

Import libraries

```
import pandas as pd # Data manipulation
import numpy as np # Linear algebra
import warnings # Ignore warnings
import seaborn as sns #plots
import matplotlib.pyplot as plt # plots
```

```
warnings.filterwarnings("ignore")
```

Load Data

```
df=pd.read_csv("../content/insurance.csv")
```

```
df.head()
```

```
Out[4]:
```

	age	sex	bmi	children	smoker	region	charges
0	19	female	27.900	0	yes	southwest	16884.92400
1	18	male	33.770	1	no	southeast	1725.55230
2	28	male	33.000	3	no	southeast	4449.46200
3	33	male	22.705	0	no	northwest	21984.47061
4	32	male	28.880	0	no	northwest	3866.85520

Data Visualization

Unique Values

```
In [137]: col_list=df.col.unique() # col in df.select_dtypes("Object")
```

```
Out[137]:
```

```
region: ['northwest', 'southeast', 'northwest', 'southeast']
sex: ['female', 'male']
smoker: ['yes', 'no']
```

```
In [137]:
```

```
class Pie():

    def __init__(self):

        self.region=df.groupby('region').size()
        self.smoker=df.groupby('smoker').size()
        self.sex=df.groupby('sex').size()

    def region_pie(self):

        self.region.plot(kind='pie',title="Region Percent", figsize=[10,8],
            colors=['#77dd77','#fdd9d9','#e4b6f4','#fddcae1'],explode=[0,0,0,0],
            autopct='%1.2f%%(%.1f%%)'.format(p,p/100)*self.region.sum())

        plt.ylabel('none')

    def smoker_pie(self):

        self.smoker.plot(kind='pie',title="Smoker Percent", figsize=[10,8],
            colors=['#77dd77','#fdd9d9'],explode=[0,0,0,0],
            autopct='%1.2f%%(%.1f%%)'.format(p,p/100)*self.region.sum())

        plt.ylabel('none')

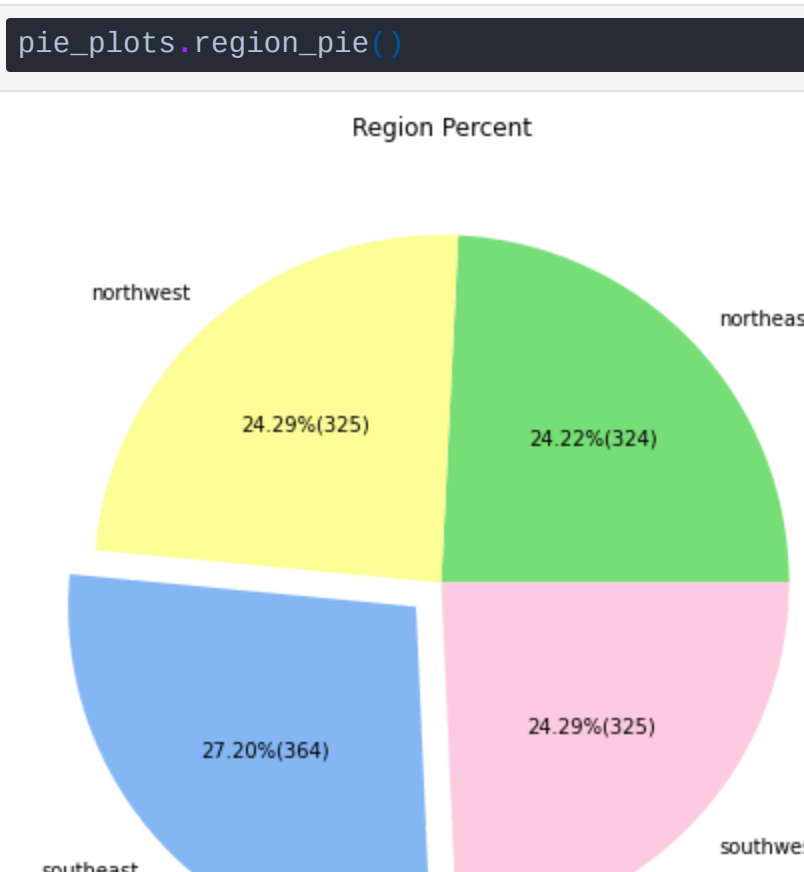
    def sex_pie(self):

        self.sex.plot(kind='pie',title="Sex Percent", figsize=[10,8],
            colors=['#FFD1DC','#2271b3'],
            autopct='%1.2f%%(%.1f%%)'.format(p,p/100)*self.sex.sum())

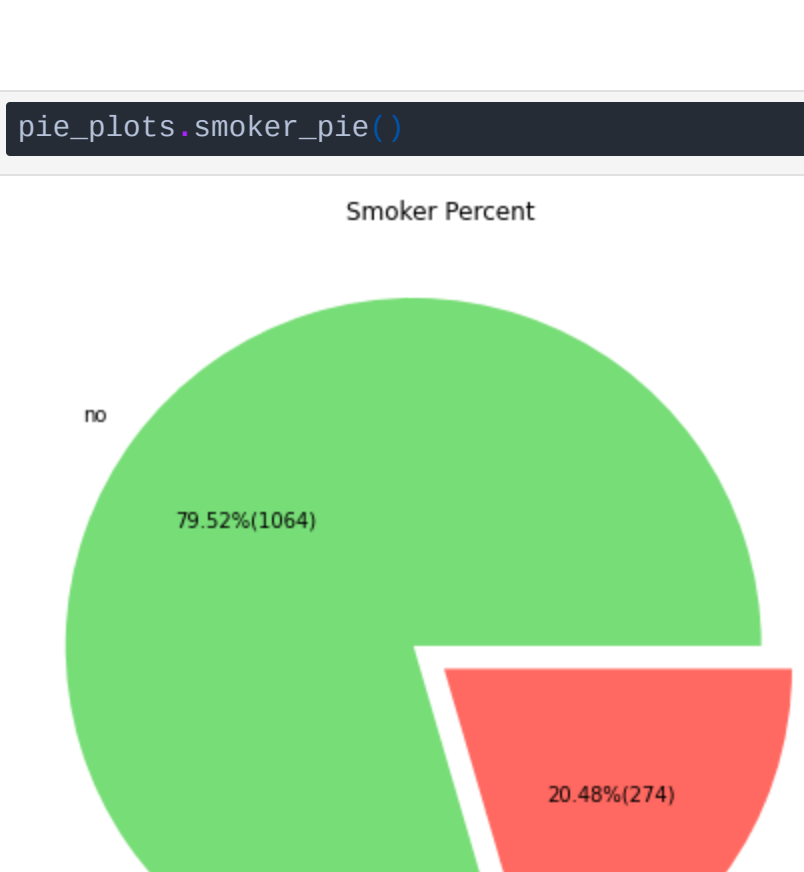
        plt.ylabel('none')
```

```
In [138]: pie_plots=Pie()
```

