

# Feature Engineering

## 1. Univariate analysis

For numerical column - histogram(Distplot | histogram)

For Categorical column - Barplot | Countplot

## 2. Bivariate analysis : Predictor vs Target Variable

Num - Num | Scatterplot, Implot

Cat - Num | Boxplot | Violinplot

Cat - Cat | pd.crosstab | Table

## 3. Missing Values : Mean and Median Imputation | Model based Imputation

```
In [1]: import matplotlib.pyplot as plt

In [2]: x1=[1,6,8,5]
        y1=[5,6,7,8]

In [3]: plt.plot(x1,y1,label="Data")
        plt.show()

Out[3]: <Figure with 1 Axes>

In [4]: x2=y1
        y2=x1

In [5]: plt.plot(x1,y1,label="Data 1",color="r",markers="o",markersize=7) #marker: to mark the position 's' = square, 'o'
        plt.plot(x2,y2,label="Data 2",color="g",linestyle="--",linewidth=4) #linestyle: to get line in diff shape
        plt.legend()
        plt.xlabel("x-axis")
        plt.ylabel("y-axis")
        plt.title("Simple Line Graph")

Out[5]: Text(0.5, 1.0, 'Simple Line Graph')

In [6]: import seaborn as sns

In [7]: d1=sns.load_dataset("tips")

In [8]: d1.head() #Head is used to get an idea abt the dataset like how the data is? To know what all feature we have

Out[8]:
   total_bill  tip    sex  smoker  day  time  size
0    16.99   1.01  Female    No   Sun  Dinner     2
1    10.34   1.66   Male    No   Sun  Dinner     3
2    21.01   3.50   Male    No   Sun  Dinner     3
3    23.68   3.31   Male    No   Sun  Dinner     2
4    24.59   3.61  Female    No   Sun  Dinner     4

In [9]: d1.tail() #Maybe this is an data of some restaurant

Out[9]:
   total_bill  tip    sex  smoker  day  time  size
239   29.03   5.92   Male    No   Sat  Dinner     3
240   27.18   2.00  Female    Yes   Sat  Dinner     2
241   22.67   2.00   Male    Yes   Sat  Dinner     2
242   17.82   1.75   Male    No   Sat  Dinner     2
243   18.78   3.00  Female    No   Thur  Dinner     2

In [10]: #To check what values are there in the data like in time
        d1['time'].unique() #So, data has lunch values also

Out[10]: array(['Dinner', 'Lunch'], dtype=object)
Categories (2, object): ['Dinner', 'Lunch']

In [11]: d1['day'].unique()

Out[11]: array(['Sun', 'Sat', 'Thur', 'Fri'], dtype=object)
Categories (4, object): ['Sun', 'Sat', 'Thur', 'Fri']

In [12]: d1.nunique() #No. of unique values in all columns

Out[12]:
total_bill    229
tip           123
sex            2
smoker         2
day            4
time           2
size          6
dtype: int64

In [13]: d1.count() #Out of tot 244 we have 229 unique values for total bill

Out[13]:
total_bill    244
tip           244
sex            244
smoker         244
day            244
time           244
size           244
dtype: int64

In [14]: #Question : Are this restaurant is prefer by more male or females?

In [15]: # Are more smokers come to restaurant? give some officer for smoker sell ciggretts

In [16]: # Are restaurant is more crowded sat and sun only? if yes give off to staff on friday

In [17]: # Was is it lunch or dinner place?

In [18]: # Is this is high cost restaurant or low cost?

In [19]: # How to attract customers who pay highest bill?
```

So in order to get answers to all these questions we need to identify the data Distributions

How the data are distributed across... We need plot and data to infer the details

We should know which plot to be used for which type of work.

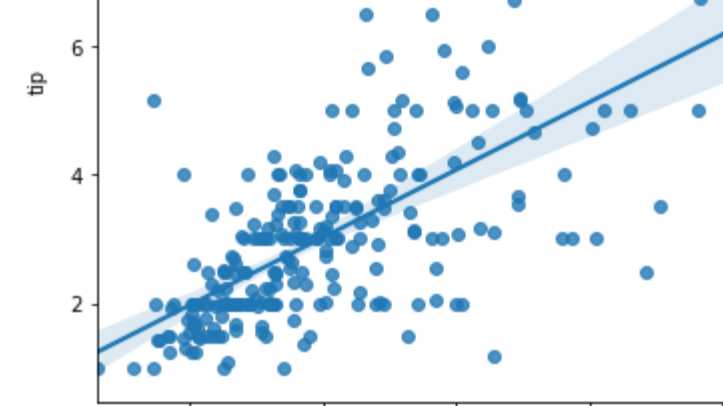
```
In [1]: ## Distribution : histogram, distplot
```

Question: Need to know the relationship b/w the total bill and tip?

```
In [14]: # So, for that we need to do bi-variate analysis : Use Scatter Plot
```

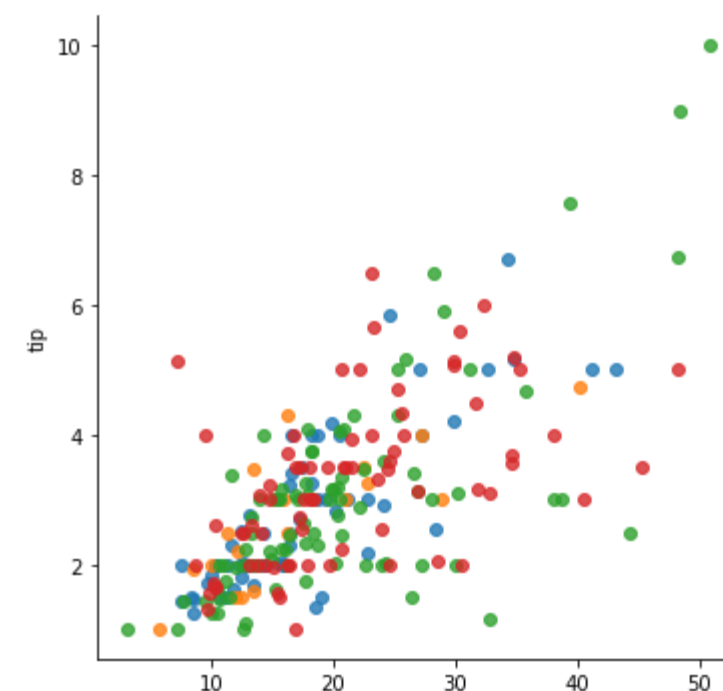
```
In [15]: plt.scatter(x="total_bill",y="tip",data=d1)
```

