

Import libraries

```
In [49]: 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
```

```
In [2]: warnings.filterwarnings("ignore")
```

Load Data

```
In [4]: df=pd.read_csv("/content/insurance.csv")
```

```
In [5]: df.head()
```

```
Out[5]:
```

	age	sex	bmi	children	smoker	region	charges
0	19	female	27.900	0	yes	southwest	16984.92400
1	18	male	33.770	1	no	southeast	1725.55200
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	3866.85620

Unique Values

```
In [6]: [col_list(df[col].unique()) for col in df.select_dtypes("object")]
```

```
Out[6]:
```

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

```
In [48]: df["sex"].value_counts()
```

```
Out[48]:
```

```
sex
female    662
male     1084
Name: sex, dtype: int64
```

```
In [46]: df["region"].value_counts()
```

```
Out[46]:
```

```
region
southeast    366
southwest    325
northwest    325
northeast    324
Name: region, dtype: int64
```

```
In [47]: df["smoker"].value_counts()
```

```
Out[47]:
```

```
smoker
yes      1684
no       274
Name: smoker, dtype: int64
```

```
In [73]: df["children"].value_counts()
```

```
Out[73]:
```

```
children
0       274
1       324
2       240
3       157
4        25
5         18
Name: children, dtype: int64
```

Data Visualization

```
In [59]: sns.set_style(style="whitegrid")
```

```
In [321]:
```

```
def histogram(feature_title):

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

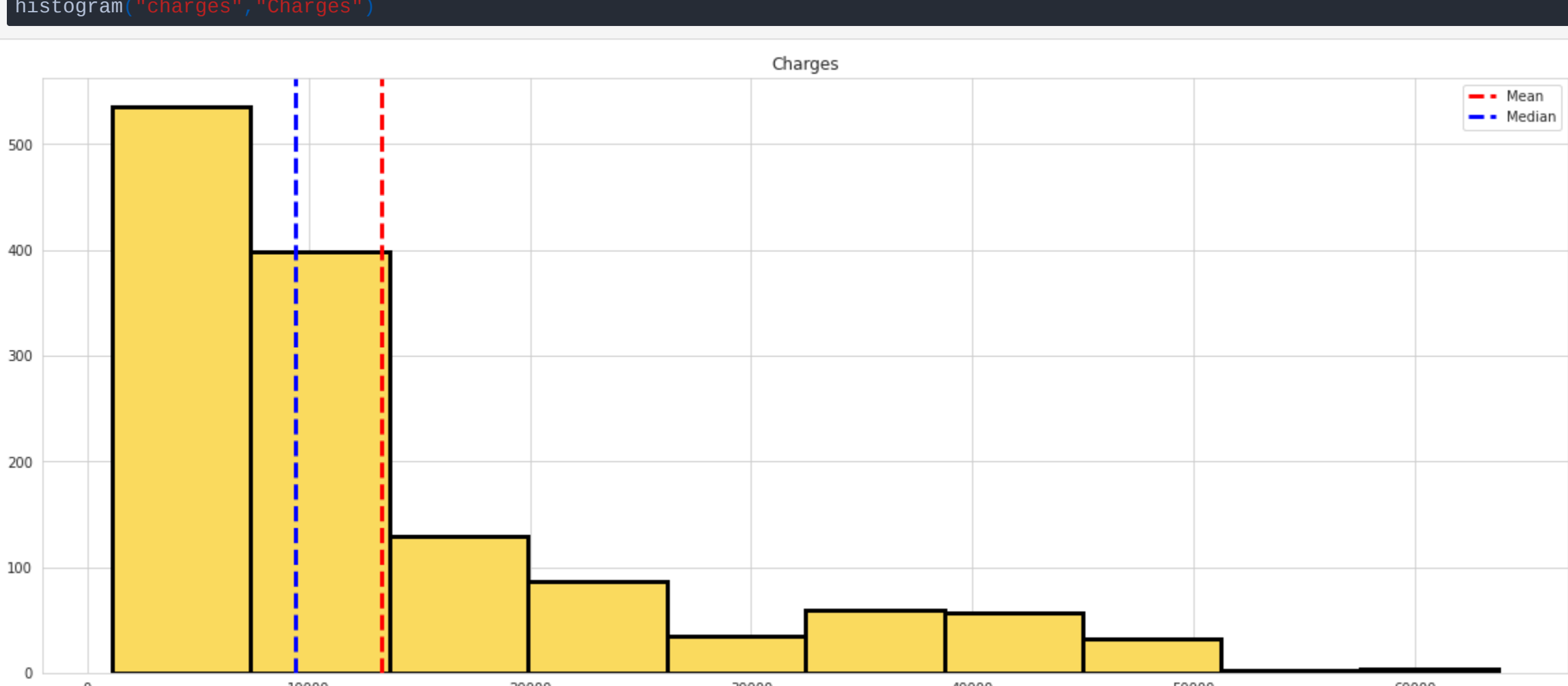
    ax=plt.subplot(1,1)
    ax.hist(df[feature], ec="k", color="#FADABE", lw=1)

    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()
```

```
In [322]: histogram("charges", "Charges")
```

