

## ▼ Business Problem

As a data scientist working at Apollo 24/7, the ultimate goal is to tease out meaningful and actionable insights from Patient-level collected data. You can help Apollo hospitals to be more efficient, to influence diagnostic and treatment processes, to map the spread of a pandemic. One of the best examples of data scientists making a meaningful difference at a global level is in the response to the COVID-19 pandemic, where they have improved information collection, provided ongoing and accurate estimates of infection spread and health system demand, and assessed the effectiveness of government policies.

Objective:

- To find which variables are significant in predicting the reason for hospitalization for different regions.
- Analysing how well some variables like viral load, smoking, Severity Level describe the hospitalization charges.

Concept used:

- Uni-Variate Analysis
- Bi-Variate Analysis
- 2-sample t-test
- ANNOVA
- Chi-square

## ▼ Basic data exploration

```
1 import numpy as np
2 import pandas as pd
3 import seaborn as sns
4 import matplotlib.pyplot as plt
5
6 from scipy.stats import norm, chi2, f # distributions
7
8 from scipy.stats import ttest_ind, ttest_rel, f_oneway, kruskal # numerical vs categorical
9 from scipy.stats import chisquare, chi2_contingency # categorical features
10 from scipy.stats import pearsonr, spearmanr # numeric vs numeric
11
12 from scipy.stats import kstest # cdf
13
14 from statsmodels.distributions.empirical_distribution import ECDF
15 # Empirical CDF
```

```
1 df = pd.read_csv('/content/scaler_apollo_hospitals.csv')
```

```
1 df
```

	Unnamed: 0	age	sex	smoker	region	viral load	severity level	hospitalization charges
0	0	19	female	yes	southwest	9.30	0	42212
1	1	18	male	no	southeast	11.26	1	4314
2	2	28	male	no	southeast	11.00	3	11124
3	3	33	male	no	northwest	7.57	0	54961
4	4	32	male	no	northwest	9.63	0	9667
...	...	...	...	...	...	...	...	...
1333	1333	50	male	no	northwest	10.32	3	26501
1334	1334	18	female	no	northeast	10.64	0	5515
1335	1335	18	female	no	southeast	12.28	0	4075
1336	1336	21	female	no	southwest	8.60	0	5020
1337	1337	61	female	yes	northwest	9.69	0	72853

1338 rows × 8 columns

```
1 df.shape
```

```
(1338, 8)
```

```
1 df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1338 entries, 0 to 1337
Data columns (total 8 columns):
 #   Column            Non-Null Count  Dtype  
--- 
 0   Unnamed: 0        1338 non-null    int64  
 1   age               1338 non-null    int64  
 2   sex               1338 non-null    object  
 3   smoker            1338 non-null    object  
 4   region            1338 non-null    object  
 5   viral load        1338 non-null    float64 
 6   severity level   1338 non-null    int64  
 7   hospitalization charges 1338 non-null  int64  
dtypes: float64(1), int64(4), object(3)
memory usage: 83.8+ KB
```

```
1 df = df.drop(df.columns[0],axis=1)
```

```
1 df.isna().sum()
```

```
age          0
sex          0
smoker       0
region       0
viral load   0
severity level 0
hospitalization charges 0
dtype: int64
```

```
1 df.nunique()
```

```
age          47
sex          2
smoker       2
region       4
viral load   462
severity level 6
hospitalization charges 1320
dtype: int64
```

```
1 df.dtypes
```

```
age          int64
sex          object
smoker       object
region       object
viral load   float64
severity level int64
hospitalization charges int64
dtype: object
```

```
1 df['severity level']= df['severity level'].astype('object')
```

```
1 cat_cols = list(df.dtypes[df.dtypes == 'object'].index)
2 num_cols = list(df.dtypes[df.dtypes != 'object'].index)
3 cat_cols,num_cols
```

```
(['sex', 'smoker', 'region', 'severity level'],
 ['age', 'viral load', 'hospitalization charges'])
```

```
1 for i in cat_cols:
2   print(df[i].value_counts())
```

```
male      676
female    662
Name: sex, dtype: int64
no       1064
yes      274
Name: smoker, dtype: int64
southeast 364
southwest 325
northwest 325
northeast 324
Name: region, dtype: int64
0      574
1      324
2      240
3      157
4      25
5      18
Name: severity level, dtype: int64
```

## ▼ Univariate Data Analysis

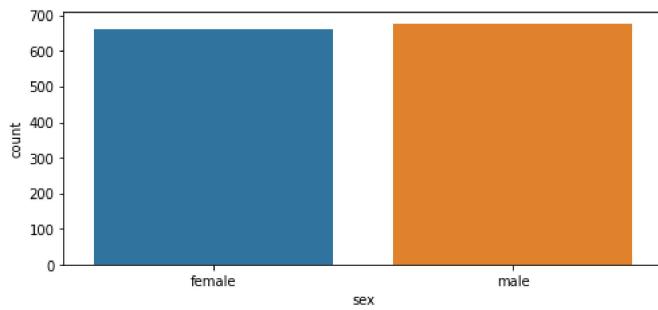
```
1 df2=df.copy()
2 cat_count = df2[cat_cols].melt().groupby(['variable', 'value'])[['value']].size().reset_index(name='Counts')
3 s = df2[cat_cols].melt().variable.value_counts()
4 cat_count['Percent'] = cat_count['Counts'].div(cat_count['variable'].map(s)).mul(100).round().astype('int')
5 cat_count.groupby(['variable', 'value','Counts','Percent']).first()
6
7 # there are just less than 5% patients with severity level 4 or 5.
```

variable	value	Counts	Percent
region	northeast	324	24
	northwest	325	24
	southeast	364	27
	southwest	325	24
severity level	0	574	43
	1	324	24
	2	240	18
	3	157	12
	4	25	2
	5	18	1
sex	female	662	49
	male	676	51
smoker	no	1064	80
	yes	274	20

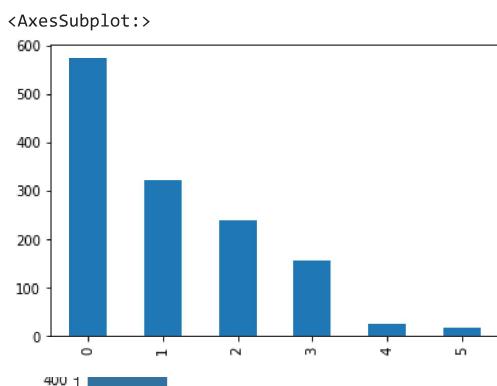
```
1 df.describe(include='object').T
```

	count	unique	top	freq
sex	1338	2	male	676
smoker	1338	2	no	1064
region	1338	4	southeast	364
severity level	1338	6	0	574

```
1 plt.figure(figsize = [8,16])
2 for i in range (len(cat_cols)):
3     plt.subplot(len(cat_cols),1, i+1)
4     sns.countplot(data=df, x=cat_cols[i])
```



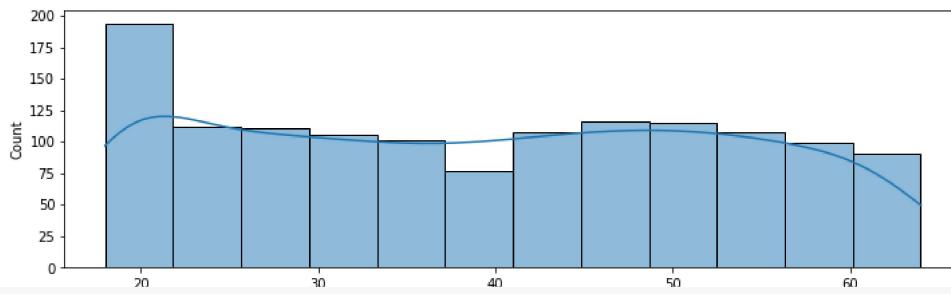
```
1 df['severity level'].value_counts().plot(kind='bar')
```



```
1 df.describe().T
```

	count	mean	std	min	25%	50%	75%	max	🔗
age	1338.0	39.207025	14.049960	18.00	27.0000	39.00	51.0000	64.00	
viral load	1338.0	10.221233	2.032796	5.32	8.7625	10.13	11.5675	17.71	
hospitalization charges	1338.0	33176.058296	30275.029296	2805.00	11851.0000	23455.00	41599.5000	159426.00	

