

```
import numpy as np
import pandas as pd
import seaborn as sns
import math
import random
import matplotlib.pyplot as plt
from sklearn import preprocessing
from sklearn.linear_model import LinearRegression
from sklearn.model_selection import train_test_split
from sklearn import metrics
from sklearn.metrics import classification_report
from sklearn.metrics import product
import matplotlib inline
import pylab
from pylab import rcParams
import statmodels.api as sm

import statistics
from scipy import stats
import sklearn
import warnings
warnings.filterwarnings('ignore')
import matplotlib as mpl
fig, ax = plt.subplots()
# Import data to dataframe
medical_df = pd.read_csv('C:/Users/MichaelRupert/Downloads/e8d8am5uf6d75k650df/medical_raw_data.csv')
list(medical_df.columns)
print(medical_df.columns)
```

```
Index(['Unnamed: 0', 'CaseOrder', 'Customer_id', 'Interaction', 'UID', 'City',
      'State', 'County', 'Zip', 'Lat', 'Lng', 'Population', 'Area',
      'Timezone', 'Job', 'Children', 'Age', 'Education', 'Employment',
      'Income', 'Marital', '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', 'Reflex_esophagitis', 'Asthma', 'Services',
      'Allergidays', 'TotalCharge', 'Additional_charges', 'Item1', 'Item2',
      'Item3', 'Item4', 'Item5', 'Item6', 'Item7', 'Item8',
      dtype='object'])
```

```
In [647]: #Drop unneeded columns
medical_df = medical_df.drop(columns=['CaseOrder', 'Customer_id', 'Interaction',
                                     'UID', 'County', 'Timezone', 'Job', 'Education', 'Employment',
                                     'Children', 'Zip', 'Doc_visits',
                                     'Initial_days', axis=1])

medical_df.columns
```

```
Out[647]: Index(['Unnamed: 0', 'City', 'State', 'Lat', 'Lng', 'Population', 'Area',
      'HighBlood', 'Stroke', 'Complication_risk', 'Overweight', 'Arthritis',
      'Diabetes', 'Hyperlipidemia', 'BackPain', 'Anxiety',
      'Allergidays', 'TotalCharge', 'Additional_charges', 'Item1', 'Item2', 'Item3', 'Item4',
      'Item5', 'Item6', 'Item7', 'Item8',
      dtype='object'])
```

```
In [648]: #Describe the data in the dataset
medical_df.describe().transpose()
```

		count	mean	std	min	25%	50%	75%	max
Unnamed: 0	10000	5000.500000	2886.895680	1.000000	2507.500000	5000.500000	7502.500000	10000.000000	
	Lat	10000.0	38.751099	5.402855	17.967190	35.255120	39.419355	42.044175	70.560990
Lng	10000.0	-91.243080	15.059398	-174.209690	-97.352982	-88.375090	-80.448505	-65.290170	
	Population	10000.0	9965.252800	14824.758614	0.000000	694.750000	2769.000000	13945.000000	122814.000000
Age	7586.0	53.295676	20.659187	0.000000	35.000000	53.000000	71.000000	89.000000	
	Income	7536.0	40484.438268	28664.861050	154.080000	19450.792500	33942.280000	54075.235000	207249.130000
VitD_levels	10000.0	1.9412675	6.732377	9.519012	1.651317	18.080560	19.789740	53.019124	
	Full_meals_eaten	10000.0	1.0041400	1.028811	0.000000	0.000000	1.000000	2.000000	7.000000
VitD_supp	10000.0	0.398900	0.628505	0.000000	0.000000	0.000000	1.000000	5.000000	
	Overweight	9018.0	0.709137	0.454186	0.000000	0.000000	1.000000	1.000000	1.000000
Anxiety	9016.0	0.322216	0.467389	0.000000	0.000000	0.000000	1.000000	1.000000	
