

```
In [127... import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import warnings
warnings.filterwarnings('ignore')
```

```
In [108... df = pd.read_csv('insurance.csv')
```

```
In [109... df.head()
```

```
Out[109... 
```

	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

```
In [110... 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   bmi         1338 non-null   float64 
3   children    1338 non-null   int64  
4   smoker      1338 non-null   object  
5   region      1338 non-null   object  
6   charges     1338 non-null   float64 
dtypes: float64(2), int64(2), object(3)
memory usage: 73.3+ KB
```

```
In [111... df.shape
```

```
Out[111... (1338, 7)
```

```
In [112... df.columns
```

```
Out[112... Index(['age', 'sex', 'bmi', 'children', 'smoker', 'region', 'charges'], dtype
='object')
```

```
In [113... df.describe
```

```
Out[113... <bound method NDFrame.describe of
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
...      ...      ...      ...      ...      ...      ...
1333      50   male  30.970      3  no   northwest  10600.54830
1334      18  female  31.920      0  no   northeast   2205.98080
1335      18  female  36.850      0  no   southeast   1629.83350
1336      21  female  25.800      0  no   southwest   2007.94500
1337      61  female  29.070      0  yes  northwest  29141.36030

[1338 rows x 7 columns]>
```

```
In [114... df.describe
```

```
Out[114... <bound method NDFrame.describe of
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
...      ...      ...      ...      ...      ...      ...
1333      50   male  30.970      3  no   northwest  10600.54830
1334      18  female  31.920      0  no   northeast   2205.98080
1335      18  female  36.850      0  no   southeast   1629.83350
1336      21  female  25.800      0  no   southwest   2007.94500
1337      61  female  29.070      0  yes  northwest  29141.36030

[1338 rows x 7 columns]>
```

```
In [115... df.sex.unique()
```

```
Out[115... array(['female', 'male'], dtype=object)
```

```
In [116... df.isnull().sum()
```

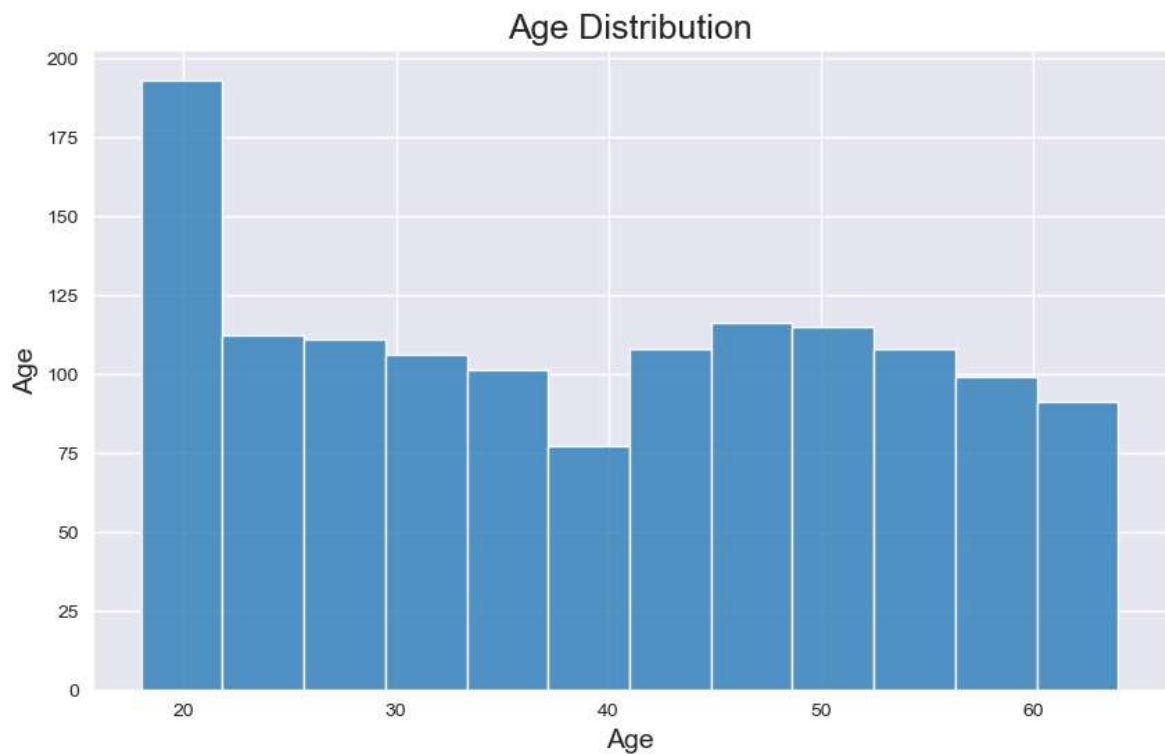
```
Out[116... age      0
sex      0
bmi      0
children 0
smoker   0
region   0
charges  0
dtype: int64
```

```
In [117... df[df.duplicated]
```