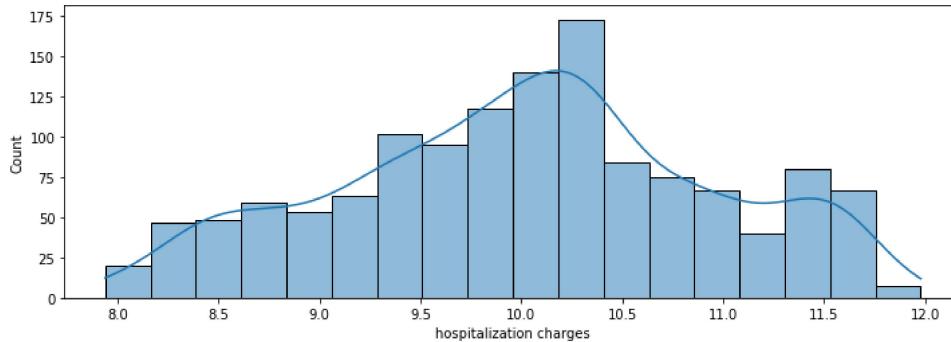
```
1 plt.figure(figsize = [12,12])
2 for i in range (len(num_cols)):
3     plt.subplot(len(num_cols),1, i+1)
4     sns.histplot(data=df, x=num_cols[i], kde=True)
```



```

1 plt.figure(figsize = [12,4])
2 sns.histplot(np.log(df['hospitalization charges']),kde=True)
3 plt.show()

```

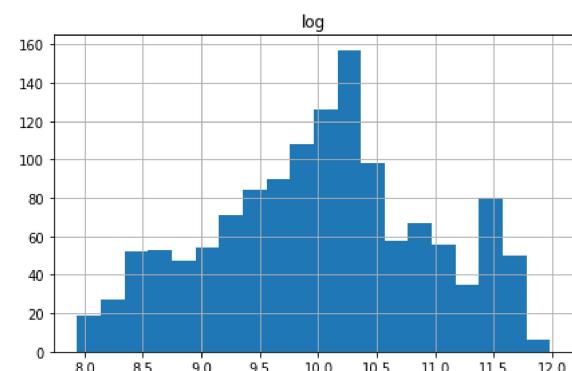
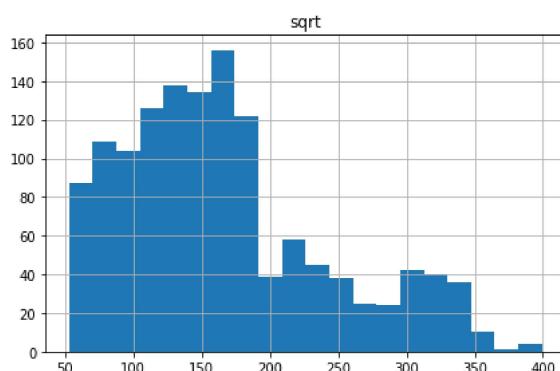


```

1 trans = df['hospitalization charges'].transform([np.sqrt, np.log])
2 trans.hist(bins=20, layout=(2,2),figsize=(16,10))
3 plt.suptitle('Transformed output')
4 plt.show()

```

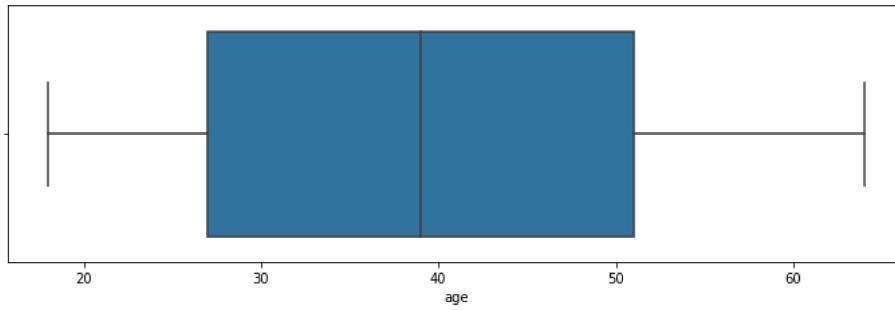
Transformed output



```

1 plt.figure(figsize = [12,12])
2 for i in range (len(num_cols)):
3     plt.subplot(len(num_cols),1, i+1)
4     sns.boxplot(data=df, x=num_cols[i])

```



```

1 for i in (num_cols):
2   print(i, round(df[i].skew(),1))
3
4 # windspeed distribution is right skewed, means it has some outliers in right.

```

```

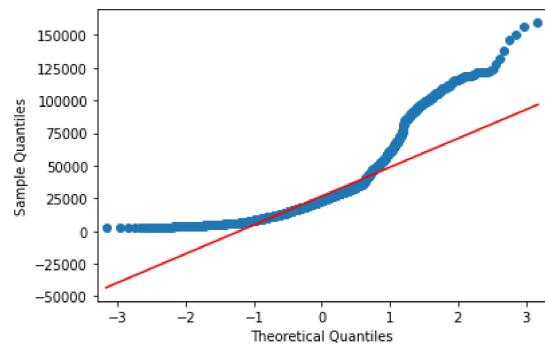
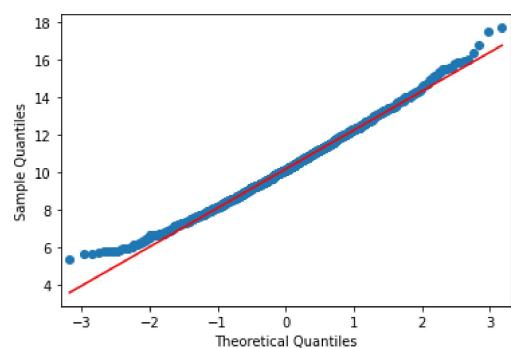
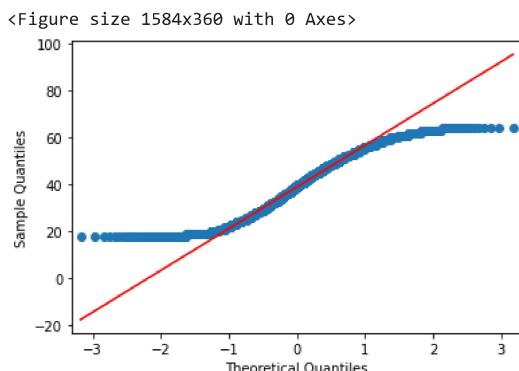
age 0.1
viral load 0.3
hospitalization charges 1.5

```

```

1 import statsmodels.api as sm
2 import matplotlib.pyplot as plt
3 plt.figure(figsize = [22,5])
4 for i in (num_cols):
5   sm.qqplot(df[i], line ='q')
6
7 # hospitalization charge values clearly do not follow the red line,
8 # which is an indication that they do not follow a normal distribution.

