

# Healthcare Unsupervised Anomaly Detection

By ⇒ Pushpendra Kumar

## Overview

Healthcare fraud diverts crucial funding from vital areas such as medication and emergency services, driven by both practitioners and patients. To combat this issue, I have initiated a project focused on healthcare data anomaly detection. Through various unsupervised techniques, we aim to uncover irregularities indicating health concerns or fraudulent behavior.

Extensive Exploratory Data Analysis (EDA) helps reveal hidden anomalies and their implications. By leveraging advanced data visualization and analytics, we can extract insights for early detection and prevention. Our ultimate goal is to develop effective anomaly detection models, which are crucial for safeguarding patient health and optimizing healthcare operations.

## Importing Libraries

```
In [3]: import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
%matplotlib inline

#filtering the warnings
import warnings
warnings.filterwarnings("ignore")
```

## Loading the dataset

```
In [25]: df = pd.read_csv('Healthcare Providers.csv')
```

```
In [26]: # Shape of the data
df.shape
```

```
Out[26]: (100000, 27)
```

```
In [27]: # First 5 rows
df.head()
```

	index	National Provider Identifier	Name/Organization Name of the Provider	Last Name of the Provider	First Name of the Provider	Middle Initial of the Provider	Credentials of the Provider	Gender of the Provider	Entity Type of the Provider	Street Address 1 of the Provider	Street Address 2 of the Provider	...	HCPCS Code
0	8774979	1891106191	UPADHYAYULA	SATYASREE	NaN		M.D.	F	I	1402 S GRAND BLVD	FDT 14TH FLOOR	...	99223
1	3354385	1346202256		JONES	WENDY	P	M.D.	F	I	2950 VILLAGE DR	NaN	...	G0202
2	3001884	1306820956		DUROCHER	RICHARD	W	DPM	M	I	WASHINGTON AVE	STE 212	...	99348
3	7594822	1770523540		FULLARD	JASPER	NaN	MD	M	I	5746 N BROADWAY ST	NaN	...	81002
4	746159	1073627758		PERROTTI	ANTHONY	E	DO	M	I	875 MILITARY TRL	SUITE 200	...	96372

5 rows × 27 columns

```
In [28]: # Last 5 rows
df.tail()
```

Out[28]:

	index	National Provider Identifier	Last Name/Organization Name of the Provider	First Name of the Provider	Middle Initial of the Provider	Credentials of the Provider	Gender of the Provider	Entity Type of the Provider	Street Address 1 of the Provider	Street Address 2 of the Provider	...	HCPSCS Code
99995	3837311	1386938868	PAPES	JOAN	NaN	PT	F	I	324 E BALTIMORE ST	NaN	...	97162
99996	2079360	1215091327	HAYNER MARGARET	S	ARNP	F	I	645 NW 4TH ST	NaN	...	99213	
99997	8927965	1902868185	VALENCIA	DANA	NaN	M.D.	M	I	3009 N BALLAS RD	SUITE 202B	...	93320 s t
99998	8854571	1891941183	GONZALEZ-LAMOS	RAFAELA	NaN	NaN	F	I	2365 BOSTON POST RD	SUITE 201	...	G0008 A ,
99999	3547535	1356772156	RAMEZANI	ELIIAN	NaN	NaN	F	I	444 MIDDLE NECK RD APT 1H	NaN	...	97112

5 rows × 27 columns

In [29]:

# Information regarding dataset  
df.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 100000 entries, 0 to 99999
Data columns (total 27 columns):
 #   Column           Non-Null Count  Dtype  
 ---  --  
 0   index            100000 non-null   int64  
 1   National Provider Identifier 100000 non-null   int64  
 2   Last Name/Organization Name of the Provider 100000 non-null   object 
 3   First Name of the Provider    95745 non-null   object 
 4   Middle Initial of the Provider 70669 non-null   object 
 5   Credentials of the Provider  92791 non-null   object 
 6   Gender of the Provider      95746 non-null   object 
 7   Entity Type of the Provider 100000 non-null   object 
 8   Street Address 1 of the Provider 100000 non-null   object 
 9   Street Address 2 of the Provider 40637 non-null   object 
 10  City of the Provider       100000 non-null   object 
 11  Zip Code of the Provider   100000 non-null   float64 
 12  State Code of the Provider 100000 non-null   object 
 13  Country Code of the Provider 100000 non-null   object 
 14  Provider Type            100000 non-null   object 
 15  Medicare Participation Indicator 100000 non-null   object 
 16  Place of Service          100000 non-null   object 
 17  HCPSCS Code              100000 non-null   object 
 18  HCPSCS Description        100000 non-null   object 
 19  HCPSCS Drug Indicator    100000 non-null   object 
 20  Number of Services        100000 non-null   object 
 21  Number of Medicare Beneficiaries 100000 non-null   object 
 22  Number of Distinct Medicare Beneficiary/Per Day Services 100000 non-null   object 
 23  Average Medicare Allowed Amount 100000 non-null   object 
 24  Average Submitted Charge Amount 100000 non-null   object 
 25  Average Medicare Payment Amount 100000 non-null   object 
 26  Average Medicare Standardized Amount 100000 non-null   object 
dtypes: float64(1), int64(2), object(24)
memory usage: 20.6+ MB
```

In [30]:

# Columns of the dataset  
df.columns

Out[30]:

```
Index(['index', 'National Provider Identifier',
       'Last Name/Organization Name of the Provider',
       'First Name of the Provider', 'Middle Initial of the Provider',
       'Credentials of the Provider', 'Gender of the Provider',
       'Entity Type of the Provider', 'Street Address 1 of the Provider',
       'Street Address 2 of the Provider', 'City of the Provider',
       'Zip Code of the Provider', 'State Code of the Provider',
       'Country Code of the Provider', 'Provider Type',
       'Medicare Participation Indicator', 'Place of Service', 'HCPSCS Code',
       'HCPSCS Description', 'HCPSCS Drug Indicator', 'Number of Services',
       'Number of Medicare Beneficiaries',
       'Number of Distinct Medicare Beneficiary/Per Day Services',
       'Average Medicare Allowed Amount', 'Average Submitted Charge Amount',
       'Average Medicare Payment Amount',
       'Average Medicare Standardized Amount'],
      dtype='object')
```

4. Continuation Description of dataset

```
In [31]: # Statistical description of dataset  
df.describe()
```

```
Out[31]:
```

	index	National Provider Identifier	Zip Code of the Provider
count	1.000000e+05	1.000000e+05	1.000000e+05
mean	4.907646e+06	1.498227e+09	4.163820e+08
std	2.839633e+06	2.874125e+08	3.082566e+08
min	2.090000e+02	1.003001e+09	6.010000e+02
25%	2.458791e+06	1.245669e+09	1.426300e+08
50%	4.901266e+06	1.497847e+09	3.633025e+08
75%	7.349450e+06	1.740374e+09	6.819881e+08
max	9.847440e+06	1.993000e+09	9.990166e+08

## Exploratory Data Analysis (EDA)

### Handling Missing Values

```
In [32]: # Checking missing values  
df.isnull().sum()
```

```
Out[32]:
```

index	0
National Provider Identifier	0
Last Name/Organization Name of the Provider	0
First Name of the Provider	4255
Middle Initial of the Provider	29331
Credentials of the Provider	7209
Gender of the Provider	4254
Entity Type of the Provider	0
Street Address 1 of the Provider	0
Street Address 2 of the Provider	59363
City of the Provider	0
Zip Code of the Provider	0
State Code of the Provider	0
Country Code of the Provider	0
Provider Type	0
Medicare Participation Indicator	0
Place of Service	0
HPCPS Code	0
HPCPS Description	0
HPCPS Drug Indicator	0
Number of Services	0
Number of Medicare Beneficiaries	0
Number of Distinct Medicare Beneficiary/Per Day Services	0
Average Medicare Allowed Amount	0
Average Submitted Charge Amount	0
Average Medicare Payment Amount	0
Average Medicare Standardized Amount	0
dtype: int64	

**Insights:** Dropping the columns:

'index', 'National Provider Identifier', 'Last Name/Organization Name of the Provider', 'First Name of the Provider', 'Middle Initial of the Provider', 'Street Address 1 of the Provider', 'Street Address 2 of the Provider'. As these doesn't add much towards detecting anomaly. Meanwhile filling rest missing columnns: ('Credentials of the Provider','Gender of the Provider') with their mode values.

```
In [33]: # Dropping columns  
DropColumns = ['index', 'National Provider Identifier', 'Last Name/Organization Name of the Provider',  
              'First Name of the Provider', 'Middle Initial of the Provider', 'Street Address 1 of the Provider',  
              'Street Address 2 of the Provider']  
df = df.drop(DropColumns, axis=1)
```

```
In [34]: df.info()
```

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 100000 entries, 0 to 99999
Data columns (total 20 columns):
 #   Column           Non-Null Count   Dtype  
--- 
 0   Credentials of the Provider    92791 non-null    object  
 1   Gender of the Provider        95746 non-null    object  
 2   Entity Type of the Provider  100000 non-null    object  
 3   City of the Provider         100000 non-null    object  
 4   Zip Code of the Provider     100000 non-null    float64 
 5   State Code of the Provider   100000 non-null    object  
 6   Country Code of the Provider 100000 non-null    object  
 7   Provider Type               100000 non-null    object  
 8   Medicare Participation Indicator 100000 non-null    object  
 9   Place of Service             100000 non-null    object  
 10  HCPCS Code                  100000 non-null    object  
 11  HCPCS Description           100000 non-null    object  
 12  HCPCS Drug Indicator        100000 non-null    object  
 13  Number of Services          100000 non-null    object  
 14  Number of Medicare Beneficiaries 100000 non-null    object  
 15  Number of Distinct Medicare Beneficiary/Per Day Services 100000 non-null    object  
 16  Average Medicare Allowed Amount 100000 non-null    object  
 17  Average Submitted Charge Amount 100000 non-null    object  
 18  Average Medicare Payment Amount 100000 non-null    object  
 19  Average Medicare Standardized Amount 100000 non-null    object  
dtypes: float64(1), object(19)
memory usage: 15.3+ MB

```

```
In [35]: # Filling missing with mode values
df["Credentials of the Provider"] = df["Credentials of the Provider"].fillna(df["Credentials of the Provider"].mode()[0])
df["Gender of the Provider"] = df["Gender of the Provider"].fillna(df["Gender of the Provider"].mode()[0])
```

```
In [36]: # Checking missing values again
df.isnull().sum()
```

```
Out[36]: Credentials of the Provider      0
Gender of the Provider                   0
Entity Type of the Provider             0
City of the Provider                    0
Zip Code of the Provider                0
State Code of the Provider              0
Country Code of the Provider            0
Provider Type                          0
Medicare Participation Indicator      0
Place of Service                        0
HCPCS Code                            0
HCPCS Description                      0
HCPCS Drug Indicator                  0
Number of Services                     0
Number of Medicare Beneficiaries      0
Number of Distinct Medicare Beneficiary/Per Day Services 0
Average Medicare Allowed Amount        0
Average Submitted Charge Amount        0
Average Medicare Payment Amount        0
Average Medicare Standardized Amount  0
dtype: int64
```

**Insights:** Converting these object columns which are meant to be numeric:

'Number of Services', 'Number of Medicare Beneficiaries', 'Number of Distinct Medicare Beneficiary/Per Day Services', 'Average Medicare Allowed Amount', 'Average Submitted Charge Amount', 'Average Medicare Payment Amount', 'Average Medicare Standardized Amount'.

Also removing any comma separators if they are present.

```
In [37]: # Data Cleaning

def RemoveComma(x):
    return str(x).replace(",","",) # Convert to string before replacing comma

numericCols = ['Number of Services', 'Number of Medicare Beneficiaries',
               'Number of Distinct Medicare Beneficiary/Per Day Services',
               'Average Medicare Allowed Amount', 'Average Submitted Charge Amount',
               'Average Medicare Payment Amount', 'Average Medicare Standardized Amount']

df[numericCols] = df[numericCols].applymap(RemoveComma).apply(pd.to_numeric, errors='ignore')
```

```
In [38]: #Checking information again
df.info()
```

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 100000 entries, 0 to 99999
Data columns (total 20 columns):
 #   Column          Non-Null Count   Dtype  
 --- 
 0   Credentials of the Provider    100000 non-null   object  
 1   Gender of the Provider        100000 non-null   object  
 2   Entity Type of the Provider  100000 non-null   object  
 3   City of the Provider         100000 non-null   object  
 4   Zip Code of the Provider     100000 non-null   float64 
 5   State Code of the Provider   100000 non-null   object  
 6   Country Code of the Provider 100000 non-null   object  
 7   Provider Type               100000 non-null   object  
 8   Medicare Participation Indicator 100000 non-null   object  
 9   Place of Service            100000 non-null   object  
 10  HCPCS Code                 100000 non-null   object  
 11  HCPCS Description          100000 non-null   object  
 12  HCPCS Drug Indicator       100000 non-null   object  
 13  Number of Services         100000 non-null   float64 
 14  Number of Medicare Beneficiaries 100000 non-null   int64  
 15  Number of Distinct Medicare Beneficiary/Per Day Services 100000 non-null   int64  
 16  Average Medicare Allowed Amount 100000 non-null   float64 
 17  Average Submitted Charge Amount 100000 non-null   float64 
 18  Average Medicare Payment Amount 100000 non-null   float64 
 19  Average Medicare Standardized Amount 100000 non-null   float64 

dtypes: float64(6), int64(2), object(12)
memory usage: 15.3+ MB

```

## Checking Duplicate Values

```
In [41]: # Checking for duplicates
df.duplicated().any()
```

```
Out[41]: True
```

```
In [43]: # Finding duplicate rows in the entire DataFrame
duplicate_rows = df[df.duplicated()]
duplicate_rows
```

	Credentials of the Provider	Gender of the Provider	Entity Type of the Provider	City of the Provider	Zip Code of the Provider	State Code of the Provider	Country Code of the Provider	Provider Type	Medicare Participation Indicator	Place of Service	HCPCS Code	HCPCS Description	HC Indic
2047	MD	M	I	TAMPA	336125513.0	FL	US	Diagnostic Radiology	Y	O	84520	Urea nitrogen level to assess kidney function	

```
In [44]: # Dropping duplicates
df = df.drop_duplicates()
```

```
In [45]: #Checking duplicates again
df.duplicated().any()
```

```
Out[45]: False
```

```
In [46]: df.head()
```

Out[46]:

	Credentials of the Provider	Gender of the Provider	Entity Type of the Provider	City of the Provider	Zip Code of the Provider	State Code of the Provider	Country Code of the Provider	Provider Type	Medicare Participation Indicator	Place of Service	HCPCS Code	HCPD Descripti
0	M.D.	F	I	SAINT LOUIS	631041004.0	MO	US	Internal Medicine	Y	F	99223	Initial hospital inpatient care typically 70
1	M.D.	F	I	FAYETTEVILLE	283043815.0	NC	US	Obstetrics & Gynecology	Y	O	G0202	Screen mammogram bilateral view study
2	DPM	M	I	NORTH HAVEN	64732343.0	CT	US	Podiatry	Y	O	99348	Establish patient history visit, typical 25 min
3	MD	M	I	KANSAS CITY	641183998.0	MO	US	Internal Medicine	Y	O	81002	Urinalysis manual test
4	DO	M	I	JUPITER	334585700.0	FL	US	Internal Medicine	Y	O	96372	Injection beneath tissue or muscle for

In [47]:

# Cleaning 'Credentials of the Provider' Column

# Removing periods '.' if present

df['Credentials of the Provider'] = df['Credentials of the Provider'].str.replace('.','')

In [48]:

df.head()

Out[48]:

	Credentials of the Provider	Gender of the Provider	Entity Type of the Provider	City of the Provider	Zip Code of the Provider	State Code of the Provider	Country Code of the Provider	Provider Type	Medicare Participation Indicator	Place of Service	HCPCS Code	HCPD Descripti
0	MD	F	I	SAINT LOUIS	631041004.0	MO	US	Internal Medicine	Y	F	99223	Initial hospital inpatient care typically 70
1	MD	F	I	FAYETTEVILLE	283043815.0	NC	US	Obstetrics & Gynecology	Y	O	G0202	Screen mammogram bilateral view study
2	DPM	M	I	NORTH HAVEN	64732343.0	CT	US	Podiatry	Y	O	99348	Establish patient history visit, typical 25 min
3	MD	M	I	KANSAS CITY	641183998.0	MO	US	Internal Medicine	Y	O	81002	Urinalysis manual test
4	DO	M	I	JUPITER	334585700.0	FL	US	Internal Medicine	Y	O	96372	Injection beneath tissue or muscle for

## Data Visualization

In [67]:

df.describe()

Out[67]:

	Zip Code of the Provider	Number of Services	Number of Medicare Beneficiaries	Number of Distinct Medicare Beneficiary/Per Day Services	Average Medicare Allowed Amount	Average Submitted Charge Amount	Average Medicare Payment Amount	Average Medicare Standardized Amount
count	9.999900e+04	99999.000000	99999.000000	99999.000000	99999.000000	99999.000000	99999.000000	99999.000000
mean	4.163828e+08	239.673621	89.810008	142.116901	101.435115	354.553846	77.359516	78.031420
std	3.082580e+08	2493.199458	1109.622428	1640.235384	257.243886	1062.613042	199.719682	200.046326
min	6.010000e+02	11.000000	11.000000	11.000000	0.010000	0.010000	0.008679	0.008679
25%	1.426300e+08	21.000000	17.000000	20.000000	24.270000	57.649074	19.335749	20.122753
50%	3.633025e+08	43.000000	32.000000	40.000000	65.100000	146.000000	47.020353	47.841582
75%	6.819881e+08	118.000000	75.000000	106.000000	113.160000	298.936222	84.895746	84.880047
max	9.990166e+08	282739.000000	190306.000000	282737.000000	20494.000000	62694.000000	16067.300000	16957.148000

By observing the standard deviation, it can be concluded that the data is widely spreaded.

## Correlation Matrix

