

Apollo Hospitals was established in 1983, renowned as the architect of modern healthcare in India. As the nation's first corporate hospital, Apollo Hospitals is acclaimed for pioneering the private healthcare revolution in the country.

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.

Lets Analyse the data

```
## Import libraries
import warnings
warnings.filterwarnings("ignore")
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from IPython.core.display import display, HTML
display(HTML("<style>.container { width:100% !important; }</style>"))
from IPython.core.display import display, HTML
display(HTML("<style>.container { width:100% !important; }</style>"))
import matplotlib_inline
matplotlib_inline.backend_inline.set_matplotlib_formats('svg')
from scipy.stats import kstest, ks_2samp
from sklearn.preprocessing import LabelEncoder
from scipy import stats
from statsmodels.graphics.gcfplots import qqplot_2samples
plot_cmmap = "gist_heat"
plot_color = "#B3C2E6"
```

```
In [64]: ## reading dataset
df = pd.read_csv("C:\Users\shahil.bansal\Desktop\scaler_apollo_hospitals.csv")
```

```
In [3]: ## getting the shape of dataset
df.shape
```

```
Out[3]: (1338, 7)
```

There are 1338 rows and 7 columns in the data

```
In [ ]:
```

```
In [4]: ## checking info of the data
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1338 entries, 0 to 1337
Data columns (total 7 columns):
 #   Column      Non-Null Count  Dtype
---  --
 0   age         1338 non-null   int64
 1   sex         1338 non-null   object
 2   smoker      1338 non-null   object
 3   region      1338 non-null   object
 4   viral load  1338 non-null   float64
 5   severity level  1338 non-null   int64
 6   hospitalization charges  1338 non-null   int64
dtypes: float64(1), int64(3), object(3)
memory usage: 73.3+ KB
```

We can see that there are 1338 data points and 7 features and there are no null values and there are both categorical and numerical values in the data

```
In [ ]:
```

```
In [5]: ## Describe the data
df.describe()
```

```
Out[5]:
```

	age	viral load	severity level	hospitalization charges
count	1338.000000	1338.000000	1338.000000	1338.000000
mean	39.207025	10.221233	1.094918	33176.056296
std	14.049960	8.202796	1.205493	30275.029296
min	18.000000	5.320000	0.000000	2895.000000
25%	27.000000	6.762500	0.000000	11851.000000
50%	39.000000	10.130000	1.000000	23455.000000
75%	51.000000	11.567500	2.000000	41999.000000
max	64.000000	17.710000	5.000000	159426.000000

we can see there are not much outliers there are some outliers in hospitalization charges that we will handle in further analysis.

```
In [ ]:
```

```
In [6]: ## checking for null values
df.isnull().sum()
```

```
Out[6]:
age          0
sex          0
smoker       0
region       0
viral load   0
severity level  0
hospitalization charges  0
dtype: int64
```

We can see there are no null values in the data

```
In [ ]:
```

```
In [7]: ## value count of data
cat_cols = ['sex', 'smoker', 'region', 'severity level']
df[cat_cols].apply(lambda x: x.value_counts()).stack()
```

```
Out[7]:
sex          female    662.0
           male      676.0
smoker      no      1064.0
           yes       274.0
region      southeast  324.0
           northwest  325.0
           southeast  364.0
           southwest  325.0
           southwest  324.0
severity level  1       324.0
              2       240.0
              3       197.0
              4       25.0
              5       18.0
dtype: float64
```

