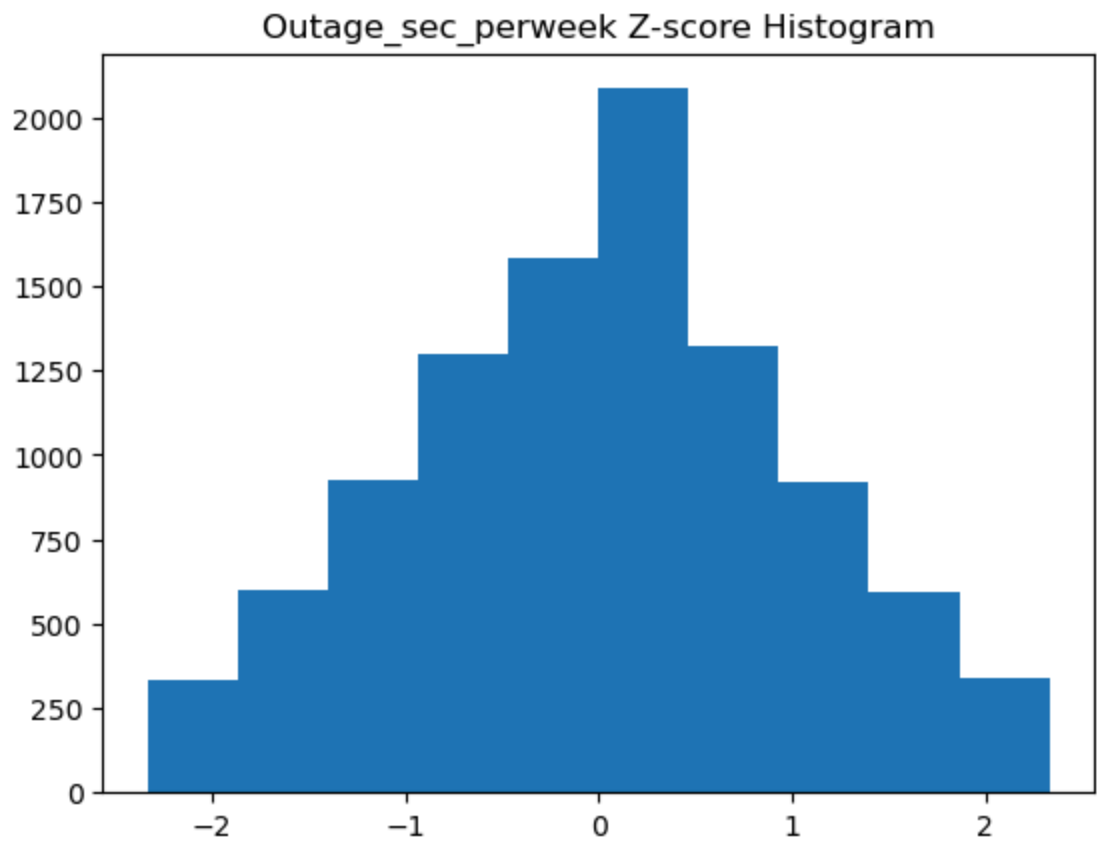
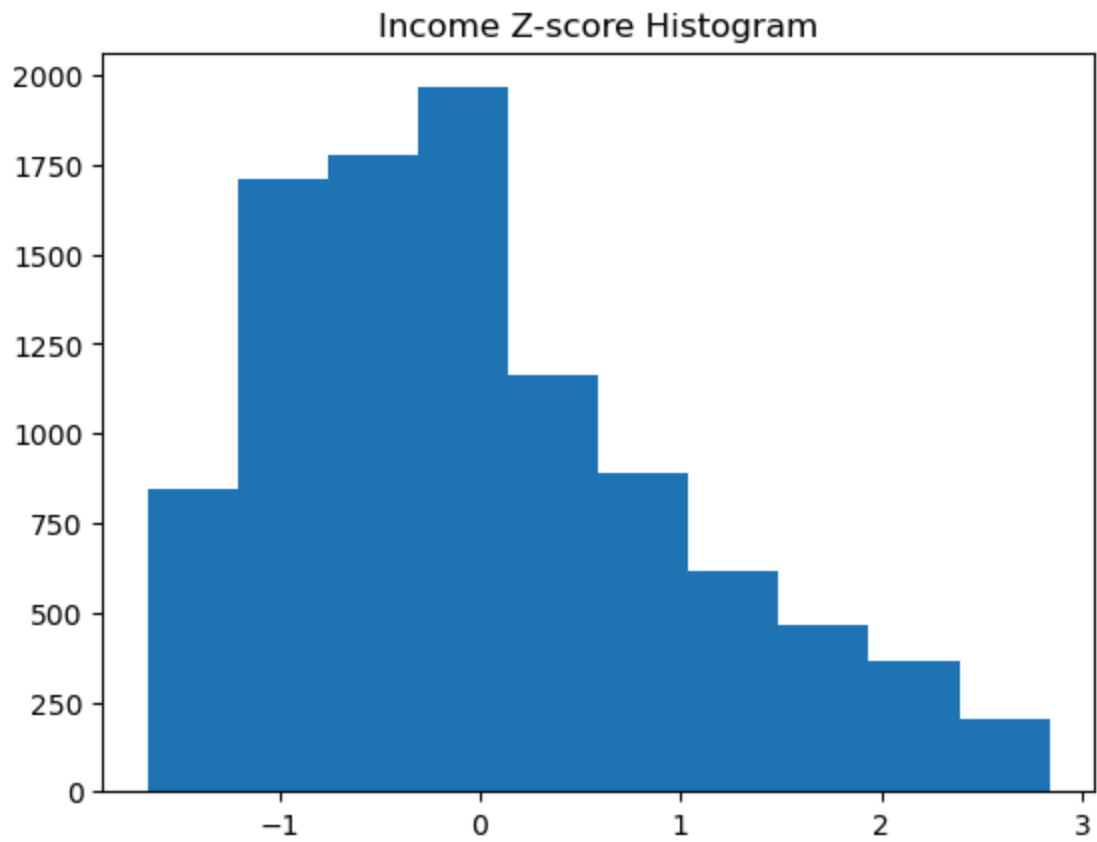


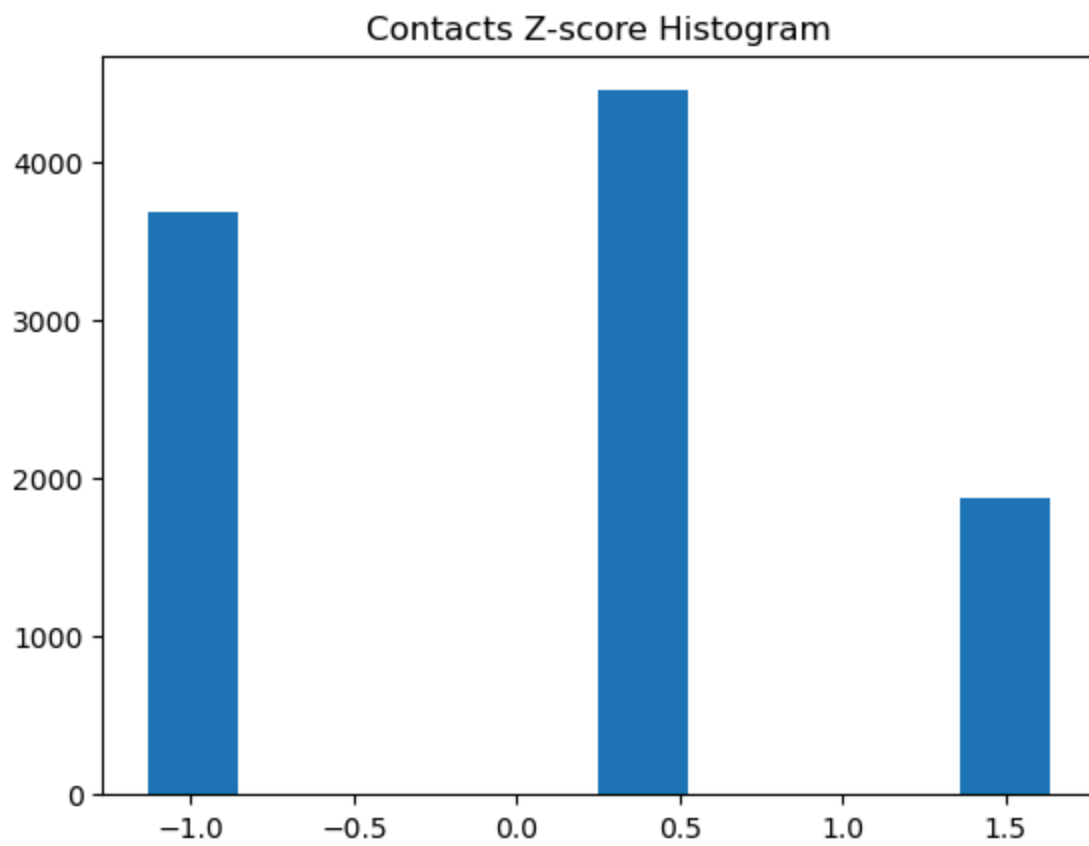
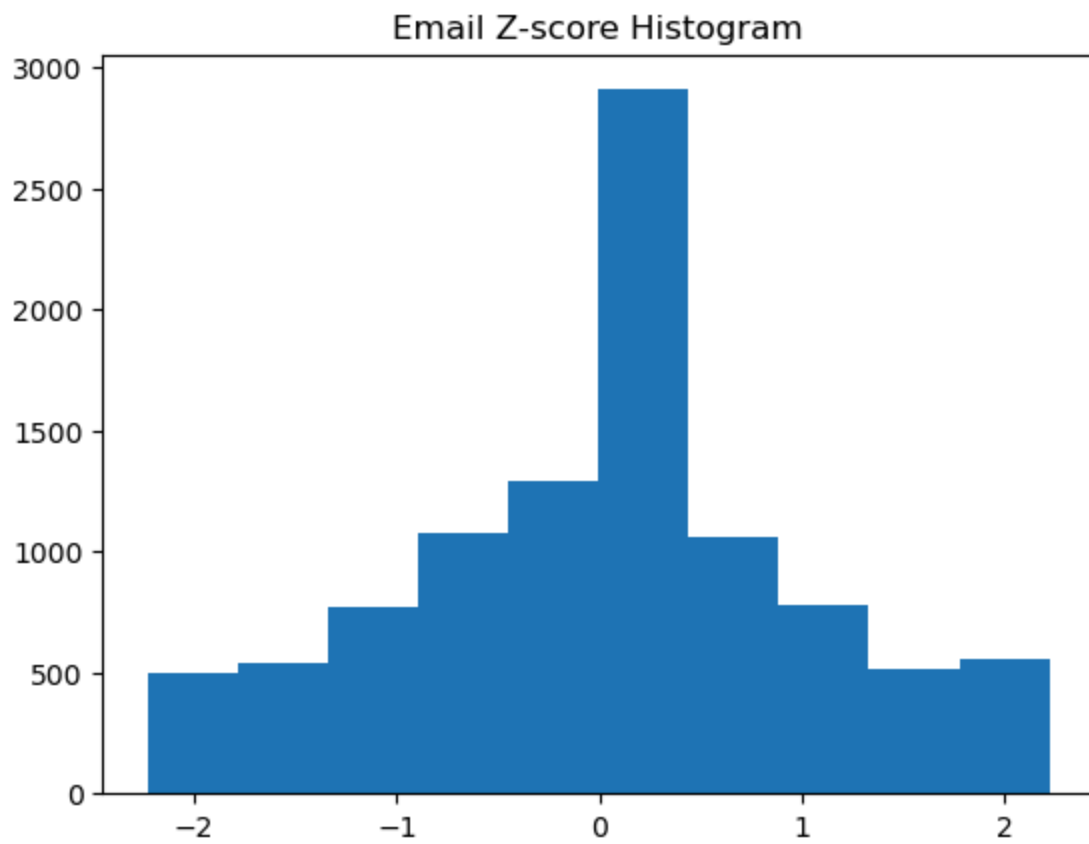
```
In [18]: for column in dfq_c:
          dfq['zscore'] = stats.zscore(dfq[column])
          plt.hist(dfq['zscore'])
          plt.title(column + ' Z-score Histogram')
          plt.show()
```



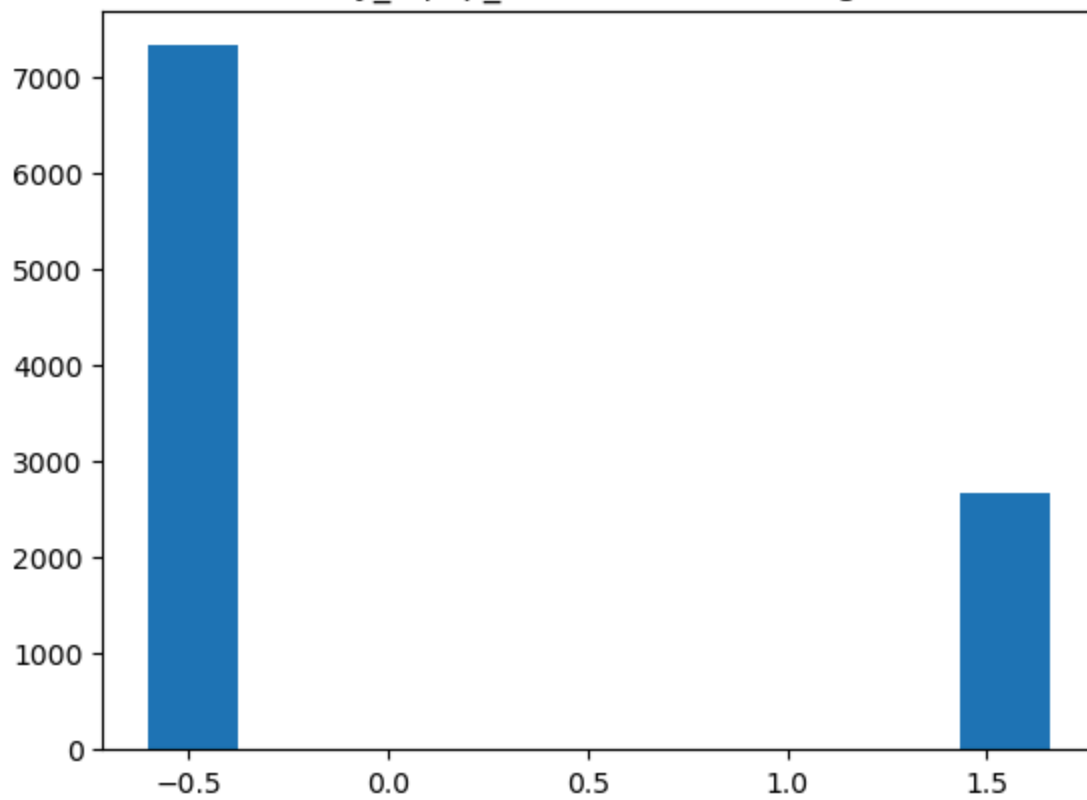




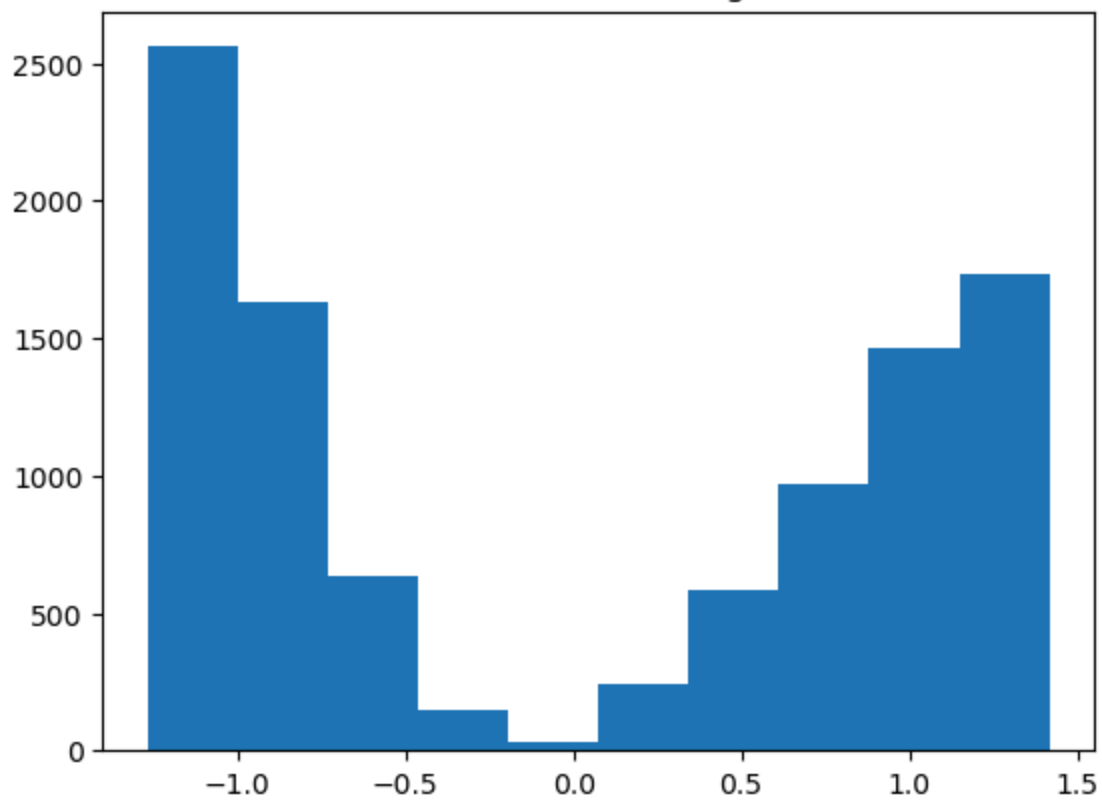


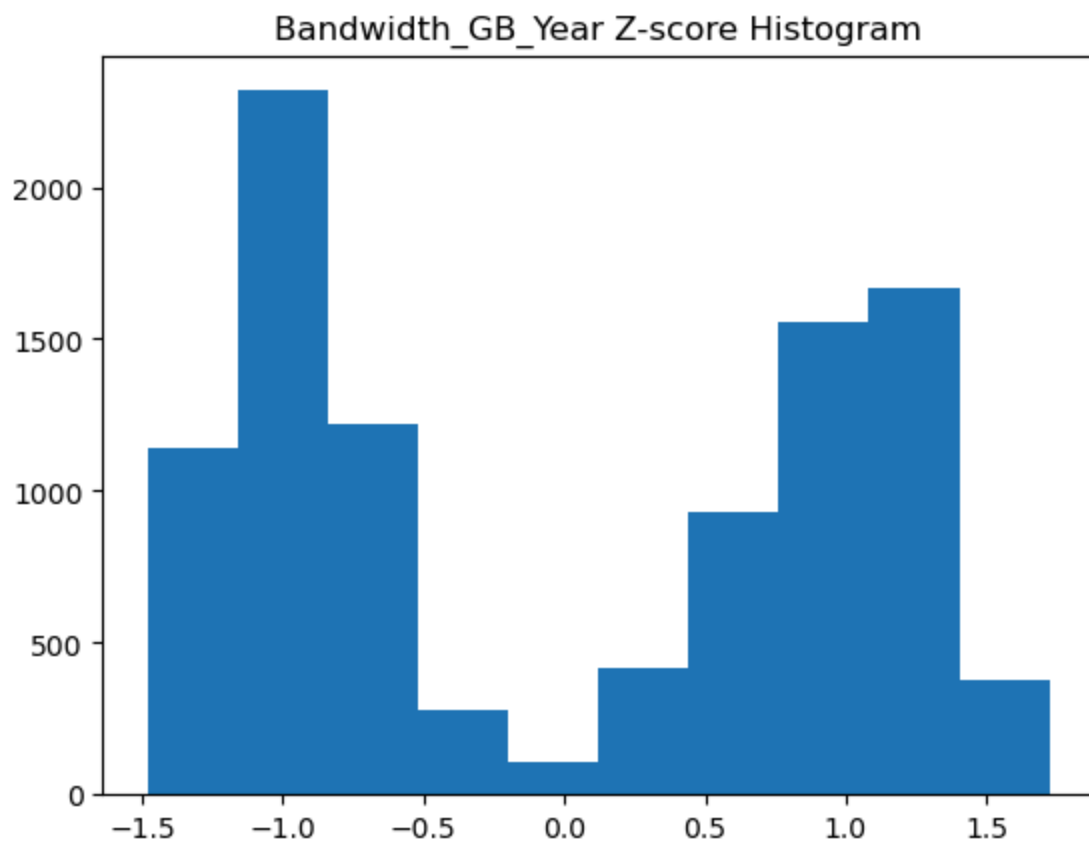
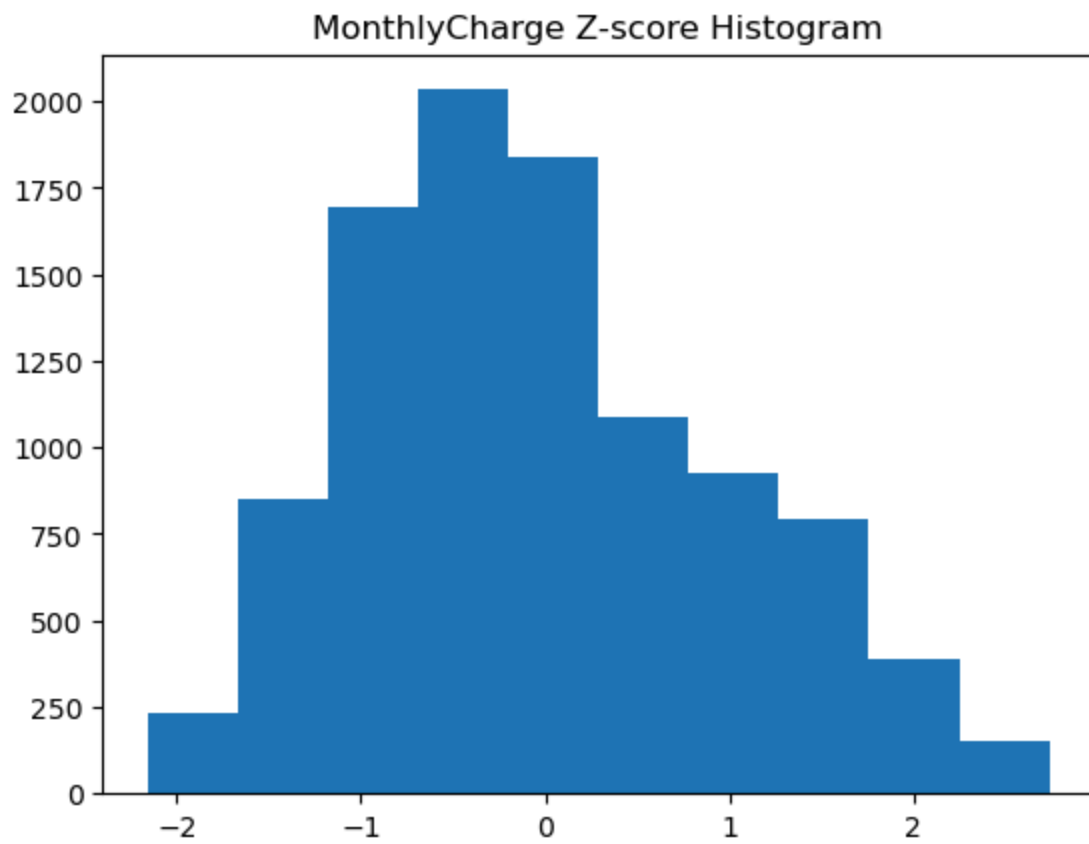


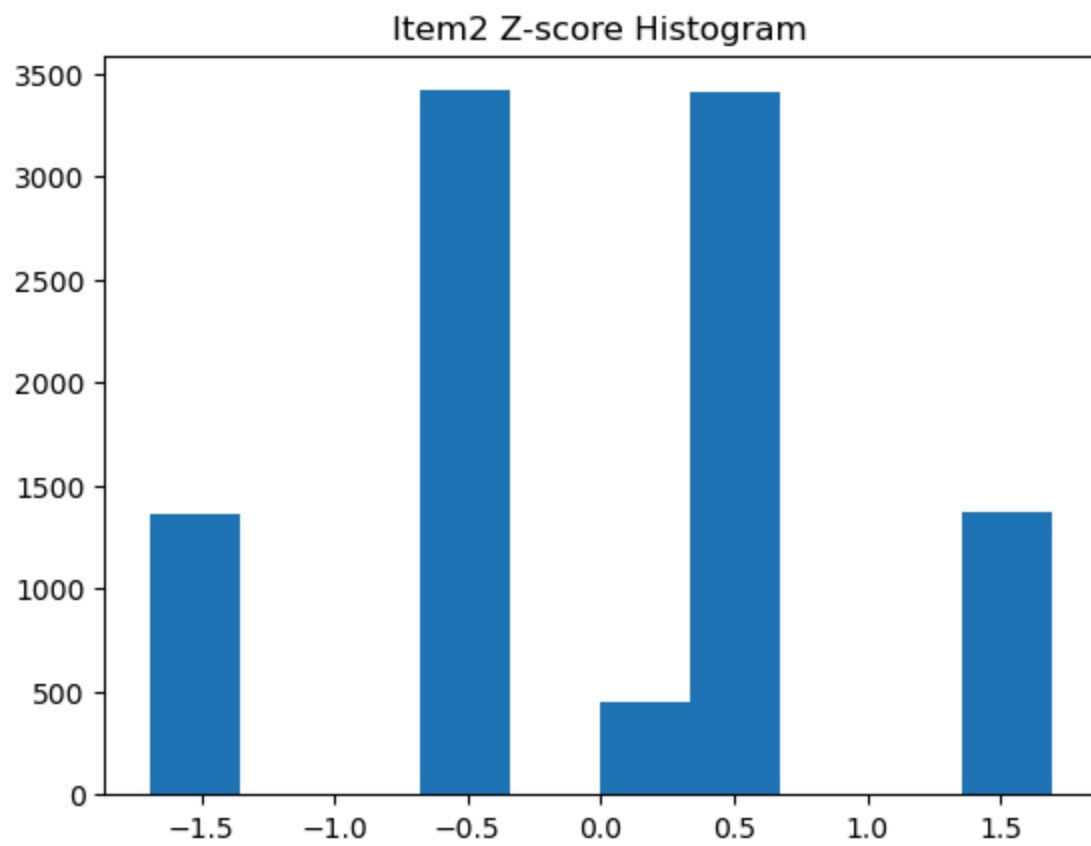
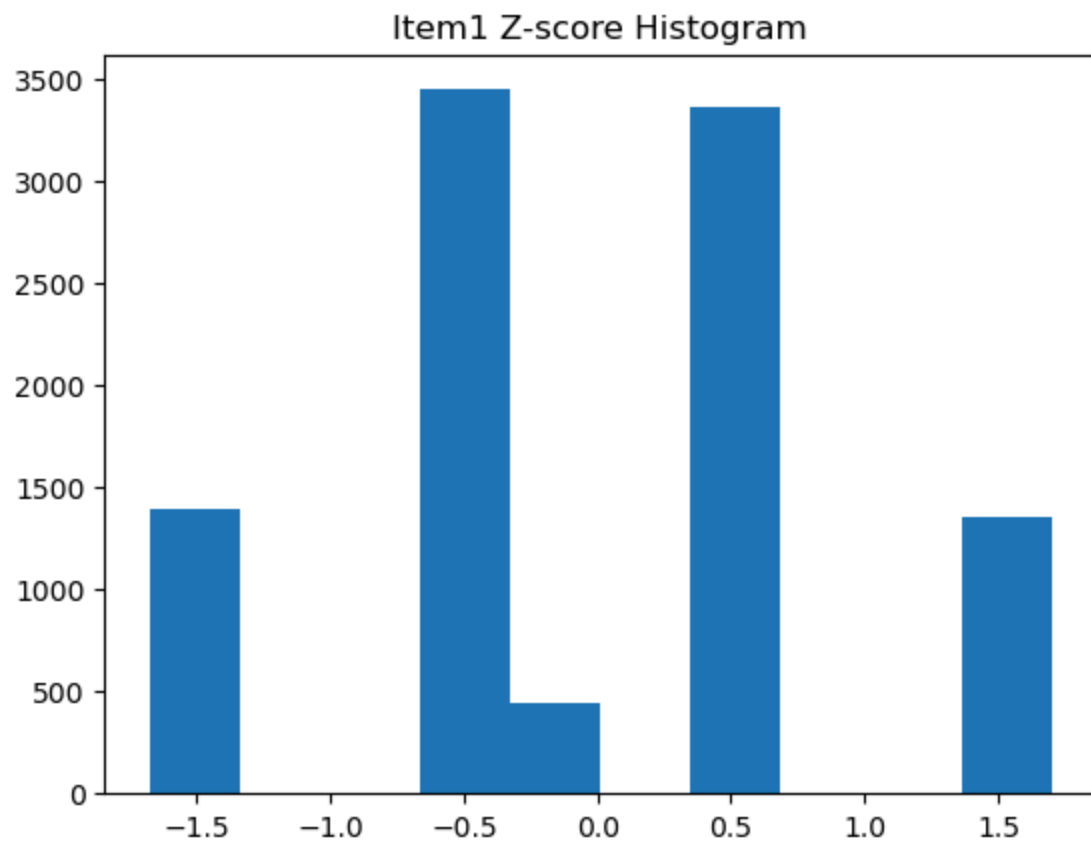
Yearly equip\_failure Z-score Histogram