```
In [50]: df_corr = df.select_dtypes(include=['number']).corr()
df_corr
```

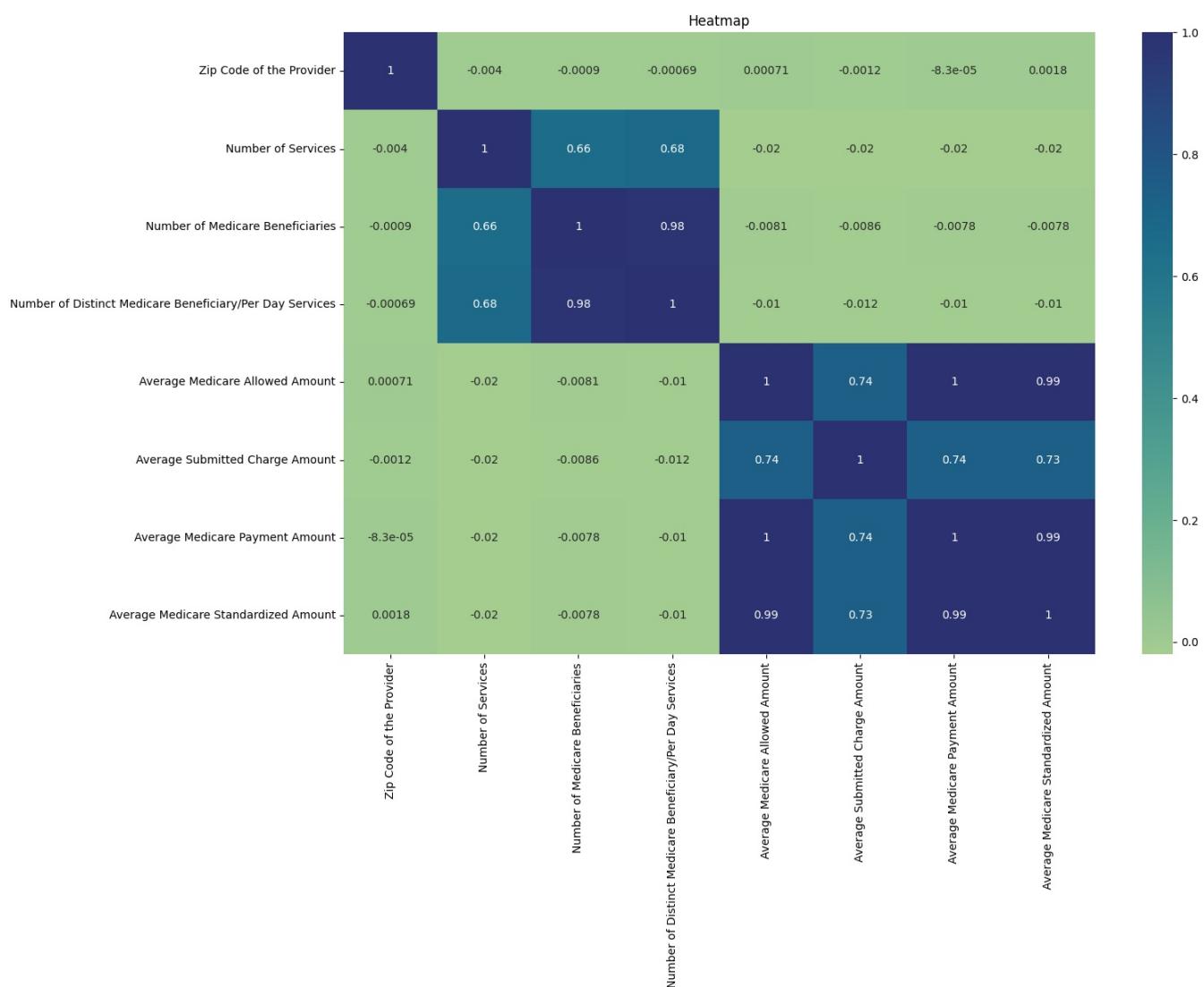
Out[50]:

	Zip Code of the Provider	Number of Services	Number of Medicare Beneficiaries	Number of Distinct Medicare Beneficiary/Per Day Services	Average Medicare Allowed Amount	Average Submitted Charge Amount	Average Medicare Payment Amount	Average Medicare Standardized Amount
Zip Code of the Provider	1.000000	-0.003997	-0.000900	-0.000686	0.000706	-0.001193	-0.000083	0.001842
Number of Services	-0.003997	1.000000	0.658923	0.675076	-0.020413	-0.019698	-0.019995	-0.020093
Number of Medicare Beneficiaries	-0.000900	0.658923	1.000000	0.981072	-0.008088	-0.008557	-0.007797	-0.007815
Number of Distinct Medicare Beneficiary/Per Day Services	-0.000686	0.675076	0.981072	1.000000	-0.010471	-0.011829	-0.010246	-0.010265
Average Medicare Allowed Amount	0.000706	-0.020413	-0.008088	-0.010471	1.000000	0.739766	0.998704	0.994831
Average Submitted Charge Amount	-0.001193	-0.019698	-0.008557	-0.011829	0.739766	1.000000	0.738814	0.733502
Average Medicare Payment Amount	-0.000083	-0.019995	-0.007797	-0.010246	0.998704	0.738814	1.000000	0.994649
Average Medicare Standardized Amount	0.001842	-0.020093	-0.007815	-0.010265	0.994831	0.733502	0.994649	1.000000

## Heatmap

```
In [80]: plt.figure(figsize=(15,10))
sns.heatmap(df_corr, annot=True, cmap='crest')
plt.title('Heatmap')
```

Out[80]:

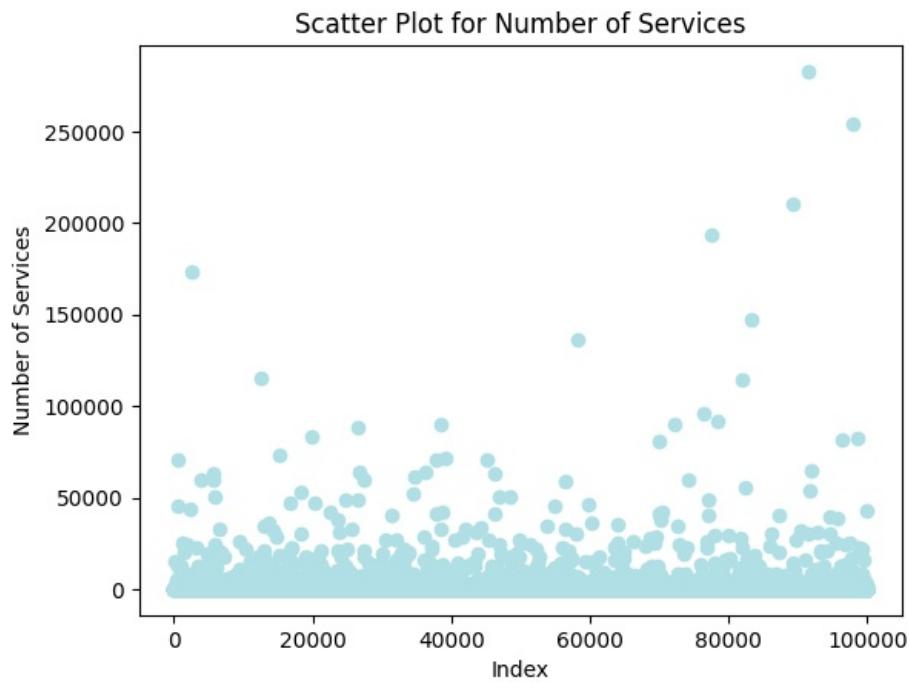
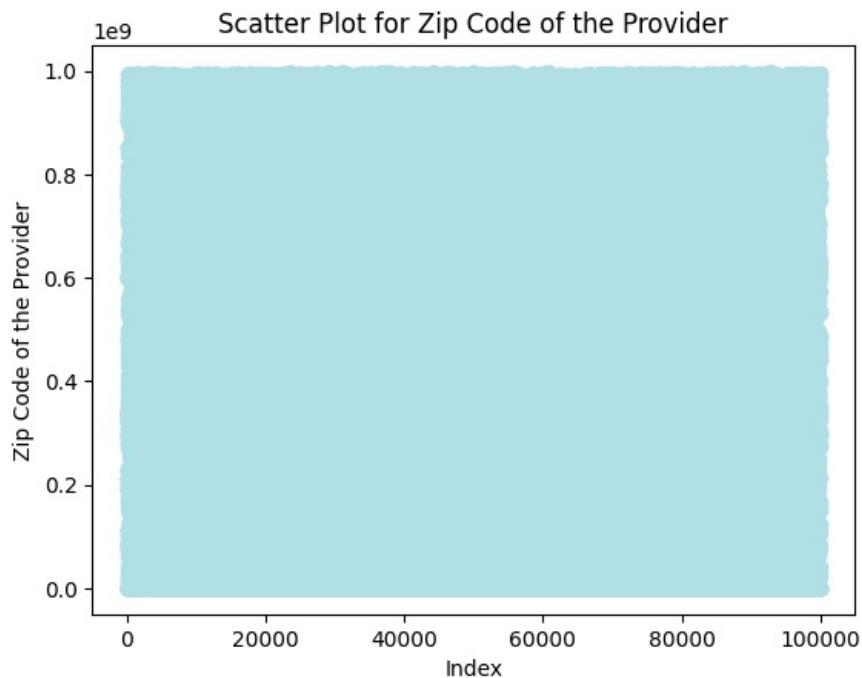


Insights: The 'Average Medicare Allowed Amount' and 'Average Medicare Payment Amount' shows a very strong correlation having value=1, suggesting that higher allowed amounts tend to result in higher payments.

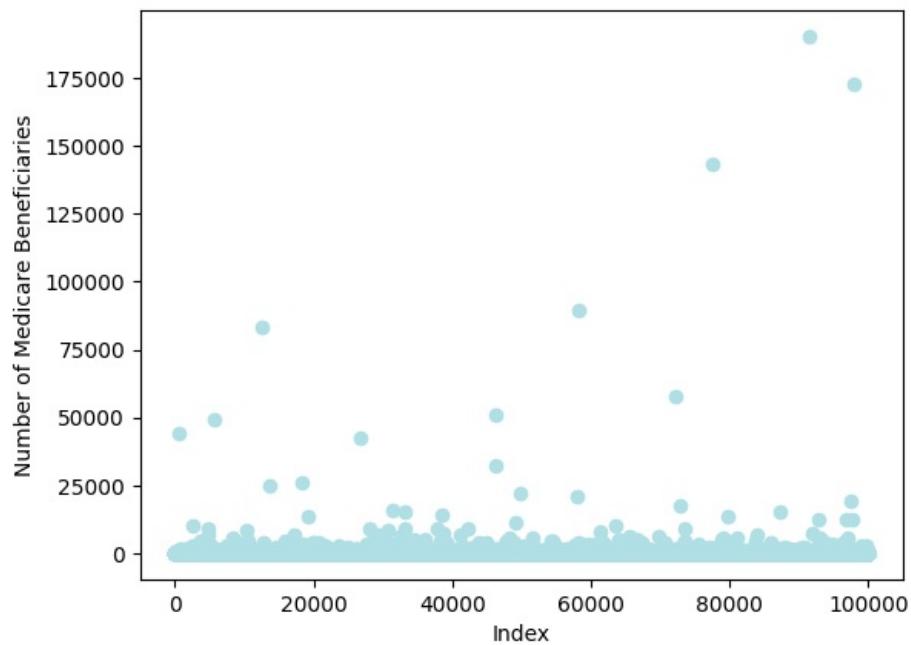
```
In [87]: numerical_cols = df.select_dtypes(include=[np.number])

for col in numerical_cols.columns:
    plt.figure()
    plt.scatter(df.index, df[col], c='powderblue')
    plt.title(f"Scatter Plot for {col}")
    plt.xlabel("Index")
    plt.ylabel(col)

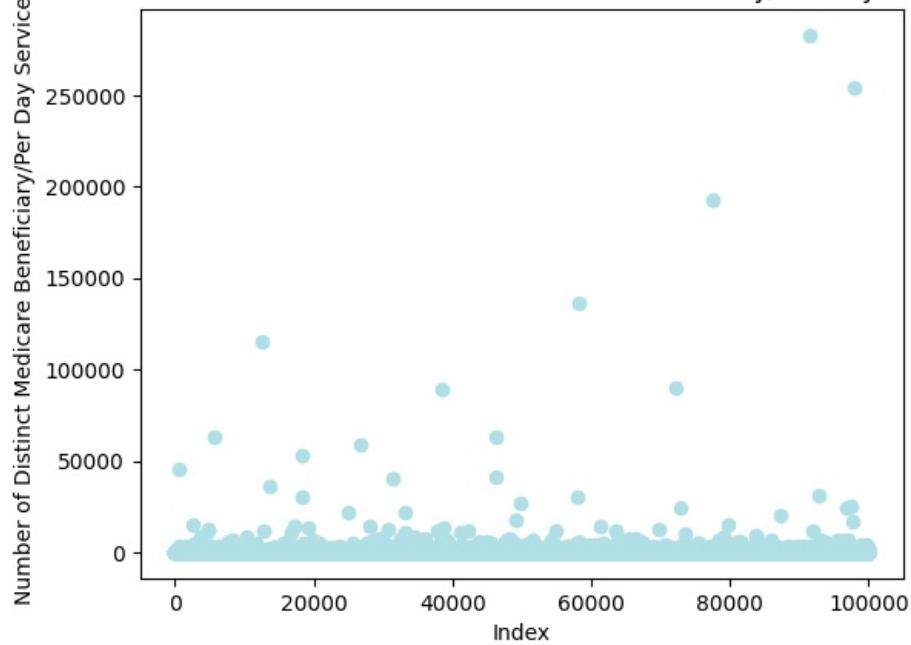
plt.show()
```



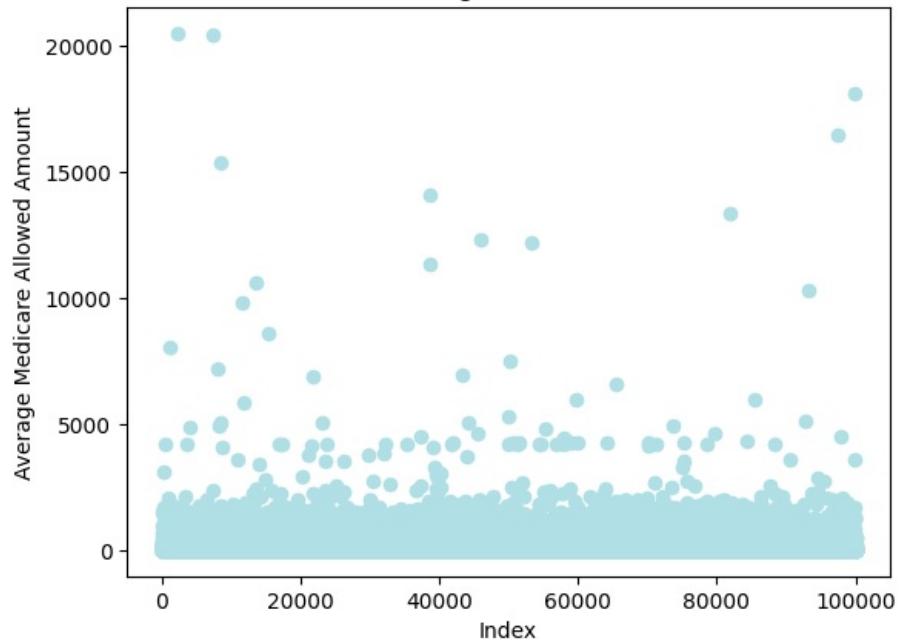
Scatter Plot for Number of Medicare Beneficiaries



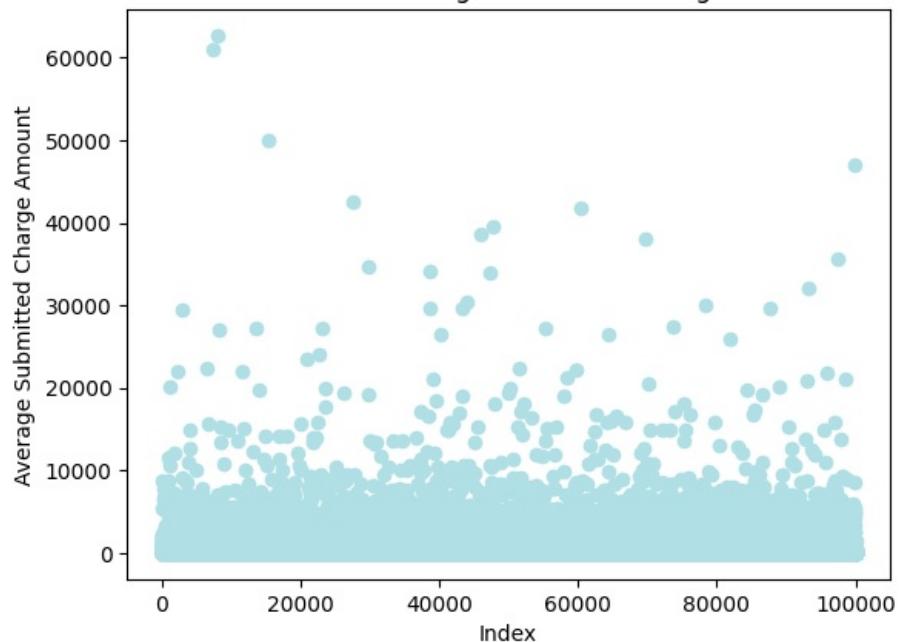
Scatter Plot for Number of Distinct Medicare Beneficiary/Per Day Services

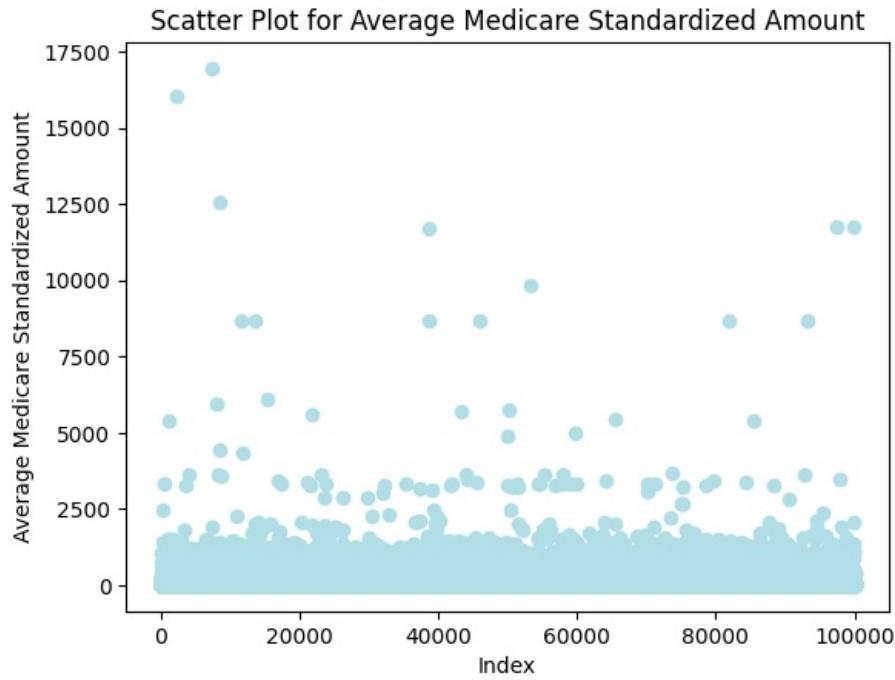
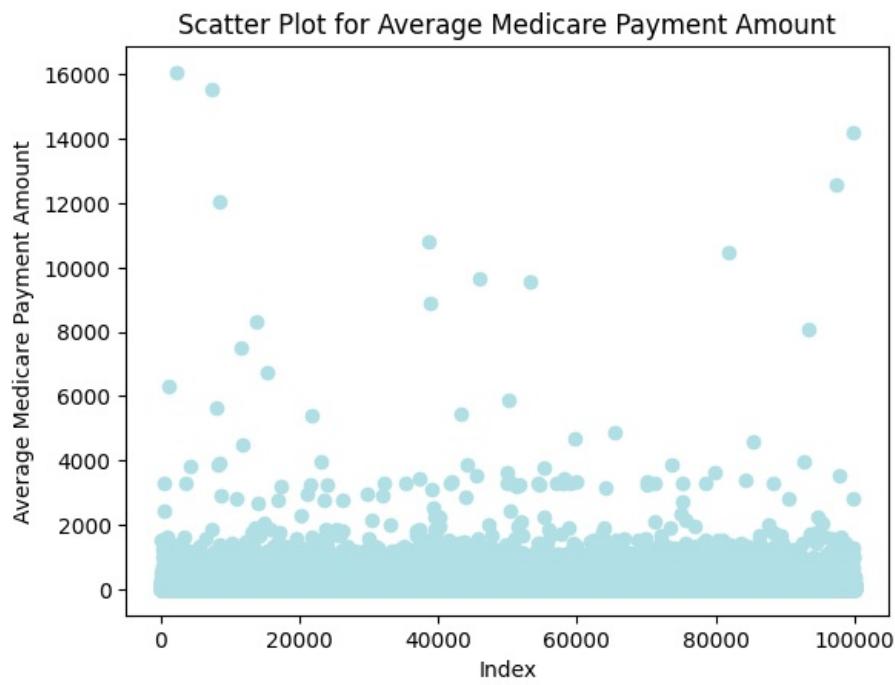


Scatter Plot for Average Medicare Allowed Amount



Scatter Plot for Average Submitted Charge Amount





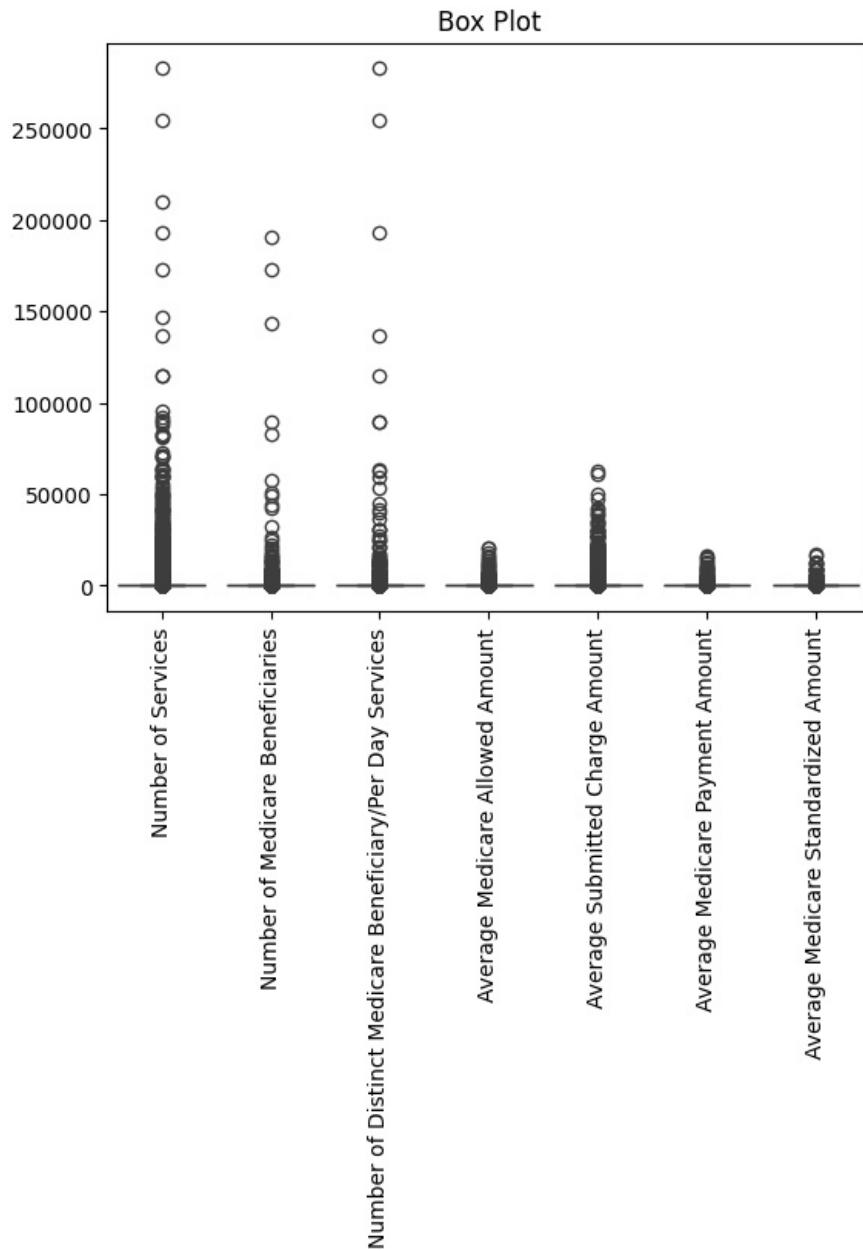
**Observation:** The 'Zip Code of the Provider' column has uniform distribution so it can be dropped while training the model for detecting outliers.

In [91]: `# Dropping Zip Code  
df.drop(columns=['Zip Code of the Provider'], inplace=True)`

Box Plot Representation