```



```

1 for i in num_cols:
2   R_whisker = np.percentile(df[i],75)+(np.percentile(df[i],75)-np.percentile(df[i],25))*1.5

```

```
3 outliers_percentage = df[df[i]>R_whisker].index.size/df.index.size*100
4 print(f'Percentage of outliers in {i} column = ', round(outliers_percentage,2))
Percentage of outliers in age column =  0.0
Percentage of outliers in viral load column =  0.67
Percentage of outliers in hospitalization charges column =  10.39
```

```
1 for i in num_cols:
2     q1=df[i].quantile(0.25)
3     q3=df[i].quantile(0.75)
4     IQR=q3-q1
5     outliers_percentage = df[((df[i]<(q1-1.5*IQR)) | (df[i]>(q3+1.5*IQR)))].index.size/df.index.size*100
6     print(f'Percentage of outliers in {i} column = ', round(outliers_percentage,2))

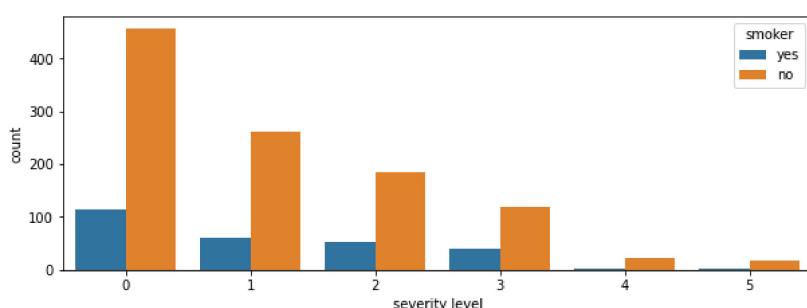
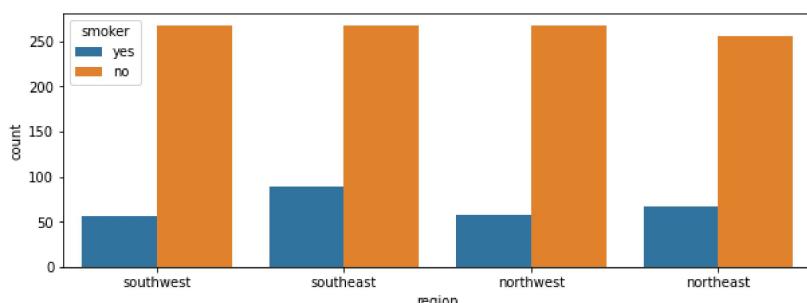
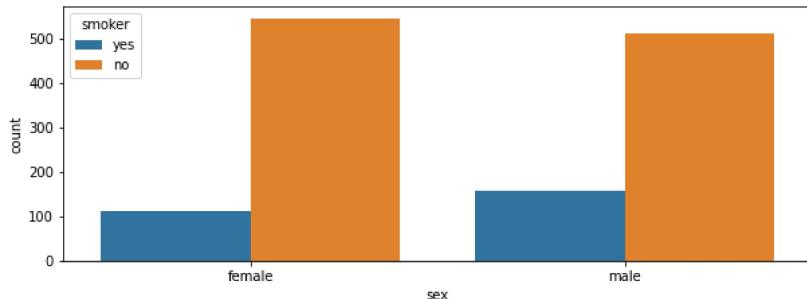
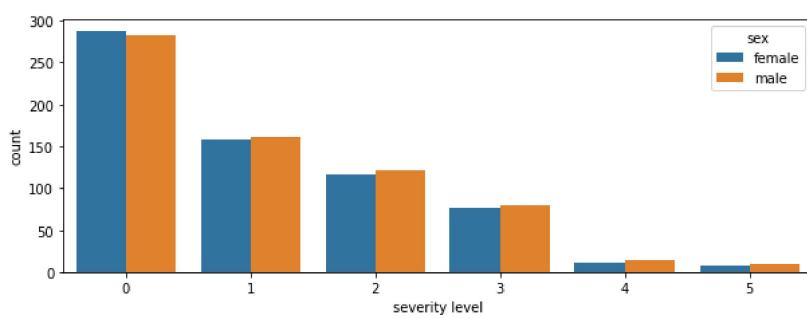
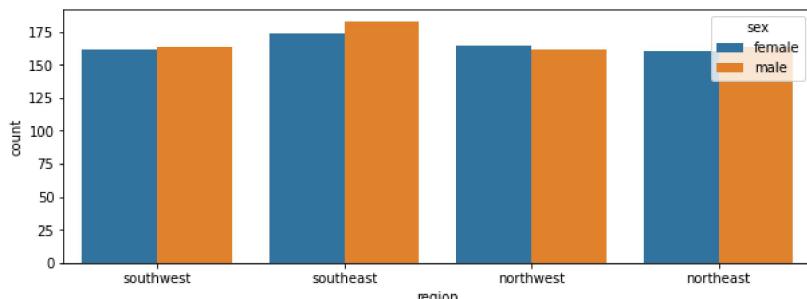
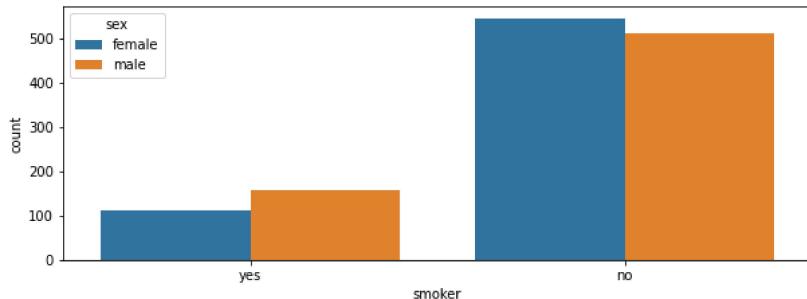
Percentage of outliers in age column =  0.0
Percentage of outliers in viral load column =  0.67
Percentage of outliers in hospitalization charges column =  10.39
```

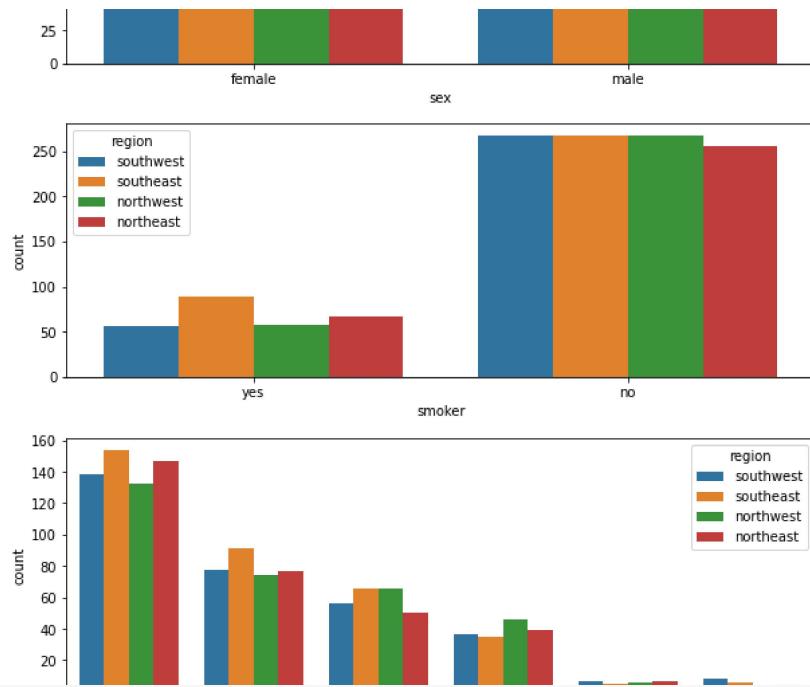
```
1 q1=df['viral load'].quantile(0.25)
2 q3=df['viral load'].quantile(0.75)
3 IQR=q3-q1
4 df = df[~((df['viral load']<(q1-1.5*IQR)) | (df['viral load']>(q3+1.5*IQR)))]
5
6 #·viral·load·column·has·less·than·1%·outliers so we will remove them using IQR method.
```

```
1 p99=df['hospitalization charges'].quantile(0.95)
2 df = df[~((df['hospitalization charges']<0) | (df['hospitalization charges']>p99))]
3
4 # hospitalization charges column has more than 10% outliers so we will just remove extreme 1% of them using percentile method.
```

## ▼ Bivariate Data Analysis

```
1 for j in range (len(cat_cols)):
2     for i in range (len(cat_cols)):
3         if i!=j:
4             plt.figure(figsize = [10,16])
5             plt.subplot(len(cat_cols), 1, i+1)
6             sns.countplot(data=df, x=cat_cols[i],hue=cat_cols[j])
7             plt.show()
```