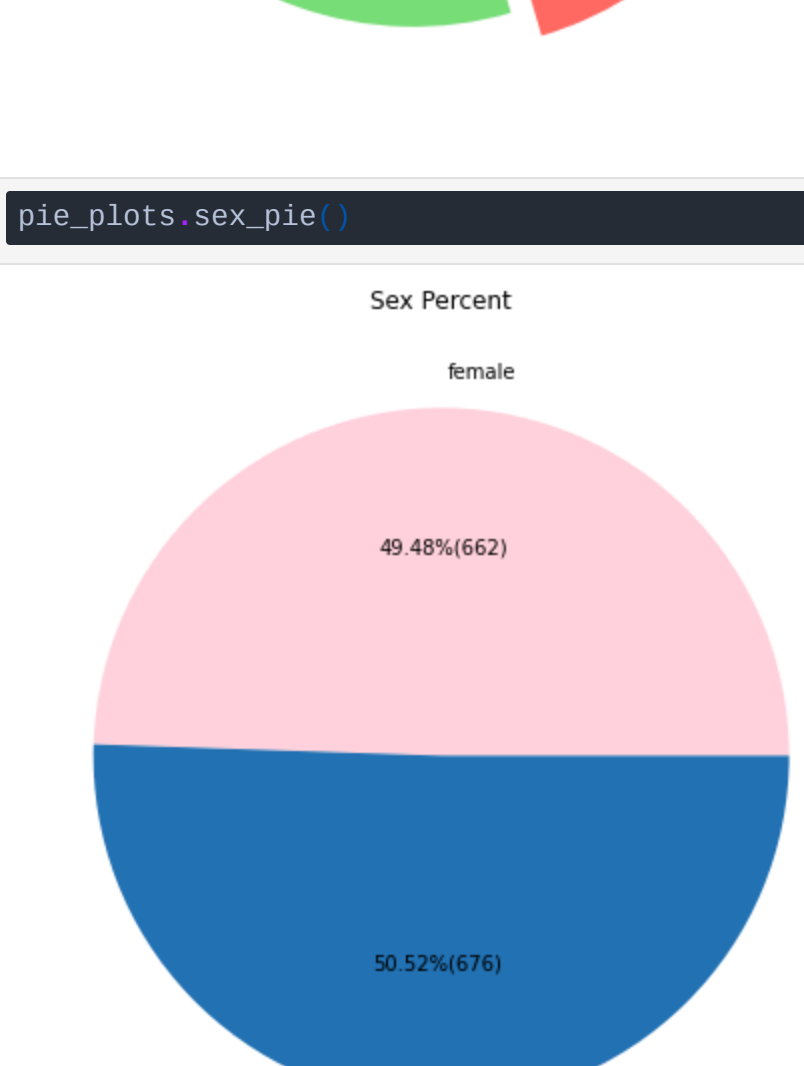
```
In [138]: pie_plots.region_pie()
```



```
In [134]: pie_plots.smoker_pie()
```



```
In [139]: pie_plots.sex_pie()
```



```
In [140]: df.children.value_counts()
```

```
Out[140]:
```

```
0    313
1    324
2    248
3    167
4     25
5     18
dtype: int64
```

```
In [141]: sns.set_style(style='whitegrid')
```

```
In [142]:
```

```
def histogram(feature,title):

    fig=plt.subplots(1,1,figsize=(10,8))