```
In [95]: sns.boxplot(data=df)
plt.title('Box Plot')
plt.xticks(rotation=90)
```

```
Out[95]: ([0, 1, 2, 3, 4, 5, 6],
[Text(0, 0, 'Number of Services'),
Text(1, 0, 'Number of Medicare Beneficiaries'),
Text(2, 0, 'Number of Distinct Medicare Beneficiary/Per Day Services'),
Text(3, 0, 'Average Medicare Allowed Amount'),
Text(4, 0, 'Average Submitted Charge Amount'),
Text(5, 0, 'Average Medicare Payment Amount'),
Text(6, 0, 'Average Medicare Standardized Amount')])
```



Insights: The columns 'Number of Services', 'Number of Medicare Beneficiaries', and 'Number of Distinct Medicare Beneficiary/Per Day Services' contains more number of outliers.

## Visualizing using AutoViz Library

### Univariate and Bivariate Analysis

```
In [97]: %matplotlib inline
from autoviz import AutoViz_Class
AV = AutoViz_Class()

dft = AV.AutoViz(df, chart_format="png")
```

Shape of your Data Set loaded: (99999, 19)

##### CLASSIFYING VARIABLES #####

Classifying variables in data set...

Number of Numeric Columns = 5  
 Number of Integer-Categorical Columns = 2  
 Number of String-Categorical Columns = 3  
 Number of Factor-Categorical Columns = 0  
 Number of String-Boolean Columns = 5  
 Number of Numeric-Boolean Columns = 0  
 Number of Discrete String Columns = 0  
 Number of NLP String Columns = 4  
 Number of Date Time Columns = 0  
 Number of ID Columns = 0  
 Number of Columns to Delete = 0  
 19 Predictors classified...

No variables removed since no ID or low-information variables found in data set

To fix these data quality issues in the dataset, import FixDQ from autoviz...

There are 1 duplicate rows in your dataset

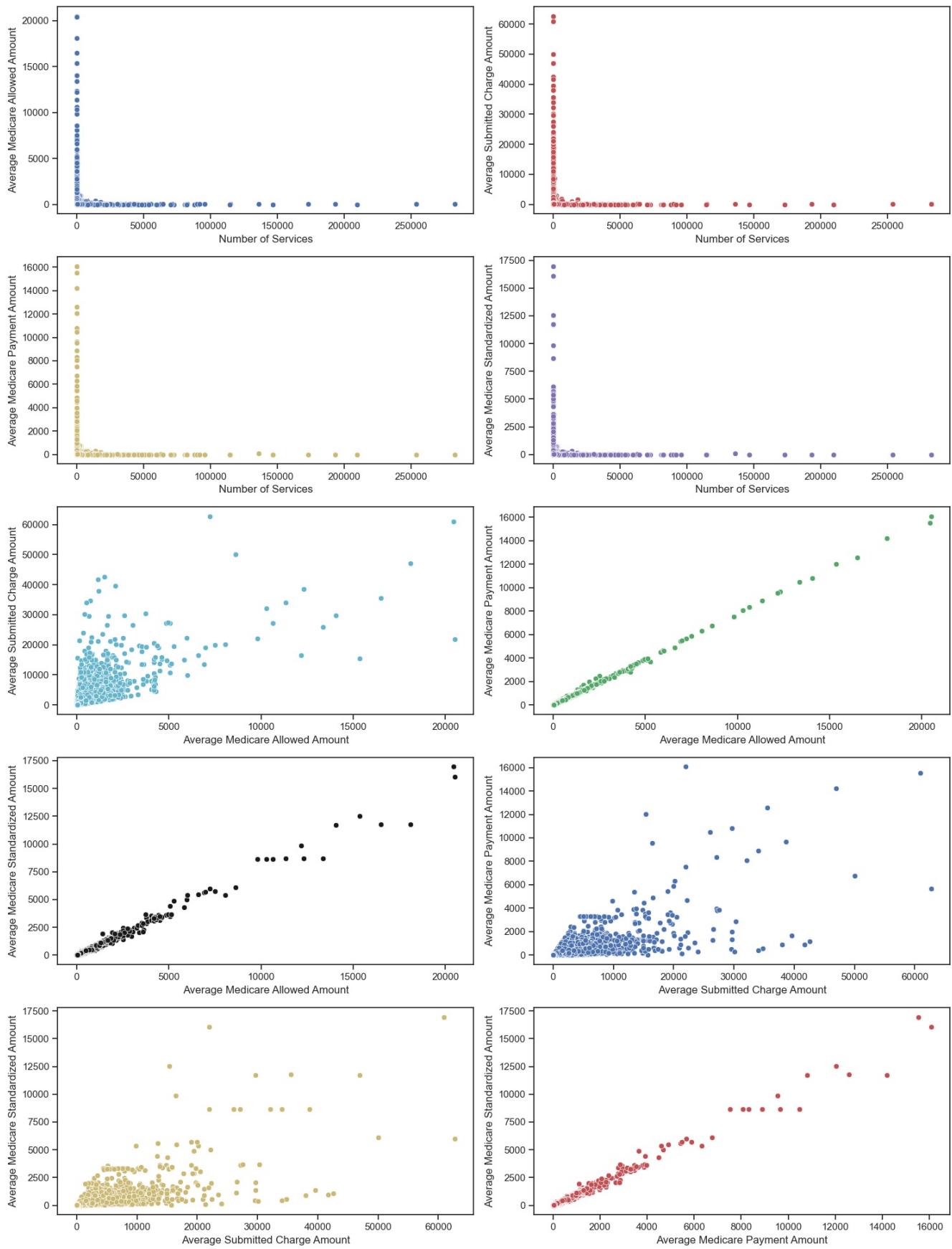
Alert: Dropping duplicate rows can sometimes cause your column data types to change to object!

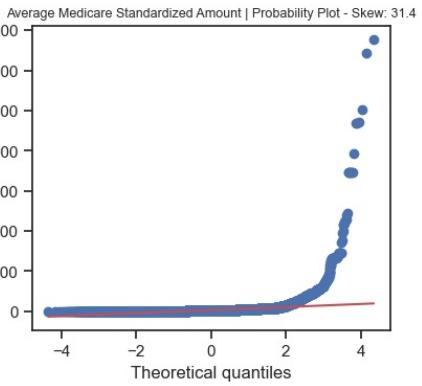
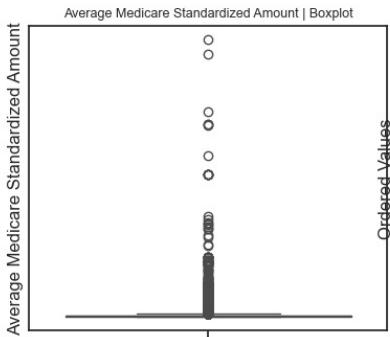
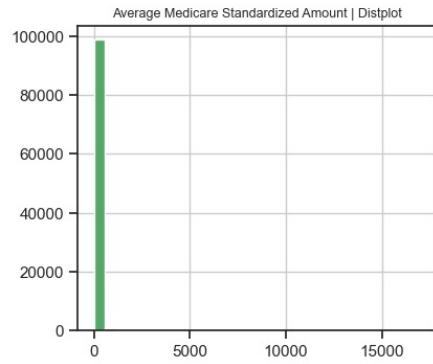
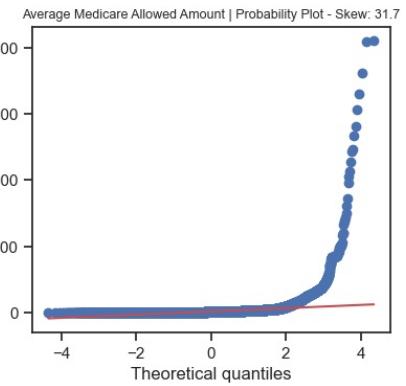
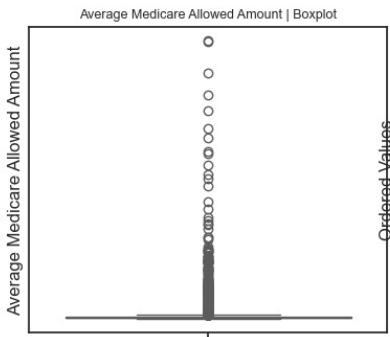
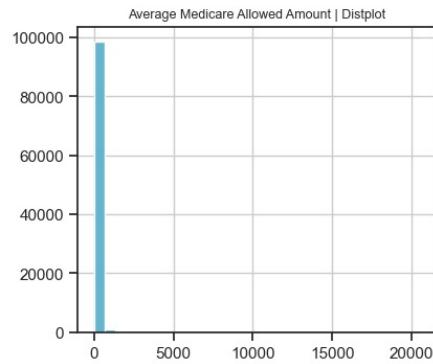
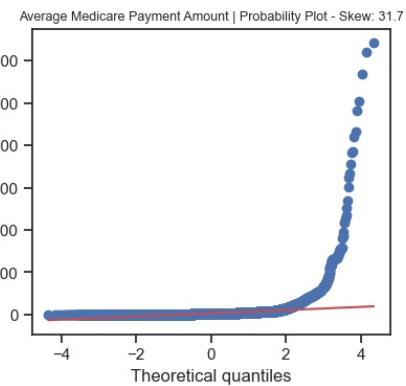
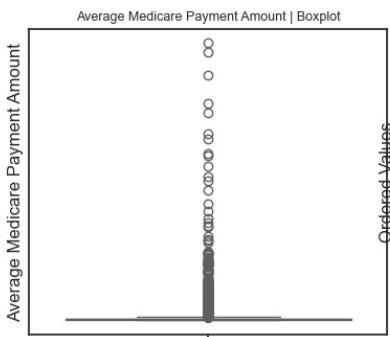
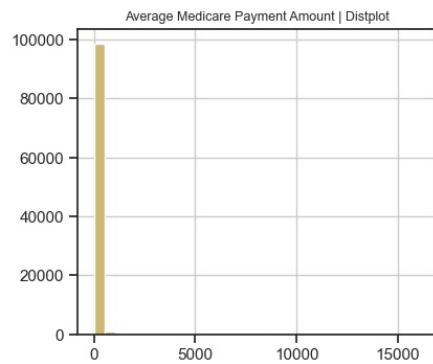
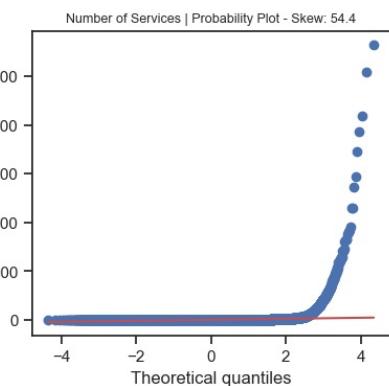
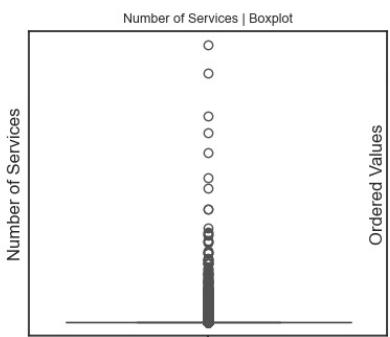
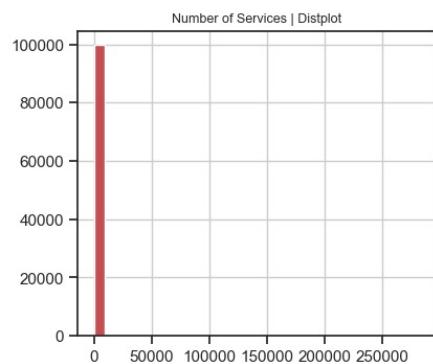
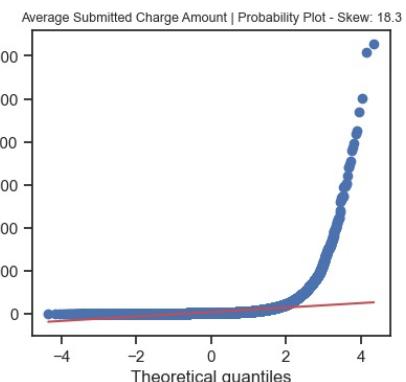
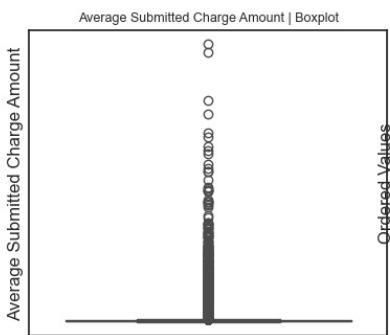
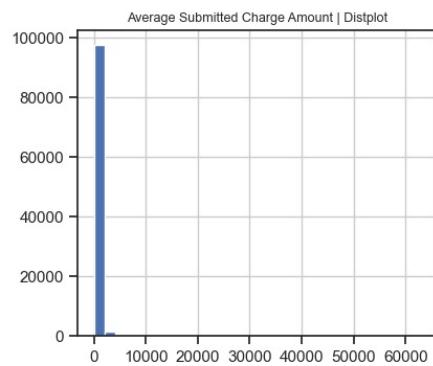
All variables classified into correct types.

	Data Type	Missing Values%	Unique Values%	Minimum Value	Maximum Value	DQ Issue
Credentials of the Provider	object	0.000000	1			No issue
Gender of the Provider	object	0.000000	0			No issue
Entity Type of the Provider	object	0.000000	0			No issue
City of the Provider	object	0.000000	5			No issue
State Code of the Provider	object	0.000000	0		25 rare categories: Too many to list. Group them into a single category or drop the categories.	