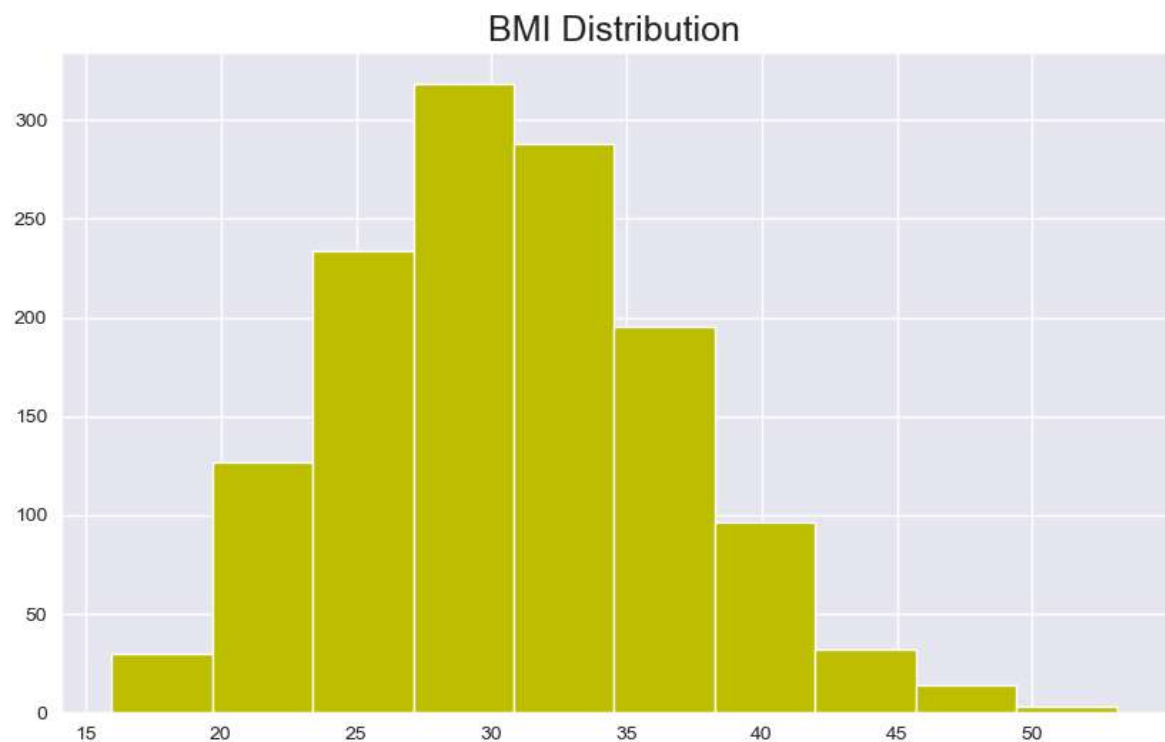
```
Out[117...    age  sex  bmi  children  smoker  region  charges
581   19  male  30.59         0     no  northwest  1639.5631
```

```
In [118... df.drop_duplicates(keep = 'first', inplace = True)
```