    ax=plt.subplot(1,1)
    ax.set_title(title)
    ax.hist(df[feature],ec='k',color='#FA0A5E',lw=1)

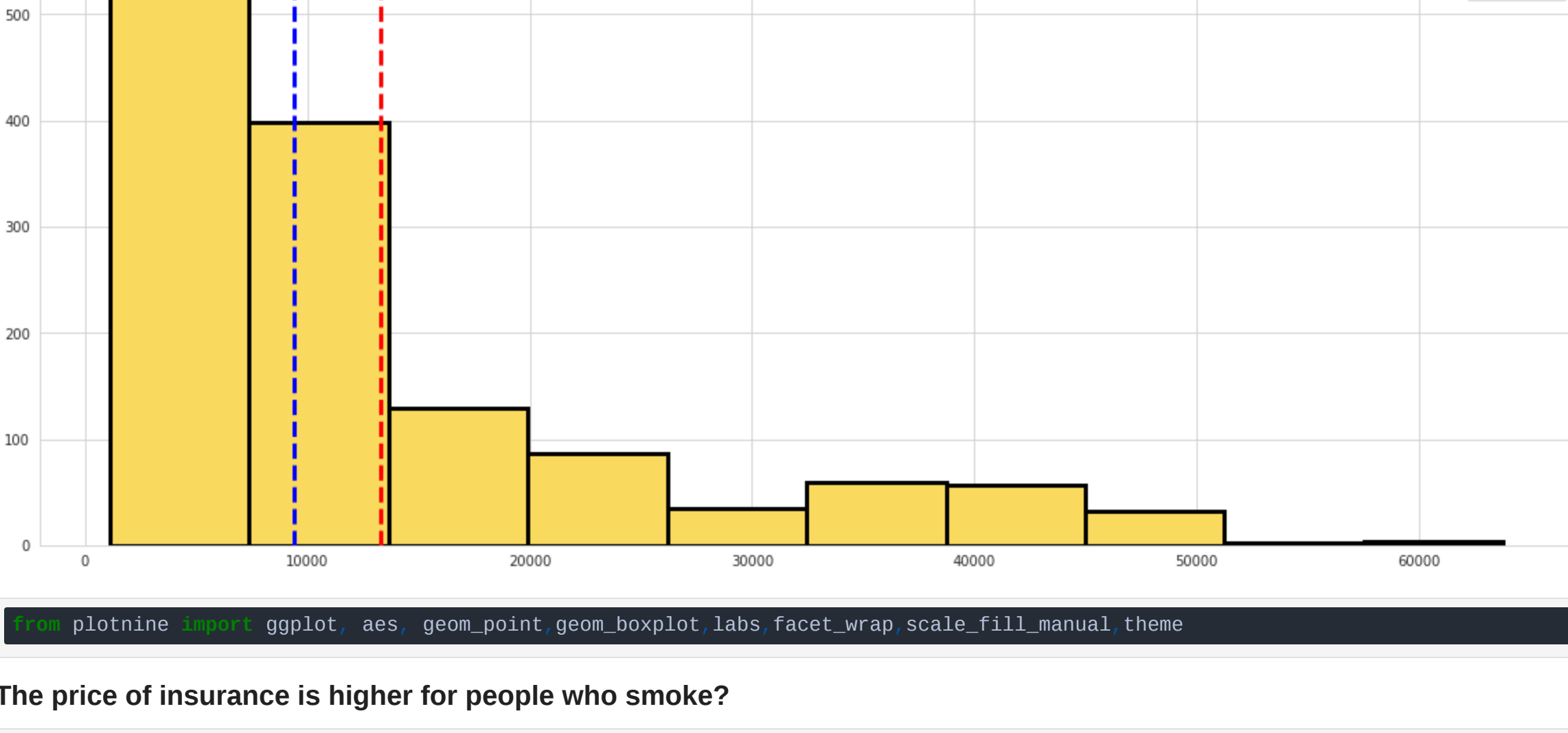
    ax.axvline(df[feature].mean(),
        color='red',
        linestyle='--',
        lw=1,label='Mean')

    ax.axvline(df[feature].median(),
        color='blue',
        linestyle='--',
        lw=1,label='Median')

    ax.legend()

    plt.show()
```

```
In [143]: histogram("charges","Charges")
```