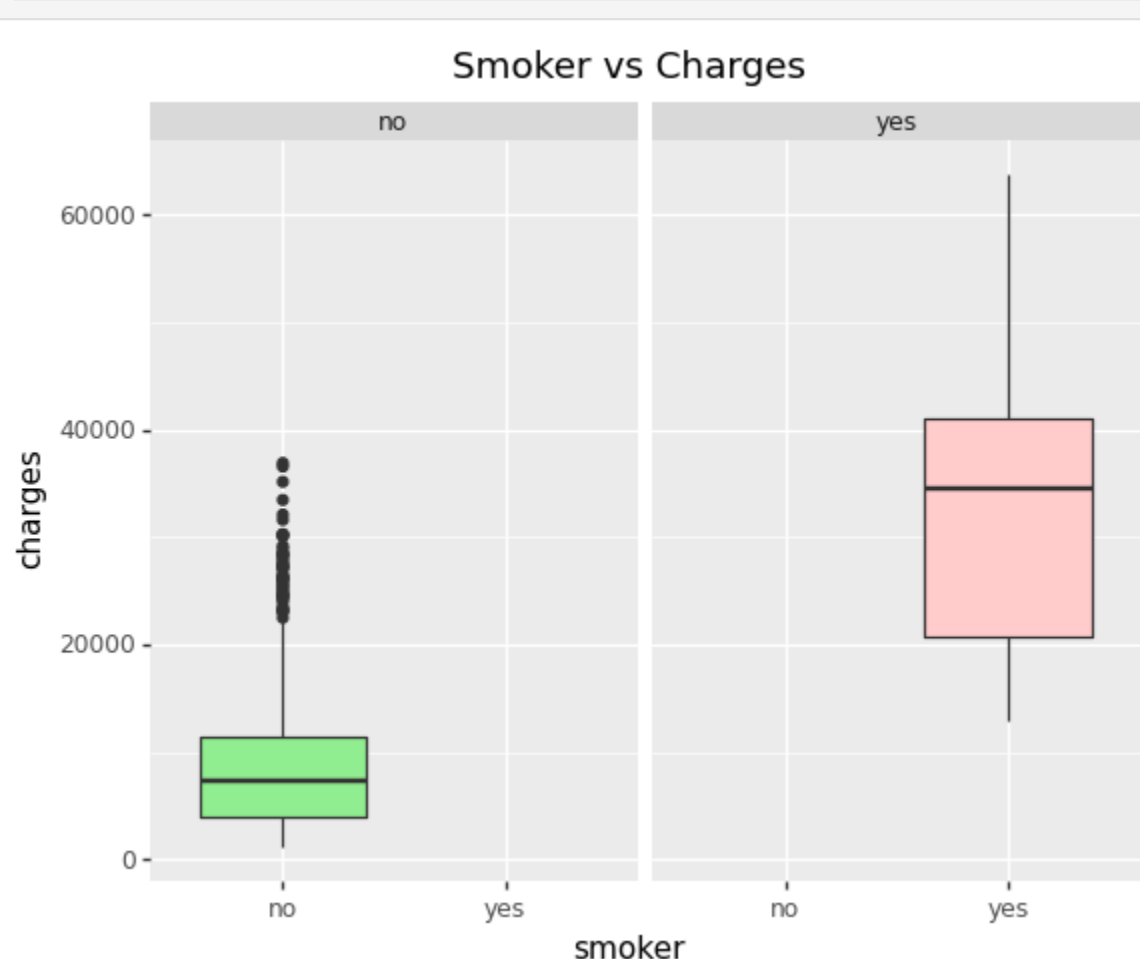


```
In [81]: plt.plotnine.theme('ggplot', aes=geom_point, geom=geom_boxplot)
```

