

```
In [31]: import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
```

```
In [32]: #Importing CSV File
#Since we are checking our model will be properly able to detect the room will sold or not we have to use cila
df = pd.read_csv('House-Price.csv', header=0)
```

```
In [33]: df.head()
```

	price	resid_area	air_qual	room_num	age	dist1	dist2	dist3	dist4	teachers	poor_prop	airport	n_hos_beds	n_hot_rooms	waterbod
0	240	32.31	0.538	6.575	65.2	4.35	3.81	4.18	4.01	24.7	4.98	YES	5.480	11.1920	Riv
1	21.6	37.07	0.469	6.421	78.9	4.99	4.70	5.12	5.06	22.2	9.14	NO	7.332	12.1728	Lak
2	34.7	37.07	0.469	7.185	61.1	5.03	4.86	5.01	4.97	22.2	4.03	NO	7.394	10.1200	Nor
3	33.4	32.18	0.458	6.998	45.8	6.21	5.93	6.16	5.96	21.3	2.94	YES	9.268	11.2672	Lak
4	36.2	32.18	0.458	7.147	54.2	6.16	5.86	6.37	5.86	21.3	5.33	NO	8.824	11.2896	Lak

## DATA PRE-PROCESSING AND CLEANING

```
In [34]: # Checking Extended Dictionary (EDD)
#EDD will help to us check the data variation (i.e mean, std, count) and it will also help us to find,
# Missing value and outliers
# we can also find outliers with help of (df.info)--> Gives information about dataframe
df.describe()
```

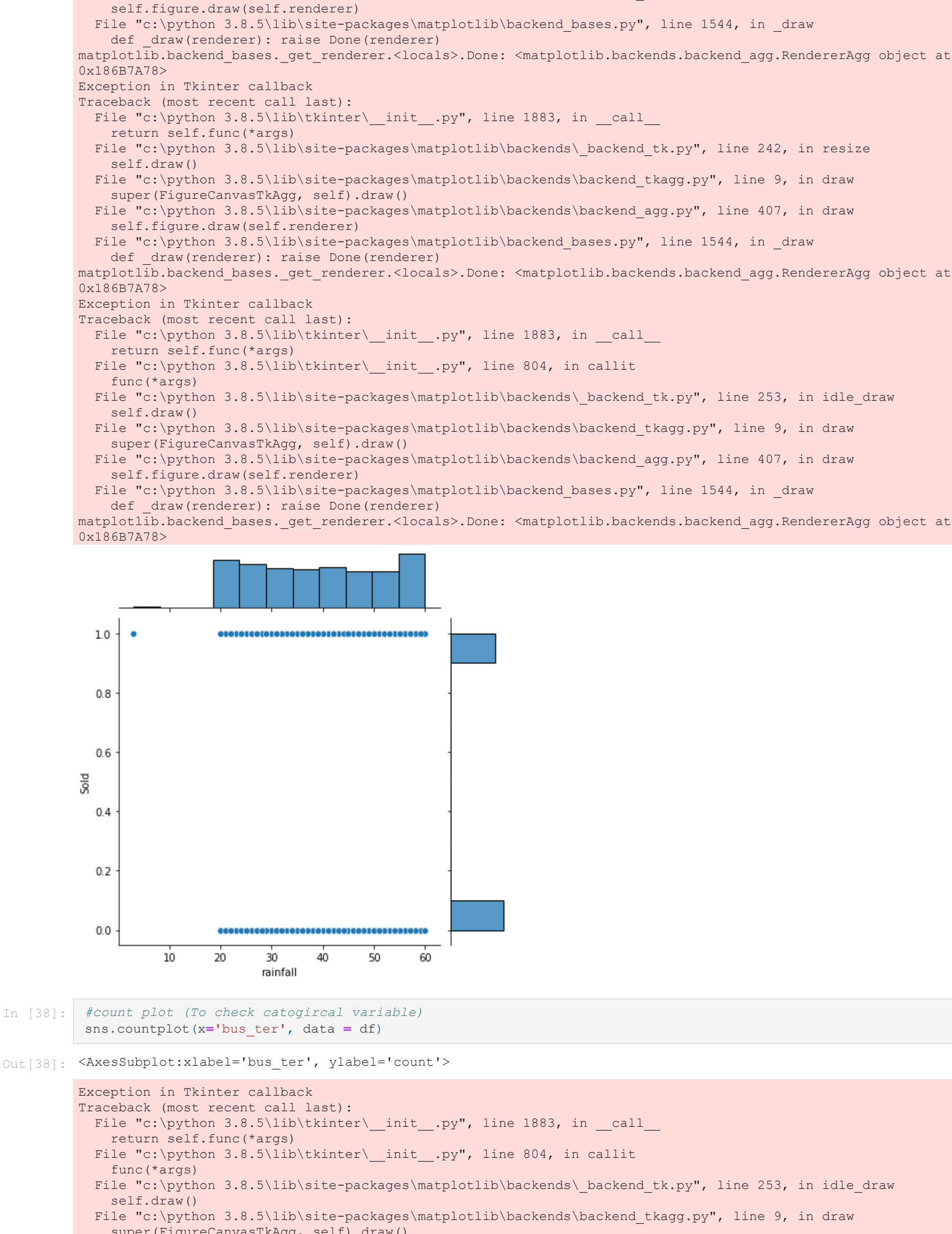
	price	resid_area	air_qual	room_num	age	dist1	dist2	dist3	dist4	teachers	poor_prop
count	506.000000	506.000000	506.000000	506.000000	506.000000	506.000000	506.000000	506.000000	506.000000	506.000000	506.000000
mean	2252854	41.136779	0.554695	6.284634	68.574901	3.971996	3.628775	3.960672	3.618972	21.544466	12.653063
std	9.182176	6.860353	0.115878	0.702617	28.148861	2.108532	2.108580	2.119797	2.099203	2.164946	7.141062
min	5.000000	30.460000	0.138500	3.561000	29.000000	1.130000	0.920000	1.150000	0.730000	18.000000	1.730000
25%	17.025000	35.190000	0.449000	5.885500	45.025000	2.270000	1.940000	2.232500	1.940000	19.800000	6.950000
50%	21.200000	39.690000	0.538000	6.208500	77.500000	3.385000	3.010000	3.375000	3.070000	20.950000	11.360000
75%	25.000000	48.100000	0.624000	6.623500	94.075000	5.367500	4.992500	5.407500	4.985000	22.600000	16.955000
max	50.000000	57.740000	0.871000	8.780000	100.000000	12.320000	11.930000	12.320000	11.940000	27.400000	37.970000

```
In [35]: #Checking Values with help of graph
```

```
In [36]: #scatter plot
sns.jointplot(x='n_hot_rooms', y='Sold', data=df)
```



```
In [37]: #scatter plot to check missing value
sns.jointplot(x='rainfall', y='Sold', data = df)
```



```
In [38]: #count plot (To check catogircal variable)
sns.countplot(x='bus_ter', data = df)
```



```
In [39]: #Observation (can be predicted by graph too)
#1) Missing value in n_hos_beds
#2) Outlier present in n_hot_rooms and rainfall
#3) Bus ha only single catogical variable
```

```
In [40]: # Dealing with missing value----> Giving missing values a value which is mean value of that particular column
df.n_hos_beds = df.n_hos_beds.fillna(df.n_hos_beds.mean())
```

```
In [41]: df.info() ## after dealing with missing value we have same number of observation in every column
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 506 entries, 0 to 505
Data columns (total 19 columns):
 #   Column        Non-Null Count  Dtype  
---  -
 0   price         506 non-null    float64
 1   resid_area    506 non-null    float64
 2   air_qual      506 non-null    float64
 3   room_num      506 non-null    float64
 4   age           506 non-null    float64
 5   dist1         506 non-null    float64
 6   dist2         506 non-null    float64
 7   dist3         506 non-null    float64
 8   dist4         506 non-null    float64
 9   teachers      506 non-null    float64
10   poor_prop     506 non-null    float64
11   airport       506 non-null    object  
12   n_hos_beds    506 non-null    float64
13   n_hot_rooms   506 non-null    float64
14   waterbody     506 non-null    object  
15   rainfall      506 non-null    int64  
16   bus_ter       506 non-null    object  
17   parks         506 non-null    float64
18   Sold          506 non-null    int64  
dtypes: float64(14), int64(2), object(3)
memory usage: 69.2+ KB
```

```
In [42]: # dealing with outliers
# from graph we came to know that in right most part of n_hot_room outlier is presnet
upper_limit = np.percentile(df.n_hot_rooms, [99])[0]
upper_limit
```

Out[42]: 15.399519999999999

```
In [43]: #Now Using Capping and flooring technique to deal with outlier
#checking for upper limit values
df[df.n_hot_rooms > upper_limit]
```

Out[43]:

	price	resid_area	air_qual	room_num	age	dist1	dist2	dist3	dist4	teachers	poor_prop	airport	n_hos_beds	n_hot_rooms	waterb
2	34.7	37.07	0.4690	7.185	61.1	5.03	4.86	5.01	4.97	22.2	4.03	NO	7.394	10.12	N
166	50.0	49.58	0.6050	7.929	96.2	2.11	1.91	2.31	1.86	25.3	3.70	YES	8.300	15.40	F
204	50.0	32.68	0.4161	8.034	31.9	5.41	4.80	5.28	4.99	25.3	2.88	YES	8.300	15.40000	F
267	50.0	33.97	0.5750	8.297	67.0	2.60	2.13	2.43	2.52	27.0	7.44	YES	8.000	15.40	N
369	50.0	48.10	0.6310	6.683	96.8	1.55	1.28	1.65	0.94	19.8	3.73	YES	6.700	15.40	F
423	13.4	48.10	0.6140	6.103	85.1	2.08	1.80	2.34	1.87	19.8	23.29	NO	8.268	81.12	N

```
In [44]: df.n_hot_rooms[(df.n_hot_rooms> 3*upper_limit)] = 3*upper_limit
```

<ipython-input-44-6cc413d964c9>:1: SettingWithCopyWarning:  
A value is trying to be set on a copy of a slice from a DataFrame  
See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user\_guide/indexing.html#returning-a-view-versus-a-copy  
df.n\_hot\_rooms[(df.n\_hot\_rooms> 3\*upper\_limit)] = 3\*upper\_limit

```
In [45]: df[df.n_hot_rooms > upper_limit] #Now values has been reduced
```

Out[45]:

	price	resid_area	air_qual	room_num	age	dist1	dist2	dist3	dist4	teachers	poor_prop	airport	n_hos_beds	n_hot_rooms	waterb
2	34.7	37.07	0.4690	7.185	61.1	5.03	4.86	5.01	4.97	22.2	4.03	NO	7.394	46.19856	N
166	50.0	49.58	0.6050	7.929	96.2	2.11	1.91	2.31	1.86	25.3	3.70	YES	8.300	15.40000	F
204	50.0	32.68	0.4161	8.034	31.9	5.41	4.80	5.28	4.99	25.3	2.88	YES	8.300	15.40000	F
267	50.0	33.97	0.5750	8.297	67.0	2.60	2.13	2.43	2.52	27.0	7.44	YES	8.000	15.40000	N
369	50.0	48.10	0.6310	6.683	96.8	1.55	1.28	1.65	0.94	19.8	3.73	YES	6.700	15.40000	F
423	13.4	48.10	0.6140	6.103	85.1	2.08	1.80	2.34	1.87	19.8	23.29	NO	8.268	46.19856	N

```
In [47]: #similary dealing with rainfall coulmn.. but on lower side
lower_limit = np.percentile(df.rainfall, [1])[0]
lower_limit
```

Out[47]: 20.0

```
In [50]: #Checking values less than 20
df[df.rainfall < lower_limit]
```

Out[50]:

	price	resid_area	air_qual	room_num	age	dist1	dist2	dist3	dist4	teachers	poor_prop	airport	n_hos_beds	n_hot_rooms	waterb
213	28.1	40.59	0.489	6.375	32.3	4.11	3.92	4.18	3.57	21.4	9.38	YES	7.562	10.2248	N

```
In [51]: #changing the value of outlier by flooring and capping method
df.rainfall[(df.rainfall < 0.3*lower_limit)] = 0.3*lower_limit
```

<ipython-input-51-78c2414f55f1>:2: SettingWithCopyWarning:  
A value is trying to be set on a copy of a slice from a DataFrame  
See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user\_guide/indexing.html#returning-a-view-versus-a-copy  
df.rainfall[(df.rainfall < 0.3\*lower\_limit)] = 0.3\*lower\_limit

```
In [52]: df[df.rainfall < lower_limit]
```

Out[52]:

	price	resid_area	air_qual	room_num	age	dist1	dist2	dist3	dist4	teachers	poor_prop	airport	n_hos_beds	n_hot_rooms	waterb
213	28.1	40.59	0.489	6.375	32.3	4.11	3.92	4.18	3.57	21.4	9.38	YES	7.562	10.2248	N

```
In [54]: #Deleting Unwanted columns(i.e columns(data) that will not effect model and which generally conveys similar
df['dist_avg'] = (df.dist1 + df.dist2 + df.dist3 + df.dist4)/4
```

```
In [55]: df.describe()
```

Out[55]:

	price	resid_area	air_qual	room_num	age	teachers	poor_prop	airport	n_hos_beds	n_hot_rooms	waterbody	rainfall	parks	Sold
count	506.000000	506.000000	506.000000	506.000000	506.000000	506.000000	506.000000	506.000000	506.000000	506.000000	506.000000	506.000000	506.000000	506.000000
mean	2252854	41.136779	0.554695	6.284634	68.574901	3.971996	3.628775	3.960672	3.618972	21.544466	12.653063			
std	9.182176	6.860353	0.115878	0.702617	28.148861	2.108532	2.108580	2.119797	2.099203	2.164946	7.141062			
min	5.000000	30.460000	0.138500	3.561000	29.000000	1.130000	0.920000	1.150000	0.730000	18.000000	1.730000			
25%	17.025000	35.190000	0.449000	5.885500	45.025000	2.270000	1.940000	2.232500	1.940000	19.800000	6.950000			
50%	21.200000	39.690000	0.538000	6.208500	77.500000	3.385000	3.010000	3.375000	3.070000	20.950000	11.360000			
75%	25.000000	48.100000	0.624000	6.623500	94.075000	5.367500	4.992500	5.407500	4.985000	22.600000	16.955000			
max	50.000000	57.740000	0.871000	8.780000	100.000000	12.320000	11.930000	12.320000	11.940000	27.400000	37.970000			

```
In [56]: del df['dist1']
del df['dist2']
del df['dist3']
del df['dist4']
del df['bus_ter']
```

```
In [58]: df.head()
```