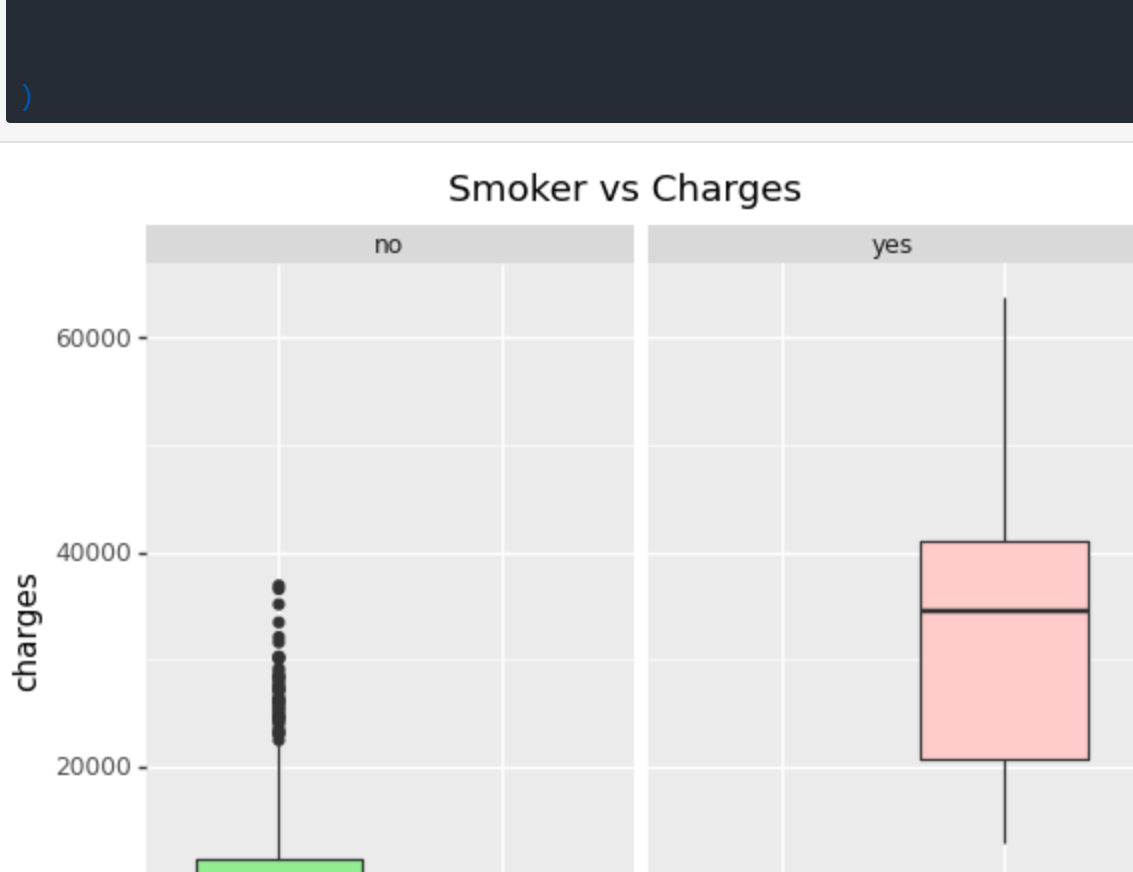


```
In [144]: plotnine==ggplot,aes=geom_point,geom=boxplot,labs=facet_wrap,scale_fill_manual,theme
```