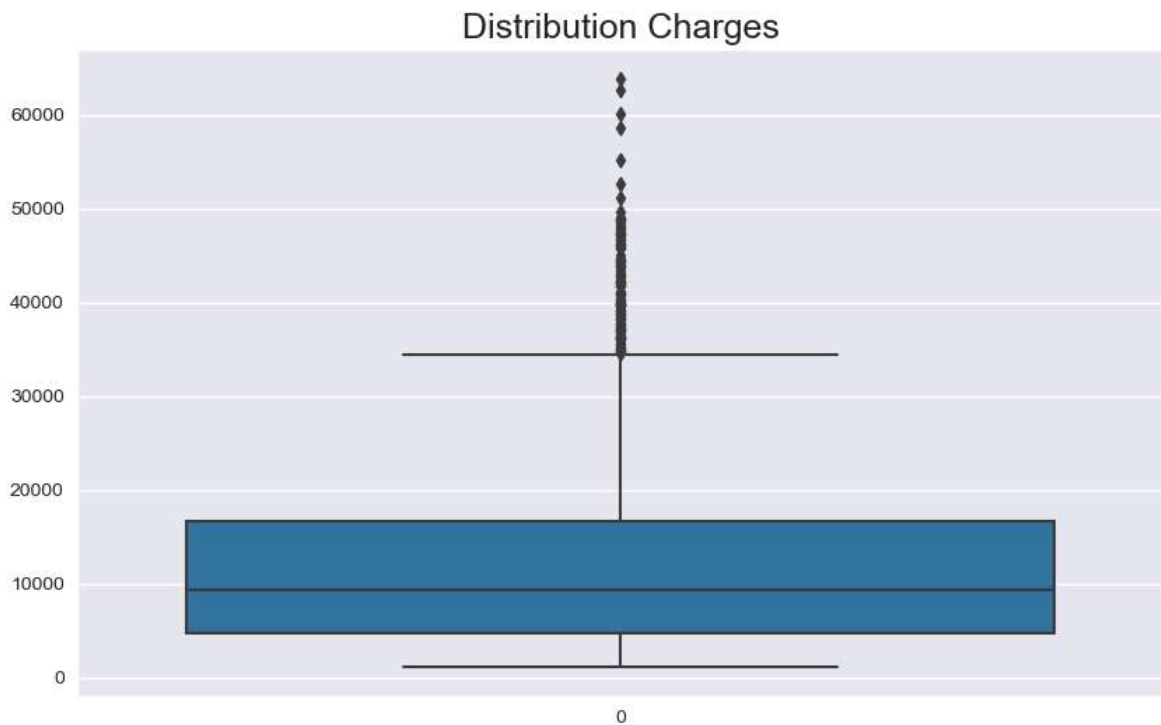
```
In [128... plt.figure(figsize=(10,6))
sns.histplot(df.age)
plt.title('Age Distribution', size=18)
plt.xlabel('Age', size=14)
plt.ylabel('Age', size=14)
plt.show()
```



```
In [120... plt.figure(figsize=(10,6))
plt.hist(df.bmi,color='y')
plt.title('BMI Distribution',size=18)
plt.show()
```



```
In [121... plt.figure(figsize = (10,6))
sns.boxplot(df.charges)
plt.title('Distribution Charges',size=18)
plt.show()
```



```
In [122... Q1 = df['charges'].quantile(0.25)
Q3 = df['charges'].quantile(0.75)
IQR = Q3 - Q1
print(IQR)
```

11911.37345

```
In [101... df[(df['charges'] < Q1 - 1.5 * IQR) | (df['charges'] > Q3 + 1.5 * IQR)]
```