data looks good there are no miscellaneous values

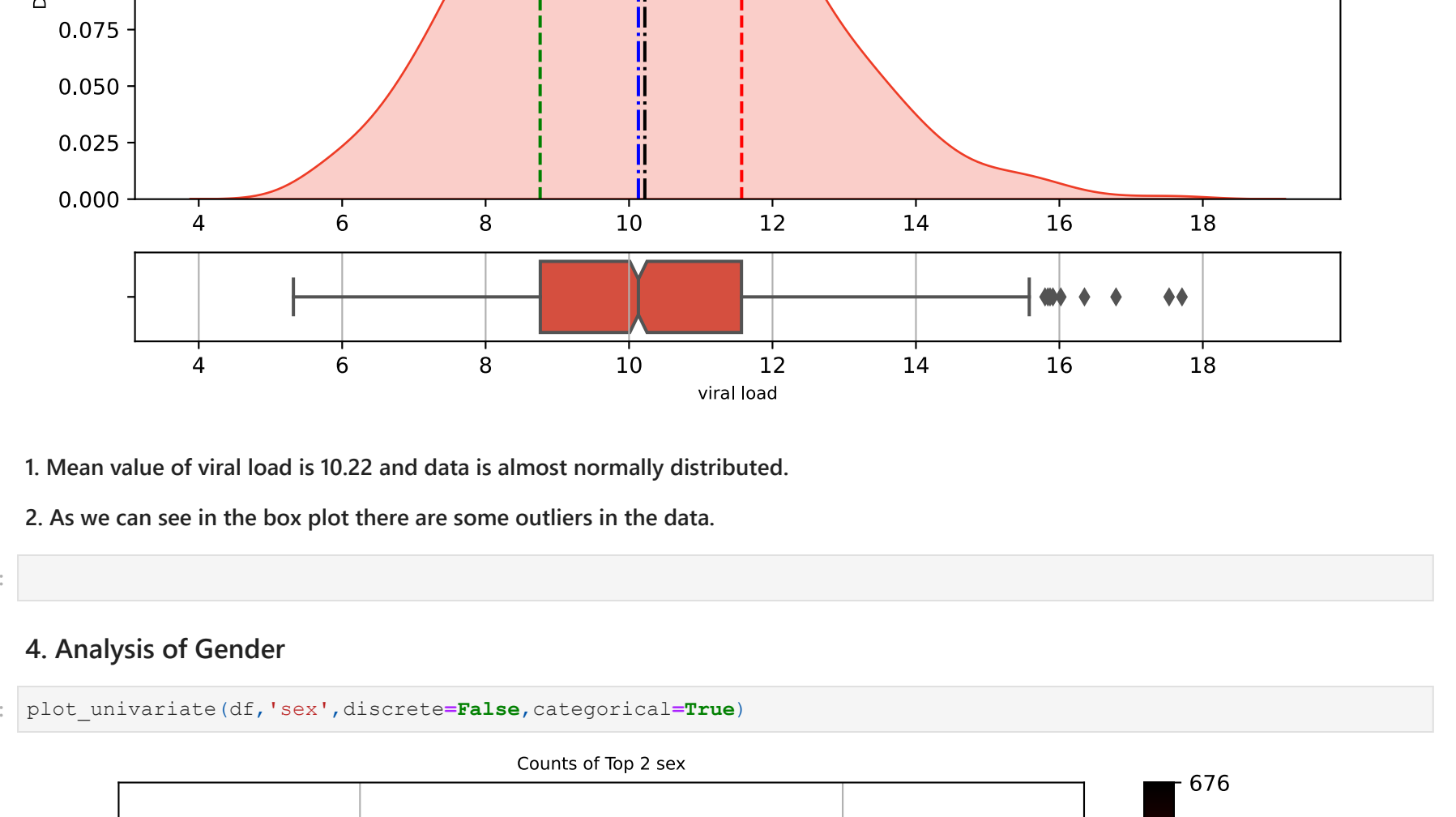
```
In [ ]:
```

Univariate Analysis

```
In [8]: ## generic function to plot the graphs
def plot_univariate(df, column, discrete=False, categorical=False, top="all"):
    plot_cmmap = "gist_heat"
    plot_color = "#B3C2E6"
    if categorical:
        # get value counts
        df_value_counts = df[column].value_counts()
        # if all, set top as length of df
        if top == "all":
            top = len(df_value_counts)
        df_value_counts = df_value_counts[:top]
        # draw count plot
        plt.figure(figsize=(10, 5))
        sns.countplot(x=column, data=df, palette=plot_cmmap, order=df_value_counts.index)
        # annotate values
        for p in ax.patches:
            # if value of bar is less than 10% of max value, change the alignment
            ax.annotate(f'{v}\n(p.get_height())', (p.get_x()+p.get_width()*0.4, p.get_height()*1.01),
                        ha="center", va="bottom", color=p.get_facecolor(), rotation=90, size=8)
        # increase ylin to accomodate annotations
        ax.set_ylim(0, ax.get_ylim()[1]*1.11)
        # add colorbar manually
        norm = plt.Normalize(df_value_counts.min(), df_value_counts.max())
        cmap = plt.get_cmap("plasma_r")
        sm = plt.cm.ScalarMappable(cmap=cmap, norm=norm)
        sm.set_array([])
        ax.figure.colorbar(sm)
        ax.set_xticklabels(ax.get_xticklabels(), rotation=90)
        # add title and x-axis grid
        title = f"Counts of Top {top} {column}"
        plt.title(title, fontsize=8)
        ax.xaxis.grid()
        plt.show()
    else:
        # get statistical parameters
        stat_params = df[column].describe().round(2)
        fig, axes = plt.subplots(2, 1, figsize=(10, 5), gridspec_kw={'height_ratios': [5, 1]})
        # plot only histogram if discrete variable else histogram with kde
        if discrete:
            cm = plt.cm.get_cmap("plasma_r")
            bins = np.histogram(bin_edges(df[column]), bins="auto")
            n, bins, patches = axes[0].hist(df[column], bins=bins, linewidth=1.2, edgecolor="black")
            bin_centers = 0.5 * (bins[1:] + bins[1:])
            # https://stackoverflow.com/questions/23061657/plot-histogram-with-colors-taken-from-colormap
            col = bin_centers - min(bin_centers)
            col /= max(col)
            for c, p in zip(col, patches):
                plt.setp(p, 'facecolor', cm(c))
            else:
                sns.kdeplot(data=df, x=column, ax=axes[0], color=plot_color, shade=True)
        # add lines showing stat params mean, median, 25% quant, 75% quant
        axes[0].axline(x=stat_params['mean'], color='b', linestyle="--", label=f"Mean-{stat_params['mean']}")
        axes[0].axline(x=stat_params['50%'], color='b', linestyle="--", label=f"Median-{stat_params['50%']}")
        axes[0].axline(x=stat_params['25%'], color='g', linestyle="--", label=f"Lower Quartile-{stat_params['25%']}")
        axes[0].axline(x=stat_params['75%'], color='r', linestyle="--", label=f"Upper Quartile-{stat_params['75%']}")
        axes[0].legend(fancybox=True, shadow=True, prop={'size': 10})
        # titles and labels
        axes[0].set_xlabel("", fontsize=8)
        axes[0].set_ylabel('Density', fontsize=8)
        axes[0].set_title(f"KDE & Box Plot-{column}", fontsize=10)
        axes[0].xaxis.grid(True)
        # plot box plot
        sns.boxplot(x=column, data=df, ax=axes[1], color=plot_color, notch=True)
        # set edge and face colours
        for i, box in enumerate(axes[1].artists):
            fc = box.get_facecolor()
            box.set_edgecolor(plot_color)
            box.set_facecolor(plt.colors.to_rgba(fc, 0.3))
        for i in range(64, 6*(i+1)):
            axes[1].lines[i].set_color(plot_color)
        # set labels of boxplot same as KDE plot
        axes[1].set_xlim(axes[0].get_xlim())
        axes[1].set_xlabel(column, fontsize=8)
        axes[1].xaxis.grid(True)
        plt.show()
```