Country Code of the Provider	object	0.000000	0		3 rare categories: ['DE', 'JP', 'TR']. Group them into a single category or drop the categories.	
Provider Type	object	0.000000	0		64 rare categories: Too many to list. Group them into a single category or drop the categories.	
Medicare Participation Indicator	object	0.000000	0			No issue
Place of Service	object	0.000000	0			No issue
HCPCS Code	object	0.000000	2			No issue
HCPCS Description	object	0.000000	2			No issue
HCPCS Drug Indicator	object	0.000000	0			No issue
Number of Services	float64	0.000000	NA 11.000000 282739.000000		Column has 12316 outliers greater than upper bound (263.50) or lower than lower bound(-124.50). Cap them or remove them.	
Number of Medicare Beneficiaries	int64	0.000000	1 11.000000 190306.000000		Column has 10473 outliers greater than upper bound (162.00) or lower than lower bound(-70.00). Cap them or remove them.	
Number of Distinct Medicare Beneficiary/Per Day Services	int64	0.000000	1 11.000000 282737.000000		Column has 11895 outliers greater than upper bound (235.00) or lower than lower bound(-109.00). Cap them or remove them., Column has a high correlation with ['Number of Medicare Beneficiaries']. Consider dropping one of them.	
Average Medicare Allowed Amount	float64	0.000000	NA 0.010000 20494.000000		Column has 5079 outliers greater than upper bound (246.49) or lower than lower bound(-109.06). Cap them or remove them.	
Average Submitted Charge Amount	float64	0.000000	NA 0.010000 62694.000000		Column has 10673 outliers greater than upper bound (660.88) or lower than lower bound(-304.29). Cap them or remove them.	
Average Medicare Payment Amount	float64	0.000000	NA 0.008679 16067.300000		Column has 6093 outliers greater than upper bound (183.24) or lower than lower bound(-79.01). Cap them or remove them., Column has a high correlation with ['Average Medicare Allowed Amount']. Consider dropping one of them.	
Average Medicare Standardized Amount	float64	0.000000	NA 0.008679 16957.148000		Column has 6078 outliers greater than upper bound (182.02) or lower than lower bound(-77.01). Cap them or remove them., Column has a high correlation with ['Average Medicare Allowed Amount', 'Average Medicare Payment Amount']. Consider dropping one of them.	

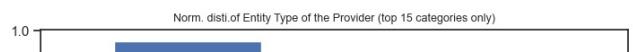
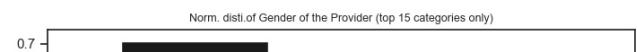
**Number of All Scatter Plots = 15**

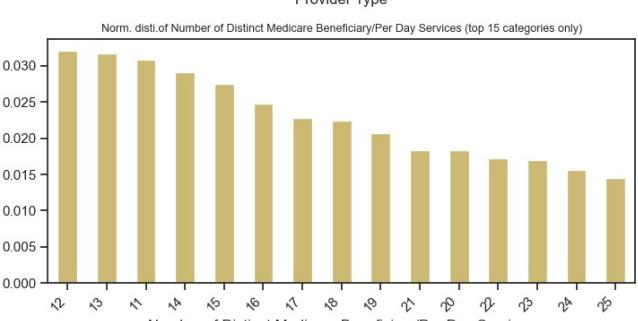
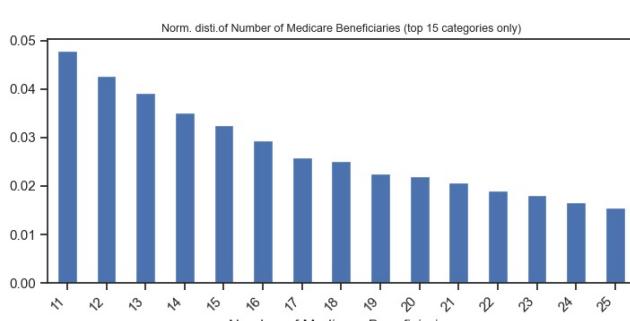
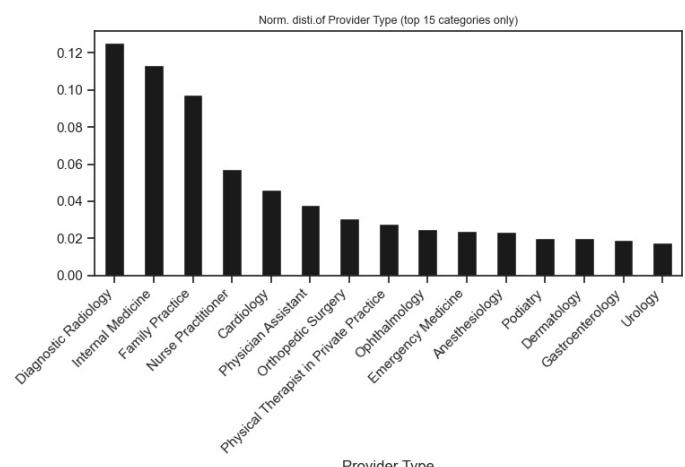
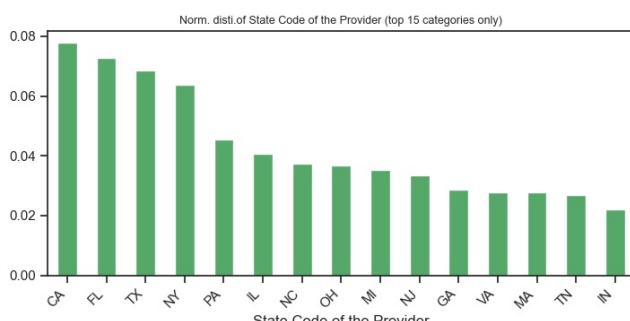
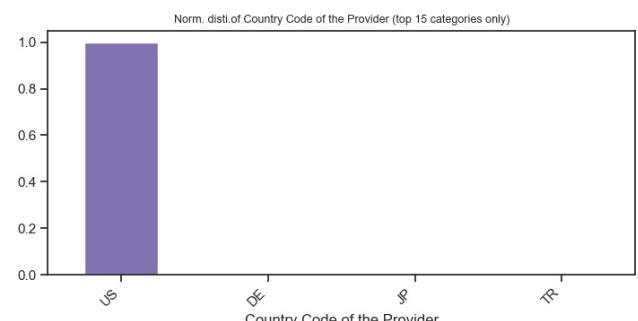
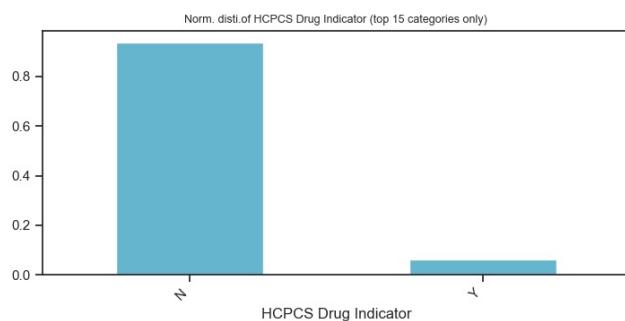
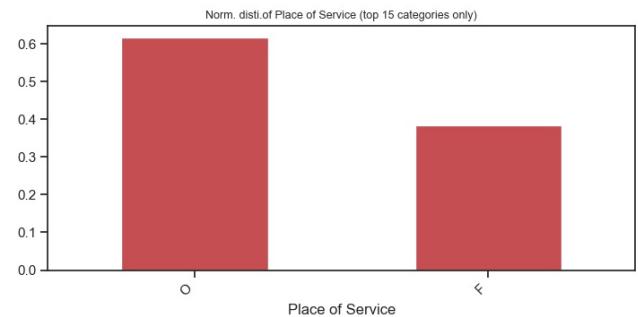
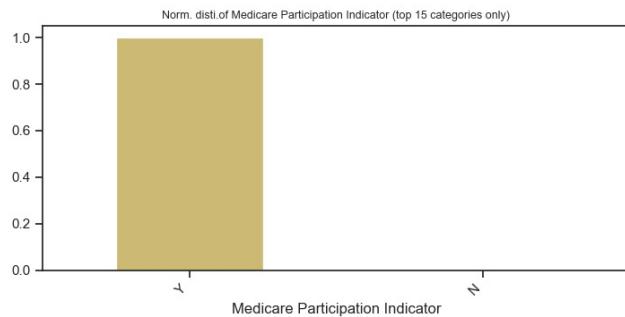
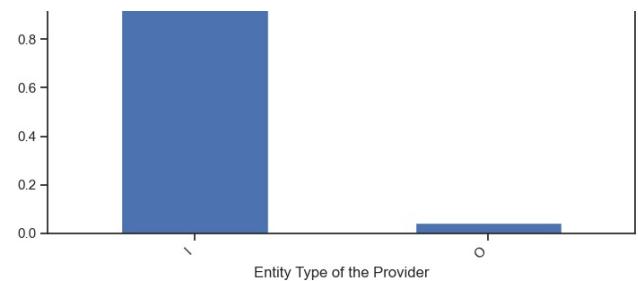
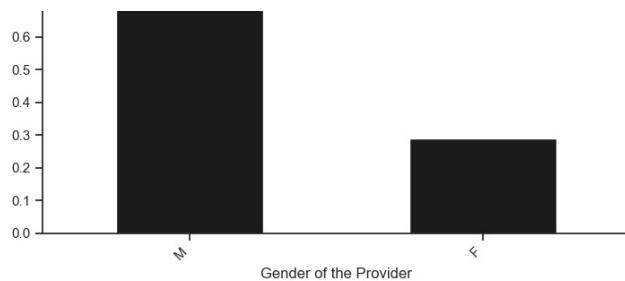
Pair-wise Scatter Plot of all Continuous Variables