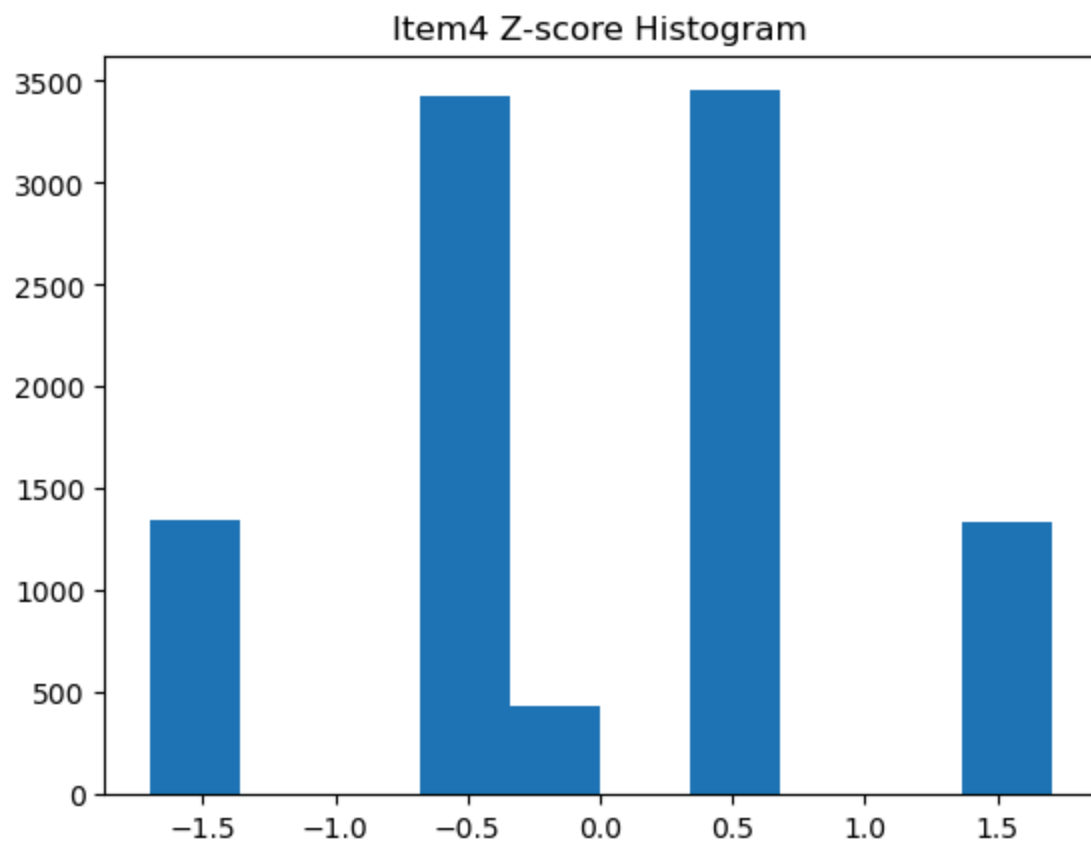
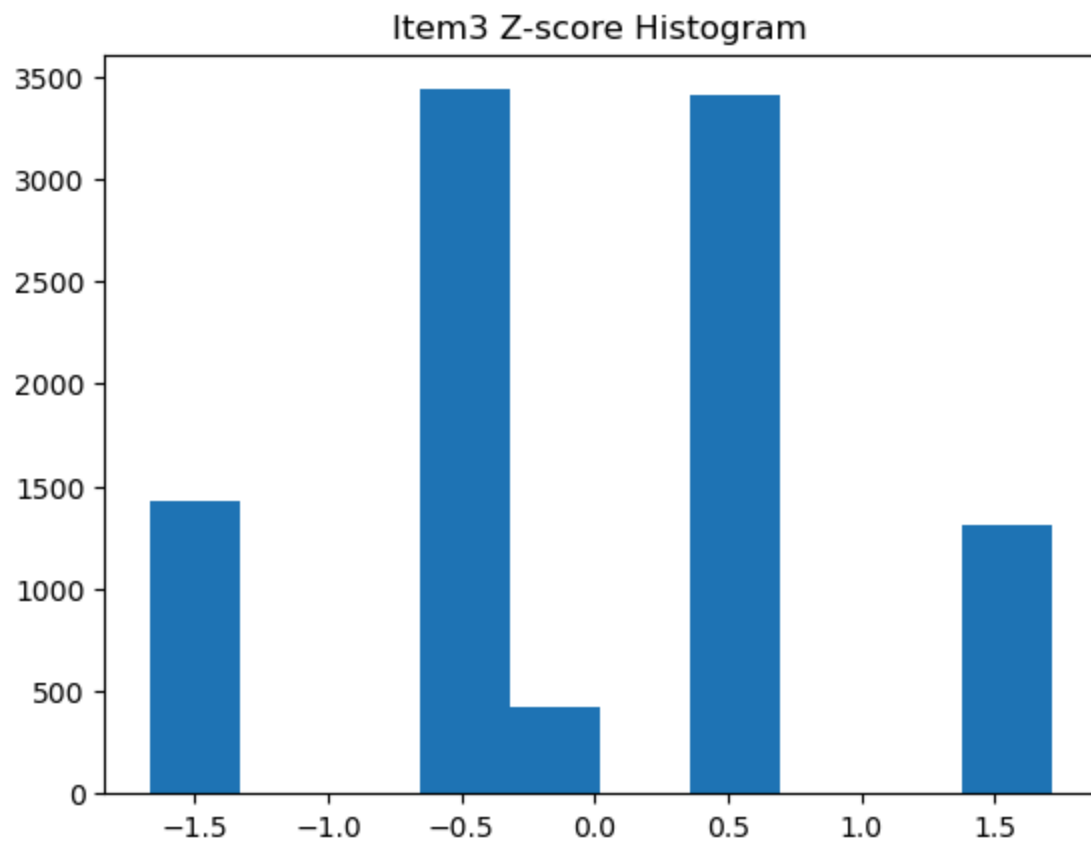


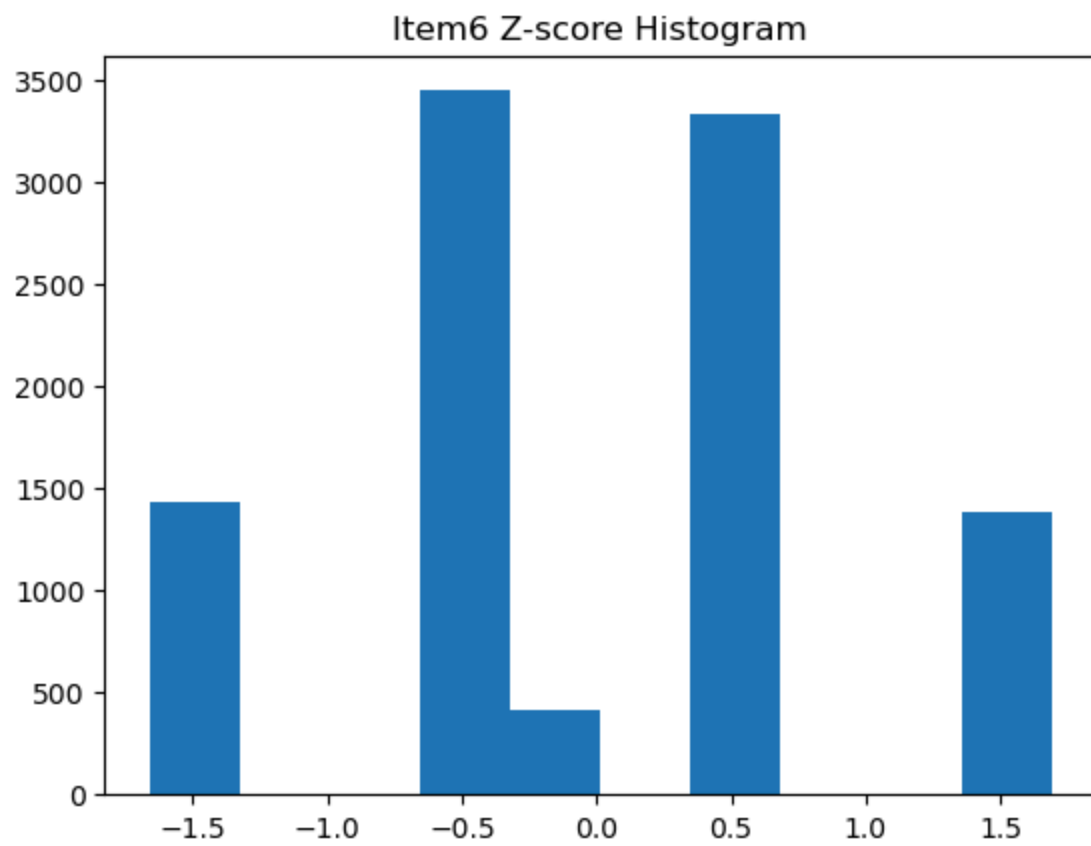
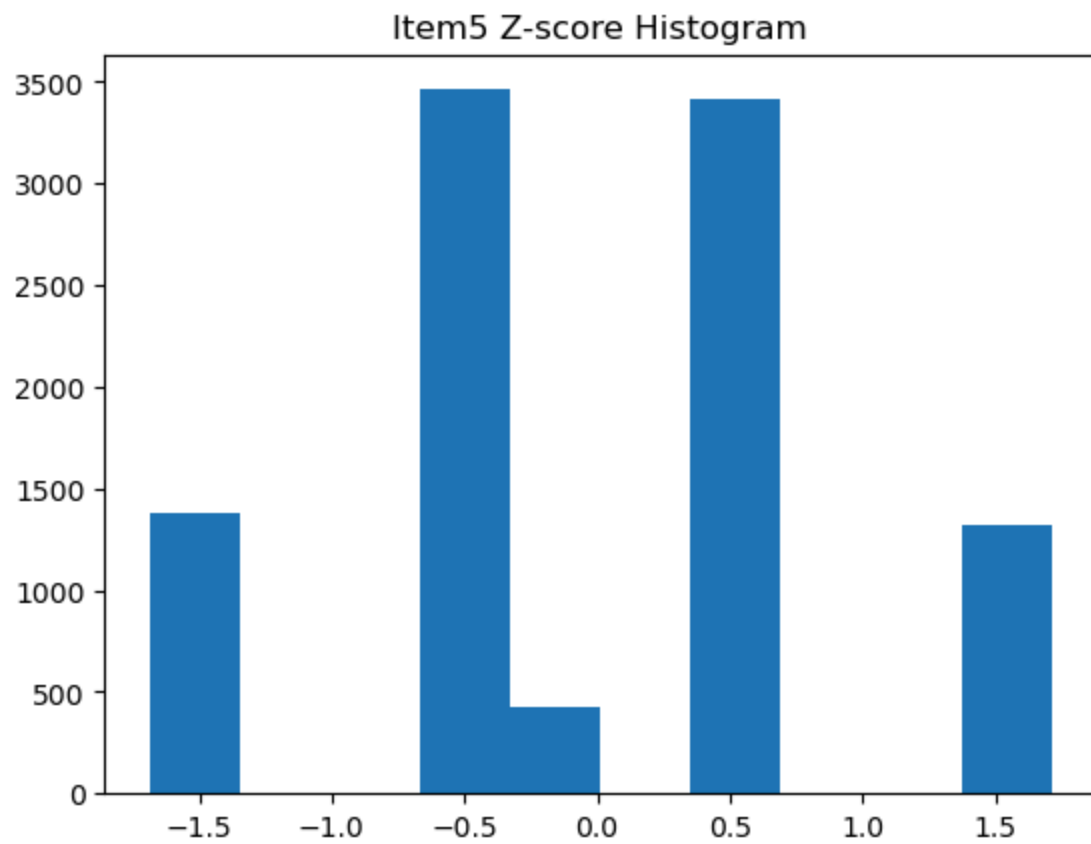
Tenure Z-score Histogram

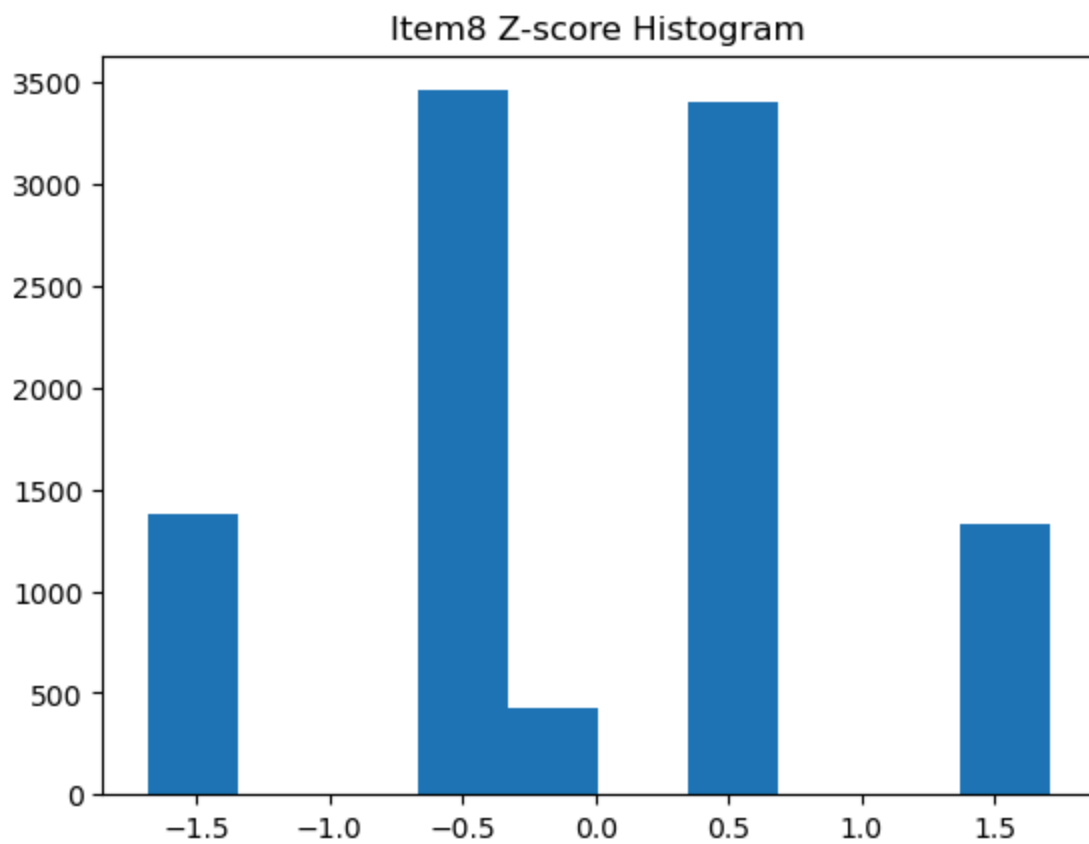
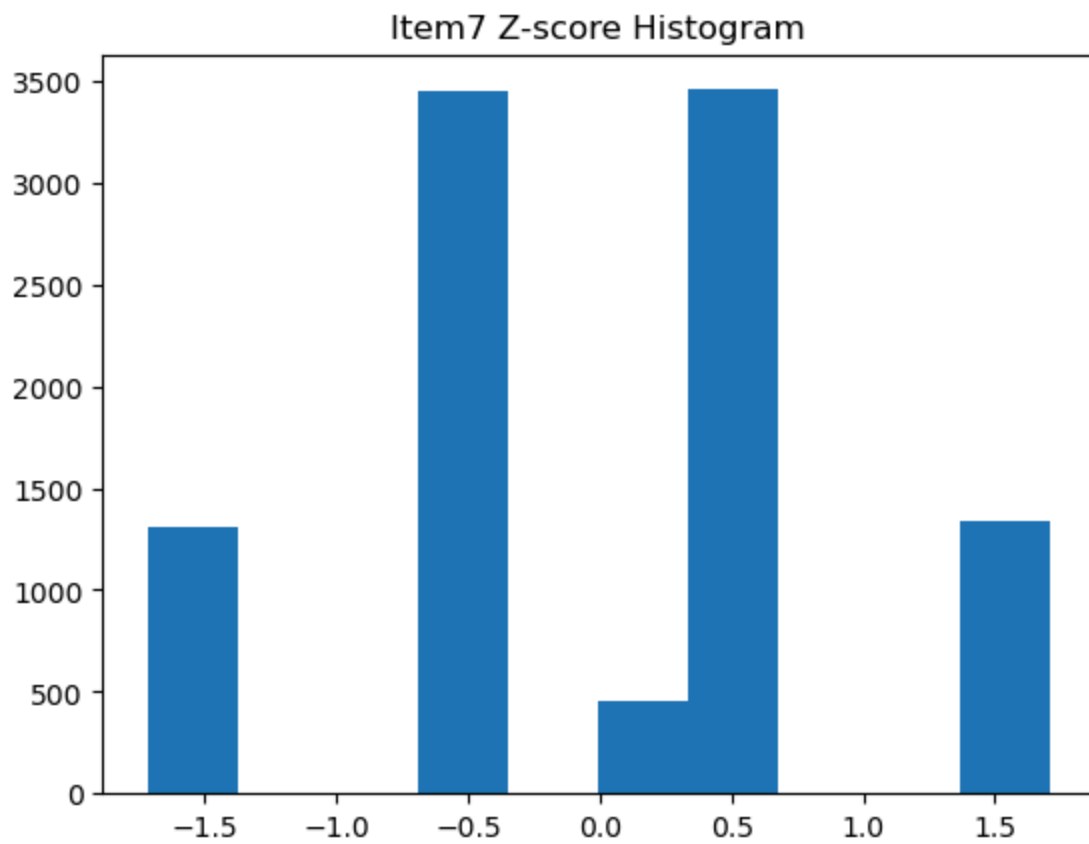












## C2. Data Exploration (EDA)



To get out summary statistics, we can use the describe function. I use a for loop to run the .describe() function and get summary statistics of all the variables. Its also important to note that all outliers have been removed.

CaseOrder: For CaseOrder, our summary statistics make sense. Each case is assigned a number and it goes from 1 to 10,000. The summary statistics match what one would expect to see for a list of number from 1 to 10,000.

Population: The average population is 6817, the minimum is 2, and the maximum is 38,597. We can once again see the zeroes have been removed as per our treatment of nulls earlier in the analysis.

Children: The mean, which is calculated by adding up all the values and dividing by the n amount, is around 1.7 children. The most children is 6, and the least is 0.

Age: The average age is around 53. This means many of the customers tend to be older. The youngest customer is 18 and the oldest is 89.

Income: The average income is around 35,688 dollars. The lowest is 348 dollars and the highest is 96,190.

Outage\_sec\_perweek: On average there is an average outage time of 10 seconds per week. The minimum or lowest time is 4 seconds, and the highest is 15. This is interesting as we learn that there is never a point in time where there is a week without outages.

Email: The average number of emails is 12. The minimum is 6 and the maximum is 18. This gives us insight that depending on the customer, different amounts of emails are sent. This may be because the business segments its customers or because some customers joined at different times and thus were not included in previous emails.

Contacts: The average for this variable is .8. The minimum is 0 and the max is 2. This shows us that customers do not frequently contact customer support, with the most a customer contacting them being 2 recorded times.

Yearly\_equip\_failure: The average for this variable is .3. The minimum is 0 and the maximum is 1. This shows us that its not frequent for a customer's equipment to fail, and that it will most likely not occur more than once according to our recorded history.

Tenure: The average tenure is around 34.5 months. The minimum is 1 and the max is 72. This shows us that tenure of the customer does generally not last for more than a few years according to our data.

MonthlyCharge: The average monthly charge is about 172 dollars a month. The minimum is 80 and max is 290. This could be due to different customers having different plans, customizable services, and offers.

Bandwidth\_gb\_year: This is the amount of gb a customer uses per year. On average it is 3392, with the lowest being 155.5 and the highest being 7158.98 gb.

Item1 through Item8: Items 1 through 8 should all have a minimum of 2 and a max of 5. This could be a result of how we cleaned the data, removing any outliers. On average the amount is about

```
In [19]: # run a for loop that goes through and uses .describe()
for column in dfq_c:
    print('Variable: ', column, '\n', dfq[column].describe(), '\n')
```

Variable: CaseOrder  
count 10000.00000  
mean 5000.50000  
std 2886.89568  
min 1.00000  
25% 2500.75000  
50% 5000.50000  
75% 7500.25000  
max 10000.00000  
Name: CaseOrder, dtype: float64

Variable: Population  
count 10000.000000  
mean 6817.325500  
std 9188.402721  
min 2.000000  
25% 782.000000  
50% 2610.000000  
75% 8810.000000  
max 38597.000000  
Name: Population, dtype: float64

Variable: Children  
count 10000.00000  
mean 1.66700  
std 1.50424  
min 0.00000  
25% 0.00000  
50% 1.00000  
75% 3.00000  
max 6.00000  
Name: Children, dtype: float64

Variable: Age  
count 10000.000000  
mean 53.078400  
std 20.698882  
min 18.000000  
25% 35.000000  
50% 53.000000  
75% 71.000000  
max 89.000000  
Name: Age, dtype: float64

Variable: Income  
count 10000.000000  
mean 35688.400354  
std 21324.318547  
min 348.670000  
25% 19224.717500  
50% 31716.910000  
75% 48598.185000  
max 96190.740000  
Name: Income, dtype: float64

Variable: Outage\_sec\_perweek  
count 10000.000000  
mean 10.001559  
std 2.545498  
min 4.065560

```
25%      8.254991
50%     10.001559
75%     11.747969
max      15.931970
Name: Outage_sec_perweek, dtype: float64
```

```
Variable: Email
count    10000.000000
mean     12.014786
std      2.695926
min      6.000000
25%     10.000000
50%     12.000000
75%     14.000000
max     18.000000
Name: Email, dtype: float64
```

```
Variable: Contacts
count    10000.000000
mean     0.819200
std      0.722886
min      0.000000
25%      0.000000
50%      1.000000
75%      1.000000
max      2.000000
Name: Contacts, dtype: float64
```

```
Variable: Yearly_equip_failure
count    10000.000000
mean     0.267000
std      0.442414
min      0.000000
25%      0.000000
50%      0.000000
75%      1.000000
max      1.000000
Name: Yearly_equip_failure, dtype: float64
```

```
Variable: Tenure
count    10000.000000
mean     34.526188
std      26.443063
min      1.000259
25%      7.917694
50%     35.430507
75%     61.479795
max     71.999280
Name: Tenure, dtype: float64
```

```
Variable: MonthlyCharge
count    10000.000000
mean     172.624816
std      42.943094
min      79.978860
25%     139.979239
50%     167.484700
75%     200.734725
max     290.160419
Name: MonthlyCharge, dtype: float64
```

```
Variable: Bandwidth_GB_Year
count    10000.000000
mean     3392.341550
std      2185.294852
min       155.506715
25%      1236.470827
50%      3279.536903
75%      5586.141370
max       7158.981530
Name: Bandwidth_GB_Year, dtype: float64
```

```
Variable: Item1
count    10000.000000
mean      3.489956
std       0.888444
min       2.000000
25%       3.000000
50%       3.489956
75%       4.000000
max       5.000000
Name: Item1, dtype: float64
```

```
Variable: Item2
count    10000.000000
mean      3.501099
std       0.885750
min       2.000000
25%       3.000000
50%       3.501099
75%       4.000000
max       5.000000
Name: Item2, dtype: float64
```

```
Variable: Item3
count    10000.000000
mean      3.481319
std       0.886958
min       2.000000
25%       3.000000
50%       3.481319
75%       4.000000
max       5.000000
Name: Item3, dtype: float64
```

```
Variable: Item4
count    10000.000000
mean      3.498798
std       0.881051
min       2.000000
25%       3.000000
50%       3.498798
75%       4.000000
max       5.000000
Name: Item4, dtype: float64
```

```
Variable: Item5
count    10000.000000
mean      3.488724
std       0.882727
```

```

min            2.000000
25%            3.000000
50%            3.488724
75%            4.000000
max            5.000000
Name: Item5, dtype: float64

```

```

Variable: Item6
count      10000.000000
mean        3.487118
std         0.895207
min         2.000000
25%         3.000000
50%         3.487118
75%         4.000000
max         5.000000
Name: Item6, dtype: float64

```

```

Variable: Item7
count      10000.000000
mean        3.504609
std         0.876074
min         2.000000
25%         3.000000
50%         3.504609
75%         4.000000
max         5.000000
Name: Item7, dtype: float64

```

```

Variable: Item8
count      10000.000000
mean        3.490077
std         0.884270
min         2.000000
25%         3.000000
50%         3.490077
75%         4.000000
max         5.000000
Name: Item8, dtype: float64

```

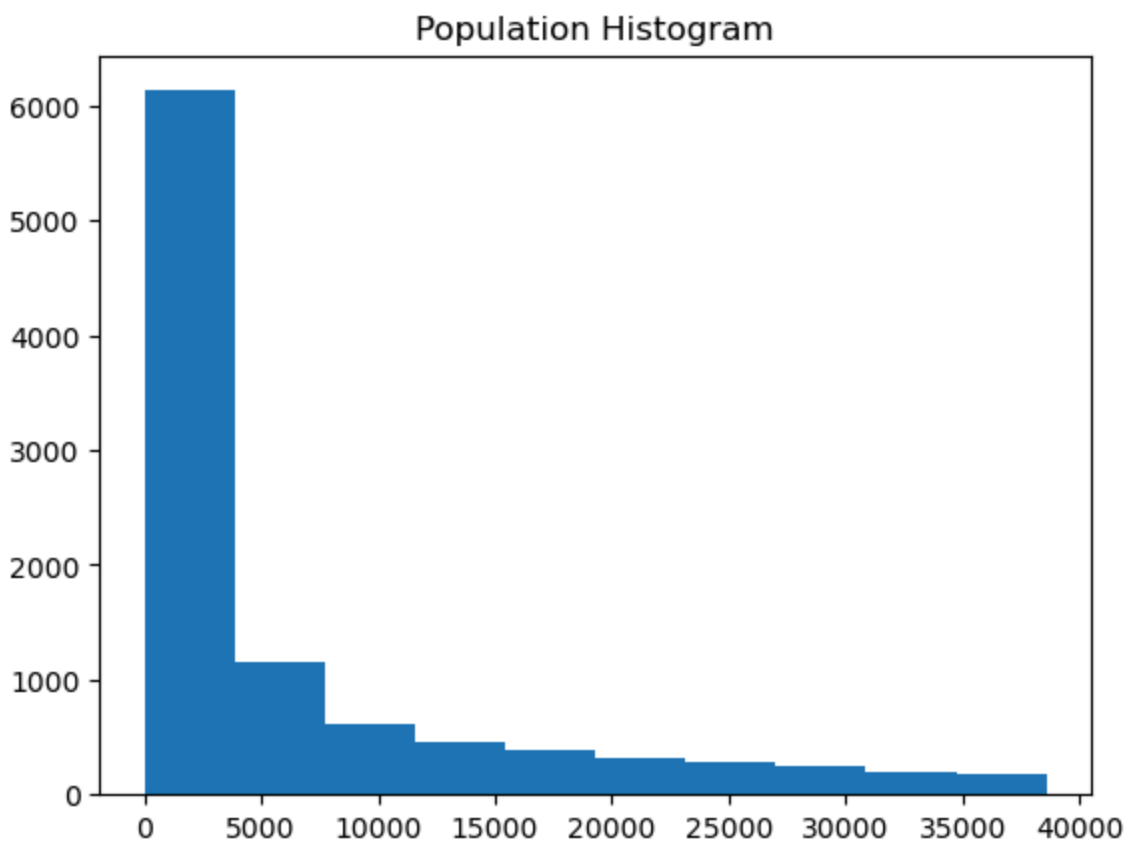
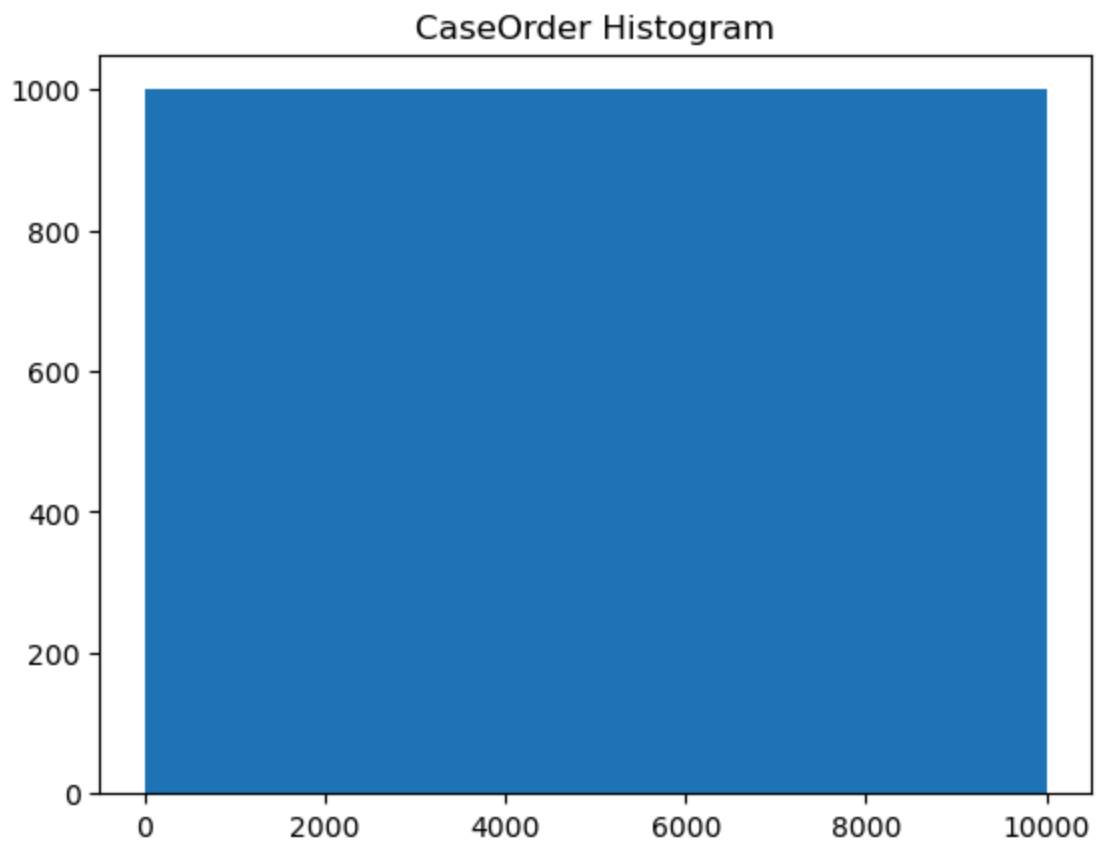
## C3. Visualizations