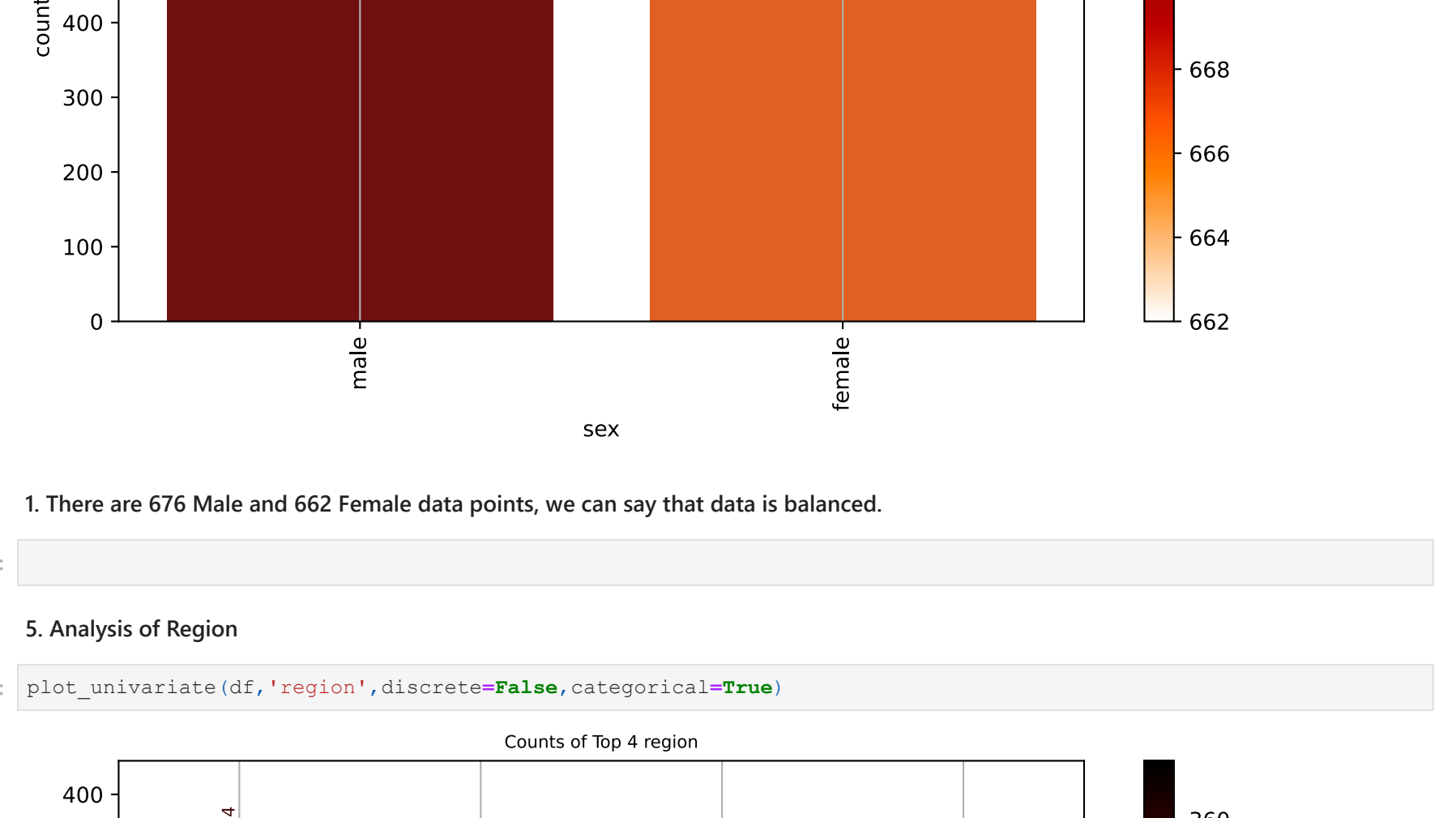
```
In [ ]:
```

1. Analysis of Age



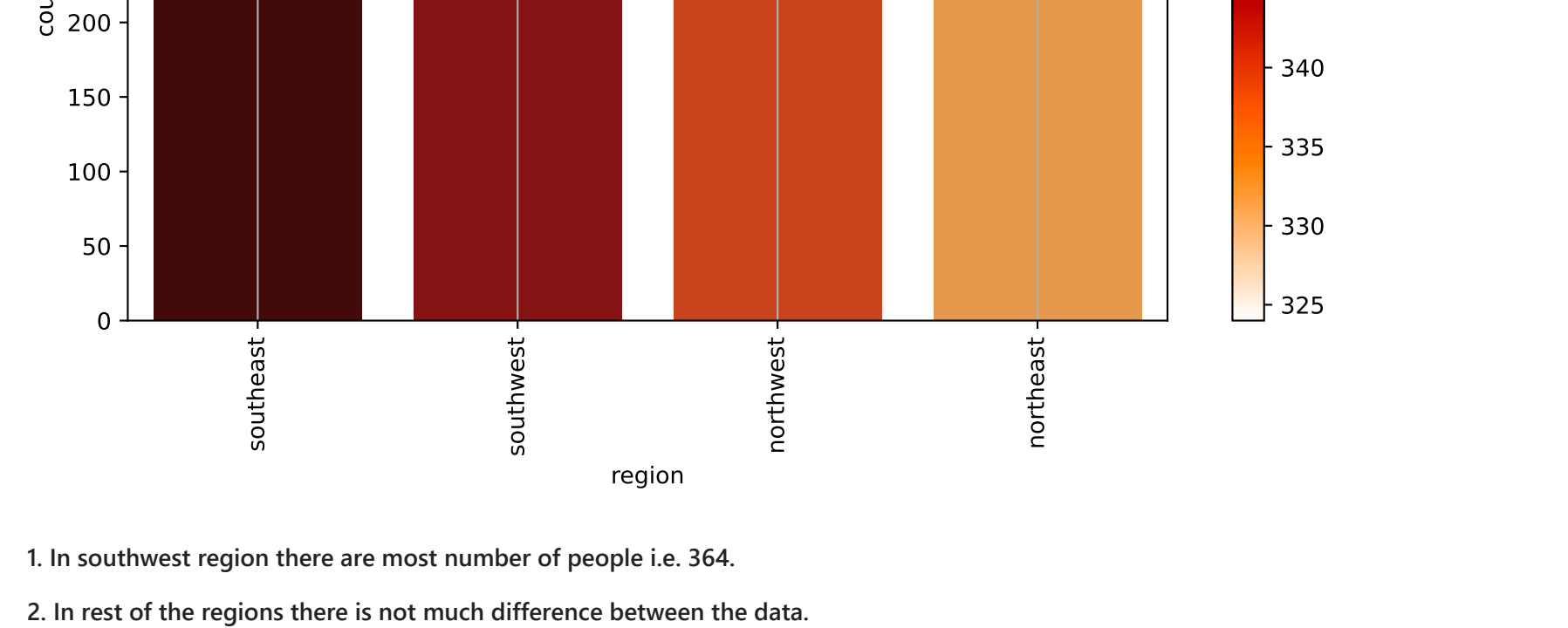
1. Mean value of age in the data is around 39.21 and it is almost normally distributed.
2. There are no outliers in the data as shown in the box plot.