Detection outliers

```
In [82]:
```

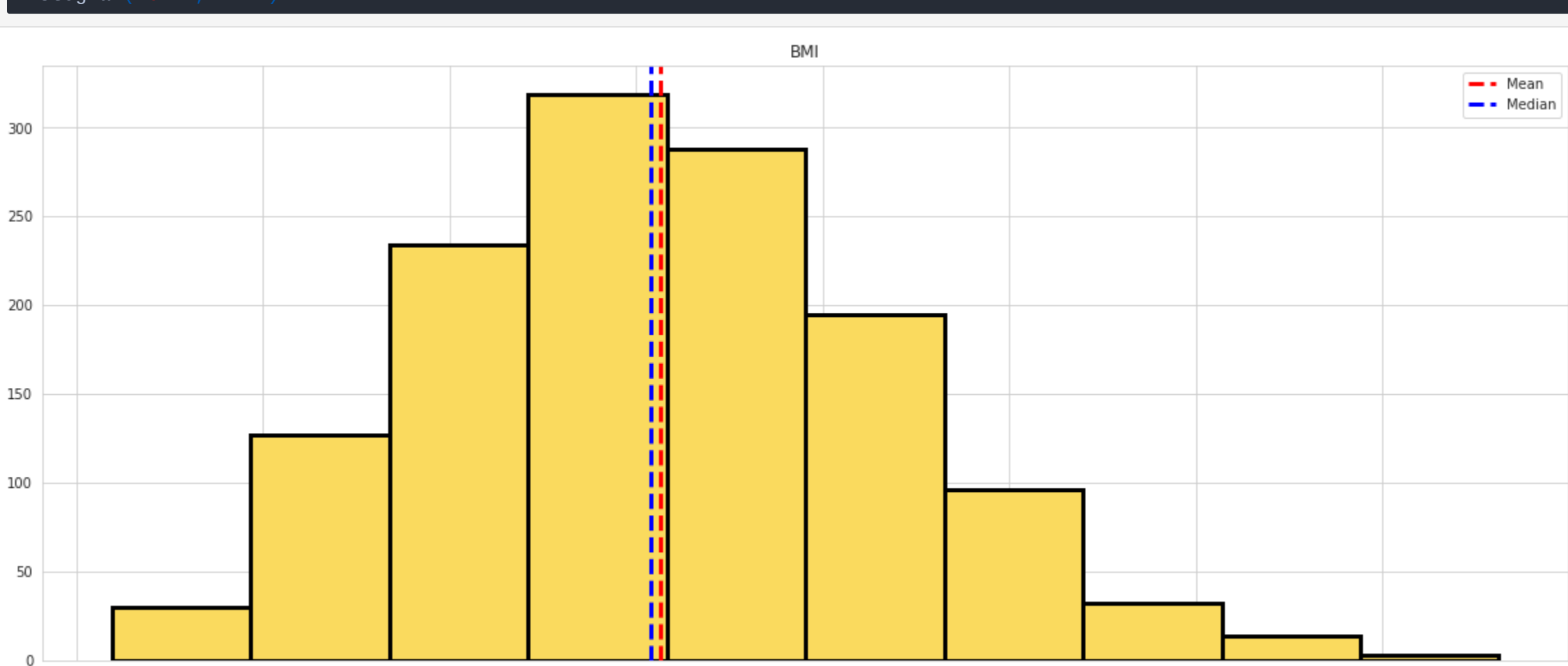
```
ggplot(df)
+ aes(x="smoker", y="charges", fill="smoker")
+ geom_boxplot() + labs(title="Smoker vs Charges")
+ facet_wrap("smoker")
+ theme(legend_position="none")
+ scale_fill_manual(values=["#99ee99", "#ffcc99"])
```



```
Out[82]: plt.plotnine.theme('ggplot', aes=geom_point)
```