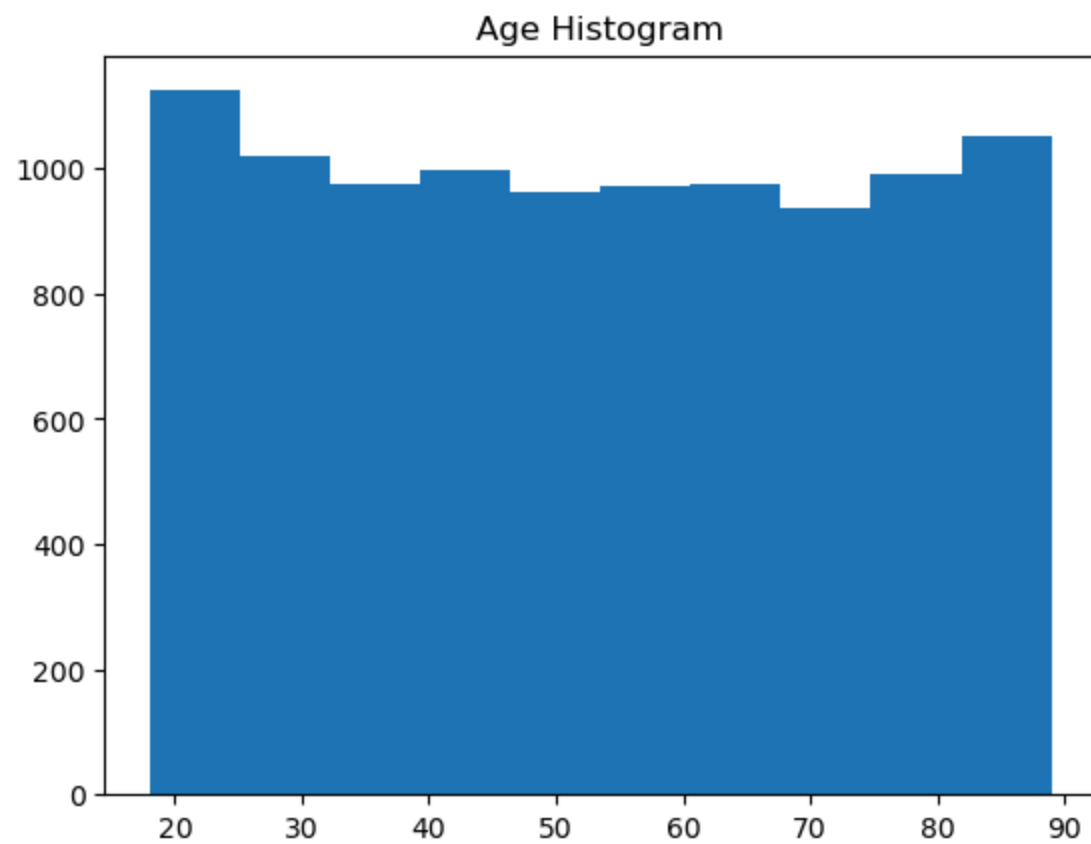
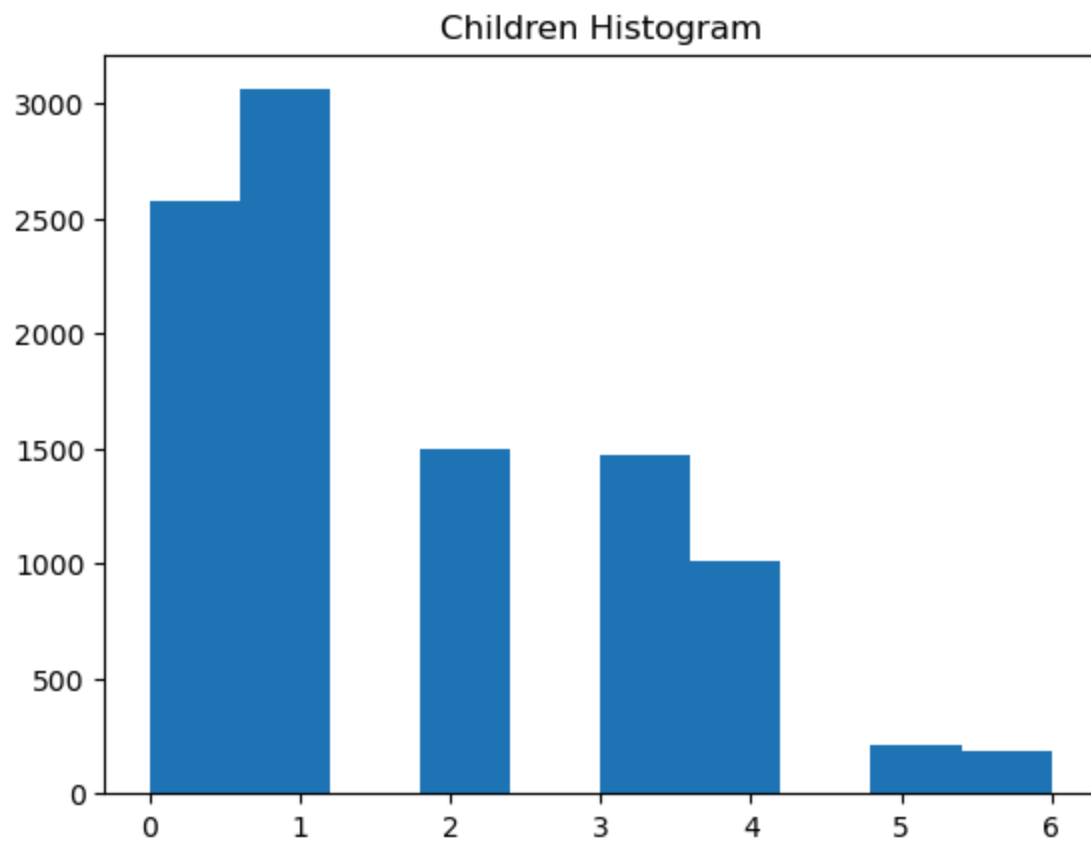
For my univariate visualizations, we can look at the distributions of all of the variables using the actual values this time rather than the z-scores. Since we are using quantitative variables for the multiple linear regression analysis, I generate scatterplots for each of the variables for my bivariate visualizations

```

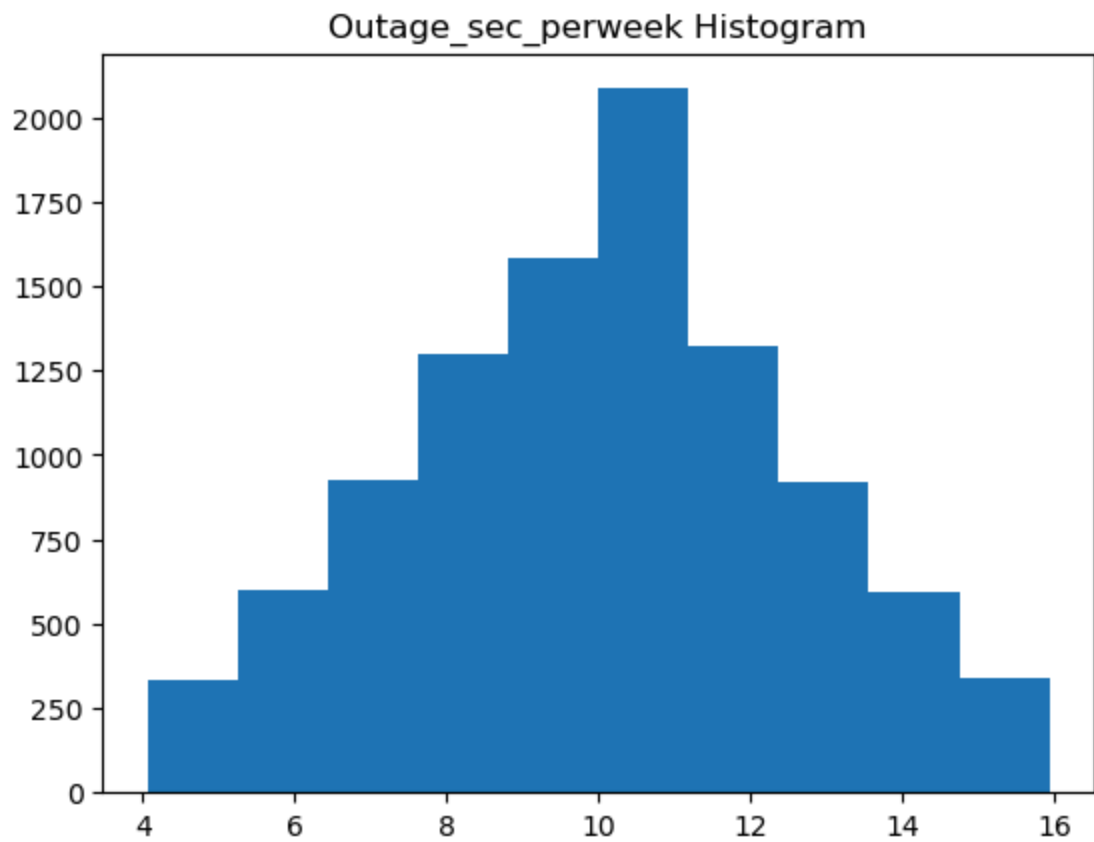
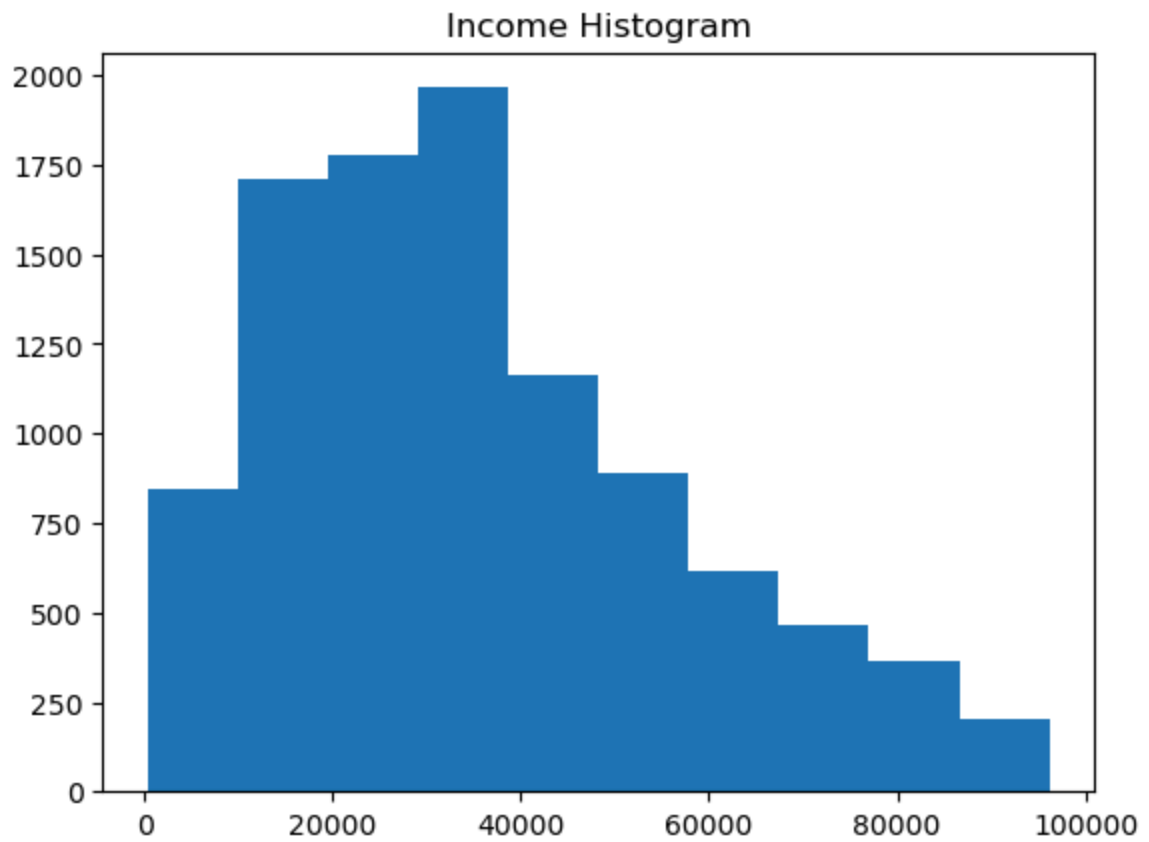
In [20]: # graph our univariate visualizations
for column in dfq_c:
    plt.hist(dfq[column])
    plt.title(column+' Histogram')
    plt.show()

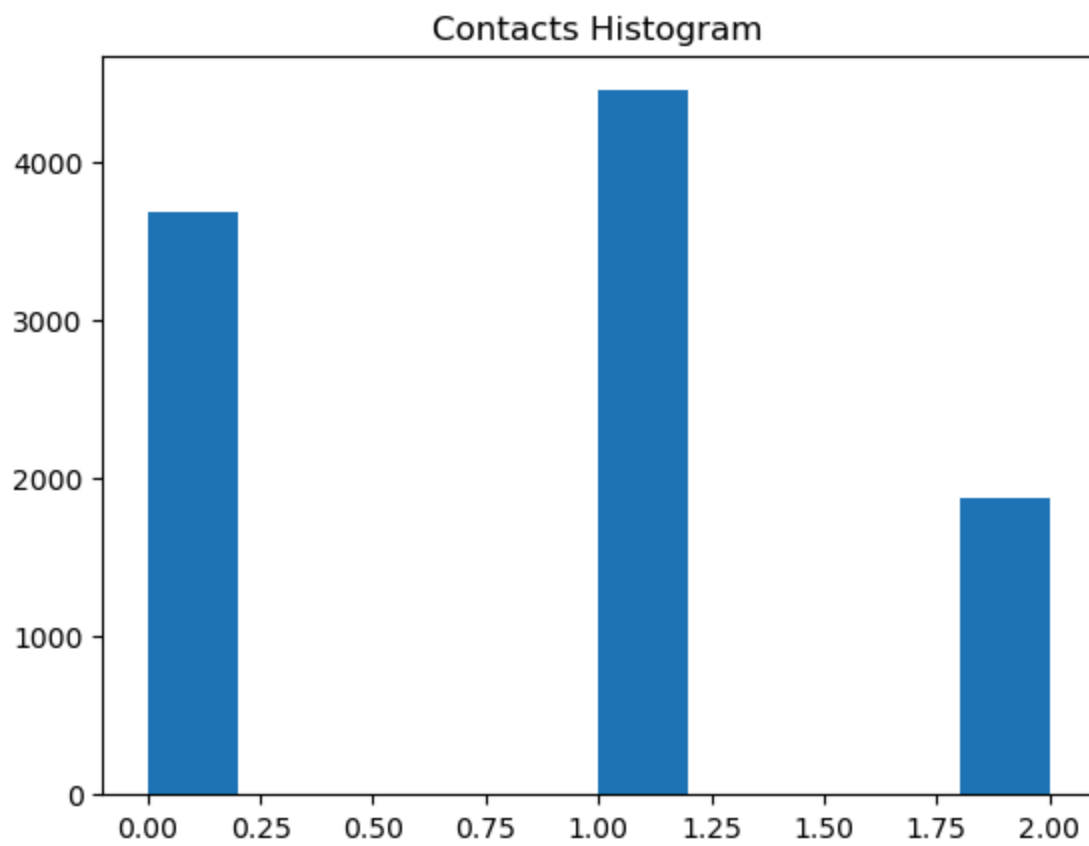
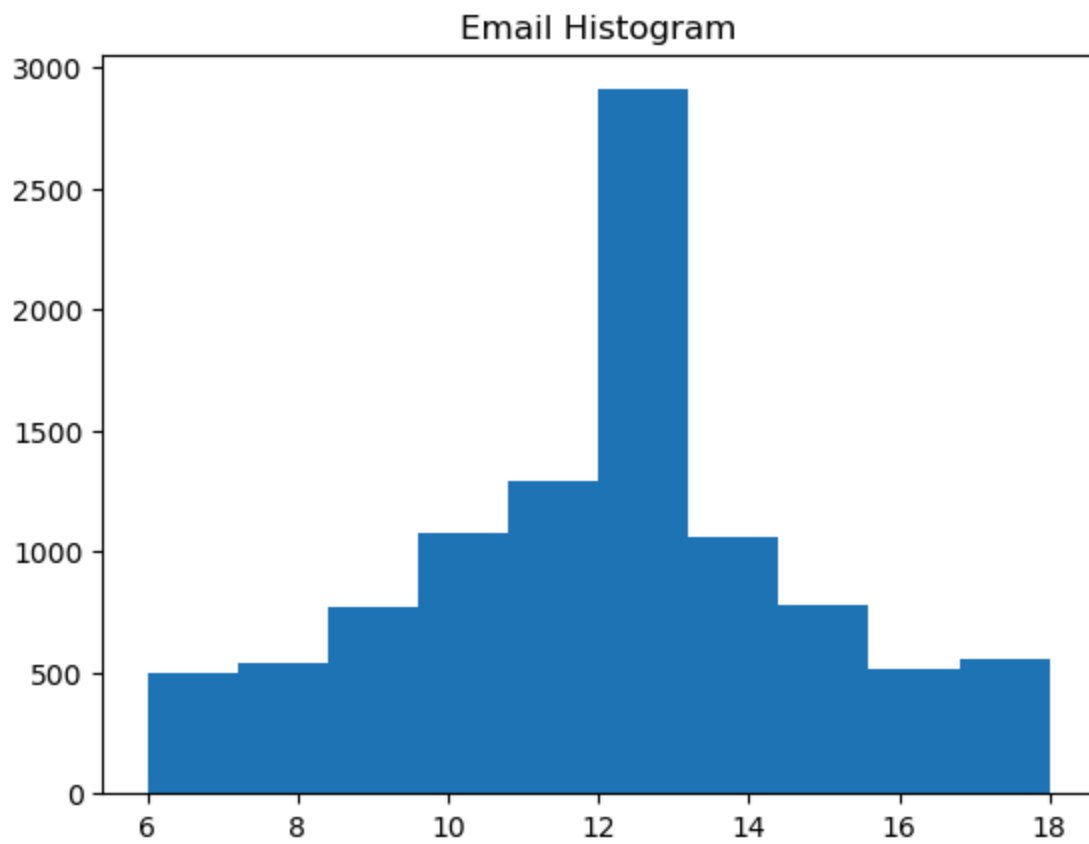
```

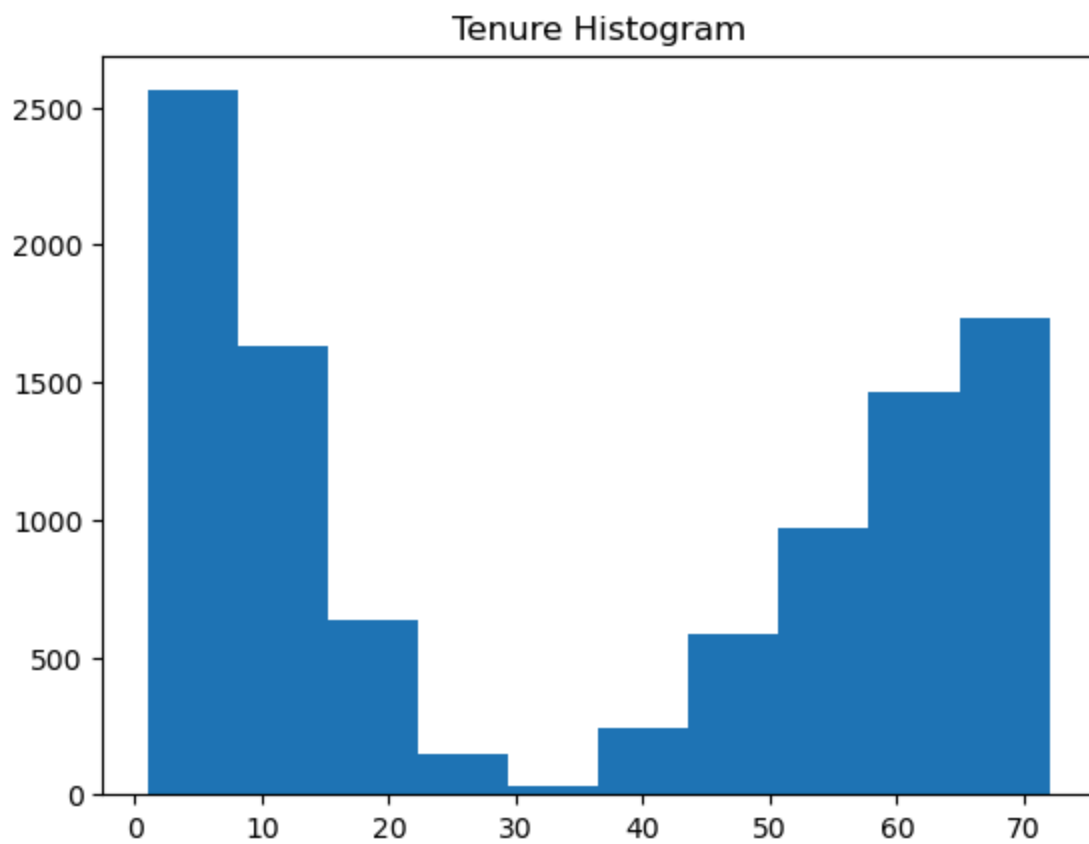
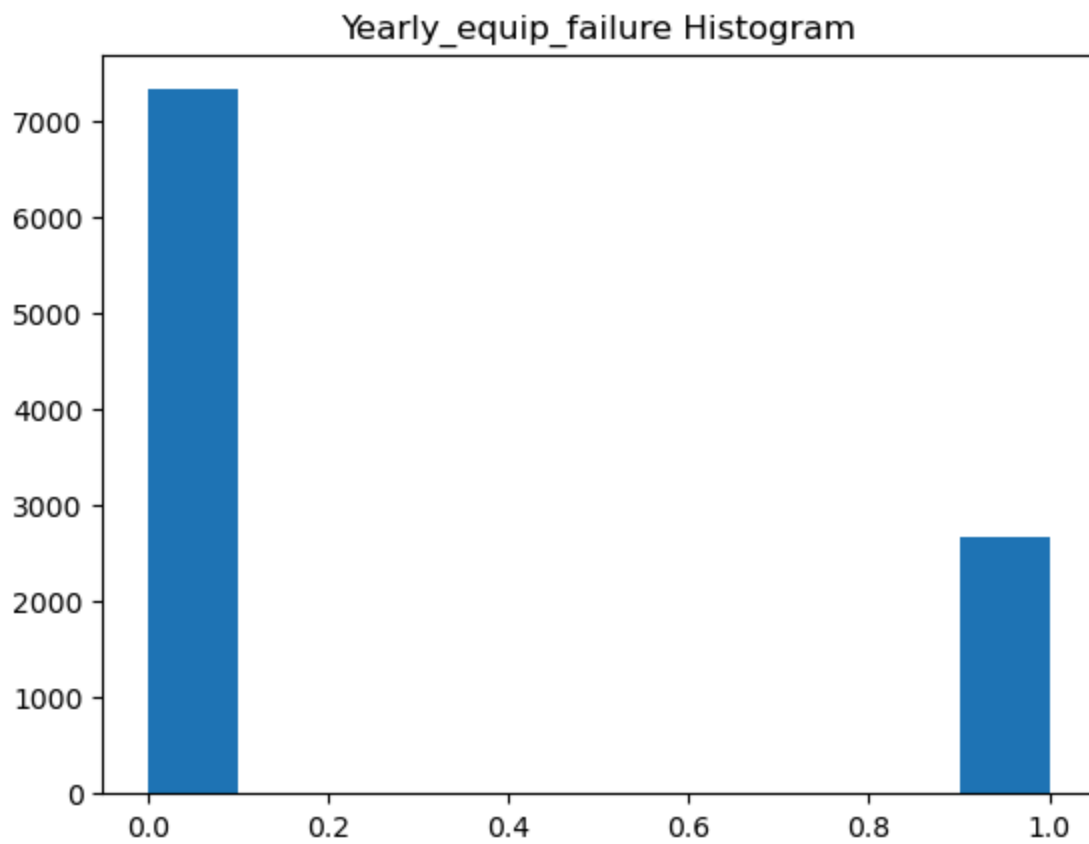