2. Analysis of Hospitalization Charges



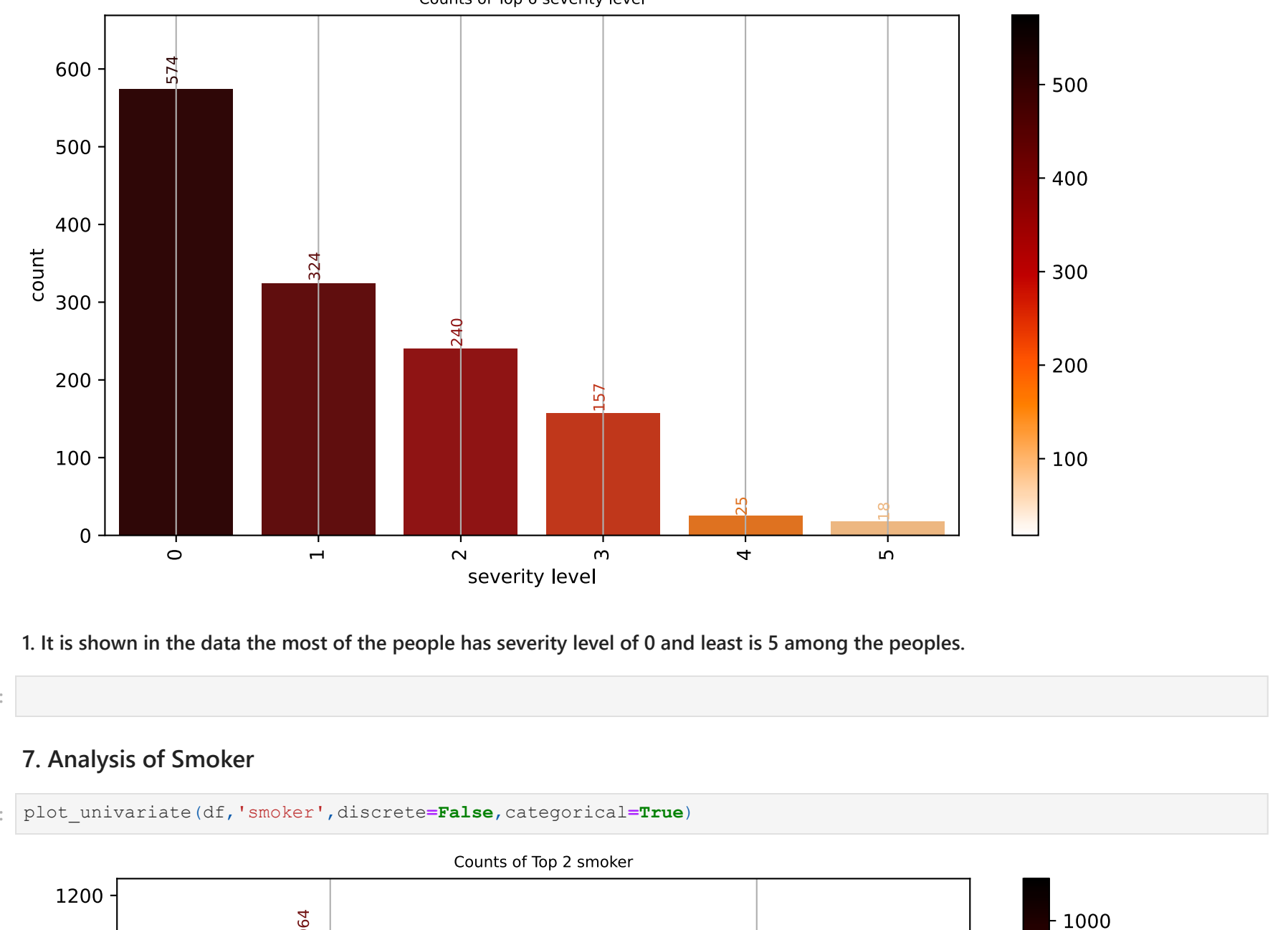
1. Mean value of Hospitalization charges is 33176 and there are outliers in the data that we will handle.
2. Data is not normally distributed, it is left skewed.

3. Analysis of Viral Load



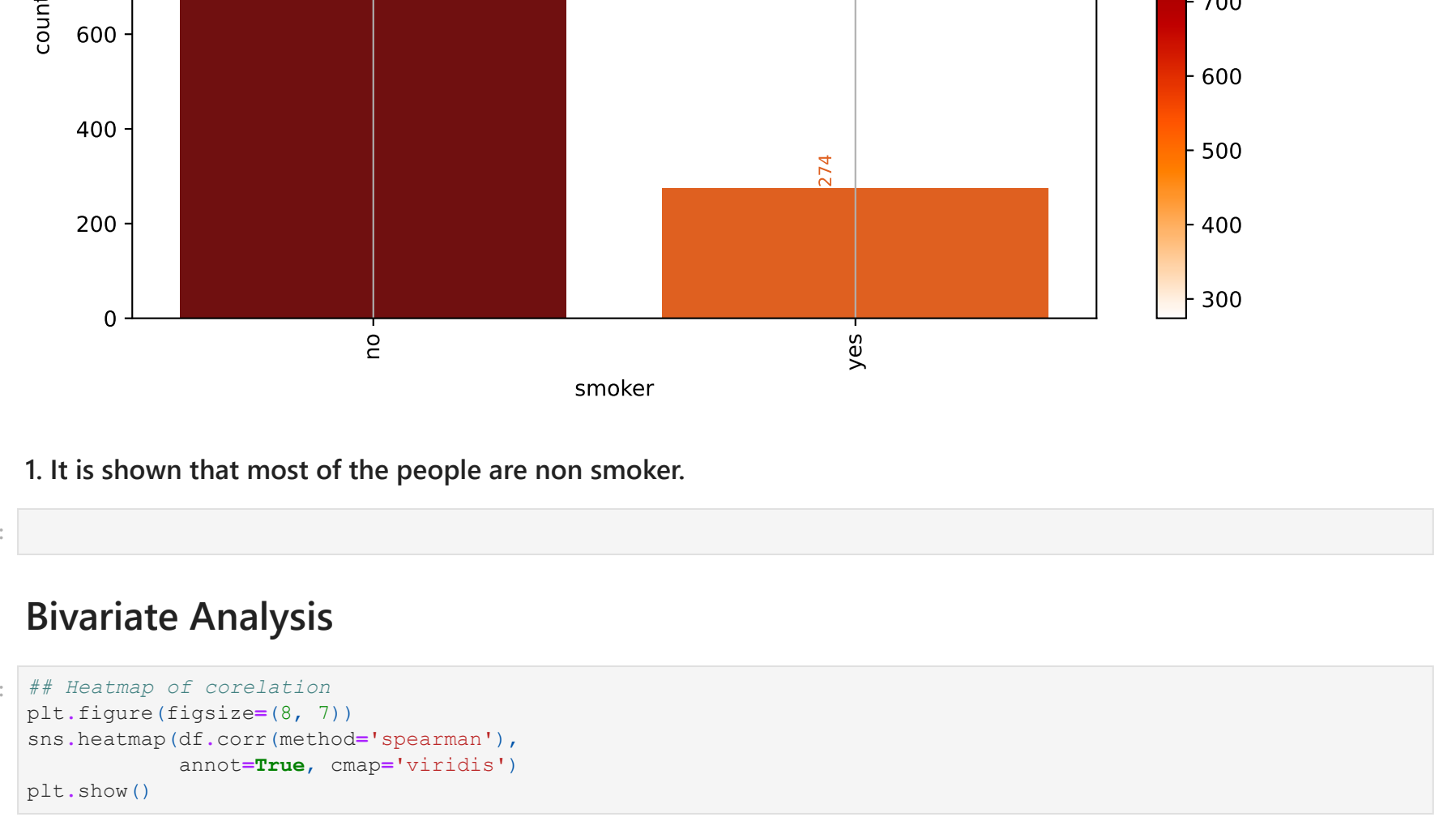
1. Mean value of viral load is 10.22 and data is almost normally distributed.
2. As we can see in the box plot there are some outliers in the data.