The price of insurance is higher for people who smoke?

```
In [145]:
```

```
ggplot(df)
+ aes(x='smoker',y='charges',fill='smoker')
+ geom_boxplot() + labs(title='Smoker vs Charges')
+ facet_wrap("smoker")
+ theme(legend_position='none')
+ scale_fill_manual(values=["#99ee99","#ffcccb"])
```



```
Out[145]:
```

We observe a strong presence of outliers, for the category of non-smokers.

```
In [146]: (df.groupby("smoker"))["charges"].mean()
```

```
Out[146]:
```

```
smoker    8434.265298
yes       32958.231832
dtype: float64
```

The average price of insurance is considerably much higher than non-smokers. Since they usually have poorer physical health and need more treatments.

```
In [147]: histogram("bmi","BMI")
```



Most of the BMI data is within a normal distribution. But even so, it is possible to appreciate outlier values in the upper range.

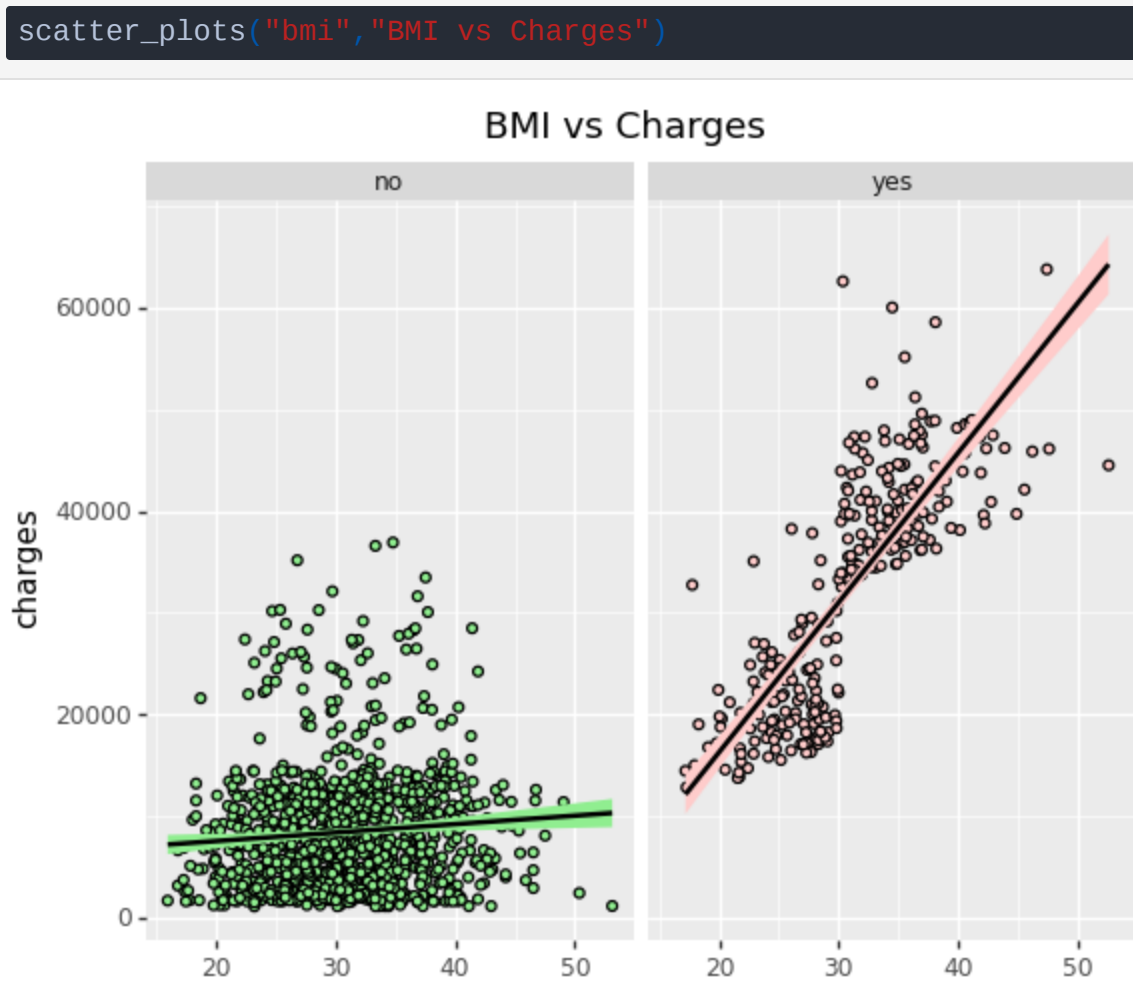
```
In [148]:
```

```
plotnine.facets==ggplot,facet_grid
plotnine.geoms==ggplot,geom_smooth
scatter_plots=feature,title:

def ggplot(df):
    + aes(x=feature,y='charges',fill='smoker',alpha=0.3)
    + geom_point() + labs(title=title,x=feature)
    + facet_wrap("smoker")
    + theme(legend_position='none')
    + scale_fill_manual(values=["#99ee99","#ffcccb"])
    + geom_smooth(method='lm')
```

People with a high BMI the insurance charge is higher?

```
In [149]: scatter_plots("bmi","BMI vs Charges")
```