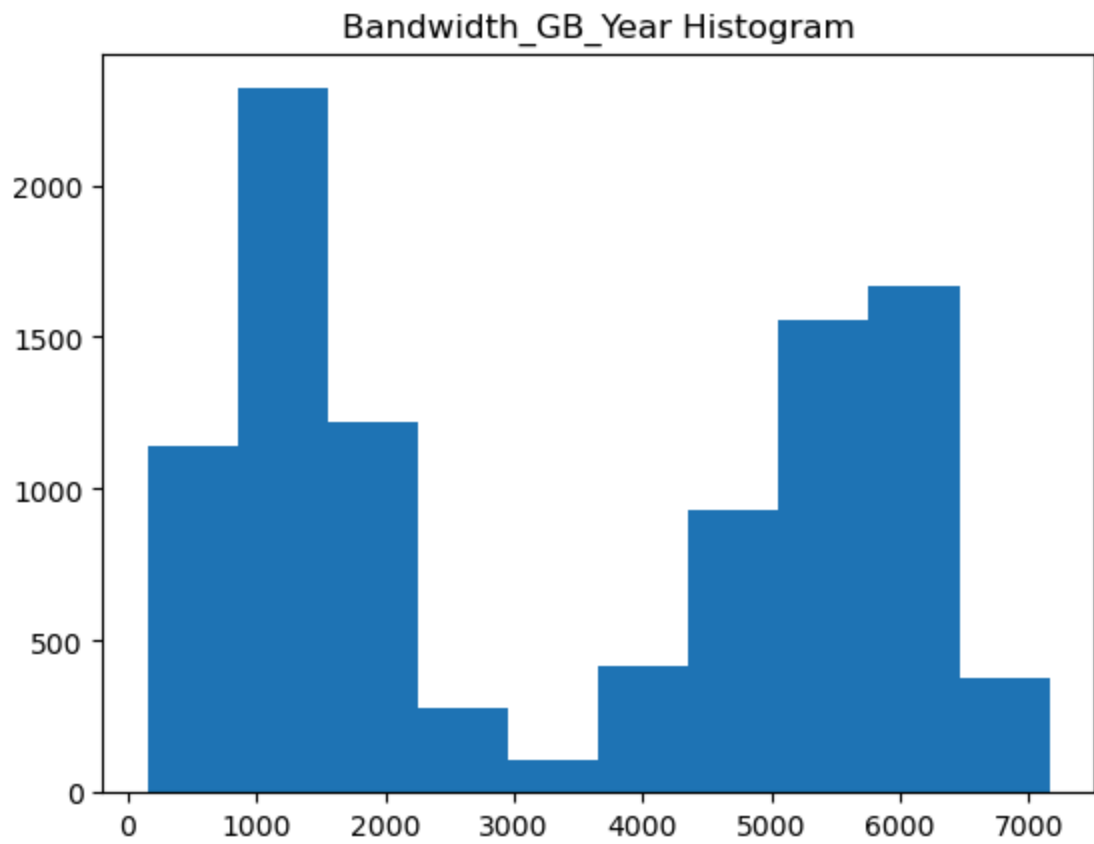
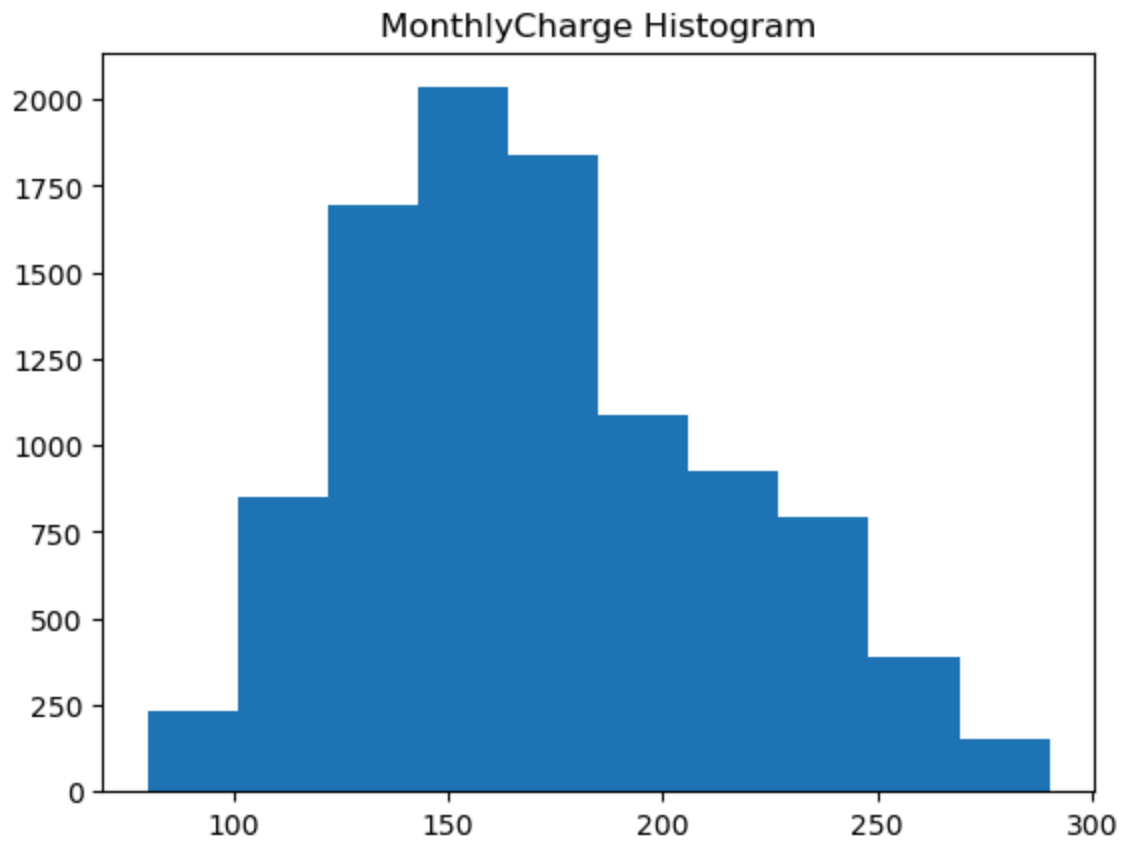


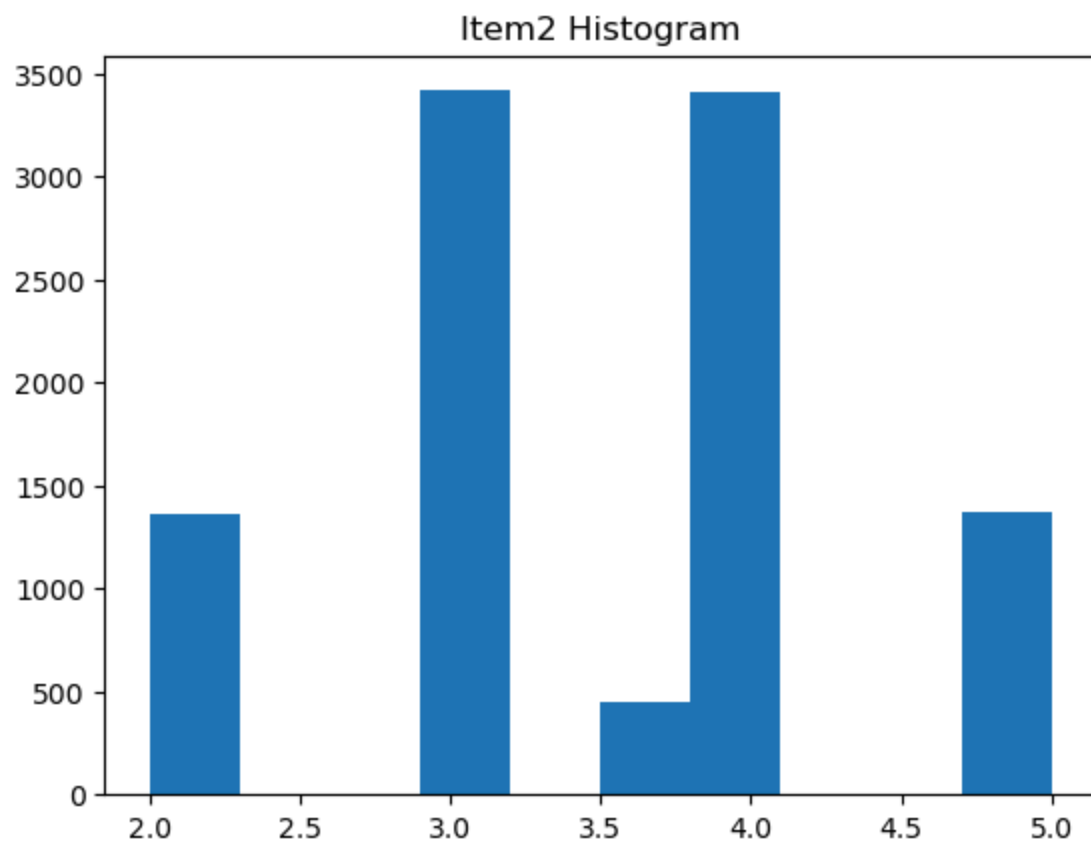
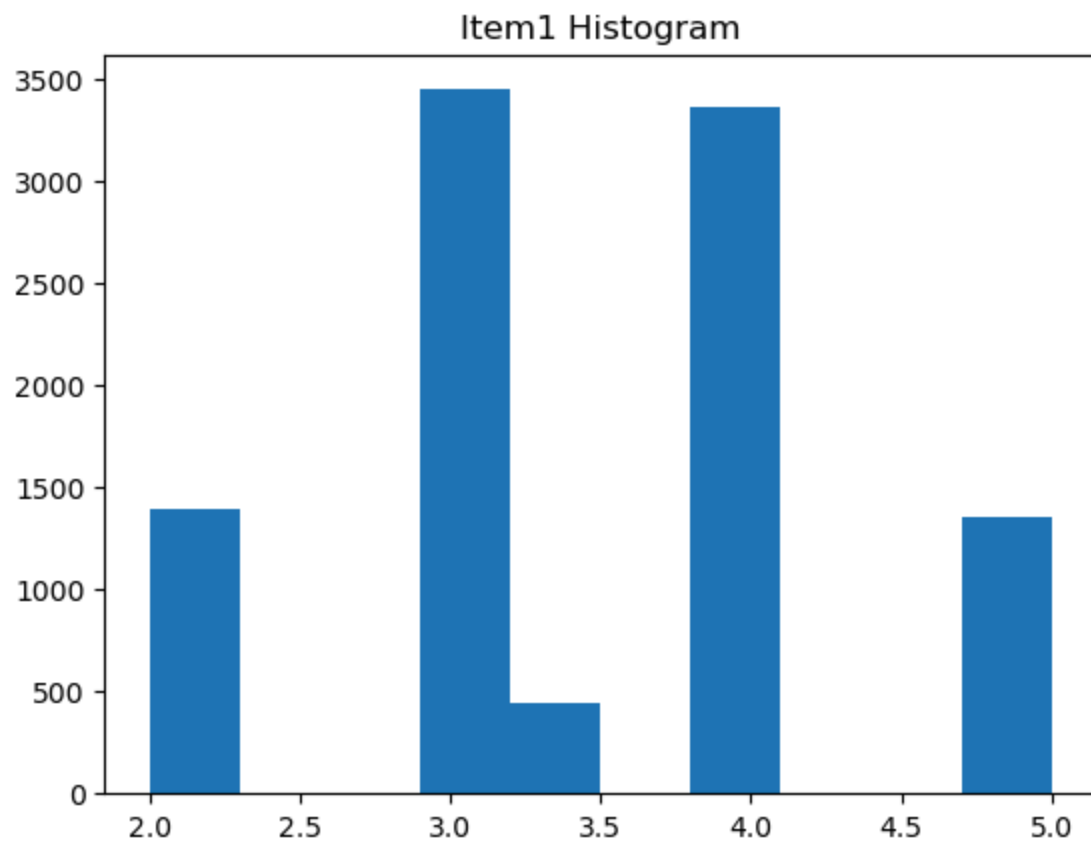


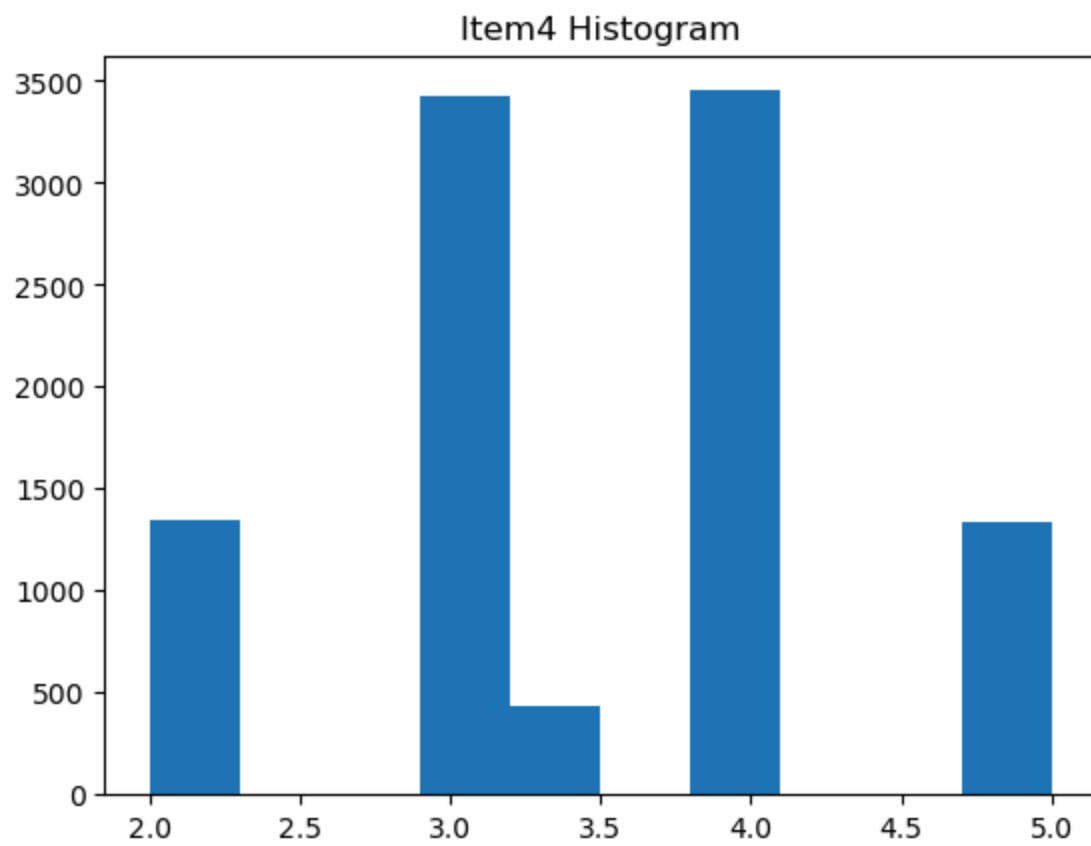
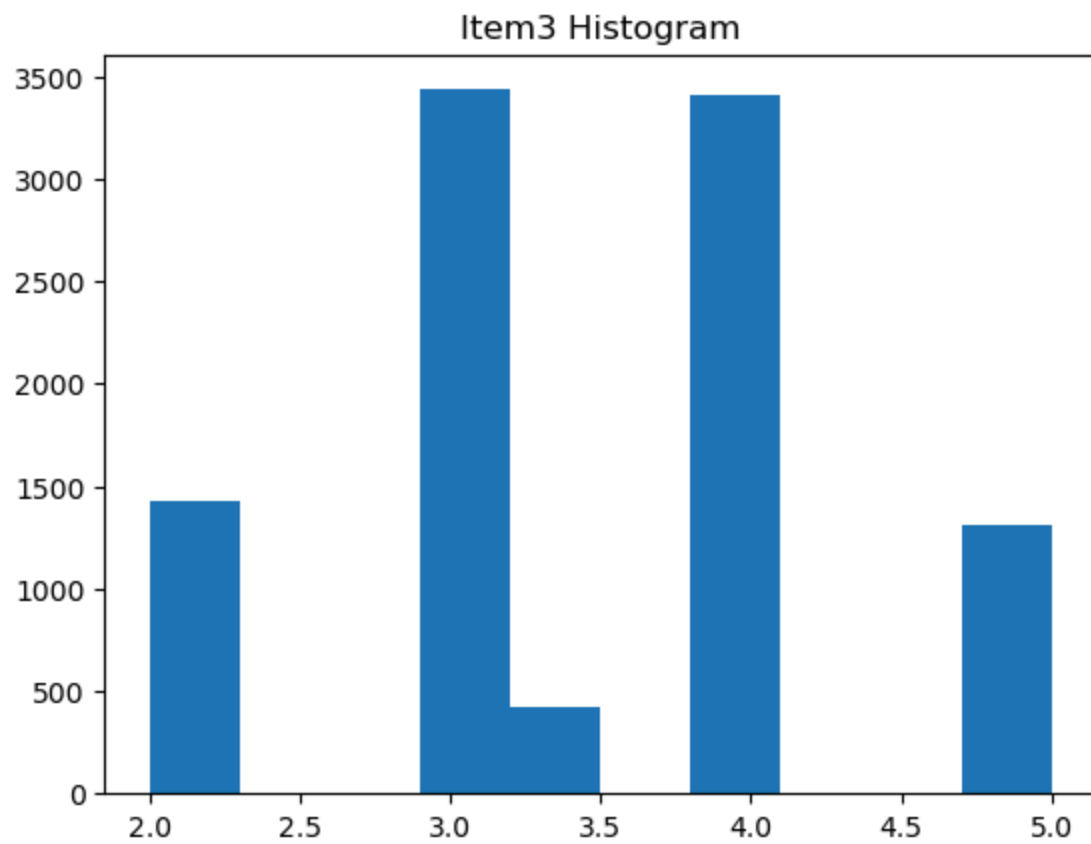


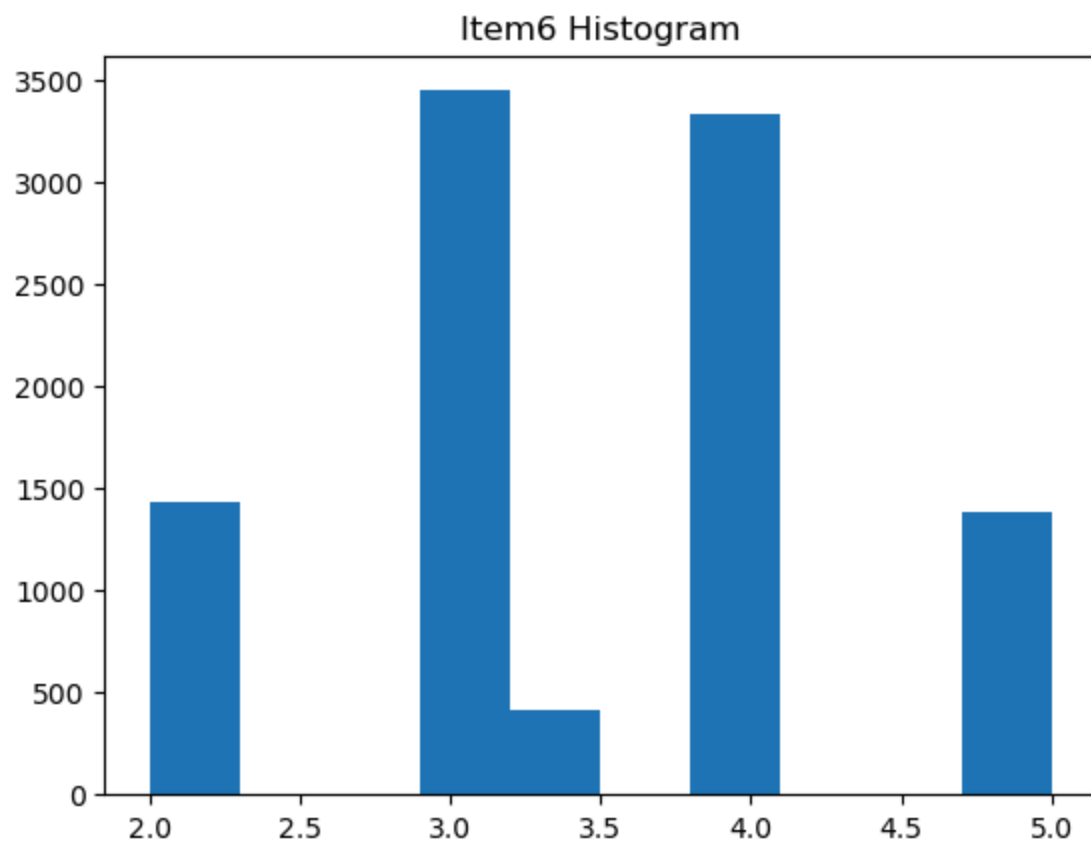
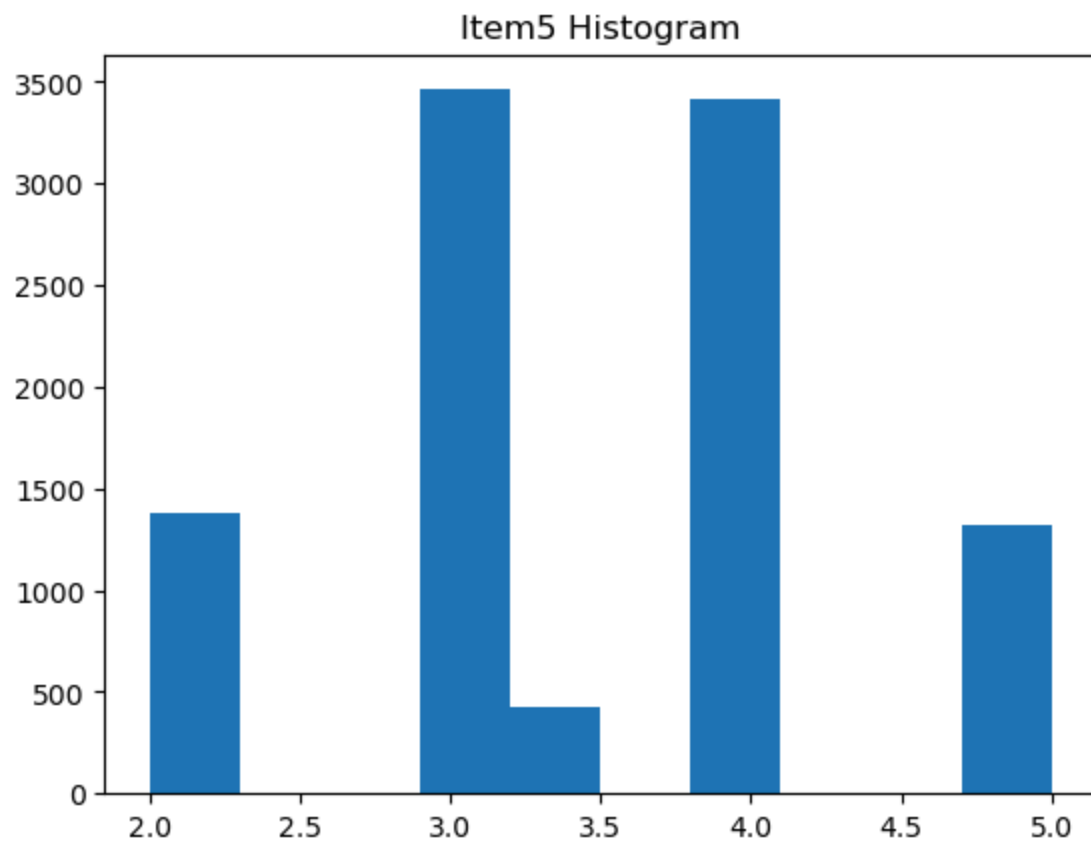


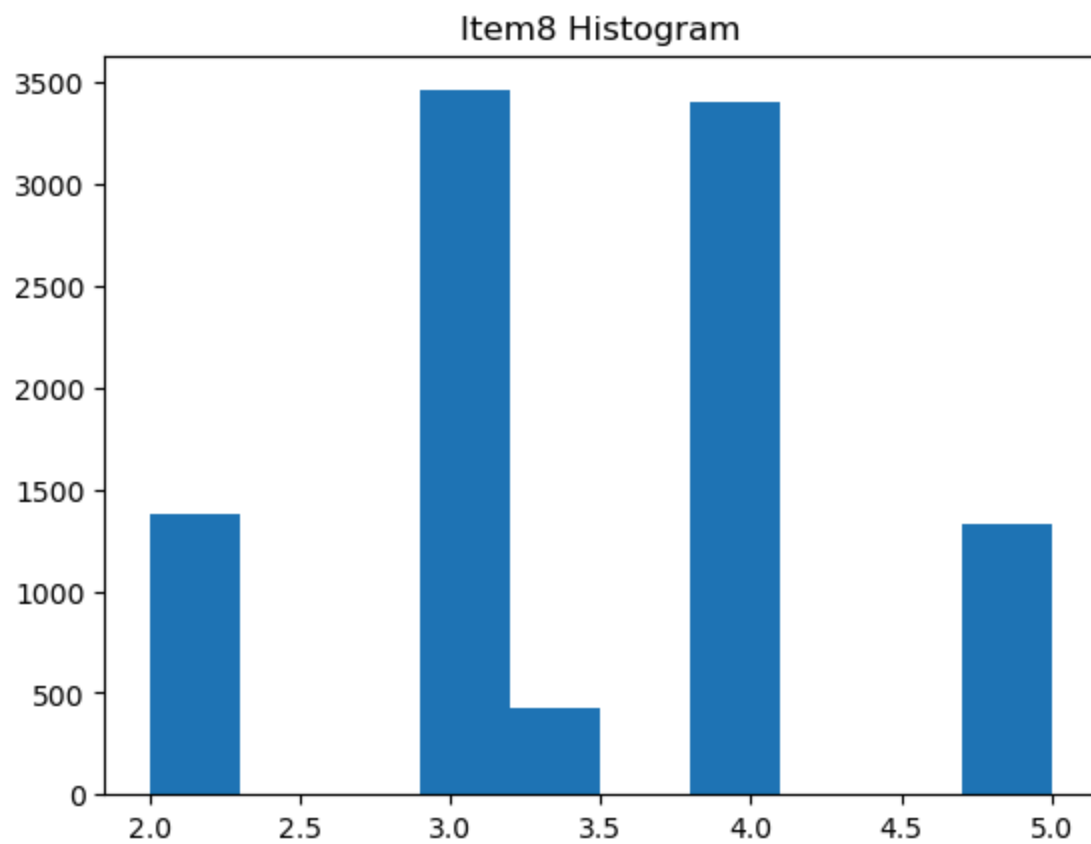
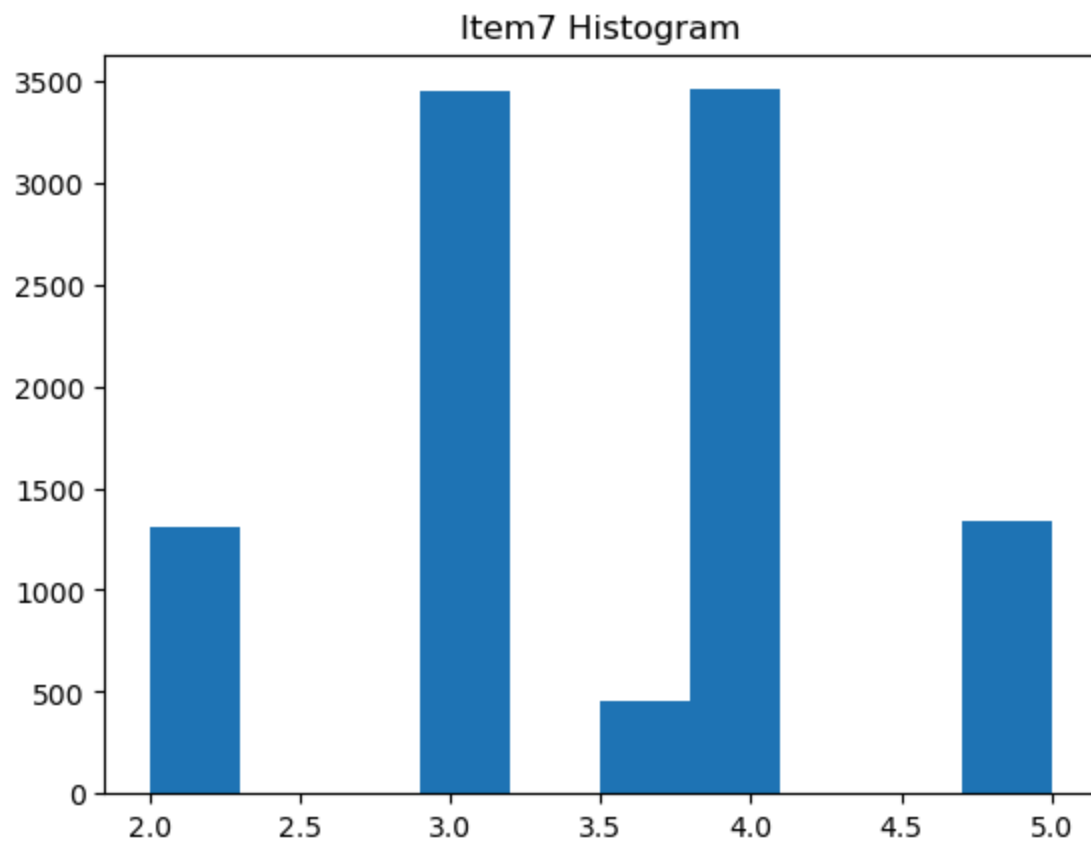