Out[15]: <matplotlib.collections.PathCollection at 0x2367da94220>



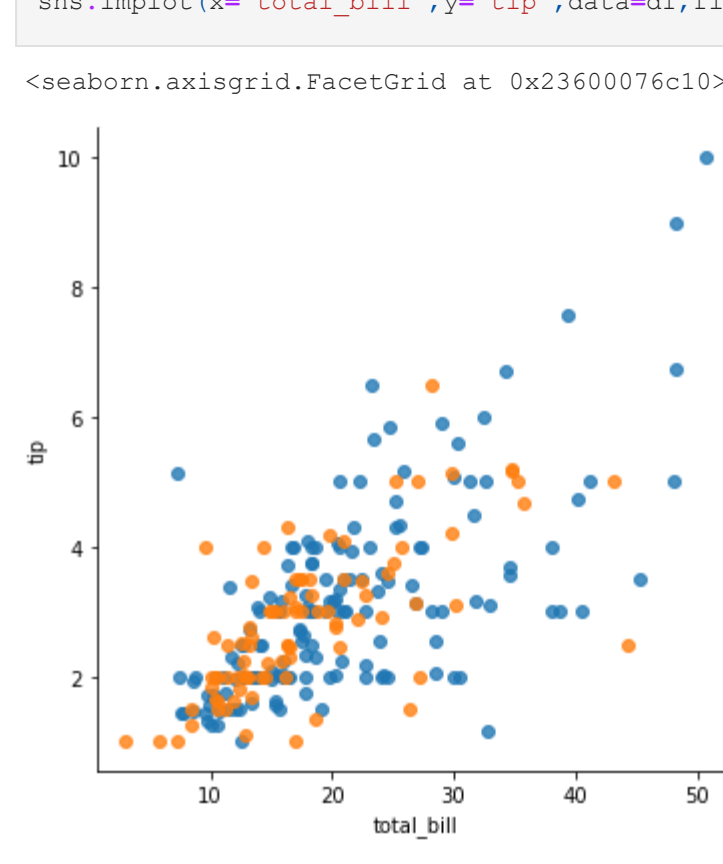
```
In [16]: sns.lmplot(x="total_bill",y="tip",data=d1,fit_reg=True) #Anap used fr ML purposes.
        #lin stands for 'Linear Model' plot, 'fit_reg' = fitting the regression line

Out[16]: <seaborn.axisgrid.FacetGrid at 0x2367dab3340>
```



```
In [17]: # To classify the data based upon which day is more traffic in the restaurant?
        sns.lmplot(x="total_bill",y="tip",data=d1,fit_reg=False,hue="day")

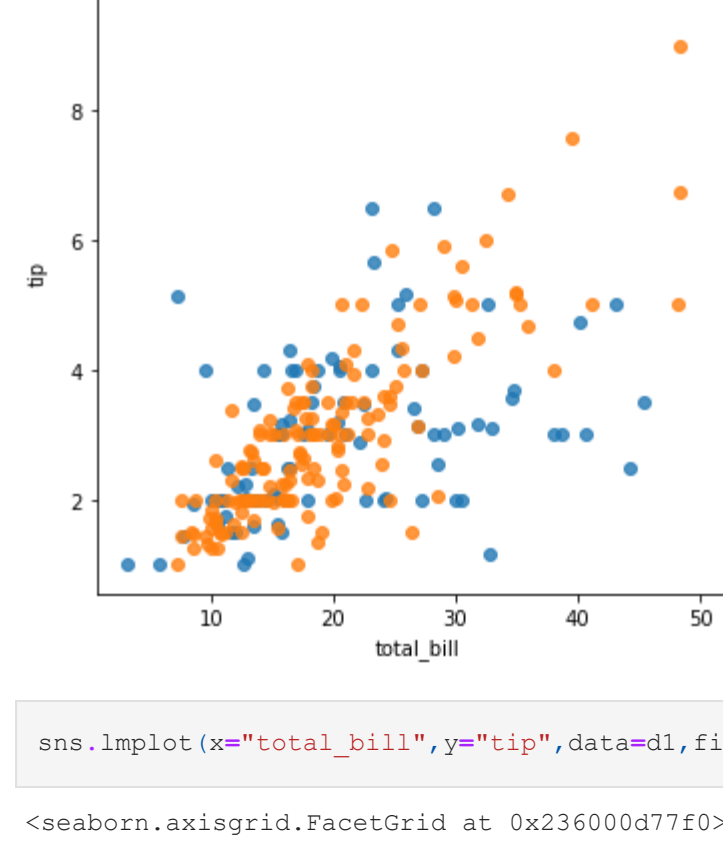
Out[17]: <seaborn.axisgrid.FacetGrid at 0x2367da34e0>
```