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

```
In [323]: 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.

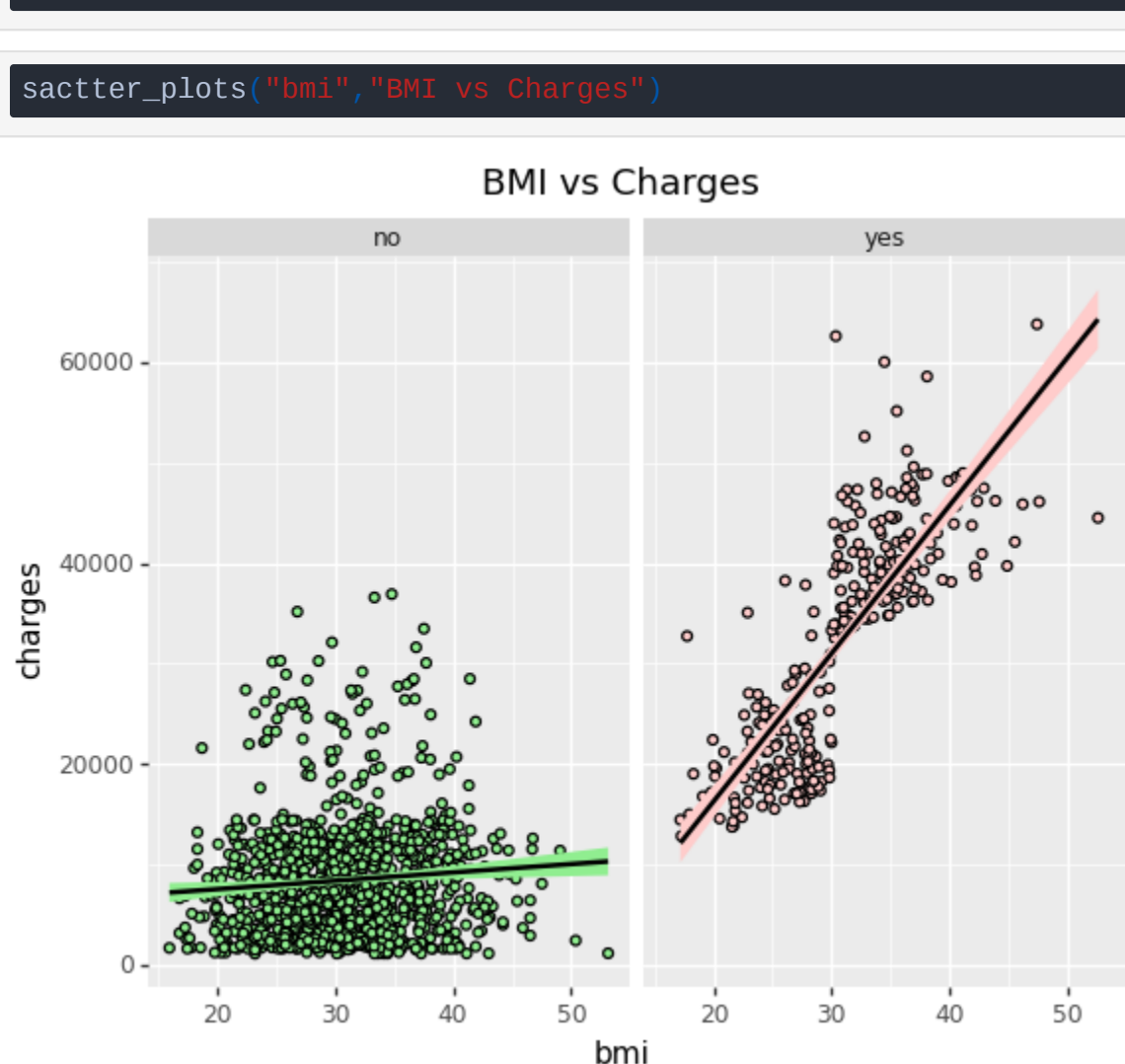
```
In [330]: plt.plotnine.facets(scales="y", facet_grid)
```

