

## Context

Machine Learning with R by Brett Lantz is a book that provides an introduction to machine learning using R. As far as I can tell, Packt Publishing does not make its datasets available online unless you buy the book and create a user account which can be a problem if you are checking the book out from the library or borrowing the book from a friend. All of these datasets are in the public domain but simply needed some cleaning up and recoding to match the format in the book.

Content Columns

age: age of primary beneficiary

sex: insurance contractor gender, female, male

bmi: Body mass index, providing an understanding of body, weights that are relatively high or low relative to height, objective index of body weight (kg / m <sup>2</sup> ) using the ratio of height to weight, ideally 18.5 to 24.9

children: Number of children covered by health insurance / Number of dependents

smoker: Smoking

region: the beneficiary's residential area in the US, northeast, southeast, southwest, northwest.

charges: Individual medical costs billed by health insurance

Acknowledgements The dataset is available on GitHub here.

Inspiration Can you accurately predict insurance costs?

```
In [1]: #We will try to analyse and visualize the dataset and try to predict the insurance costs of different i
        individuals
        Fitting regression models and try to understand different variable selection techniques like
        Backward elimination and forward selection
```

## Importing Libraries

```
In [2]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import plotly.express as px
import plotlylib inline
```

## Getting Dataset

```
In [3]: data_file=r'insurance.csv'
```

```
In [4]: dataset=pd.read_csv(data_file)
```

## Getting insights from Dataset

```
In [5]: dataset
```

```
Out[5]:
```

	age	sex	bmi	children	smoker	region	charges
0	19	female	27.900	0	yes	southwest	16884.92400
1	18	male	33.770	1	no	southwest	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	2007.84500
1337	61	female	29.070	0	yes	northwest	29141.36030

1338 rows x 7 columns

```
In [6]: dataset.head()
```

```
Out[6]:
```

	age	sex	bmi	children	smoker	region	charges
0	19	female	27.900	0	yes	southwest	16884.92400
1	18	male	33.770	1	no	southwest	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 [7]: dataset.tail()
```

```
Out[7]:
```

	age	sex	bmi	children	smoker	region	charges
1333	50	male	30.97	3	no	northwest	10600.5483
1334	18	female	31.92	0	no	northeast	2205.9808
1335	18	female	36.85	0	no	southeast	1629.8335
1336	21	female	25.80	0	no	southwest	2007.8450
1337	61	female	29.07	0	yes	northwest	29141.3603

```
In [8]: dataset.info()
```

```
Out[8]:
```

```
<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 [9]: dataset.describe('all')
```

```
Out[9]:
```

	age	sex	bmi	children	smoker	region	charges
count	1338.000000	1338	1338.000000	1338.000000	1338	1338	1338.000000
unique	NaN	male	NaN	NaN	no	southeast	NaN
top	NaN	male	NaN	NaN	no	southeast	NaN
freq	39.207025	NaN	30.863397	1.094918	NaN	NaN	13270.422265
min	14.049960	NaN	6.068187	1.205493	NaN	NaN	12110.011237
std	18.000000	NaN	15.960000	0.000000	NaN	NaN	1121.873900
25%	27.000000	NaN	26.096250	0.000000	NaN	NaN	4740.287150
50%	30.000000	NaN	30.400000	1.000000	NaN	NaN	9382.033000
75%	51.000000	NaN	34.693750	2.000000	NaN	NaN	16639.912515
max	64.000000	NaN	53.130000	5.000000	NaN	NaN	63770.428010

```
In [10]: dataset.isnull().sum()
```

```
Out[10]:
```

age	0
sex	0
bmi	0
children	0
smoker	0
region	0
charges	0
dtype:	int64

```
In [11]: dataset.columns
```

```
Out[11]: Index(['age', 'sex', 'bmi', 'children', 'smoker', 'region', 'charges'], dtype='object')
```

## Visualization

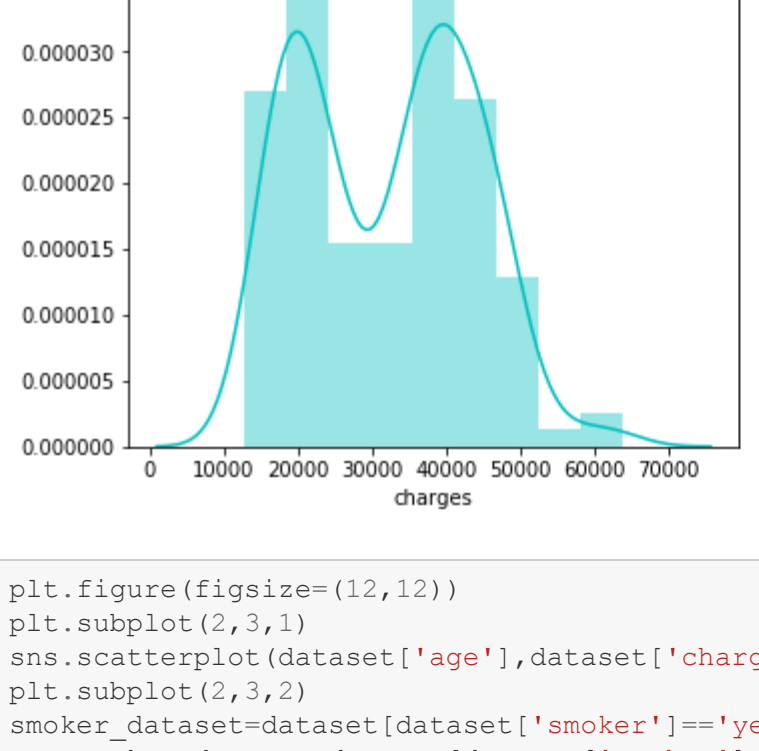
```
In [12]: correlation_plot=dataset.corr()
mask = np.triu(np.ones_like(correlation_plot, dtype=np.bool))
sns.heatmap(correlation_plot,mask=mask,annot=True,cmap='YlOrRd',linewidth=0.5)
```

