

```
In [7]: import pandas as pd
df = pd.read_csv("E:/tips.csv")

In [8]: df.head()

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]: df.tail()

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]:

Out[10]:
   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
...
...
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

244 rows x 7 columns

In [11]: df.shape

Out[11]:
(244, 7)

In [14]: df.columns

Out[14]:
Index(['total_bill', 'tip', 'sex', 'smoker', 'day', 'time', 'size'], dtype='object')

In [22]: df["sex"].value_counts()

Out[22]:
Male      157
Female     87
Name: sex, dtype: int64

In [24]: df["tip"].value_counts()

Out[24]:
2.00     33
4.00     23
4.99     12
5.00     10
2.50     10
...
4.34      1
1.56      1
5.20      1
2.60      1
1.75      1
Name: tip, Length: 123, dtype: int64

In [28]: df["total_bill"].value_counts()

Out[28]:
13.42      3
13.91      2
15.98      2
17.82      2
18.97      2
...
24.71      1
21.16      1
28.97      1
22.49      1
18.78      1
Name: total_bill, Length: 229, dtype: int64

In [5]: import pandas as pd
df = pd.read_csv("E:/tips.csv")

In [6]: df

Out[6]:
   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
...
...
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

244 rows x 7 columns

In [8]: df["size"].value_counts()

Out[8]:
2     156
3      38
4      37
1       4
6       4
Name: size, dtype: int64

In [9]: df["size"]

Out[9]:
0      2
1      3
2      3
3      2
4      4
...
239     3
240     2
241     2
242     2
243     2
Name: size, Length: 244, dtype: int64

In [16]: import matplotlib.pyplot as plt
plt

Out[16]:
<module 'matplotlib.pyplot' from 'E:\anaconda\lib\site-packages\matplotlib\pyplot.py">

In [51]: import matplotlib.pyplot as plt
z=plt.pie(df["size"].value_counts())

In [58]: import matplotlib.pyplot as plt
z=plt.pie(df["size"].value_counts(),autopct="%0.2f%")

In [48]: import matplotlib.pyplot as plt
z=plt.pie(df["size"].value_counts(),labels=['2','3','4','5','1','6'],autopct="%0.2f%")

In [39]: df["day"].value_counts()

Out[39]:
Sat      87
Sun      76
Thur     62
Fri      19
Name: day, dtype: int64

In [47]: import matplotlib.pyplot as plt
z=plt.pie(df["day"].value_counts(),labels=['sat','3sun','thur','fri'],autopct="%0.2f%")

In [39]: df["smoker"].value_counts()

Out[39]:
No      151
Yes      93
Name: smoker, dtype: int64

In [48]: import matplotlib.pyplot as plt
z=plt.pie(df["smoker"].value_counts(),labels=['no','yes'],autopct="%0.2f%")

In [49]: import matplotlib.pyplot as plt
z=plt.pie(df["time"].value_counts(),labels=['lunch','dinner'],autopct="%0.2f%")

In [43]: df["total_bill"].max()

Out[43]:
50.81

In [45]: df["total_bill"].min()

Out[45]:
3.07

In [53]: df["tip"].value_counts()

Out[53]:
2.00     33
3.00     23
4.00     12
4.99     10
2.50     10
...
4.34      1
1.56      1
5.20      1
2.60      1
1.75      1
Name: tip, Length: 123, dtype: int64

In [54]: df["tip"].max()

Out[54]:
10.0

In [55]: df["tip"].min()

Out[55]:
1.0

In [5]: import pandas as pd
df = pd.read_csv("E:/tips.csv")

In [6]: df

Out[6]:
   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
...
...