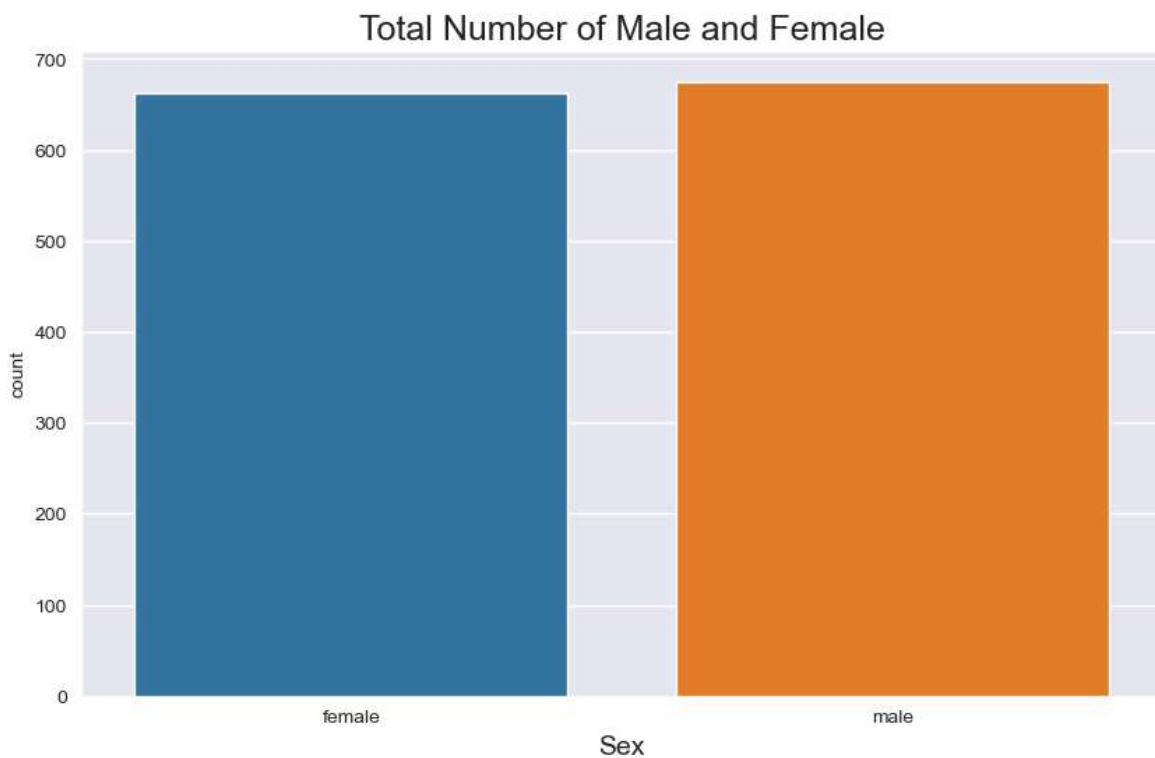
Out[101...

	age	sex	bmi	children	smoker	region	charges
<b>14</b>	27	male	42.130	0	yes	southeast	39611.75770
<b>19</b>	30	male	35.300	0	yes	southwest	36837.46700
<b>23</b>	34	female	31.920	1	yes	northeast	37701.87680
<b>29</b>	31	male	36.300	2	yes	southwest	38711.00000
<b>30</b>	22	male	35.600	0	yes	southwest	35585.57600
...	...	...	...	...	...	...	...
<b>1300</b>	45	male	30.360	0	yes	southeast	62592.87309
<b>1301</b>	62	male	30.875	3	yes	northwest	46718.16325
<b>1303</b>	43	male	27.800	0	yes	southwest	37829.72420
<b>1313</b>	19	female	34.700	2	yes	southwest	36397.57600
<b>1323</b>	42	female	40.370	2	yes	southeast	43896.37630

139 rows × 7 columns

In [102...

```
plt.figure(figsize=(10,6))
sns.countplot(x = 'sex', data = df)
plt.title('Total Number of Male and Female',size=18)
plt.xlabel('Sex',size=14)
plt.show()
```



In [103...

```
df.smoker.value_counts()
```

```
Out[103... smoker  
no      1063  
yes      274  
Name: count, dtype: int64
```

```
In [132... plt.figure(figsize = (10,6))  
sns.set_style('darkgrid')  
sns.boxplot(x='smoker',y='charges',data=df)  
plt.title('Smoker vs Charges',size=18);
```