```
Out[12]: <matplotlib.axes._subplots.AxesSubplot at 0x1cd8f5a550>
```



```
In [85]: correlation_plot=dataset.corr()
mask = np.triu(np.ones_like(correlation_plot, dtype=np.bool))
sns.heatmap(correlation_plot,mask=mask,annot=True,cmap='YlOrRd',linewidth=0.8)
```

```
Out[85]: <matplotlib.axes._subplots.AxesSubplot at 0x1cd8f56760b8>
```



```
In [13]: correlation_plot['charges'].sort_values()
```

```
Out[13]:
```

children	0.067998
bmi	0.198341
age	0.299008
charges	1.000000

Name: charges, dtype: float64

```
In [86]: correlation_plot['charges'].sort_values()
```

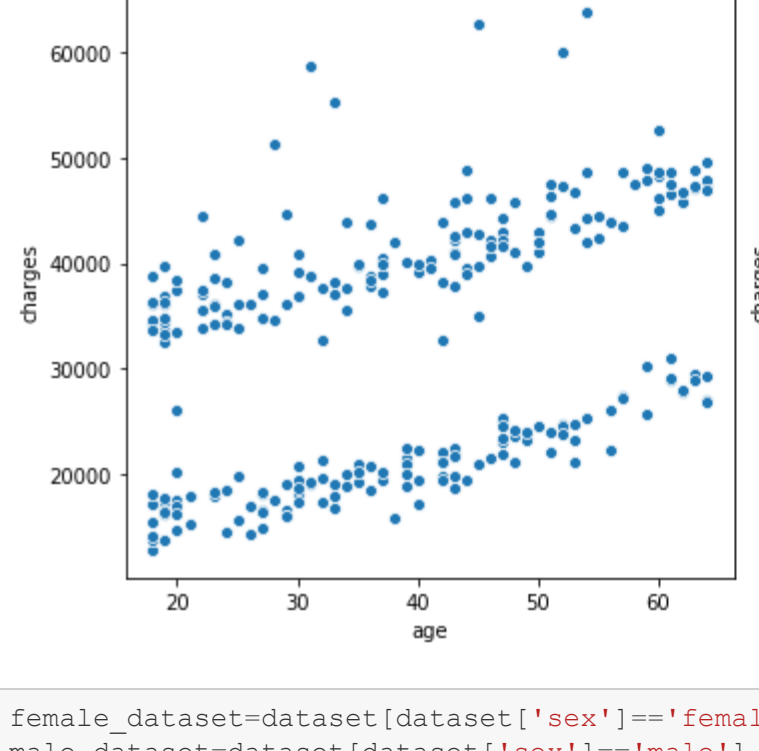
```
Out[86]:
```

region	-0.006208
sex	0.057292
children	0.067998
bmi	0.198341
age	0.299008
smoker	0.787251
charges	1.000000

Name: charges, dtype: float64

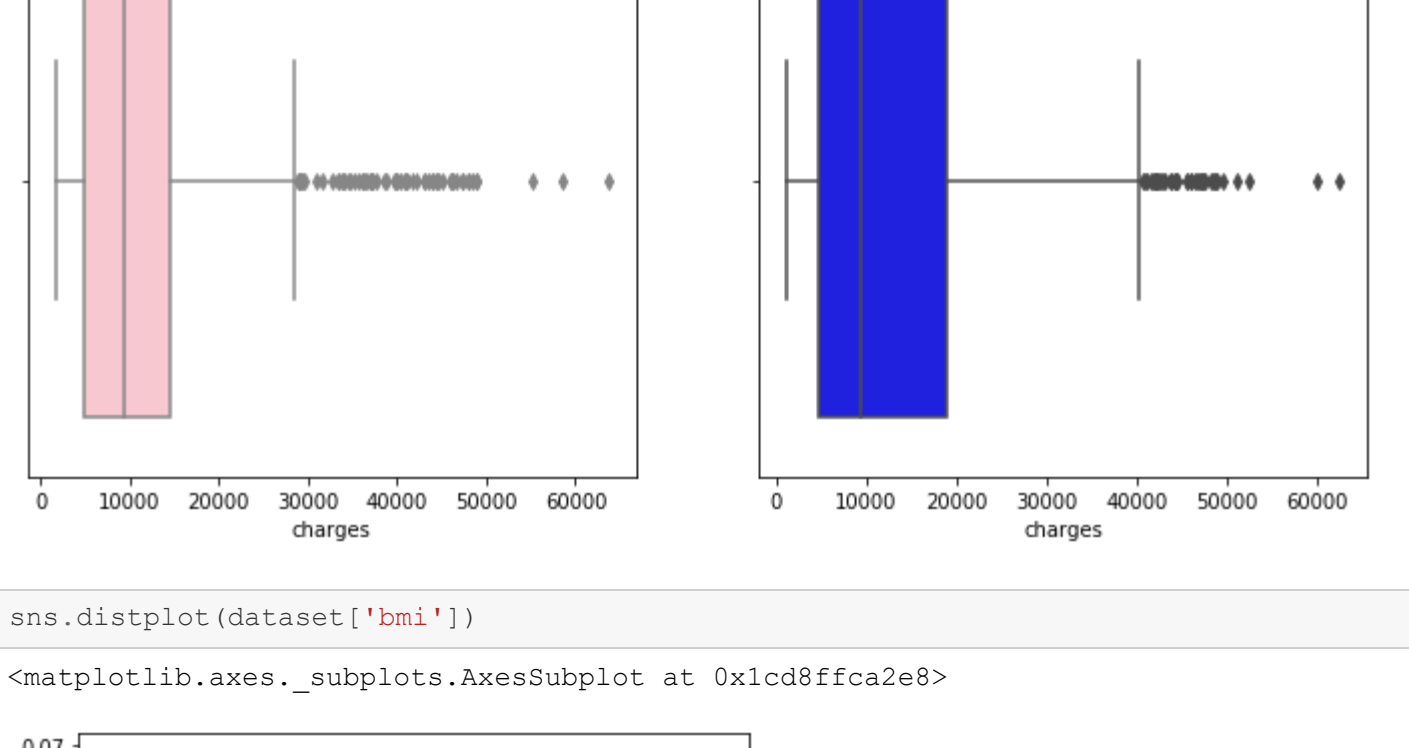
```
In [14]: sns.distplot(dataset['charges'])
```