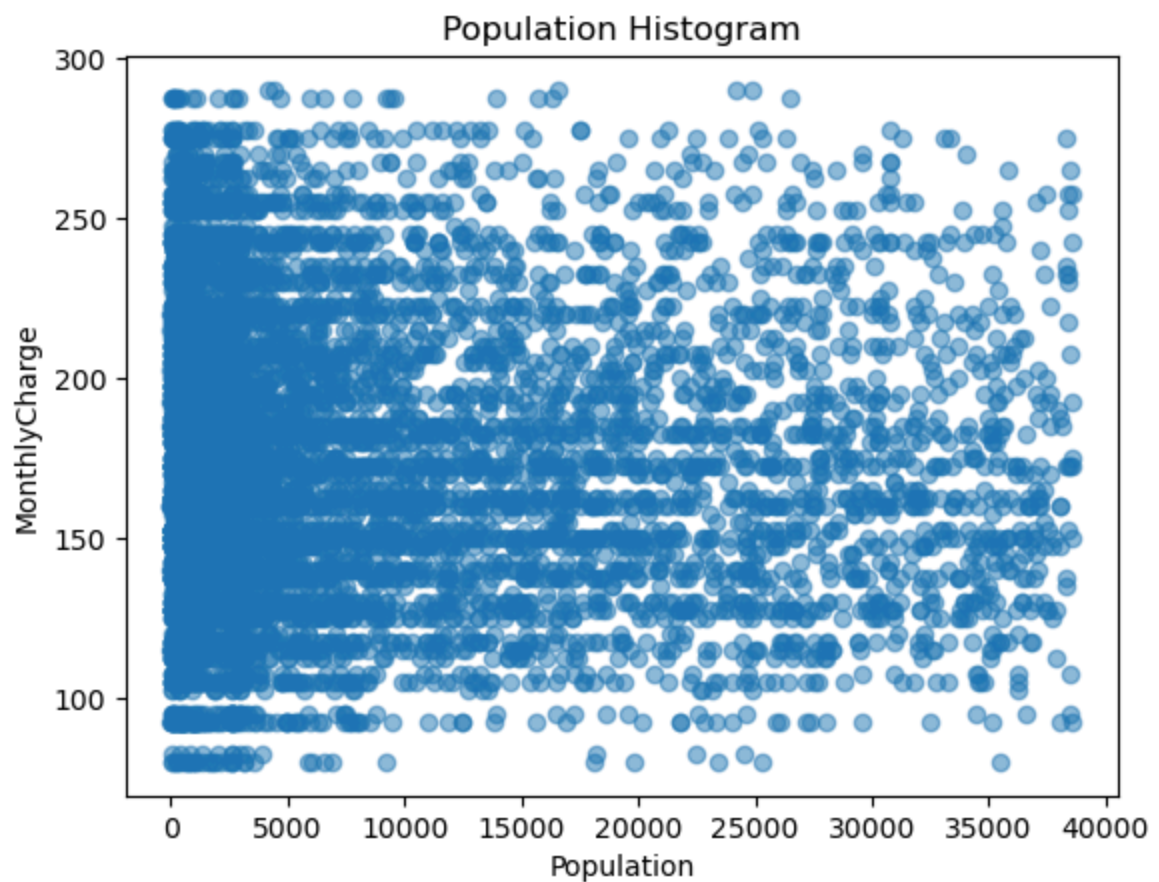
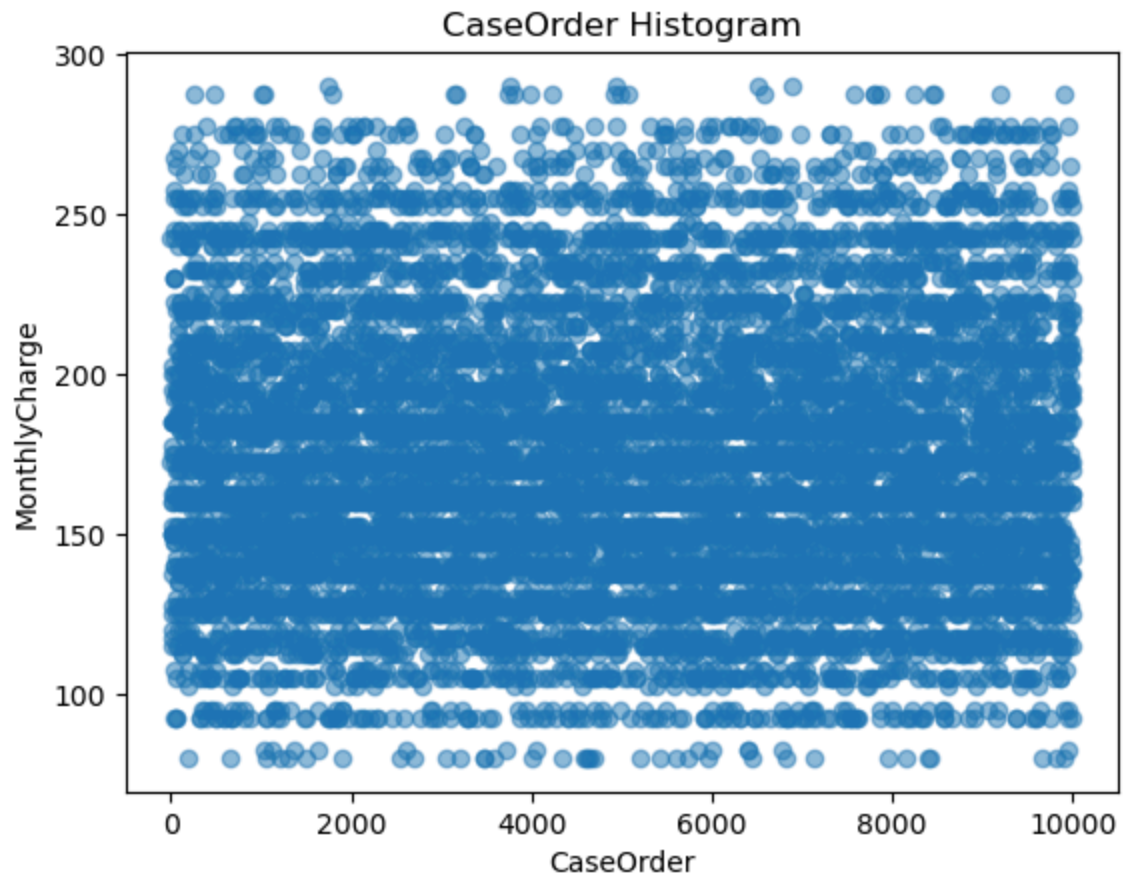


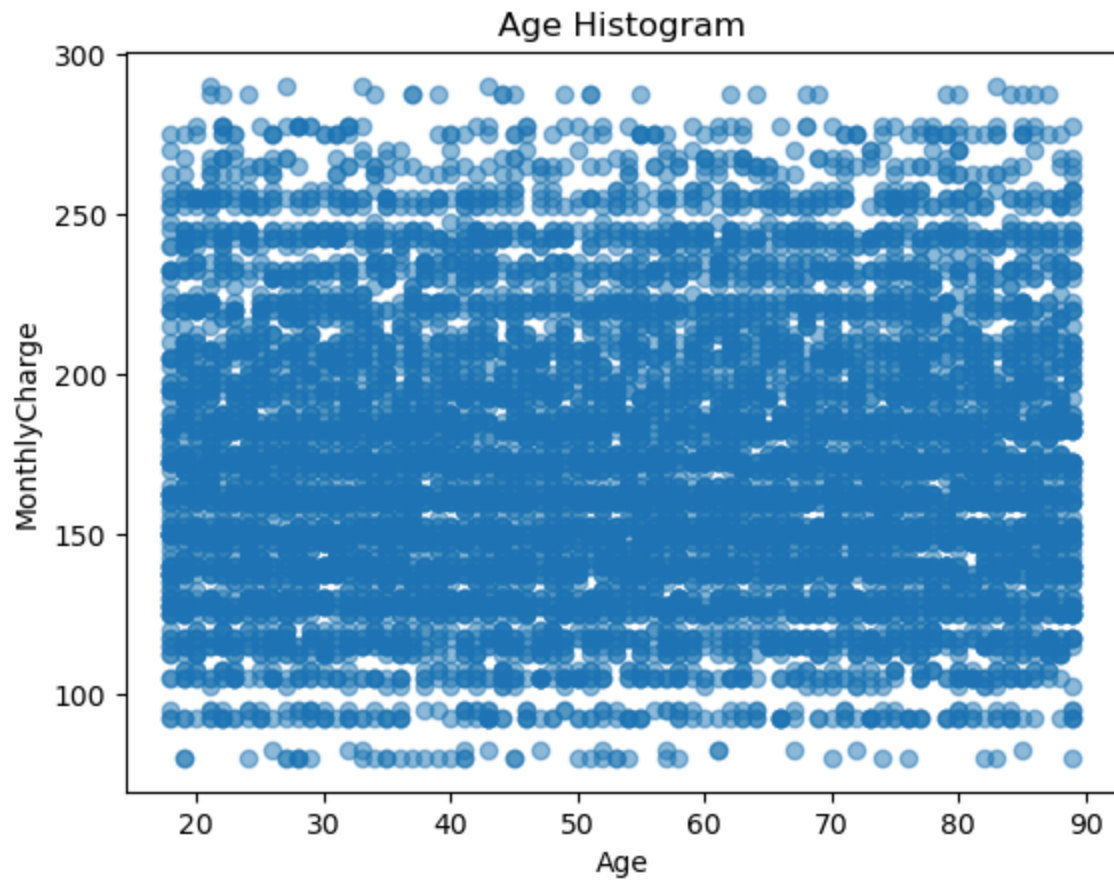
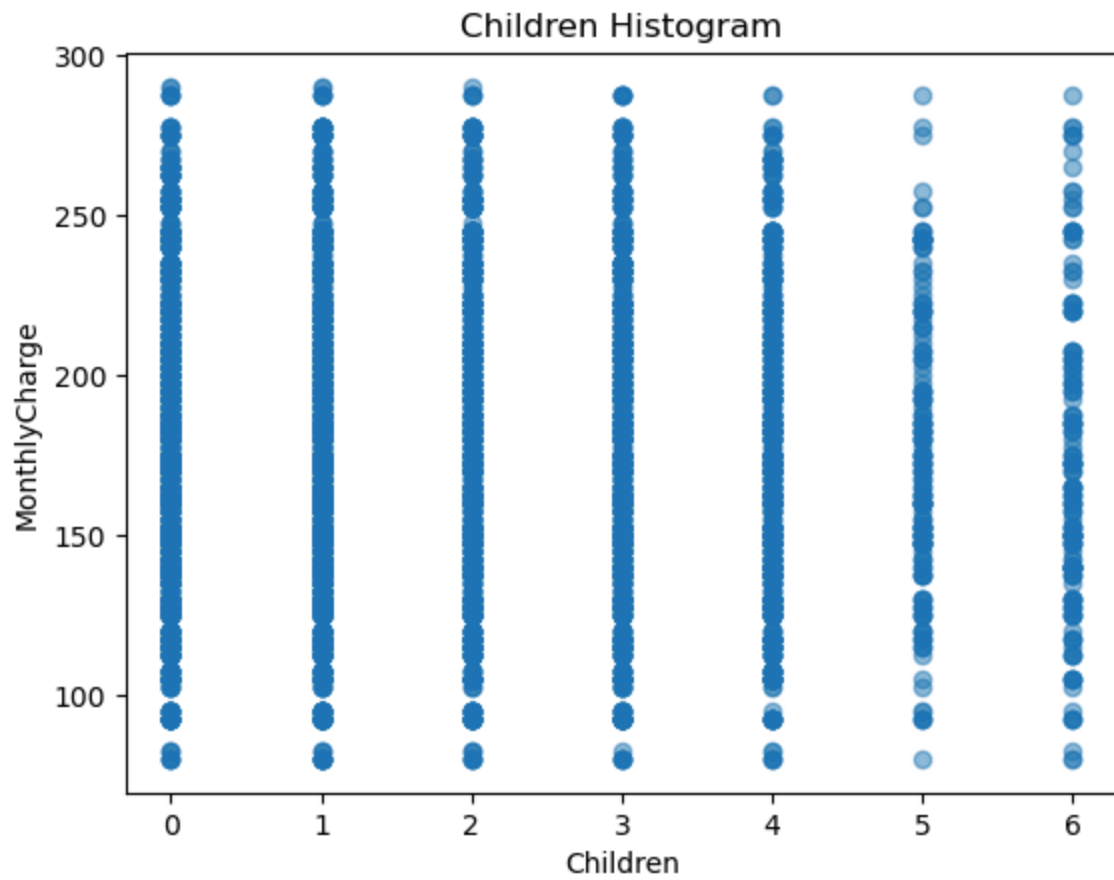


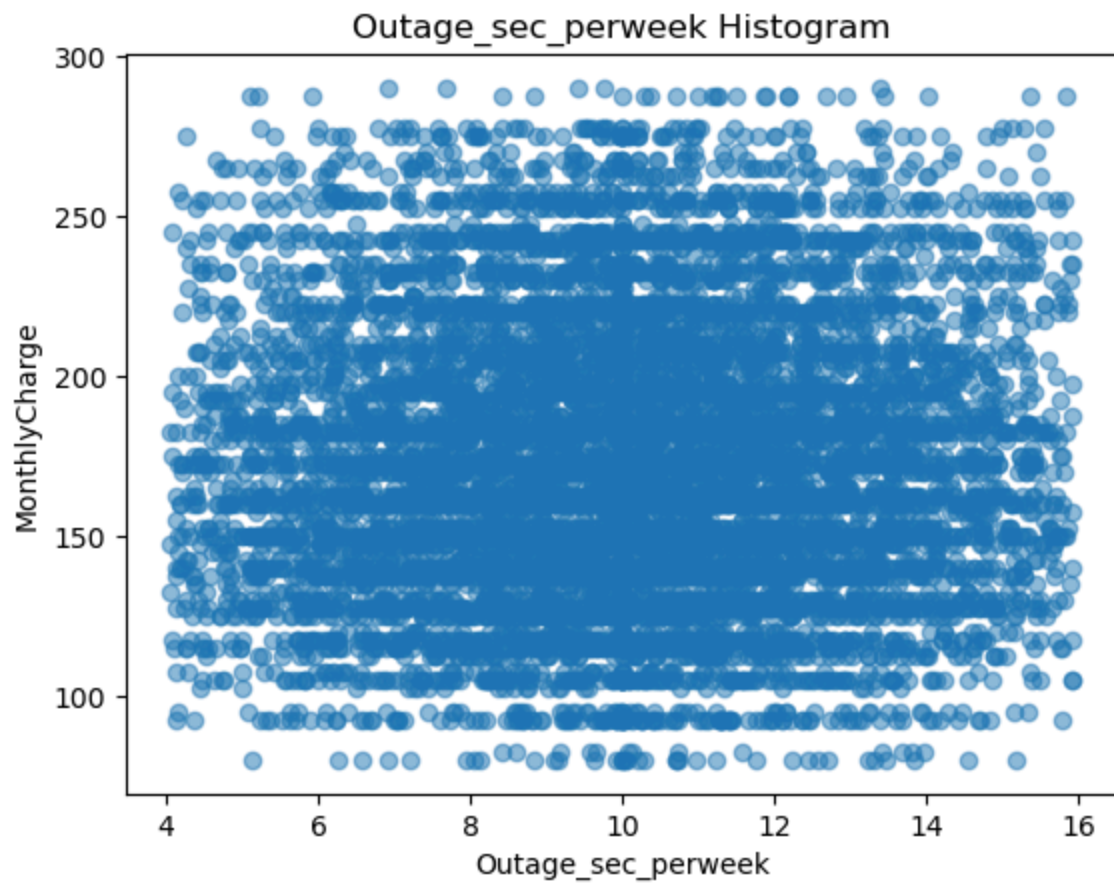
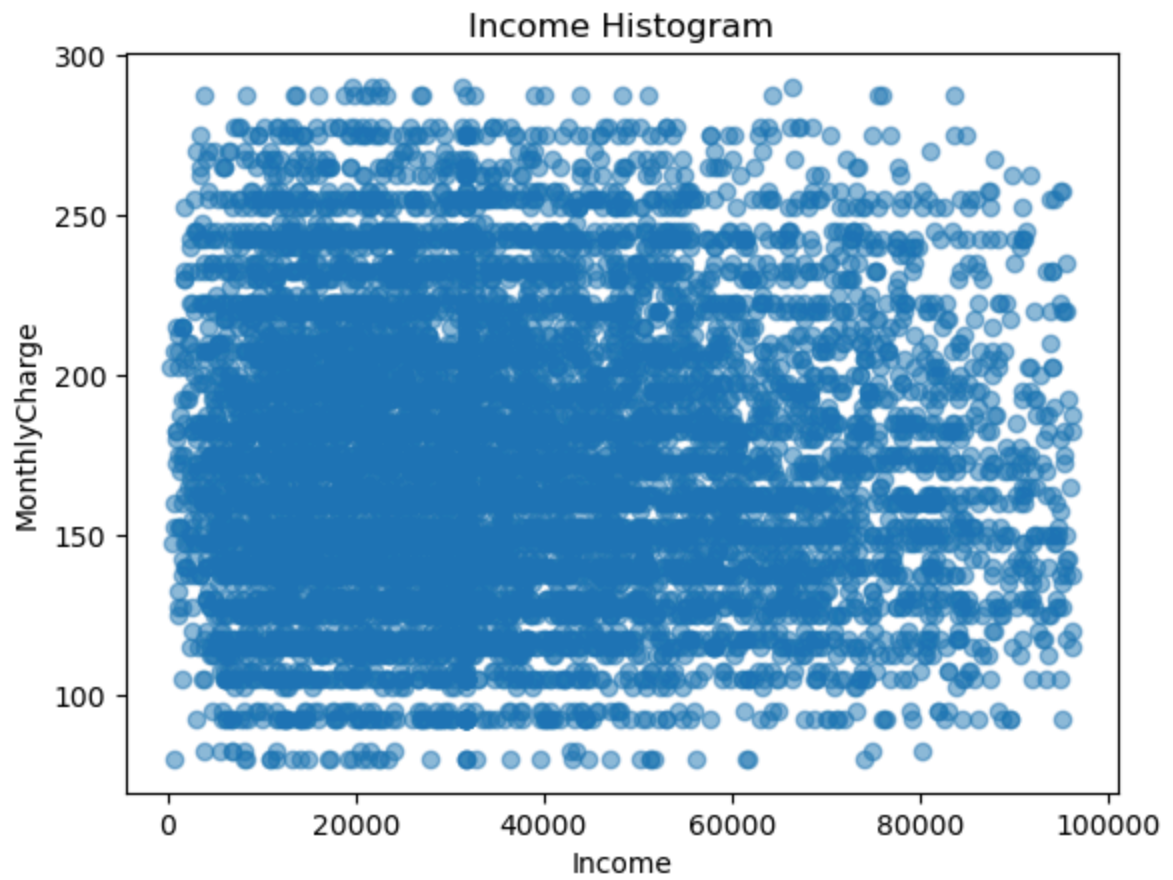
```
In [21]: for column in dfq_c:
          x = dfq[column]
          y = dfq['MonthlyCharge']
          plt.scatter(x, y, alpha=0.5)
          plt.title(column + ' Histogram')
```

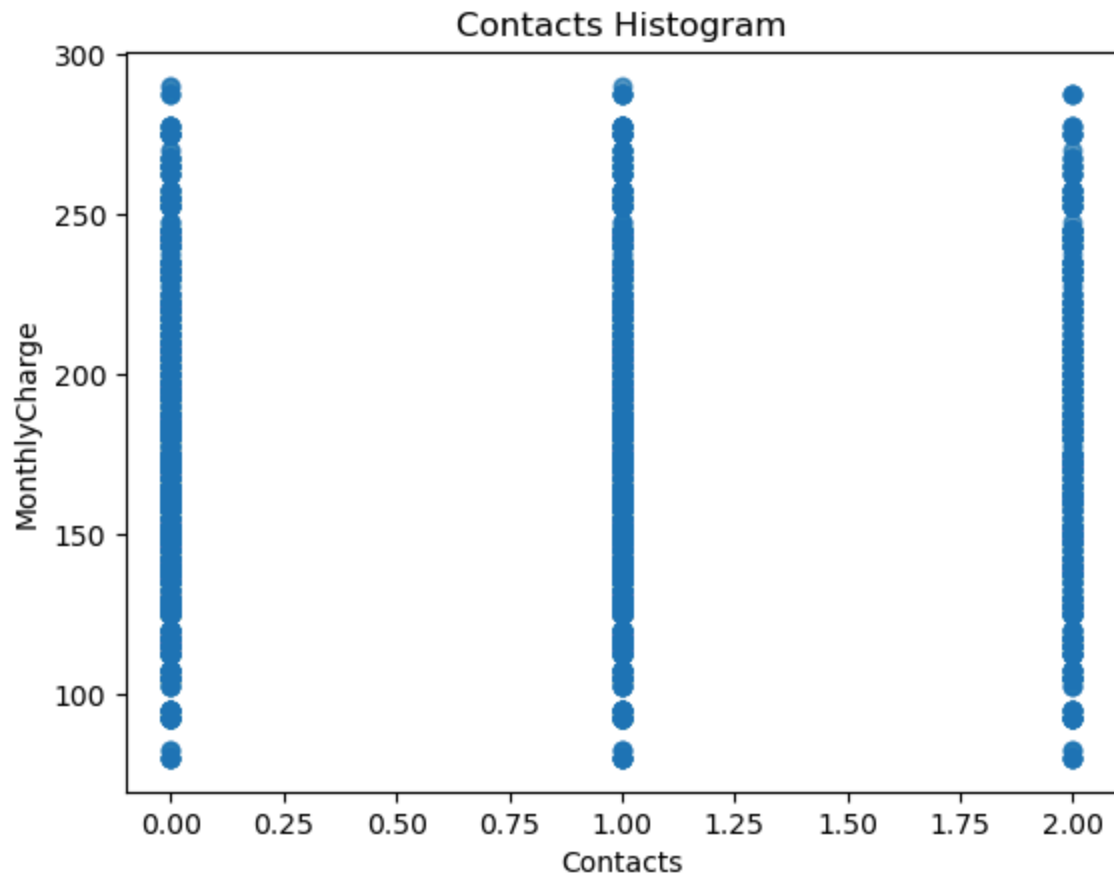
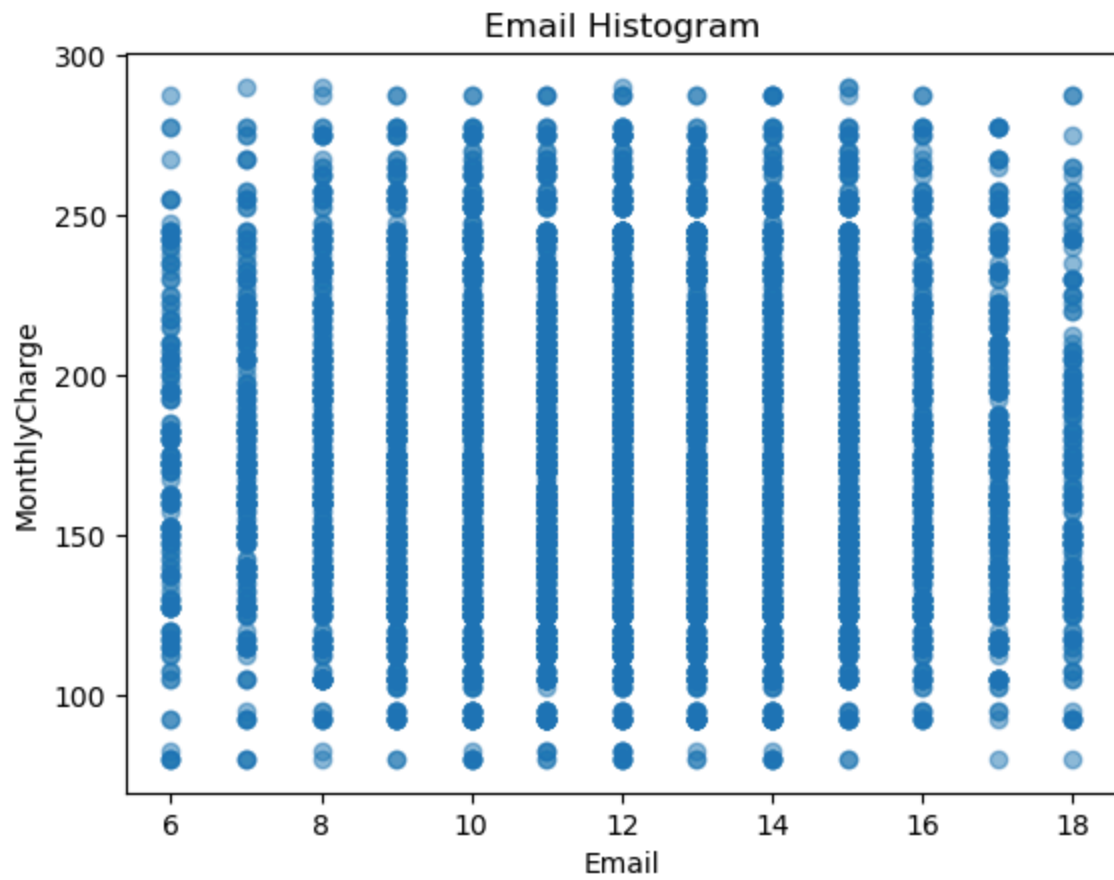


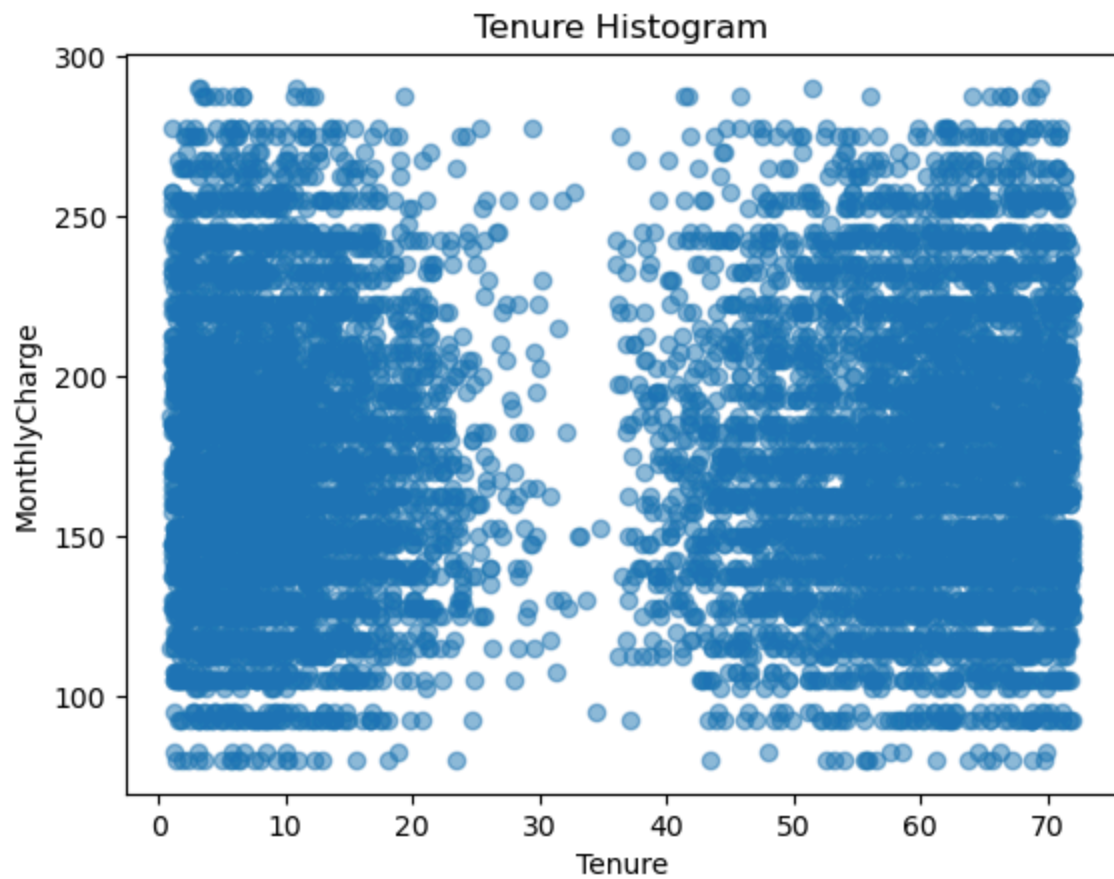
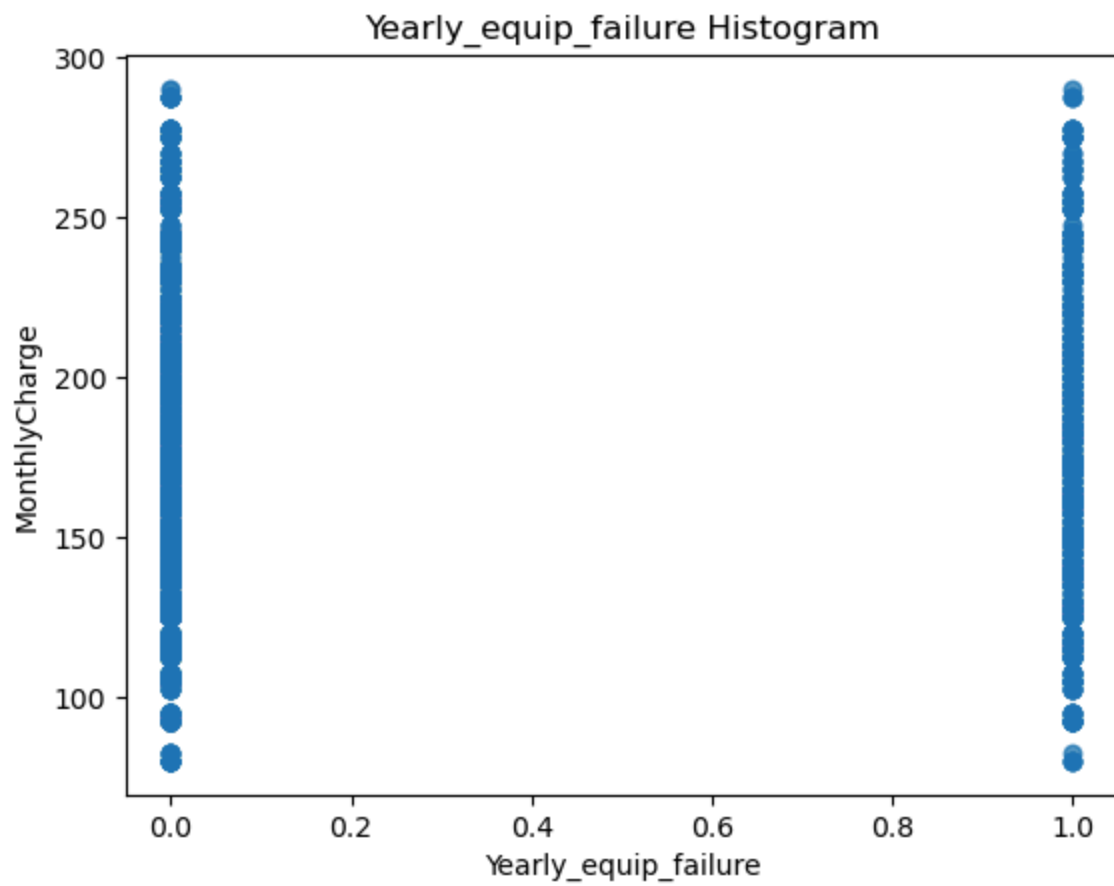
```
plt.xlabel(column)
plt.ylabel("MonthlyCharge")
plt.show()
```