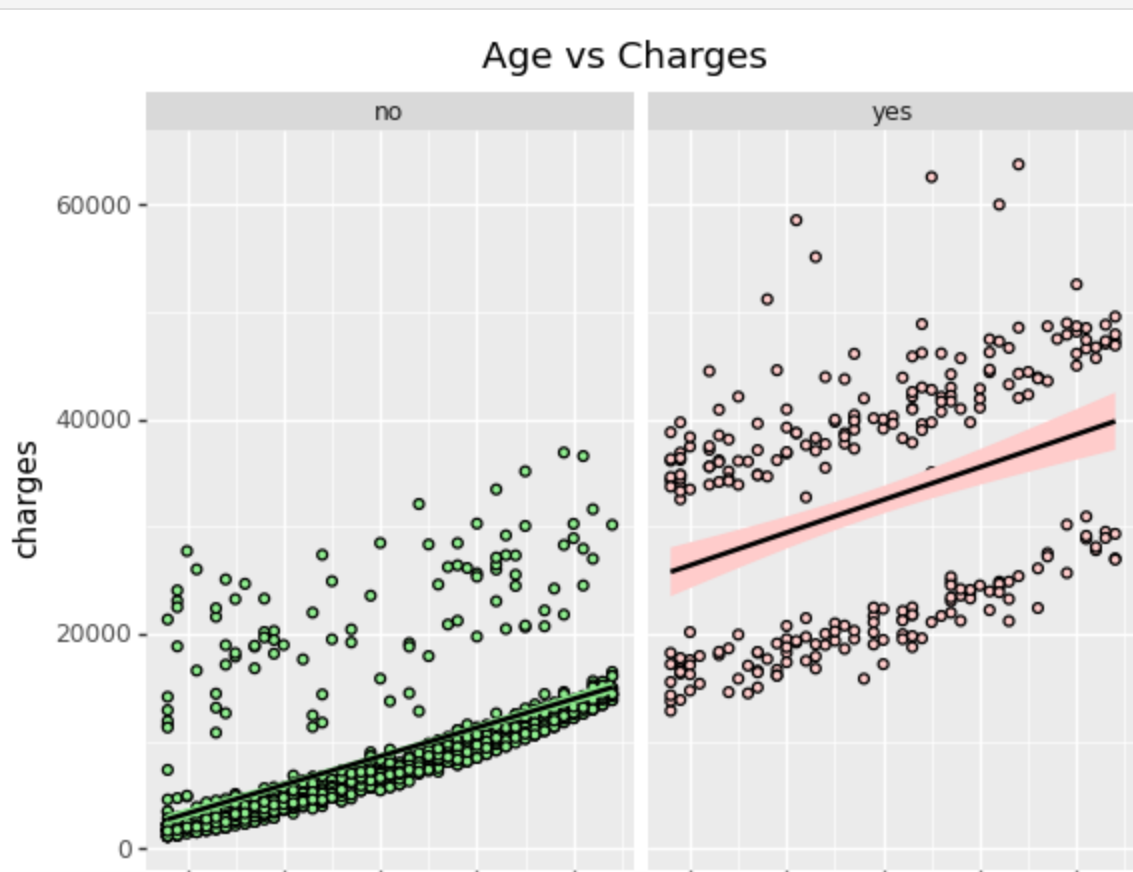


```
Out[149]:
```

- For non-smokers, the data trend remains constant.
- While for smokers the trend line is linear, that is, one value increases proportionally with another.

Does age influence the price of insurance?

```
In [150]: scatter_plots("age","Age vs Charges")
```



```
Out[150]:
```

- For non-smokers, the trend of the data remains linear with respect to age. So we can create a linear regression model to model the data for non-smokers. But with the disadvantage that it would not explain so well for the opposite case.
- For smokers, the relationship between age seems non-existent.

```
In [151]:
```

```
corr_pearson=df.corr(method='pearson')
corr_spearman=df.corr(method='spearman')

corr_matrix=corr_pearson

plt.figure(figsize=(10,8))
sns.heatmap(corr_matrix,annot=True)
```

Correlation Matrix

It is used to establish possible relationships between variables

Pearson

```
In [152]: correlation_matrix.corr_pearson
```



Spearman

```
In [153]: correlation_matrix.corr_spearman
```



The correlation is measured from 0 to 1 if it is positive. There does not appear to be a strong relationship between the variable of interest. It is still too early to start ruling out variables, since these variables can complement the predictions.

Conclusion

- The **smoker** variable **positively** affects the price of health insurance.
- The **age** has a **better relationship** for the group of **non smokers**. While the BMI maintains more importance when the person smokes.
- The variables that describe the characteristics of the individual regarding their state of health. Like age, the BMI and whether the **person is a smoker increase considerably** if the individual is in poor health, while for **healthy people** the charge is **not usually high**. Especially if the person does **not smoke**, the charge increases proportionally to the number of years of the individual.