```
Out[14]: <matplotlib.axes._subplots.AxesSubplot at 0x1cd8f53668>
```



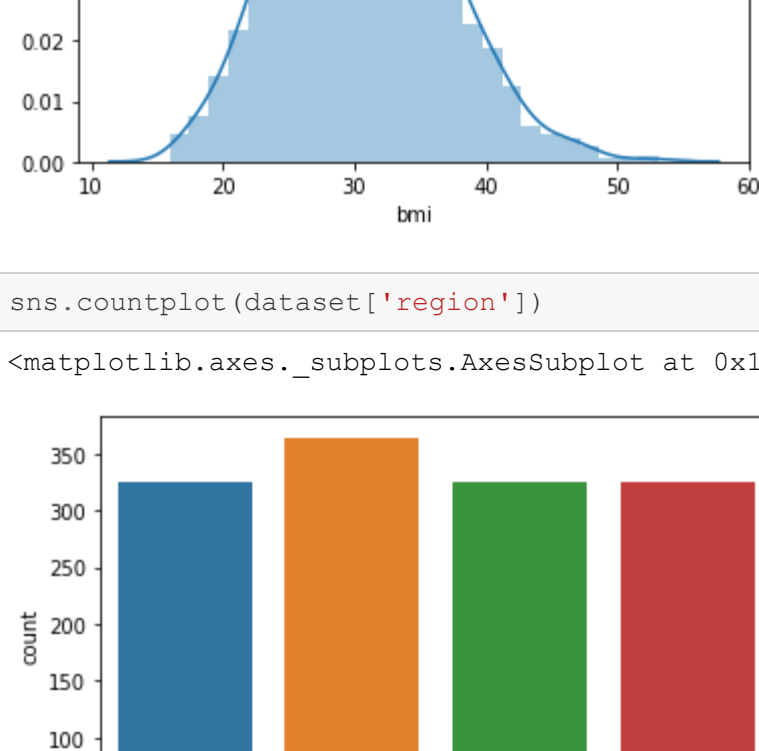
```
In [15]: plt.figure(figsize=(12,12))
plt.subplot(2,2,1)
explode=(0.1,0)
color='blue','orange'
label=['Non-Smoker','Smoker']
dataset['smoker'].value_counts().plot.pie(autopct='%1.2f%%',shadow=True,explode=explode,color=color,label=label)
plt.title('Smokers vs Non-Smokers')
plt.subplot(2,2,2)
color='blue','orange'
dataset['smoker'].value_counts().plot.bar(color=color)
plt.title('Smokers vs Non-Smokers')
```

```
Out[15]: Text(0.5, 1.0, 'Smokers vs Non-Smokers')
```



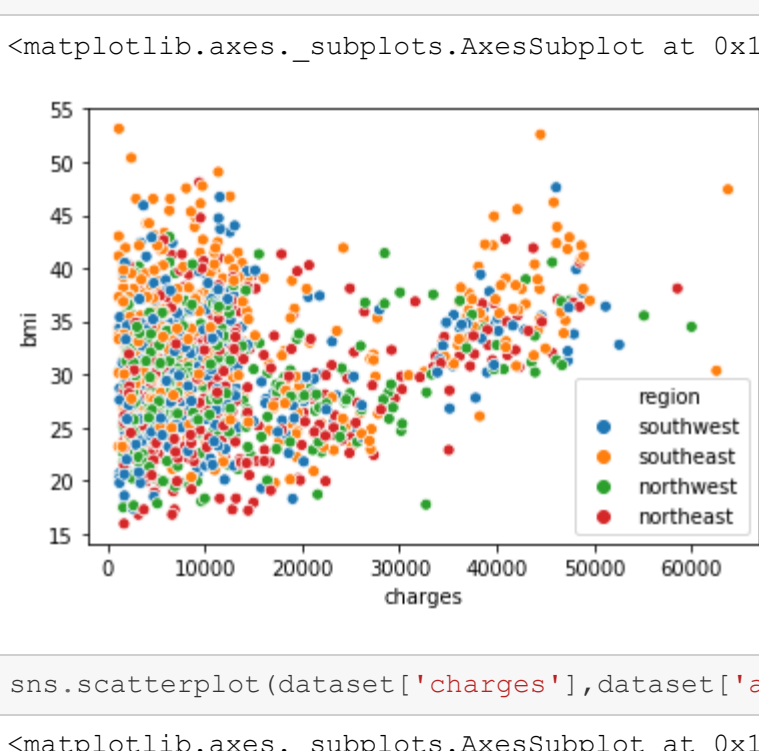
```
In [16]: sns.scatterplot(dataset['charges'],dataset['smoker'])
```