```
plt.plotnine.geoms(scales="y", geom_smooth)

def scatter_plots(feature_title):

    ggplot(df)
    + aes(x=feature, y="charges", fill="smoker", alpha=.1)
    + geom_point() + labs(title=title, x=feature)
    + facet_wrap("smoker")
    + theme(legend_position="none")
    + scale_fill_manual(values=["#99ee99", "#ffcc99"])
    + geom_smooth(method="lm")
```

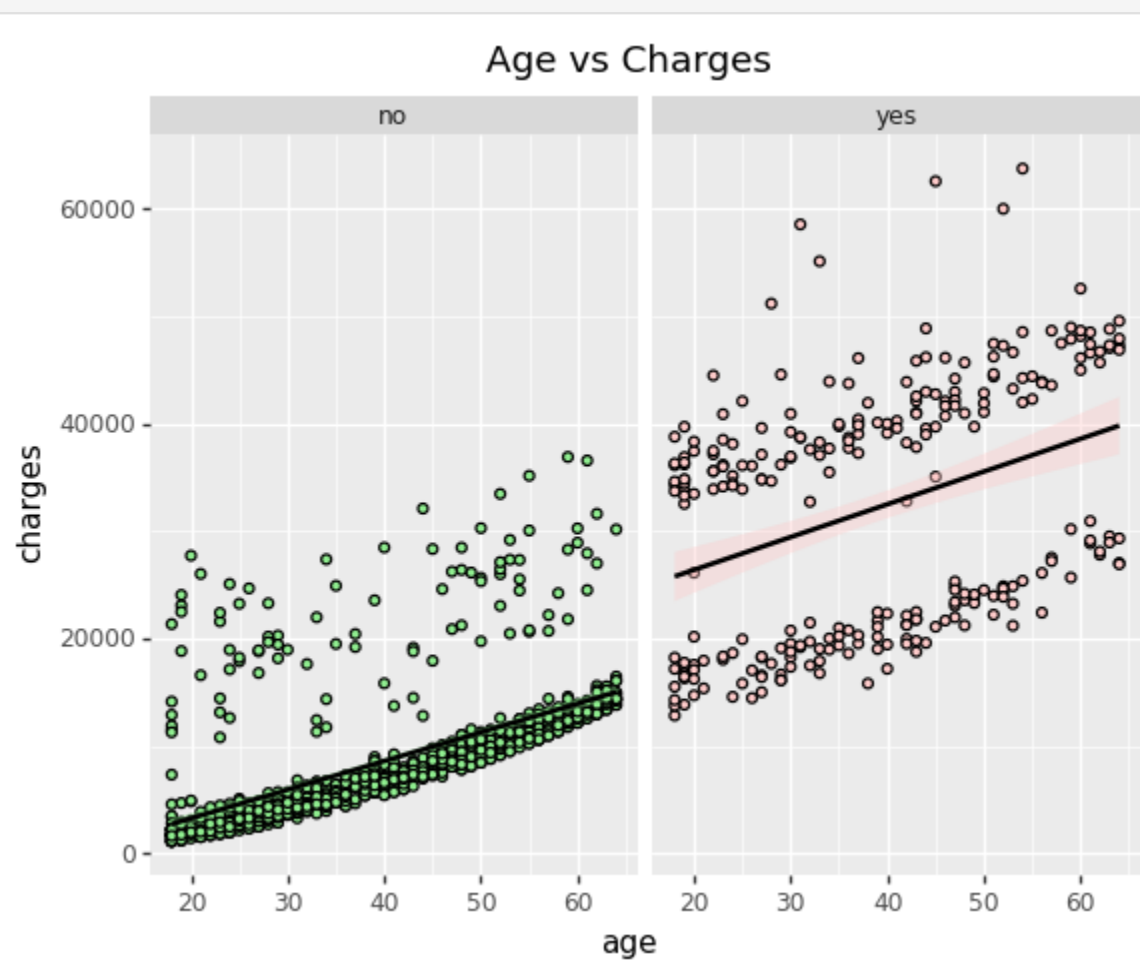
```
In [331]: scatter_plots("bmi", "BMI vs Charges")
```



```
Out[331]: plt.plotnine.theme('ggplot', aes=geom_point)
```

- For non-smokers, the data trend remains constant.
- While for smokers the trend line is linear.