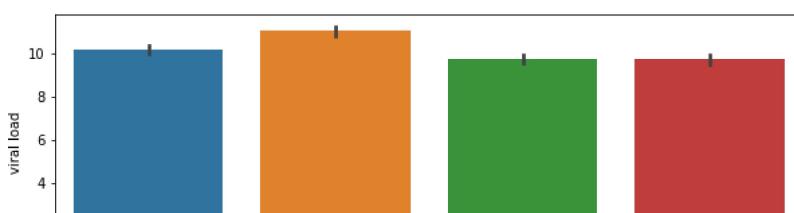
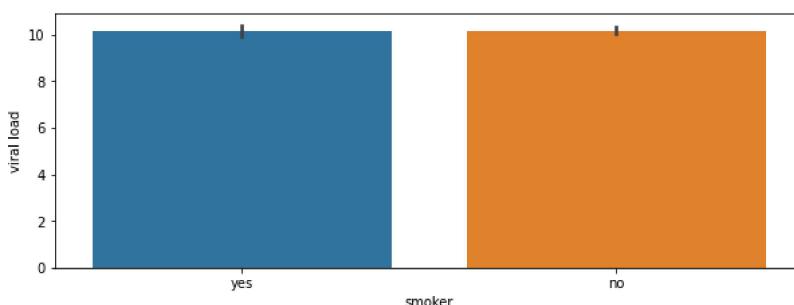
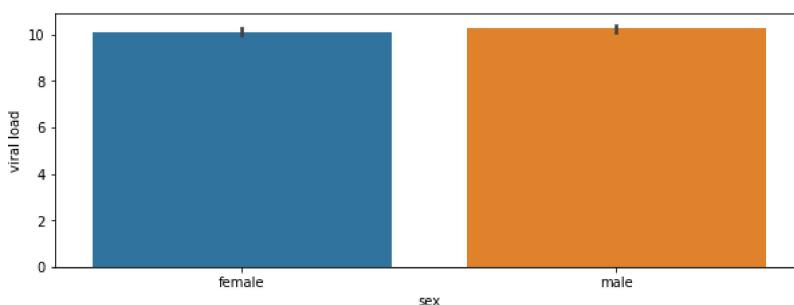
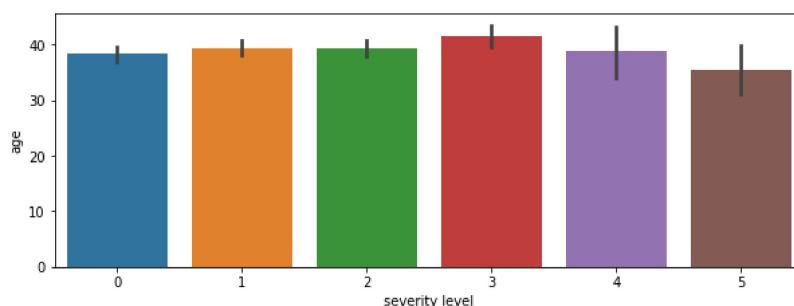
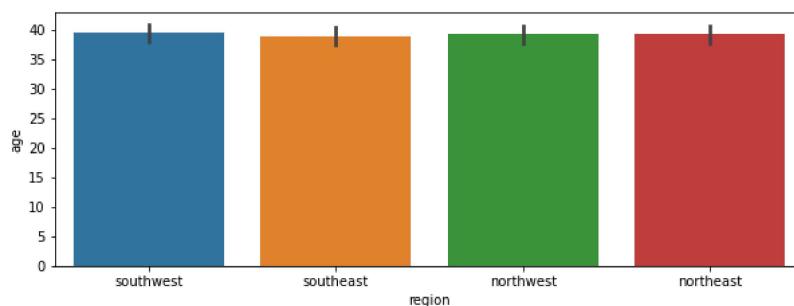
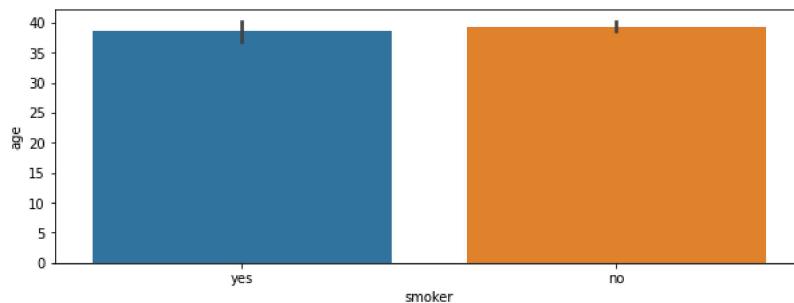
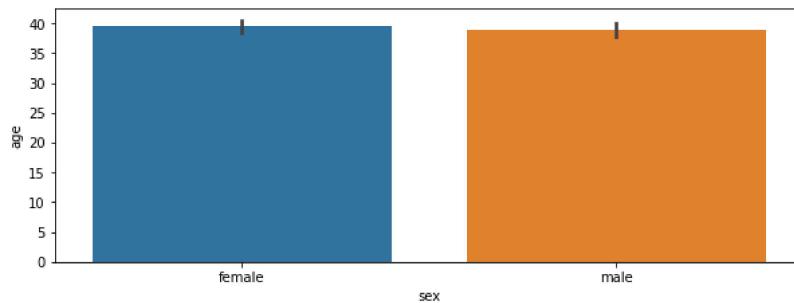


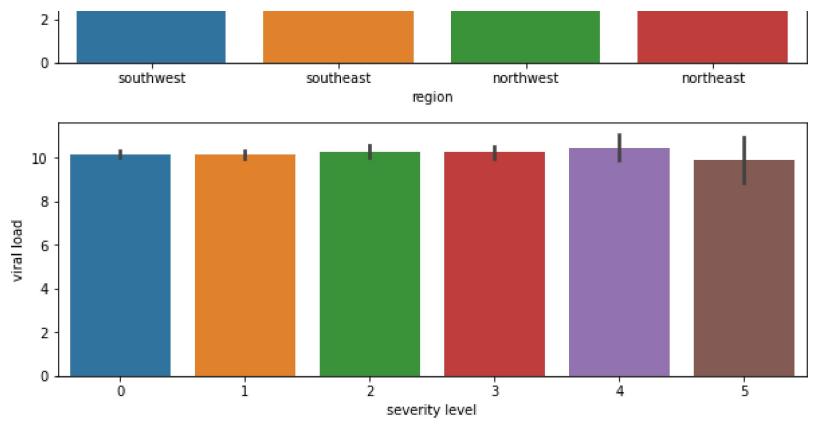


```

1 for j in range(len(num_cols)):
2     for i in range(len(cat_cols)):
3         plt.figure(figsize = [10,16])
4         plt.subplot(len(cat_cols), 1, i+1)
5         sns.barplot(data=df, x=cat_cols[i],y=num_cols[j])
6         plt.show()