```
Out[16]: <matplotlib.axes._subplots.AxesSubplot at 0x1cd8f5b205d0>
```



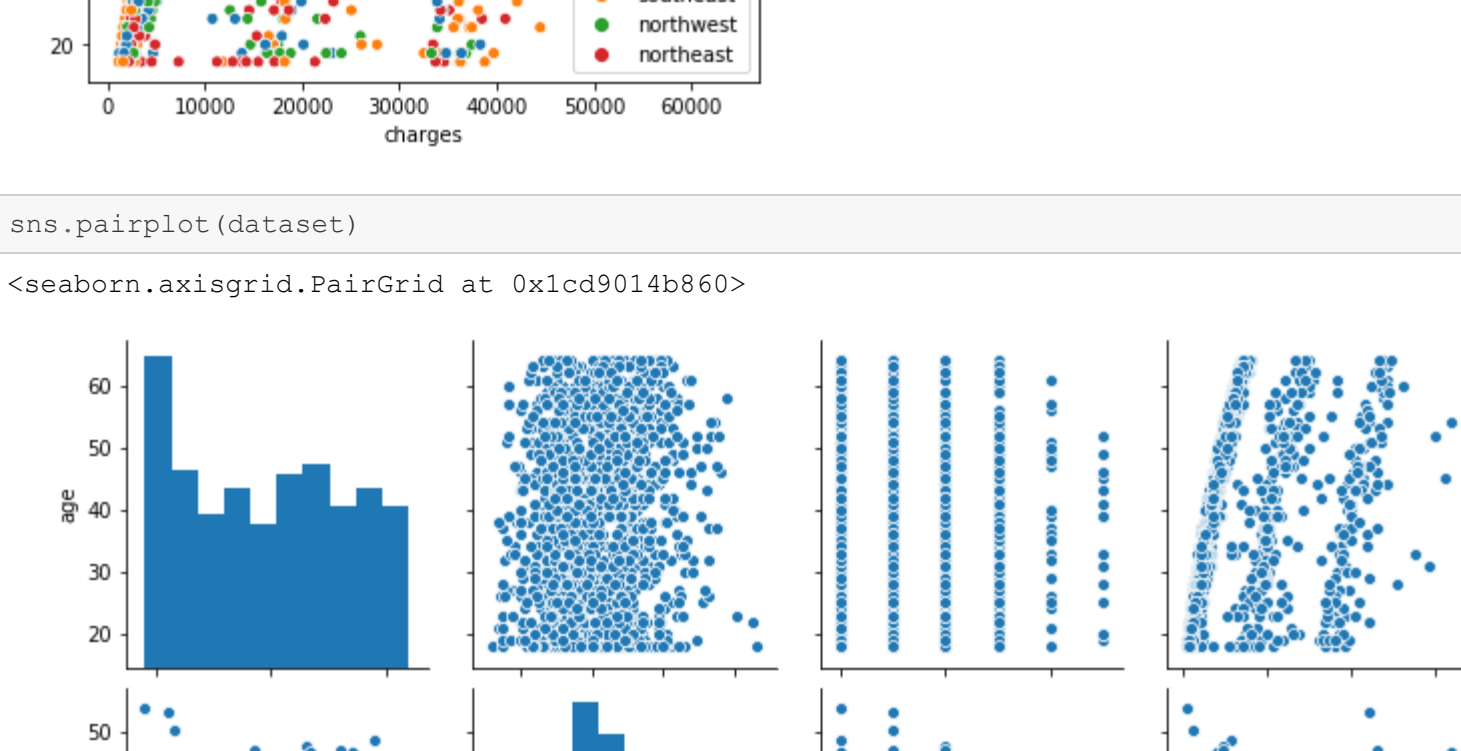
```
In [17]: sns.violinplot(dataset['charges'],dataset['smoker'],hue=dataset['sex'])
```

```
Out[17]: <matplotlib.axes._subplots.AxesSubplot at 0x1cd8f5b4ad30>
```



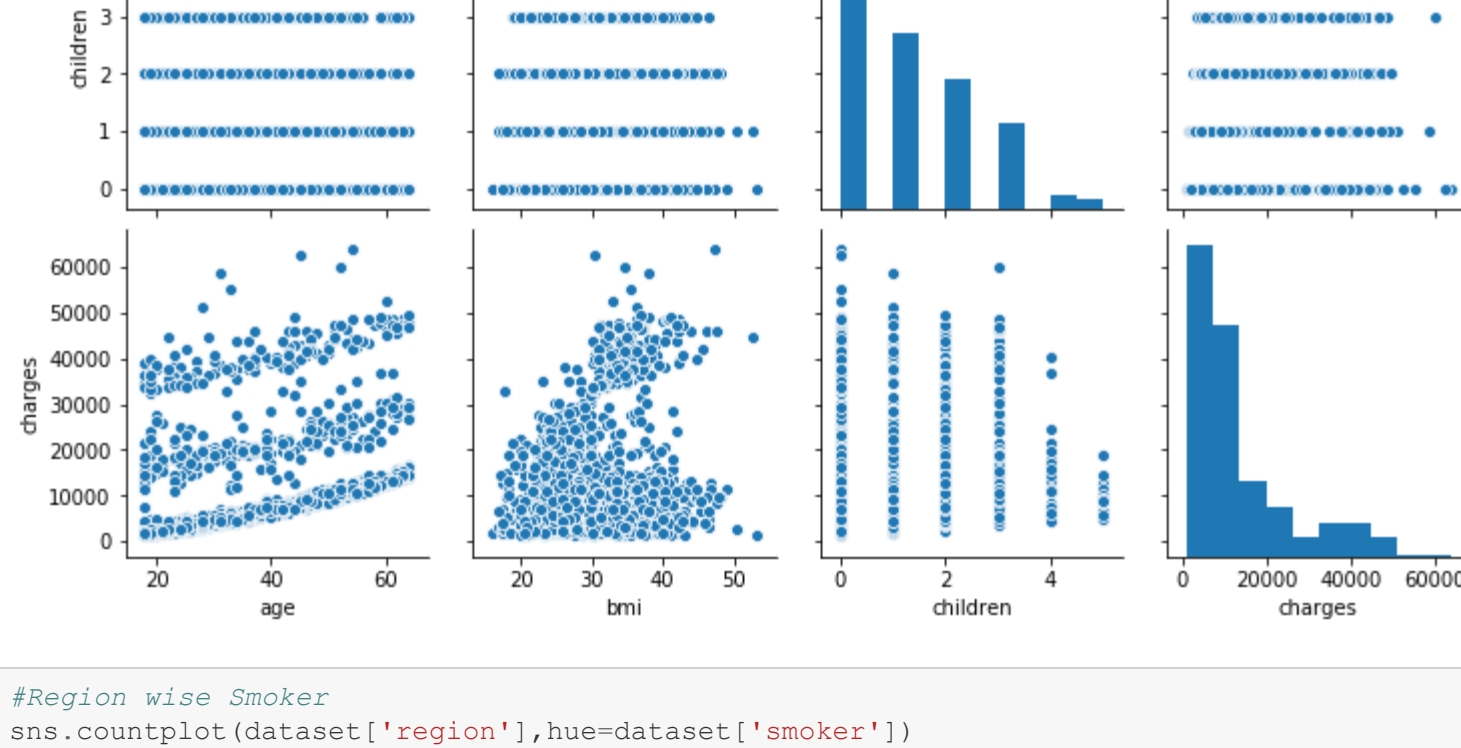
```
In [18]: f=plt.figure(figsize=(12,5))
ax=sns.boxplot(dataset['charges'],dataset['smoker']=='yes')
ax.add_subplot(121)
ax=f.add_subplot(122)
sns.distplot(dataset['charges'],dataset['smoker']=='no')
```

