

Part I: Research Question

A. Question

Can we identify the variables most influential in the readmittance of patients following a visit at our hospital?

B. Variables

CaseOrder: Integer that maintains initial order of raw data Customer_id: String that provides unique key to patient Interaction: String that provides unique key for patient transactions, procedures, and admissions UID: String that provides unique key for patient transactions, procedures, and admissions City: String that indicates patient city of residence on billing statement State: String that indicates patient state of residence on billing statement County: String that indicates patient county of residence on billing statement Zip: Integer that indicates patient zip code of residence on billing statement Lat: Float that indicates Latitude based on billing statement address Lng: Float that indicates Longitude based on billing statement address Population: Integer that indicates population based on census data within one-mile radius of address Area: String that indicates area type based on census Timezone: String that indicates timezone based on patient provided residence Job: String that indicates patient job as provided by admissions information Children: Float that indicates number of children provided by admissions information Age: Float that indicates age of patient provided by admissions information Education: String that indicates highest earned level of education provided by admissions information Employment: String that indicates current employment status provided by admissions information Income: Float that indicates annual income of patient (or primary insurance holder) provided by admissions information Marital: String that indicates current marital status of patient (or primary insurance holder) provided by admissions information Gender: String that indicates self-identified gender ReAdmis: String that indicates whether the patient was readmitted within a month of relevant visit VitD_levels: Float that indicates patient vitamin D levels measured in ng/ml Doc_visits: Integer that indicates the number of times that the PCP (Primary Care Physician) visited the patient during their initial stay Full_meals_eaten: Integer that indicates the number of full meals the patient ate during hospitalization (partial = 0) VitD_supp: Integer that indicates the number of times vitamin D supplements were administered to patient Soft_drink: String that indicates if a patient drinks 3 or more sodas in a day frequently Initial_admin: String that indicates the route in which a patient was admitted into the hospital HighBlood: String that indicates if the patient has high blood pressure Stroke: String that indicates if the patient has had a stroke Complication_risk: String that indicates the level of risk for complication associated with the patient Overweight: Float that indicates if the patient is overweight Arthritis: String that indicates if the patient has arthritis Diabetes: String that indicates if the patient has diabetes Hyperlipidemia: String that indicates if the patient has hyperlipidemia BackPain: String that indicates if the patient has chronic back pain Anxiety: Float that indicates if the patient has an anxiety disorder Allergic_rhinitis: String that indicates if the patient has allergic rhinitis Reflux_esophagitis: String that indicates if the patient has reflux esophagitis Asthma: String that indicates if the patient if the patient has asthma Services: String that indicates the primary service a

patient received during their stay Initial_days: Float that indicates the number of days the patient stayed during the initial visit TotalCharge: Float that indicates the average cost per day (Total Cost/# of days) Additional_charges: Float that indicates the average cost of miscellaneous services received during stay Item1: Integer that indicates survey answer about the importance of "Timely Admission" (Scale of 1 most - 8 least) Item2: Integer that indicates survey answer about the importance of "Timely Treatment" (Scale of 1 most - 8 least) Item3: Integer that indicates survey answer about the importance of "Timely Visits" (Scale of 1 most - 8 least) Item4: Integer that indicates survey answer about the importance of "Reliability" (Scale of 1 most - 8 least) Item5: Integer that indicates survey answer about the importance of "Options" (Scale of 1 most - 8 least) Item6: Integer that indicates survey answer about the importance of "Hours of Treatment" (Scale of 1 most - 8 least) Item7: Integer that indicates survey answer about the importance of "Courteous Staff" (Scale of 1 most - 8 least) Item8: Integer that indicates survey answer about the importance of "Evidence of Active Listening from Doctor" (Scale of 1 most - 8 least)

Part II: Data-Cleaning Plan

C. Plan Explanation

C1.

- Orient Myself to the data
- Value Counts
- Missing Values
- Standardized numeric columns
- Categorical variables to Numeric
- Histograms and Boxplots (Outlier Detection)
- PCA

C2.

This data cleaning plan is based up on the Data Preparation Phase explained in the text "Data Science Using Python and R." I have added a few pieces to be more specific about my process that is based on my experience as a full-time Data Analyst.

The first time is orienting myself to the data. In general - this is just the step where I am going to look at the raw data itself. What is my initial input and what stands out? What are the data types I am working with? etc.

The second step is value counts. In general - this will help me to identify unique identifier columns that may not be beneficial in a model. This will also give me a high level view of distributions for both categorical and numeric values.

Missing values is the step in which I will utilize multiple imputation to address missing values. ML models typically do not account for missing data (there is nuance in this), so you need to assign values in every row of every column. I will use IterativeImputer from sklearn to fill in our missing measures. IterativeImputer is similar to MICE within R. It is going to iterate through the dataset multiple times and establish estimations of the missing value based upon the other variables. (Scikit-learn)

Standardized numeric columns will be two fold - there will be some mapping and then standardizing the columns with sklearn's StandardScaler. This module will center distributions around 0 as the mean. This will be necessary for PCA and will be easier to evaluate outliers and distribution since the scale will be standardized. (Larose, 2019)

Converting categorical variables to numeric is necessary on a couple levels. We need this for our IterativeImputer so that categorical variables can be evaluated as a part of that model. Also - again this will be necessary in future ML applications. (Larose, 2019)

Histograms and boxplots will allow me to visualize the distributions and outliers in a graphic format. This will help to inform decisions on how to deal with outliers. Do certain observations need to be excluded?

PCA is really the end goal of this project. PCA is a form of dimensionality reduction. "Principal components analysis (PCA) seeks to account for the correlation structure of a set of predictor variables, using a smaller set of uncorrelated linear combinations of these variables, called components," (Larose, 2019). I find that the easiest way to conceptualize this is as the vector transformation along a single line to explain the maximal variance. With this method - we can reduce our necessity for all variables and explain most of the variance in a dataset.

C3.

I am utilizing Python as my programming language of choice in this project. Python is a highly flexible programming language that handles data really well. It is consistently one of the most used programming languages in the world (Eastwood, 2020). In thinking of long-term strategy - it is much simpler to put a Python built model into production within an existing application environment. This allows the data analyst to not just be putting together presentations and visuals, but to truly contribute in production. R is great at statistical programming - it is very often used in research, because it is built for that very implementation. So, there is certainly in argument for its utilization in any project like this. But, I prefer Python for its flexibility. You can create web applications, desktop applications, api's and the list goes on. So, if we ever needed to pivot this analysis into a different domain - it would be far easier to do in Python than R.

As for packages - I am using several. Pandas is one of the most popular packages in Python and it deals with tabular data. It effectively allows you to interact with the data in a format we are traditionally used to - with column headers and row numbers. Numpy is a package that I will use for some calculations and also interacting with the data at certain points as an array. Matplotlib is the package that will be used for visualizations. Sklearn is a scientific package for Python meant to implement the various stages of Machine Learning. So, sklearn will be used in imputation,

transformation, and ultimately PCA. Bioinfokit is an additional package that I will be using for some visualization of PCA.

```
In [397...
#Package imports
import matplotlib.pyplot as plt
import pandas as pd
import numpy as np
from matplotlib import pyplot
from sklearn import preprocessing
from sklearn.experimental import enable_iterative_imputer
from sklearn.impute import IterativeImputer
from sklearn.preprocessing import StandardScaler
from sklearn.decomposition import PCA
```

```
In [398...
#Read in dataset and set index to first column that already has indexes. It also 0 index.
Med_df = pd.read_csv(r'C:\Users\jacob.colp.UNITY\Downloads\Medical Data Raw\medical_raw.csv')
Med_df = Med_df.reset_index(drop = True)
```

```
In [399...
#Show relevant statistical information about numeric columns.

Med_df.describe()
```

```
Out[399...
      CaseOrder      Zip      Lat      Lng      Population      Children      Age
count  10000.000000  10000.000000  10000.000000  10000.000000  10000.000000  7412.000000  7586.000000
mean     5000.500000  50159.323900    38.751099   -91.243080    9965.253800    2.098219   53.295676
std     2886.89568   27469.588208     5.403085    15.205998   14824.758614    2.155427   20.659182
min        1.000000    610.000000    17.967190   -174.209690     0.000000    0.000000   18.000000
25%     2500.750000  27592.000000    35.255120   -97.352982    694.750000    0.000000   35.000000
50%     5000.500000  50207.000000    39.419355   -88.397230   2769.000000    1.000000   53.000000
75%     7500.250000  72411.750000    42.044175   -80.438050  13945.000000    3.000000   71.000000
max    10000.000000  99929.000000    70.560990   -65.290170  122814.000000   10.000000   89.000000
```

8 rows × 25 columns

```
In [400...
#Convert yes/no columns to binary. I also converted Complication Risk to ordinal values

for x in Med_df.columns:
    if Med_df[x].dtype == object:
        Med_df[x].replace({'Yes':1, 'yes':1, 'No':0, 'no':0}, inplace = True)

Med_df.Complication_risk.replace({'High':3, 'Medium':2, 'Low':1}, inplace = True)

Med_df.info()
```

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 10000 entries, 0 to 9999
Data columns (total 52 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                             7412 non-null   float64
15  Age                                   7586 non-null   float64
16  Education                            10000 non-null  object
17  Employment                           10000 non-null  object
18  Income                               7536 non-null   float64
19  Marital                              10000 non-null  object
20  Gender                               10000 non-null  object
21  ReAdmis                              10000 non-null  int64
22  VitD_levels                          10000 non-null  float64
23  Doc_visits                           10000 non-null  int64
24  Full_meals_eaten                     10000 non-null  int64
25  VitD_supp                            10000 non-null  int64
26  Soft_drink                           7533 non-null   float64
27  Initial_admin                        10000 non-null  object
28  HighBlood                            10000 non-null  int64
29  Stroke                               10000 non-null  int64
30  Complication_risk                    10000 non-null  int64
31  Overweight                           9018 non-null   float64
32  Arthritis                            10000 non-null  int64
33  Diabetes                             10000 non-null  int64
34  Hyperlipidemia                       10000 non-null  int64
35  BackPain                             10000 non-null  int64
36  Anxiety                              9016 non-null   float64
37  Allergic_rhinitis                    10000 non-null  int64
38  Reflux_esophagitis                   10000 non-null  int64
39  Asthma                               10000 non-null  int64
40  Services                             10000 non-null  object
41  Initial_days                          8944 non-null   float64
42  TotalCharge                           10000 non-null  float64
43  Additional_charges                   10000 non-null  float64
44  Item1                                10000 non-null  int64
45  Item2                                10000 non-null  int64
46  Item3                                10000 non-null  int64
47  Item4                                10000 non-null  int64
48  Item5                                10000 non-null  int64
49  Item6                                10000 non-null  int64
50  Item7                                10000 non-null  int64
51  Item8                                10000 non-null  int64
dtypes: float64(12), int64(25), object(15)
memory usage: 4.0+ MB

```

```
#This will tell me the number of null values in each column
```

```
Med_df.isnull().sum()
```

```
Out[401...] CaseOrder      0
Customer_id    0
Interaction     0
UID             0
City            0
State          0
County         0
Zip            0
Lat            0
Lng            0
Population     0
Area           0
Timezone       0
Job            0
Children       2588
Age            2414
Education      0
Employment     0
Income         2464
Marital        0
Gender         0
ReAdmis        0
VitD_levels    0
Doc_visits     0
Full_meals_eaten 0
VitD_supp      0
Soft_drink     2467
Initial_admin  0
HighBlood      0
Stroke         0
Complication_risk 0
Overweight     982
Arthritis      0
Diabetes       0
Hyperlipidemia 0
BackPain       0
Anxiety        984
Allergic_rhinitis 0
Reflux_esophagitis 0
Asthma         0
Services       0
Initial_days   1056
TotalCharge    0
Additional_charges 0
Item1          0
Item2          0
Item3          0
Item4          0
Item5          0
Item6          0
Item7          0
Item8          0
dtype: int64
```

```
In [402...] #Find unique columns
for x in Med_df:
```

```
print(x+'': '+str(Med_df[x].is_unique))
```

```
CaseOrder: True
Customer_id: True
Interaction: True
UID: True
City: False
State: False
County: False
Zip: False
Lat: False
Lng: False
Population: False
Area: False
Timezone: False
Job: False
Children: False
Age: False
Education: False
Employment: False
Income: False
Marital: False
Gender: False
ReAdmis: False
VitD_levels: True
Doc_visits: False
Full_meals_eaten: False
VitD_supp: False
Soft_drink: False
Initial_admin: False
HighBlood: False
Stroke: False
Complication_risk: False
Overweight: False
Arthritis: False
Diabetes: False
Hyperlipidemia: False
BackPain: False
Anxiety: False
Allergic_rhinitis: False
Reflux_esophagitis: False
Asthma: False
Services: False
Initial_days: False
TotalCharge: True
Additional_charges: False
Item1: False
Item2: False
Item3: False
Item4: False
Item5: False
Item6: False
Item7: False
Item8: False
```

In [403...

```
#I am dropping Job, Income, and Marital Status because those all may not relate to the
Med_df.drop(['CaseOrder', 'Customer_id', 'Interaction', 'UID', 'Job', 'Income', 'Marita
#Rename all Item survey columns to their question topic.
```

```
Med_df.rename({'Item1':'Survey_Timely_Admission', 'Item2':'Survey_Timely_Treatment', 'I
              'Item7':'Survey_Courteous_Staff', 'Item8':'Survey_Evidence_of_Active_Listeni

Med_df.columns
```

Out[403...

```
Index(['City', 'State', 'County', 'Zip', 'Lat', 'Lng', 'Population', 'Area',
      'Timezone', 'Children', 'Age', 'Education', 'Employment', 'Gender',
      'ReAdmis', 'VitD_levels', 'Doc_visits', 'Full_meals_eaten', 'VitD_supp',
      'Soft_drink', 'Initial_admin', 'HighBlood', 'Stroke',
      'Complication_risk', 'Overweight', 'Arthritis', 'Diabetes',
      'Hyperlipidemia', 'BackPain', 'Anxiety', 'Allergic_rhinitis',
      'Reflux_esophagitis', 'Asthma', 'Services', 'Initial_days',
      'TotalCharge', 'Additional_charges', 'Survey_Timely_Admission',
      'Survey_Timely_Treatment', 'Survey_Timely_Visits', 'Survey_Reliability',
      'Survey_Options', 'Survey_Hours_of_Treatment', 'Survey_Courteous_Staff',
      'Survey_Evidence_of_Active_Listening'],
      dtype='object')
```

In [404...