```

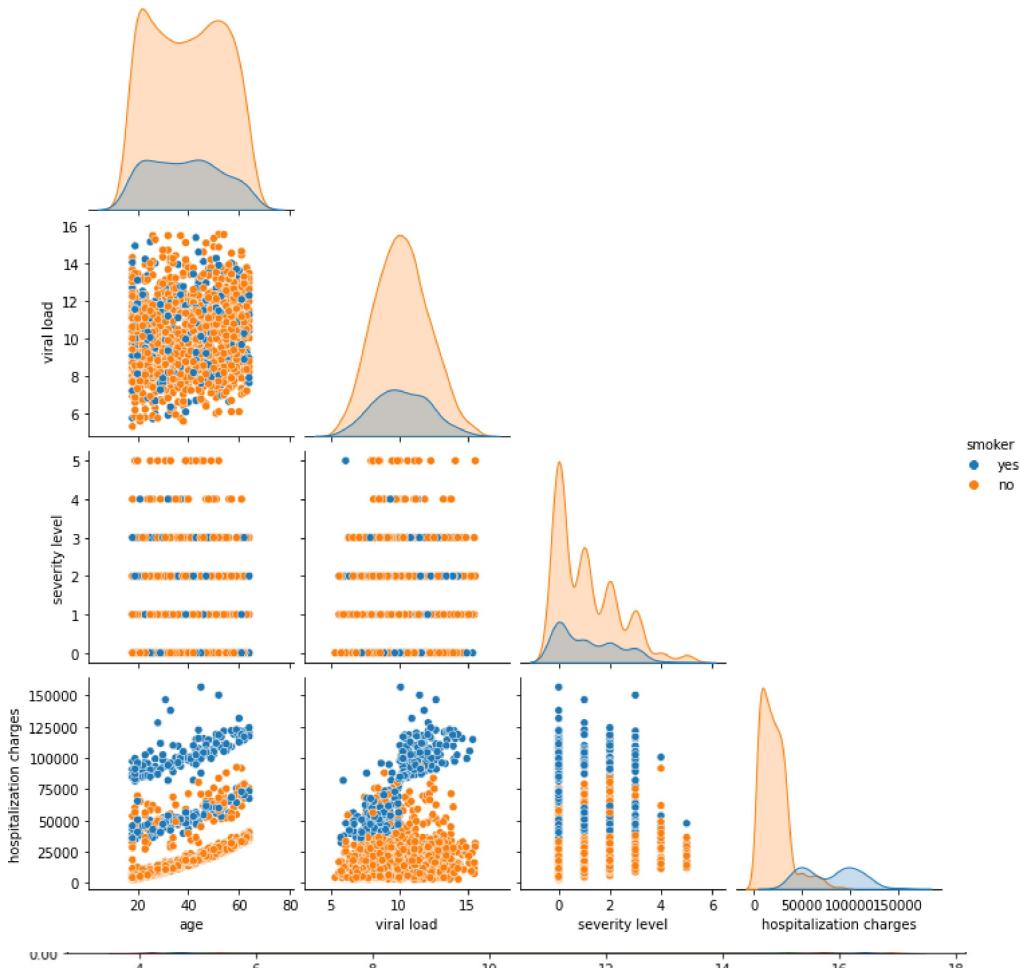




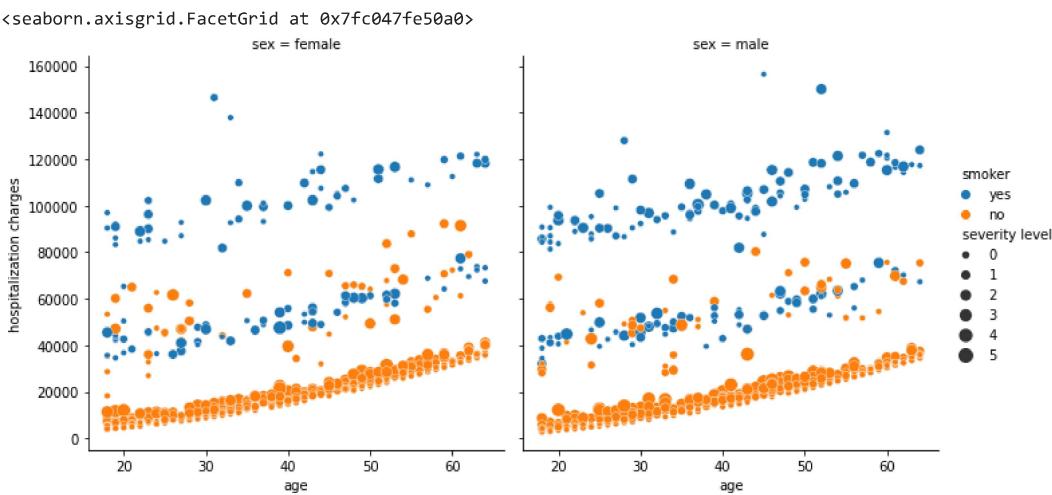
```
1 plt.figure(figsize = [12,20])
2 for j in range(len(cat_cols)):
3     plt.subplot(len(cat_cols), 1,j+1)
4     sns.kdeplot(data=df,x='viral load',hue=cat_cols[j],fill=True)
```

```
0.10
```

```
1 sns.pairplot(data = df, hue= 'smoker', corner=True)
2 plt.show()
```



```
1 sns.relplot(data=df, y="hospitalization charges", x="age", hue='smoker', size="severity level", sizes=(20, 100), col="sex")
2
3 # There are more male who is smoker.
```



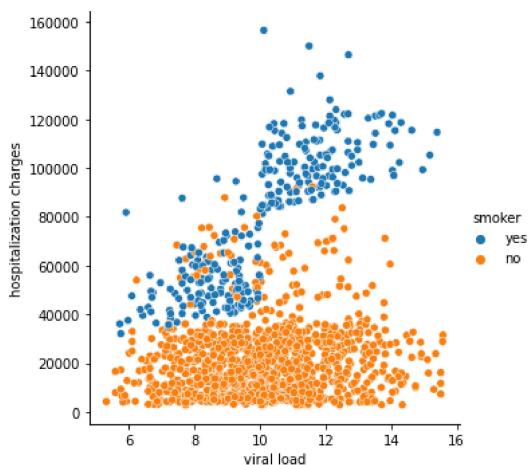
```
1 df[df['smoker']=='yes'].describe().T
```

	count	mean	std	min	25%	50%	75%	max
age	271.0	38.523985	13.932057	18.00	27.000	38.00	49.000	64.0
viral load	271.0	10.167970	2.013258	5.73	8.685	10.12	11.675	15.4
hospitalization charges	271.0	79588.354244	28461.379883	32074.00	51899.500	85758.00	102342.000	156482.0

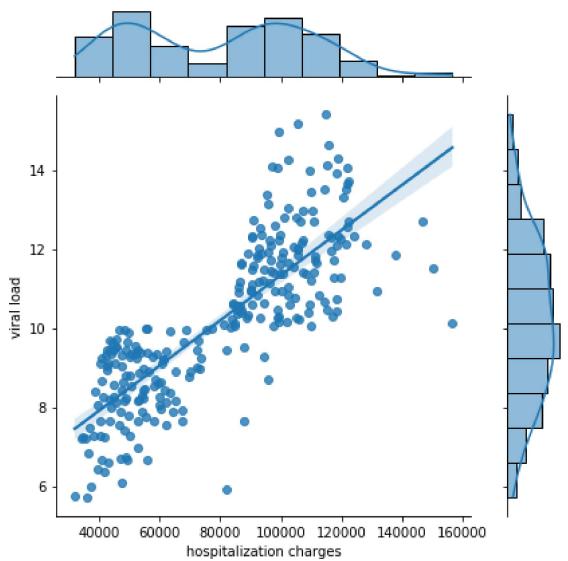
```
1 df[df['smoker']=='no'].describe().T
```

	count	mean	std	min	25%	50%	75%	max	🔗
<b>age</b>	1058.0	39.378072	14.078600	18.00	27.00	40.0	52.000	64.00	
<b>viral load</b>	1058.0	10.182108	1.964027	5.32	8.77	10.1	11.455	15.58	
<b>hospitalization charges</b>	1058.0	21105.423440	15007.155103	2805.00	9971.25	18353.5	28407.750	92277.00	

```
1 sns.relplot(data=df, y="hospitalization charges", x="viral load", hue='smoker')
2 plt.show()
```



```
1 sns.jointplot(x='hospitalization charges', y='viral load', data=df[df['smoker']=='yes'], kind='reg')
2 plt.tight_layout()
3 plt.show()
```