1. More traffic on Saturday and Sunday
2. More tip on Saturday and Sunday

```
In [18]: sns.lmplot(x="total_bill",y="tip",data=d1,fit_reg=False,hue="sex")

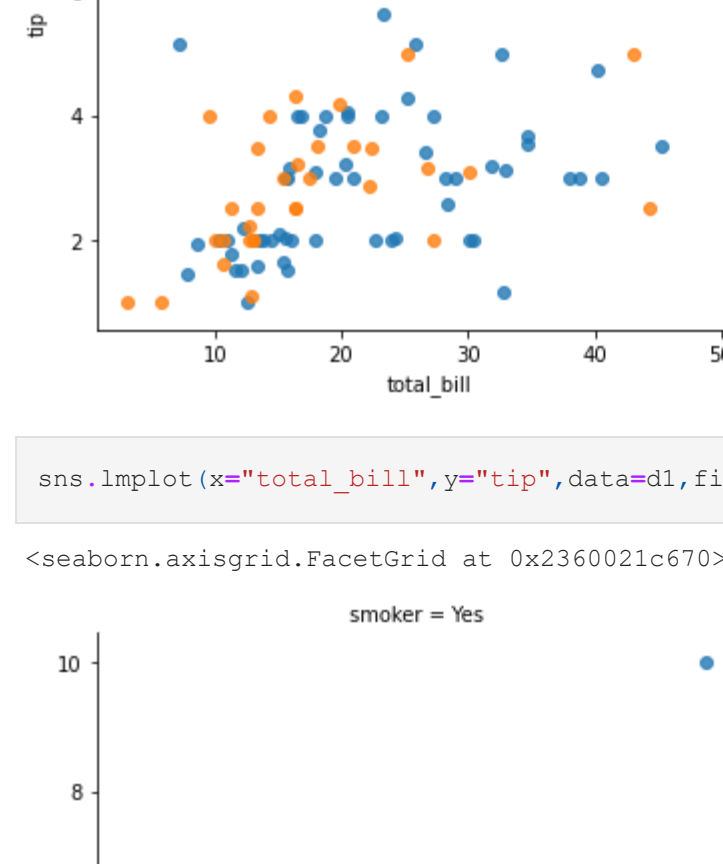
Out[18]: <seaborn.axisgrid.FacetGrid at 0x23600076e10>
```



1. Highest tip was given by 'Male'

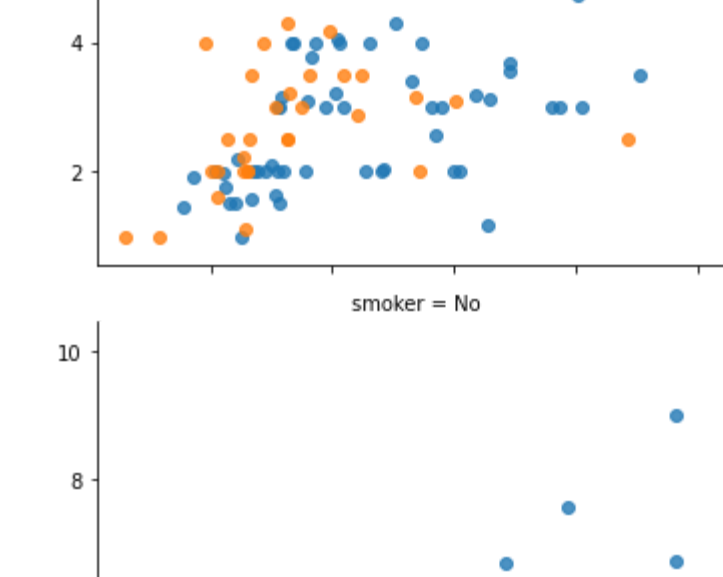
```
In [19]: sns.lmplot(x="total_bill",y="tip",data=d1,fit_reg=False,hue="smoker")

Out[19]: <seaborn.axisgrid.FacetGrid at 0x236000feb50>
```



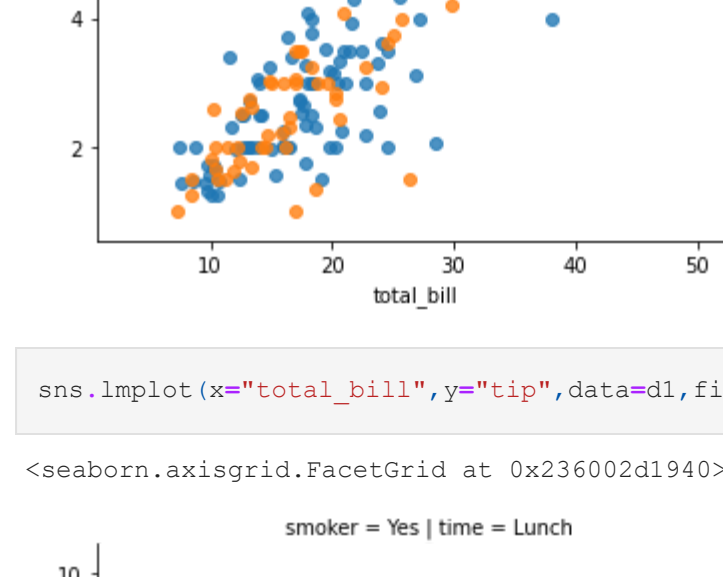
```
In [20]: sns.lmplot(x="total_bill",y="tip",data=d1,fit_reg=False,hue="sex",col="smoker")

Out[20]: <seaborn.axisgrid.FacetGrid at 0x236000d77f0>
```