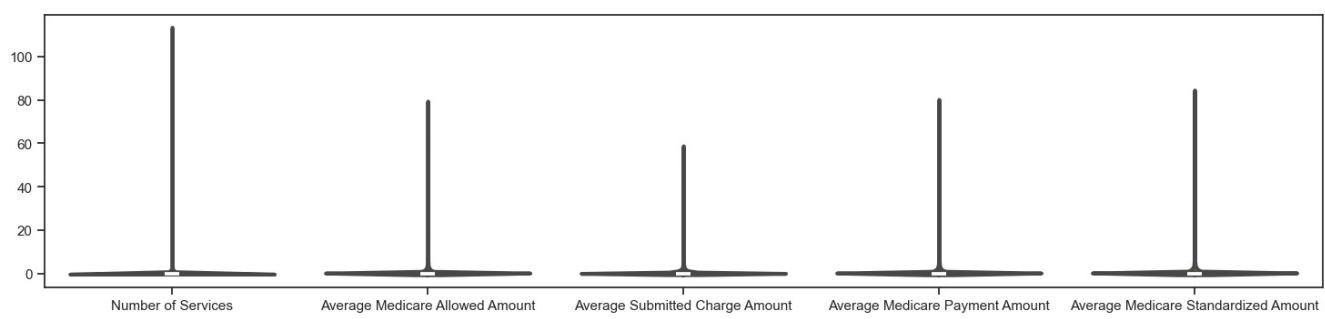


Histograms and Normalized distributions of all variables

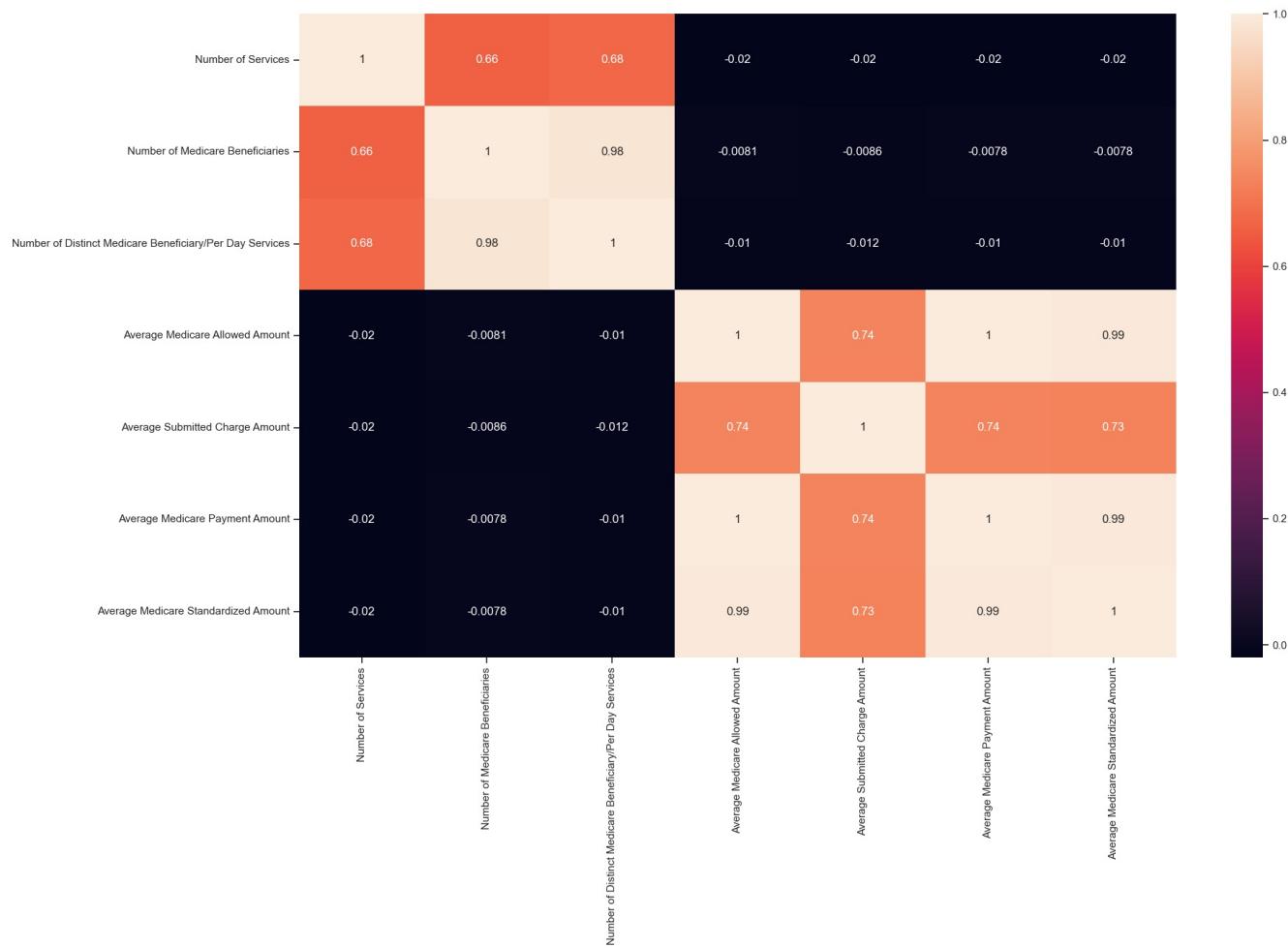




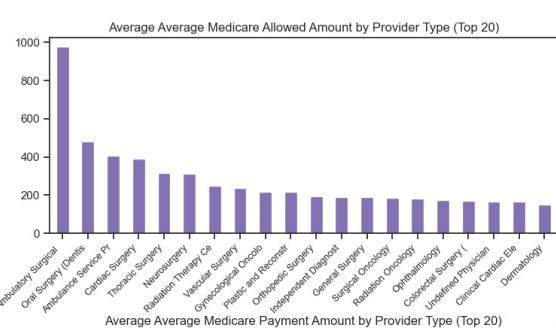
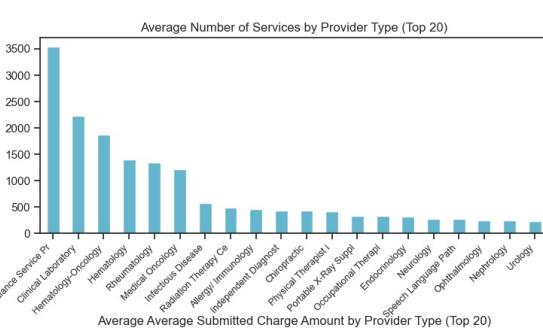
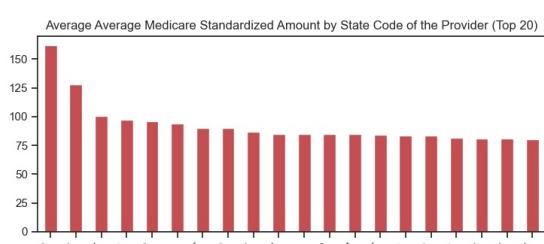
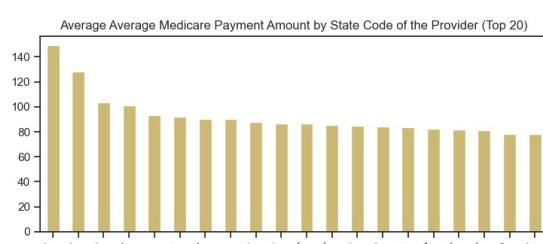
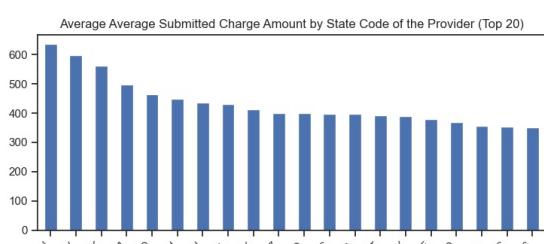
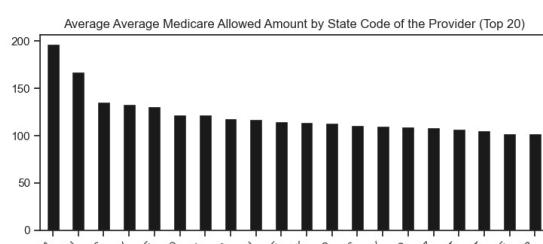
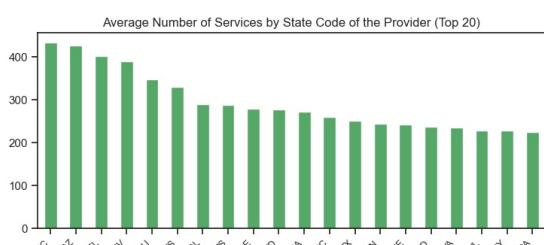
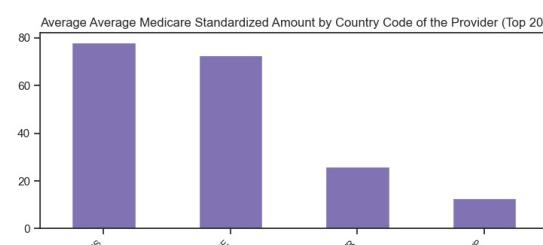
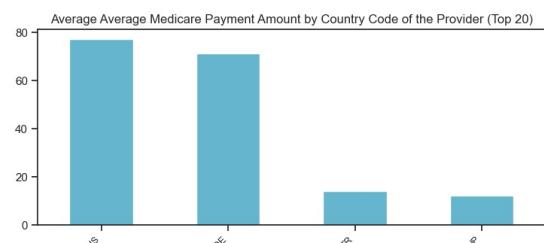
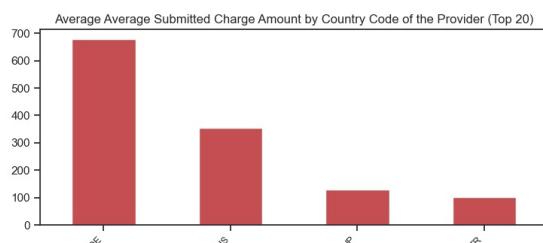
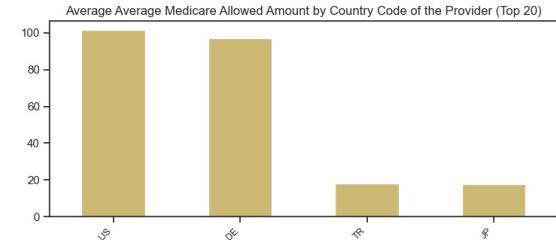
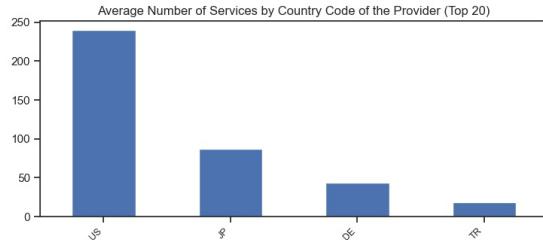
Violin Plot of all Continuous Variables

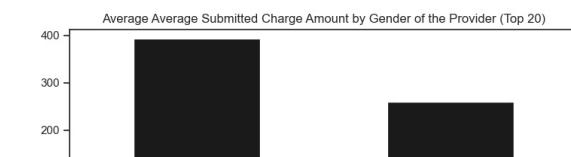
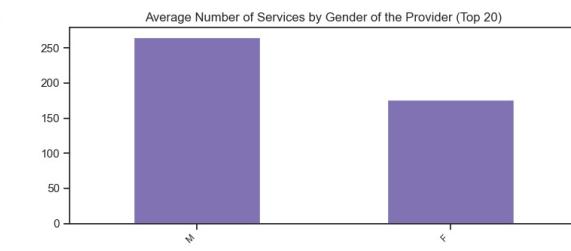
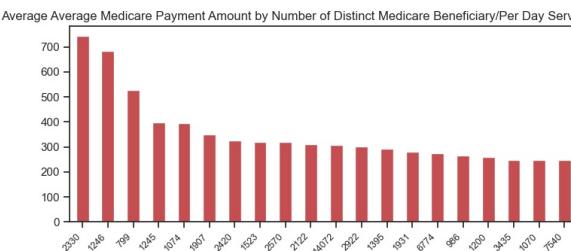
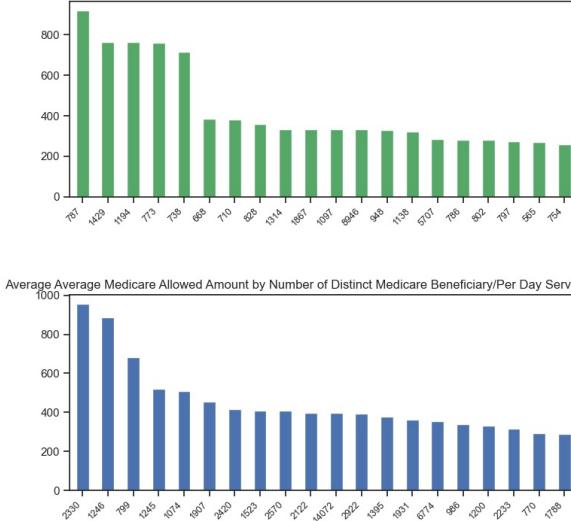
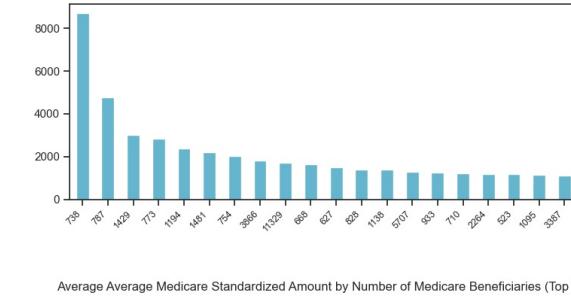
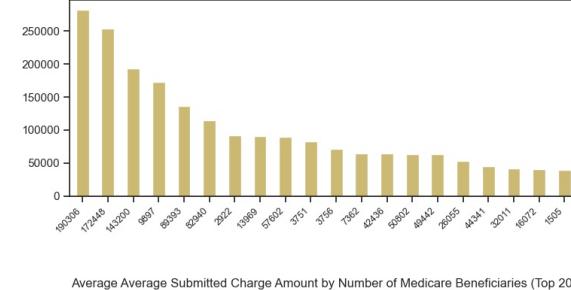
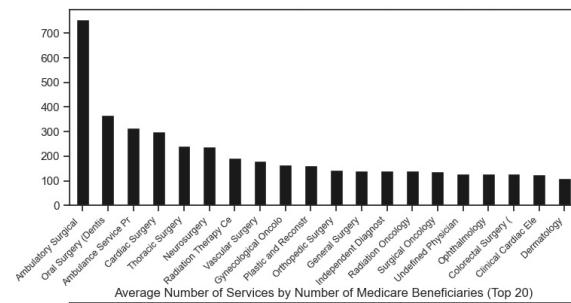
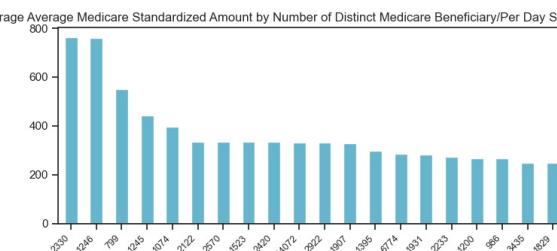
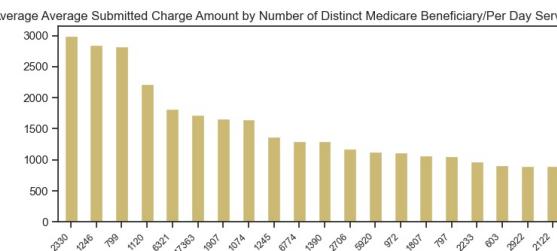
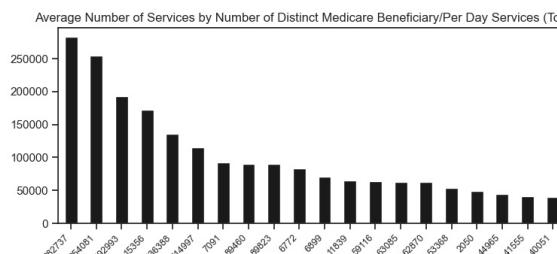
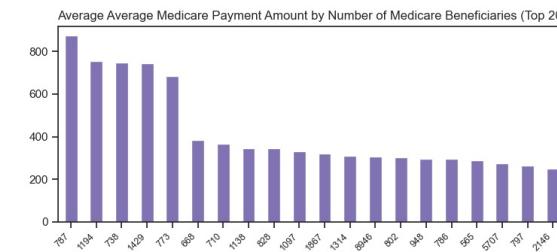
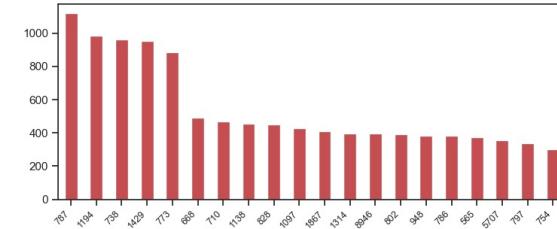
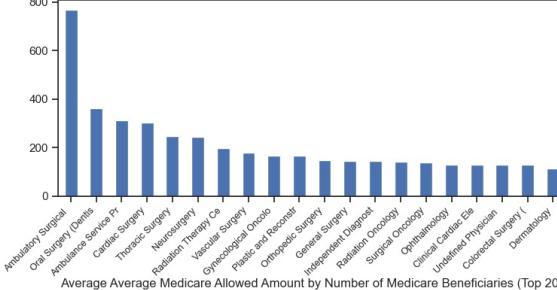
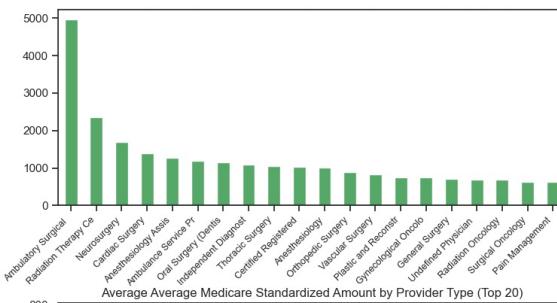


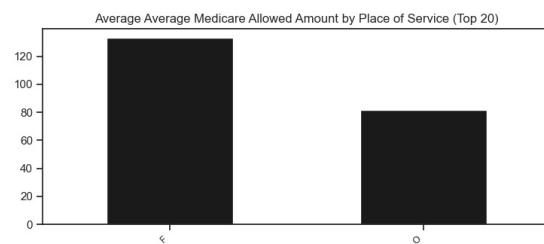
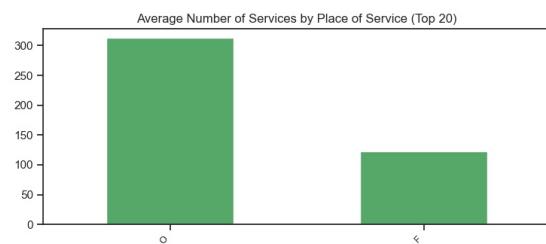
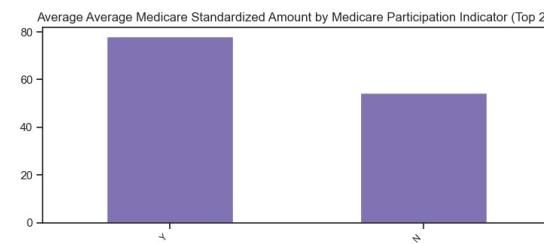
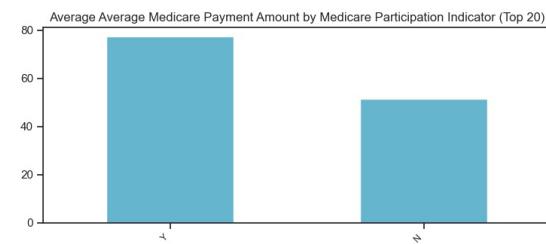
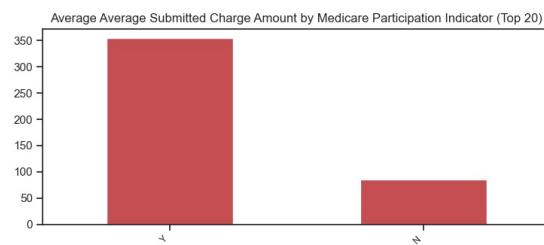
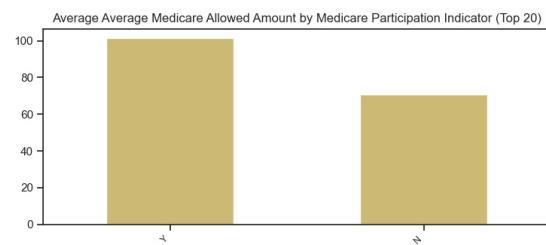
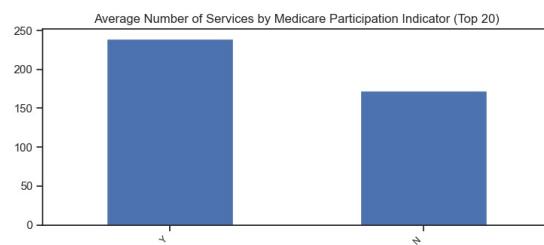
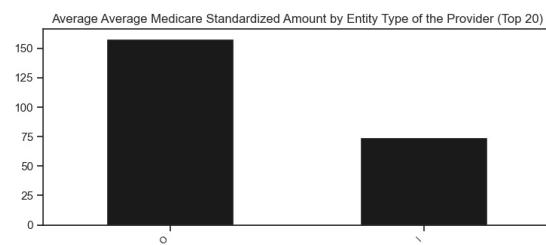
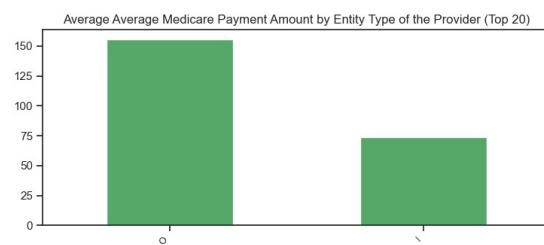
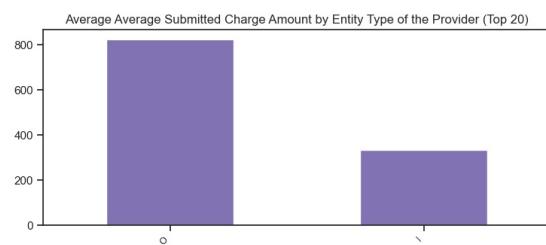
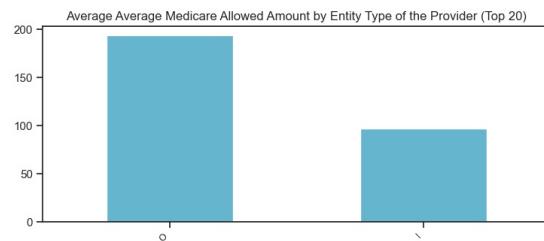
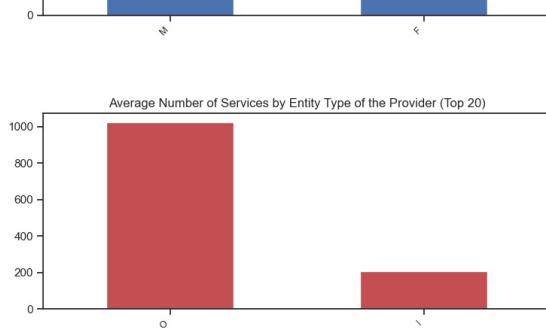
Heatmap of all Numeric Variables including target:

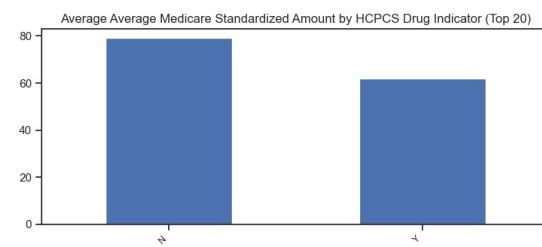
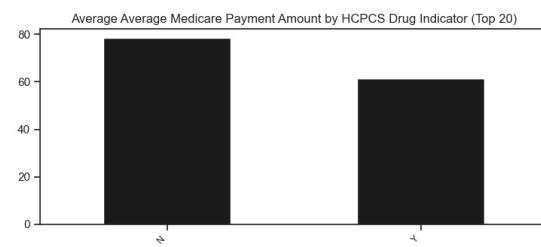
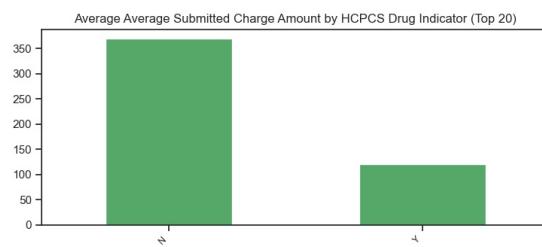
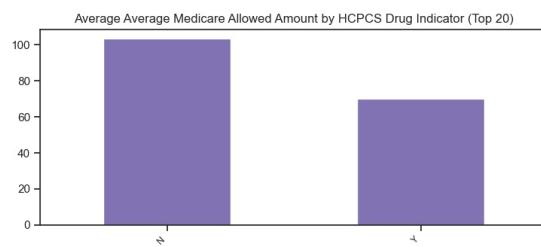
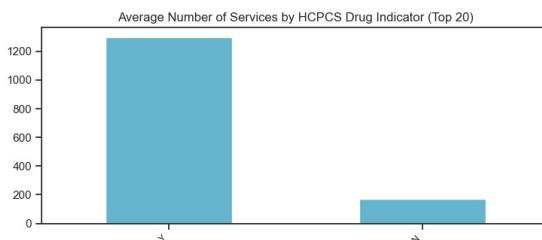
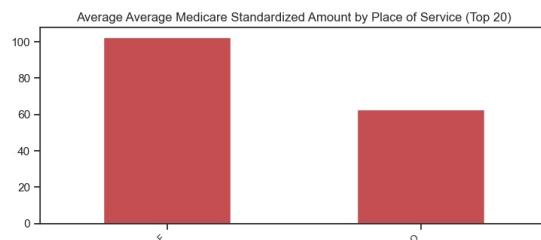
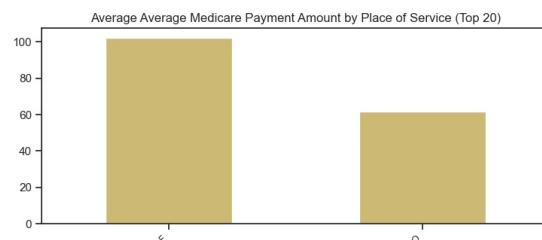
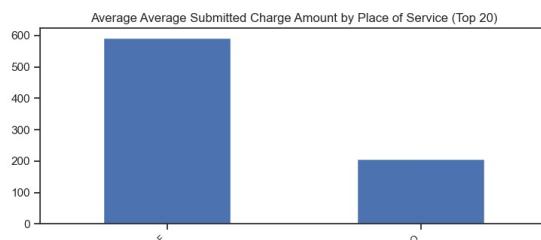


Bar plots for each Continuous by each Categorical variable



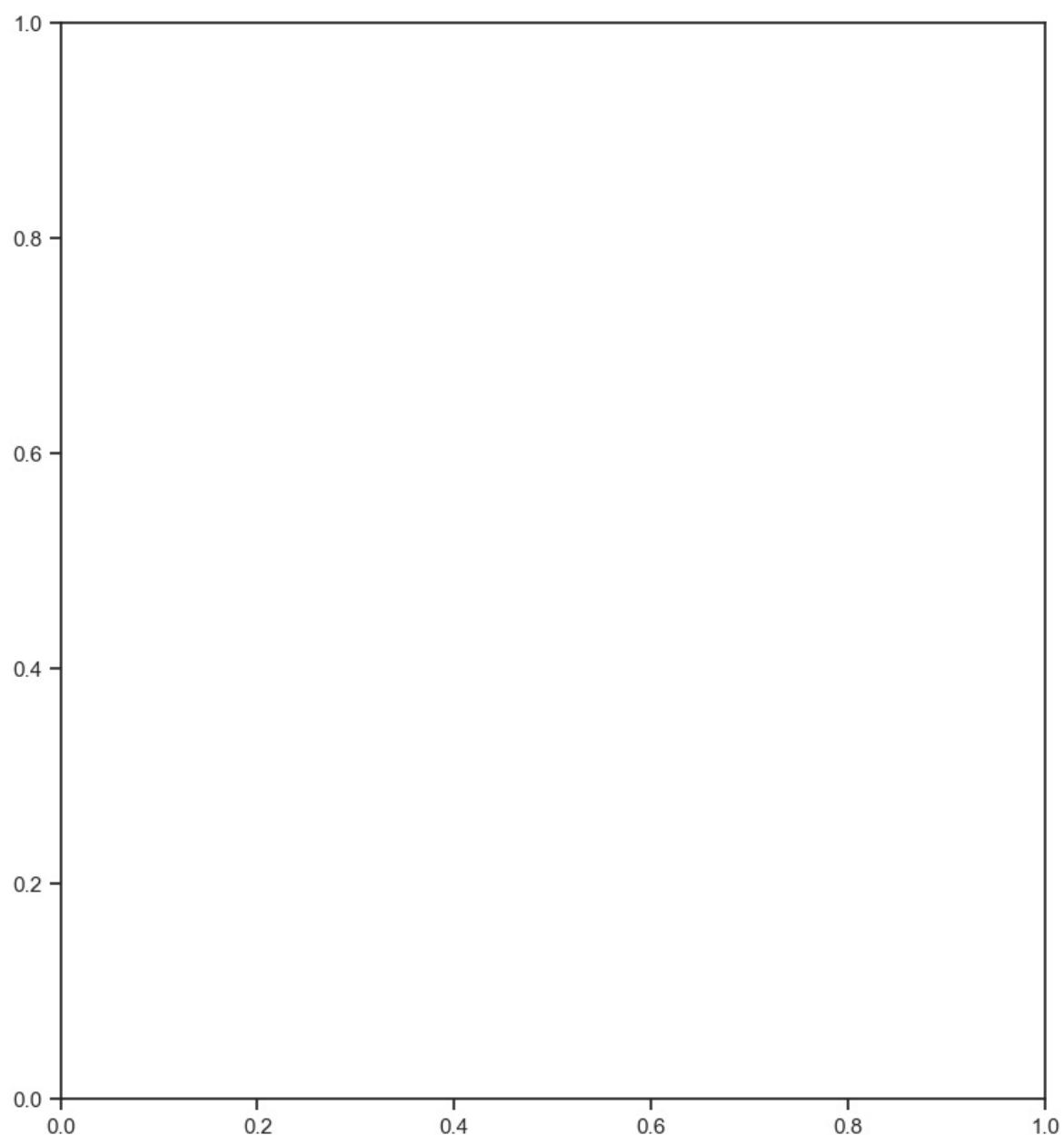




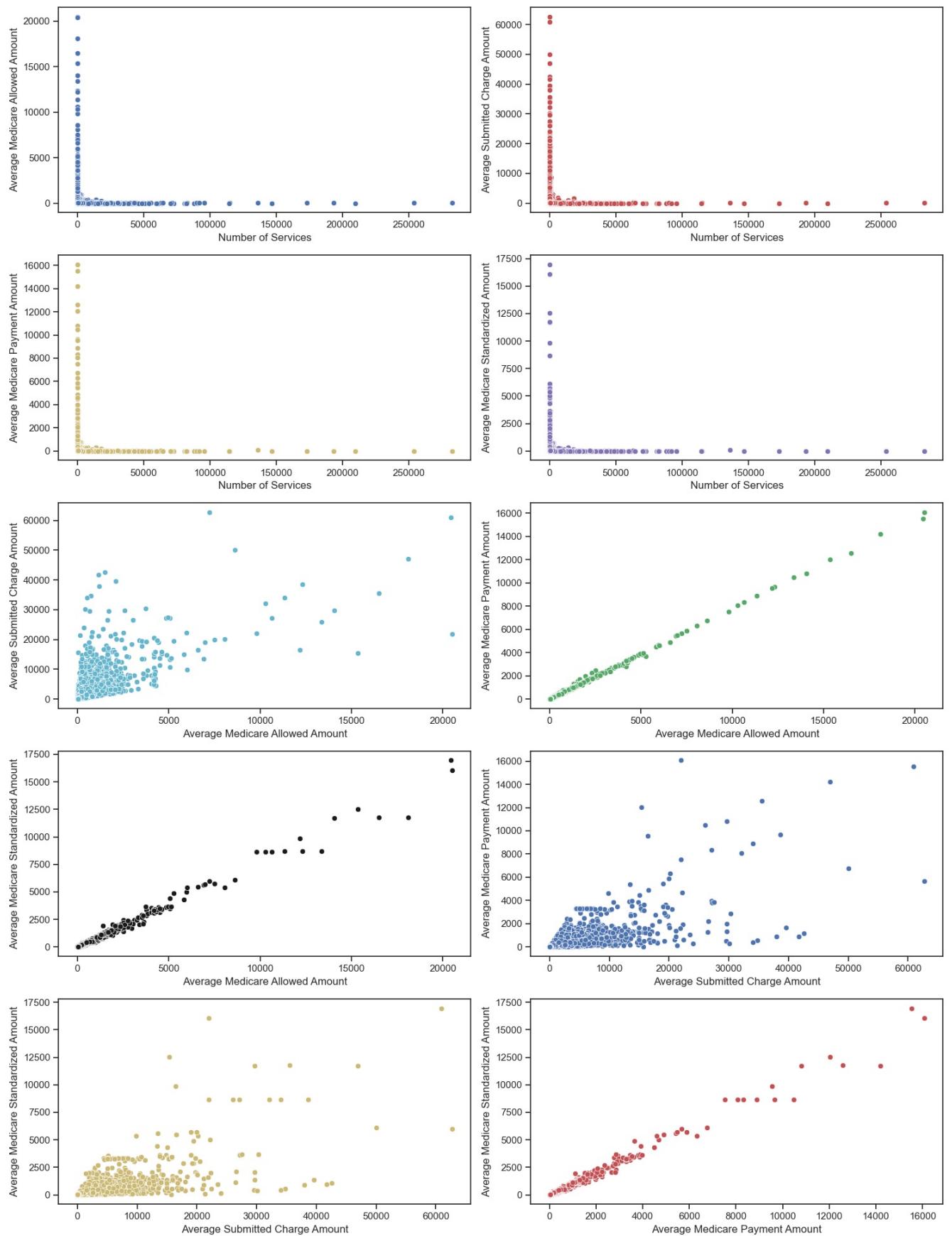


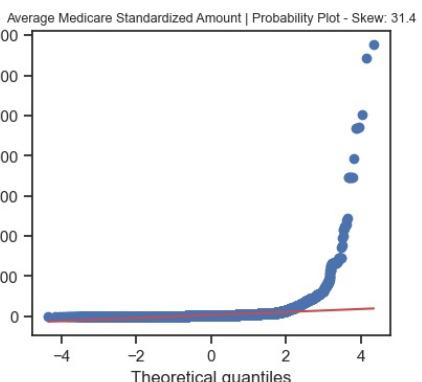
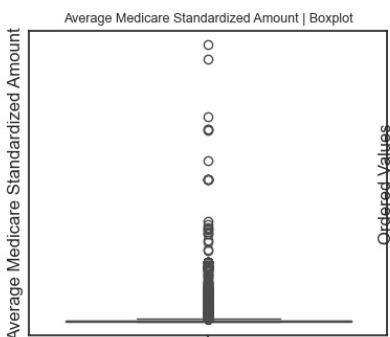
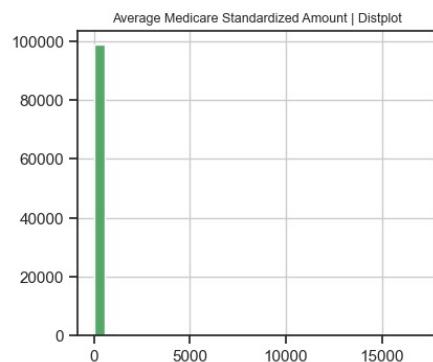
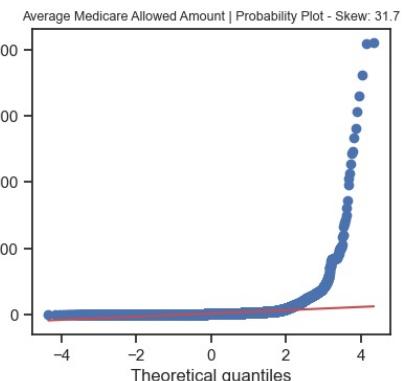
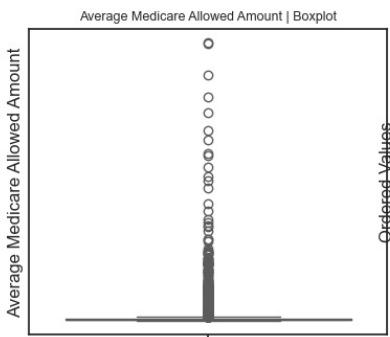
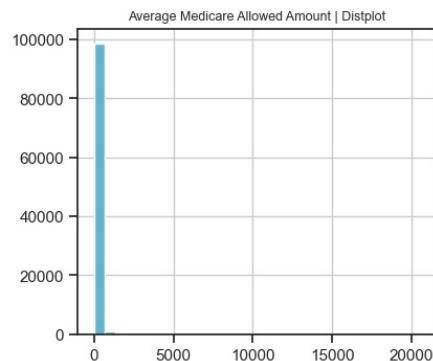
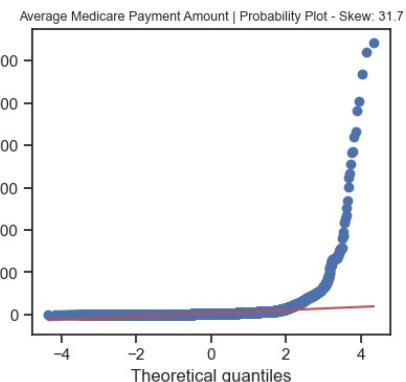
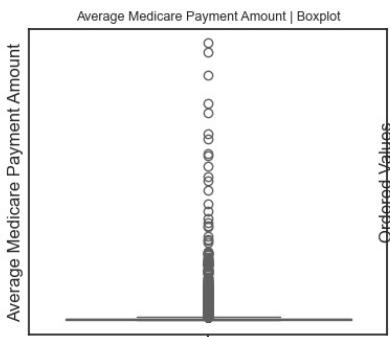
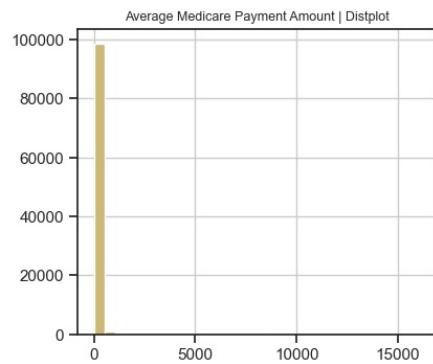
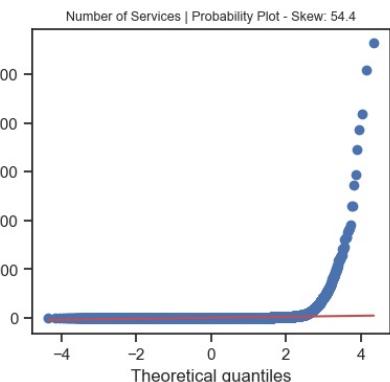
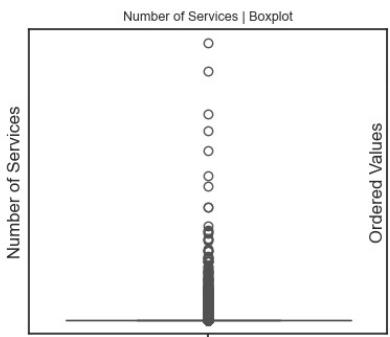
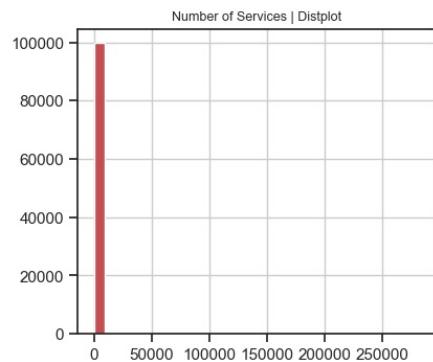
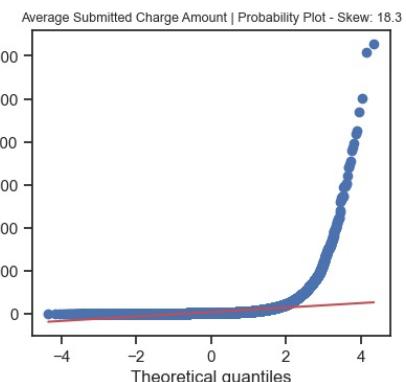
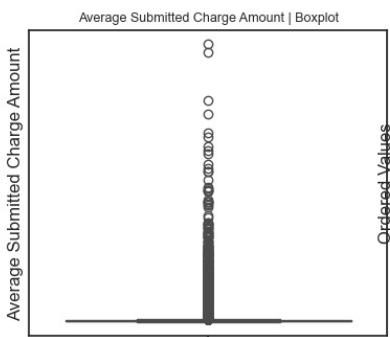