	TotalCharge	10000.0	5891.538261	6342.601344	1256.751699	3253.239465	5852.250564	7614.989701	21524.224210
AdditionalCharges	10000.0	12945.652856	6542.601544	3125.702716	7996.487642	11581.979365	15626.601495	30566.073130	
	Item1	10000.0	3.518800	1.031966	1.000000	3.000000	4.000000	4.000000	8.000000
Item2	10000.0	3.506700	1.034825	1.000000	3.000000	3.000000	4.000000	7.000000	
	Item3	10000.0	3.511100	1.022755	1.000000	3.000000	4.000000	4.000000	8.000000
Item4	10000.0	3.515100	1.036282	1.000000	3.000000	4.000000	4.000000	7.000000	
	Item5	10000.0	3.496600	1.030192	1.000000	3.000000	3.000000	4.000000	7.000000
Item6	10000.0	3.522500	1.032376	1.000000	3.000000	4.000000	4.000000	7.000000	
	Item7	10000.0	3.494000	1.021405	1.000000	3.000000	3.000000	4.000000	7.000000
	Item8	10000.0	3.509700	1.042312	1.000000	3.000000	3.000000	4.000000	7.000000

```
In [649]: #check for duplicates
is_dups_bool = medical_df.duplicated()
print("Are there any duplicates? " + str(is_dups_bool.value_counts()))
```

```
Are there any duplicates? False      10000
dtype: int64
```

```
In [650]: #check for zeros
zeros = (medical_df["Population"] == 0).sum()
print(zeros)
```

```
109
```

```
In [651]: #drop records with zeros in population field
count = 0
secolist = []
for x in medical_df["Population"]:
```

```
    if x ==0:
        secolist.append(count)
        count = count + 1
medical_df.drop(secolist, axis=0, inplace=True)
medical_df.reset_index(inplace=True)
medical_df.describe().transpose()
```

```
zeros = (medical_df["Population"] == 0).sum()
print(zeros)
```

```
0
```

```
In [652]: #encode the data with numerics
my_dict3={"yes": 1, "No": 0}
my_dict4={"Male": 1, "Female": 2, "Prefer not to answer": 0}
my_dict5={"Low":1, "Medium":2, "High":3}
```

```
medical_df_V2 = medical_df.replace({"ReAdmis": my_dict3, "Gender": my_dict4, "Soft_drink": my_dict3,
                                   "HighBlood": my_dict3,
                                   "Stroke": my_dict3, "Complication_risk": my_dict5, "Arthritis": my_dict3,
                                   "Hyperlipidemia": my_dict3, "BackPain": my_dict3, "Diabetes": my_dict3,
                                   "Allergidays": my_dict3, "Reflex_esophagitis": my_dict3, "Asthma": my_dict3})
medical_df_V2.describe().transpose()
```

		count	mean	std	min	25%	50%	75%	max
Unnamed: 0	9891.0	5004.544839	2885.109958	1.000000	2507.500000	5004.000000	7502.500000	10000.000000	
	Lat	9891.0	38.761388	5.402855	17.967190	35.268175	39.440770	42.049850	70.560990
Lng	9891.0	-91.215769	15.176267	-174.209690	-97.319560	-88.375090	-80.444955	-65.290170	
	Population	9891.0	10075.072086	14869.065326	1.000000	730.000000	2859.000000	14161.000000	122814.000000
Age	7596.0	53.331601	20.664683	0.000000	35.000000	53.000000	71.000000	89.000000	
	Income	7466.0	40485.233793	28659.104817	154.080000	19454.922500	33947.305000	54107.657500	207249.130000
Gender	9891.0	1.480942	0.540864	0.000000	1.000000	2.000000	2.000000	2.000000	
	ReAdmis	9891.0	0.366899	0.628505	0.000000	0.000000	0.000000	1.000000	1.000000
VitD_levels	9891.0	1.9415365	6.732357	9.519012	1.6514963	18.078917	19.792348	53.019124	