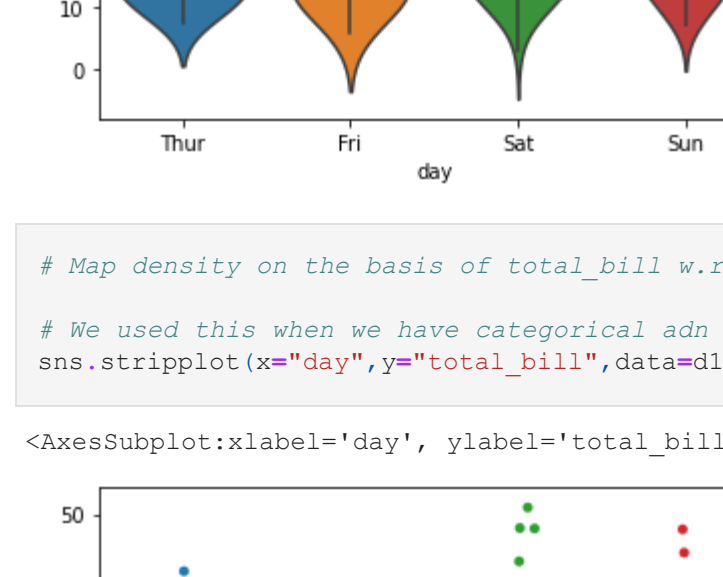
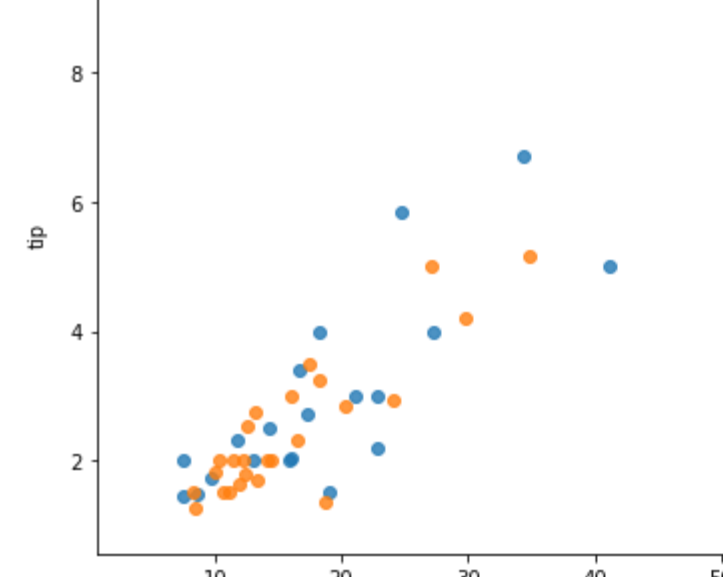
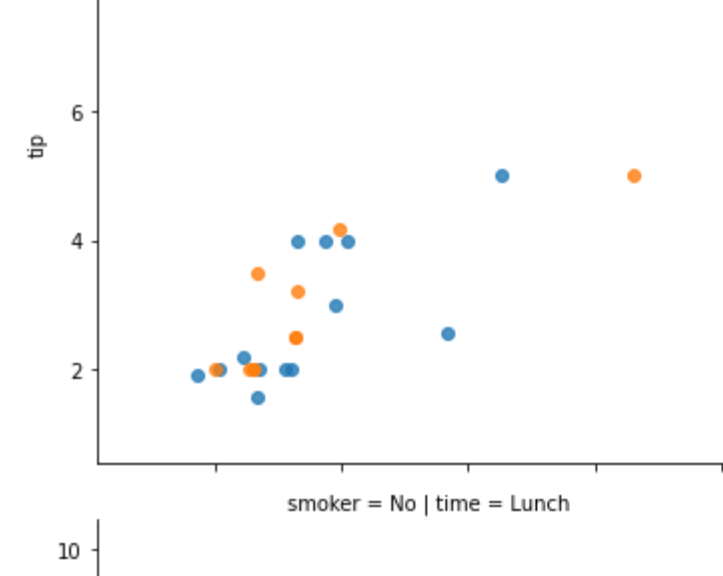
```
In [21]: sns.lmplot(x="total_bill",y="tip",data=d1,fit_reg=False,hue="sex",row="smoker")

Out[21]: <seaborn.axisgrid.FacetGrid at 0x2360021c670>
```

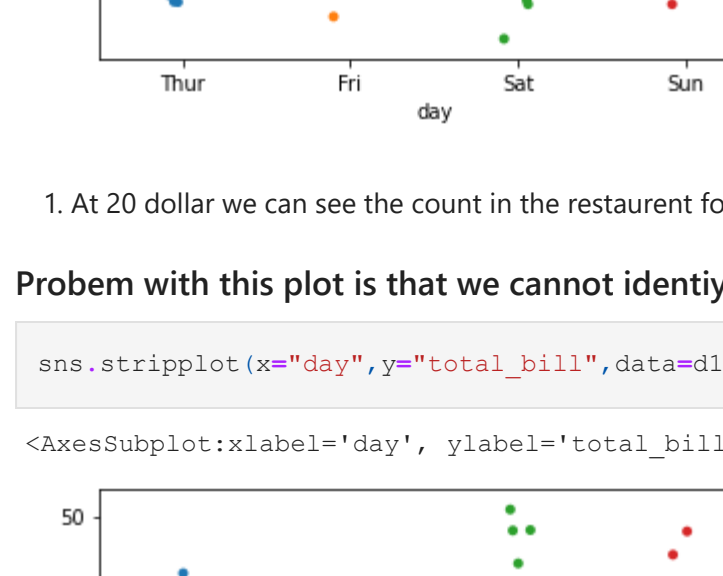


```
In [22]: sns.lmplot(x="total_bill",y="tip",data=d1,fit_reg=False,hue="sex",row="smoker",col="time")

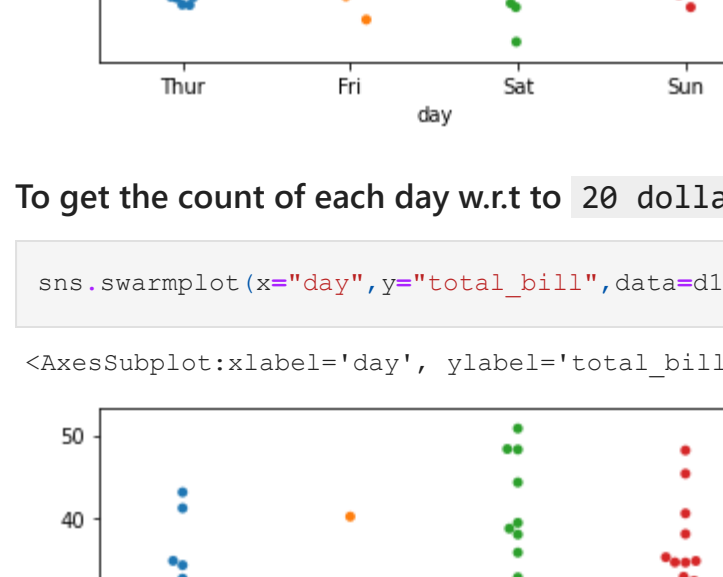
Out[22]: <seaborn.axisgrid.FacetGrid at 0x236002d1940>
```



```
In [23]: # A violin plot plays a similar role as a box and whisker plot.
        # It shows the distribution of quantitative data across several levels of one (or more) categorical variables
        sns.violinplot(x="day",y="total_bill",data=d1)
```