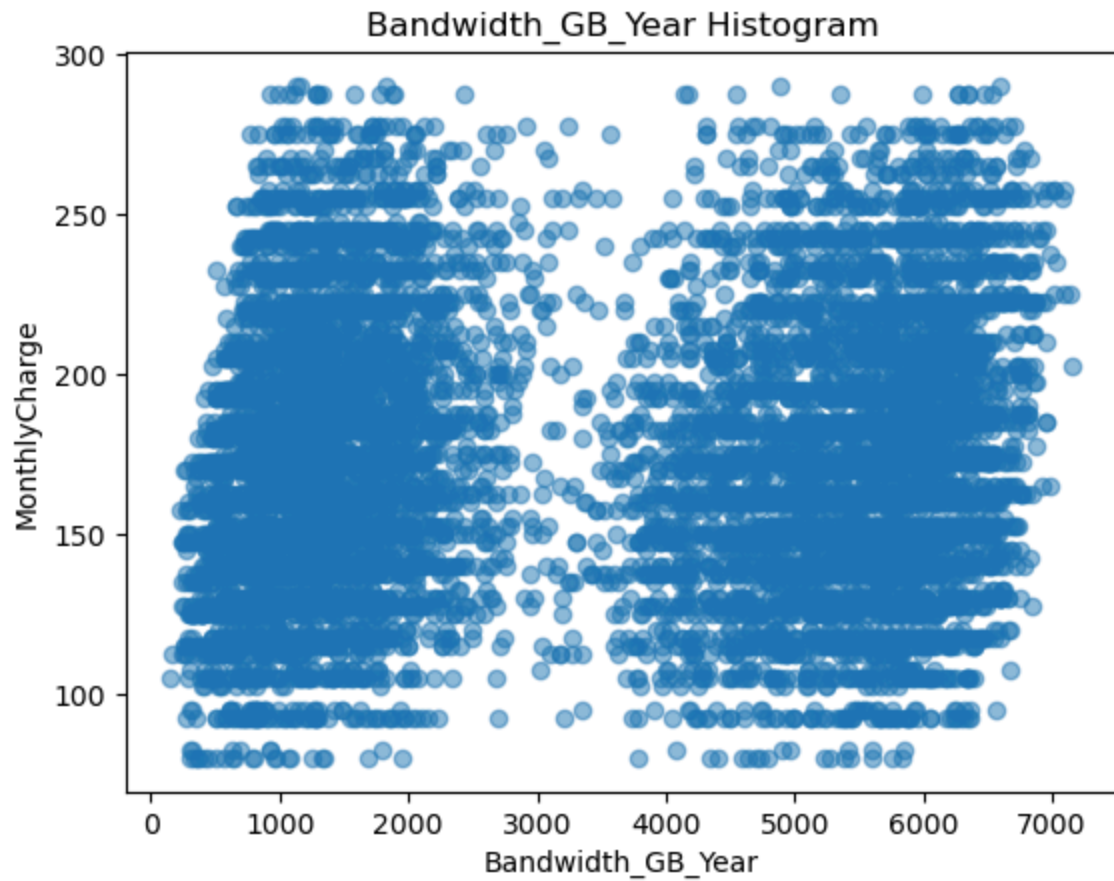
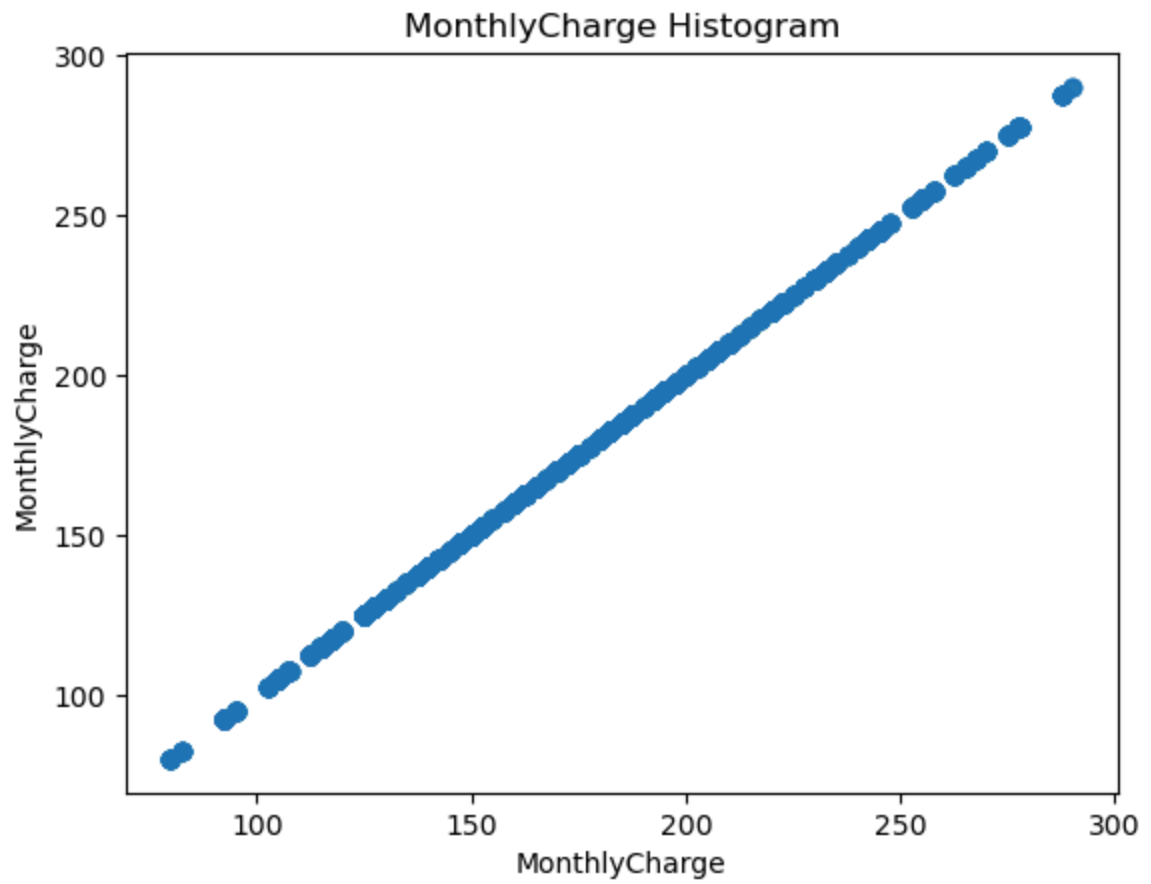


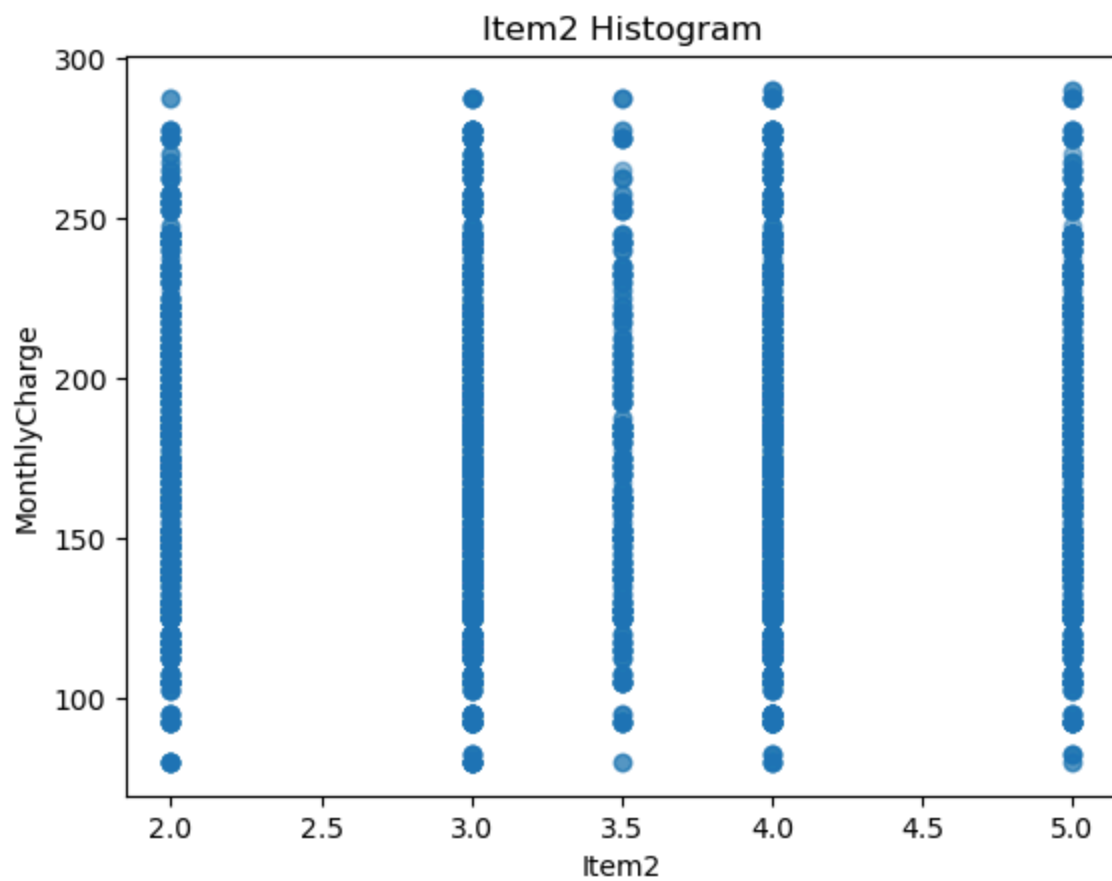
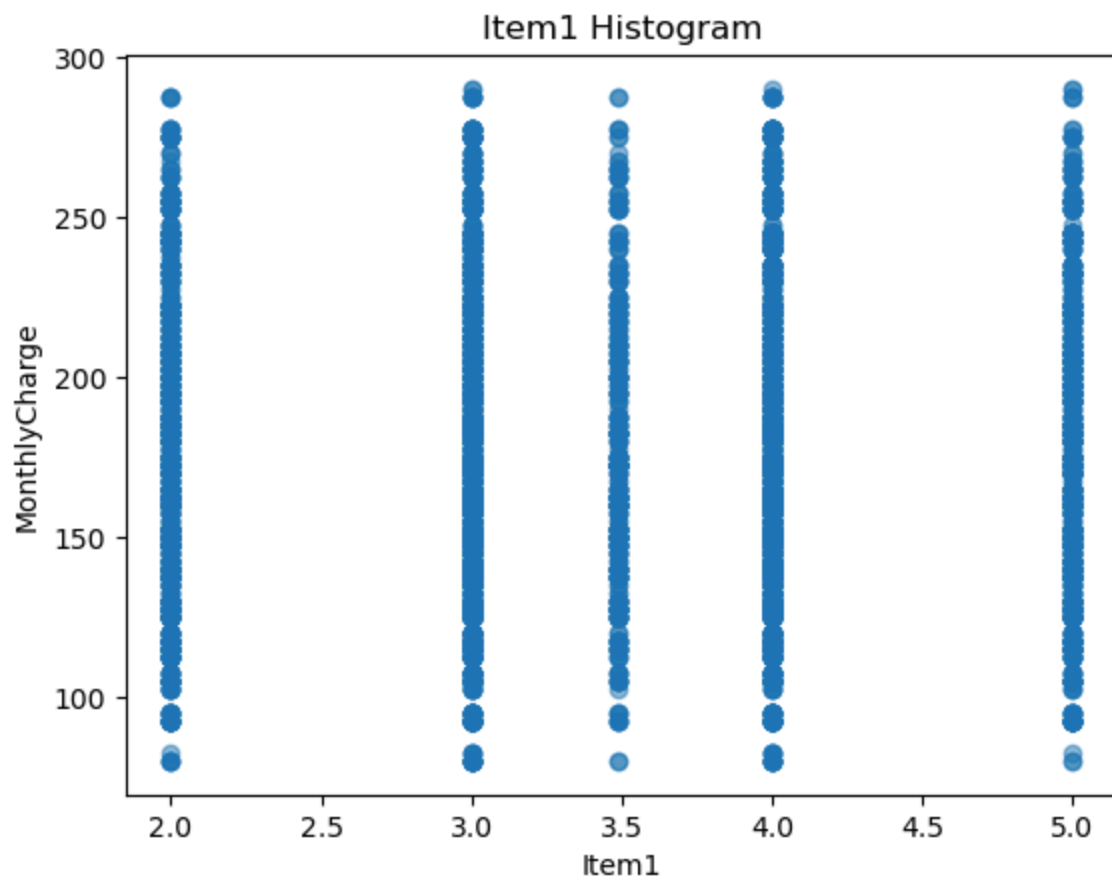


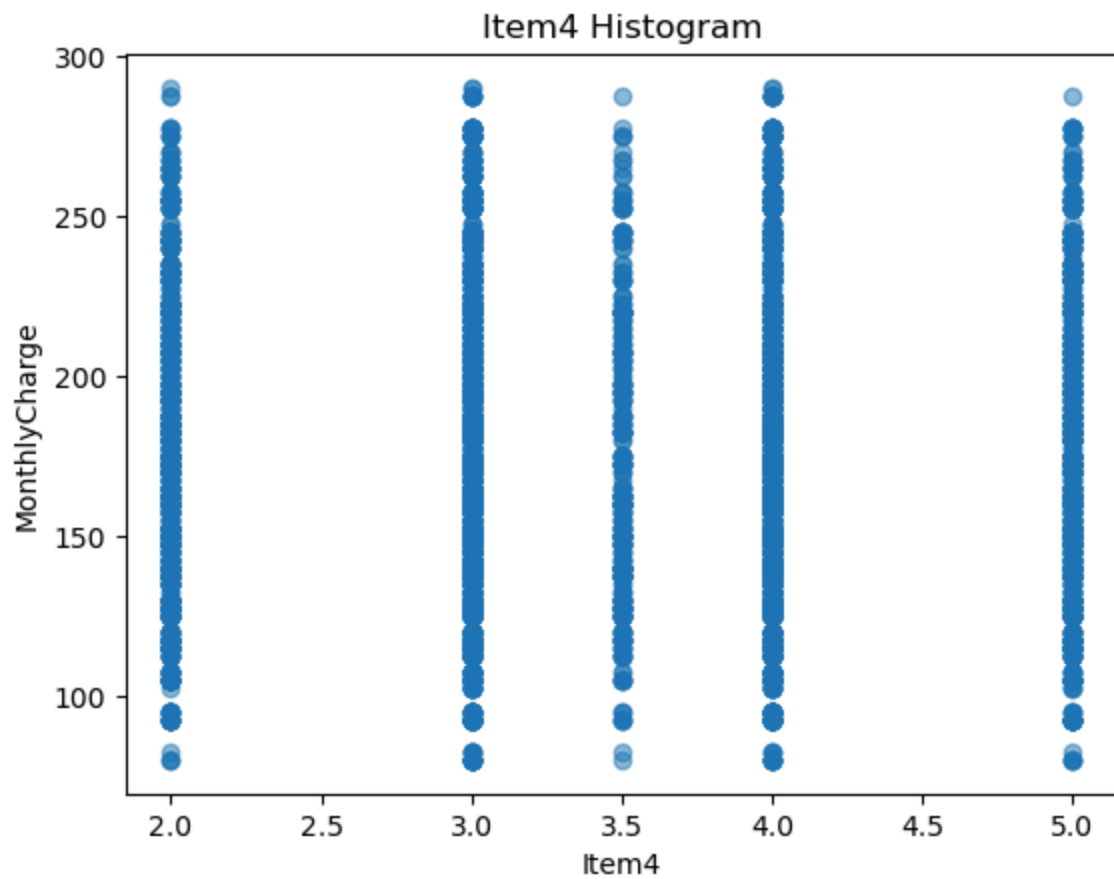
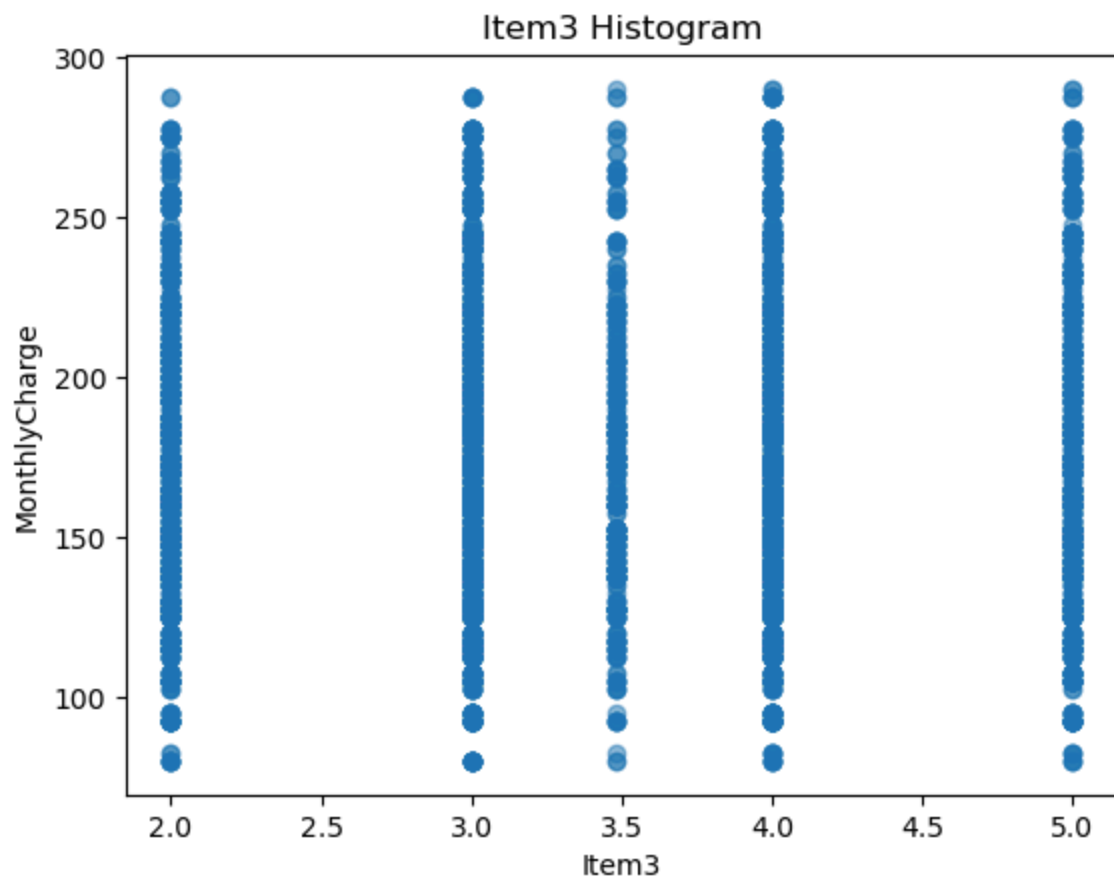




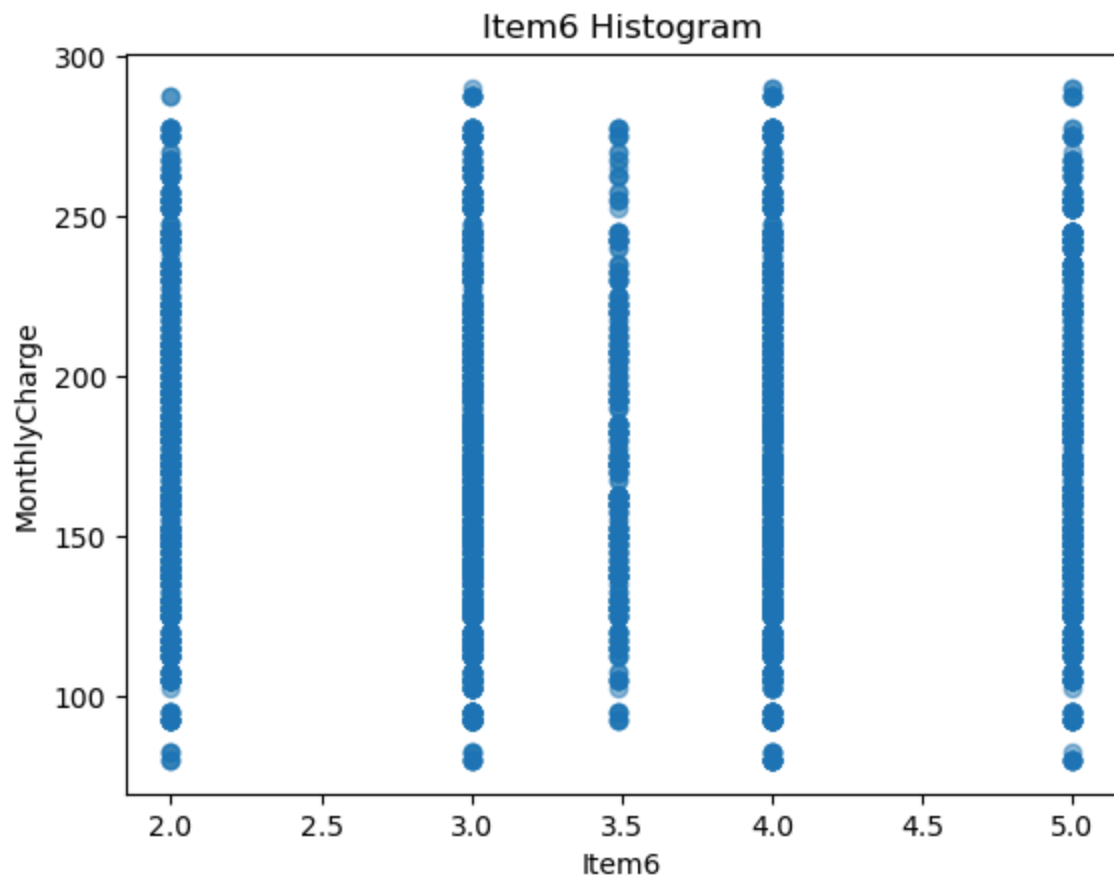
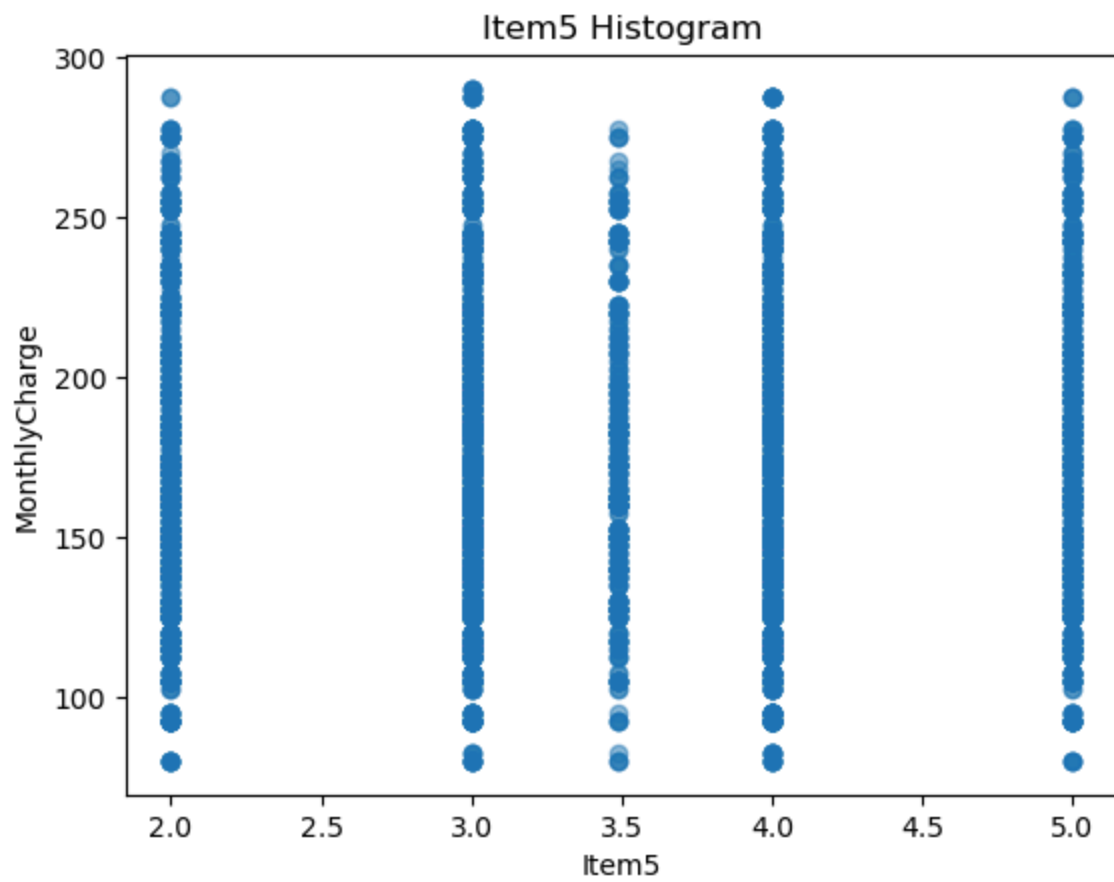