4. Analysis of Gender



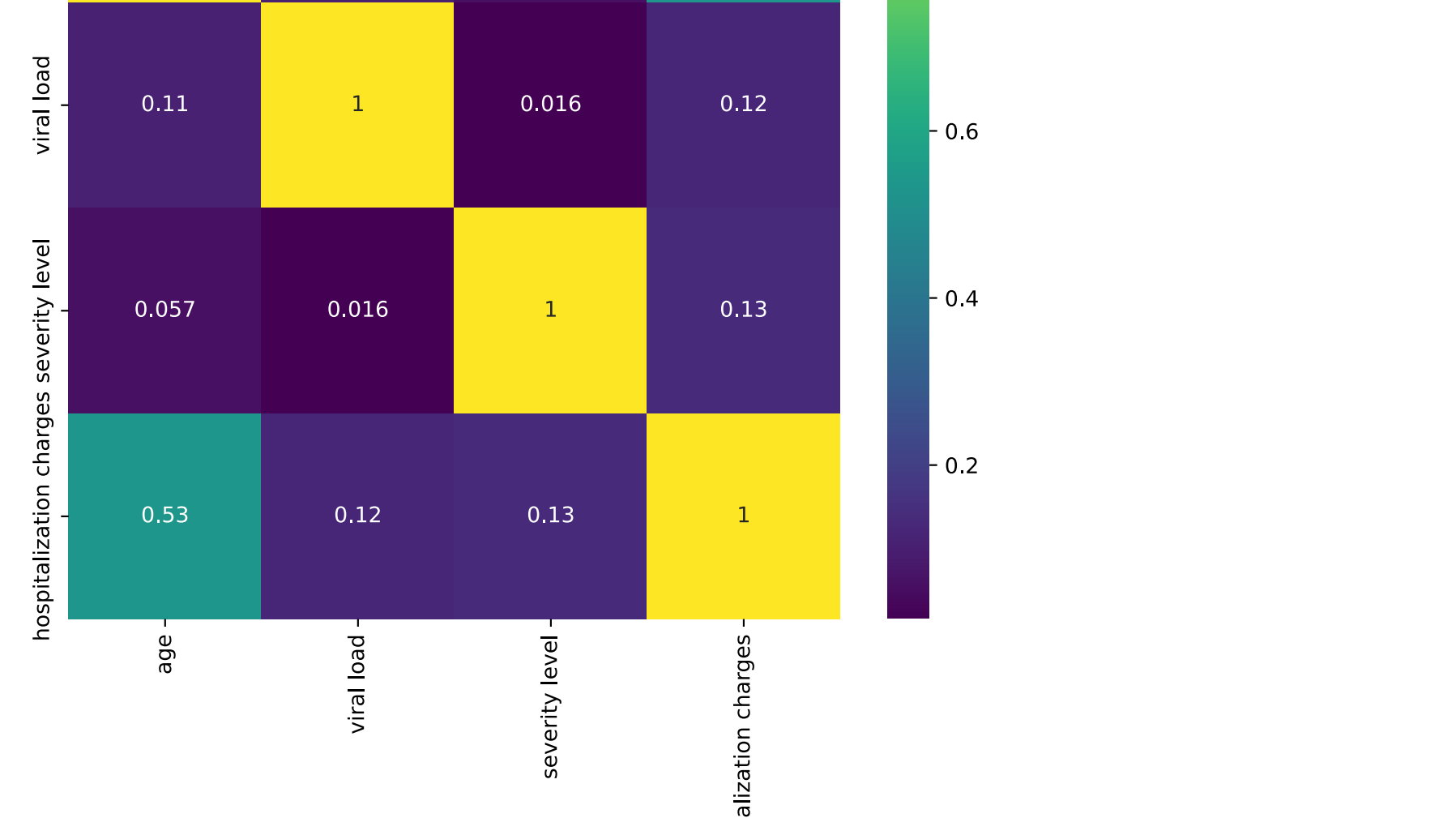
1. There are 676 Male and 662 Female data points, we can say that data is balanced.

5. Analysis of Region



1. In southwest region there are most number of people i.e. 364.
2. In rest of the regions there is not much difference between the data.