	Full_meals_eaten	9891.0	1.000506	1.008006	0.000000	0.000000	1.000000	2.000000	7.000000
VitD_supp	9891.0	0.398847	0.628892	0.000000	0.000000	0.000000	1.000000	5.000000	
	Soft_drink	7446.0	0.258125	0.437633	0.000000	0.000000	0.000000	1.000000	1.000000
HighBlood	9891.0	0.409362	0.491741	0.000000	0.000000	0.000000	1.000000	1.000000	
	Stroke	9891.0	0.199070	0.399321	0.000000	0.000000	0.000000	0.000000	1.000000
Complication_risk	9891.0	2.123142	0.730667	1.000000	2.000000	2.000000	3.000000	3.000000	
	Overweight	8521.0	0.709562	0.453990	0.000000	0.000000	1.000000	1.000000	1.000000
Arthritis	9891.0	0.358103	0.479467	0.000000	0.000000	0.000000	1.000000	1.000000	
	Diabetes	9891.0	0.272975	0.445511	0.000000	0.000000	0.000000	1.000000	1.000000
Hyperlipidemia	9891.0	0.337883	0.473012	0.000000	0.000000	0.000000	1.000000	1.000000	
	BackPain	9891.0	0.411586	0.492146	0.000000	0.000000	0.000000	1.000000	1.000000
Anxiety	9916.0	0.321781	0.467186	0.000000	0.000000	0.000000	1.000000	1.000000	
	Allergic_rhinitis	9891.0	0.393287	0.488504	0.000000	0.000000	0.000000	1.000000	1.000000
Reflex_esophagitis	9891.0	0.413305	0.492452	0.000000	0.000000	0.000000	1.000000	1.000000	
	Asthma	9891.0	0.289859	0.453720	0.000000	0.000000	0.000000	1.000000	1.000000
TotalCharge	9891.0	5893.668940	6378.780141	1256.751699	3256.755457	5855.357730	7615.923547	21524.224210	
	Additional_charges	9891.0	12945.165184	6541.957948	3125.702716	7997.923737	11581.934090	15626.601495	30566.073130
Item1	9891.0	3.519968	1.032378	1.000000	3.000000	4.000000	4.000000	8.000000	
	Item2	9891.0	3.506723	1.035678	1.000000	3.000000	3.000000	4.000000	8.000000
Item3	9891.0	3.512587	1.033082	1.000000	3.000000	3.000000	4.000000	7.000000	
	Item4	9891.0	3.514811	1.037058	1.000000	3.000000	4.000000	4.000000	7.000000
Item5	9891.0	3.497826	1.029823	1.000000	3.000000	3.000000	4.000000	7.000000	
	Item6	9891.0	3.523607	1.032399	1.000000	3.000000	4.000000	4.000000	7.000000
Item7	9891.0	3.494692	1.023311	1.000000	3.000000	3.000000	4.000000	7.000000	
	Item8	9891.0	3.508543	1.042476	1.000000	3.000000	3.000000	4.000000	7.000000

```
In [651]: #find fields and columns with nan values
nan_values = medical_df_V2.isnull()
nan_columns = nan_values.any()
```

```
columns_with_nan = medical_df_V2.columns[nan_columns].tolist()
print(columns_with_nan)
```

```
['Age', 'Income', 'Soft_drink', 'Overweight', 'Anxiety']
```

```
In [654]: # Replace this set of columns with a trinary value
medical_df_V2[['Soft_drink', 'Overweight', 'Anxiety']] = medical_df_V2[['Soft_drink', 'Overweight', 'Anxiety']].fillna(value = 2)
```

```
In [655]: #recheck for nulls in those columns
nan_values = medical_df_V2.isnull()
nan_columns = nan_values.any()
```

```
columns_with_nan = medical_df_V2.columns[nan_columns].tolist()
print(columns_with_nan)
```

```
['Age', 'Income']
```

```
In [656]: #replace continuous values with randomly chosen values of similar tupe either integer or floats
count = 0
for x in medical_df_V2["Age"]:
```

```
    randage = random.randint(18,89)
    if math.isnan(x):
        medical_df_V2.loc[count,medical_df_V2.columns.get_loc("Age")] = randage
        count = count + 1
    randage = 0
    count = 0
```

```
for x in medical_df_V2["Income"]:
```

```
    randinc = random.uniform(200.0,208000.0)
    if math.isnan(x):
        medical_df_V2.loc[count,medical_df_V2.columns.get_loc("Income")] = randinc
        count = count + 1
    randage = 0
    count = 0
```

```
nan_values = nan_values.any()
columns_with_nan = medical_df_V2.columns[nan_columns].tolist()
print(columns_with_nan)
medical_df_V2.describe().transpose()
```

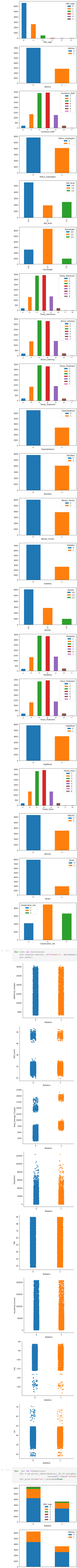
```

'Hyperlipidemia', 'BackPain', 'Anxiety', 'Allergic_rhinitis',
'Reflux_esophagitis', 'Asthma', 'Timely_Admission', 'Timely_Treatment',
'Timely_Visits', 'Reliability', 'Mourne_Treatment', 'Cautious_Staff',
'Active_Listening')

for x in Continuous:
    medical_df_V2[[x]].hist()

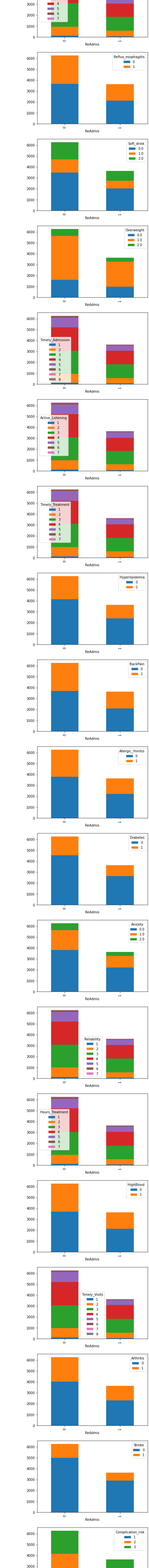
```





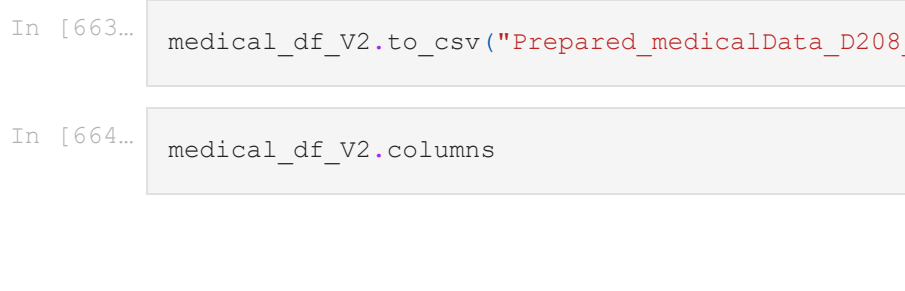
In [66]:  

```
for cont in Continuous:  
    sns.catplot(y=cont, x="ReAdmis", data=medical_df_V2)  
    plt.show()
```



In [66]:  

```
cat = cat in Categorical:  
tbl = pd.pivot_table(medical_df_V2.groupby(['ReAdmis', cat]).size().reset_index(),  
                     values=0, index='ReAdmis', columns=cat, aggfunc=np.sum)  
tbl.plot(kind='bar', stacked=True)
```



In [66]:  

```
medical_df_V2.to_csv("Prepared_medicalData_D208_Task2.csv", index=False)
```

In [66]:  

```
medical_df_V2.columns
```

```
Index(['Unnamed: 0', 'Lat', 'Lng', 'Population', 'Age', 'Household_Income',
'Readmis', 'VitD_levels', 'VitD_supp', 'Soft_drink', 'HighBlood',
'Stroke', 'Complication_risk', 'Overweight', 'Arthritis', 'Diabetes',
'Hyperlipidemia', 'BackPain', 'Anxiety', 'Allergic_rhinitis',
'Reflux_esophagitis', 'Asthma', 'Daily_Average_Charges',
'Additional_charges', 'Timely_Admission', 'Timely_Treatment',
'Timely_Visits', 'Reliability', 'Hours_Treatment', 'Courteous_Staff',
'Active_Listening'],
dtype='object')
```

```
In [665]: #build an initial model using logit
model_Readmis_V1 = sm.Logit(medical_df_V2[['Readmis'],medical_df_V2[['Lat', 'Lng', 'Population', 'Age', 'Household_Income',
'VitD_levels', 'VitD_supp', 'Soft_drink', 'HighBlood',
'Stroke', 'Complication_risk', 'Overweight', 'Arthritis', 'Diabetes',
'Hyperlipidemia', 'BackPain', 'Anxiety', 'Allergic_rhinitis',
'Reflux_esophagitis', 'Asthma', 'Daily_Average_Charges',
'Additional_charges', 'Timely_Admission', 'Timely_Treatment',
'Timely_Visits', 'Reliability', 'Hours_Treatment', 'Courteous_Staff',
'Active_Listening']]).fit()
print(model_Readmis_V1.params)
print(model_Readmis_V1.summary())
```

```
Optimization terminated successfully.
Current function value: 0.213818
Iterations: 9

Lat                0.003802
Lng                0.000435
Population         0.000003
Age               -0.001739
Household_Income  0.000001
VitD_levels       -0.627233
VitD_supp         -0.054598
Soft_drink        0.031811
HighBlood         -0.172891
Stroke            0.123099
Complication_risk -0.274730
Overweight        0.017370
Arthritis         -0.219184
Diabetes          -0.165793
Hyperlipidemia   -0.148815
BackPain         -0.106898
Anxiety           -0.125126
Allergic_rhinitis -0.263822
Reflux_esophagitis -0.165699
Asthma           -0.058190
Daily_Average_Charges 0.001862
Additional_charges 0.000007