```
#I am converting categorical values into the "category" datatype from pandas. This will

categorical_variables = ['City', 'State', 'County', 'Zip', 'Area', 'Timezone', 'Educati

for x in categorical_variables:
    Med_df[x] = Med_df[x].astype('category')

Med_df.dtypes
```

Out[404...

City	category
State	category
County	category
Zip	category
Lat	float64
Lng	float64
Population	int64
Area	category
Timezone	category
Children	float64
Age	float64
Education	category
Employment	category
Gender	category
ReAdmis	int64
VitD_levels	float64
Doc_visits	int64
Full_meals_eaten	int64
VitD_supp	int64
Soft_drink	float64
Initial_admin	category
HighBlood	int64
Stroke	int64
Complication_risk	int64
Overweight	float64
Arthritis	int64
Diabetes	int64
Hyperlipidemia	int64
BackPain	int64
Anxiety	float64
Allergic_rhinitis	int64
Reflux_esophagitis	int64

Asthma	int64
Services	category
Initial_days	float64
TotalCharge	float64
Additional_charges	float64
Survey_Timely_Admission	int64
Survey_Timely_Treatment	int64
Survey_Timely_Visits	int64
Survey_Reliability	int64
Survey_Options	int64
Survey_Hours_of_Treatment	int64
Survey_Courteous_Staff	int64
Survey_Evidence_of_Active_Listening	int64
dtype:	object

In [405...

#This code gives me the value counts for each value of each column. I can roughly see d

```
for x in Med_df:
    print(Med_df[x].value_counts())
```

```
Houston      36
San Antonio  26
Springfield  22
Miami        21
New York     21
..
Hollenberg   1
Hollandale   1
Holland Patent  1
Holcombe     1
Zumbro Falls 1
Name: City, Length: 6072, dtype: int64
TX      553
CA      550
PA      547
NY      514
IL      442
OH      383
MO      328
FL      304
VA      287
IA      276
MI      273
MN      267
NC      254
GA      247
KS      220
WI      214
KY      210
WV      207
OK      207
IN      195
TN      194
AL      194
WA      191
AR      190
NE      185
CO      179
NJ      176
LA      173
```

MA	149
MS	134
MD	131
SC	128
SD	123
OR	122
ME	122
MT	112
NM	110
ID	109
ND	108
AZ	108
CT	80
NH	79
UT	72
AK	70
VT	60
NV	51
WY	51
PR	43
HI	34
DE	17
RI	14
DC	13

Name: State, dtype: int64

Jefferson	118
Washington	100
Franklin	93
Los Angeles	88
Montgomery	80

...

Churchill	1
Republic	1
Cimarron	1
St. Martin	1
Lamoille	1

Name: County, Length: 1607, dtype: int64

24136	4
88345	4
77663	4
38330	4
37324	4

..

37146	1
37144	1
37138	1
37134	1
99929	1

Name: Zip, Length: 8612, dtype: int64

36.06702	4
33.34798	4
35.25512	4
39.38610	4
37.86890	4

..

41.00911	1
39.20560	1
46.36035	1
34.96563	1
40.49998	1

Name: Lat, Length: 8588, dtype: int64

```

-121.28753    4
-82.35159    4
-85.99134    4
-105.68001   4
-89.03658    4
..
-74.87894    1
-99.17911    1
-91.81854    1
-106.83727   1
-80.19959    1
Name: Lng, Length: 8601, dtype: int64
0          109
195         14
115         11
178         11
285         11
...
8092         1
11147        1
27175        1
7371         1
41524        1
Name: Population, Length: 5951, dtype: int64
Rural        3369
Suburban     3328
Urban        3303
Name: Area, dtype: int64
America/New_York      3889
America/Chicago      3771
America/Los_Angeles   937
America/Denver        612
America/Detroit       262
America/Indiana/Indianapolis  151
America/Phoenix       100
America/Boise         86
America/Anchorage     50
America/Puerto_Rico   43
Pacific/Honolulu      34
America/Menominee     14
America/Nome          12
America/Indiana/Vincennes  8
America/Sitka         6
America/Kentucky/Louisville  6
America/Toronto       5
America/Indiana/Tell_City  3
America/Indiana/Marengo  3
America/North_Dakota/Beulah  2
America/Indiana/Winamac  1
America/Indiana/Vevay   1
America/North_Dakota/New_Salem  1
America/Indiana/Knox    1
America/Yakutat        1
America/Adak           1
Name: Timezone, dtype: int64
0.0      1880
1.0      1858
3.0      1113
2.0      1094
4.0       739
8.0       157

```

```

7.0      154
6.0      145
5.0      126
9.0       83
10.0     63
Name: Children, dtype: int64
30.0     126
47.0     124
74.0     123
38.0     123
40.0     122
...
75.0     90
82.0     90
63.0     90
51.0     89
36.0     85
Name: Age, Length: 72, dtype: int64
Regular High School Diploma      2444
Bachelor's Degree                 1724
Some College, 1 or More Years, No Degree  1484
9th Grade to 12th Grade, No Diploma      832
Associate's Degree                797
Master's Degree                   701
Some College, Less than 1 Year        642
Nursery School to 8th Grade          552
GED or Alternative Credential         389
Professional School Degree           208
No Schooling Completed              133
Doctorate Degree                    94
Name: Education, dtype: int64
Full Time      6029
Student        1017
Part Time      991
Unemployed     983
Retired        980
Name: Employment, dtype: int64
Female          5018
Male            4768
Prefer not to answer    214
Name: Gender, dtype: int64
0      6331
1      3669
Name: ReAdmis, dtype: int64
17.802330      1
18.423248      1
15.954743      1
19.566698      1
19.221626      1
..
18.107325      1
17.331743      1
49.013013      1
18.292722      1
20.421883      1
Name: VitD_levels, Length: 10000, dtype: int64
5      3823
6      2436
4      2385
7       634
3       595

```

```

8      61
2      58
1      6
9      2
Name: Doc_visits, dtype: int64
0      3715
1      3615
2      1856
3      612
4      169
5      25
6      6
7      2
Name: Full_meals_eaten, dtype: int64
0      6702
1      2684
2      544
3      64
4      5
5      1
Name: VitD_supp, dtype: int64
0.0      5589
1.0      1944
Name: Soft_drink, dtype: int64
Emergency Admission      5060
Elective Admission      2504
Observation Admission    2436
Name: Initial_admin, dtype: int64
0      5910
1      4090
Name: HighBlood, dtype: int64
0      8007
1      1993
Name: Stroke, dtype: int64
2      4517
3      3358
1      2125
Name: Complication_risk, dtype: int64
1.0      6395
0.0      2623
Name: Overweight, dtype: int64
0      6426
1      3574
Name: Arthritis, dtype: int64
0      7262
1      2738
Name: Diabetes, dtype: int64
0      6628
1      3372
Name: Hyperlipidemia, dtype: int64
0      5886
1      4114
Name: BackPain, dtype: int64
0.0      6110
1.0      2906
Name: Anxiety, dtype: int64
0      6059
1      3941
Name: Allergic_rhinitis, dtype: int64
0      5865
1      4135

```

```

Name: Reflux_esophagitis, dtype: int64
0    7107
1    2893
Name: Asthma, dtype: int64
Blood Work    5265
Intravenous   3130
CT Scan       1225
MRI           380
Name: Services, dtype: int64
10.585770     1
64.630142     1
48.772686     1
67.036508     1
63.334689     1
..
9.216747      1
1.021594      1
10.261690     1
17.170461     1
70.850592     1
Name: Initial_days, Length: 8944, dtype: int64
3191.048774    1
7329.393066    1
8498.290160    1
8451.833926    1
7530.770634    1
..
2065.518265    1
3409.593273    1
15289.590000    1
4383.419018    1
8700.856021    1
Name: TotalCharge, Length: 10000, dtype: int64
8013.787149     5
22000.064780    4
3241.339760     4
11303.682330    4
8755.123303     4
..
20461.526600    1
13357.949060    1
5316.329223     1
24412.109160    1
11643.189930    1
Name: Additional_charges, Length: 8888, dtype: int64
4    3455
3    3404
5    1377
2    1315
6     225
1     213
7      10
8       1
Name: Survey_Timely_Admission, dtype: int64
3    3439
4    3351
5    1421
2    1360
1     213
6     204
7      12

```

Name: Survey_Timely_Treatment, dtype: int64

```
4    3464
3    3379
5    1358
2    1356
6     220
1     211
7       11
8         1
```

Name: Survey_Timely_Visits, dtype: int64

```
3    3422
4    3394
5    1388
2    1346
6     231
1     207
7       12
```

Name: Survey_Reliability, dtype: int64

```
4    3446
3    3423
2    1380
5    1308
6     219
1     211
7       13
```

Name: Survey_Options, dtype: int64

```
4    3464
3    3371
5    1403
2    1319
6     220
1     213
7       10
```

Name: Survey_Hours_of_Treatment, dtype: int64

```
4    3487
3    3456
2    1345
5    1274
1     215
6     212
7       11
```

Name: Survey_Courteous_Staff, dtype: int64

```
3    3401
4    3337
5    1429
2    1391
6     221
1     209
7       12
```

Name: Survey_Evidence_of_Active_Listening, dtype: int64

In [406...

```
#Encode categorical variables for multiple imputation in IterativeImputer. Original var

label_encode = preprocessing.LabelEncoder()

for x in Med_df[categorical_variables]:
    encode_column = str(x)+'_Encoded'
    Med_df[encode_column] = label_encode.fit_transform(Med_df[x])

Med_df_cat = Med_df[categorical_variables]
```

```
Med_df.drop(categorical_variables, axis = 1, inplace = True)
```

In [407...

```
#Use IterativeImputer from Sklearn to impute missing values in the dataset. Again, this
```

```
It_Imp = IterativeImputer(skip_complete = True, min_value = 0)
```

```
Med_df_fill = np.round(It_Imp.fit_transform(Med_df))
```

```
Med_df_fill = pd.DataFrame(Med_df_fill)
```

```
Med_df_fill.isna().sum()
```

Out[407...

```
0      0
1      0
2      0
3      0
4      0
5      0
6      0
7      0
8      0
9      0
10     0
11     0
12     0
13     0
14     0
15     0
16     0
17     0
18     0
19     0
20     0
21     0
22     0
23     0
24     0
25     0
26     0
27     0
28     0
29     0
30     0
31     0
32     0
33     0
34     0
35     0
36     0
37     0
38     0
39     0
40     0
41     0
42     0
43     0
44     0
dtype: int64
```



```
In [408... #This renames all of my columns in my newly formed DF to match their original names for

for x, y in enumerate(It_Imp.feature_names_in_):
    Med_df_fill.rename(columns= {x:y}, inplace = True)

for x in categorical_variables:
    cat_var = str(x)+'_Encoded'
    Med_df_fill[cat_var] = Med_df_fill[cat_var].astype('category')

Med_df_fill.info()
```

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

```
RangeIndex: 10000 entries, 0 to 9999
```

```
Data columns (total 45 columns):
```

#	Column	Non-Null Count	Dtype
0	Lat	10000 non-null	float64
1	Lng	10000 non-null	float64
2	Population	10000 non-null	float64
3	Children	10000 non-null	float64
4	Age	10000 non-null	float64
5	ReAdmis	10000 non-null	float64
6	VitD_levels	10000 non-null	float64
7	Doc_visits	10000 non-null	float64
8	Full_meals_eaten	10000 non-null	float64
9	VitD_supp	10000 non-null	float64
10	Soft_drink	10000 non-null	float64
11	HighBlood	10000 non-null	float64
12	Stroke	10000 non-null	float64
13	Complication_risk	10000 non-null	float64
14	Overweight	10000 non-null	float64
15	Arthritis	10000 non-null	float64
16	Diabetes	10000 non-null	float64
17	Hyperlipidemia	10000 non-null	float64
18	BackPain	10000 non-null	float64
19	Anxiety	10000 non-null	float64
20	Allergic_rhinitis	10000 non-null	float64
21	Reflux_esophagitis	10000 non-null	float64
22	Asthma	10000 non-null	float64
23	Initial_days	10000 non-null	float64
24	TotalCharge	10000 non-null	float64
25	Additional_charges	10000 non-null	float64
26	Survey_Timely_Admission	10000 non-null	float64
27	Survey_Timely_Treatment	10000 non-null	float64
28	Survey_Timely_Visits	10000 non-null	float64
29	Survey_Reliability	10000 non-null	float64
30	Survey_Options	10000 non-null	float64
31	Survey_Hours_of_Treatment	10000 non-null	float64
32	Survey_Courteous_Staff	10000 non-null	float64
33	Survey_Evidence_of_Active_Listening	10000 non-null	float64
34	City_Encoded	10000 non-null	category
35	State_Encoded	10000 non-null	category
36	County_Encoded	10000 non-null	category
37	Zip_Encoded	10000 non-null	category
38	Area_Encoded	10000 non-null	category
39	Timezone_Encoded	10000 non-null	category
40	Education_Encoded	10000 non-null	category
41	Employment_Encoded	10000 non-null	category
42	Gender_Encoded	10000 non-null	category
43	Initial_admin_Encoded	10000 non-null	category
44	Services_Encoded	10000 non-null	category

dtypes: category(11), float64(34)
memory usage: 3.3 MB

```
In [409... #Create hold df for our target variable and then drop it from our working DF

Target = Med_df_fill.ReAdmis

Med_df_fill.drop(columns='ReAdmis', axis = 1, inplace=True)

Target
```

```
Out[409...
0      0.0
1      0.0
2      0.0
3      0.0
4      0.0
...
9995   0.0
9996   1.0
9997   1.0
9998   1.0
9999   1.0
Name: ReAdmis, Length: 10000, dtype: float64
```

```
In [410... #Temporarily add back in categorical values for cleaned data set

Med_df_fill = Med_df_fill.merge(Med_df_cat, left_index=True, right_index=True)
```

```
In [411... #Export clean dataset and then drop non-numeric categorical values

Med_df_fill.to_csv(r'C:\Users\jacob.colp.UNITY\Downloads\Medical Data Raw\medical_raw_d

for x in Med_df_cat:
    Med_df_fill.drop(x, axis=1, inplace=True)
```

```
In [412... #This scales all of the numeric values relative to their mean.

StandardScaler = StandardScaler()

Scaled_Med_df = StandardScaler.fit_transform(Med_df_fill)

Scaled_Med_df = pd.DataFrame(Scaled_Med_df)

for x, y in enumerate(StandardScaler.feature_names_in_):
    Scaled_Med_df.rename(columns = {x:y}, inplace = True)

#Getting rid of qualitative variables

Scaled_Med_df = Scaled_Med_df.drop(columns=['City_Encoded', 'County_Encoded', 'State_En