```
Out[18]: <matplotlib.axes._subplots.AxesSubplot at 0x1cd8f5b6d300>
```



```
In [19]: plt.figure(figsize=(12,12))
plt.subplot(2,3,1)
smoker_dataset=dataset[dataset['smoker']=='yes']
non_smoker_dataset=dataset[dataset['smoker']=='no']
sns.scatterplot(smoker_dataset['age'],smoker_dataset['charges'])
plt.subplot(2,3,2)
sns.scatterplot(non_smoker_dataset['age'],non_smoker_dataset['charges'],color='orange')
```

```
Out[19]: <matplotlib.axes._subplots.AxesSubplot at 0x1cd8f5b82400>
```



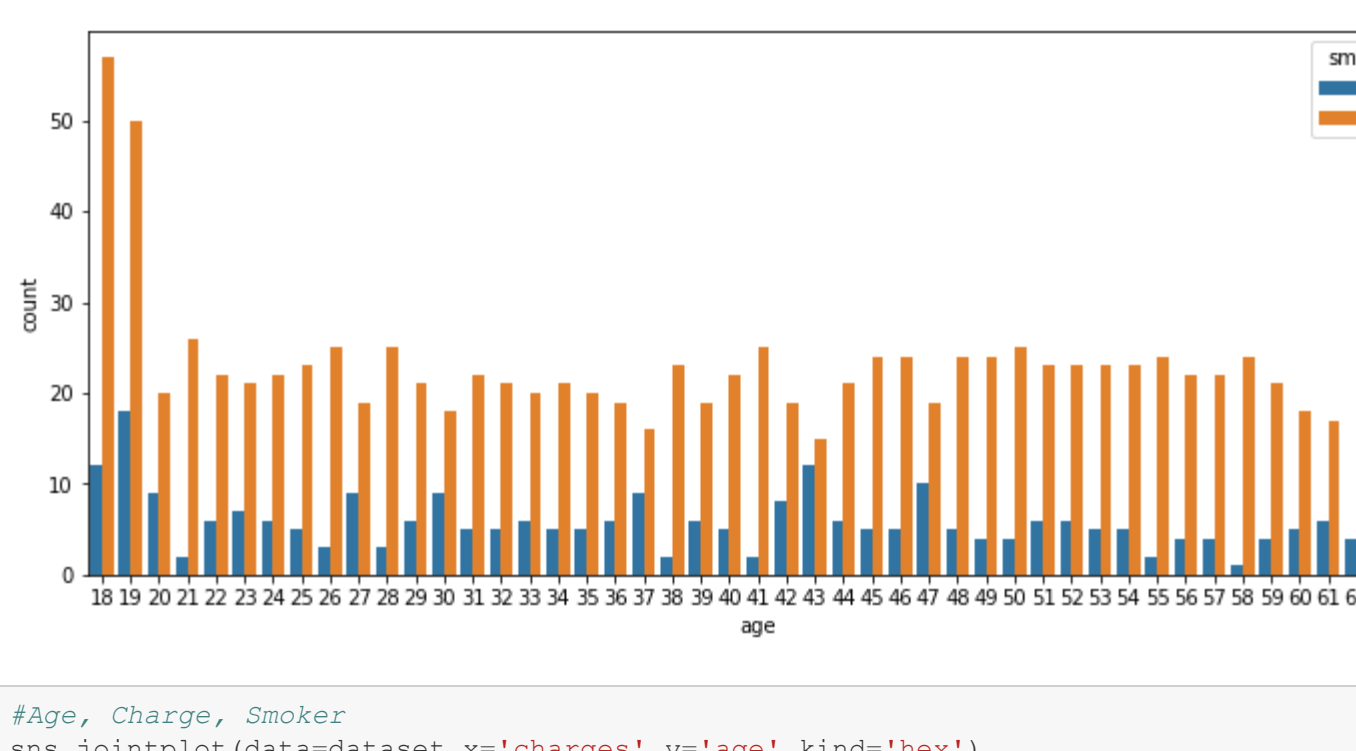
```
In [20]: plt.figure(figsize=(12,12))
smoker_dataset=dataset[dataset['smoker']=='yes']
non_smoker_dataset=dataset[dataset['smoker']=='no']
sns.scatterplot(smoker_dataset['age'],smoker_dataset['charges'])
plt.subplot(2,3,2)
sns.scatterplot(non_smoker_dataset['age'],non_smoker_dataset['charges'],color='orange')
```