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

244 rows x 7 columns

In [7]: df["day"].unique()

Out[7]:
array(['Sun', 'Sat', 'Thur', 'Fri'], dtype=object)

In [7]: df["day"].value_counts()

Out[7]:
Sat      87
Sun      76
Thur     62
Fri      19
Name: day, dtype: int64

In [10]: df["time"].unique()

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

In [12]: df["time"].value_counts()

Out[12]:
Dinner    176
Lunch     68
Name: time, dtype: int64

In [13]: which day more collection
g=df.groupby(df["day"])

In [18]: g.max()

Out[18]:
   total_bill  tip  sex  smoker  time  size
day
Fri    40.17  4.73  Male    Yes  Lunch     4
Sat    50.81 10.00  Male    Yes  Dinner     5
Sun    48.17  6.50  Male    Yes  Dinner     6
Thur   43.11  6.70  Male    Yes  Lunch     6

In [19]: which day less collection
g.min()

Out[19]:
   total_bill  tip  sex  smoker  time  size
day
Fri     5.75  1.00  Female    No  Dinner     1
Sat     3.07  1.00  Female    No  Dinner     1
Sun     7.25  1.01  Female    No  Dinner     2
Thur    7.51  1.25  Female    No  Dinner     1

In [28]: g.max().sum()

Out[28]:
total_bill      182.26
tip              27.93
sex             MaleMaleMale
smoker           YesYesYesYes
time             LunchDinnerDinnerLunch
size             21

In [21]: df["size"].max()

Out[21]:
6

In [22]: df.loc[df["size"]==6]

Out[22]:
   total_bill  tip    sex  smoker  day  time  size
125   29.80  4.2  Female    No  Thur  Lunch     6
141   34.30  6.7    Male    No  Thur  Lunch     6
143   27.05  5.0  Female    No  Thur  Lunch     6
156   48.17  5.0    Male    No  Sun  Dinner     6

In [1]: import pandas as pd
df = pd.read_csv("E:/tips.csv")

In [2]: df

Out[2]:
   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
...
...
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

244 rows x 7 columns

In [3]: df["total_bill"].max()

Out[3]:
50.81

In [5]: df.loc[df["total_bill"]==50.81]

Out[5]:
   total_bill  tip    sex  smoker  day  time  size
170   50.81 10.0  Male    Yes  Sat  Dinner     3

In [7]: df["tip"].max()

Out[7]:
10.0

In [8]: df.loc[df["tip"]==10.0]

Out[8]:
   total_bill  tip    sex  smoker  day  time  size
170   50.81 10.0  Male    Yes  Sat  Dinner     3

In [18]: # only female data
df.loc[df["sex"]=="Female"]

Out[18]:
   total_bill  tip    sex  smoker  day  time  size
4    24.59  3.61  Female    No  Sun  Dinner     4
11   35.26  5.00  Female    No  Sun  Dinner     4
14   14.83  3.02  Female    No  Sun  Dinner     2
16   10.33  1.67  Female    No  Sun  Dinner     3
...
...
226   10.09  2.00  Female    Yes  Fri  Lunch     2
229   22.12  2.88  Female    Yes  Sat  Dinner     2
238   35.83  4.67  Female    No  Sat  Dinner     3
240   27.18  2.00  Female    Yes  Sat  Dinner     2
243   18.78  3.00  Female    No  Thur  Dinner     2

87 rows x 7 columns

In [12]: # only male data
df.loc[df["sex"]=="Male"]

Out[12]:
   total_bill  tip    sex  smoker  day  time  size
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
5    25.29  4.71    Male    No  Sun  Dinner     4
6    8.77  2.00    Male    No  Sun  Dinner     2
...
...
236   12.60  1.00    Male    Yes  Sat  Dinner     2
237   32.83  1.17    Male    Yes  Sat  Dinner     3
239   29.03  5.92    Male    No  Sat  Dinner     3
241   22.67  2.00    Male    Yes  Sat  Dinner     2
242   17.82  1.75    Male    No  Sat  Dinner     2

157 rows x 7 columns

In [21]: # machine learning
from sklearn.preprocessing import LabelEncoder
le1=LabelEncoder()
le2=LabelEncoder()
le3=LabelEncoder()
le4=LabelEncoder()
df["sex"]=le1.fit_transform(df["sex"])
df["smoker"]=le2.fit_transform(df["smoker"])
df["day"]=le3.fit_transform(df["day"])
df["time"]=le4.fit_transform(df["time"])

In [22]: df

Out[22]:
   total_bill  tip  sex  smoker  day  time  size
0    16.99  1.01    0    0    2    0    2
1    10.34  1.66    1    0    2    0    3
2    21.01  3.50    1    0    2    0    3
3    23.68  3.31    1    0    2    0    2
4    24.59  3.61    0    0    2    0    4
...
...
239   29.03  5.92    1    0    1    0    3
240   27.18  2.00    0    1    1    0    2
241   22.67  2.00    1    1    1    0    2
242   17.82  1.75    1    0    1    0    2
243   18.78  3.00    0    0    3    0    2

244 rows x 7 columns

In [31]: # separate input and output
X=df.drop(columns="tip")
Y=df["tip"]

In [33]: # split the data for training and testing
from sklearn.model_selection import train_test_split
X_train,X_test,Y_train,Y_test=train_test_split(X,Y,test_size=0.2)

In [35]: X_train

Out[35]:
   total_bill  sex  smoker  day  time  size
104    20.92    0    0    1    0    2
82     10.07    0    0    3    1    1
239    10.07    1    0    1    0    2
141    34.30    1    0    3    1    6
174    16.82    1    1    2    0    2
...
...
119    24.08    0    0    3    1    4
31     18.35    1    0    1    0    4
208    24.27    1    1    1    0    2
216    28.15    1    1    1    0    5
3     23.68    1    0    2    0    2

195 rows x 6 columns

In [36]: Y_train

Out[36]:
184    4.08
82     1.83
239    1.25
141     6.70
174     4.00
...
119     2.92
31     2.50
208     2.03
216     3.00
3      3.31
Name: tip, Length: 195, dtype: float64

In [42]: # create a ML model
from sklearn.neighbors import KNeighborsRegressor
k=KNeighborsRegressor(n_neighbors=5)
# train the model
K.fit(X_train,Y_train)

Out[42]: KNeighborsRegressor()

In [43]: # test the model
Y_pred=K.predict(X_test)

In [44]: Y_pred

Out[44]:
array([3.46 , 3.07, 3.29 , 2.882, 5.05 , 2.644, 3.616, 4.168, 4.518,
       2.11 , 1.972, 3.732, 2.244, 1.924, 2.638, 2.112, 2.138, 2.918,
       1.872, 3.826, 4.256, 2.55 , 4.142, 4.132, 1.584, 2.802, 3.337,
       1.77 , 3.146, 3.942, 2.448, 2.778, 2.486, 2.540, 2.638, 2.632,
       2.666, 3.234, 3.124, 2.344, 1.64 , 1.792, 2.892, 5.446, 1.584,
       3.364, 3.086, 3.45 , 4.45 ]])

In [46]: Y_test.values

Out[46]:
array([3.5 , 3.07, 3.5 , 3. , 5.17, 2.64, 3.48, 1.17, 4.73, 1.68, 1. ,
       3.75, 2.5 , 1.5 , 5.15, 2. , 1.44, 3.88, 1.47, 3.27, 5. , 1.64,
       5.07, 3.08, 2. , 3. , 3.21, 1.87, 2.71, 6.5 , 3.71, 4.3 , 3.48,
       1.48, 6.5 , 3.5 , 1.02, 3.48, 3.5 , 2. , 1.58, 3.5 , 1.88, 2.5 ,
       1.48, 2.5 , 3.5 , 1.5 , 3. , 3. ])

In [52]: # find mean squared error
from sklearn.metrics import mean_squared_error
mse=mean_squared_error(Y_test,Y_pred)
mse

Out[52]:
1.3855133661224494

In [53]: # r 2 score
from sklearn.metrics import r2_score
rse=r2_score(Y_test,Y_pred)

Out[53]:
0.208385126501247

In [5]:
```