```
In [327]: scatter_plots("age", "Age vs Charges")
```



```
Out[327]: plt.plotnine.theme('ggplot', aes=geom_point)
```

- For non-smokers, the trend of the data remains linear with respect to age.
- For smokers, the relationship between age remains non-existent.

```
In [349]: corr_pearson=df.corr(method="pearson")
corr_spearman=df.corr(method="spearman")
```

```
def correlation_matrix(corr_matrix):

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

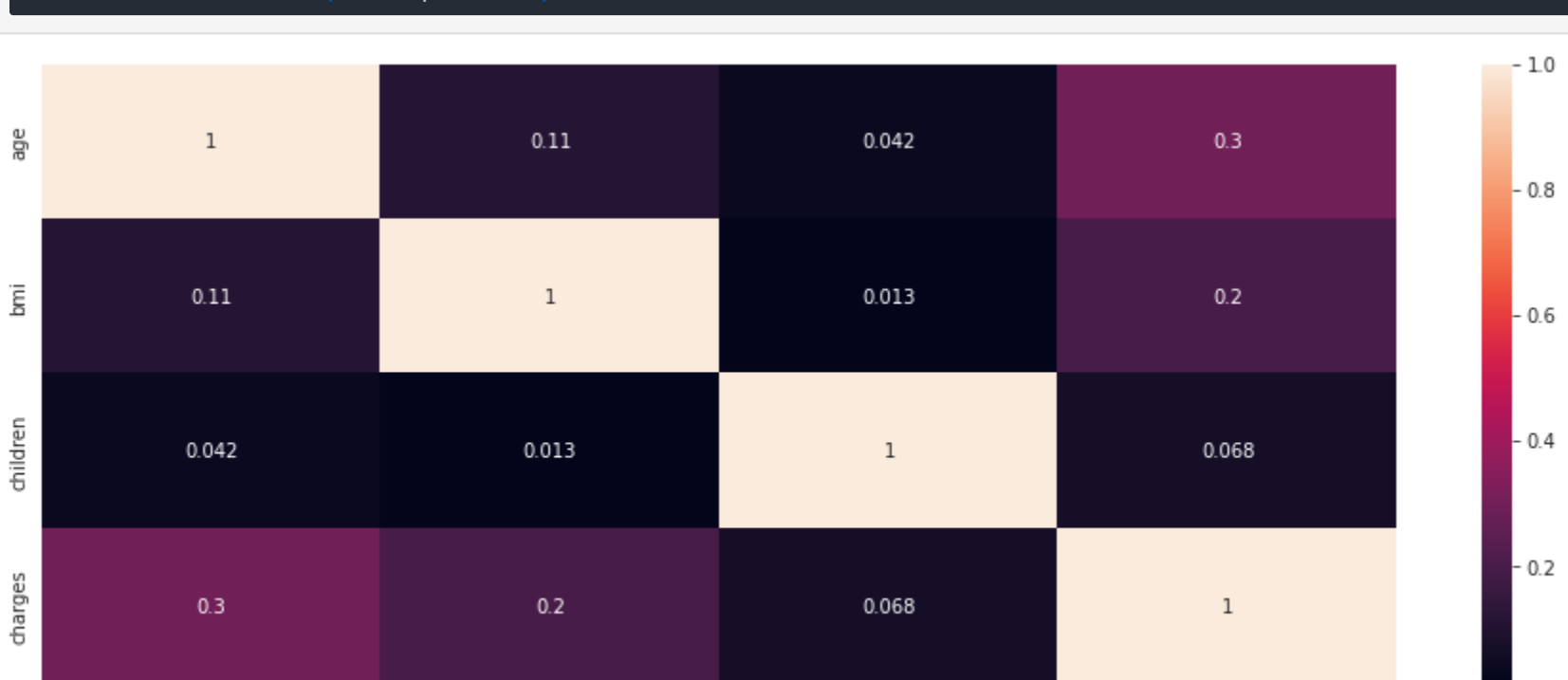
    plt.show()
```

Correlation Matrix