```
Out[20]: <matplotlib.axes._subplots.AxesSubplot at 0x1cd8f5b99200>
```



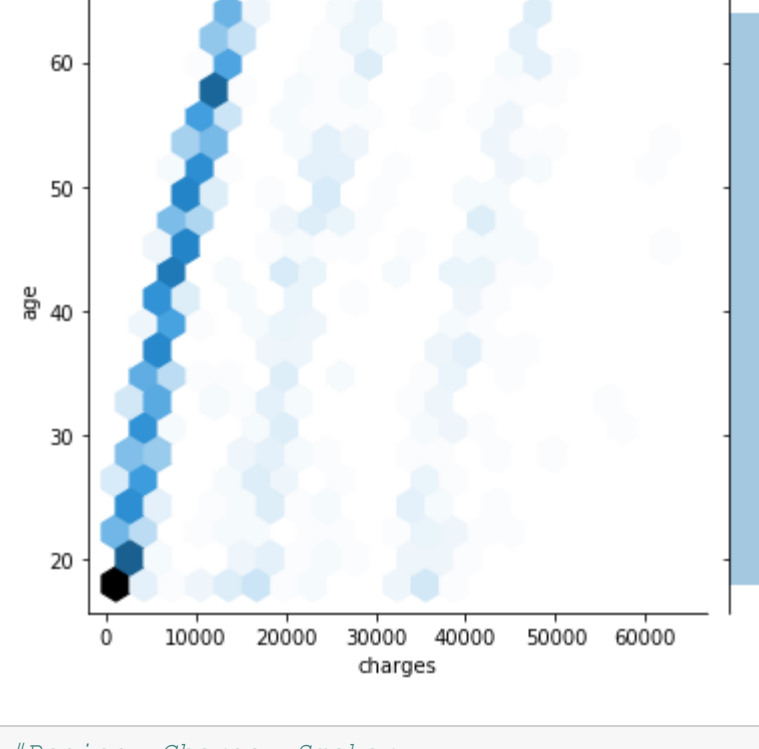
```
In [21]: female_dataset=dataset[dataset['sex']=='female']
male_dataset=dataset[dataset['sex']=='male']
plt.figure(figsize=(12,12))
plt.subplot(2,2,1)
sns.boxplot(female_dataset['charges'],color='pink')
plt.subplot(2,2,2)
sns.boxplot(male_dataset['charges'],color='blue')
```

```
Out[21]: <matplotlib.axes._subplots.AxesSubplot at 0x1cd8f5b85c00>
```



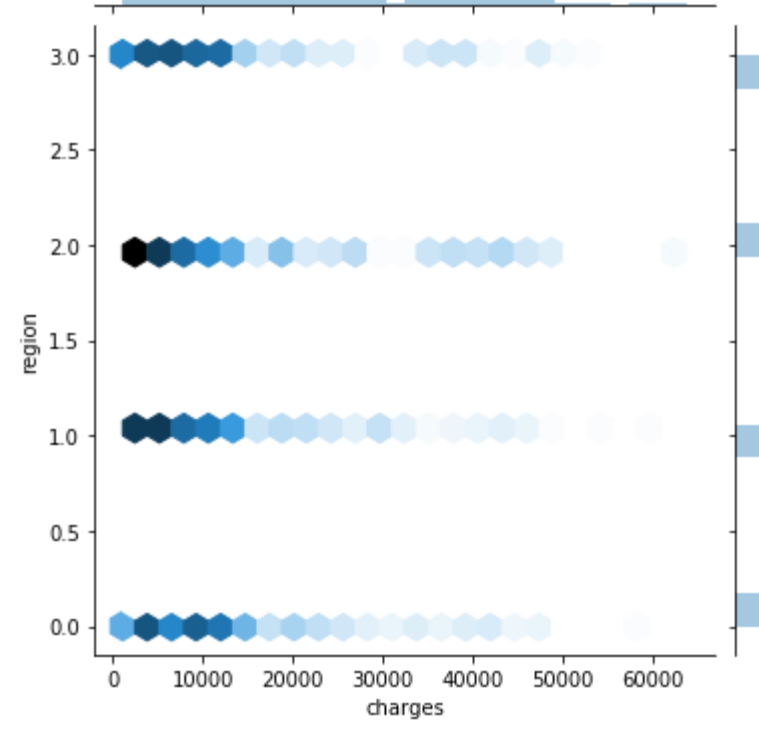
```
In [22]: sns.distplot(dataset['bmi'])
```

```
Out[22]: <matplotlib.axes._subplots.AxesSubplot at 0x1cd8f5b82e80>
```



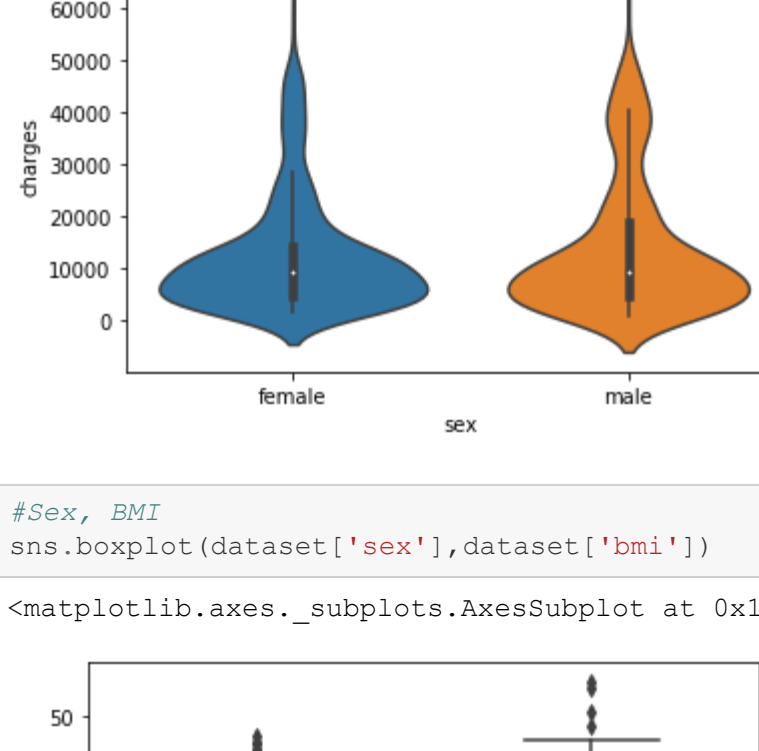
```
In [23]: sns.countplot(dataset['region'])
```