Timely_Admission  0.054570
Timely_Treatment  0.092855
Timely_Visits     -0.030254
Reliability       0.088076
Hours_Treatment  -0.058190
Courteous_Staff  -0.058190
Active_Listening  -0.066533
dtype: float64

Logit Regression Results
=====
Dep. Variable:      Readmis      No. Observations:      9891
Model:              Logit      Df Residuals:            9882
Method:             MLE      Df Model:                  28
Date:              Mon, 30 May 2022      Pseudo R-squ.:    0.6747
Time:              21:11:17      Log-Likelihood:        -2114.9
converged:          True      LL-Null:                  -6501.2
Covariance Type:    nonrobust      LRR p-value:          0.000
=====
coef      std err      z      P>|z|      [0.025      0.975]
-----
Lat                0.0038      0.006      0.615      0.538      -0.008      0.016
Lng                0.0004      0.002      0.187      0.852      -0.004      0.005
Population         0.0000      0.003      -0.067      0.943      -0.006      0.006
Age               -0.0018      0.003      -0.667      0.505      -0.007      0.003
Household_Income  0.0000      8.6e-07      1.199      0.230      -6.64e-07      2.72e-06
VitD_levels       -0.6427      0.016      -40.229      0.000      -0.674      -0.611
VitD_supp         -0.0545      0.063      -0.869      0.385      -0.177      0.068
Soft_drink        0.0318      0.047      0.680      0.497      -0.040      0.124
HighBlood         -0.1729      0.126      -1.370      0.171      -0.420      0.074
Stroke            0.1231      0.098      1.315      0.189      -0.063      0.322
Complication_risk -0.2747      0.054      -5.132      0.000      -0.380      -0.170
Overweight        0.0174      0.066      0.261      0.794      -0.113      0.148
Arthritis         -0.2192      0.082      -2.685      0.007      -0.379      -0.059
Diabetes          -0.1658      0.087      -1.898      0.058      -0.337      0.005
Hyperlipidemia   -0.1488      0.084      -1.778      0.075      -0.313      0.015
BackPain         -0.1069      0.079      -1.345      0.179      -0.263      0.049
Anxiety           -0.1251      0.058      -2.141      0.032      -0.240      -0.011
Allergic_rhinitis -0.2638      0.080      -3.289      0.001      -0.421      -0.107
Reflux_esophagitis -0.1657      0.080      -2.076      0.038      -0.322      -0.009
Asthma           -0.0582      0.087      -1.272      0.203      -0.281      0.060
Daily_Average_Charges 0.0019      4.6e-05      42.265      0.000      0.002      0.002
Additional_charges 0.0018      1.12e-05      0.651      0.515      -1.46e-05      2.92e-05
Timely_Admission  0.0546      0.056      0.972      0.331      -0.055      0.165
Timely_Treatment  0.0929      0.052      1.782      0.075      -0.009      0.195
Timely_Visits     -0.0303      0.047      -0.641      0.522      -0.123      0.062
Reliability       0.0881      0.039      1.750      0.080      -0.008      0.144
Hours_Treatment  -0.0582      0.047      -1.248      0.212      -0.150      0.033
Courteous_Staff  -0.0264      0.042      -0.625      0.532      -0.109      0.056
Active_Listening  -0.0665      0.041      -1.626      0.104      -0.147      0.014
=====
```