6. Analysis of Severity Level



1. It is shown in the data that most of the people have severity level of 0 and least is 5 among the peoples.

7. Analysis of Smoker



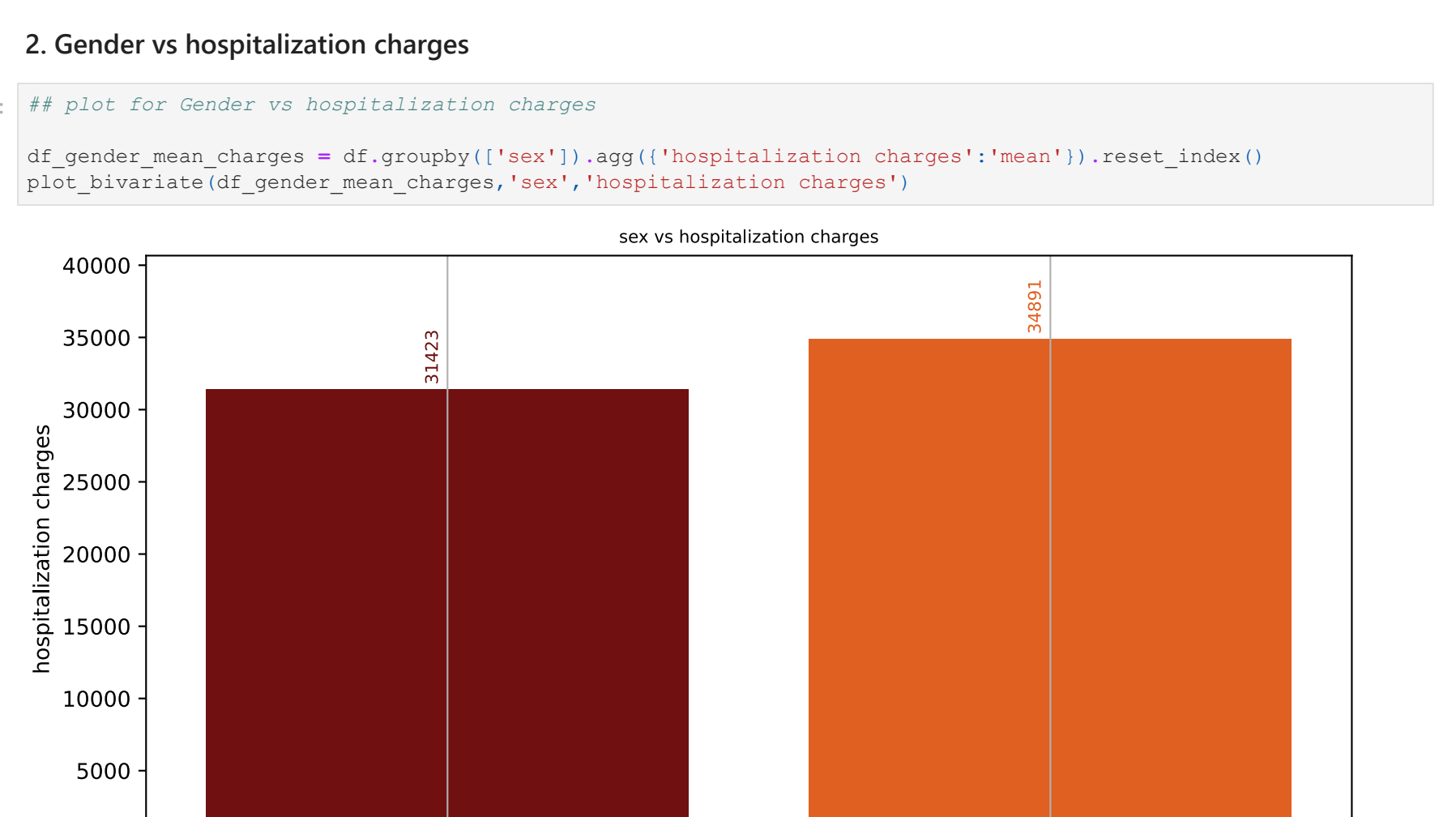
1. It is shown that most of the people are non smoker.

Bivariate Analysis



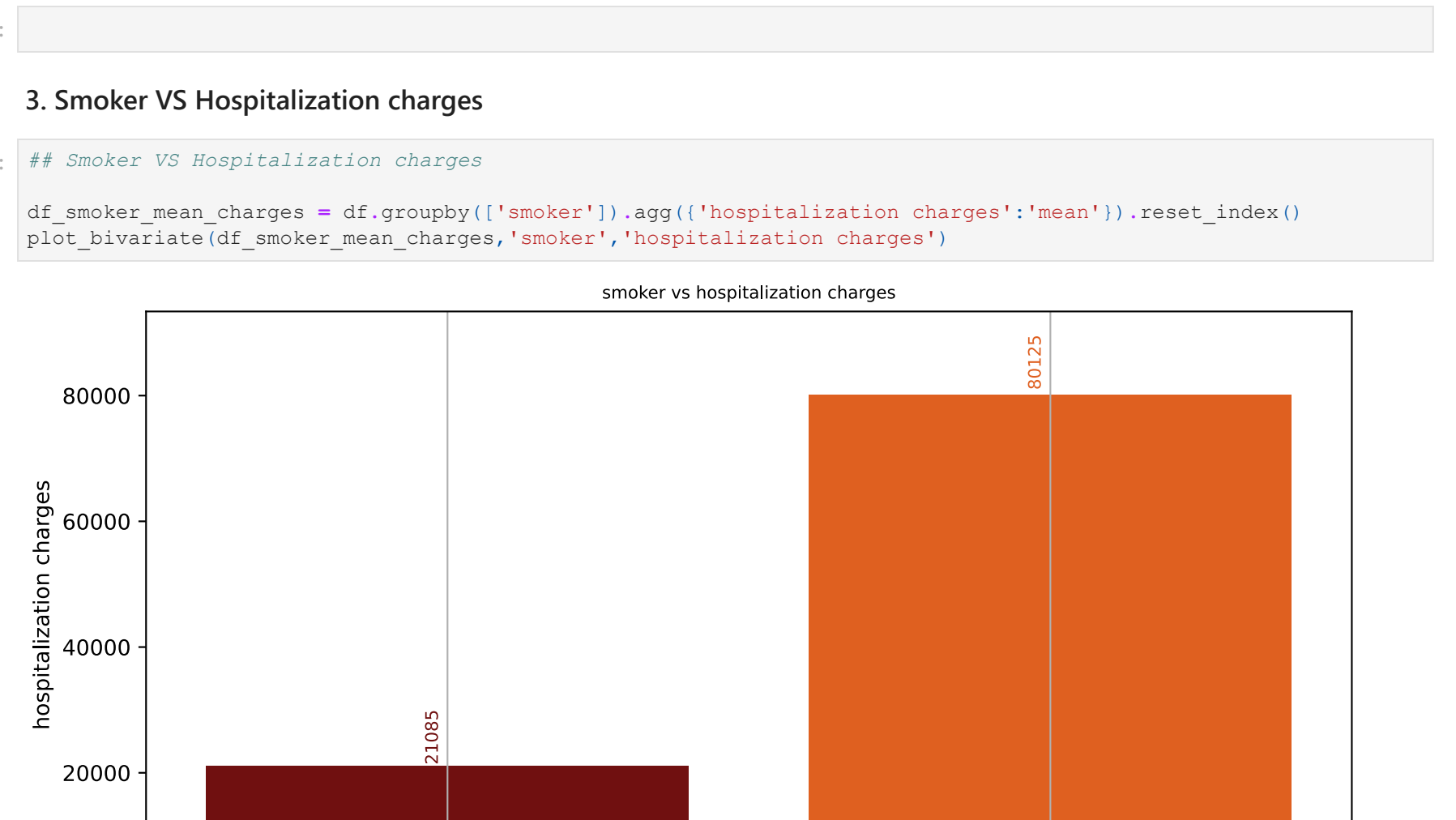
We can see that there is a strong positive correlation between age and hospitalization charges, and some relation between severity level and hospitalization charges and viral load

