

## Import libraries

```
import pandas as pd # data manipulation
import numpy as np # linear algebra
import matplotlib.pyplot as plt # plots
import seaborn as sns # plots
import warnings # ignore warning messages
```

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

## Load Data

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

Using the histogram and the box plot. We confirm the presence of outliers. So we have to give it special processing.

Technically we can give outliers the same treatment as missing values.

- Delete those values.
- Replace them with a statistical measure.
- Use a model. So that it generates values closer to the real ones.

## Feature engineering

```
class lower_upper_limits():
    def __init__(self,dataset,feature,limit):
        self.dataset=dataset
        self.mean=dataset[feature].mean()
        self.std=dataset[feature].std()
        self.limit=limit
    def upper_limit(self):
        return self.mean+self.limit*self.std
    def lower_limit(self):
        return self.mean-self.limit*self.std
```

```
df.groupby("sex").describe()["bmi"]
```

	count	mean	std	min	25%	50%	75%	max
sex								
female	662.0	30.377749	6.046023	16.815	26.125	30.1075	34.31375	48.07
male	676.0	30.943129	6.140435	15.960	26.410	30.6875	34.99250	53.13

Selected from the best interval

```
def selected_interval_upper_dataframe(feature,limits):
    limit in limits
    upper_limit=lower_upper_limits(dataframe,feature,limit).upper_limit()
    print(limit,upper_limit)
def selected_interval_lower_dataframe(feature,limits):
    limit in limits
    upper_limit=lower_upper_limits(df,"bmi",limit).lower_limit()
    print(limit,upper_limit)
```

## Upper limit

```
upper_intervals=np.arange(0,5,step=1)
upper_intervals=upper_intervals.round(1)
```

```
selected_interval_upper_df["bmi"]=upper_intervals
```

0.0	46.463457169443460
0.2	44.8794888686837
0.4	42.29884444805118
0.6	40.51888281135198
0.8	47.738326213687786
0.9	48.48787786682889
1.2	50.177594778559354
1.4	51.3872323686952
1.6	52.618287743921
1.8	53.618877125868

With an interval of 3 it gives a good value. To be able to replace outliers.

Replaces values greater than the upper range

```
df["bmi"]=np.where(df["bmi">>upper_intervals["bmi"],df["bmi"]
```

## Lower limit

Replaces values less than the lower range

```
lower_intervals=np.arange(0,5,step=1)
lower_intervals=upper_intervals.round(1)
```

```
selected_interval_lower_df["bmi"]=lower_intervals
```

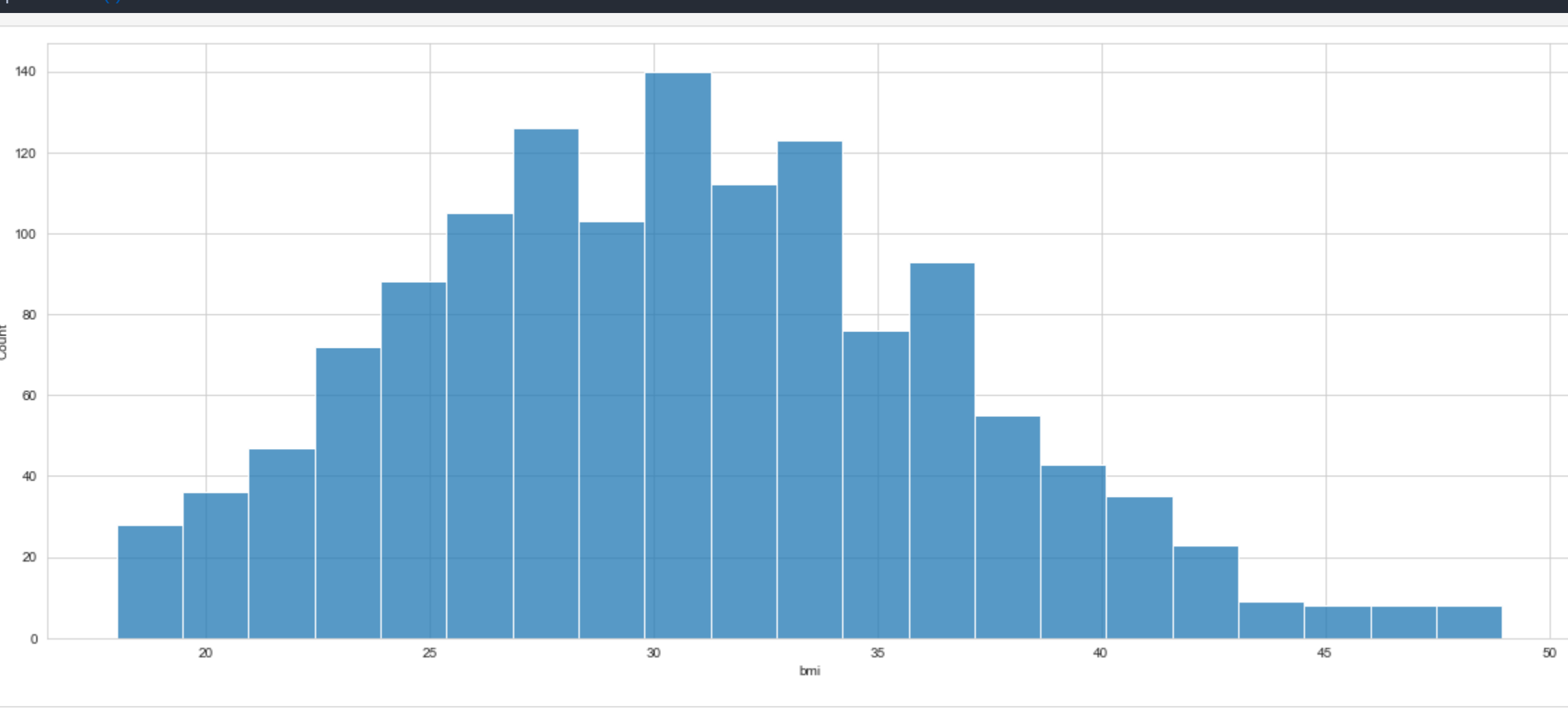
0.0	16.86888277913177
0.2	17.82488921947952
0.4	18.879813758742888
0.6	14.86888489543622
0.8	13.6458522928444
0.9	12.43881398118252
1.2	11.21586707184867
1.4	10.88884482748883
1.6	8.7882328241899
1.8	7.570789244348514