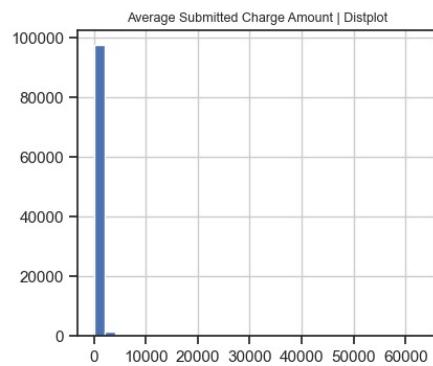
## Wordcloud for Credentials of the Provider



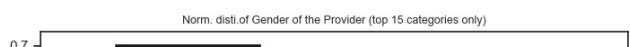


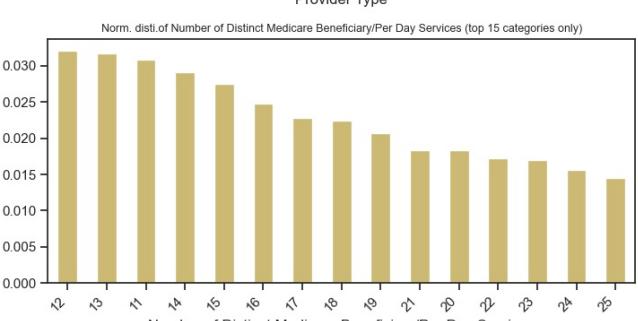
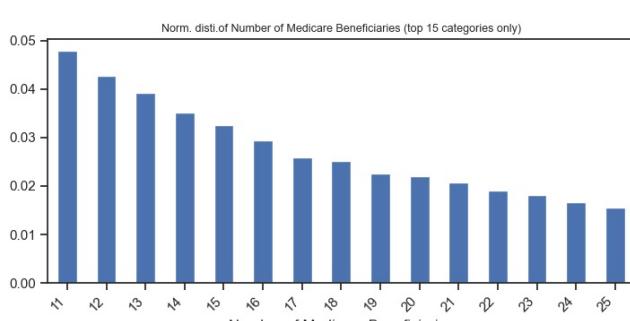
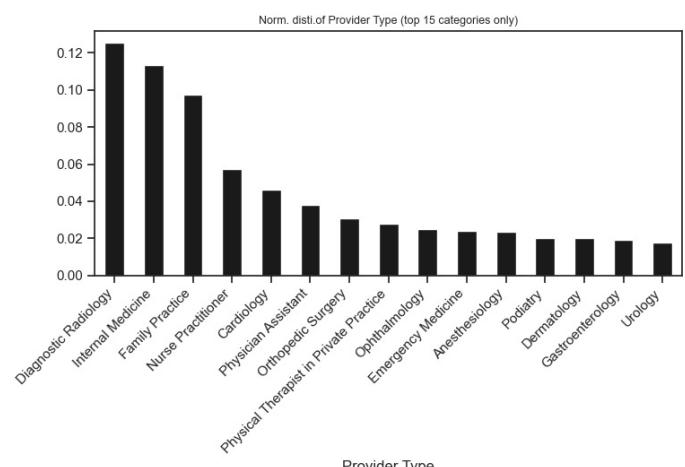
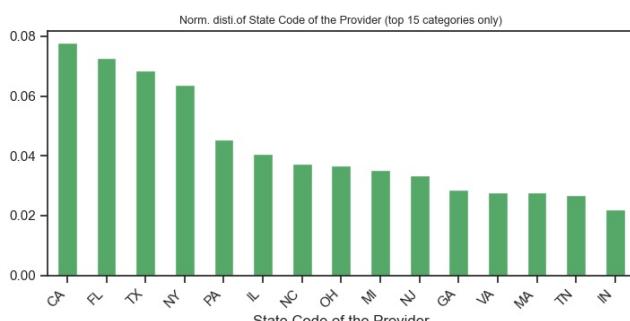
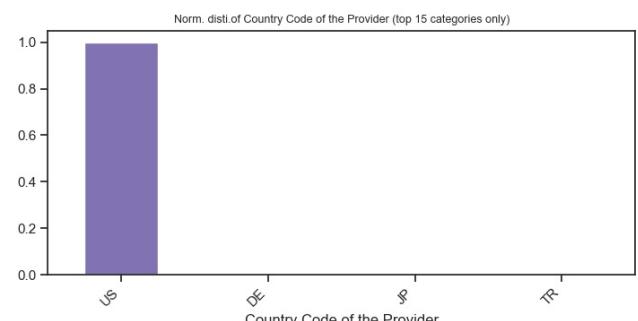
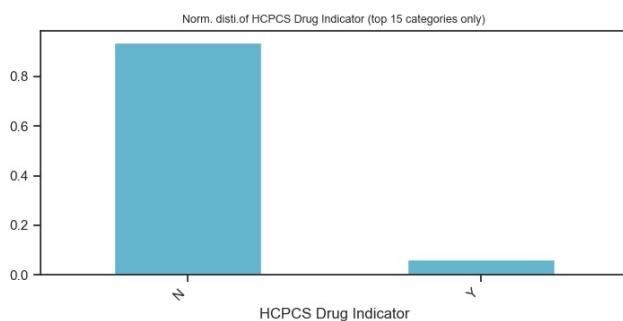
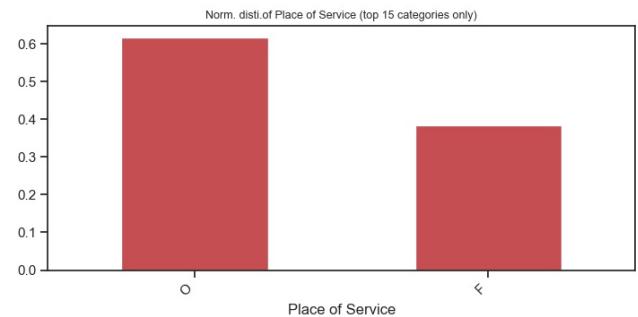
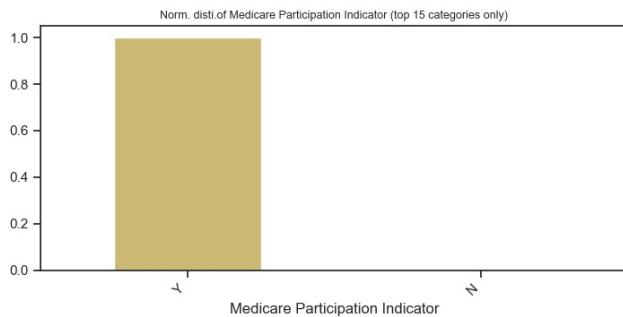
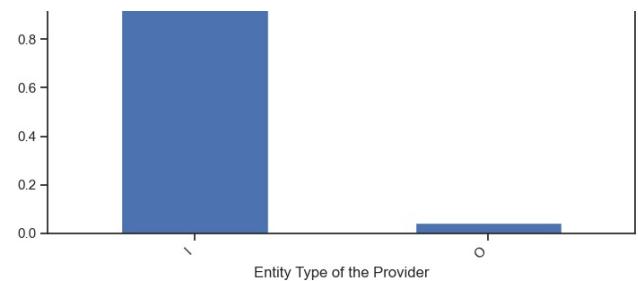
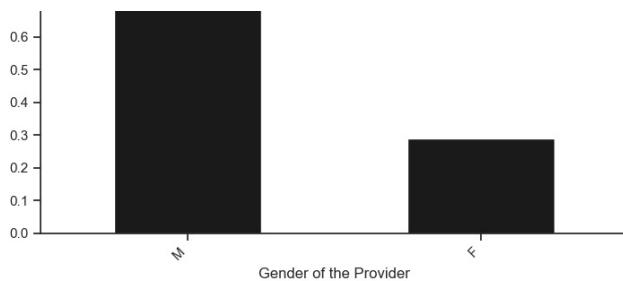
Pair-wise Scatter Plot of all Continuous Variables



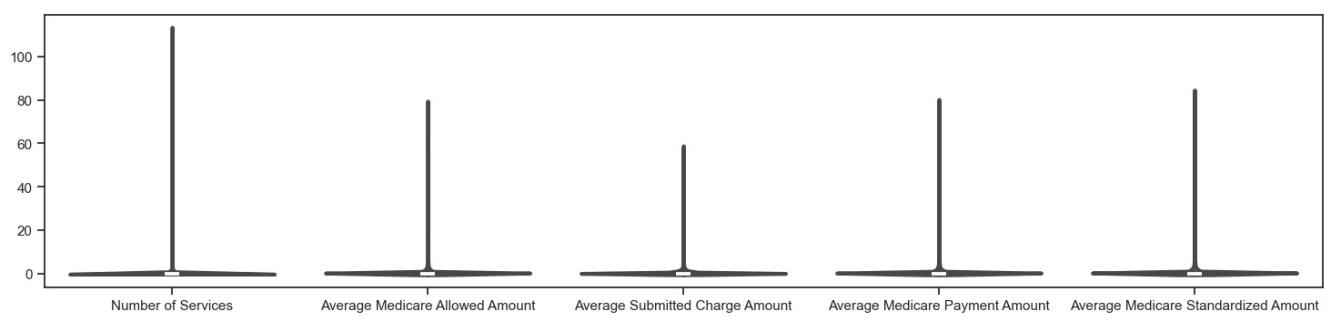


Histograms and Normalized distributions of all variables

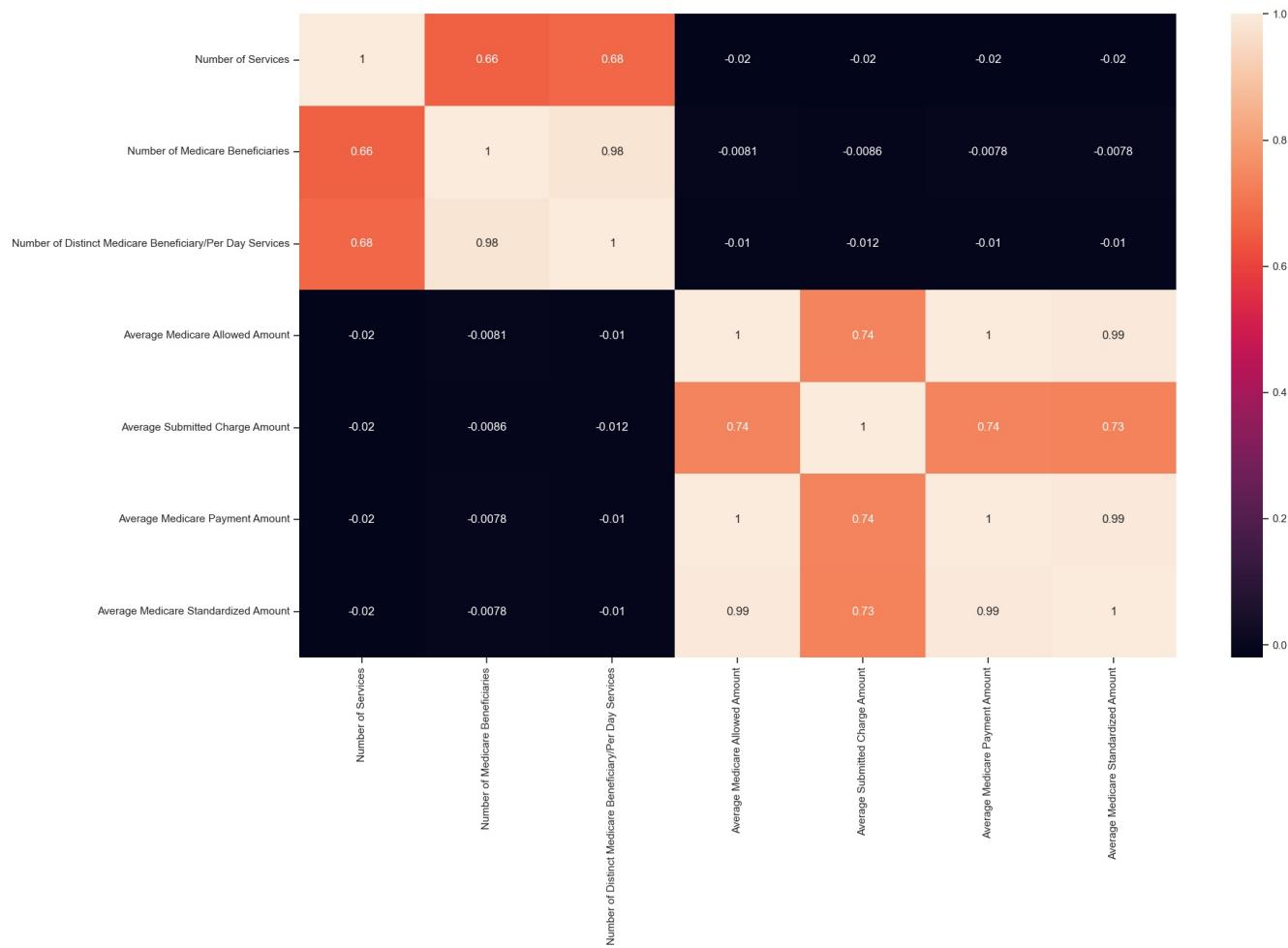




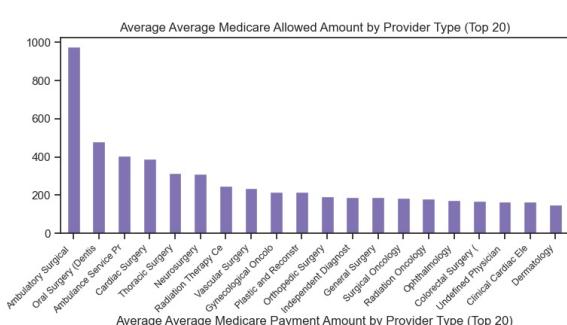
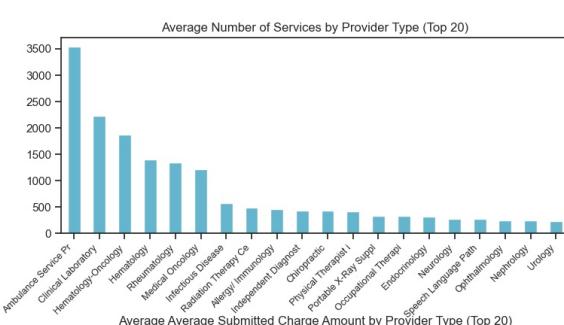
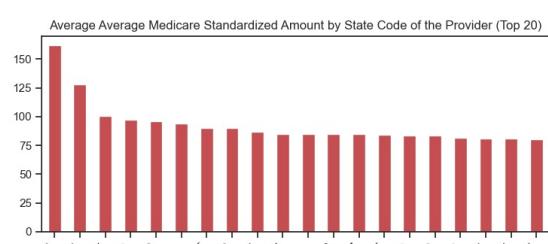
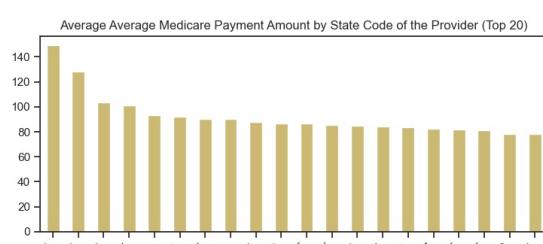
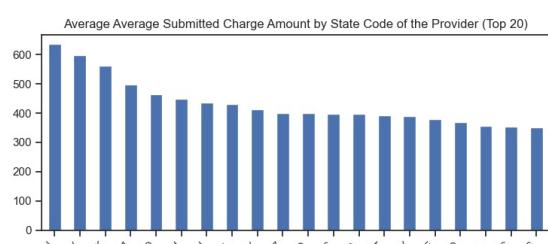
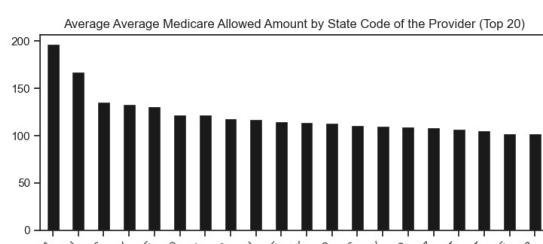
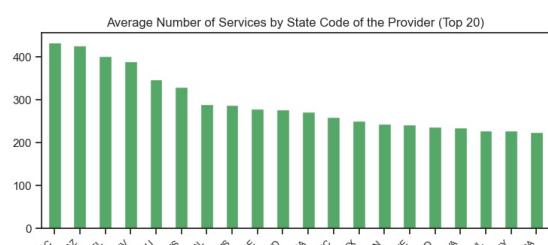
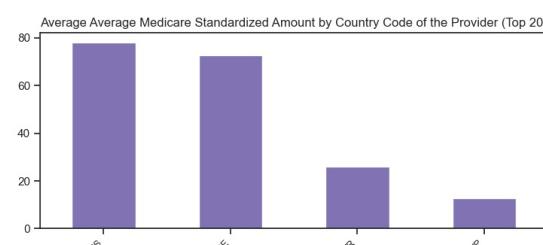
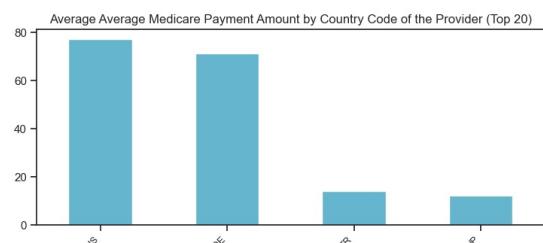
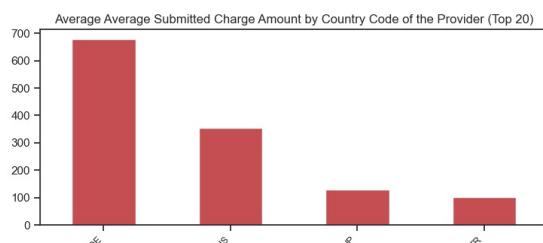
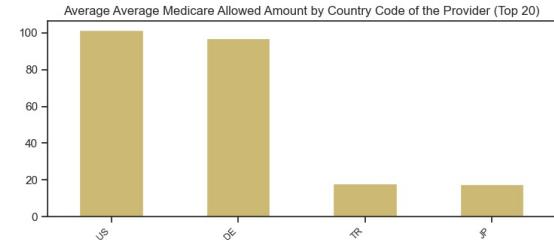
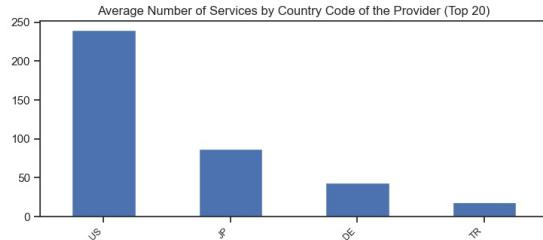
Violin Plot of all Continuous Variables

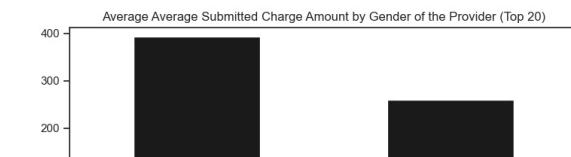
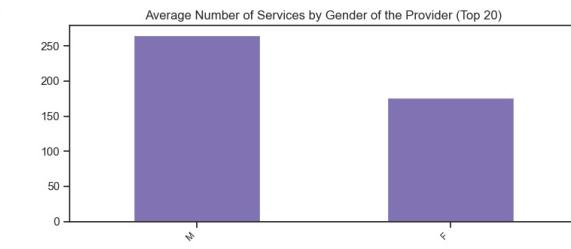
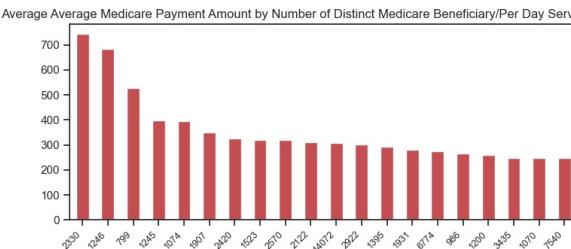
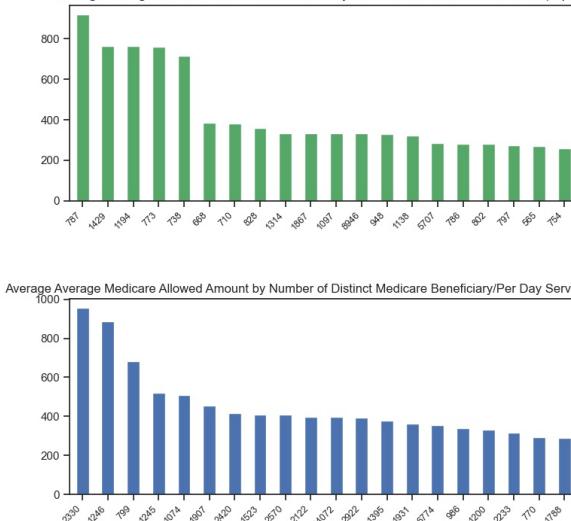
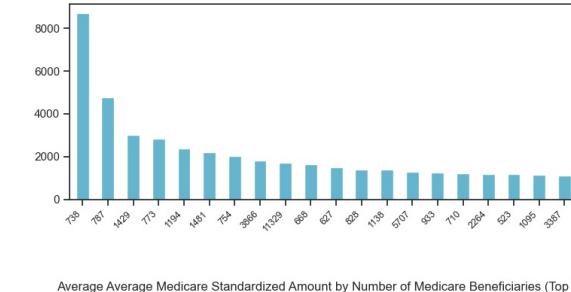
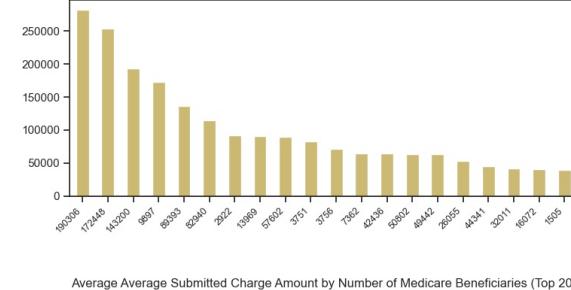
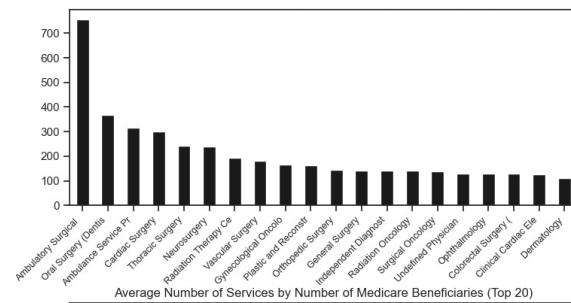
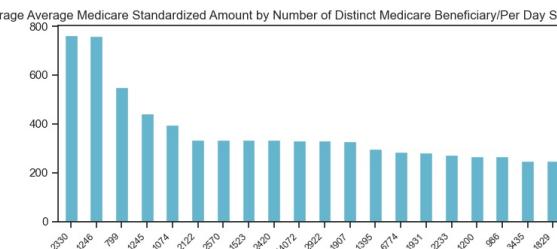
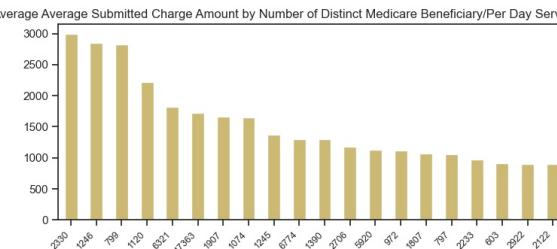
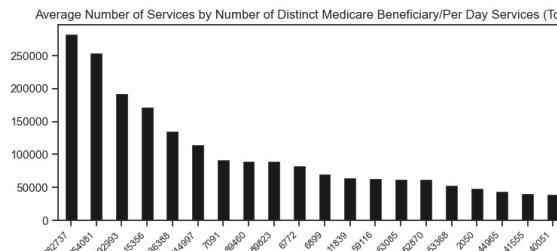
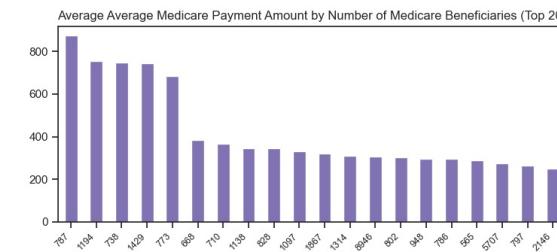
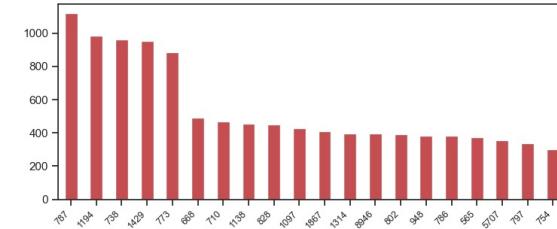
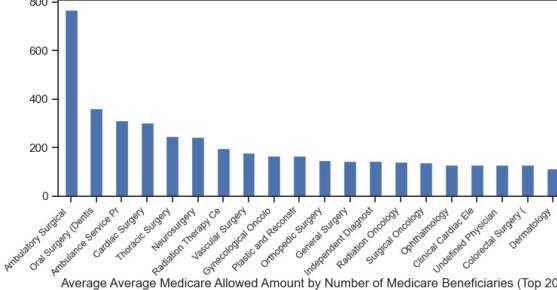
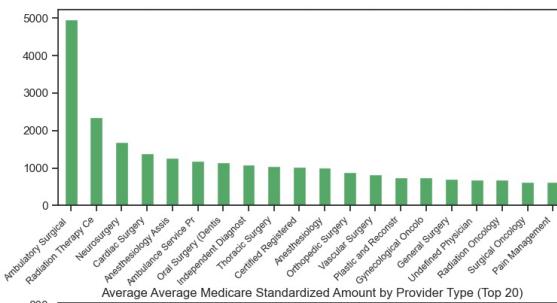


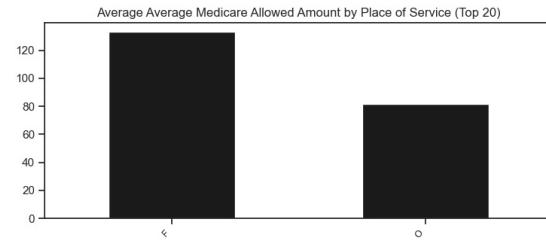
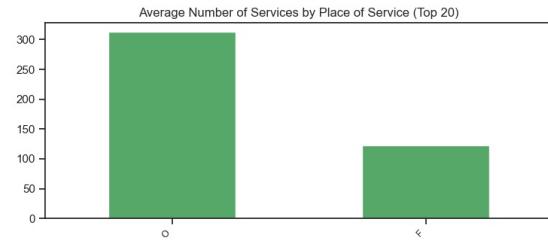
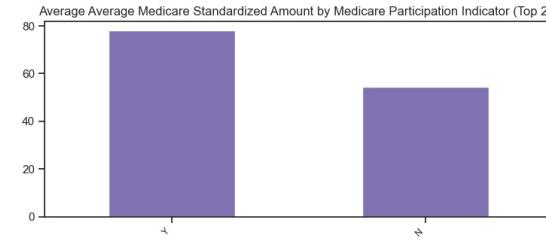
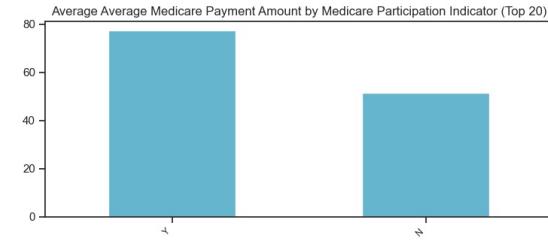
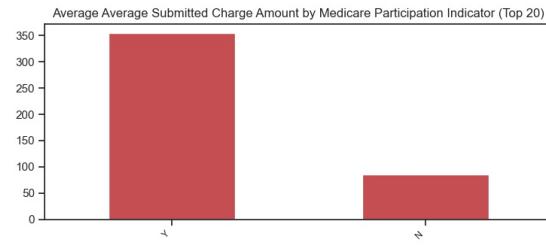
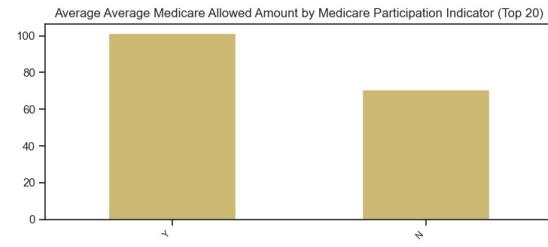
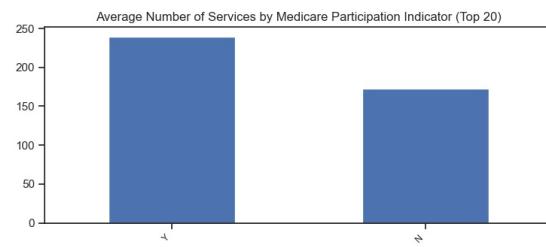
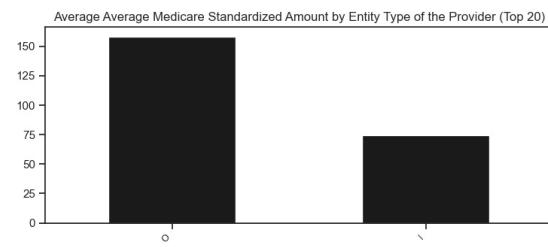
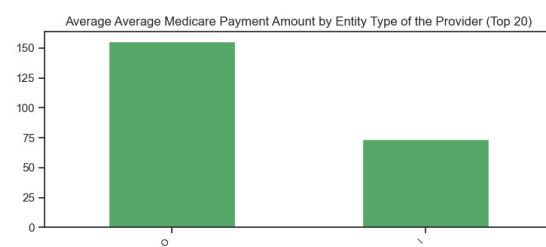
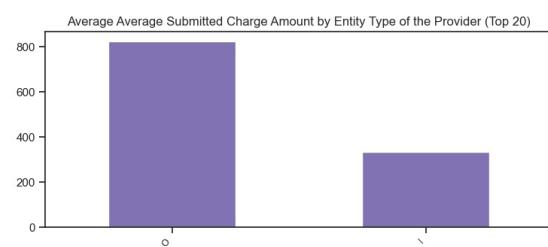
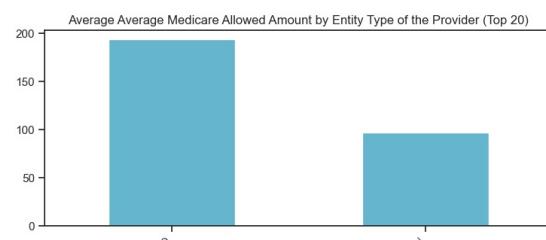
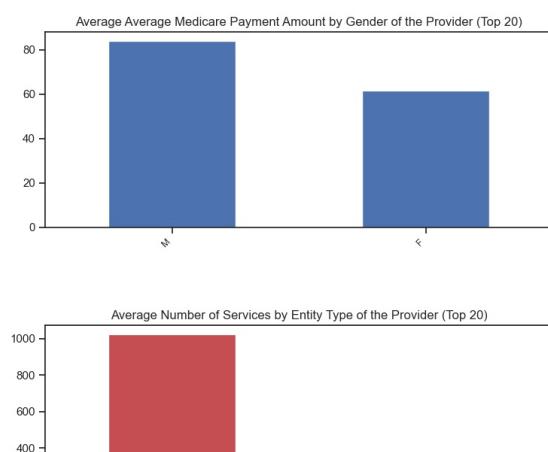
Heatmap of all Numeric Variables including target:

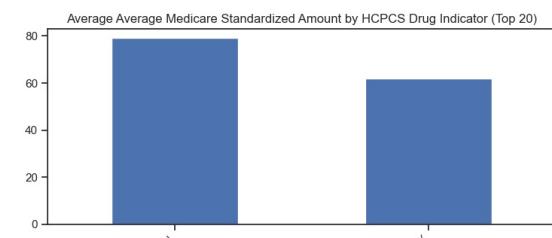
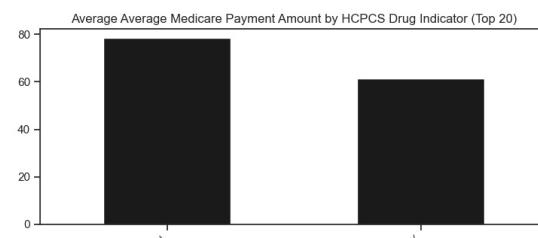
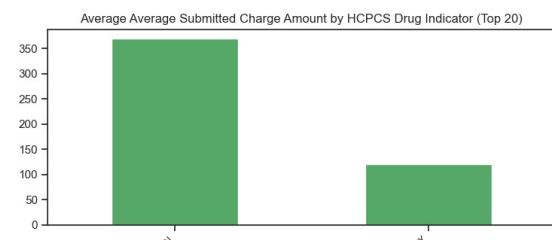
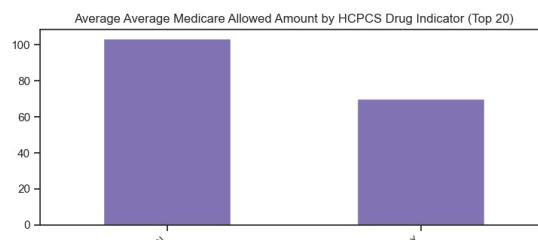
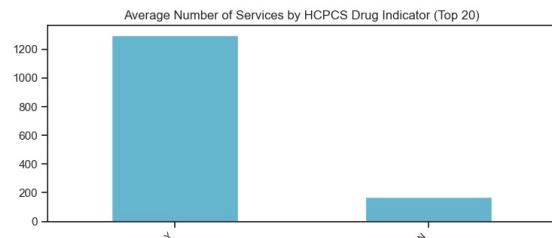
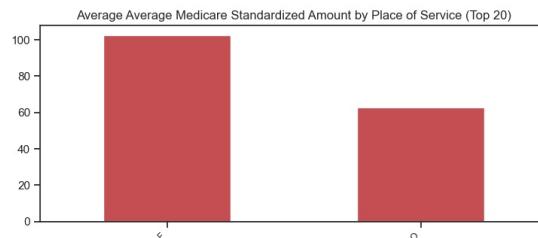
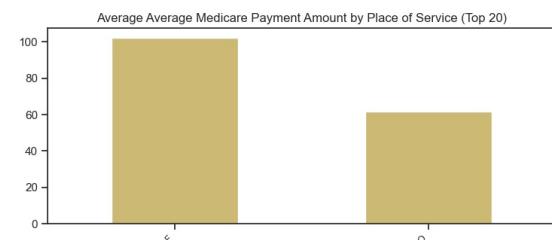
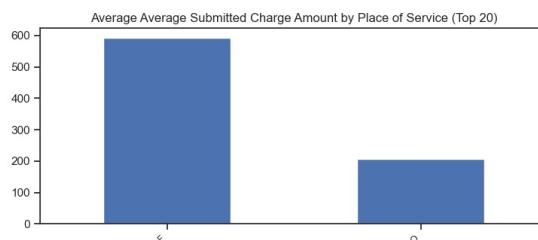


Bar plots for each Continuous by each Categorical variable





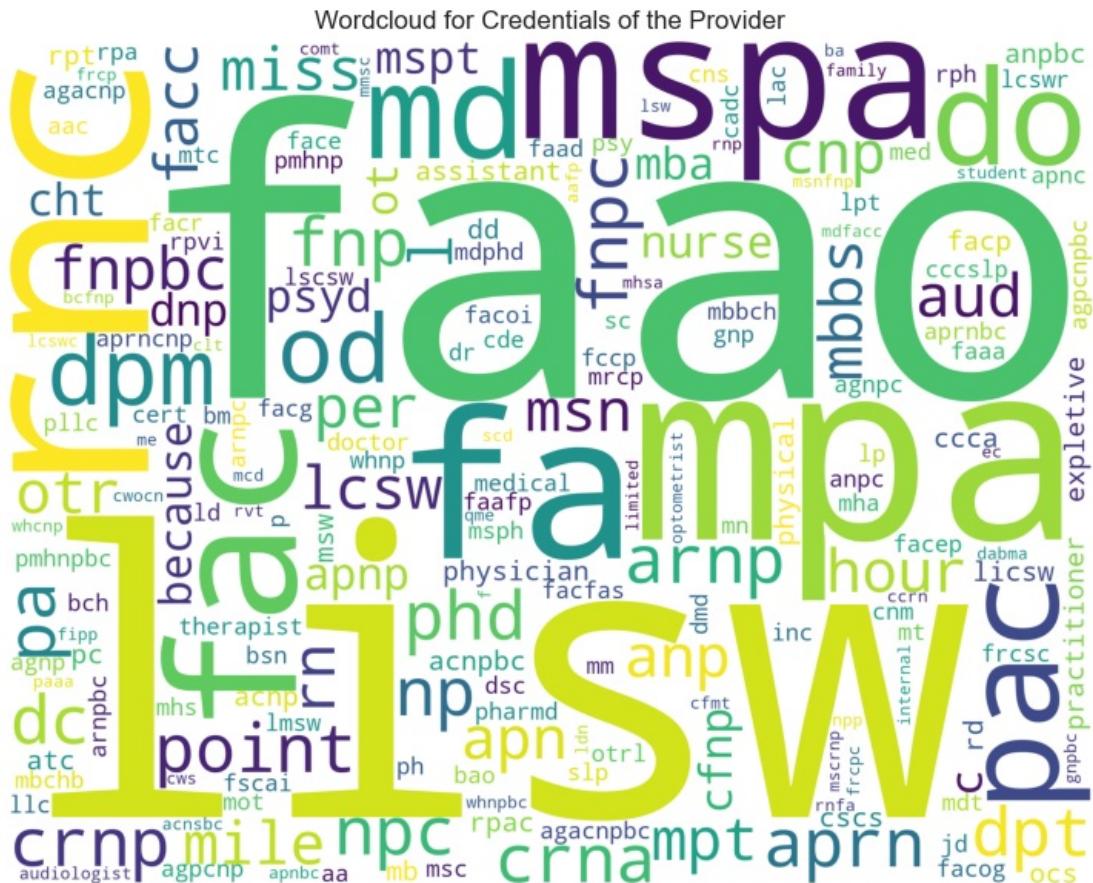