1. Smoker VS Gender



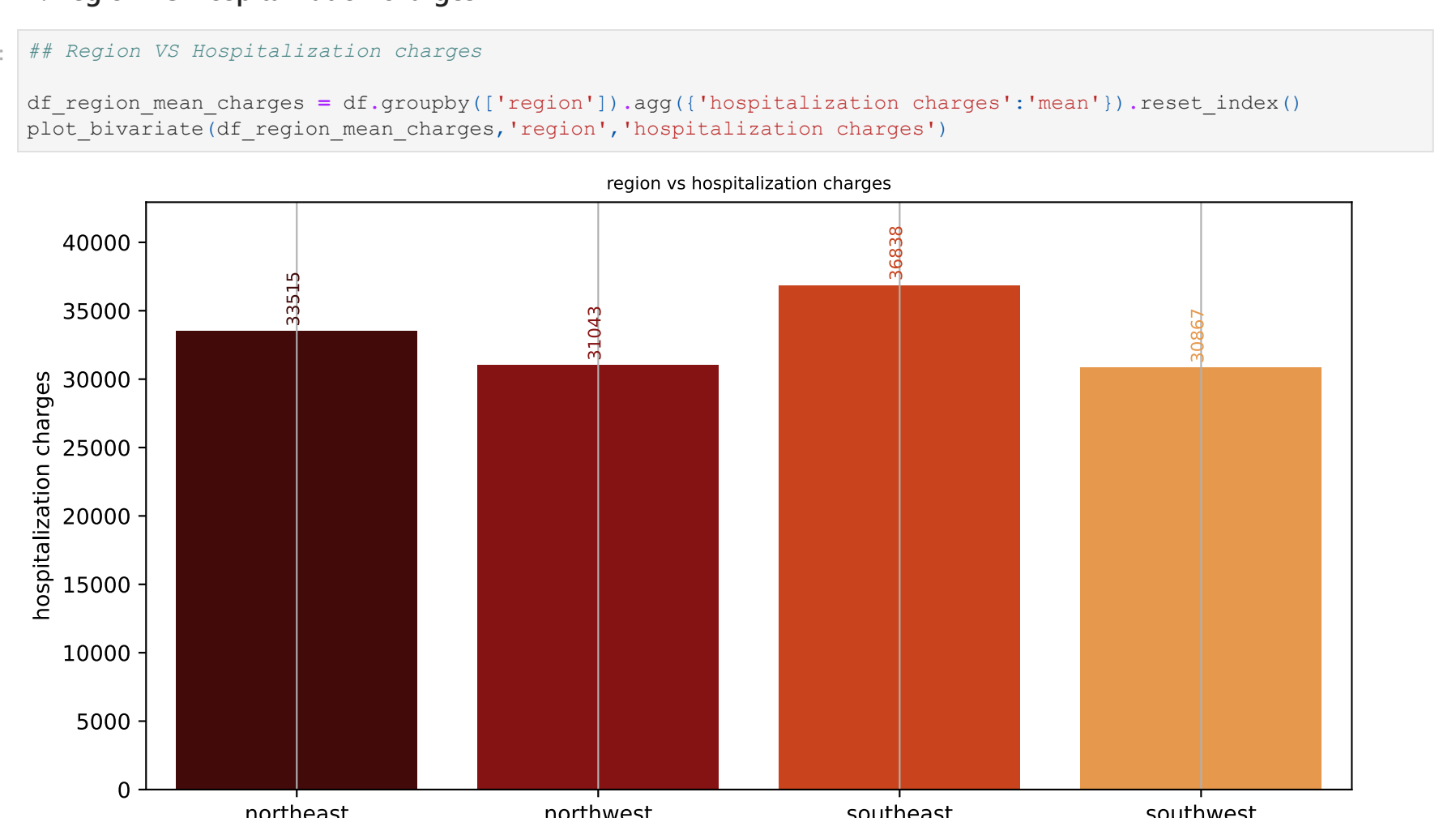
There are 115 females who smoke and 159 males who smoke