Scaled_Med_df.describe()
```

```
Out[412...
      Population  Children  Age  VitD_levels  Doc_visits  Full_meals_eaten
count  1.000000e+04  1.000000e+04  1.000000e+04  1.000000e+04  1.000000e+04  1.0
```

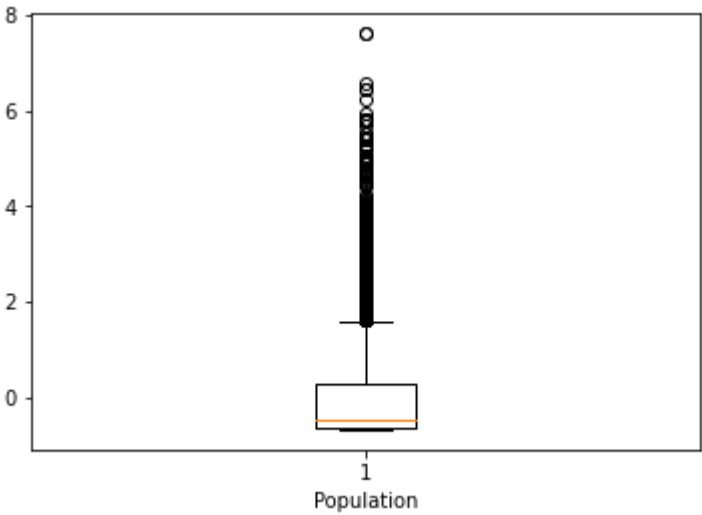
	Population	Children	Age	VitD_levels	Doc_visits	Full_meals_eaten	
mean	-1.207923e-17	1.065814e-17	-1.353584e-16	2.405187e-16	3.161915e-17	-7.815970e-17	3.6
std	1.000050e+00	1.000050e+00	1.000050e+00	1.000050e+00	1.000050e+00	1.000050e+00	1.0
min	-6.722371e-01	-1.116784e+00	-2.276538e+00	-1.397478e+00	-3.836921e+00	-9.933869e-01	-
25%	-6.253705e-01	-5.780035e-01	-8.081510e-01	-3.576545e-01	-9.679806e-01	-9.933869e-01	-
50%	-4.854456e-01	-3.922320e-02	-2.501152e-02	-2.091084e-01	-1.166703e-02	-1.388797e-03	-
75%	2.684661e-01	4.995571e-01	8.193108e-01	8.798388e-02	9.446465e-01	9.906093e-01	9.5
max	7.612562e+00	4.271019e+00	2.275461e+00	4.990006e+00	3.813587e+00	5.950600e+00	7.3

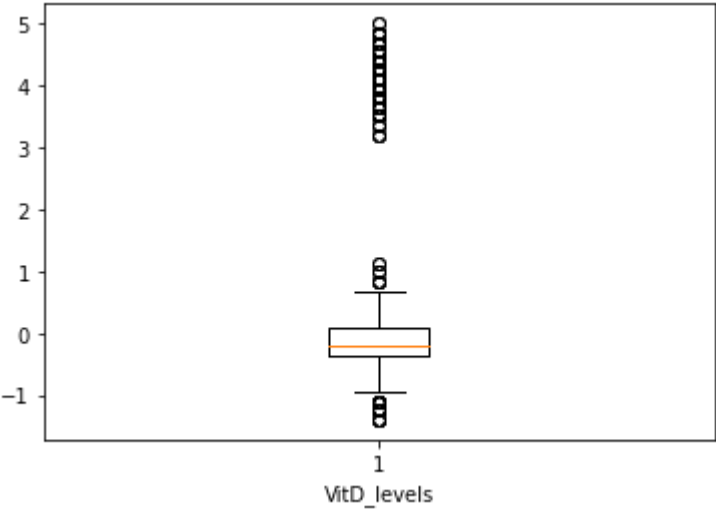
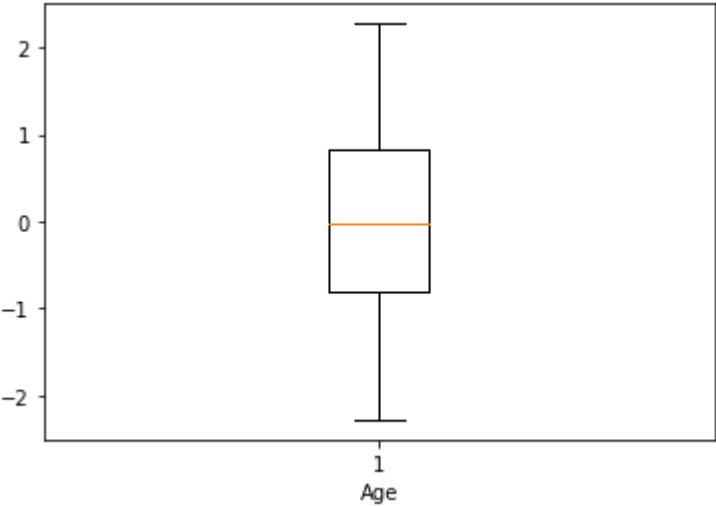
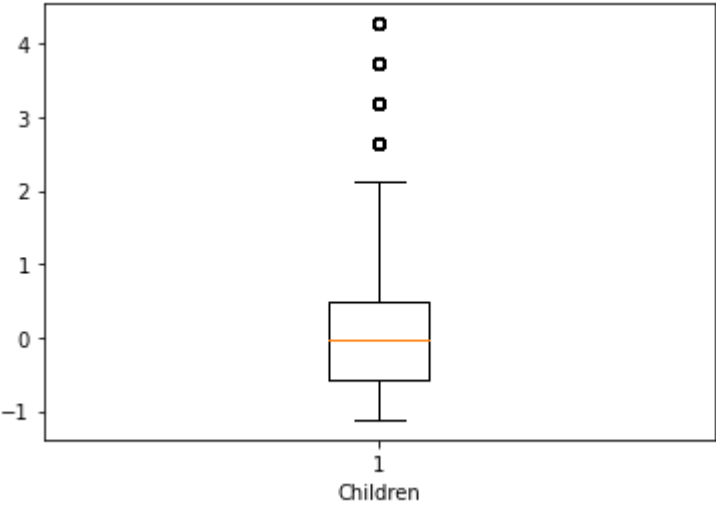
8 rows × 31 columns

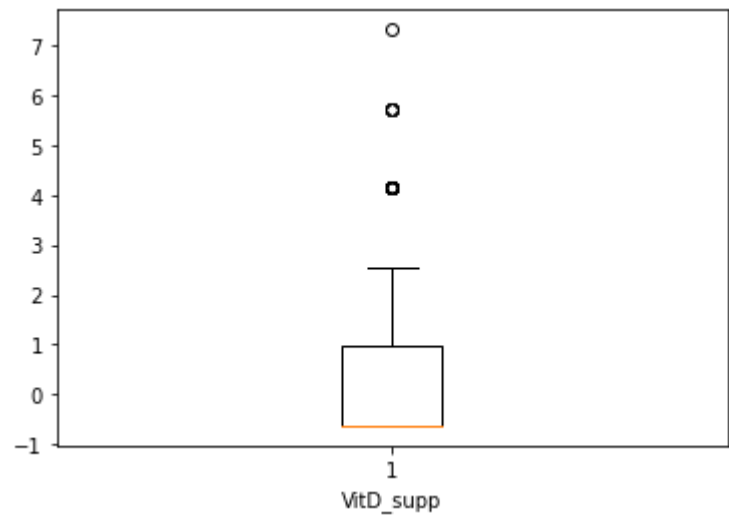
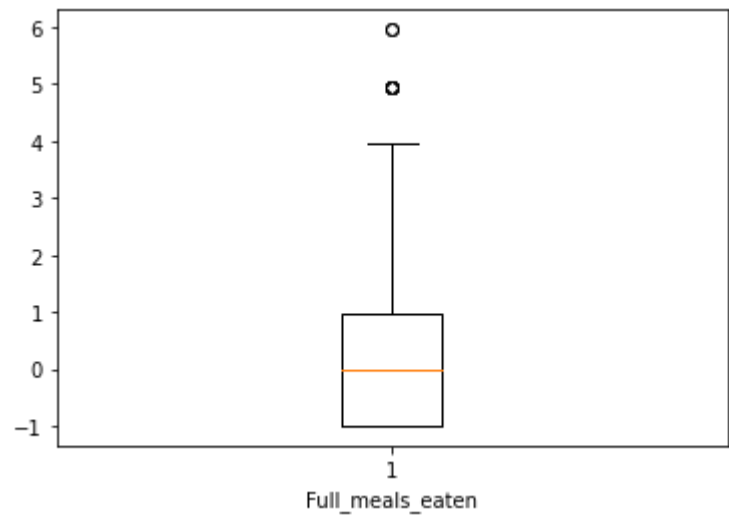
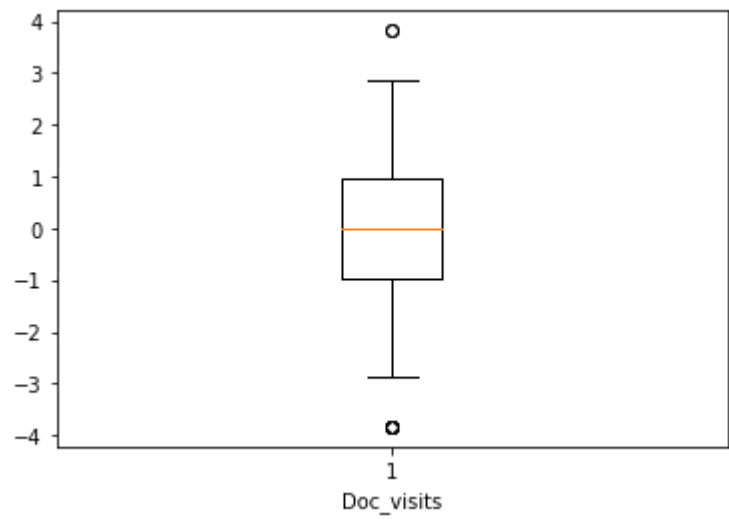
In [413...

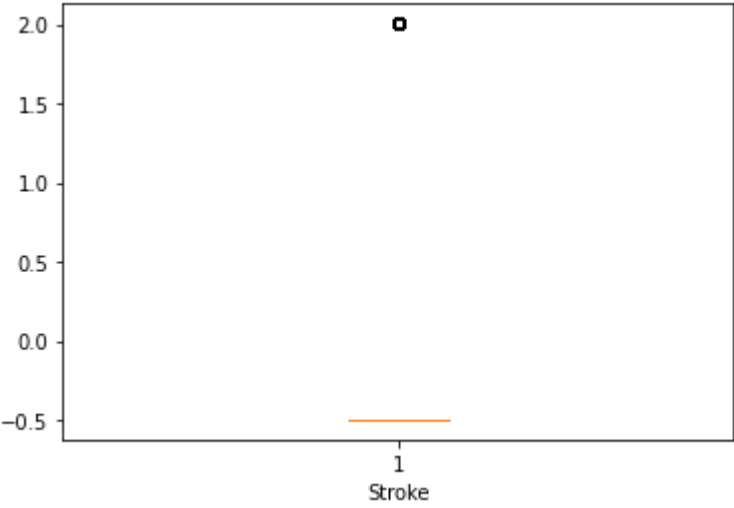
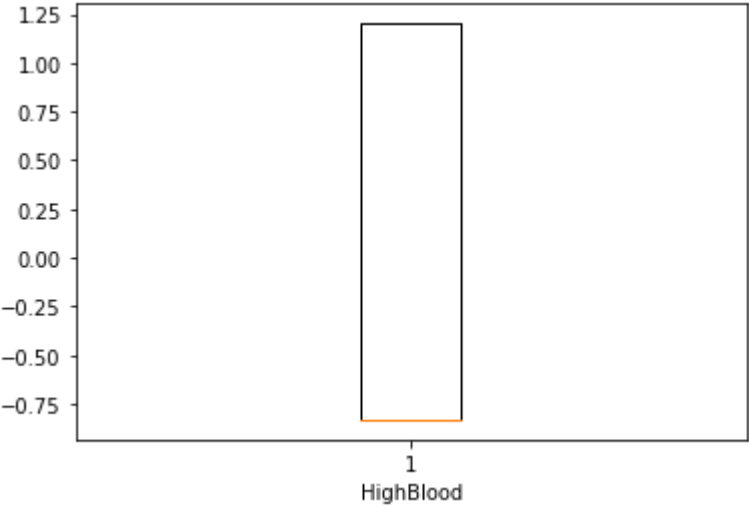
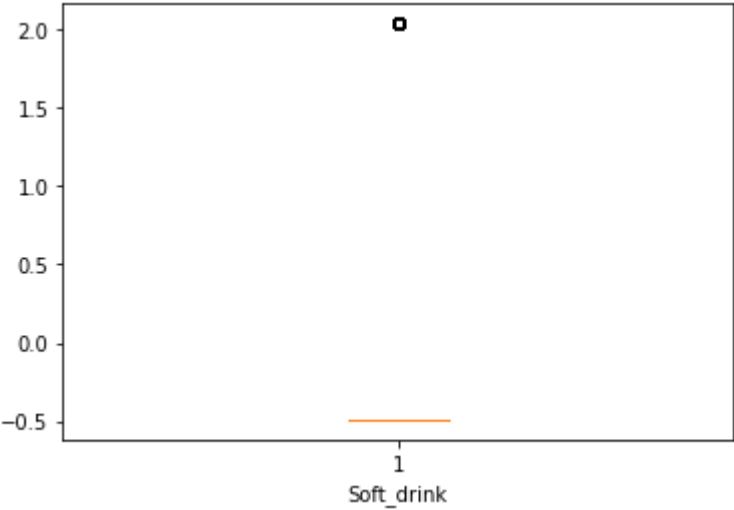
```
#Print boxplots for all non-categorical variables