```
Out[23]: <matplotlib.axes._subplots.AxesSubplot at 0x1cd900f93200>
```



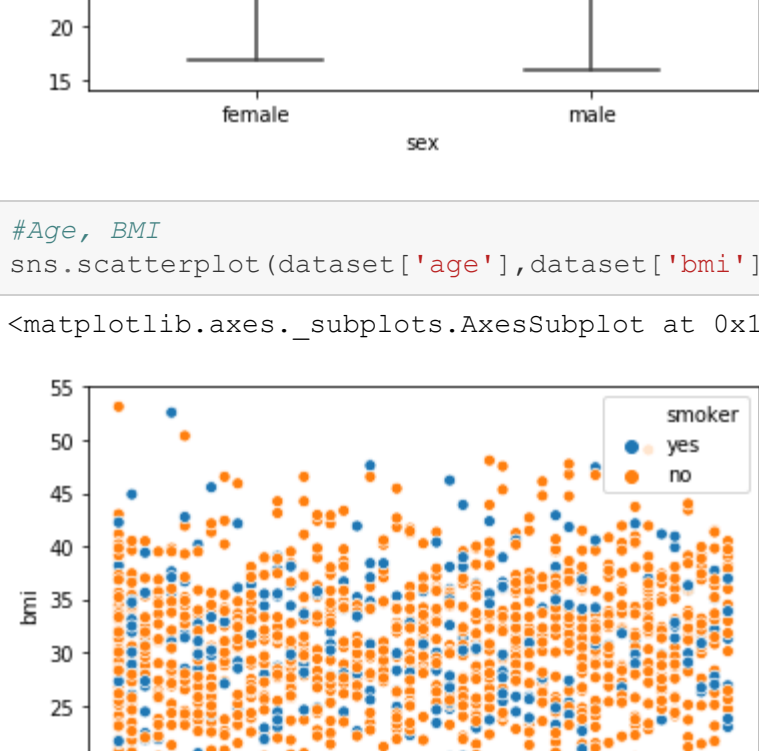
```
In [24]: sns.scatterplot(dataset['charges'],dataset['bmi'],hue=dataset['region'])
```

```
Out[24]: <matplotlib.axes._subplots.AxesSubplot at 0x1cd900f93200>
```

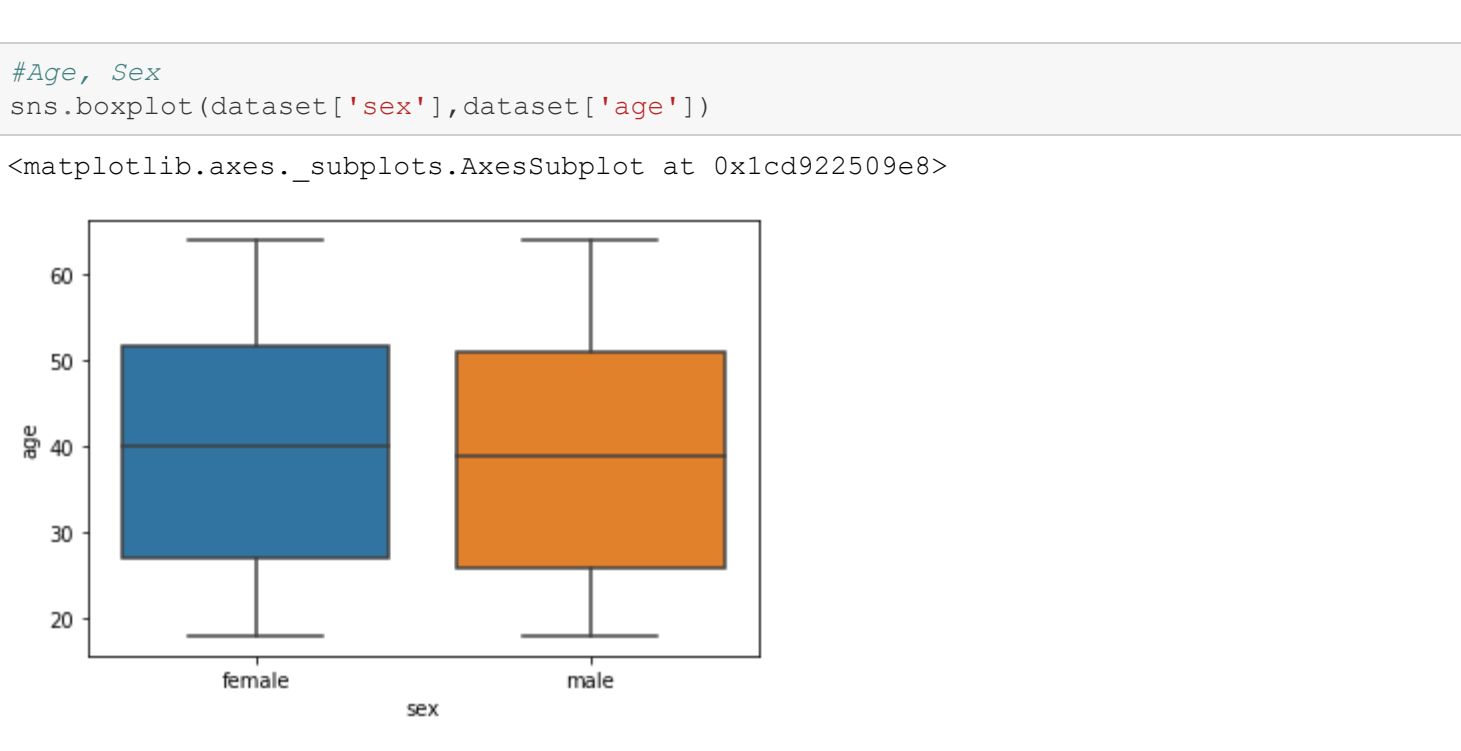