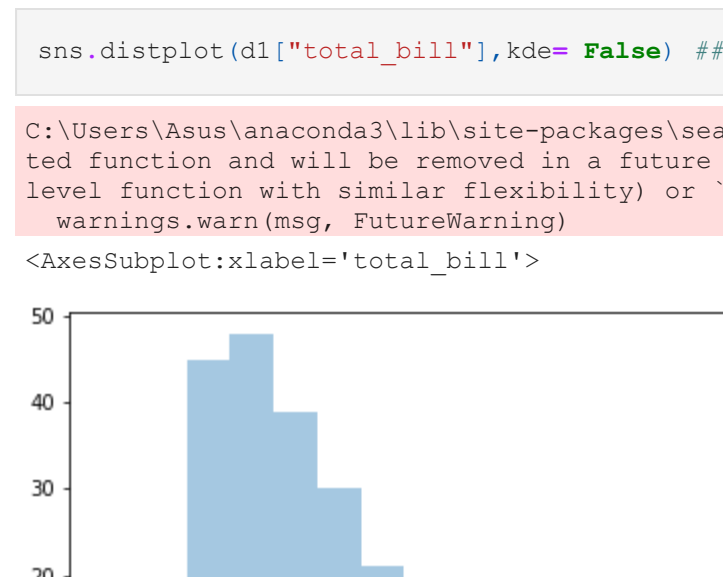
```
In [24]: # Map density on the basis of total_bill w.r.t to each day as in on which denomination each entry are there
        # We used this when we have categorical and discrete data
        sns.stripplot(x="day",y="total_bill",data=d1,jitter=True)
```



1. At 20 dollar we can see the count in the restaurant for each day.

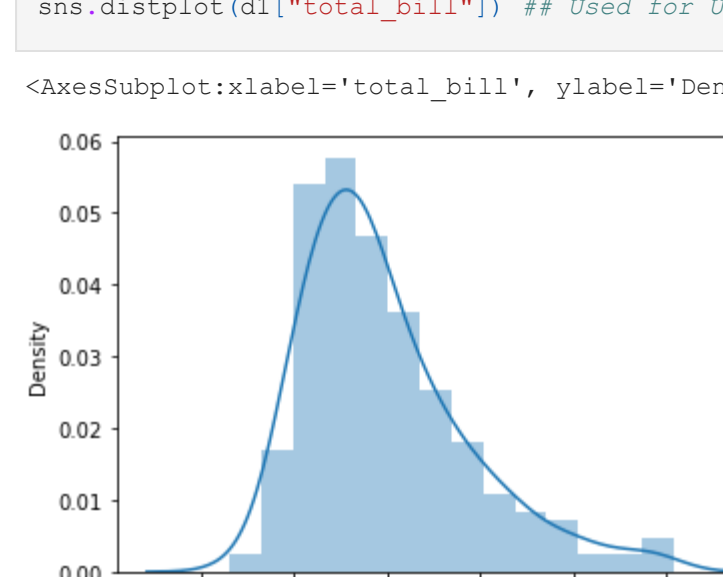
Problem with this plot is that we cannot identify which day has more count of 20 dollar bill

```
In [25]: sns.stripplot(x="day",y="total_bill",data=d1)
```



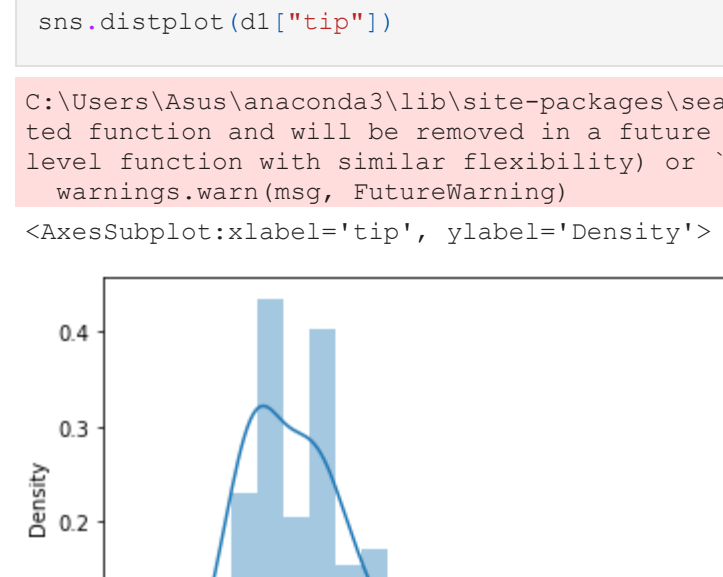
To get the count of each day w.r.t to 20 dollar bill

```
In [26]: sns.swarmplot(x="day",y="total_bill",data=d1)
```