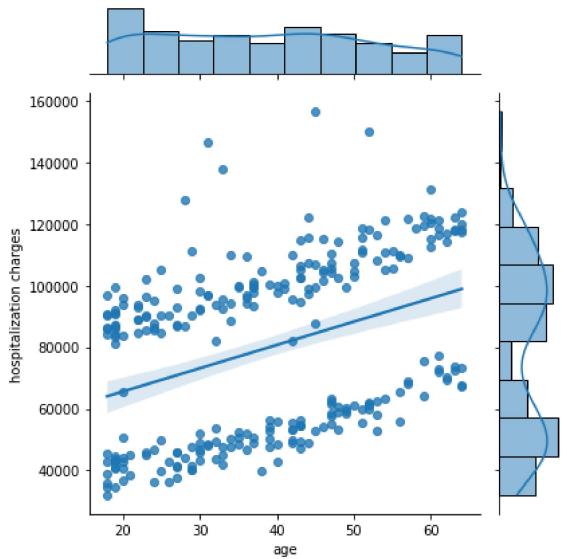


```
1 sns.jointplot(x='hospitalization charges', y='viral load', data=df[df['smoker']=='no'], kind='reg')
2 plt.tight_layout()
3 plt.show()
```

```

1 sns.jointplot(y='hospitalization charges', x='age', data=df[df['smoker']=='yes'], kind='reg')
2 plt.tight_layout()
3 plt.show()

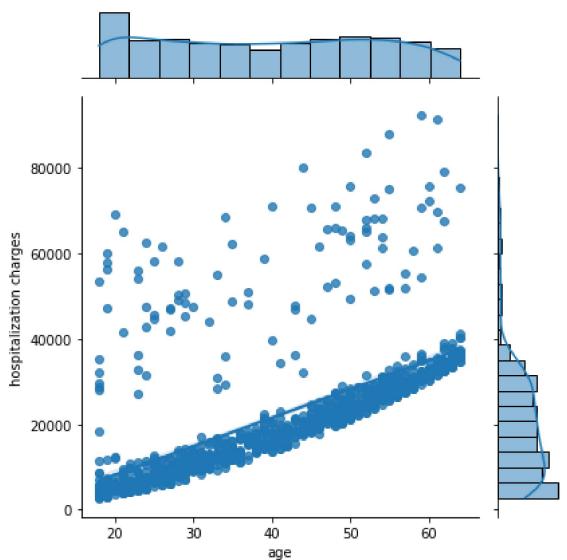
```



```

1 sns.jointplot(y='hospitalization charges', x='age', data=df[df['smoker']=='no'], kind='reg')
2 plt.tight_layout()
3 plt.show()

```



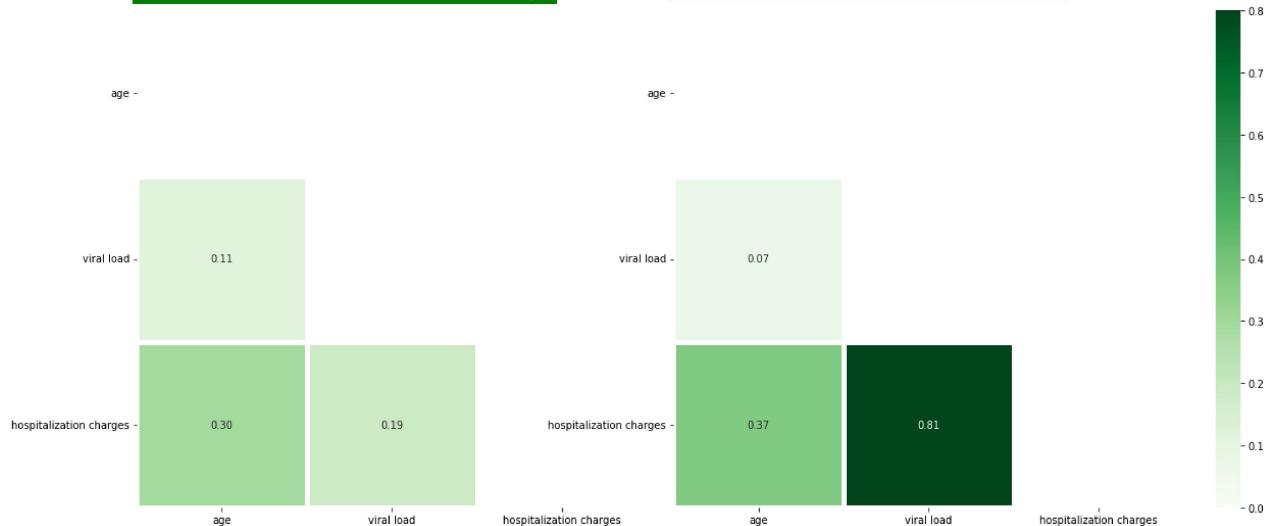
```

1 fig=plt.figure(figsize = [18,8])
2 df_smoker = df.copy()[df['smoker']=='yes']
3
4 plt.subplot(1, 2, 1)
5 mask = np.triu(np.ones_like(df.corr()))
6 sns.heatmap(df.corr().round(2), cmap= "Greens", annot=True, mask=mask, fmt=".2f", linewidths=5, vmin=0, vmax=0.8)
7 plt.tight_layout()
8 plt.title('correlation heatmap', color='w',fontsize=40, fontweight = 'normal',backgroundcolor = 'g', pad = 20, loc='left')
9 plt.yticks(rotation = 0)
10
11 plt.subplot(1, 2, 2)
12 mask = np.triu(np.ones_like(df_smoker.corr()))
13 sns.heatmap(df_smoker.corr().round(2), cmap= "Greens", annot=True, mask=mask, fmt=".2f", linewidths=5, vmin=0, vmax=0.8)
14 plt.tight_layout()
15 plt.title('correlation heatmap for smokers=Yes', color='w',fontsize=20, fontweight = 'normal',backgroundcolor = 'r', pad = 20, loc='left')
16 plt.yticks(rotation = 0)
17 plt.show()
18
19 # In general there is no correlation greater than 0.5 percent between any two variables,
20 # But if we see correlation for onle smoker patients then,
21 # we can observe that there is highly correlation between viral load and hospitalization charge
22 # and also there is correlation between age and hospitalization charge

```

# correlation heatmap

correlation heatmap for smokers=Yes



## ▼ Missing values treatment & Outlier treatment

```
1 # dataset has no missing values, checked in first section.  
2 # Outlier has treated in second section with IQR method (for few outliers) and with percentile methods (for many outliers).
```