2. Gender vs hospitalization charges



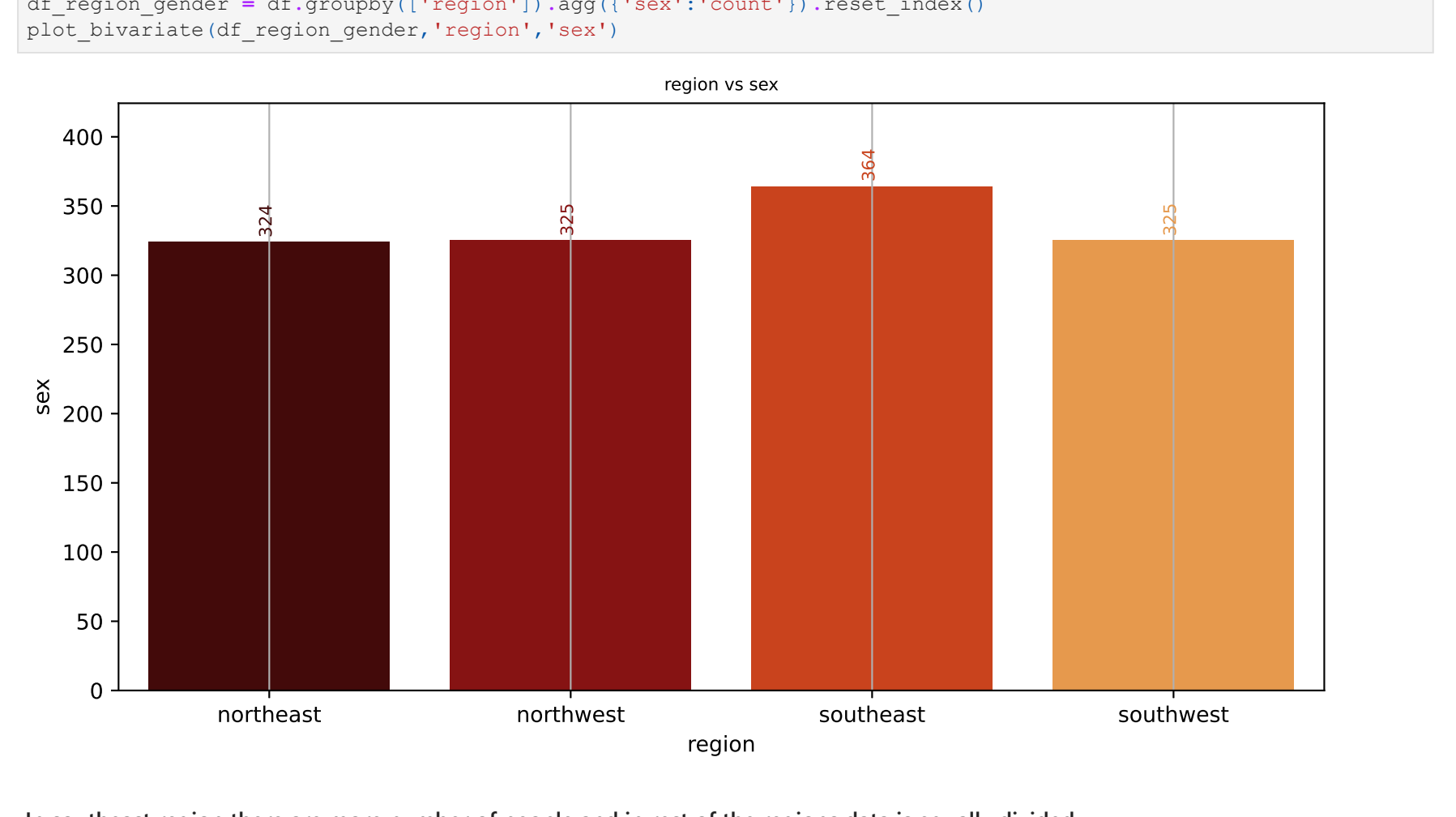
Mean value of hospitalization charges of females is 31423 and that of male is 34891. We can see that there is little difference between the means of charges for male and females.

3. Smoker VS Hospitalization charges



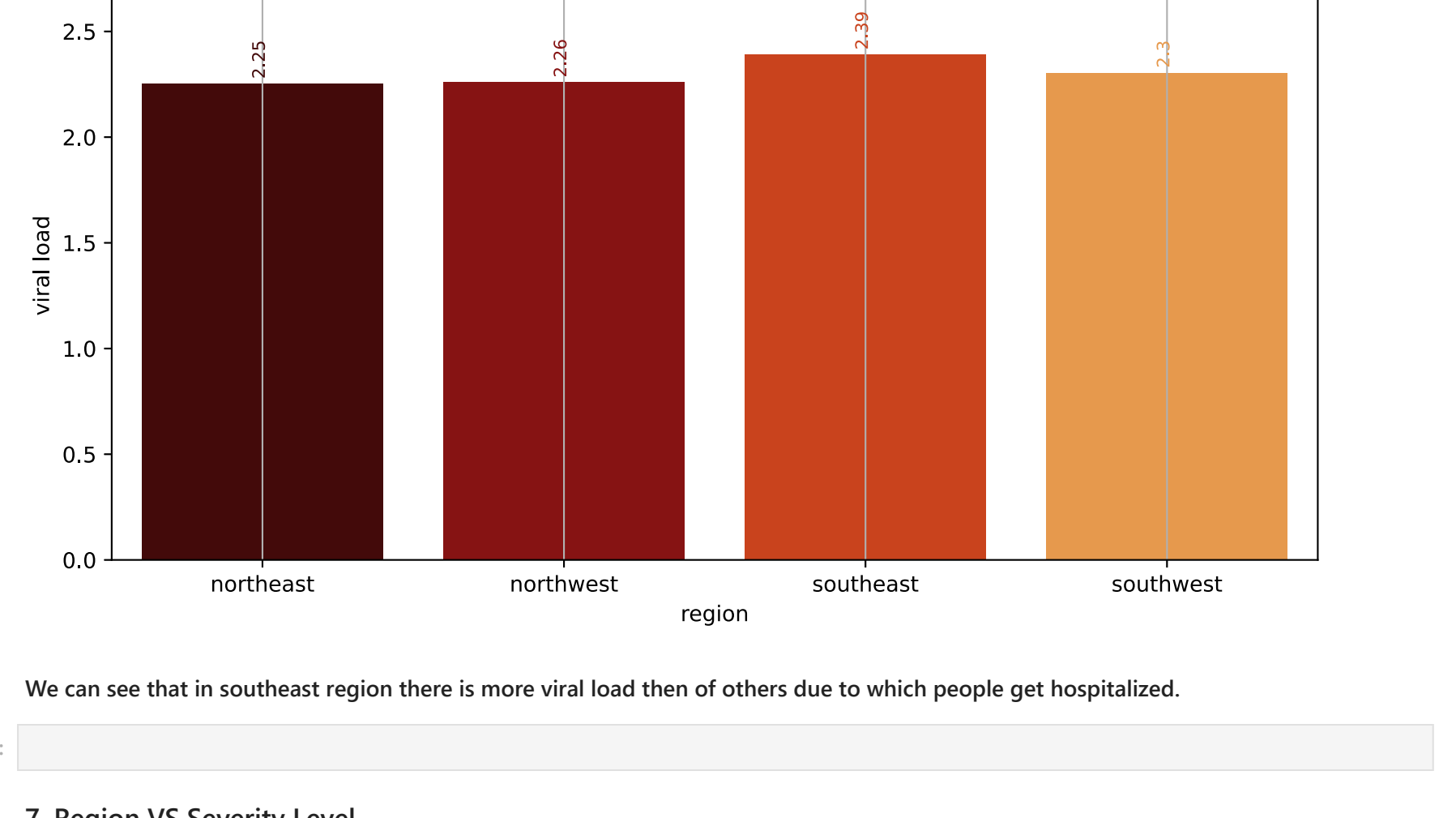
We can see that the mean charges of a smoker has more than the non smoker, because of the region that the people who smoke are more infected to the virus due to which the cost is high.

4. Region VS Hospitalization charges



As we have already seen that the people in the southeast are more and the hospitalization charges is also more in this region and in rest of the region there is no more difference.

5. Region VS Gender



In southeast region there are more number of people and in rest of the regions data is equally divided.

6. Region VS Viral Load



We can see that in southeast region there is more viral load than others due to which people get hospitalized.

7. Region VS Severity Level