Points are not getting overlapped i.e. on '20 dollar' we can map how many entries are there so that we can see which day can win as 'Saturday' has more width as compare to other days

```
In [27]: sns.distplot(d1["total_bill"],kde=False) ## Used for Un-variate data
```



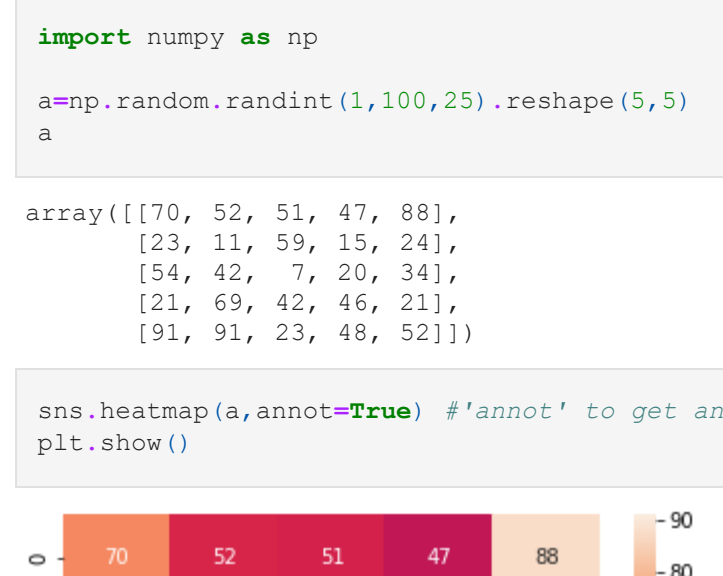
Tells us about the shape of the data or how the data is distributed

Line in the plot is made w.r.t KDE function i.e. Kernel Density Estimator

KDE Area Under the Curve Probability

```
In [29]: sns.distplot(d1["tip"])

Out[29]: <seaborn.axisgrid.FacetGrid at 0x236002d1940>
```



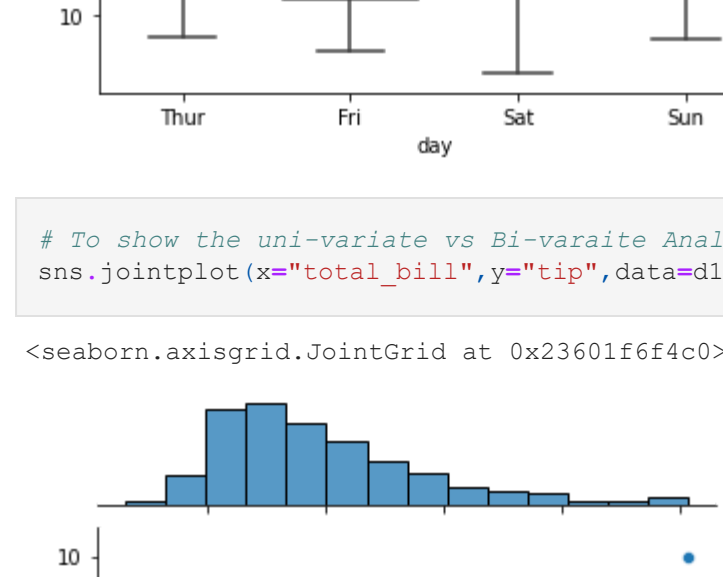
```
In [30]: d1.skew()
```

```
Out[30]:
total_bill    1.139213
tip           1.465493
size          1.447882
dtype: float64
```

Heat Graph

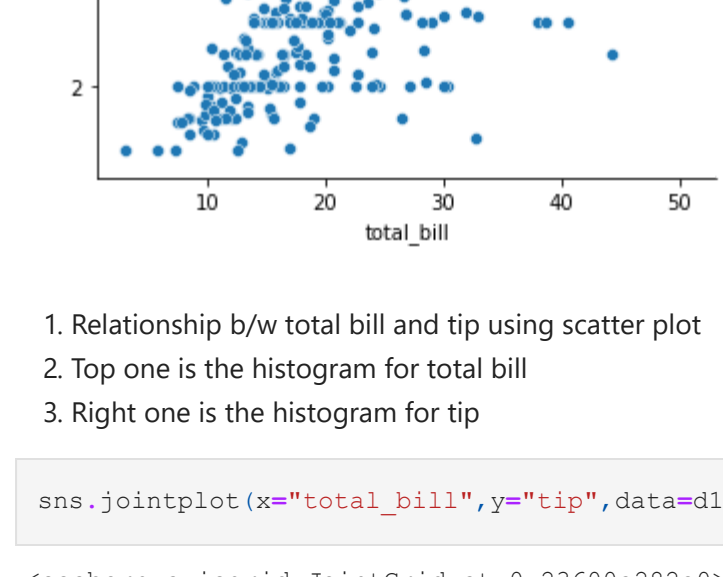
Heat Map is used where used to get the correlation in the data

```
In [31]: import numpy as np
        a=np.random.randint(1,100,25).reshape(5,5)
```