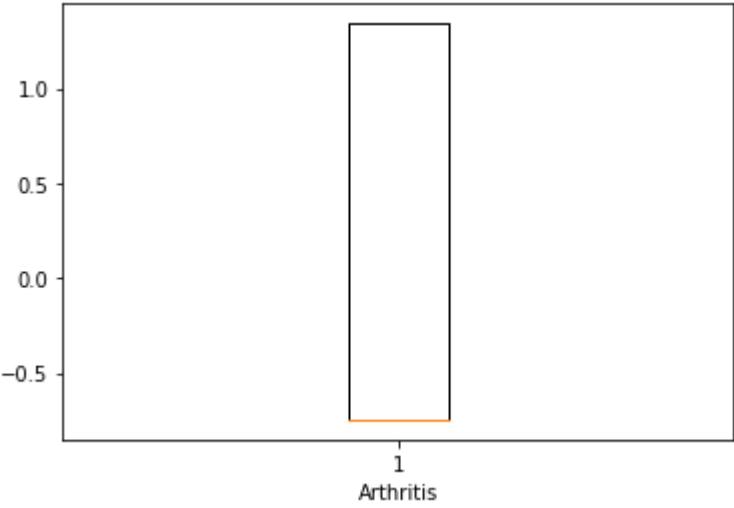
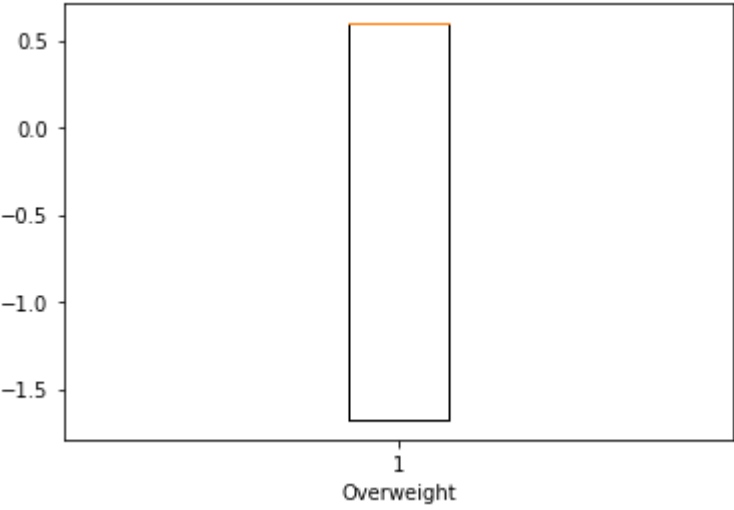
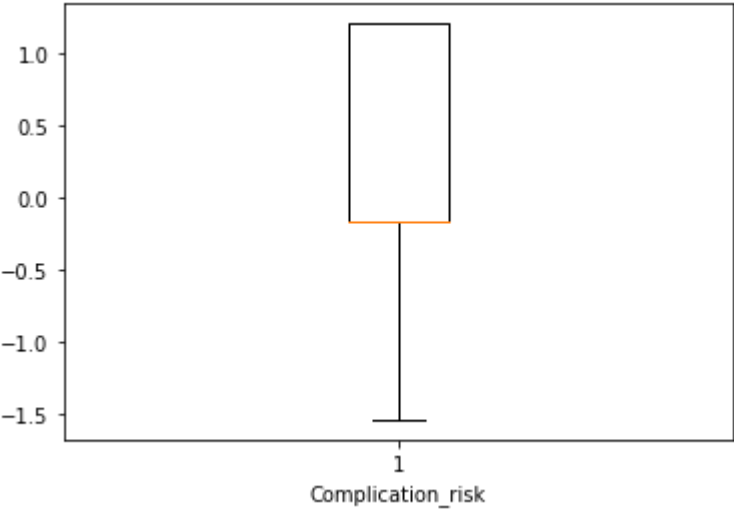
for x in Scaled_Med_df:
    pyplot.boxplot(Scaled_Med_df[x])
    pyplot.xlabel(x)
    pyplot.show()
```

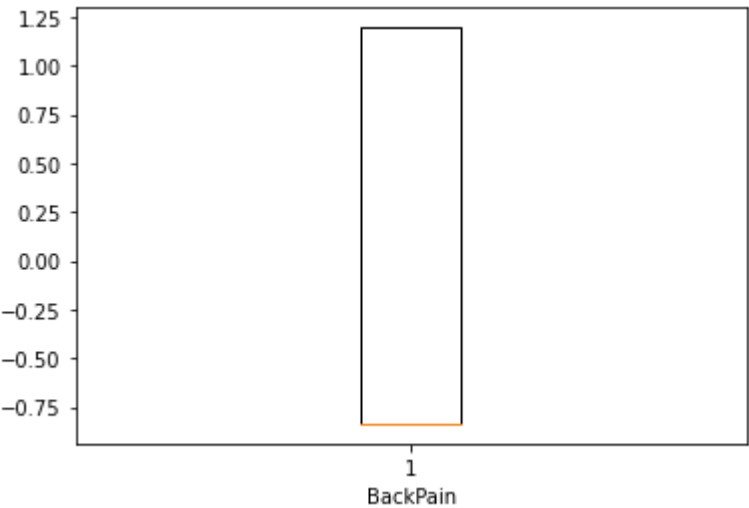
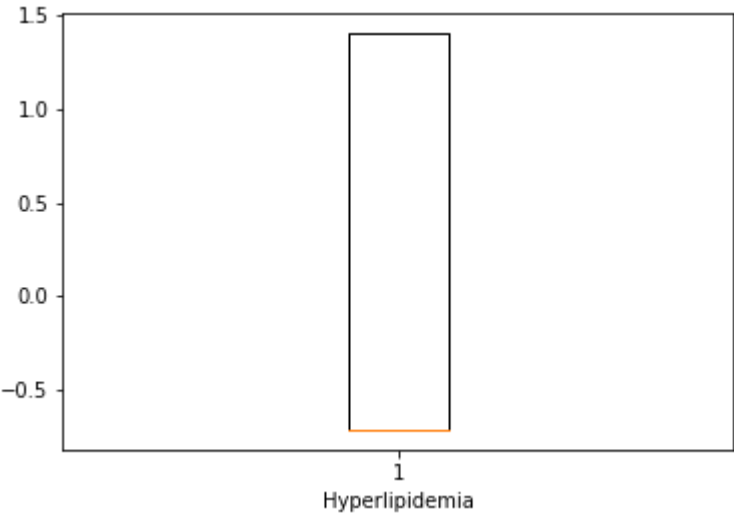
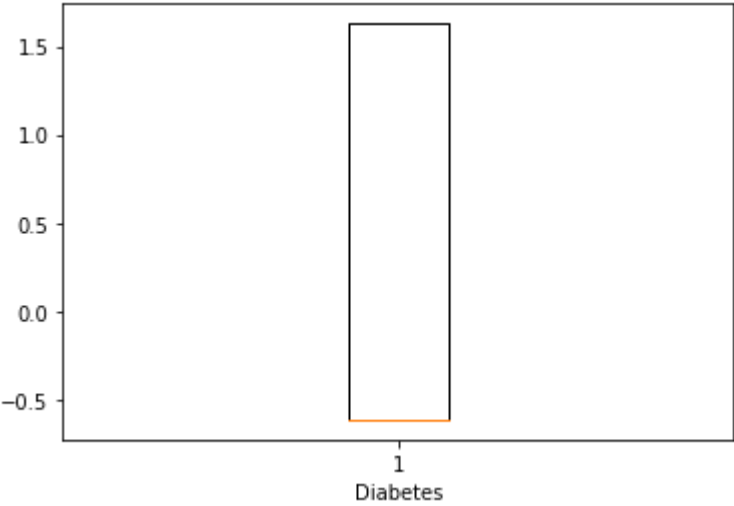


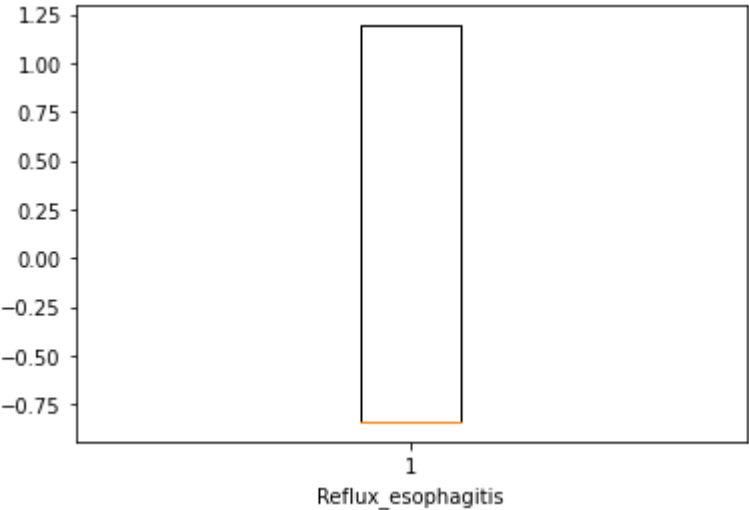
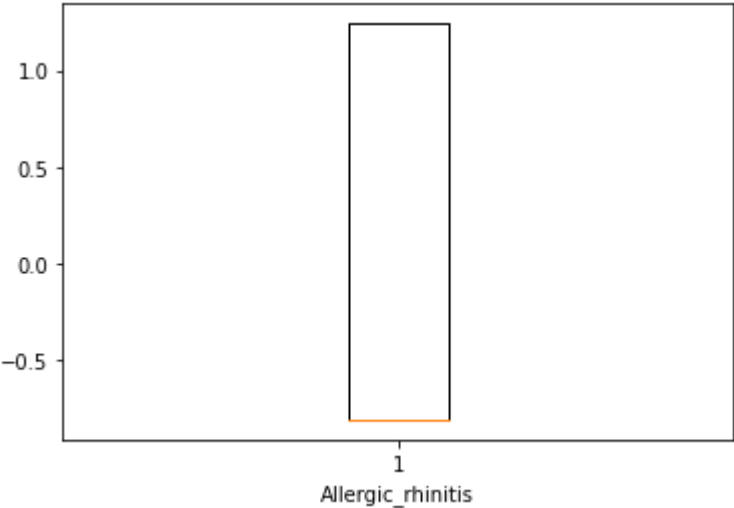


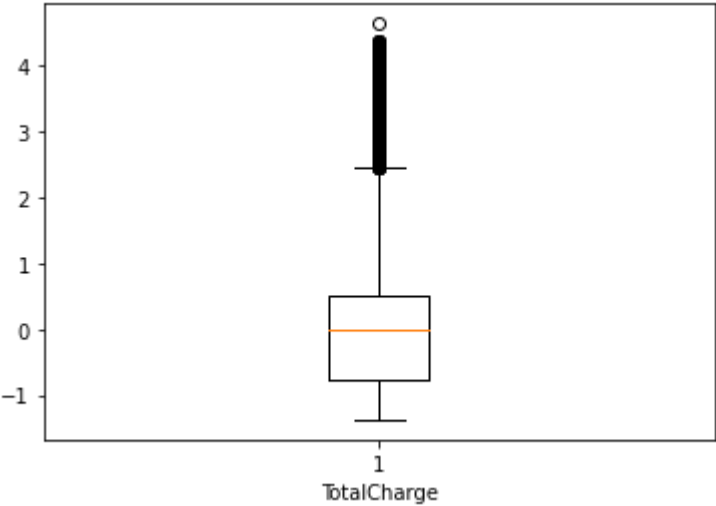
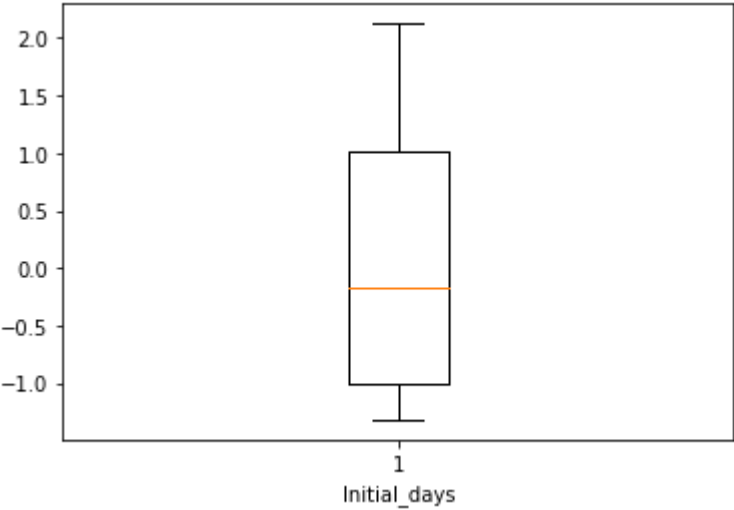
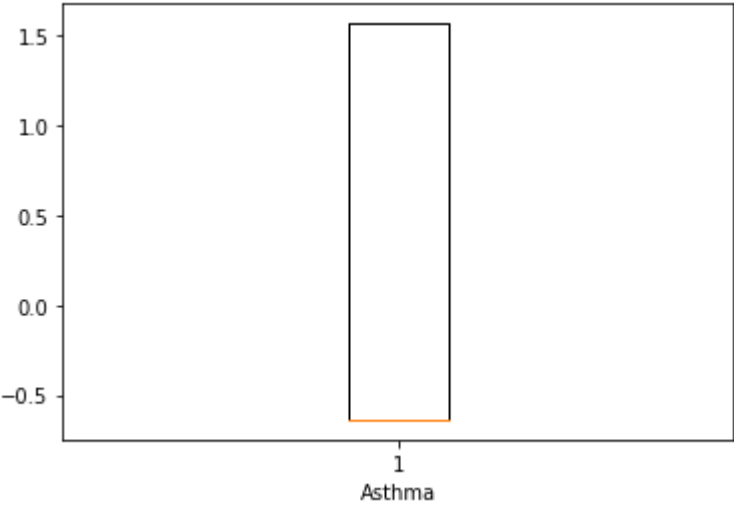


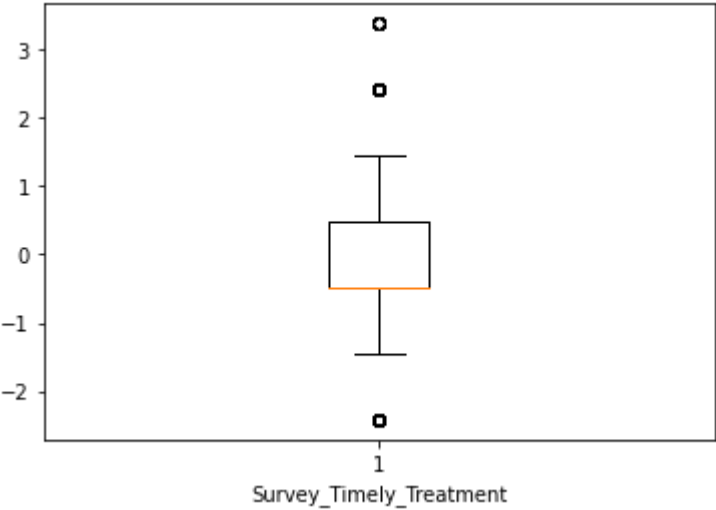
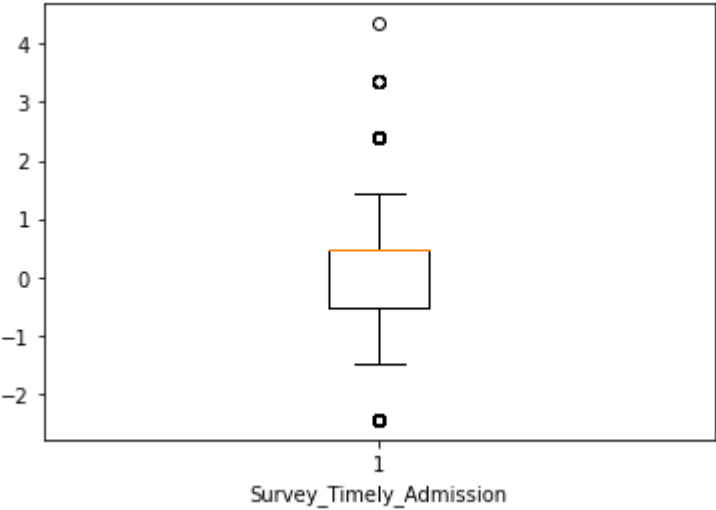
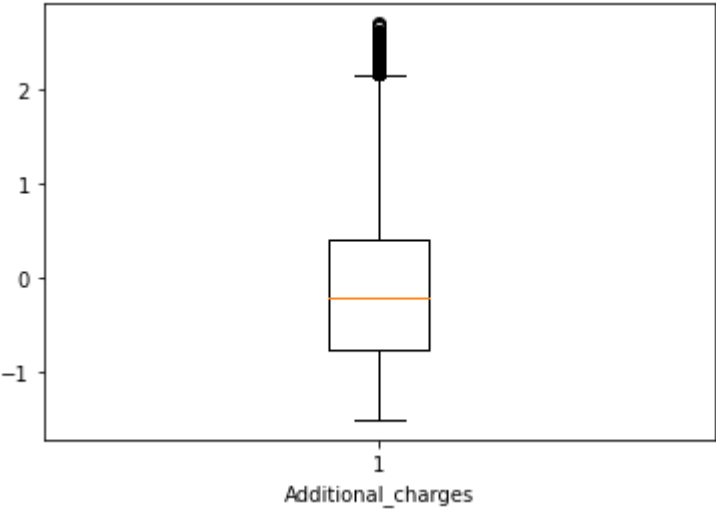


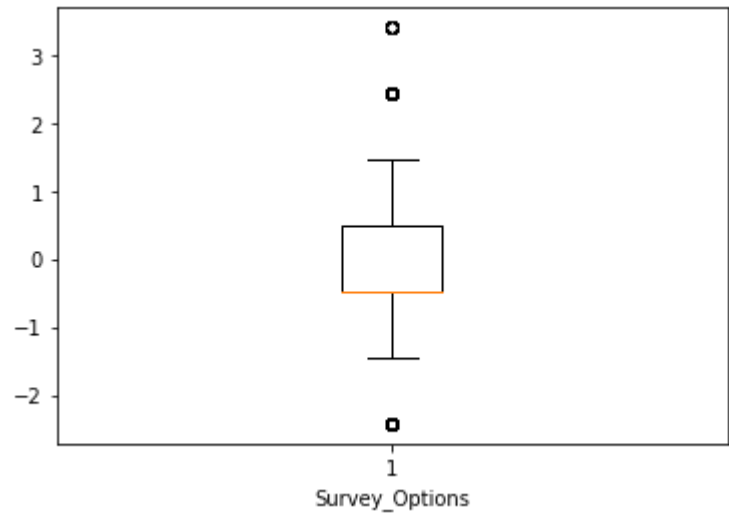
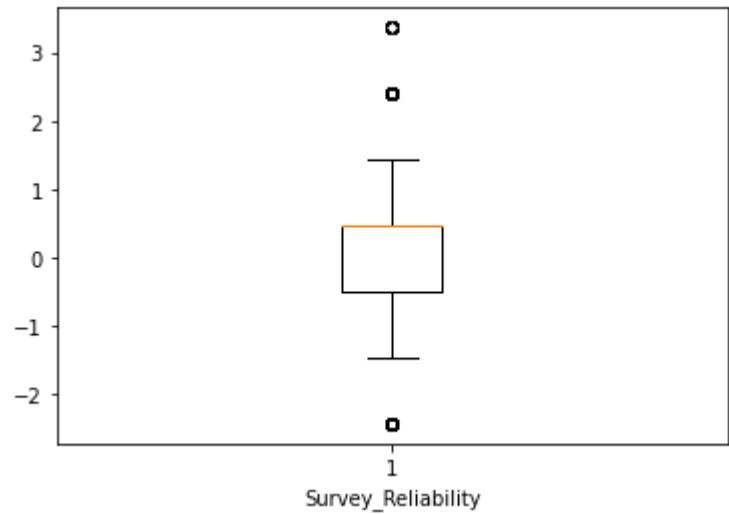
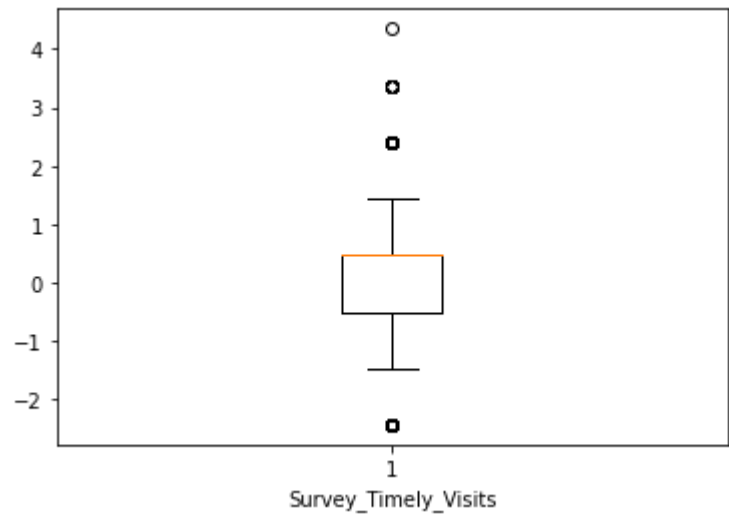


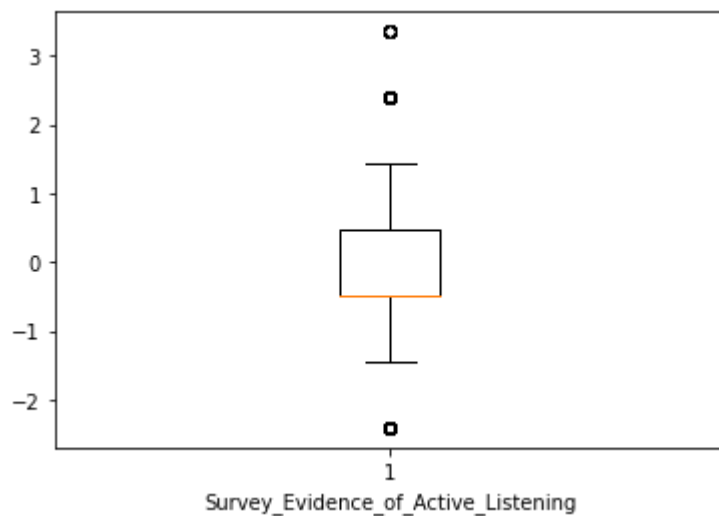
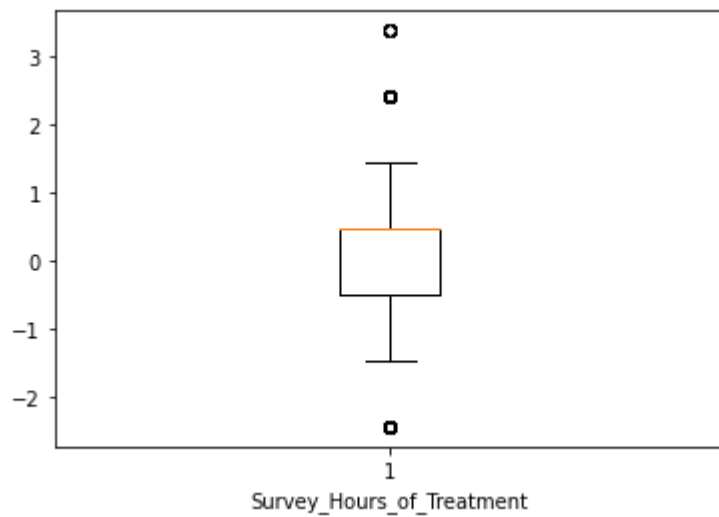








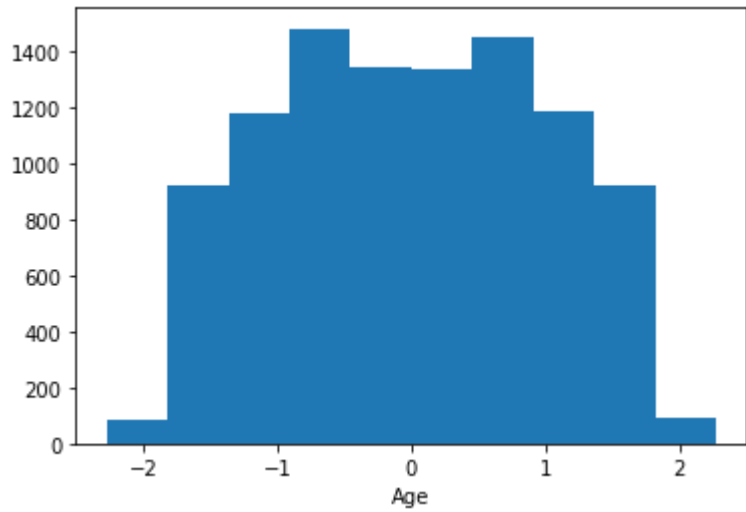
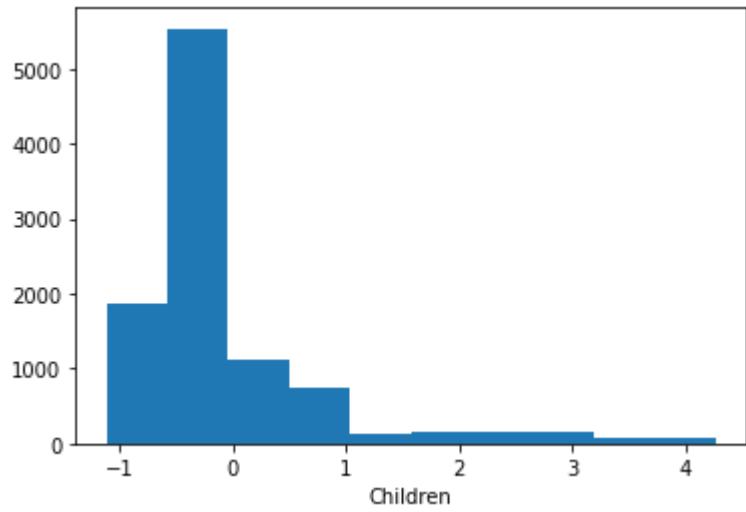
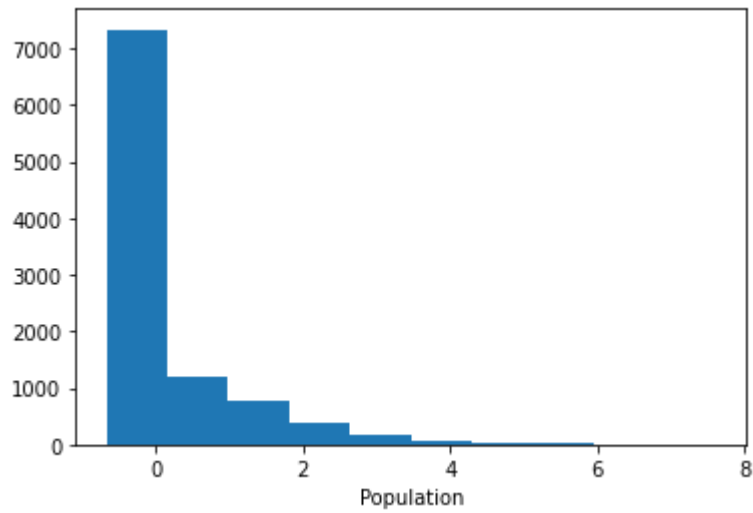


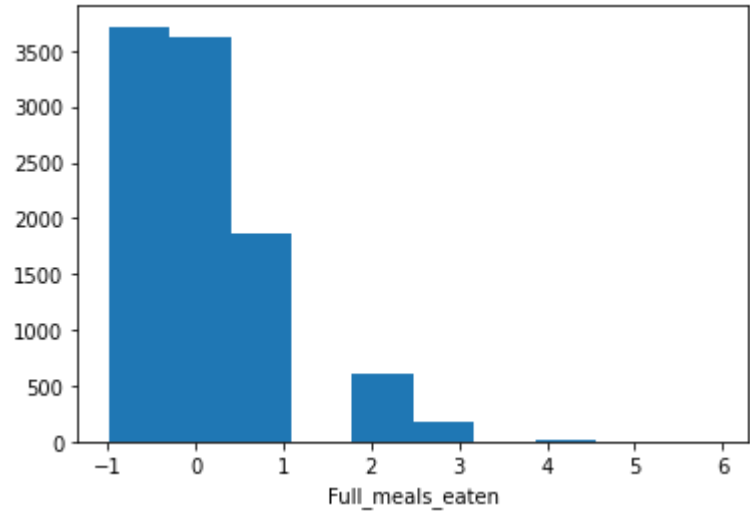
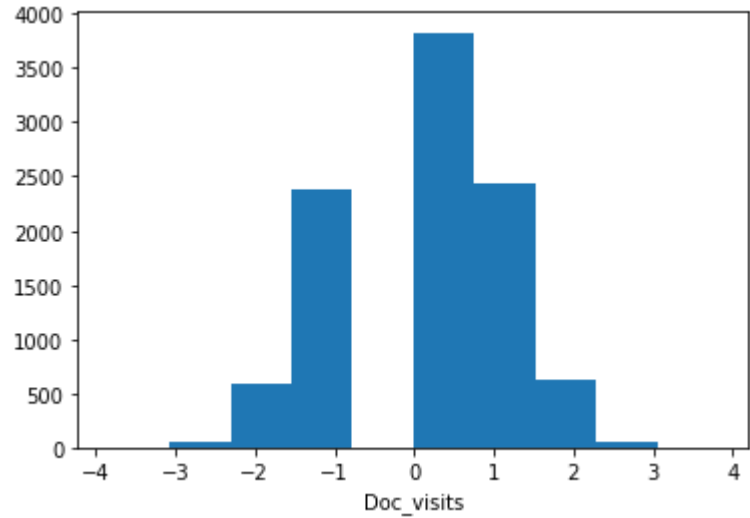
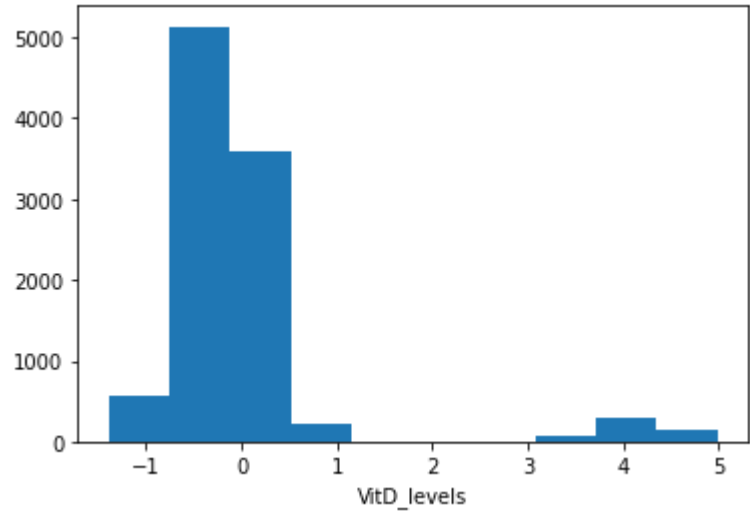


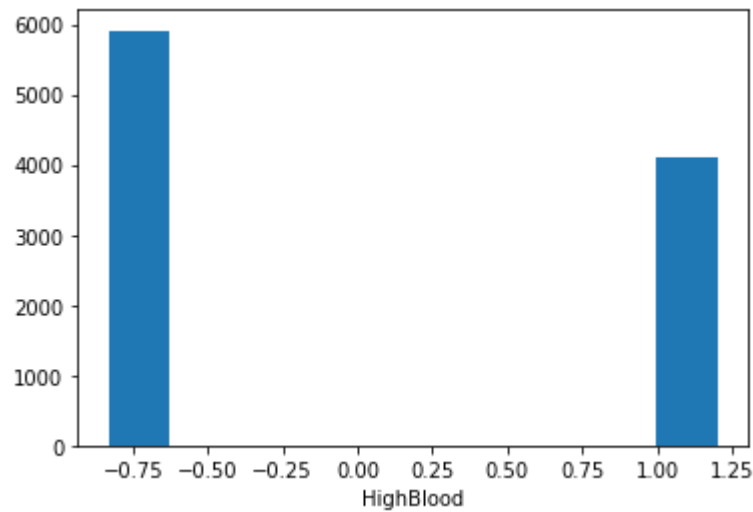
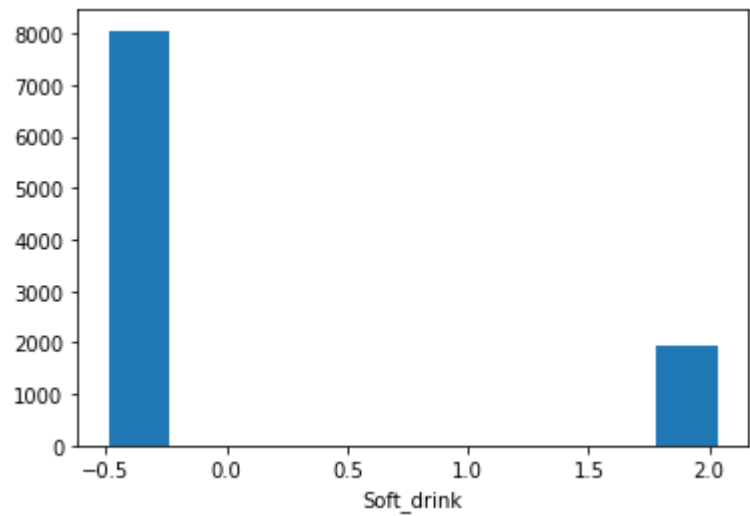
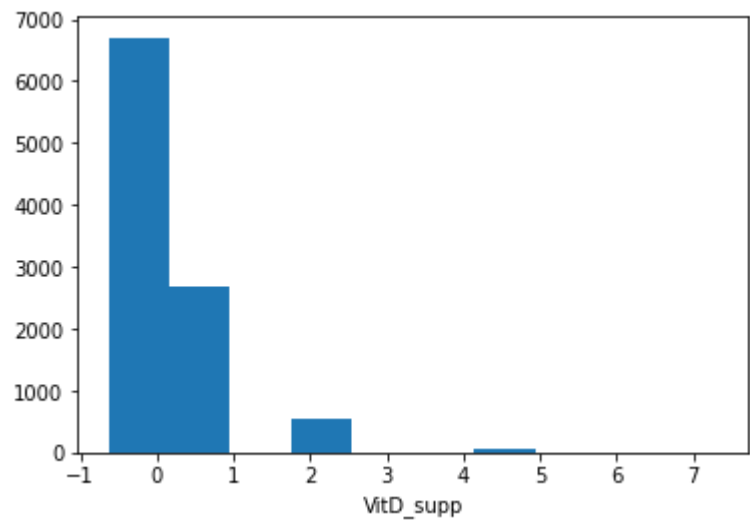
In [414...

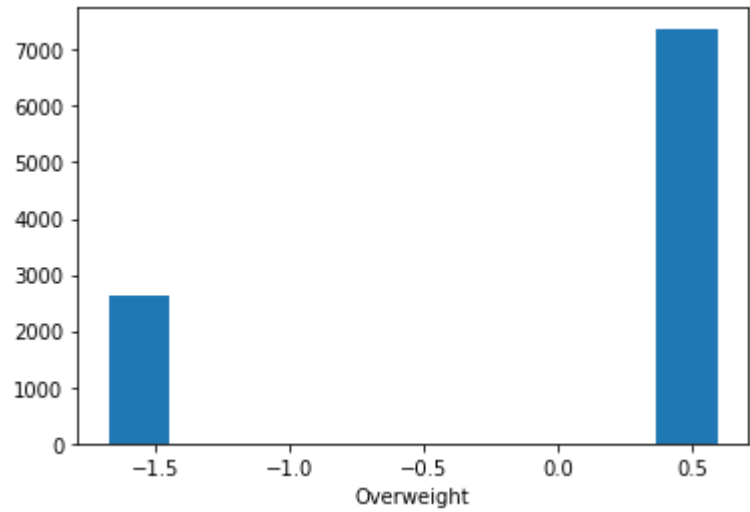
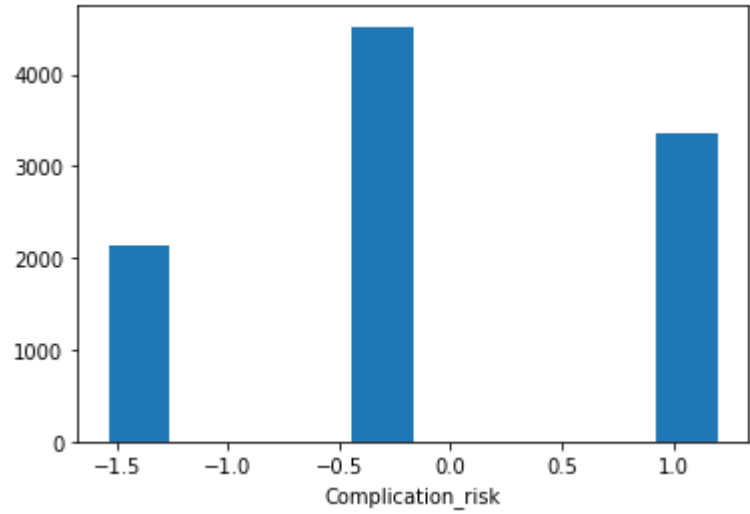
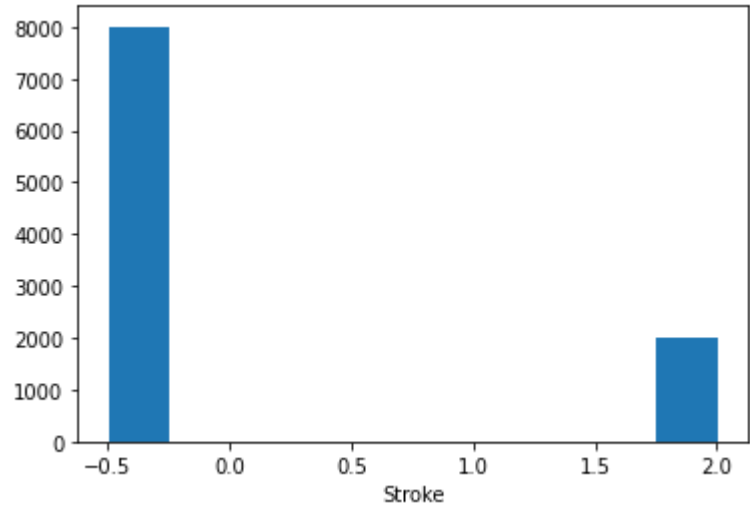
#Print histograms for all non-categorical variables

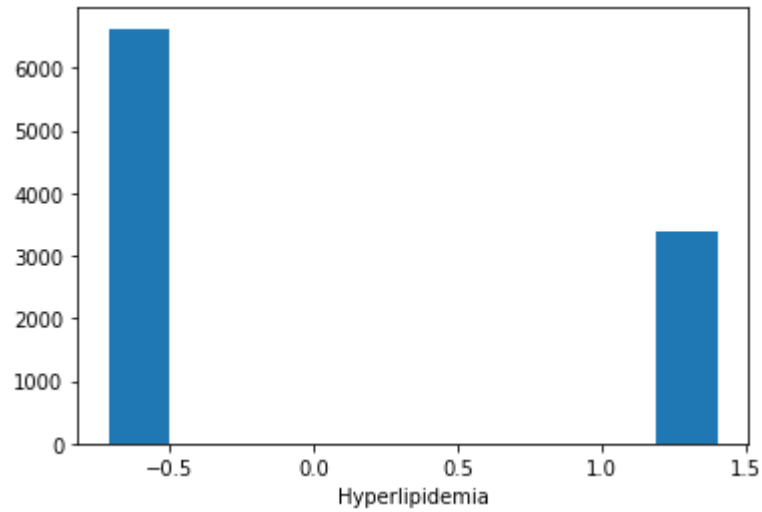
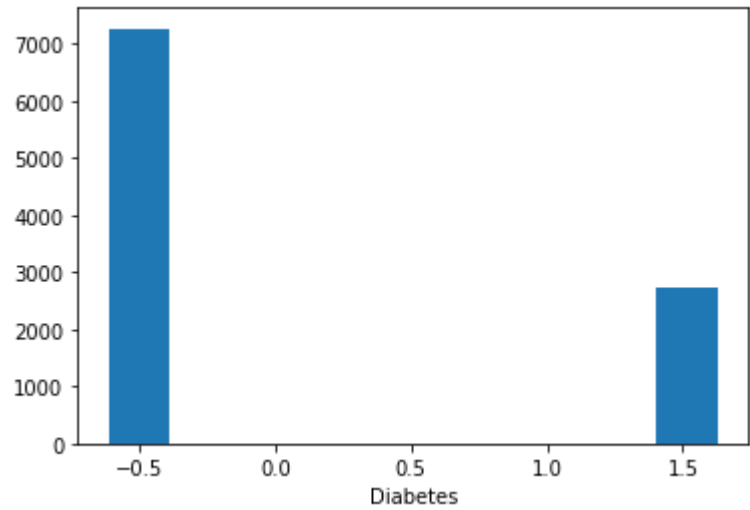
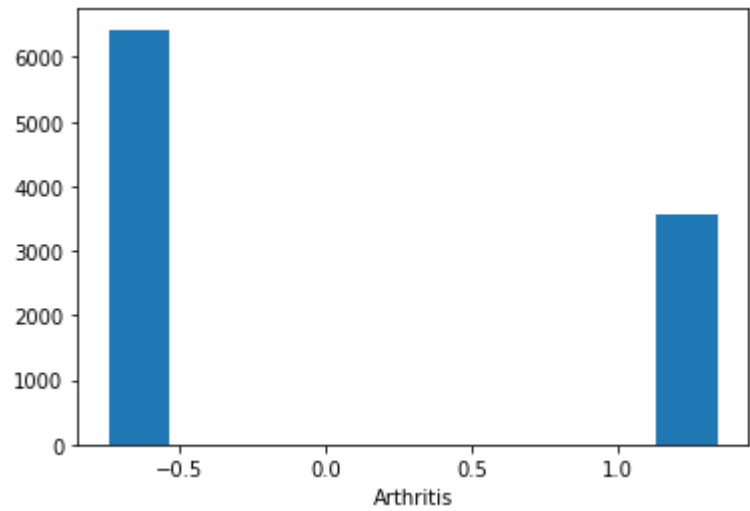
```
for x in Scaled_Med_df:  
    pyplot.hist(Scaled_Med_df[x])  
    pyplot.xlabel(x)  
    pyplot.show()
```

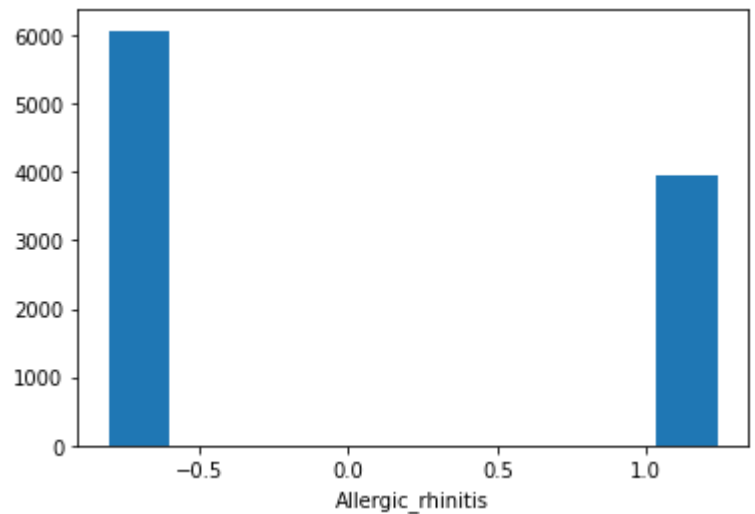
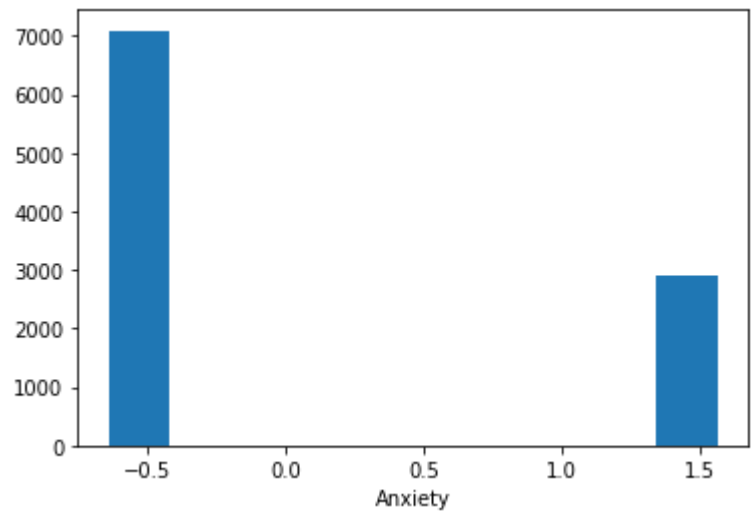
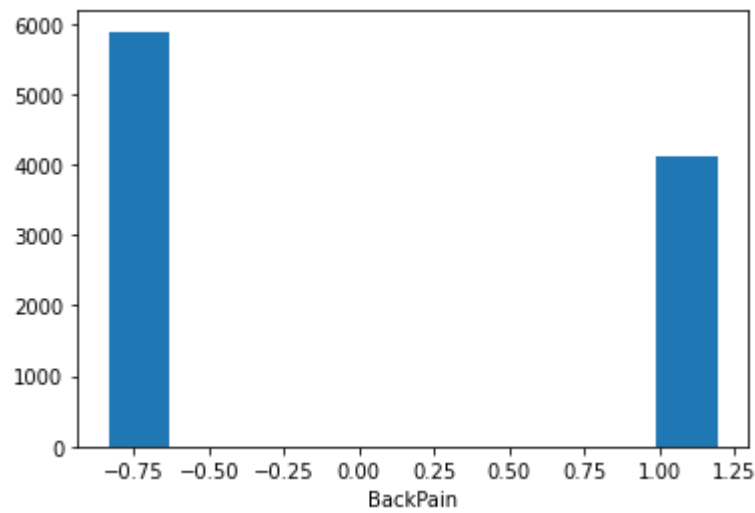


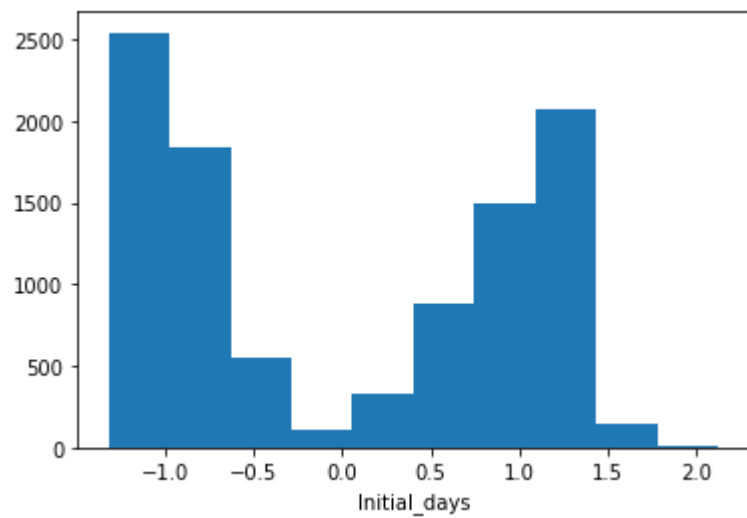
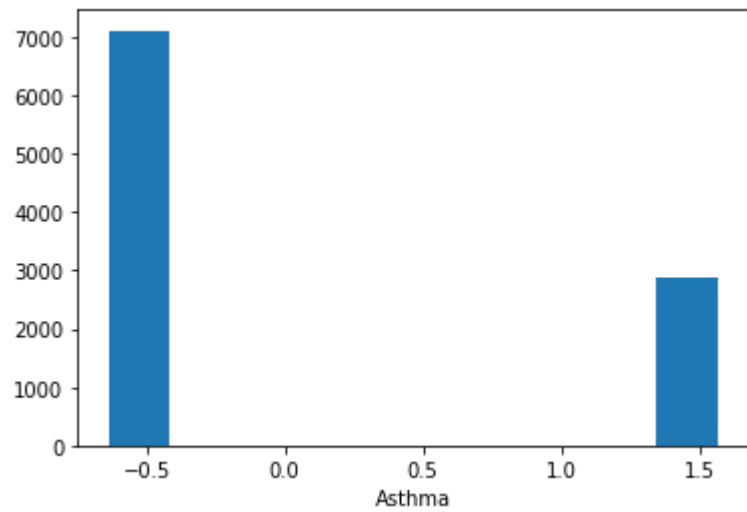
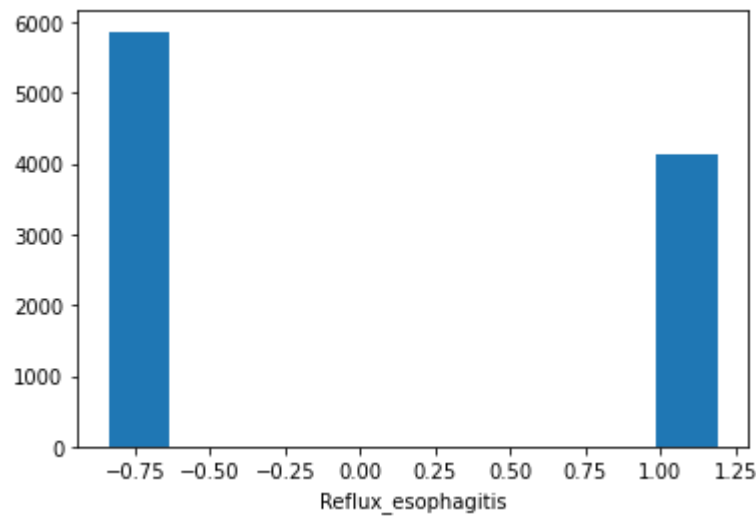


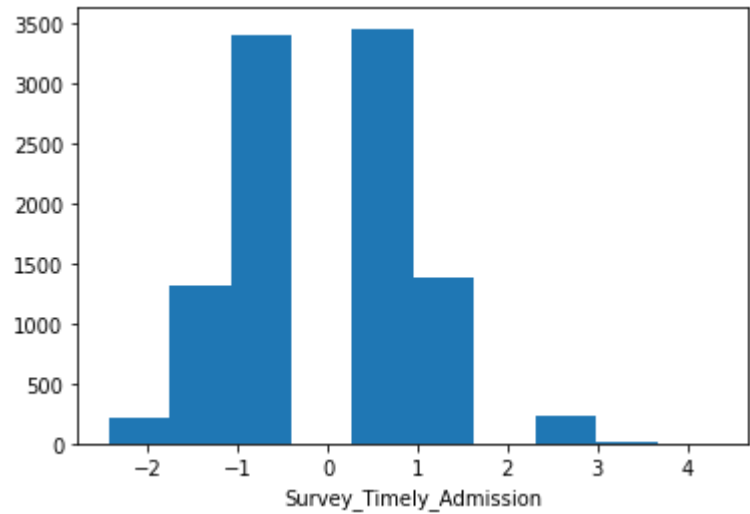
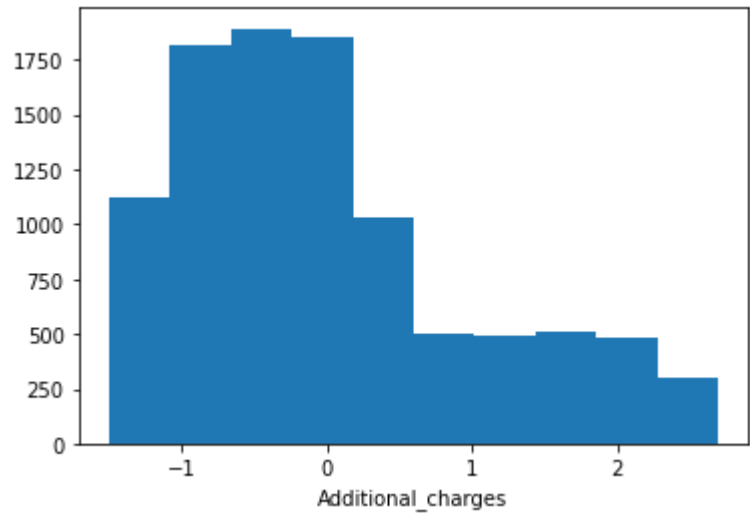
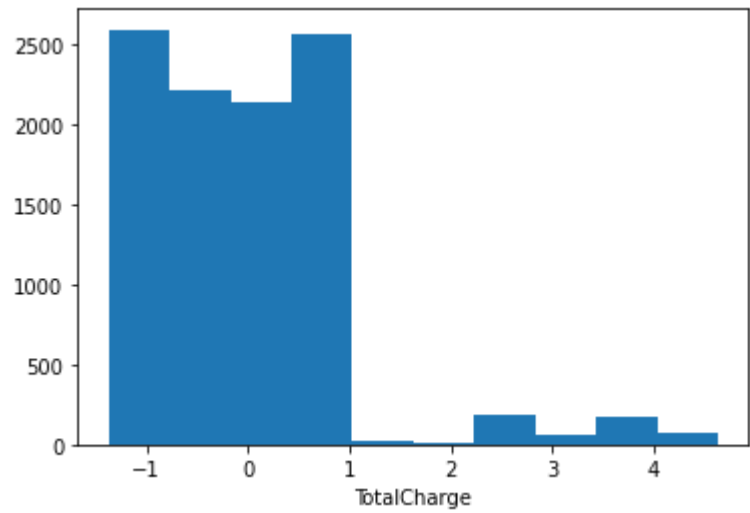


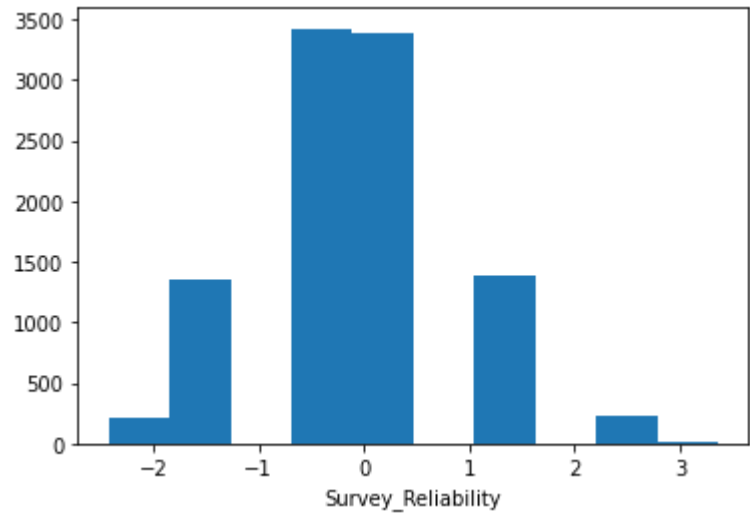
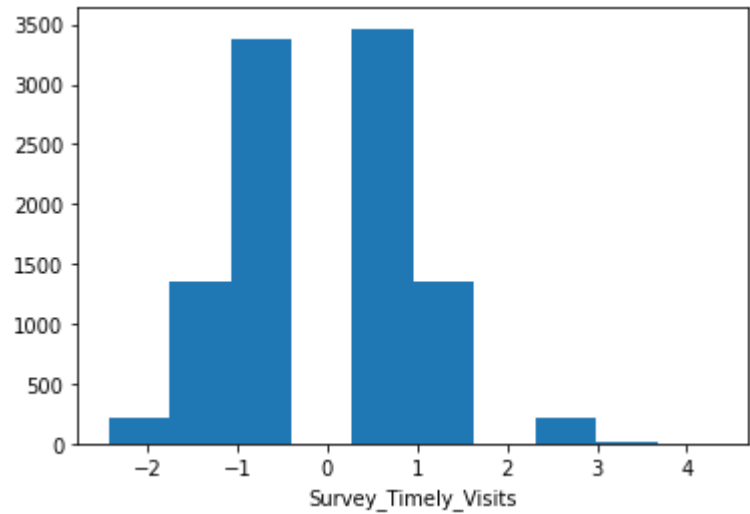
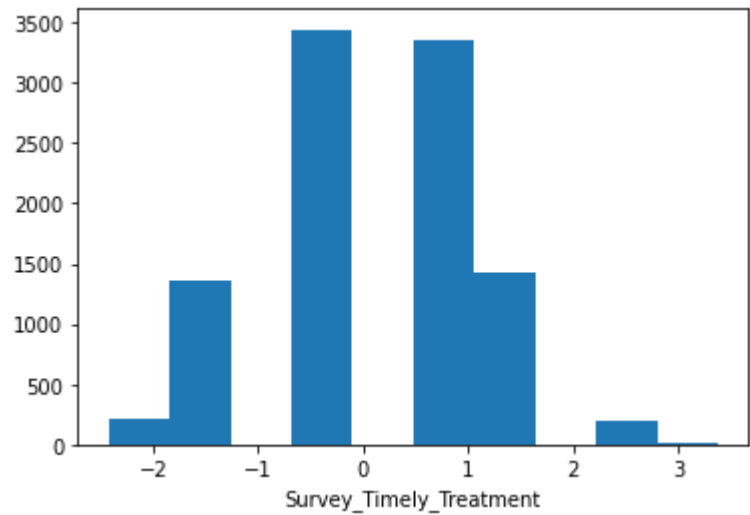


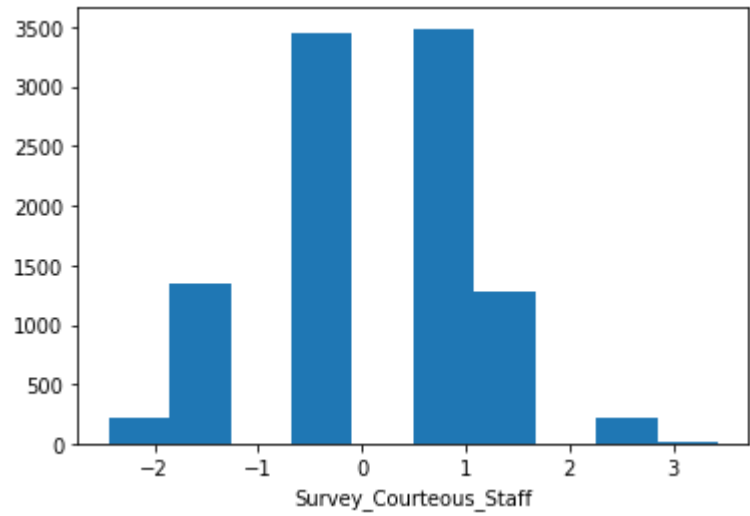
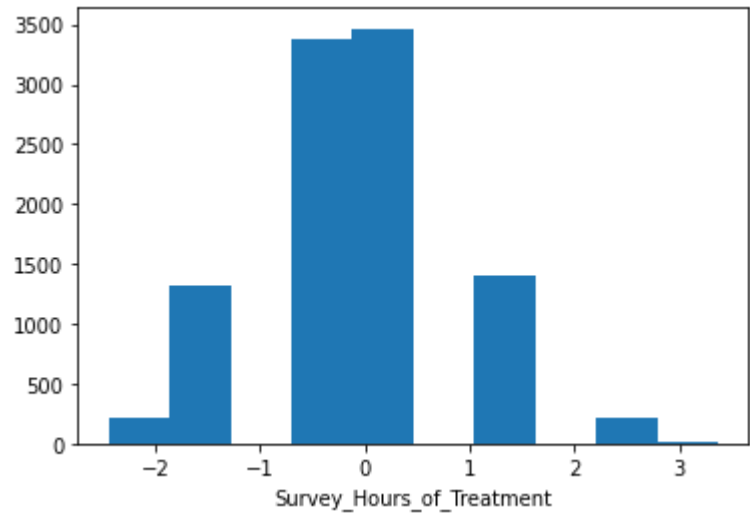
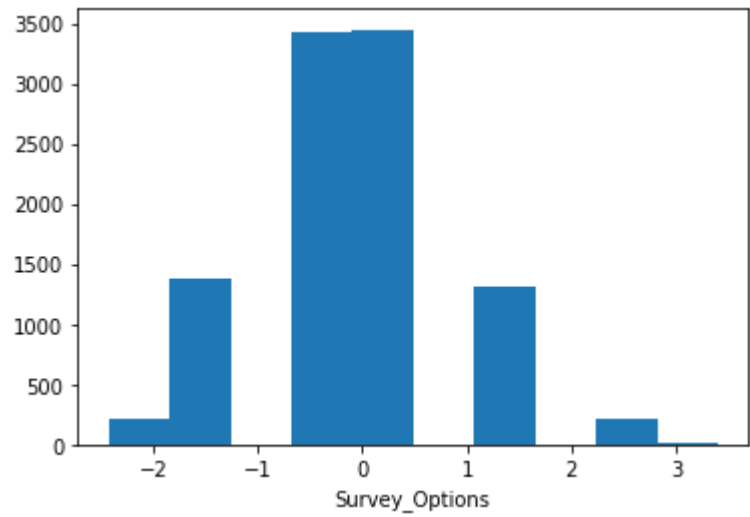


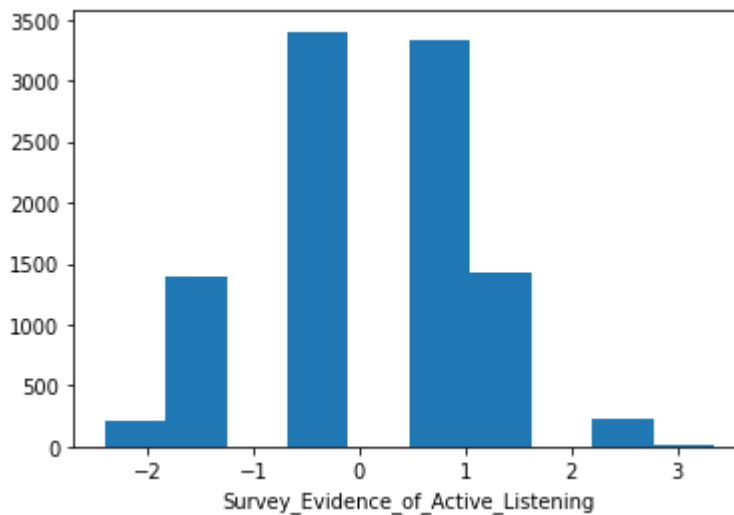












Part III: Data-Cleaning

D. Data-Cleaning Process Summary

- Orient Myself to the data
- Value Counts
- Missing Values
- Standardized numeric columns
- Categorical variables to Numeric