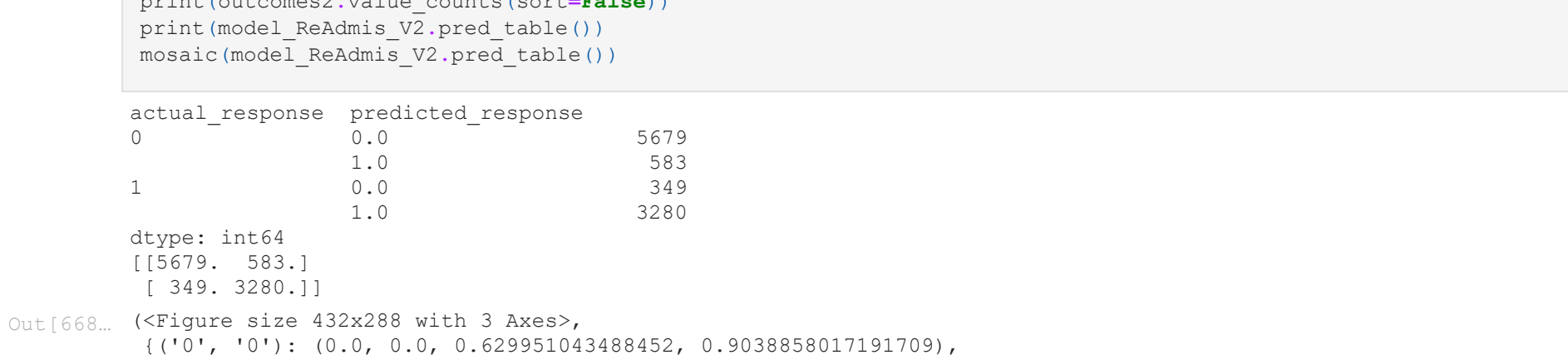
```
In [666]: from sklearn.linear_model import LogisticRegression
from sklearn.metrics import classification_report, confusion_matrix
from statsmodels.graphics.mosaicplot import mosaic
# Split the dataset into a training and testing set. Using an 80/20 testing/training split.
actual_response1 = medical_df_V2['Readmis']

predicted_response1 = np.round(model_Readmis_V1.predict())

outcome1 = pd.DataFrame({'actual_response': actual_response1, 'predicted_response': predicted_response1})
print(outcome1.value_counts(sort=False))
print(model_Readmis_V1.pred_table())
mosaic(model_Readmis_V1.pred_table())
TN1 = model_Readmis_V1.pred_table()[0,0]
FP1 = model_Readmis_V1.pred_table()[1,1]
FN1 = model_Readmis_V1.pred_table()[1,0]
PP1 = model_Readmis_V1.pred_table()[0,1]

sensitivity1 = TN1/(FN1+TN1)
accuracy1 = (TN1 + FP1)/(TN1 + FN1 + PP1 + TP1)
specificity1 = TN1/(TN1 + TP1)
print("The sensitivity of this model is: " + str(sensitivity1))
print("The accuracy of this model is: " + str(accuracy1))
print("The specificity of this model is: " + str(specificity1))
```

```
actual_response predicted_response
0                0.0                5676
1                1.0                586
1                0.0                353
dtype: int64
[[5676, 586,]
 [ 353, 3276,]]
The sensitivity of this model is: 0.900862107976949
The accuracy of this model is: 0.900862107976949
The specificity of this model is: 0.6340482937376541
```



```
In [667]: model_Readmis_V2 = sm.Logit(medical_df_V2[['Readmis'],medical_df_V2[['VitD_levels', 'Complication_risk',
'Arthritis', 'Anxiety', 'Allergic_rhinitis',
'Reflux_esophagitis',
'Additional_charges']]).fit()
print(model_Readmis_V2.params)
print(model_Readmis_V2.summary())
```

```
Optimization terminated successfully.
Current function value: 0.215313
Iterations: 9

VitD_levels       -0.636669
Complication_risk -0.259786
Arthritis         -0.204376
Anxiety           -0.110393
Allergic_rhinitis -0.248485
Reflux_esophagitis -0.154311
Daily_Average_Charges 0.001852
dtype: float64

Logit Regression Results
=====
Dep. Variable:      Readmis      No. Observations:      9891
Model:              Logit      Df Residuals:            9884