The appropriate interval value is 2. Since considered with the minimum BMI values regardless of sex is 18.

```
df["bmi"]=np.where(df["bmi"]<2,18,df["bmi"])
```

```
sns.set_style(style="whitegrid")
```

```
plt.subplots(figsize=(10,1))
sns.histplot(data=df.query("smoker=='no'")
```



```
smoker_no_split=df.query("smoker=='no'")
smoker_yes_split=df.query("smoker=='yes'")
```

We still see outliers. Therefore, I will decide to divide the dataset based on the age of the user. For better cleaning.

## Smoker no split

### Upper limit

```
lower_upper_limits(smoker_no_split,"charges").upper_limit()
```

```
smoker_no_split["charges"]
```

Using an interval of 2 gives good results. In addition, for this case, the values that are out of the normal, we are going to transform them by null values and later we are going to replace said value, using a linear model. Since there are variables that are highly correlated.

Transform the outliers to null

```
smoker_no_split["charges"]=smoker_no_split["charges"].apply(lambda x: np.nan if x>upper_limit else x)
```

Calculate the percentage of null values

```
(smoker_no_split.isnull().sum()/len(df))*100).sort_values(ascending=False)
```

```
charges
```

```
age
```

```
sex
```

```
smoker
```

```
region
```

```
children
```

```
smoker
```

```
region
```

```
children
```

```
smoker
```

```
region
```

```
children
```

```
smoker
```

```
region
```

```
children
```

```
smoker
```

```
region
```

```
children
```

```
smoker
```

```
region
```

```
children
```

```
smoker
```

```
region
```

```
children
```

```
smoker
```

```
region
```

```
children
```

```
smoker
```

```
region
```

```
children
```

```
smoker
```

```
region
```

```
children
```

```
smoker
```

```
region
```

```
children
```

```
smoker
```

```
region
```

```
children
```

```
smoker
```

```
region
```

```
children
```

```
smoker
```

```
region
```

```
children
```

```
smoker
```

```
region
```

```
children
```

```
smoker
```

```
region
```

```
children
```

```
smoker
```

```
region
```

```
children
```

```
smoker
```

```
region
```

```
children
```

```
smoker
```

```
region
```

```
children
```

```
smoker
```

```
region
```

```
children
```

```
smoker
```

```
region
```

```
children
```

```
smoker
```

```
region
```

```
children
```

```
smoker
```

```
region
```

```
children
```

```
smoker
```

```
region
```

```
children
```

```
smoker
```

```
region
```

```
children
```

```
smoker
```

```
region
```

```
children
```

```
smoker
```

```
region
```

```
children
```

```
smoker
```

```
region
```

```
children
```

```
smoker
```

```
region
```

```
children
```

```
smoker
```

```
region
```

```
children
```

```
smoker
```

```
region
```

```
children
```

```
smoker
```

```
region
```

```
children
```

```
smoker
```

```
region
```

```
children
```

```
smoker
```

```
region
```

```
children
```

```
smoker
```

```
region
```

```
children
```

```
smoker
```

```
region
```

```
children
```

```
smoker
```

```
region
```

```
children
```

```
smoker
```

```
region
```

```
children
```

```
smoker
```

```
region
```

```
children
```

```
smoker
```

```
region
```

```
children
```

```
smoker
```

```
region
```

```
children
```

```
smoker
```

```
region
```

```
children
```

```
smoker
```

```
region
```

```
children
```

```
smoker
```

```
region
```

```
children
```

```
smoker
```

```
region
```

```
children
```

```
smoker
```

```
region
```

```
children
```

```
smoker
```

```
region
```

```
children
```