- Histograms and Boxplots (Outlier Detection)

D1.

I am choosing not to remove any outliers. I believe that all anomalies are a normal part of our studied population. As noted in "Guidelines for Removing and Handling Outliers in Data" - anomalies should be removed when they are either the result of an error, or they cease to represent our studied population. I do not believe any outliers to be indicative of either of those things. Therefore - I am maintaining all values. To be fair - anomalies do exist within the data. To name a few - income, education, and population. As we are evaluating a hospital - it is reasonable to expect that you are going to have broad variations on these factors. Were our research question to be centered around the readmission of low income/low education individuals, or low population areas - we could certainly scrub our data of extremely high anomalies, but as this is an evaluation of an inclusive hospital population - these will be retained.

D2.

As I have chosen to leave all data intact - there are no mitigation methods to be justified. However, these could have been dealt with using some kind of normalization. Perhaps reassigning these outliers to more representative values, or dropping them all together.

D3.

The first data cleaning step was familiarizing myself with the data. The ultimate outcome of this was me being able to identify data types and their general statistical values.

The second step was value counts. This really helps me be able to identify the unique columns that will not aid in the ultimate goal of dimensionality reduction, and gives me a general idea of what exists in each column.

Third was dealing with missing values. I utilized IterativeImputer from the sklearn package. This iterates over the data multiple times to train a model that will impute missing values based upon that variables relationship to other columns. Ultimately - this helps us to infer missing data points and provide a complete dataset to future models.

Standardized numeric columns was the fourth step that normalized all of our values around a mean of 0. This has multiple functions - it does help in scale for viewing outliers. It is also a necessary step in PCA that will be performed below. It gives the model values to recognize variance across the whole dataset.

Fifth was converting categorical variables to numeric. This is a necessity for our imputation performed above. The iterativeimputer module requires numeric values. It would also help if they were needed for another future model down the line where we wanted to evaluate their impact. Of course we would need to define that there is no ordinal value so that a model would not infer that, but again it was a necessary part of what we were doing now as well.