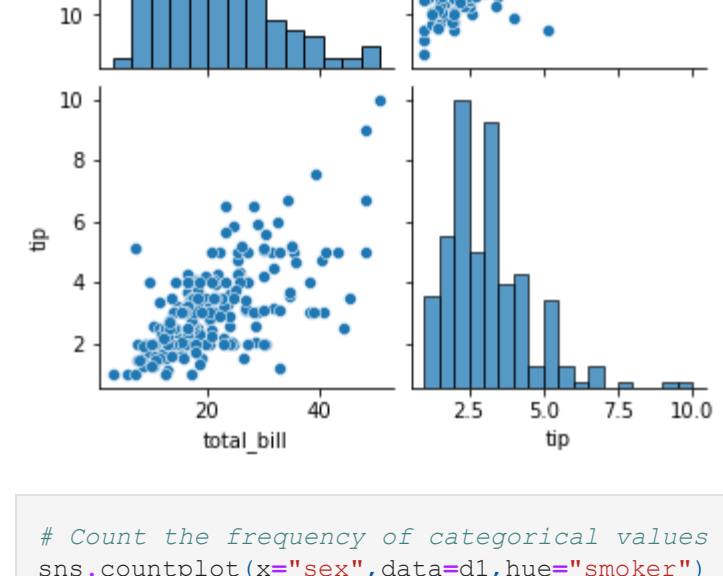
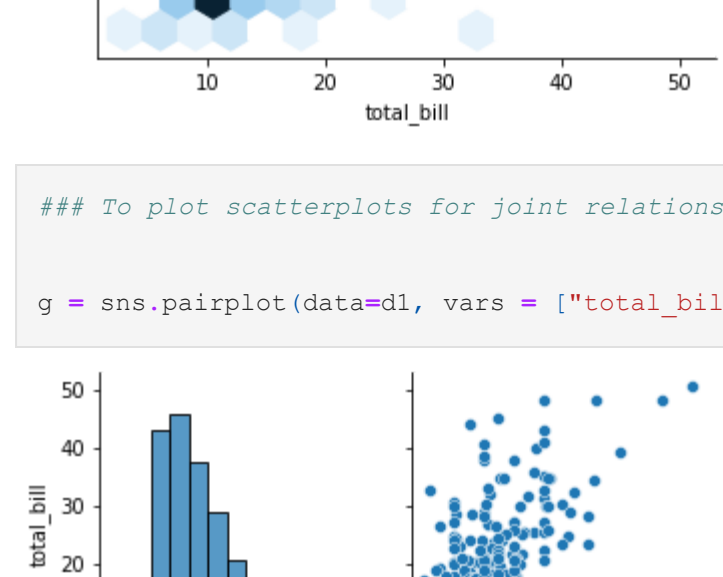
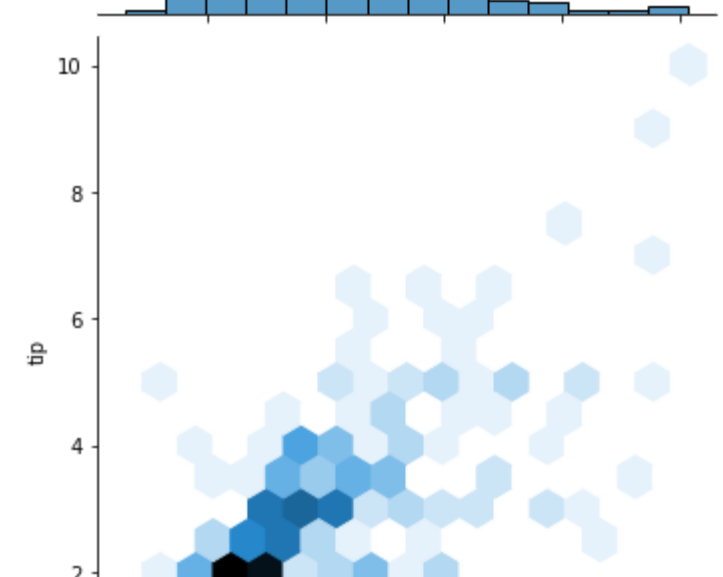


Here color intensity depends upon the value of the cell i.e. if the value is lowest then color is darkest

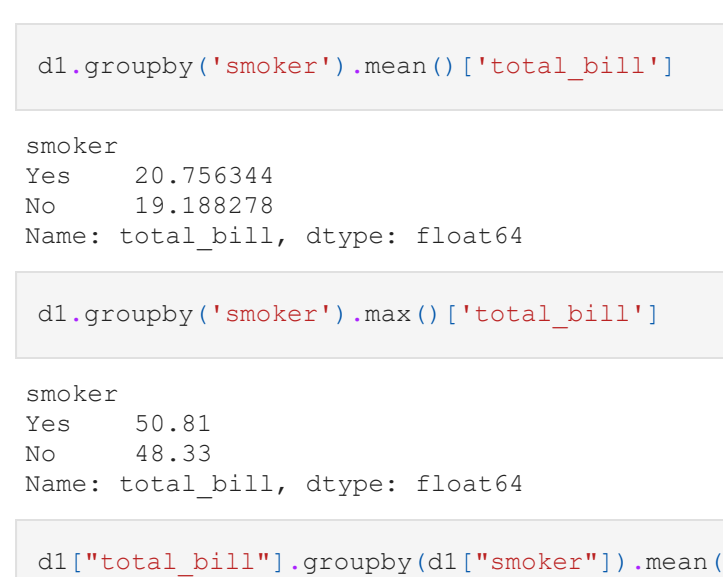
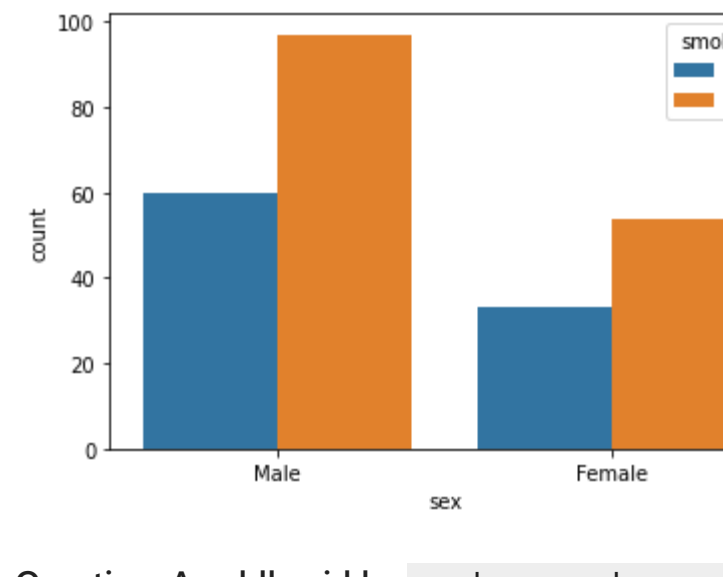
```
In [33]: sns.boxplot(y="total_bill",x="day",data=d1)
```