## ▼ Hypothesis Testing

```
1 def Hypothesis_testing(Samples,alpha,alternative,mu,crosstab):  
2  
3     if alternative=='two-sided':  
4         H0 = 'means are equal'  
5     elif alternative== 'less':  
6         H0 = 'Sample1>=Sample2'  
7     elif alternative=='greater':  
8         H0 = 'Sample1<=Sample2'  
9     else:  
10        print('Check the alternative input')  
11        exit()  
12    H0_chi_stat = 'obs and exp values are same, proportions are similar'  
13    H0_chi2_contingency = 'variables are independent, likelihoods are similar'  
14    if not all(Samples):  
15        print('Check Samples detail input')  
16        exit()  
17  
18    if type(Samples[0][0])==list:  
19  
20        if sum(Samples[0][0])-sum(Samples[0][1])<0.001:  
21            print('Chi square test for goodness of fit (obs,exp)')  
22            from scipy.stats import chisquare  
23            test, p_val = chisquare(*Sample1)  
24            print("p-value is: " + str(p_val))  
25            print('null hypothesis: ',H0_chi_stat)  
26            if p_val <= alpha:  
27                print('We can reject the null hypothesis')  
28            else:  
29                print('We can accept the null hypothesis')  
30  
31    else:  
32        print('Chi square test for independence (contingency)')  
33        from scipy.stats import chi2_contingency  
34        test, p_val, dof, expected_val = chi2_contingency(Sample1)  
35        print("p-value is: " + str(p_val))  
36        print('null hypothesis: ',H0_chi2_contingency)  
37        if p_val <= alpha:  
38            print('We can reject the null hypothesis')  
39        else:  
40            print('We can accept the null hypothesis')  
41
```

```

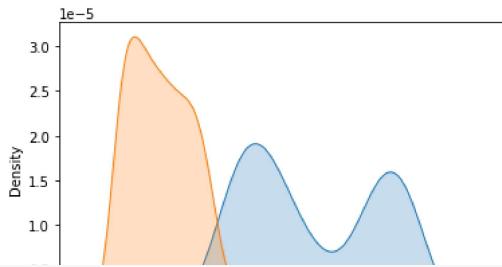
42 elif len(Samples)>2:
43     print(len(Samples), 'Sample', 'Anova f test')
44     from scipy.stats import f_oneway
45     f_test, p_val = f_oneway(*Samples)
46     print("p-value is: " + str(p_val))
47     print('null hypothesis: ',H0)
48     if p_val < alpha:
49         print("We can reject the null hypothesis")
50     else:
51         print("We can accept the null hypothesis")
52
53 elif len(Samples)>1:
54     print(len(Samples), 'Sample', 't-test')
55     from scipy.stats import ttest_ind
56     t_test,p_val = ttest_ind(*Samples)
57     print("p-value of independent t-test is: ", p_val)
58     print('null hypothesis: ',H0)
59     if p_val < alpha:
60         print("We can reject the null hypothesis")
61     else:
62         print("We can accept the null hypothesis")
63
64 if len(Samples[0])==len(Samples[1]):
65     from scipy.stats import ttest_rel
66     t_test, p_val = ttest_rel(*Samples, alternative=alternative)
67     print("p-value of the paired t-test is: ", p_val)
68     print('null hypothesis: ',H0)
69     if p_val < 0.05:
70         print("We can reject the null hypothesis")
71     else:
72         print("We can accept the null hypothesis")
73
74
75 elif len(Samples)==1:
76     print(len(Samples), 'Sample', 't test')
77     from scipy.stats import ttest_1samp
78     t_stat, p_val = ttest_1samp(Samples[0],mu, alternative=alternative)
79     print("P-value is: ", p_val)
80     print('null hypothesis: ',H0)
81     if p_val < 0.05:
82         print(" We can reject the null hypothesis")
83     else:
84         print("We can accept the null hypothesis")
85
86 import matplotlib.pyplot as plt
87 import seaborn as sns
88 if type(Sample1[0])==list:
89     crosstab.plot(kind="bar",figsize=(8,5))
90     print(crosstab)
91 else:
92     for samp in (Samples):
93         sns.kdeplot(samp,fill=True)
94     plt.show()
95
96 import numpy as np
97 for i in range(len(Samples)):
98     if alternative=='less':
99         print(f'mean of Sample{i+1}:',np.mean(Samples[i]).round(2), '      95% Range [5%-100%] :', [np.percentile(Samples[i],5).round(2),np.percentile(Samples[i],100).round(2)])
100    elif alternative=='greater':
101        print(f'mean of Sample{i+1}:',np.mean(Samples[i]).round(2), '      95% Range [0%-95%] :', [np.min(Samples[i]).round(2),np.percentile(Samples[i],95).round(2)])
102    else:
103        print(f'mean of Sample{i+1}:',np.mean(Samples[i]).round(2), '      95% Range [2.5%-97.7%] :', [np.percentile(Samples[i],2.5).round(2),np.percentile(Samples[i],97.7).round(2)])
1 # null hypothesis = hospitalization charge of people who do smoking is lesser than those who don't.
2 # alternative hypothesis = hospitalization charge of people who do smoking is greater than those who don't.
3
4 Sample1=df[df['smoker']=='yes']['hospitalization charges'].tolist()
5 Sample2=df[df['smoker']=='no']['hospitalization charges'].tolist()
6 Samples=[Sample1,Sample2]
7 alpha = 0.05
8 alternative = 'greater' # 'two-sided' # 'less' , 'greater'
9 mu = None
10 crosstab=None
11 Hypothesis_testing(Samples,alpha,alternative,mu,crosstab)
12
13 # conclusion: hospitalization charge of people who do smoking is greater than those who don't.

```

```

2 Sample t-test
p-value of independent t-test is: 5.209104236115724e-207
null hypothesis: Sample1<=Sample2
We can reject the null hypothesis

```



```

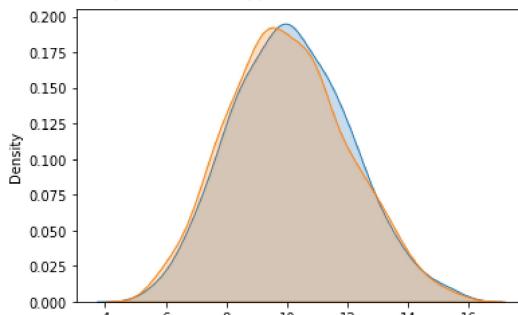
1 # null hypothesis = viral load of females is same as males.
2 # null hypothesis = viral load of females is different as males.
3
4 Sample1=df[df['sex']=='male']['viral load'].tolist()
5 Sample2=df[df['sex']=='female']['viral load'].tolist()
6 Samples=[Sample1,Sample2]
7 alpha = 0.05
8 alternative = 'two-sided' # 'two-sided' # 'less' , 'greater'
9 mu = None
10 crosstab=None
11 Hypothesis_testing(Samples,alpha,alternative,mu,crosstab)
12
13 # conclusion: viral load of females is same as males

```

```

2 Sample t-test
p-value of independent t-test is: 0.3144637521964963
null hypothesis: means are equal
We can accept the null hypothesis

```



```

mean of Sample1: 10.13      95% Range [2.5%-97.7%] : [6.64, 14.07]
mean of Sample2: 10.02      95% Range [2.5%-97.7%] : [6.47, 13.98]

1 # null hypothesis = proportion of smoking significantly similar across different regions.
2 # alternative hypothesis = proportion of smoking significantly different across different regions.
3
4 import pandas as pd
5 crosstab = pd.crosstab(index=df['smoker'], columns=df['region'])    # Sample1 = crosstab.values.tolsit()
6 Sample1=crosstab.values.tolist()
7 Samples=[Sample1]
8 mu = None
9 alpha = 0.05
10 alternative = 'two-sided' # 'less' , 'greater'
11 Hypothesis_testing(Samples,alpha,alternative,mu,crosstab)
12
13 # Conclusion: proportion of smoking significantly similar across different regions.

```

```

Chi square test for indepedence (contingency)
p-value is: 0.558169723528082
null hypothesis: variables are independet, liklihoods are similar
We can accept the null hypothesis
region northeast northwest southeast southwest
smoker

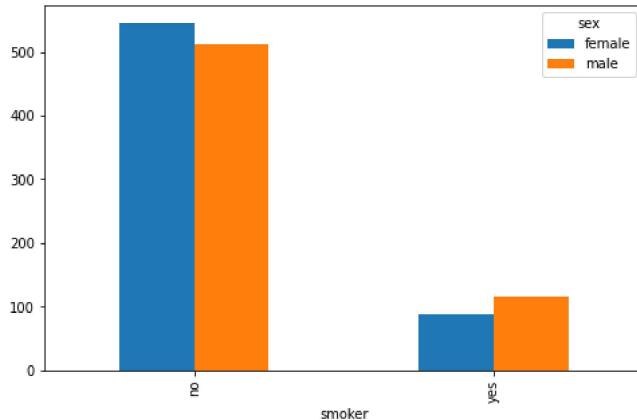
1 # null hypothesis = proportion of smoking significantly similar across different gender.
2 # alternative hypothesis = proportion of smoking significantly different across different gender.
3
4 import pandas as pd
5 crosstab = pd.crosstab(index=df['smoker'], columns=df['sex'])    # Sample1 = crosstab.values.tolsit()
6 Sample1=crosstab.values.tolist()
7 Samples=[Sample1]
8 mu = None
9 alpha = 0.05
10 alternative = 'two-sided' # 'less' , 'greater'
11 Hypothesis_testing(Samples,alpha,alternative,mu,crosstab)
12
13 # conclusion: we can say that there are more males who are smoker.

```

```

Chi square test for indepedence (contingency)
p-value is: 0.03245297207395062
null hypothesis: variables are independet, liklihoods are similar
We can reject the null hypothesis
sex      female   male
smoker
no        546     512
yes       88      116

```



```

1 # null hypothesis = proportion of smoking significantly similar across different severity level.
2 # alternative hypothesis = proportion of smoking significantly different across different severity level.
3
4 import pandas as pd
5 crosstab = pd.crosstab(index=df['smoker'], columns=df['severity level'])    # Sample1 = crosstab.values.tolsit()
6 Sample1=crosstab.values.tolist()
7 Samples=[Sample1]
8 mu = None
9 alpha = 0.05
10 alternative = 'two-sided' # 'less' , 'greater'
11 Hypothesis_testing(Samples,alpha,alternative,mu,crosstab)
12
13 # conclusion: we can not say that people who smoke has higher severity level.