Finally - we have histograms and boxplots. This helps the user to be able to visually identify distribution and outliers. To me - distributions look as to be expected. And, outliers are believed to be a normal part of our studied population.

D4.

Code can be found above this section of commentary.

D5.

Cleaned data was exported above this section and will be found in project submission.

D6.

There are certainly limitations within my approach to data cleaning in this project. Ideally - with missing values, I would have preferred to have conversations with the client around potential reasons as to why they were missing in the first place. I had to make the assumption that there was nothing to correct in the process of data gathering and those values needed to be imputed. I also chose to leave anomalies as a part of my analysis moving forward. This of course could be misinterpreted by a model the line, but I again could not have a conversation with the client to determine if these were true values or the result of an error. So, I chose to err on the side of caution and include under the assumption that they were true observations. I do believe that this is ethically correct based upon the intended outcome of the project.

D7.

So, how could the limitations of my approach impact analysis? Imputation could certainly lead to a misrepresentation of the truth. Is there any 100% guarantee that the imputed values are going to be indicative of reality? No. I cannot say with 100% certainty that the imputed values are correct, however that is the nature of imputation and the nature of dealing with missing values. There is no way of getting back to truth unless you go back and confirm that observation. I believe that the ultimate impact of this is mitigated, because the total number of missing values is small relative to the total number of observations, but it is certainly context to be aware of when communicating this to stakeholders.

As for the non-action on anomalies - this could lead to the outliers impacting our model by providing extremes. However, as they are assumed-to-be-true values - they need to be included as to have accurate information relative to our population. So, while this may lead to some differences in the model compared to their exclusion - it does fall in line with best practices.

In [415...