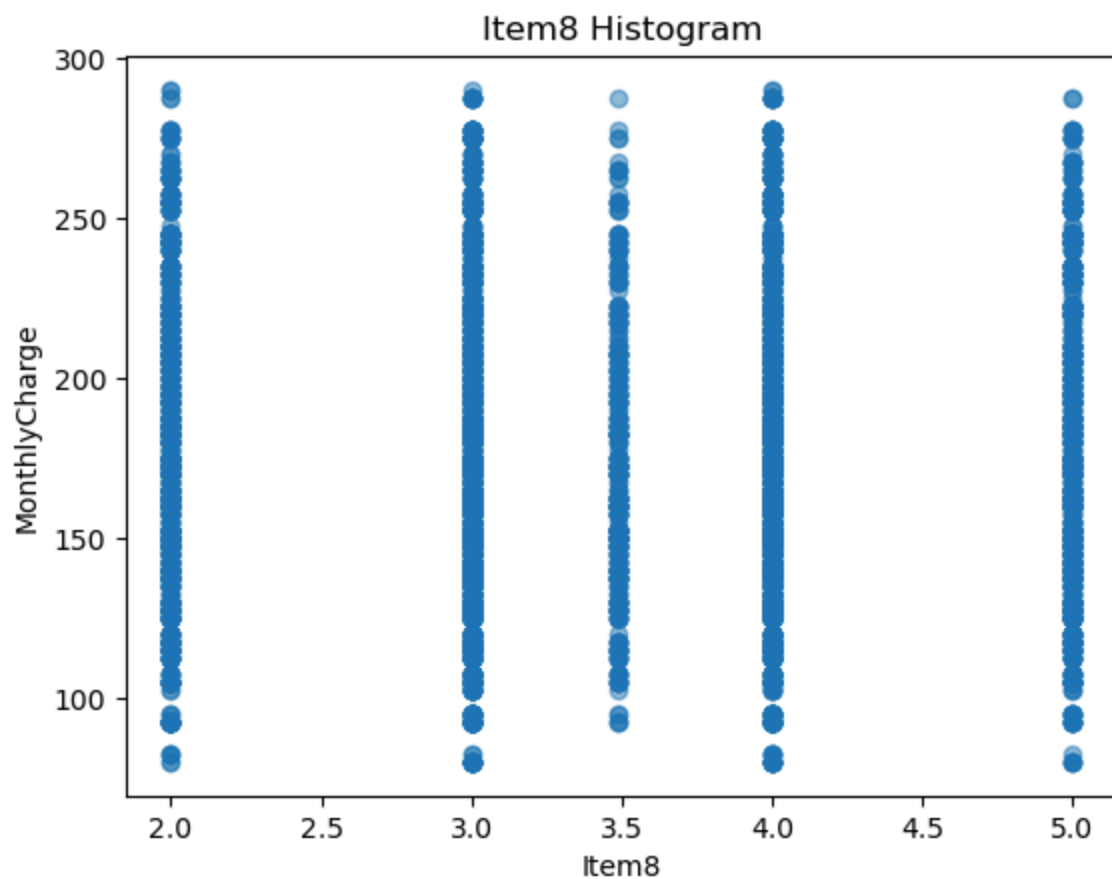
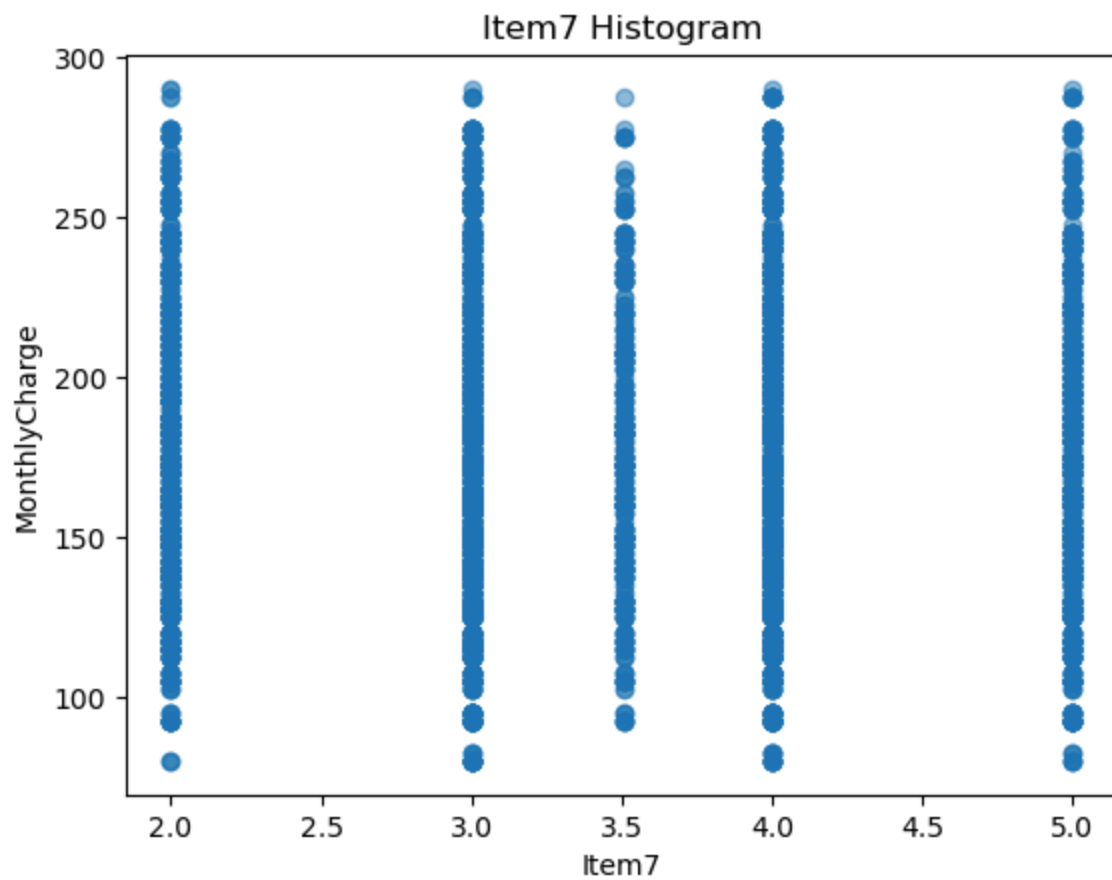












## C4. Data Transformation (Data Wrangling)

Due to the variables that were selected, no re-expression of categorical variables was necessary. I kept my analysis strictly to the quantitative variables for this test.

For the data transformations, I did data cleaning earlier in the analysis. This included removing nulls from population and removing outliers from all other variables.

## C5. Prepared Dataset

Here, I export the prepared dataset as a .csv file

```
In [22]: dfq.to_csv('prepared_data_task1.csv')
```

## D1. Initial Model

Here I construct an initial multiple linear regression model

```
In [23]: # set the dependent variable as monthlycharge and the independent variables as
y = dfq['MonthlyCharge']
X = dfq[['Population', 'Children', 'Income', 'Contacts', 'Yearly equip_failure

# fit the model to x and Y
model = sm.OLS(y, X)
results = model.fit()

#print the results
print(results.summary())
```

## OLS Regression Results

```

=====
Dep. Variable:          MonthlyCharge   R-squared (uncentered):
0.935
Model:                  OLS           Adj. R-squared (uncentered):
0.935
Method:                 Least Squares   F-statistic:
9588.
Date:                   Wed, 08 May 2024   Prob (F-statistic):
0.00
Time:                   16:21:46         Log-Likelihood:
-52328.
No. Observations:       10000           AIC:
1.047e+05
Df Residuals:           9985           BIC:
1.048e+05
Df Model:               15
Covariance Type:        nonrobust
=====

```

```

=====
                                coef      std err          t      P>|t|      [0.025
0.975]
-----

```

```

Population              0.0001    4.93e-05      2.473      0.013      2.53e-05
0.000
Children                0.9796      0.300        3.266      0.001        0.392
1.568
Income                  9.042e-05  2.11e-05      4.290      0.000      4.91e-05
0.000
Contacts                2.0475      0.626        3.272      0.001        0.821
3.274
Yearly_equip_failure    2.4305      1.024        2.373      0.018        0.423
4.438
Outage_sec_perweek      2.2284      0.168       13.280      0.000        1.899
2.557
Email                   2.0620      0.155       13.278      0.000        1.758
2.366
Item1                   1.8084      0.658        2.747      0.006        0.518
3.099
Item2                   1.6390      0.630        2.603      0.009        0.405
2.873
Item3                   1.4500      0.590        2.457      0.014        0.293
2.607
Item4                   7.3810      0.500       14.750      0.000        6.400
8.362
Item5                   9.9448      0.466       21.349      0.000        9.032
10.858
Item6                   3.5439      0.562        6.302      0.000        2.442
4.646
Item7                   3.6323      0.551        6.597      0.000        2.553
4.712
Item8                   3.8973      0.530        7.352      0.000        2.858
4.936

```

```

=====
Omnibus:                261.682    Durbin-Watson:           1.959
Prob(Omnibus):           0.000    Jarque-Bera (JB):        234.454
Skew:                    0.321    Prob(JB):                1.23e-51
Kurtosis:                2.611    Cond. No.                9.49e+04
=====

```

**Notes:**

- [1]  $R^2$  is computed without centering (uncentered) since the model does not contain a constant.
- [2] Standard Errors assume that the covariance matrix of the errors is correctly specified.
- [3] The condition number is large,  $9.49e+04$ . This might indicate that there are strong multicollinearity or other numerical problems.

## D2. & D3. Model Reduction Method and Justification

- Discuss and describe what statistically based feature selection procedure or a model metric procedure you will use to reduce the initial model in a way that aligns with the research question. (Hint: Discussions regarding model reduction methods can be found in Getting Started with D208 Part I). (\*)
- Explain (justify) why this method is being used to reduce your model.

For my model reduction method, I am going to use Backwards Stepwise Elimination (BSE). In the process of BSE, you iteratively remove the predictors with p-value greater than .05. This method repeats until there are no values greater than .05. We can start by checking which variables have greater than .05 p-value. This includes "Population", "Yearly\_equip\_failure", and "Item3". I will drop these from the table and run the model again.

After the second iteration, there are no longer variables with a p-value of greater than .05. This means we have finished our backwards stepwise elimination.

I chose to implement backwards stepwise elimination because it is an efficient method for helping us to determine the most important predictors of the dependent variable.`

```
In [39]: # set the dependent variable as monthlycharge and the independent variables as
y = dfq['MonthlyCharge']
X = dfq[['Children', 'Income', 'Contacts', 'Outage_sec_perweek', 'Email', 'Item']

# fit the model to x and Y
model = sm.OLS(y, X)
results = model.fit()

#print the results
print(results.summary())
```

## OLS Regression Results

```

=====
Dep. Variable:          MonthlyCharge   R-squared (uncentered):
0.935
Model:                  OLS             Adj. R-squared (uncentered):
0.935
Method:                 Least Squares    F-statistic:
1.197e+04
Date:                   Wed, 08 May 2024  Prob (F-statistic):
0.00
Time:                   16:52:19         Log-Likelihood:
-52337.
No. Observations:      10000            AIC:
1.047e+05
Df Residuals:          9988             BIC:
1.048e+05
Df Model:               12
Covariance Type:       nonrobust
=====

```

```

=====
                                coef      std err          t      P>|t|      [0.025
0.975]
-----

```

```

Children                0.9882        0.300        3.292      0.001      0.400
1.577
Income                  9.139e-05    2.11e-05     4.334      0.000      5e-05
0.000
Contacts                2.0792        0.626        3.321      0.001      0.852
3.306
Outage_sec_perweek      2.2700        0.168       13.544      0.000      1.941
2.599
Email                  2.1066        0.155       13.591      0.000      1.803
2.410
Item1                   2.3235        0.631        3.683      0.000      1.087
3.560
Item2                   1.9474        0.617        3.157      0.002      0.738
3.157
Item4                   7.4463        0.500       14.880      0.000      6.465
8.427
Item5                  10.1755        0.462       22.026      0.000      9.270
11.081
Item6                   3.7071        0.560        6.622      0.000      2.610
4.804
Item7                   3.7903        0.549        6.906      0.000      2.715
4.866
Item8                   4.0240        0.529        7.608      0.000      2.987
5.061

```

```

=====
Omnibus:                259.757    Durbin-Watson:           1.959
Prob(Omnibus):           0.000    Jarque-Bera (JB):       234.122
Skew:                    0.322    Prob(JB):               1.45e-51
Kurtosis:                 2.616    Cond. No.                6.98e+04
=====

```

## Notes:

- [1]  $R^2$  is computed without centering (uncentered) since the model does not contain a constant.
- [2] Standard Errors assume that the covariance matrix of the errors is correct

ly specified.

[3] The condition number is large,  $6.98e+04$ . This might indicate that there are strong multicollinearity or other numerical problems.

## E1. Model Comparison

I will be comparing the adjusted r squared. This is because it is a model evaluation metrics that allows us to take into consideration a larger multitude of variables compared to the default r squared. The initial model's adjusted r-squared is 0.935. The reduced model's adjusted r squared is also 0.935. There is not much of a difference because the initial model's independent variables were already significant enough to explain the dependent, and the independent variables that were removed were not that much greater of a p-value from .05. We can see this in how the prob. f-statistic of the initial model is already 0.00 according to the summary. This means that the overall regression was meaningful from the start, and the fact that the prob. f-statistic is still 0.00 in the reduced model means its continues to remain significantly meaningful

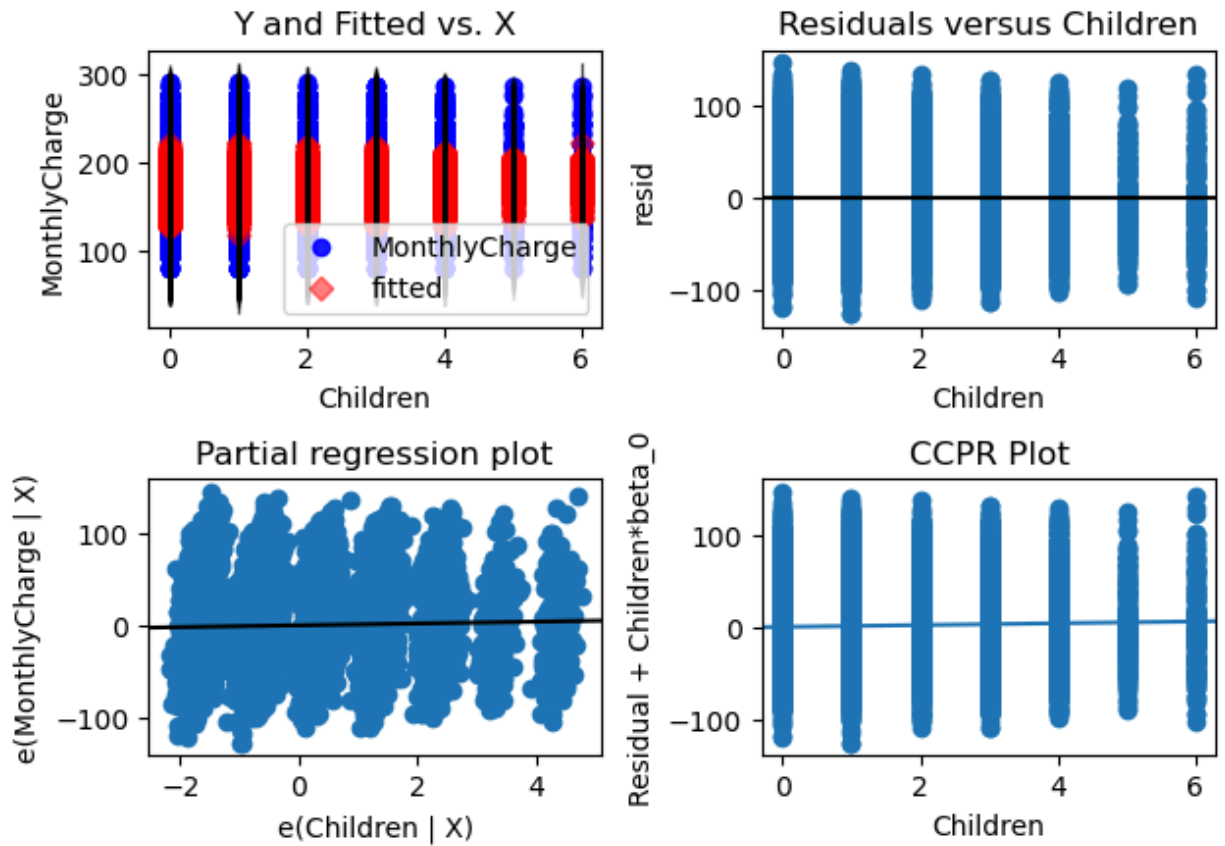
## E2. & E3. Output and Calculations

For the regression plots, I create them for each of the independent variables used in the model. This is done using a for loop running it across all of the variables in the model.

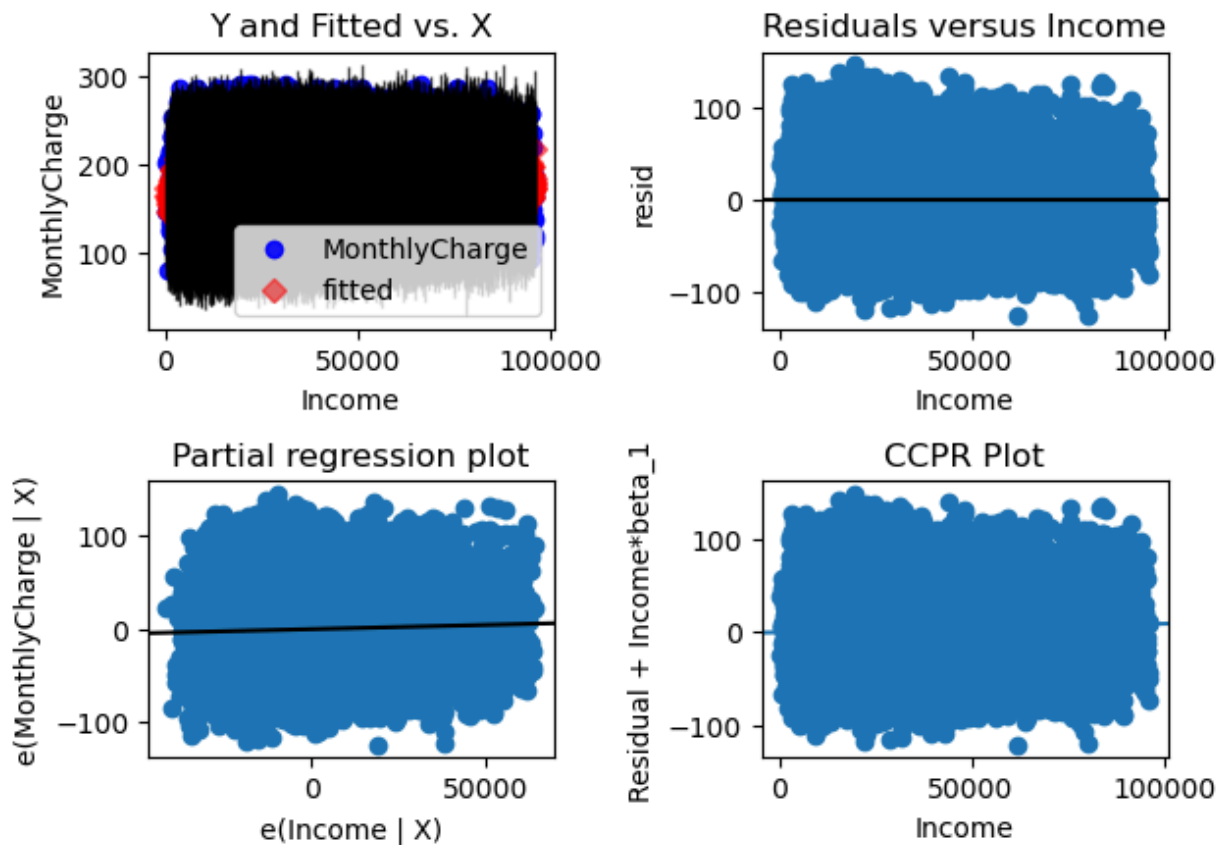
```
In [25]: result_c = ['Children', 'Income', 'Contacts', 'Outage_sec_perweek', 'Email', 'I-
for column in result_c:
    fig = sm.graphics.plot_regress_exog(results, column)
    fig.tight_layout(pad=1.0)
```

```
eval_env: 1
eval_env: 1
eval_env: 1
eval_env: 1
eval_env: 1
eval_env: 1
eval_env: 1
eval_env: 1
eval_env: 1
eval_env: 1
eval_env: 1
eval_env: 1
```

## Regression Plots for Children

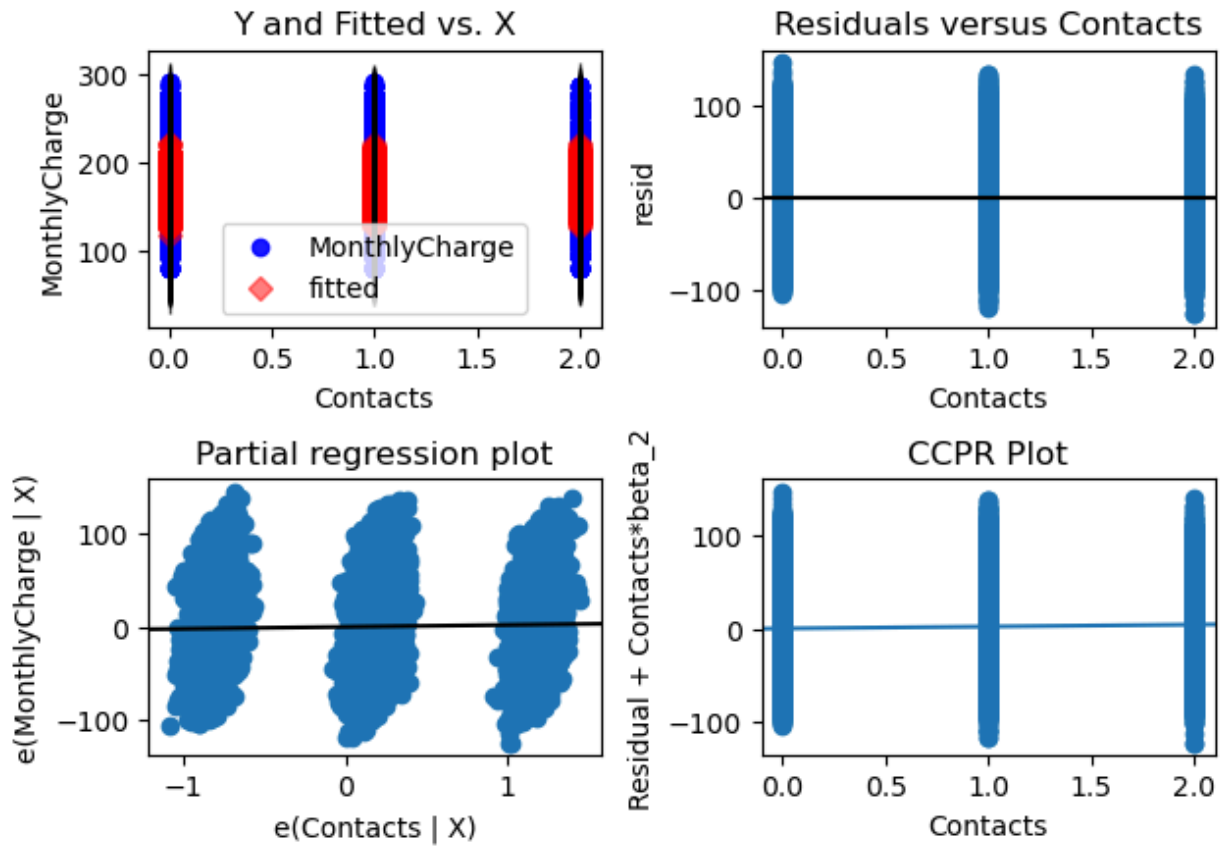


## Regression Plots for Income

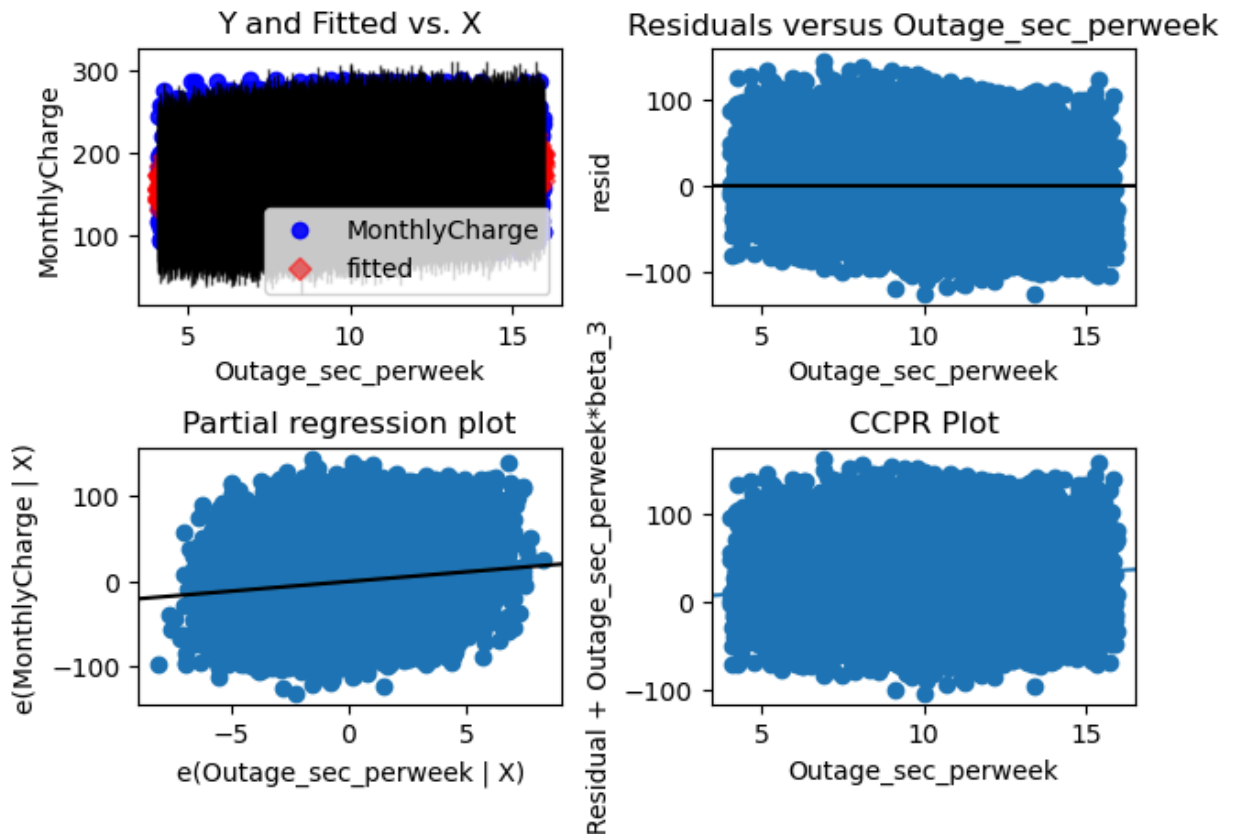




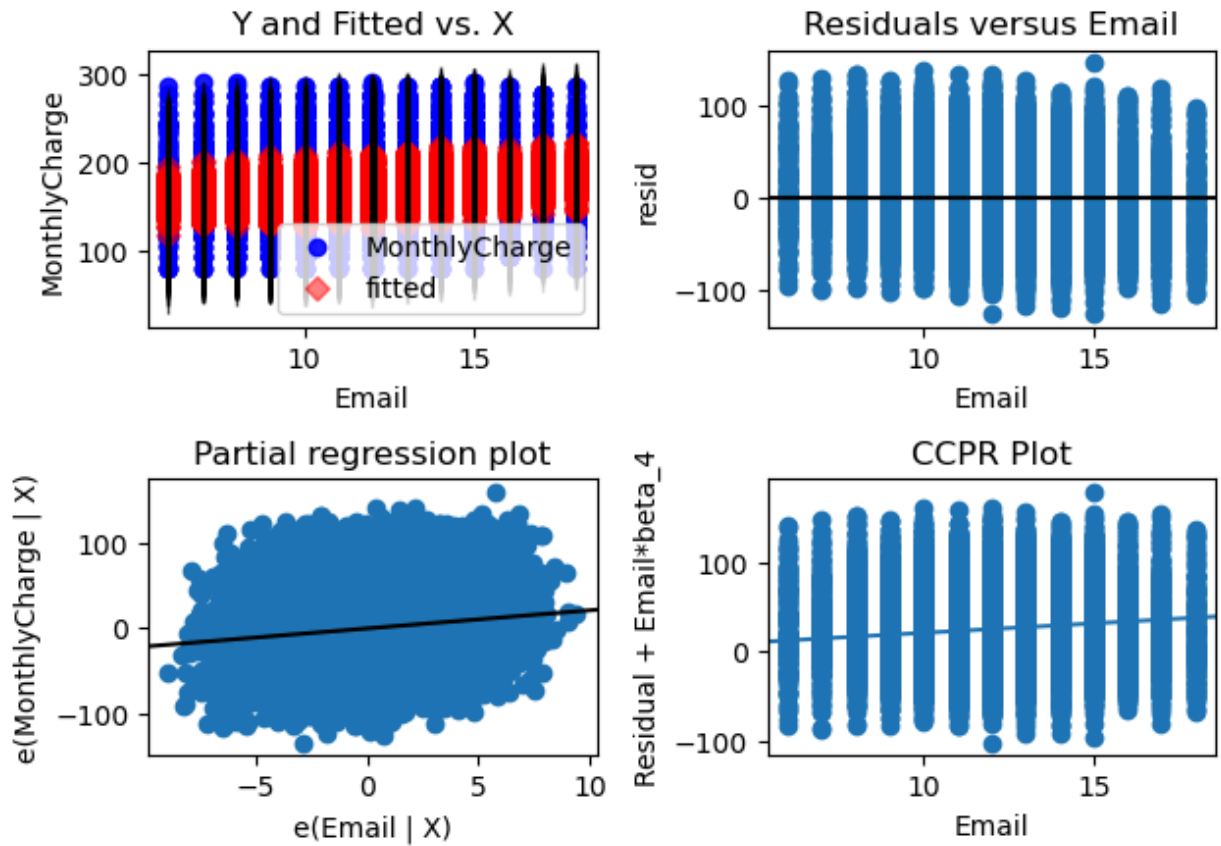
## Regression Plots for Contacts



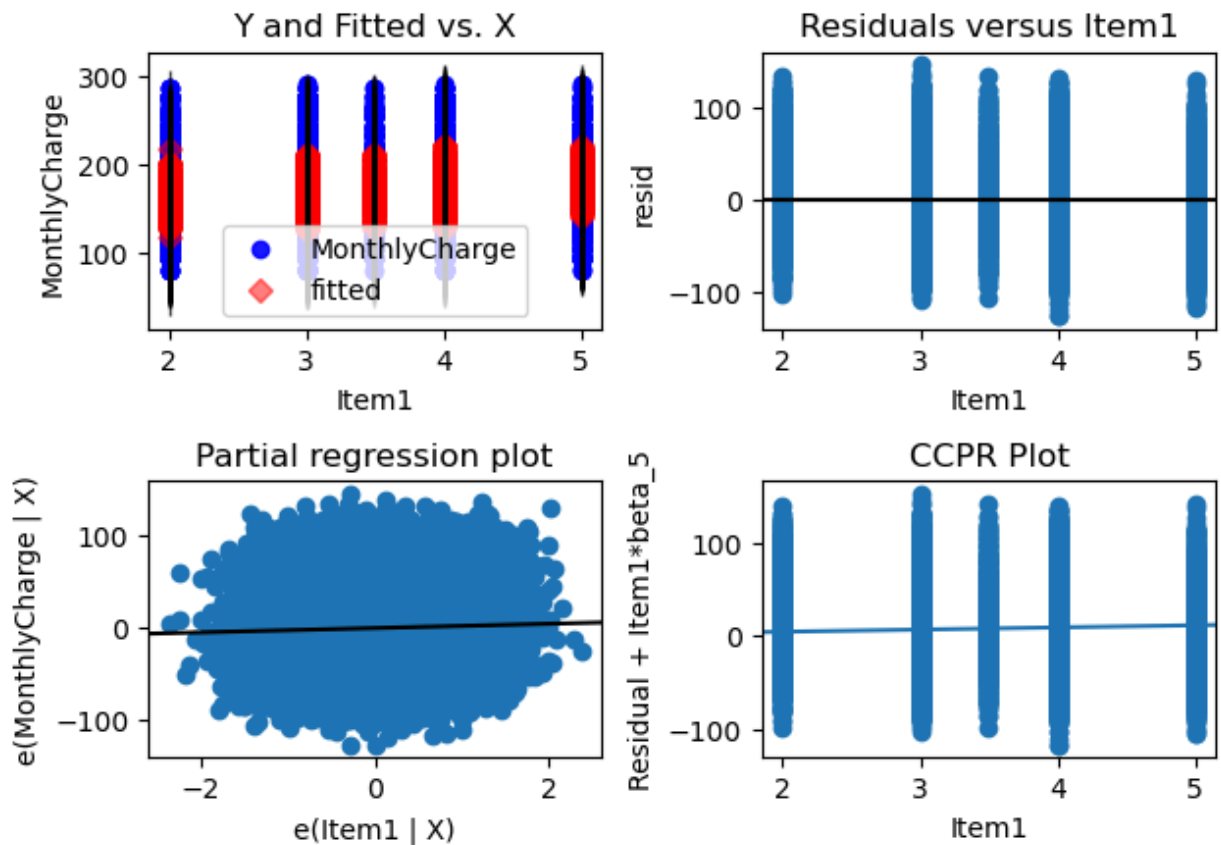
## Regression Plots for Outage\_sec\_perweek