It is used to establish possible relationships between variables

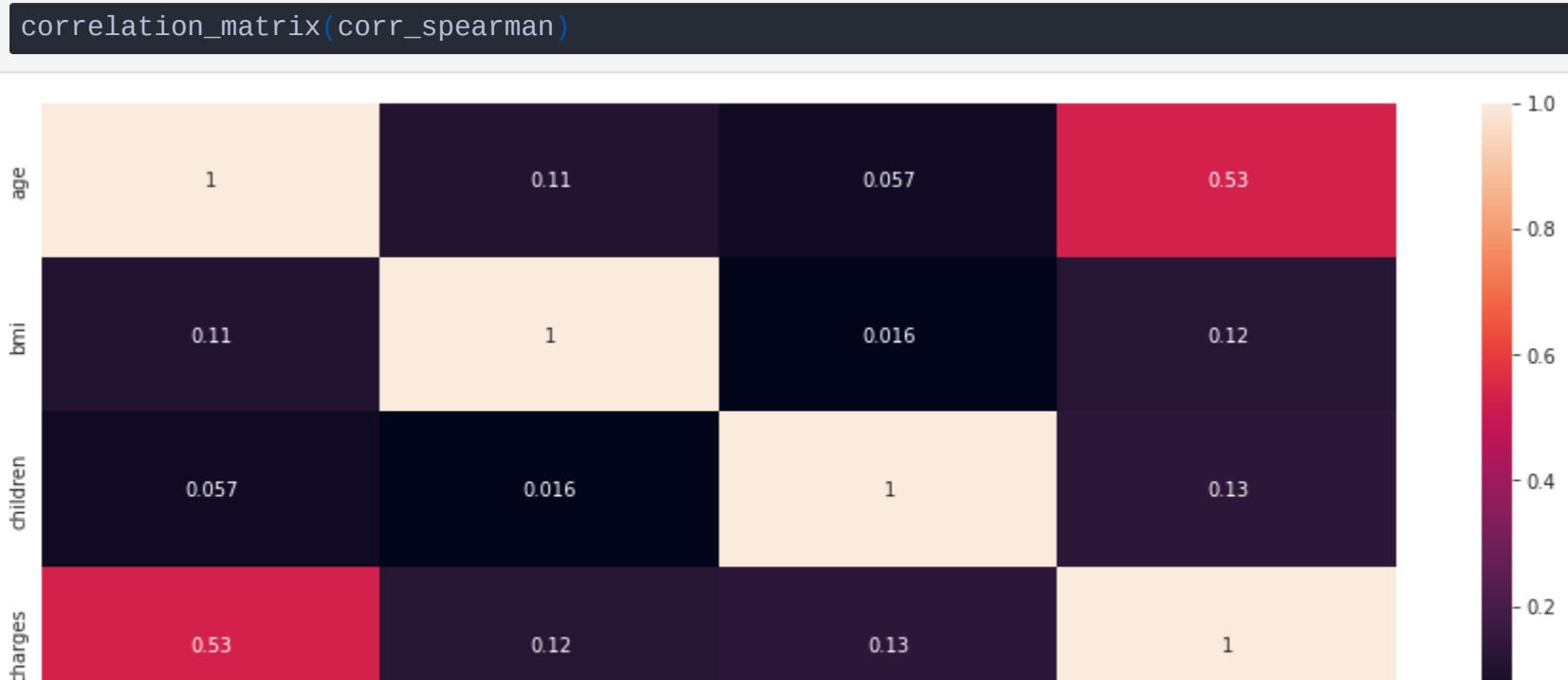
Pearson

```
In [350]: correlation_matrix(corr_pearson)
```



Spearman

```
In [352]: correlation_matrix(corr_spearman)
```



Negative r values indicate a negative correlation, in which the values of one variable tend to increase while the values of the other variable decrease. The values 1 and -1 represent positive and negative "perfect" equivalence, respectively. Although the correlations are not so insignificant. They can serve as complementary variables; to generate more accurate predictions.