Method:             MLE      Df Model:                  6
Date:              Mon, 30 May 2022      Pseudo R-squ.:    0.6724
Time:              21:11:17      Log-Likelihood:        -2125.7
converged:          True      LL-Null:                  -6501.2
Covariance Type:    nonrobust      LRR p-value:          0.000
=====
coef      std err      z      P>|z|      [0.025      0.975]
-----
VitD_levels       -0.6367      0.016      -40.924      0.000      -0.667      -0.606
Complication_risk -0.2598      0.045      -5.783      0.000      -0.348      -0.172
Arthritis         -0.2044      0.080      -2.548      0.011      -0.362      -0.047
Anxiety           -0.1104      0.058      -1.915      0.055      -0.223      0.003
Allergic_rhinitis -0.2485      0.079      -3.150      0.002      -0.403      -0.094
Reflux_esophagitis -0.1543      0.039      -1.970      0.048      -0.307      -0.001
Daily_Average_Charges 0.0019      4.3e-05      42.670      0.000      0.002      0.002
=====
```

```
In [668]: from sklearn.linear_model import LogisticRegression
from sklearn.metrics import classification_report, confusion_matrix
from statsmodels.graphics.mosaicplot import mosaic
# Split the dataset into a training and testing set. Using an 80/20 testing/training split.
actual_response2 = medical_df_V2['Readmis']

predicted_response2 = np.round(model_Readmis_V2.predict())

outcome2 = pd.DataFrame({'actual_response': actual_response2, 'predicted_response': predicted_response2})
print(outcome2.value_counts(sort=False))
print(model_Readmis_V2.pred_table())
mosaic(model_Readmis_V2.pred_table())
TN2 = model_Readmis_V2.pred_table()[0,0]
FP2 = model_Readmis_V2.pred_table()[1,1]
FN2 = model_Readmis_V2.pred_table()[1,0]
PP2 = model_Readmis_V2.pred_table()[0,1]

sensitivity2 = TN2/(FN2+TN2)
accuracy2 = (TN2 + FP2)/(TN2 + FN2 + PP2 + TP2)
specificity2 = TN2/(TN2 + TP2)
print("The sensitivity of this model is: " + str(sensitivity2))
print("The accuracy of this model is: " + str(accuracy2))
print("The specificity of this model is: " + str(specificity2))
```

```
The sensitivity of this model is: 0.9038302562689446
The accuracy of this model is: 0.9057723248812051
The specificity of this model is: 0.633887710619958
```



```
In [ ]: 
```