```
[nltk_data]  Downloading collection 'popular'
[nltk_data]
[nltk_data]    |   Downloading package cmudict to
[nltk_data]    |   C:\Users\Pushpendra.PUSHPENDRA-
[nltk_data]    |   PC.001\AppData\Roaming\nltk_data...
[nltk_data]    |   Package cmudict is already up-to-date!
[nltk_data]    |   Downloading package gazetteers to
[nltk_data]    |   C:\Users\Pushpendra.PUSHPENDRA-
[nltk_data]    |   PC.001\AppData\Roaming\nltk_data...
[nltk_data]    |   Package gazetteers is already up-to-date!
[nltk_data]    |   Downloading package genesis to
[nltk_data]    |   C:\Users\Pushpendra.PUSHPENDRA-
[nltk_data]    |   PC.001\AppData\Roaming\nltk_data...
[nltk_data]    |   Package genesis is already up-to-date!
[nltk_data]    |   Downloading package gutenberg to
[nltk_data]    |   C:\Users\Pushpendra.PUSHPENDRA-
[nltk_data]    |   PC.001\AppData\Roaming\nltk_data...
[nltk_data]    |   Package gutenberg is already up-to-date!
[nltk_data]    |   Downloading package inaugural to
[nltk_data]    |   C:\Users\Pushpendra.PUSHPENDRA-
[nltk_data]    |   PC.001\AppData\Roaming\nltk_data...
[nltk_data]    |   Package inaugural is already up-to-date!
[nltk_data]    |   Downloading package movie_reviews to
[nltk_data]    |   C:\Users\Pushpendra.PUSHPENDRA-
[nltk_data]    |   PC.001\AppData\Roaming\nltk_data...
[nltk_data]    |   Package movie_reviews is already up-to-date!
[nltk_data]    |   Downloading package names to
[nltk_data]    |   C:\Users\Pushpendra.PUSHPENDRA-
[nltk_data]    |   PC.001\AppData\Roaming\nltk_data...
[nltk_data]    |   Package names is already up-to-date!
[nltk_data]    |   Downloading package shakespeare to
[nltk_data]    |   C:\Users\Pushpendra.PUSHPENDRA-
[nltk_data]    |   PC.001\AppData\Roaming\nltk_data...
[nltk_data]    |   Package shakespeare is already up-to-date!
[nltk_data]    |   Downloading package stopwords to
[nltk_data]    |   C:\Users\Pushpendra.PUSHPENDRA-
[nltk_data]    |   PC.001\AppData\Roaming\nltk_data...
[nltk_data]    |   Package stopwords is already up-to-date!
[nltk_data]    |   Downloading package treebank to
[nltk_data]    |   C:\Users\Pushpendra.PUSHPENDRA-
[nltk_data]    |   PC.001\AppData\Roaming\nltk_data...
[nltk_data]    |   Package treebank is already up-to-date!
```