```
#This is our correlation matrix to show the relationship between variables
```

```
Scaled_Med_df.corr()
```

Out[415...

	Population	Children	Age	VitD_levels	Doc_visits	Full_meals_
Population	1.000000	0.007205	-0.018371	0.002124	0.012646	-0.0
Children	0.007205	1.000000	0.005393	-0.002391	-0.004734	-0.0
Age	-0.018371	0.005393	1.000000	0.019033	0.005166	0.0
VitD_levels	0.002124	-0.002391	0.019033	1.000000	0.001367	0.0
Doc_visits	0.012646	-0.004734	0.005166	0.001367	1.000000	-0.0
Full_meals_eaten	-0.025608	-0.000856	0.010050	0.009170	-0.002767	1.0
VitD_supp	0.009781	-0.000463	0.008860	0.009991	0.005681	-0.0
Soft_drink	0.004115	0.007961	0.004361	-0.000697	0.017951	0.0
HighBlood	0.009764	0.004411	0.010891	0.004970	0.008967	0.0
Stroke	-0.001690	0.004573	0.013436	-0.009912	-0.002230	0.0
Complication_risk	0.015936	-0.003008	-0.000490	0.005316	0.012306	0.0
Overweight	0.000367	-0.021073	-0.007507	-0.007787	0.001087	-0.0
Arthritis	0.000055	0.004476	0.008995	-0.000469	-0.000719	0.0
Diabetes	-0.009975	0.018689	0.005136	-0.023462	0.012781	0.0
Hyperlipidemia	-0.006222	-0.009856	0.003964	0.000824	-0.026730	0.0
BackPain	0.006437	-0.023047	0.019881	-0.003450	0.008514	-0.0
Anxiety	-0.012899	0.005274	0.008150	0.014533	-0.002834	0.0
Allergic_rhinitis	0.007681	-0.019174	0.014716	-0.002394	0.002920	0.0
Reflux_esophagitis	0.014340	0.004483	-0.016896	-0.007717	-0.005330	-0.0
Asthma	-0.001510	0.005987	0.009301	0.011450	-0.017989	0.0

	Population	Children	Age	VitD_levels	Doc_visits	Full_meals_
Initial_days	0.019087	0.010640	0.016488	0.008069	-0.007615	-0.0
TotalCharge	0.013751	0.003058	0.024185	0.727561	-0.004515	-0.0
Additional_charges	-0.004820	0.009551	0.728847	0.016425	0.008072	0.0
Survey_Timely_Admission	0.014312	0.005446	0.004909	-0.004335	0.003680	0.0
Survey_Timely_Treatment	0.023612	0.010784	0.002969	-0.017683	0.006024	-0.0
Survey_Timely_Visits	-0.001248	0.002291	0.004855	-0.012496	-0.002718	0.0
Survey_Reliability	-0.004660	0.005824	0.005639	0.012671	-0.006538	-0.0
Survey_Options	0.008705	0.006133	-0.007860	-0.012255	-0.009434	0.0
Survey_Hours_of_Treatment	0.008159	-0.002160	-0.000156	0.007350	0.012530	0.0
Survey_Courteous_Staff	0.010034	0.005558	0.010302	0.001992	0.008589	0.0
Survey_Evidence_of_Active_Listening	-0.000220	-0.011324	-0.003642	0.004033	0.004571	-0.0

31 rows × 31 columns

In [416...