```

```

Chi square test for indepedence (contingency)
p-value is: 0.5963357865679622
null hypothesis: variables are independet, liklihoods are similar
We can accept the null hypothesis

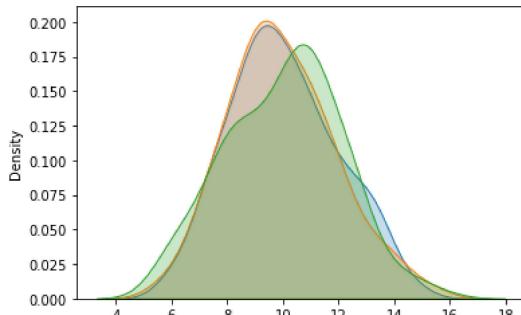
1 # null hypothesis = the mean viral load of women with 0 Severity level, 1 Severity level, and 2 Severity level are the same.
2 # null hypothesis = the mean viral load of women with 0 Severity level, 1 Severity level, and 2 Severity level are different.
3
4 cf=df['sex']=='female'
5 c0=df['severity level']==0
6 c1=df['severity level']==1
7 c2=df['severity level']==2
8
9 Sample1=df[(cf) & (c0)]['viral load'].tolist()
10 Sample2=df[(cf) & (c1)]['viral load'].tolist()
11 Sample3=df[(cf) & (c2)]['viral load'].tolist()
12 Samples=[Sample1,Sample2,Sample3]
13 alpha = 0.05
14 alternative = 'two-sided' # 'two-sided' # 'less' , 'greater'
15 mu = None
16 crosstab=None
17 Hypothesis_testing(Samples,alpha,alternative,mu,crosstab)
18
19 # conclusion: the mean viral load of women with 0 Severity level, 1 Severity level, and 2 Severity level are the same.

```

```

3 Sample Anova f test
p-value is: 0.90641731132652
null hypothesis: means are equal
We can accept the null hypothesis

```



```

mean of Sample1: 10.02      95% Range [2.5%-97.7%] : [6.73, 13.79]
mean of Sample2: 9.94      95% Range [2.5%-97.7%] : [6.68, 14.06]
mean of Sample3: 10.03      95% Range [2.5%-97.7%] : [5.97, 14.26]

```

```

1 # hospitalization charges column is not normally distributed
2 # we can perform nonparametric statistical hypothesis tests
3
4 # null hypothesis = hospitalization charge of people who do smoking is lesser than those who don't.
5 # alternative hypothesis = hospitalization charge of people who do smoking is greater than those who don't.
6
7 from scipy.stats import mannwhitneyu
8 sample1 = df[df['smoker']=='yes']['hospitalization charges'].tolist()
9 sample2 = df[df['smoker']=='no']['hospitalization charges'].tolist()
10 stat, p = mannwhitneyu(sample1, sample2)
11 print('stat=%3f, p=%3f' % (stat, p))
12 if p > 0.05:
13     print('Probably the same distribution, accept null hypothesis.')
14 else:
15     print('Probably different distributions, reject null hypothesis.')

stat=208428.500, p=0.000
Probably different distributions, reject null hypothesis.

```

```

1 # hospitalization charges column is not normally distributed
2 # we can perform nonparametric statistical sypothesis tests
3
4 # null hypothesis = hospitalization charges are similar across different severity level.
5 # alternative hypothesis = hospitalization charges are different across different severity level.
6
7 from scipy.stats import mannwhitneyu
8 sample1 = df[df['severity level']==0]['hospitalization charges'].tolist()
9 sample2 = df[df['severity level']==1]['hospitalization charges'].tolist()
10 sample3 = df[df['severity level']==2]['hospitalization charges'].tolist()
11 sample4 = df[df['severity level']==3]['hospitalization charges'].tolist()
12 stat, p = kruskal(sample1, sample2, sample3, sample4)
13 print('stat=%3f, p=%3f' % (stat, p))
14 if p > 0.05:
15     print('Probably the same distribution, accept null hypothesis.')
16 else:
17     print('Probably different distributions, reject null hypothesis.')
18

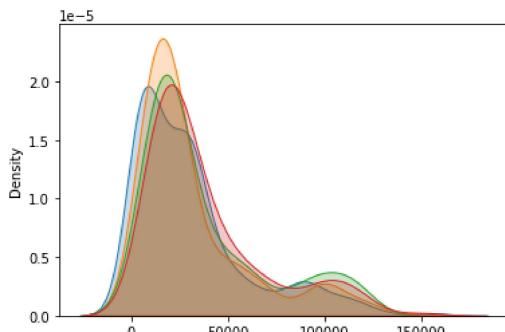
```

```
19 # Conclusion: hospitalization charges are different across different severity level.
```

```
stat=27.000, p=0.000
Probably different distributions, reject null hypothesis.
```

```
1 # parametric test is not applicable hospitalization charges,
2 # but we are just comparing the result with previous non parametric test.
3
4 sample1 = df[df['severity level']==0]['hospitalization charges'].tolist()
5 sample2 = df[df['severity level']==1]['hospitalization charges'].tolist()
6 sample3 = df[df['severity level']==2]['hospitalization charges'].tolist()
7 sample4 = df[df['severity level']==3]['hospitalization charges'].tolist()
8 Samples = [sample1, sample2, sample3, sample4]
9 alpha = 0.05
10 alternative = 'greater' # 'two-sided' # 'less' , 'greater'
11 mu = None
12 crosstab=None
13 Hypothesis_testing(Samples,alpha,alternative,mu,crosstab)
14
15 # Conclusion: hospitalization charges are different across different severity level.
```

```
4 Sample Anova f test
p-value is: 0.0035397453106568765
null hypothesis: Sample1<=Sample2
We can reject the null hypothesis
```



```
mean of Sample1: 30914.94      95% Range [0%-95%] : [2805, 97698.45]
mean of Sample2: 31827.94      95% Range [0%-95%] : [4278, 102023.6]
mean of Sample3: 37683.91      95% Range [0%-95%] : [5760, 110514.25]
mean of Sample4: 38388.31      95% Range [0%-95%] : [8608, 106125.4]
```

## Business Insights and Recommendations

```
1 # There are less patients who are smoker.
2 # 80% patients does not smoke.
3 # There are more number of males who are smoker than females.
4 # there are just less than 5% patients with higher severity level 4 or 5.
5 # there 43% patients with very low severity level 0.
6 # most patients belong to southeast region.
7 # maximum hospitalization charge has paid by patient is 159426.
8 # minimum hospitalization charge has paid by patient is 2805.
9 # age column is showing continuous uniform distribution.
10 # viral load column is showing continuous normal distribution.
11 # hospitalization charge column is showing lognormal distribution.
12 # dataset conatins no missing values.
13 # age column has no outliers present in dataset.
14 # viral load column has less than 1% outliers, we removed.
15 # hospitalization charge column has more than 10% outliers, removed 1%.
16 # In general there is no correlation between any two variables,
17 # But if we see correlation for onle smoker patients then,
18 # we can observe that there is highly correlation between viral load and hospitalization charge
19 # and also there is correlation between age and hospitalization charge
20 # using parametric test we concluded that:
21 # hospitalization charge of people who do smoking is greater than those who don't.
22 # proportion of smoking significantly similar across different regions.
23 # viral load of females is same as males.
24 # we can say that there are more males who are smoker.
25 # we can not say that people who smoke has higher severity level.
26 # the mean viral load of women with 0 Severity level, 1 Severity level, and 2 Severity level are the same.
27 # hospitalization charges are different across different severity level.
```

---

✓ 0s completed at 11:24PM