```
[nltk_data] | Downloading package twitter_samples to
[nltk_data] |   C:\Users\Pushpendra.PUSHPENDRA-
[nltk_data] |   PC.001\AppData\Roaming\nltk_data...
[nltk_data] | Package twitter_samples is already up-to-date!
[nltk_data] | Downloading package omw to
[nltk_data] |   C:\Users\Pushpendra.PUSHPENDRA-
[nltk_data] |   PC.001\AppData\Roaming\nltk_data...
[nltk_data] | Package omw is already up-to-date!
[nltk_data] | Downloading package omw-1.4 to
[nltk_data] |   C:\Users\Pushpendra.PUSHPENDRA-
[nltk_data] |   PC.001\AppData\Roaming\nltk_data...
[nltk_data] | Package omw-1.4 is already up-to-date!
[nltk_data] | Downloading package wordnet to
[nltk_data] |   C:\Users\Pushpendra.PUSHPENDRA-
[nltk_data] |   PC.001\AppData\Roaming\nltk_data...
[nltk_data] | Package wordnet is already up-to-date!
[nltk_data] | Downloading package wordnet2021 to
[nltk_data] |   C:\Users\Pushpendra.PUSHPENDRA-
[nltk_data] |   PC.001\AppData\Roaming\nltk_data...
[nltk_data] | Package wordnet2021 is already up-to-date!
[nltk_data] | Downloading package wordnet31 to
[nltk_data] |   C:\Users\Pushpendra.PUSHPENDRA-
[nltk_data] |   PC.001\AppData\Roaming\nltk_data...
[nltk_data] | Package wordnet31 is already up-to-date!
[nltk_data] | Downloading package wordnet_ic to
[nltk_data] |   C:\Users\Pushpendra.PUSHPENDRA-
[nltk_data] |   PC.001\AppData\Roaming\nltk_data...
[nltk_data] | Package wordnet_ic is already up-to-date!
[nltk_data] | Downloading package words to
[nltk_data] |   C:\Users\Pushpendra.PUSHPENDRA-
[nltk_data] |   PC.001\AppData\Roaming\nltk_data...
[nltk_data] | Package words is already up-to-date!
[nltk_data] | Downloading package maxent_ne_chunker to
[nltk_data] |   C:\Users\Pushpendra.PUSHPENDRA-
[nltk_data] |   PC.001\AppData\Roaming\nltk_data...
[nltk_data] | Package maxent_ne_chunker is already up-to-date!
[nltk_data] | Downloading package punkt to
[nltk_data] |   C:\Users\Pushpendra.PUSHPENDRA-
[nltk_data] |   PC.001\AppData\Roaming\nltk_data...
[nltk_data] | Package punkt is already up-to-date!
[nltk_data] | Downloading package snowball_data to
[nltk_data] |   C:\Users\Pushpendra.PUSHPENDRA-
[nltk_data] |   PC.001\AppData\Roaming\nltk_data...
[nltk_data] | Package snowball_data is already up-to-date!
[nltk_data] | Downloading package averaged_perceptron_tagger to
[nltk_data] |   C:\Users\Pushpendra.PUSHPENDRA-
[nltk_data] |   PC.001\AppData\Roaming\nltk_data...
[nltk_data] | Package averaged_perceptron_tagger is already up-
[nltk_data] |   to-date!
[nltk_data] Done downloading collection popular
```

Done downloading collection popular



## Wordcloud for City of the Provider



## Wordcloud for HCPCS Codes



## Wordcloud for HCPCS Description



All Plots done

Time to run AutoViz = 181 seconds

##### AUTO VISUALIZATION Completed #####

In [ ]:

In [ ]:

In [ ]:

Loading [MathJax]/extensions/Safe.js