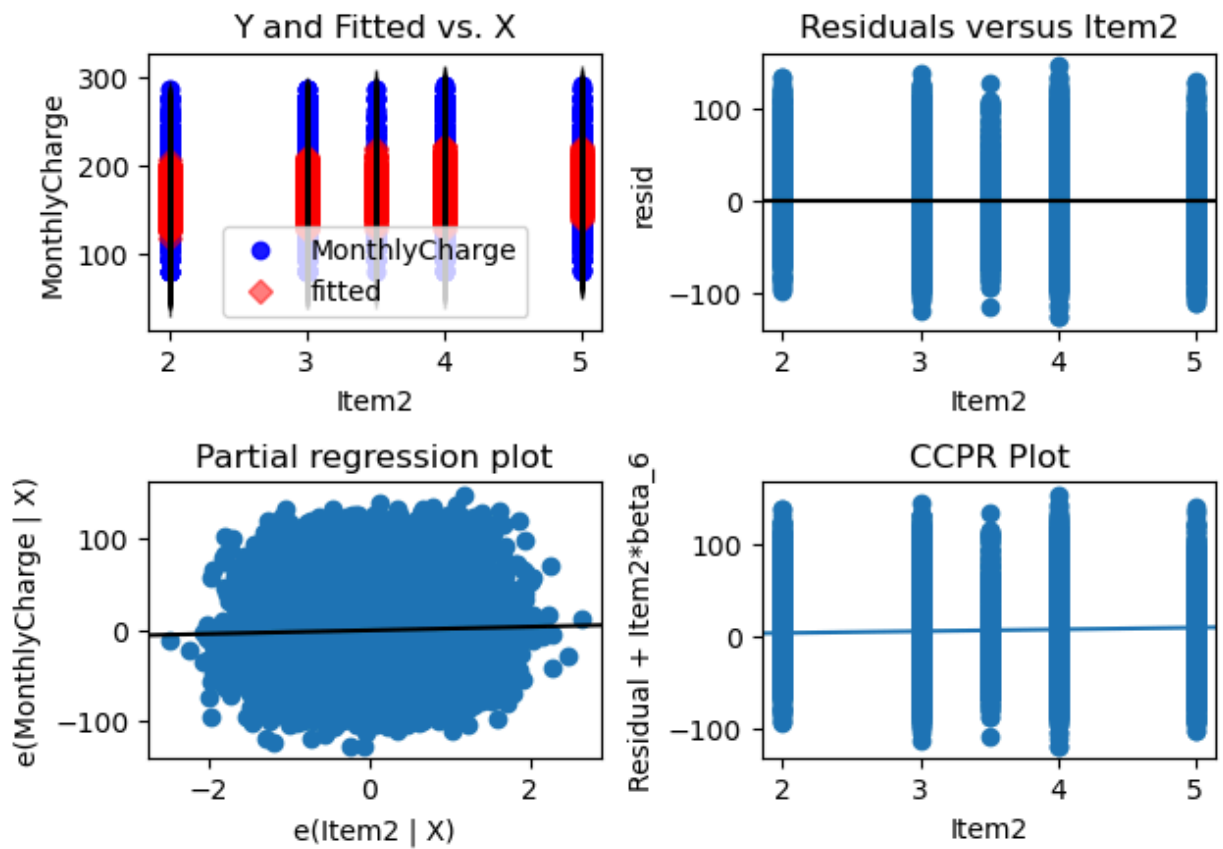
## Regression Plots for Email



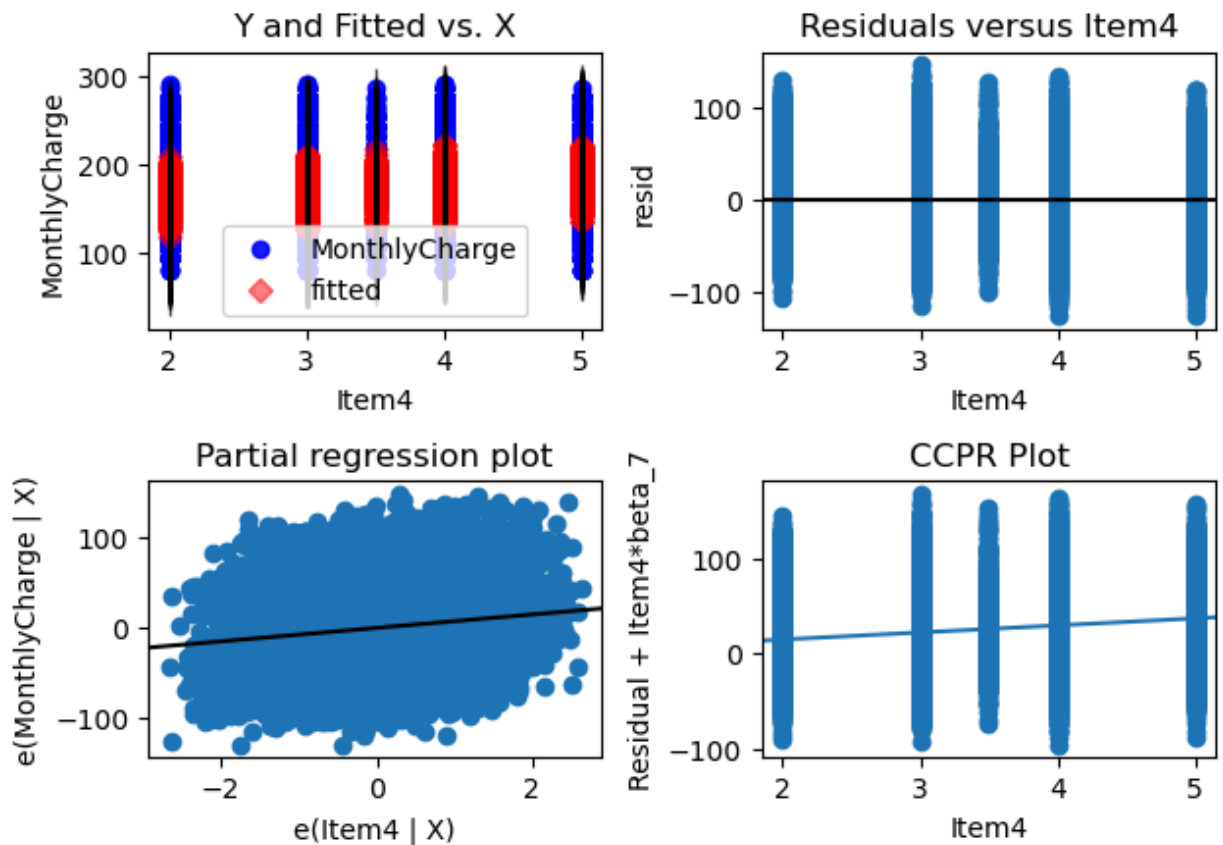
## Regression Plots for Item1



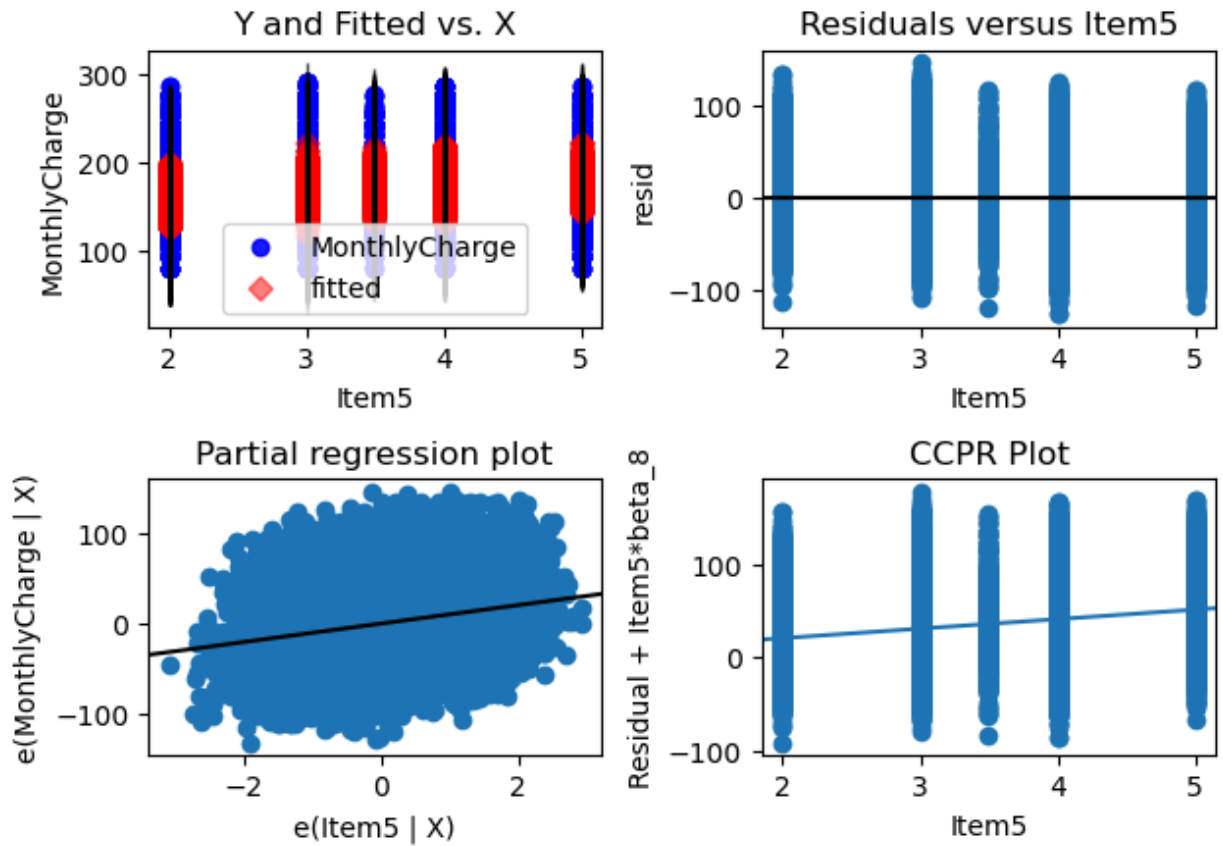
## Regression Plots for Item2



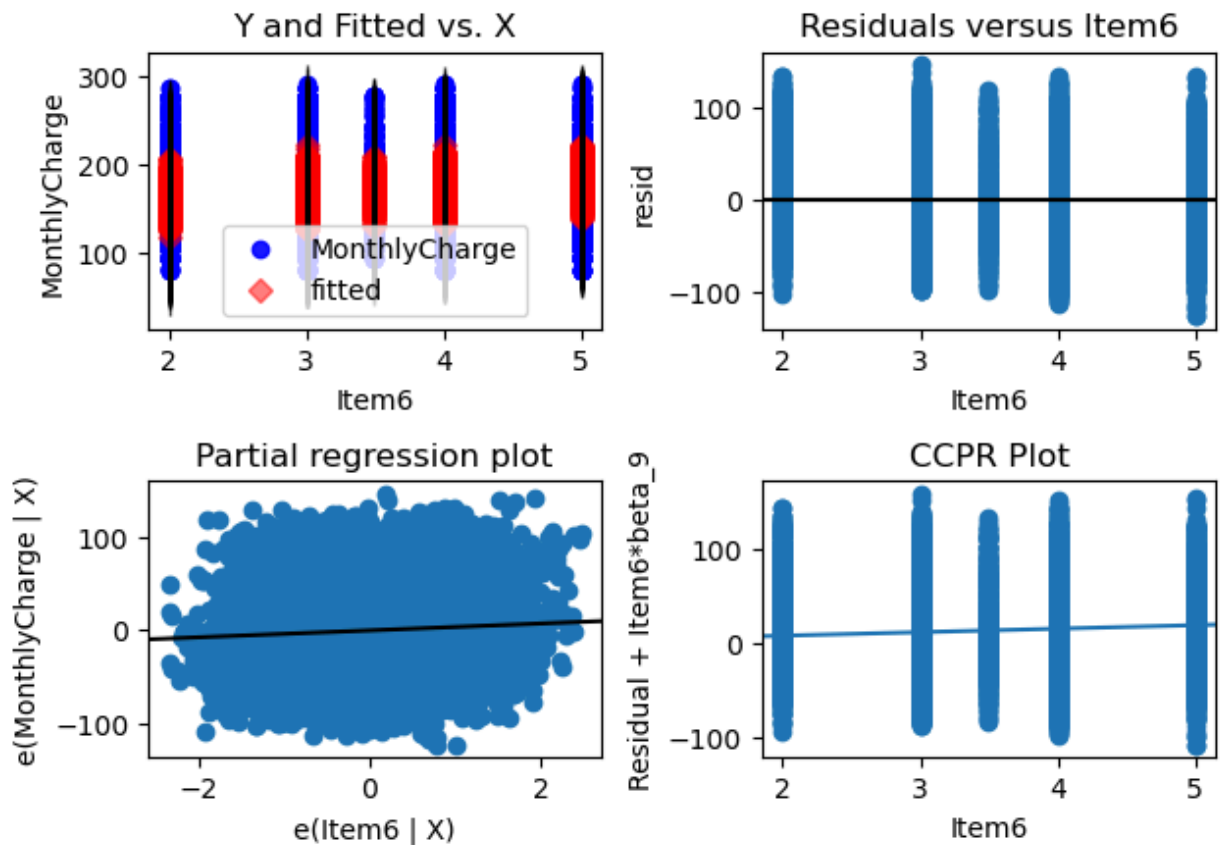
## Regression Plots for Item4



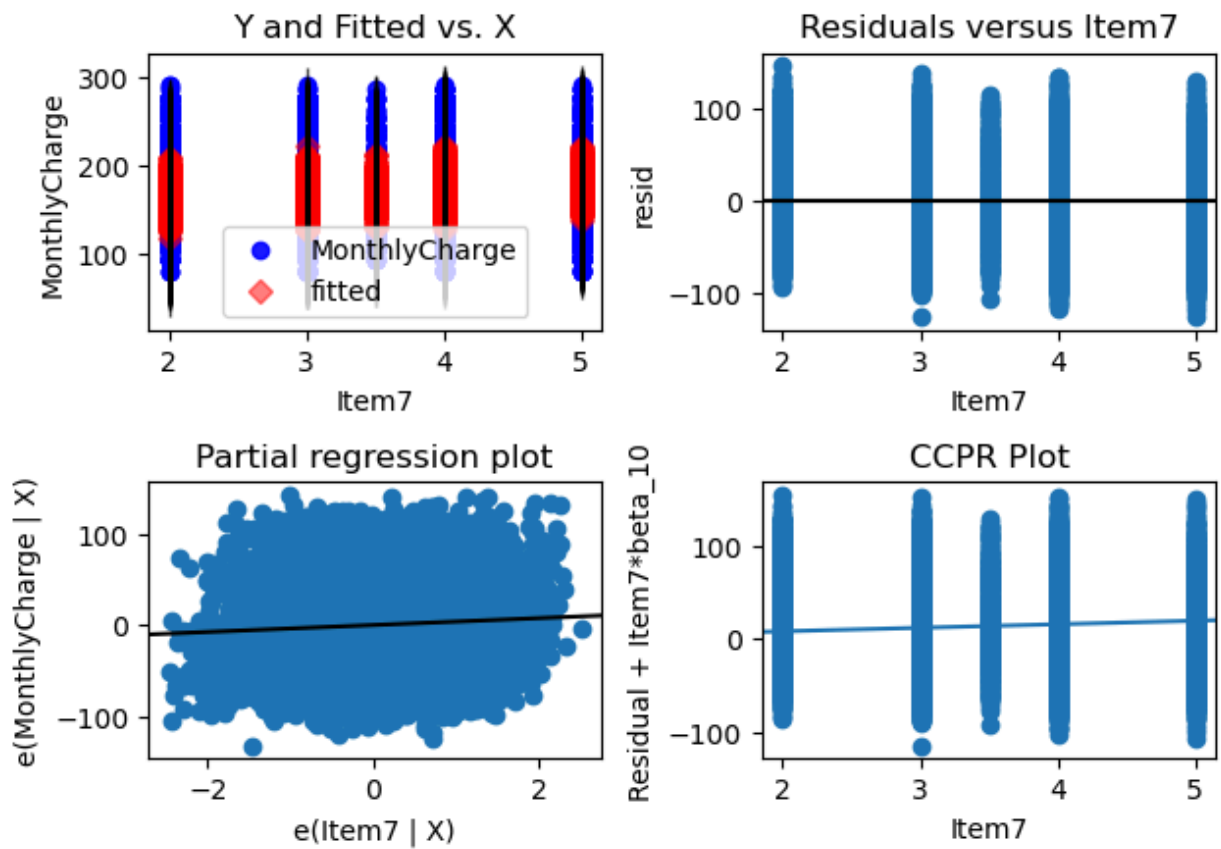
## Regression Plots for Item5



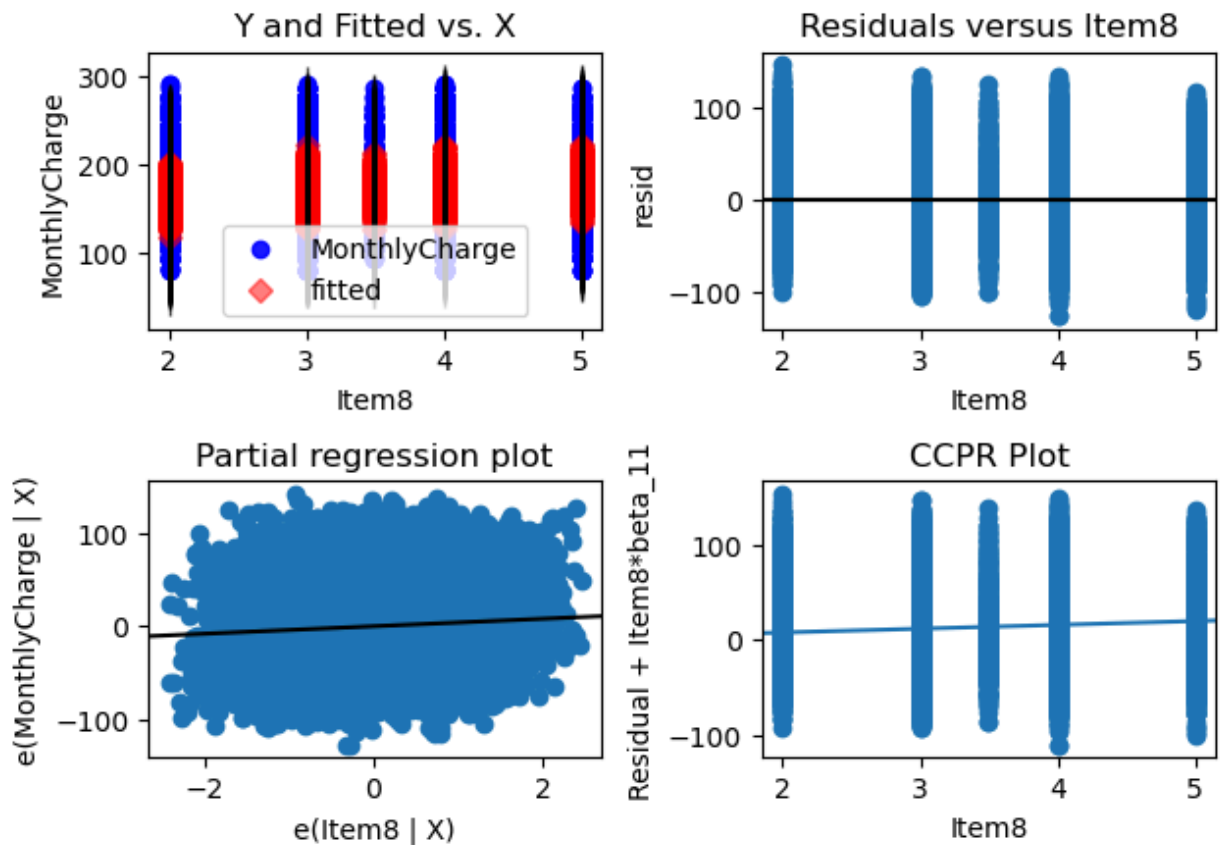
## Regression Plots for Item6



## Regression Plots for Item7



## Regression Plots for Item8



To find the residual standard error, we can use the following function: `results.bse`.

This is useful as the smaller the residual standard error, the better the regression model fits the dataset. We can see that income is the smallest and the best fit in terms of fitting for the dataset.

In [27]: `print(results.bse)`

```
Children          0.300153
Income            0.000021
Contacts          0.625991
Outage_sec_perweek 0.167605
Email             0.155004
Item1             0.630858
Item2             0.616943
Item4             0.500414
Item5             0.461986
Item6             0.559808
Item7             0.548802
Item8             0.528898
dtype: float64
```

## F1. Regression Equation, Coefficients, etc.

### Regression Equation

We can make a regression equation from the summary of the reduced model. This is the dependent variable (y) is equal to (x\_n) times the coefficients added together. This means our regression equation is:

$$\begin{aligned} \text{MonthlyCharge} = & \{x_{\text{children}}\} * 0.9882 + \{x_{\text{Income}}\} * (9.139 * 10^{-5}) + \\ & \{x_{\text{Contacts}}\} * 2.0792 + \{x_{\text{Outage\_sec\_perweek}}\} * 2.2700 + \{x_{\text{Email}}\} * 2.1066 + \\ & \{x_{\text{Item1}}\} * 2.3235 + \{x_{\text{Item2}}\} * 1.9474 + \{x_{\text{Item4}}\} * 7.4463 + \{x_{\text{Item5}}\} * 10.1755 + \\ & \{x_{\text{Item6}}\} * 3.7071 + \{x_{\text{Item7}}\} * 3.7903 + \{x_{\text{Item8}}\} * 4.0240 \end{aligned}$$

### Coefficients

We can see from the regression equation that we have a multitude of different coefficients. For children, it is the MonthlyCharge for a certain customer will be added to based on the number of children they have. Income is reduced by  $9.139 * 10^{-5}$ , and that is added to the MonthlyCharge. For contacts, it is the contacts multiplied by 2.0792. For outage it is seconds multiplied by 2.2. This is interesting because it means because the coefficient is not negative, that the more seconds out per week, the greater the monthly charge. Email is, item1, and item2 are similar in values of their coefficient, multiplying the value of that variable by around 2. Item1 and Item2 are timely response and timely fix rating, which are both time related metrics. Then we can see Item4 and Item5 are large multipliers. Looking at our data dictionary, we can see that these are Reliability and Options. So the more reliable and more options a customer feels they have, the more likely they are to have a large MonthlyCharge. Item6, Item7, and Item8 are respectful response, courteous exchange, and evidence of active listening. These are all conversation related metrics. The item

coefficients are interesting because they can be categorized by how large their impact on the equation is.

## Statistical Significance & Practical Significance of Reduced Model

The model is statistically significant because the adjusted r-squared tells us that it is good at predicting MonthlyChurn because it is close to 1, which is the maximum it can be. The number of variables has also been considered since the adjusted r-squared adjusts the score to take into consideration the number of independent variables.

In terms of being practically significant, it can help us identify what are the important items that affect MonthlyCharge. For instance, Item4, Item5, and Item8 are the items out of the survey that affect the amount of MonthlyCharge associated with a customer. That would be reliability, options, and evidence of active listening. It indicates that there should be an effort to understand why increasing customer perception of the variables leads to an increased MonthlyCharge. It also tell us that variables like how many children are not as significant in effecting the MonthlyCharge as any of the items, because it is less than 1 for a coefficient. This can lead us to explore questions like "Does customer perception matter more than customer characteristics like the number of children they have?".

## Disadvantages

There are some disadvantages we should consider when using the linear regression model. One is that we have to be careful of multicollinearity. This can occur if any of our independent variables are to closely aligned in terms of their significance with one another. This poses a significant risk for our model since we have a large number of variables. Its also a problem that the basis of a linear regression model assumes that these relationships are linear. However, many of these relationship may not be, and we would not be accurately able to represent the nuances of their relationships with one another using only a straight line on a graph like linear regression attempts to use.

## F2. Recommendations

The recommendations based off of the model that was created are to focus on understanding why the highest coefficient items play such a significant affect on MonthlyCharge. This can help us understand why some customers may be paying more than others, and perhaps allow us to target those high value customers. For instance, questions like "Does customer perception matter more than customer characteristics like the number of children they have?" allows us to focus our attention more on customer perception of items rather than targeting customers with lots of children. They may both be positively correlated, but focusing allows us to optimize for a greater revenue for the amount of effort put in.

## H. Third Party Sources of Code

No third party code used

## I. Sources

Adjusted R-squared. Corporate Finance Institute. (2023, November 21).

<https://corporatefinanceinstitute.com/resources/data-science/adjusted-r-squared/>

Grover, J. (2021, March 28). Short python code for Backward Elimination with detailed explanation. Medium. <https://groverjatin.medium.com/short-python-code-for-backward-elimination-with-detailed-explanation-52894a9a7880>

Regression plots¶. Regression Plots - statsmodels 0.15.0 (+272). (n.d.).

[https://www.statsmodels.org/devel/examples/notebooks/generated/regression\\_plots.html](https://www.statsmodels.org/devel/examples/notebooks/generated/regression_plots.html)