```
#Initiating PCA
#Reference: (Kindsonthegenius, 2020)
Med_PCA = PCA()

Med_PCA_FT = Med_PCA.fit_transform(Scaled_Med_df)

Med_PCA_FT = pd.DataFrame(Med_PCA_FT)

Med_PCA_FT
```

Out[416...

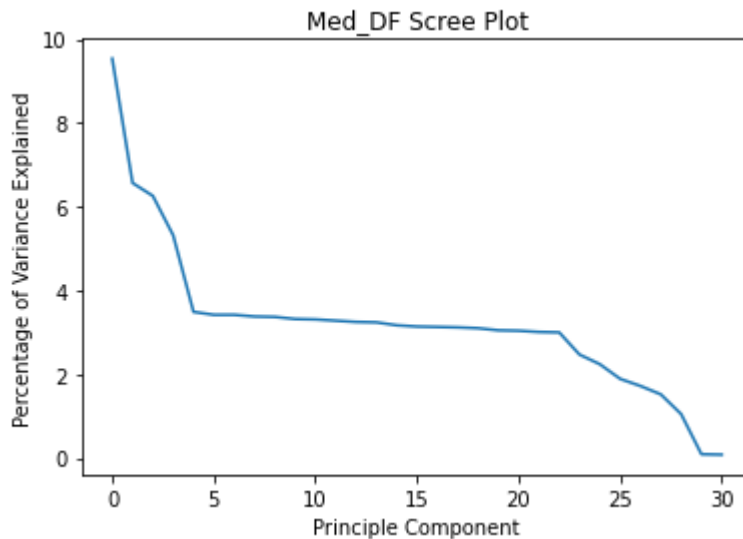
	0	1	2	3	4	5	6	7	8
0	-1.606819	0.087476	-1.603797	0.352726	1.125682	-0.824507	0.640937	-1.515131	0.697186
1	-0.325365	0.219360	-1.222288	-0.104782	-0.962116	-0.164027	-1.638954	0.928217	0.414256
2	-0.201751	-0.355915	-1.831390	-0.660291	0.278594	-0.051720	-1.148161	0.144641	0.345984
3	2.345966	-0.854095	-1.201547	0.391742	1.022480	0.349036	-0.389637	-0.557362	-0.722921
4	-2.372715	-2.752463	0.395789	-0.337430	-0.334034	1.333370	0.521480	1.905066	-1.221259
...
9995	-2.107126	-0.304723	0.550835	-0.270944	-0.520516	-0.188865	-1.047018	-0.071344	0.030868
9996	-0.703980	2.668207	-1.576351	1.724087	1.016712	0.329201	-0.810658	-1.916748	0.112562
9997	-1.880049	0.880328	0.037438	0.261843	0.846171	-0.806931	-0.263213	0.499075	0.946759
9998	0.796273	-0.095912	1.684987	0.896192	-0.497747	-0.143362	0.280589	0.528423	-0.508669
9999	0.676988	0.715358	1.253763	0.224172	-0.985252	2.139548	-1.036507	-1.679030	-0.553948

10000 rows × 31 columns

In [417...

```
percent_variance = np.round(Med_PCA.explained_variance_ratio_*100, decimals=2)

pyplot.plot(Med_PCA_FT.columns, percent_variance)
pyplot.title('Med_DF Scree Plot')
pyplot.xlabel('Principle Component')
pyplot.ylabel('Percentage of Variance Explained')
plt.show()
```



In [418...

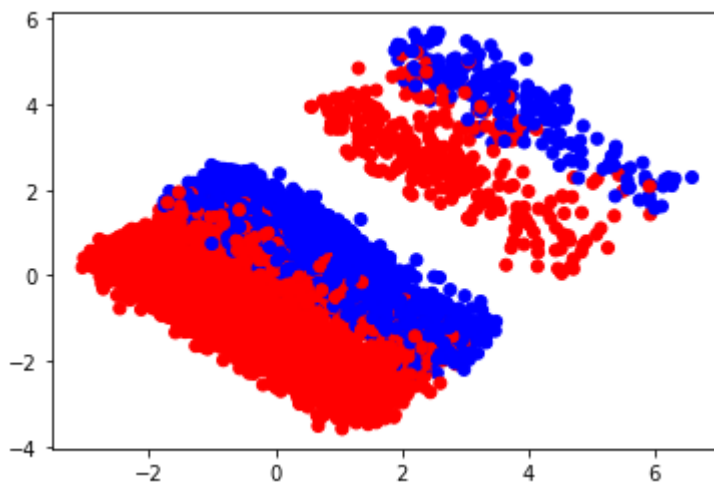
```
Med_PCA_FT = Med_PCA_FT.merge(Target, left_index=True, right_index=True)

Med_PCA_FT['Color'] = Med_PCA_FT['ReAdmis'].replace({0:'r', 1:'b'})

pyplot.scatter(Med_PCA_FT[1], Med_PCA_FT[2], c = Med_PCA_FT.Color)
```

Out[418...

<matplotlib.collections.PathCollection at 0x2218b651e50>



In [419...

```
for x, y in zip(Med_PCA_FT.columns, np.cumsum(Med_PCA.explained_variance_ratio_)):
    print(x, y)
```

```
0 0.09537300264151131
1 0.1611225273178129
```

2 0.2237484613314441
3 0.2769622917197733
4 0.3118530427783067
5 0.34609321081467875
6 0.38024913482062234
7 0.41407433387307635
8 0.44781038979305654
9 0.4810487349915523
10 0.5141275727512862
11 0.5468873589467472
12 0.5793946548583053
13 0.611750188981412
14 0.6434984437742507
15 0.6749148229595627
16 0.7061959206457349
17 0.7374189541249802
18 0.7683845215516597
19 0.7988632337291206
20 0.8292759952287545
21 0.8593506065006781
22 0.8893644401291758
23 0.9140778533385729
24 0.9364863689582283
25 0.9554105521881042
26 0.9725630552530493
27 0.9877610200006122
28 0.9983011849043533
29 0.9992245607476545
30 1.0000000000000002

In [420...

```
#This gives me all our initial variables and their loadings in each principle component

loadings = pd.DataFrame(Med_PCA.components_.T*100, index= Scaled_Med_df.columns)

loadings
```

Out[420...

	0	1	2	3	4	!
Population	1.026974	0.665143	2.180759	2.527812	-38.867162	9.32201
Children	0.252623	1.025776	-0.351872	1.369559	17.525652	40.57281
Age	0.625204	38.748640	-35.165643	2.854789	12.690066	17.42189
VitD_levels	-0.957820	36.082300	38.019565	3.446524	12.190871	-23.28975
Doc_visits	0.706793	0.613241	-1.351132	-0.533841	-13.355231	-3.77970
Full_meals_eaten	-0.051667	1.052641	-2.852055	2.213620	42.384042	-20.99512
VitD_supp	-0.467146	2.892635	1.281249	0.797140	-35.693543	18.59583
Soft_drink	0.689039	-0.001508	-0.066078	1.389166	25.409101	9.80328
HighBlood	-0.340725	33.691132	-33.268180	0.445878	-15.031272	-15.87263
Stroke	-0.235984	1.392486	-3.759940	1.544252	-2.891159	11.73140
Complication_risk	1.285496	4.839999	-0.920201	-0.953717	-15.346986	-10.70158
Overweight	0.426311	-0.103823	-2.925938	1.157815	-18.473379	-50.03268
Arthritis	-1.422671	2.196020	0.959379	-0.741561	9.220582	1.10043

	0	1	2	3	4	5
Diabetes	-0.295194	-1.237835	-1.666642	3.139808	35.350400	20.05220
Hyperlipidemia	1.707249	0.334629	1.237430	-1.475862	11.077778	21.05807
BackPain	-1.299412	2.929841	0.178080	-0.655148	-12.771244	-6.91411
Anxiety	-0.075996	3.347648	1.315573	-2.262854	20.985467	-4.94182
Allergic_rhinitis	0.481664	2.343741	-1.145508	1.778240	-2.767455	-25.23259
Reflux_esophagitis	0.632136	-0.241043	2.315297	-1.112417	-11.421476	21.44813
Asthma	-1.063100	0.931665	-1.676664	2.043914	28.676481	-29.09239
Initial_days	-1.986170	31.075724	34.955306	6.561606	-11.140513	24.07338
TotalCharge	-1.854153	48.252981	50.801088	6.605548	1.330958	-1.51388
Additional_charges	0.505671	51.079931	-48.241300	2.715812	-0.857136	2.31586
Survey_Timely_Admission	45.453775	-2.177145	-0.479634	29.500934	-0.139541	0.42481
Survey_Timely_Treatment	42.822559	-2.234316	-0.597632	29.182524	-1.840905	1.04289
Survey_Timely_Visits	39.516046	-2.408963	-0.202250	29.369536	0.238986	0.01781
Survey_Reliability	15.205224	4.398749	2.833105	-55.460947	1.262114	3.34240
Survey_Options	-18.994090	-5.919344	-2.362698	57.951343	0.370184	-0.36436
Survey_Hours_of_Treatment	41.003413	1.808836	2.106847	-16.146311	1.997520	0.88180
Survey_Courteous_Staff	35.636388	3.445278	1.404365	-16.893725	2.500745	0.73984
Survey_Evidence_of_Active_Listening	31.204330	1.960899	2.029133	-16.537068	-1.907401	-7.31687

31 rows × 31 columns

Part III: Data Cleaing (Continued)

E. Principle Component Analysis

E1., E2., E3.

I have printed out the results of PCA in the above code. I will choose to maintain 20 of the principle components as to explain 80% of the total variance in the data moving forward. You can tell this by looking at their cumulative explained variance in combination with the above scree plot. This allows me to have reasonable dimensionality reduction while maintaining the majority of what is being communicated in the data set. PC1 explains about 10% of the total variance and it's variables with highest correlation are:

Timely_Admission Timely_Treatment Hours_of_Treatment Timely_Visits

All of these variables are highly positively correlated to the variance explained in PC1.

So, how can an organization benefit from this analysis? Ultimately - these would be identified as the areas that they should focus research on. Ideally - they would want to reach out to the patients that answered to the extremes within these survey questions. Theoretically - they could define additional areas of concerns and provide that as a part of the dataset for deeper evaluation. Really what this is going to do is give them areas to focus in on for further evaluation. You can see the relationship of the first two principle components to our target in the scatter plot found above.

Works Cited

Chantal D. Larose, and Daniel T. Larose. Data Science Using Python and R. Wiley, 2019. EBSCOhost, search.ebscohost.com/login.aspx?direct=true&AuthType=sso&db=nlebk&AN=2091371&site=eds-live&scope=site.

Scikit-learn: Machine Learning in Python, Pedregosa et al., JMLR 12, pp. 2825-2830, 2011.

Eastwood, Brian. "The 10 Most Popular Programming Languages to Learn in 2021." Northeastern University Graduate Programs, 18 June 2020, <https://www.northeastern.edu/graduate/blog/most-popular-programming-languages/>.

Frost, Jim. "Guidelines for Removing and Handling Outliers in Data." Statistics By Jim, 5 Apr. 2021, <https://statisticsbyjim.com/basics/remove-outliers/>.

Kindsonthegenius. "Principal Components Analysis(Pca) in Python – Step by Step." Kindson The Genius, 10 Sept. 2020, <https://www.kindsonthegenius.com/principal-components-analysispca-in-python-step-by-step/>.

Bedre, Renesh. "Performing and Visualizing the Principal Component Analysis (PCA) from PCA Function and Scratch in Python." Renesh Bedre, 18 Apr. 2021, <https://www.reneshbedre.com/blog/principal-component-analysis.html>.