```
In [25]: sns.scatterplot(dataset['charges'],dataset['age'],hue=dataset['region'])
```

```
Out[25]: <matplotlib.axes._subplots.AxesSubplot at 0x1cd900f93200>
```

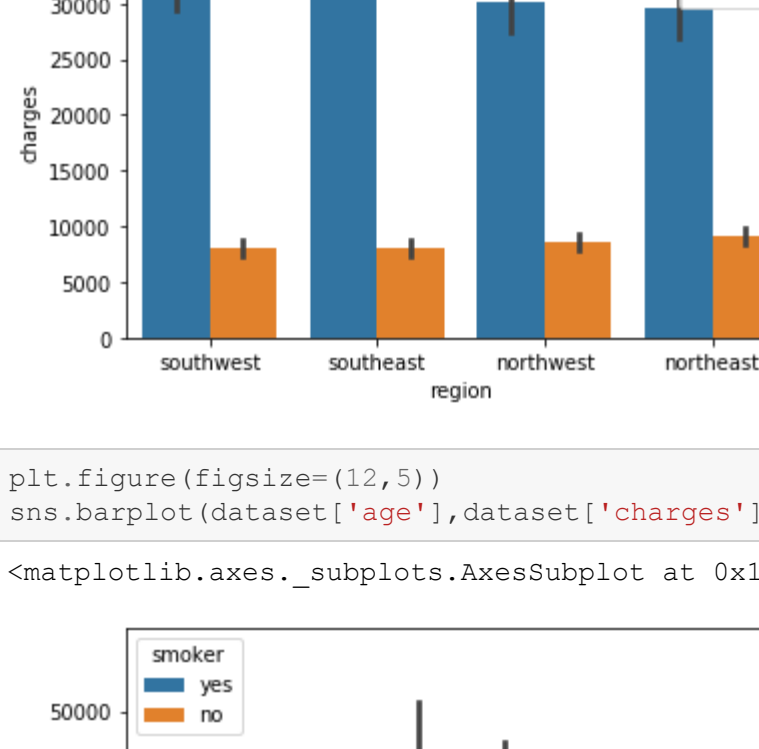


```
In [26]: sns.pairplot(dataset)
```



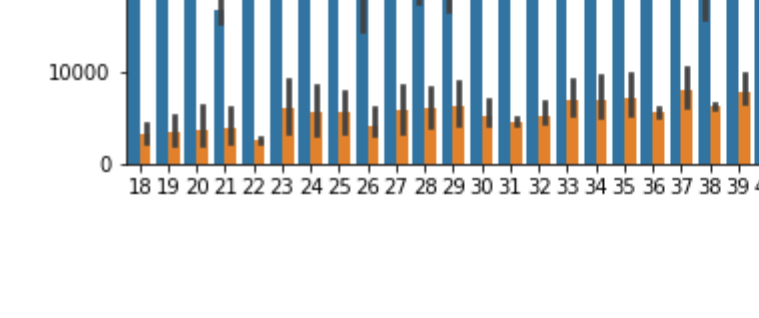
```
In [27]: #Region wise Smoker
sns.countplot(dataset['region'],hue=dataset['smoker'])
```

```
Out[27]: <matplotlib.axes._subplots.AxesSubplot at 0x1cd91a75dd80>
```



```
In [28]: #Sex wise Smoker
sns.countplot(dataset['sex'],hue=dataset['smoker'])
```

```
Out[28]: <matplotlib.axes._subplots.AxesSubplot at 0x1cd91a75dd80>
```



```
In [29]: #Age wise Smoker
plt.figure(figsize=(12,5))
sns.countplot(dataset['age'],hue=dataset['smoker'])
```

```
Out[29]: <matplotlib.axes._subplots.AxesSubplot at 0x1cd91a75dd80>
```



```
In [30]: #Age, Charge, Smoker
sns.jointplot(data=dataset,x='charges',y='age',kind='hex')
```

```
Out[30]: <seaborn.axisgrid.JointGrid at 0x1cd91d4cbb00>
```



```
In [87]: #Region, Charge, Smoker
sns.jointplot(data=dataset,x='charges',y='region',kind='hex')
```

```
Out[87]: <seaborn.axisgrid.JointGrid at 0x1cd9677f5180>
```



```
In [31]: #Sex, Charge
sns.violinplot(dataset['sex'],dataset['charges'])
```

```
Out[31]: <matplotlib.axes._subplots.AxesSubplot at 0x1cd920b85880>
```



```
In [32]: #Sex, BMI
sns.boxplot(dataset['sex'],dataset['bmi'])
```

```
Out[32]: <matplotlib.axes._subplots.AxesSubplot at 0x1cd9215cd880>
```



```
In [33]: #Age, BMI
sns.scatterplot(dataset['age'],dataset['bmi'],hue=dataset['smoker'])
```

```
Out[33]: <matplotlib.axes._subplots.AxesSubplot at 0x1cd920c70700>
```



```
In [34]: #Age, Sex
sns.boxplot(dataset['sex'],dataset['age'])
```

```
Out[34]: <matplotlib.axes._subplots.AxesSubplot at 0x1cd922509e80>
```



```
In [35]: sns.barplot(dataset['region'],dataset['charges'],hue=dataset['smoker'])
```

```
Out[35]: <matplotlib.axes._subplots.AxesSubplot at 0x1cd922b02e80>
```



```
In [36]: plt.figure(figsize=(12,5))
sns.barplot(dataset['age'],dataset['charges'],hue=dataset['smoker'])
```