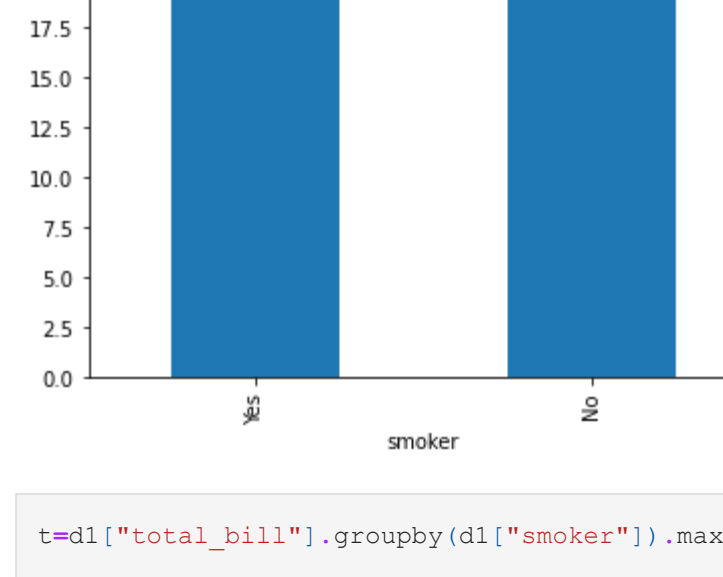
```
In [34]: # To show the uni-variate vs Bi-variate Analysis under one plot
        sns.jointplot(x="total_bill",y="tip",data=d1)
```



```
In [36]: ## To plot scatterplots for joint relationships and histograms for univariate distributions:
        g = sns.pairplot(data=d1, vars = ("total_bill","tip"))
```



```
In [37]: # Count the frequency of categorical values w.r.t to other var
        sns.countplot(x="sex",data=d1,hue="smoker")
        plt.show()
```



Question: Avg bill paid by smokers and non-smokers

```
In [38]: d1.groupby('smoker').mean()['total_bill']
```

```
Out[38]:
smoker
Yes    20.756344
No     19.188278
Name: total_bill, dtype: float64
```

```
In [39]: d1.groupby('smoker').max()['total_bill']
```

```
Out[39]:
smoker
Yes    50.81
No     48.33
Name: total_bill, dtype: float64
```

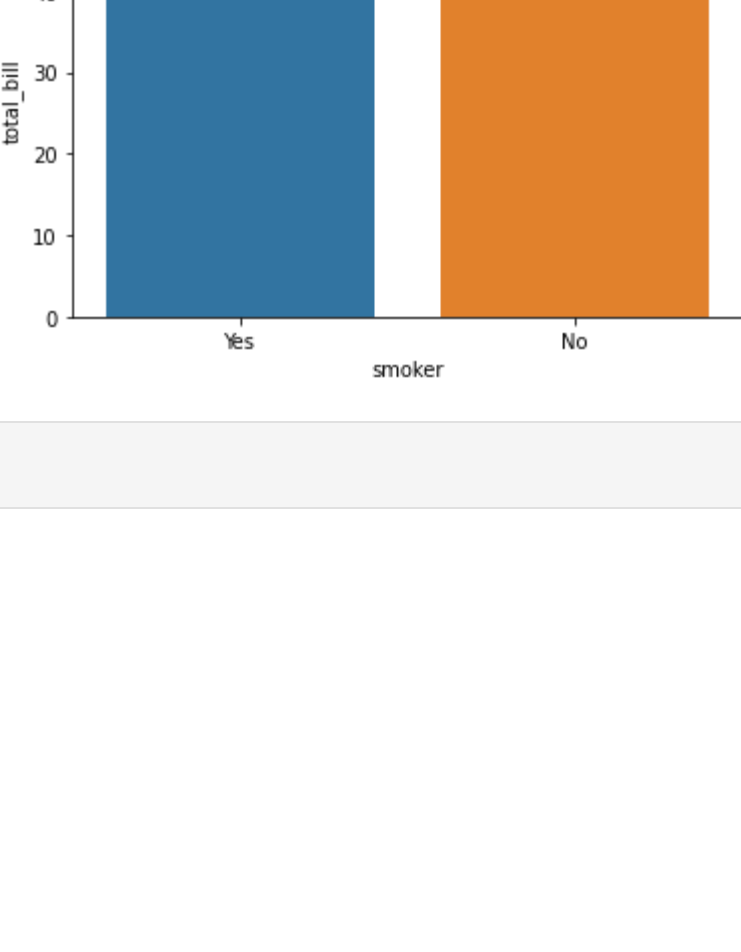
```
In [40]: d1['total_bill'].groupby(d1['smoker']).mean().plot(kind='bar')
        plt.show()
```



```
In [41]: m=d1["total_bill"].groupby(d1["smoker"]).max()
```

```
[42]: sns.barplot(r.index,t)
      plt.show()
```

C:\Users\Asus\anaconda3\lib\site-packages\seaborn\decorators.py:36: FutureWarning: Pass the following variable as keyword args: x, y. From version 0.12, the only valid positional argument will be data , and passing other arguments without an explicit keyword will result in an error or misinterpretation.



In [ ]: