

```
In [1]: import pandas as pd

import numpy as np

import matplotlib.pyplot as plt

%matplotlib inline

In [2]: housing = pd.read_csv('C:\Users\91820\Downloads\original.CSV')

In [3]: housing

Out[3]:
```

	price	crime_rate	resid_area	air_qual	room_num	age	dist1	dist2	dist3	dist4	teache
0	24.0	0.00632	32.31	0.538	6.575	65.2	4.35	3.81	4.18	4.01	24
1	21.6	0.02731	37.07	0.469	6.421	78.9	4.99	4.70	5.12	5.06	22
2	34.7	0.02729	37.07	0.469	7.185	61.1	5.03	4.86	5.01	4.97	22
3	33.4	0.03237	32.18	0.458	6.998	45.8	6.21	5.93	6.16	5.96	21
4	36.2	0.06905	32.18	0.458	7.147	54.2	6.16	5.86	6.37	5.86	21
...
501	22.4	0.06263	41.93	0.573	6.593	69.1	2.64	2.45	2.76	2.06	15
502	20.6	0.04527	41.93	0.573	6.120	76.7	2.44	2.11	2.46	2.14	15
503	23.9	0.06076	41.93	0.573	6.976	91.0	2.34	2.06	2.29	1.98	15
504	22.0	0.10959	41.93	0.573	6.794	89.3	2.54	2.31	2.40	2.31	15
505	19.0	0.04741	41.93	0.573	6.030	80.8	2.72	2.24	2.64	2.42	15
506 rows x 19 columns											

```
In [4]: housing.describe()

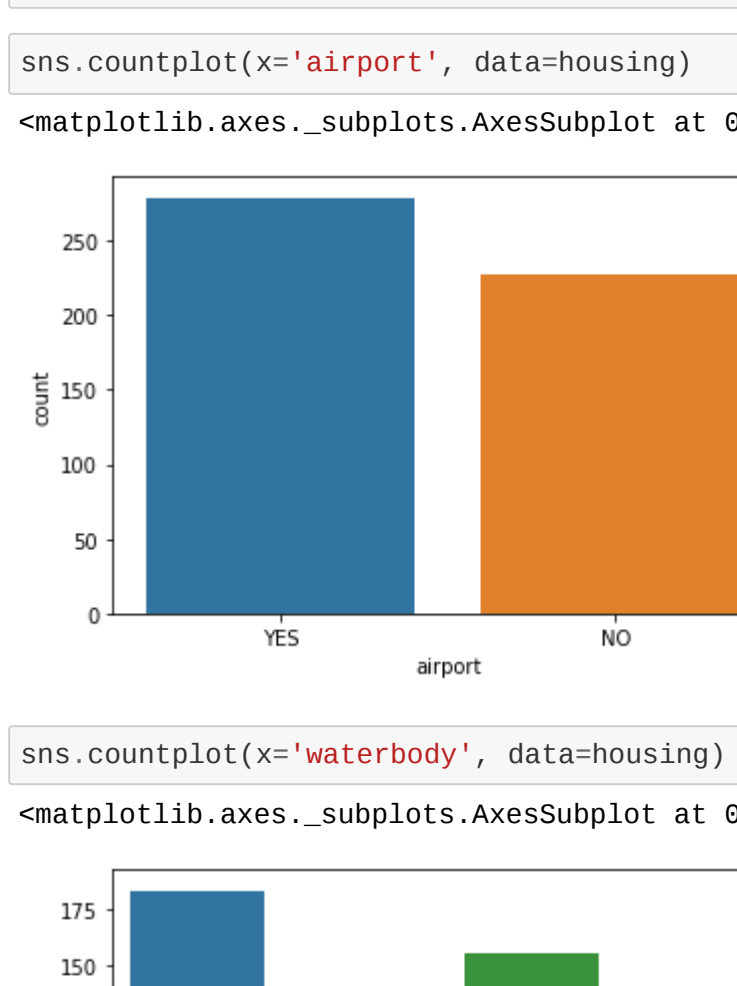
Out[4]:
```

	price	crime_rate	resid_area	air_qual	room_num	age	dist1
count	506.000000	506.000000	506.000000	506.000000	506.000000	506.000000	506.000000
mean	22.528854	3.613524	41.136779	0.554095	6.284634	68.574901	3.971996
std	9.182176	8.601545	6.860353	0.115878	0.702617	28.148861	2.108532
min	5.000000	0.006320	30.460000	0.385000	3.561000	2.900000	1.130000
25%	17.025000	0.082045	35.190000	0.449000	5.885500	45.025000	2.270000
50%	21.200000	0.256510	39.690000	0.538000	6.208500	77.500000	3.385000
75%	25.000000	3.677082	48.100000	0.624000	6.623500	94.075000	5.367500
max	50.000000	88.976200	57.740000	0.871000	8.780000	100.000000	12.320000

```
In [5]: # observation: missing values in "n_hos_beds"

In [6]: plt.scatter(x='n_hot_rooms',y='price',data=housing)

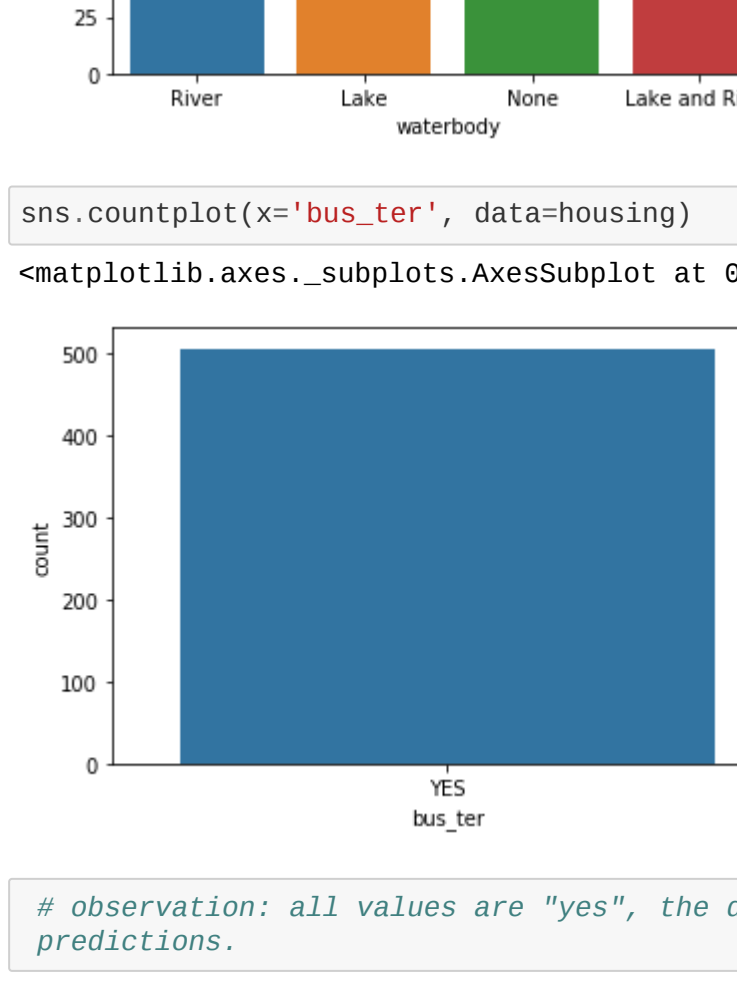
Out[6]: <matplotlib.collections.PathCollection at 0x1535c943048>
```



```
In [7]: # observation: two outliers

In [8]: plt.scatter(x='rainfall',y='price',data=housing)

Out[8]: <matplotlib.collections.PathCollection at 0x1535cd33848>
```



```
In [9]: # observation: single outlier

In [10]: housing.head()

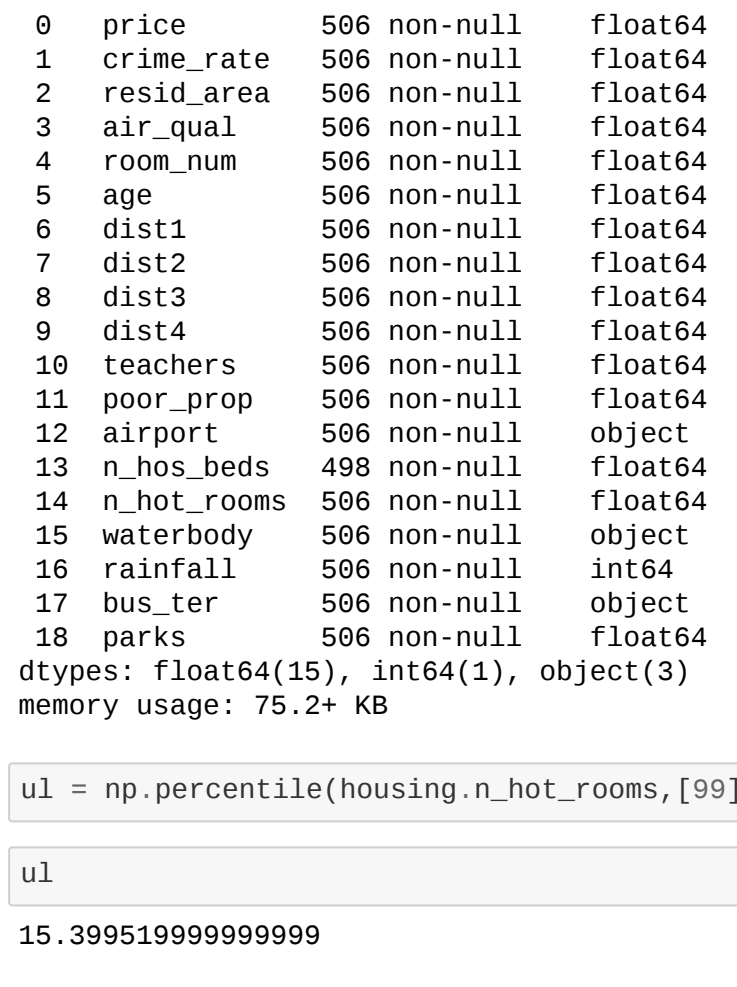
Out[10]:
```

	price	crime_rate	resid_area	air_qual	room_num	age	dist1	dist2	dist3	dist4	teachers
0	24.0	0.00632	32.31	0.538	6.575	65.2	4.35	3.81	4.18	4.01	24.7
1	21.6	0.02731	37.07	0.469	6.421	78.9	4.99	4.70	5.12	5.06	22.2
2	34.7	0.02729	37.07	0.469	7.185	61.1	5.03	4.86	5.01	4.97	22.2
3	33.4	0.03237	32.18	0.458	6.998	45.8	6.21	5.93	6.16	5.96	21.3
4	36.2	0.06905	32.18	0.458	7.147	54.2	6.16	5.86	6.37	5.86	21.3

```
In [11]: import seaborn as sns

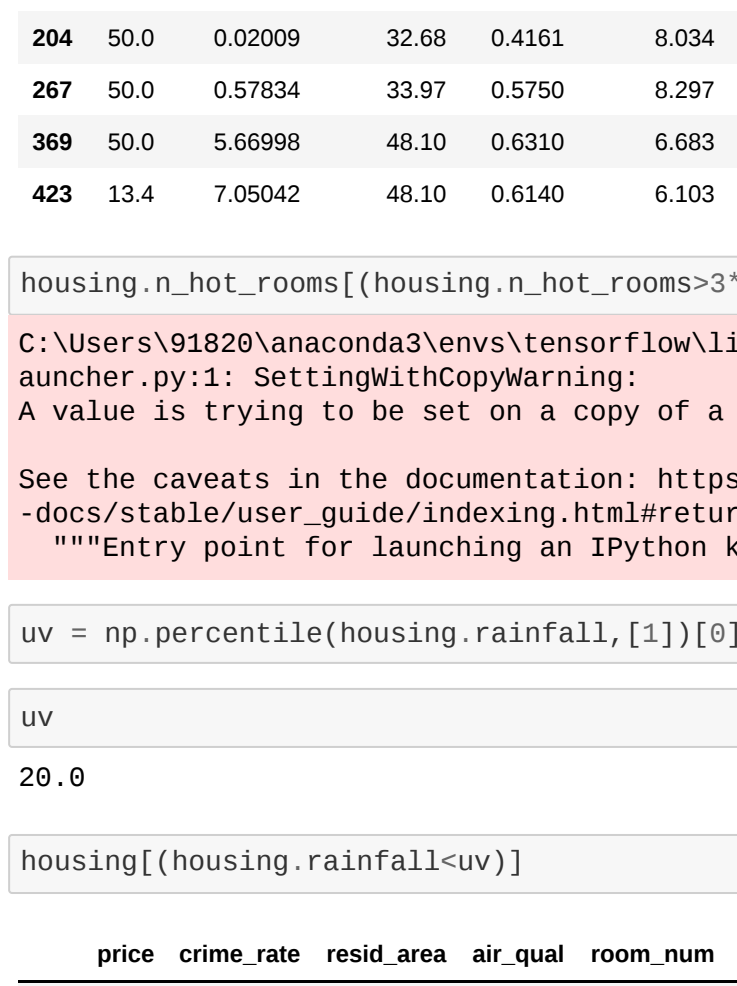
In [12]: sns.countplot(x='airport', data=housing)

Out[12]: <matplotlib.axes._subplots.AxesSubplot at 0x1535ed7d048>
```



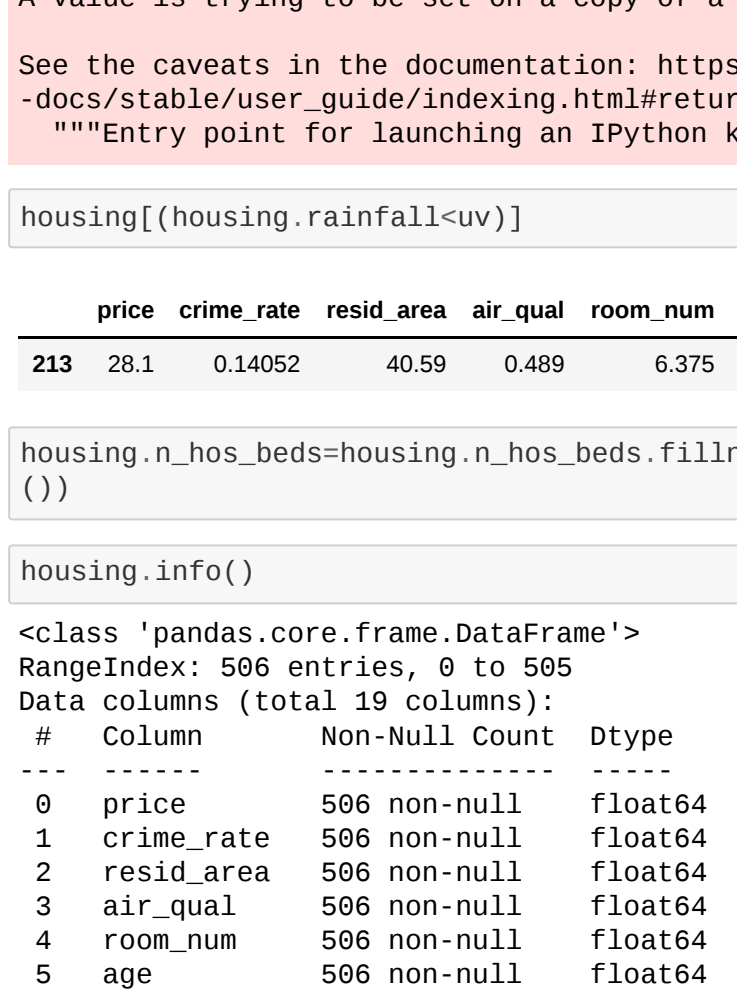
```
In [13]: sns.countplot(x='waterbody', data=housing)

Out[13]: <matplotlib.axes._subplots.AxesSubplot at 0x1535edef488>
```



```
In [14]: sns.countplot(x='bus_ter', data=housing)

Out[14]: <matplotlib.axes._subplots.AxesSubplot at 0x1535ee53148>
```



```
In [15]: # observation: all values are "yes", the data is not going to effect predictions.

In [16]: housing.info()

<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 crime_rate 506 non-null float64
2 resid_area 506 non-null float64
3 air_qual 506 non-null float64
4 room_num 506 non-null float64
5 age 506 non-null float64
6 dist1 506 non-null float64
7 dist2 506 non-null float64
8 dist3 506 non-null float64
9 dist4 506 non-null float64
10 teachers 506 non-null float64
11 poor_prop 506 non-null float64
12 airport 506 non-null object
13 n_hos_beds 498 non-null float64
14 n_hot_rooms 506 non-null float64
15 waterbody 506 non-null object
16 rainfall 506 non-null int64
17 bus_ter 506 non-null object
18 parks 506 non-null float64
dtypes: float64(15), int64(1), object(3)
memory usage: 75.2+ KB

In [17]: ul = np.percentile(housing.n_hot_rooms,[99])[0]

In [18]: ul

Out[18]: 15.399519999999999

In [19]: housing[housing.n_hot_rooms>ul]

Out[19]:
```

	price	crime_rate	resid_area	air_qual	room_num	age	dist1	dist2	dist3	dist4	teache
16	34.7	0.02729	37.07	0.4690	7.185	61.1	5.03	4.86	5.01	4.97	22
262	50.0	2.01019	49.58	0.6050	7.929	96.2	2.11	1.91	2.31	1.86	25
204	50.0	0.02009	32.68	0.4161	8.034	31.9	5.41	4.80	5.28	4.59	25
367	50.0	0.57834	33.97	0.5750	8.297	67.0	2.60	2.13	2.43	2.52	27
269	50.0	5.66998	48.10	0.6310	6.683	96.8	1.55	1.28	1.65	0.94	15
423	13.4	7.05042	48.10	0.6140	6.103	85.1	2.08	1.80	2.34	1.87	15

```
In [20]: housing.n_hot_rooms[(housing.n_hot_rooms>3*ul)] = 3*ul

C:\Users\91820\anaconda3\envs\tensorflow\lib\site-packages\ipykernel_launcher.py: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
"""Entry point for launching an IPython kernel.

In [21]: uv = np.percentile(housing.rainfall,[1])[0]

In [22]: uv

Out[22]: 20.0

In [23]: housing[housing.rainfall<uv]

Out[23]:
```

	price	crime_rate	resid_area	air_qual	room_num	age	dist1	dist2	dist3	dist4	teache
213	28.1	0.14052	40.59	0.489	6.375	32.3	4.11	3.92	4.18	3.57	21

```
In [24]: housing.rainfall[(housing.rainfall<uv)]=0.3*uv

C:\Users\91820\anaconda3\envs\tensorflow\lib\site-packages\ipykernel_launcher.py: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
"""Entry point for launching an IPython kernel.

In [25]: housing[housing.rainfall<uv]

Out[25]:
```

	price	crime_rate	resid_area	air_qual	room_num	age	dist1	dist2	dist3	dist4	teache
213	28.1	0.14052	40.59	0.489	6.375	32.3	4.11	3.92	4.18	3.57	21

```
In [26]: housing.n_hos_beds=housing.n_hos_beds.fillna(housing.n_hos_beds.mean())

In [27]: housing.info()

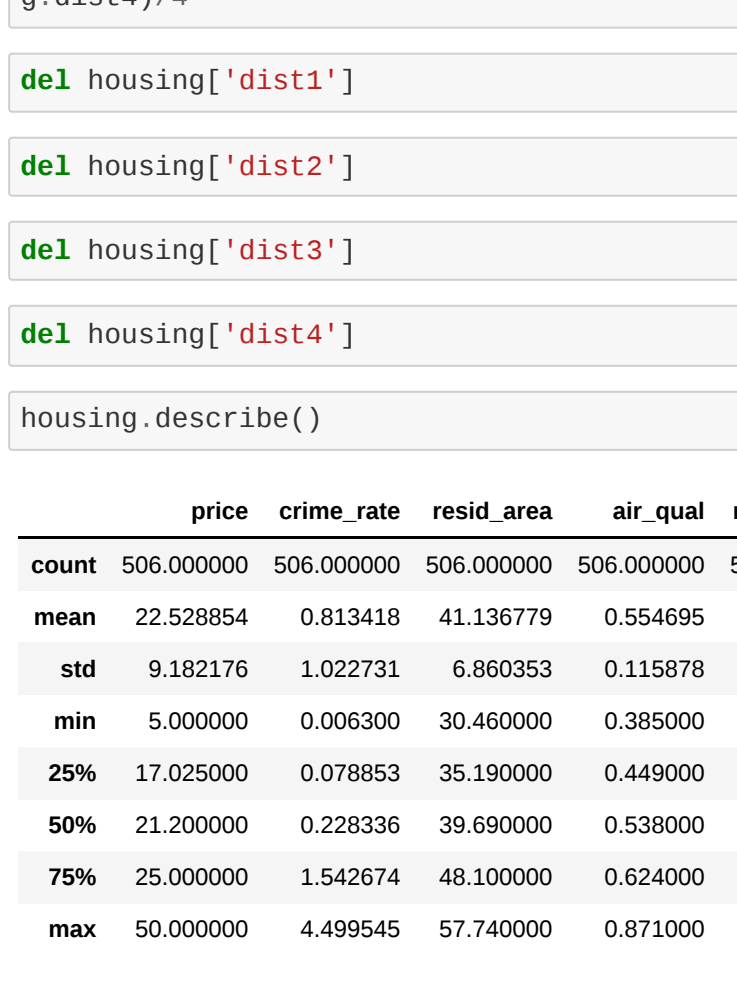
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 506 entries, 0 to 505
Data columns (total 19 columns):
 # Column Non-Null Count Dtype
---
0 price 506 non-null float64
1 crime_rate 506 non-null float64
2 resid_area 506 non-null float64
3 air_qual 506 non-null float64
4 room_num 506 non-null float64
5 age 506 non-null float64
6 dist1 506 non-null float64
7 dist2 506 non-null float64
8 dist3 506 non-null float64
9 dist4 506 non-null float64
10 teachers 506 non-null float64
11 poor_prop 506 non-null float64
12 airport 506 non-null object
13 n_hos_beds 506 non-null float64
14 n_hot_rooms 506 non-null float64
15 waterbody 506 non-null object
16 rainfall 506 non-null int64
17 bus_ter 506 non-null object
18 parks 506 non-null float64
dtypes: float64(15), int64(1), object(3)
memory usage: 75.2+ KB

In [28]: # fixed missing values

In [29]: plt.scatter(x='crime_rate',y='price',data=housing)

# it apparently has outliers but it's not linear so will be treated differently

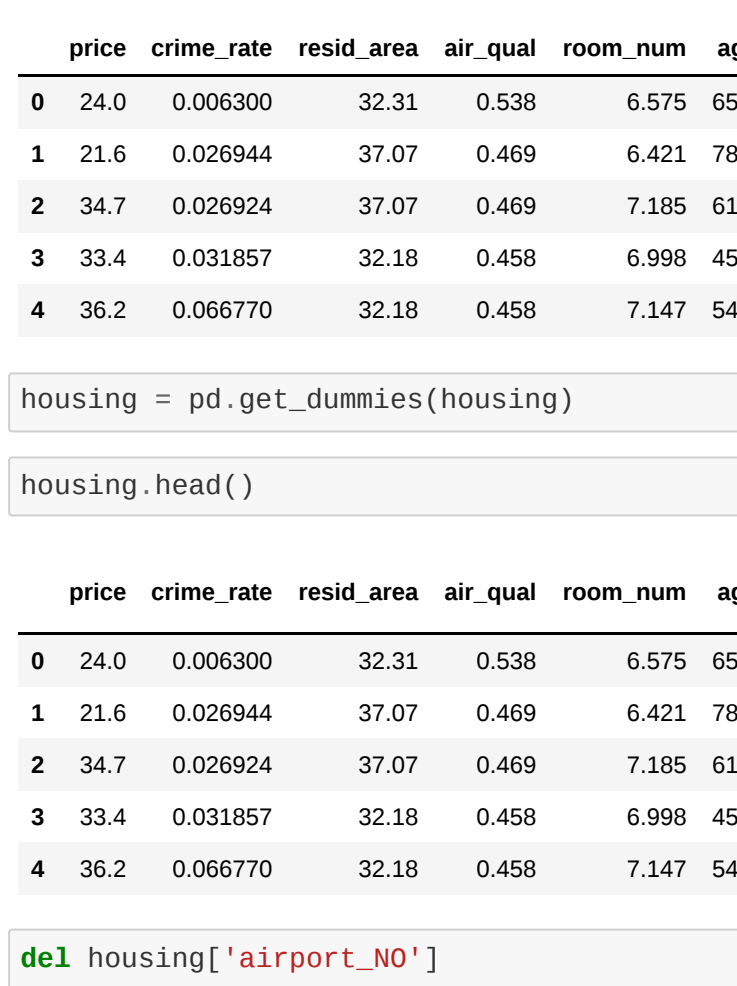
Out[29]: <matplotlib.collections.PathCollection at 0x1535ef6e248>
```



```
In [30]: housing.crime_rate = np.log(1+housing.crime_rate)

In [31]: plt.scatter(x='crime_rate',y='price',data=housing)

Out[31]: <matplotlib.collections.PathCollection at 0x1535ef5bc88>
```



```
In [32]: housing['avg_dist']=(housing.dist1+housing.dist2+housing.dist3+housing.dist4)/4

In [33]: del housing['dist1']

In [34]: del housing['dist2']

In [35]: del housing['dist3']

In [36]: del housing['dist4']

In [37]: housing.describe()

Out[37]:
```

	price	crime_rate	resid_area	air_qual	room_num	age	teachers
count	506.000000	506.000000	506.000000	506.000000	506.000000	506.000000	506.000000
mean	22.528854	0.813418	41.136779	0.554095	6.284634	68.574901	21.544466
std	9.182176	1.022731	6.860353	0.115878	0.702617	28.148861	2.164466
min	5.000000	0.006320	30.460000	0.385000	3.561000	2.900000	18.000000
25%	17.025000	0.078853	35.190000	0.449000	5.885500	45.025000	19.800000
50%	21.200000	0.228336	39.690000	0.538000	6.208500	77.500000	20.950000
75%	25.000000	1.542674	48.100000	0.624000	6.623500	94.075000	22.600000
max	50.000000	4.499545	57.740000	0.871000	8.780000	100.000000	27.400000

```
In [38]: #all locations have bus terminal no effect on prediction so

In [39]: del housing['bus_ter']

In [40]: housing.head()

Out[40]:
```

	price	crime_rate	resid_area	air_qual	room_num	age	teachers	poor_prop	airport	n_hos_beds
0	24.0	0.006300	32.31	0.538	6.575	65.2	24.7	4.98	YES	
1	21.6	0.026944	37.07	0.469	6.421	78.9	22.2	9.14	NO	
2	34.7	0.026924	37.07	0.469	7.185	61.1	22.2	4.03	NO	
3	33.4	0.031857	32.18	0.458	6.998	45.8	21.3	2.94	YES	
4	36.2	0.066770	32.18	0.458	7.147	54.2	21.3	5.33	NO	

```
In [41]: housing = pd.get_dummies(housing)

In [42]: housing.head()

Out[42]:
```

	price	crime_rate	resid_area	air_qual	room_num	age	teachers	poor_prop	n_hos_beds
0	24.0	0.006300	32.31	0.538	6.575	65.2	24.7	4.98	5.480
1	21.6	0.026944	37.07	0.469	6.421	78.9	22.2	9.14	7.332
2	34.7	0.026924	37.07	0.469	7.185	61.1	22.2	4.03	7.394
3	33.4	0.031857	32.18	0.458	6.998	45.8	21.3	2.94	9.268
4	36.2	0.066770	32.18	0.458	7.147	54.2	21.3	5.33	8.824

```
In [43]: del housing['airport_NO']

In [44]: del housing['waterbody_None']

In [45]: # dealing with categorical data

In [46]: housing.corr()

Out[46]:
```

	price	crime_rate	resid_area	air_qual	room_num	age	teachers
price	1.000000	-0.466527	-0.484754	-0.429300	0.696304	-0.377999	0.505655
crime_rate	-0.466527	1.000000	0.660283	0.707587	-0.288784	0.559591	-0.390052
resid_area	-0.484754	0.660283	1.000000	0.763551	-0.391676	0.644779	-0.383248
air_qual	-0.429300	0.707587	0.763551	1.000000	-0.302188	0.731470	-0.188933
room_num	0.696304	-0.288784	-0.391676	-0.302188	1.000000	-0.240265	0.355501
age	-0.377999	0.559591	0.644779	0.731470	-0.240265	1.000000	-0.155155
teachers	0.505655	-0.390052	-0.383248	-0.188933	0.355501	-0.261515	1.000000
poor_prop	-0.740836	0.608970	0.603800	0.590979	-0.613808	0.602339	-0.374044
n_hos_beds	0.108880	-0.040089	0.005799	-0.049553	0.032009	-0.021012	-0.008056
n_hot_rooms	0.017007	0.056870	0.005871	-0.003728	0.014583	0.013918	-0.037007
rainfall	-0.047200	0.082150	0.055845	0.091956	-0.064817	0.074684	-0.045928
parks	-0.391574	0.638951	0.707635	0.915544	-0.282817	0.673850	-0.187004
avg_dist	0.240289	-0.586371	-0.708022	-0.769247	0.205241	-0.747906	0.232452
airport_YES	0.182667	-0.134486	-0.115401	-0.073903	0.163774	0.005101	0.069437
waterbody_Lake	0.036233	-0.025390	-0.026590	-0.046439	-0.004195	0.003452	-0.046717
waterbody_Lake and River	0.037497	0.009076	0.051649	0.013893	0.010554	-0.004354	0.048811
waterbody_River	0.071751	-0.060099	-0.098976	-0.037772	0.046251	-0.088609	0.094256

```
In [47]: del housing['parks']

In [50]: from sklearn.linear_model import LinearRegression

In [51]: y = housing['price']

In [52]: X = housing[['room_num']]

In [53]: lm = LinearRegression()

In [54]: lm.fit(X,y)

Out[54]: LinearRegression(copy_X=True, fit_intercept=True, n_jobs=None, normalize=False)

In [55]: print(lm.intercept_,lm.coef_)

-34.65924312309721 [ 9.09966966]

In [56]: lm.predict(X)

Out[56]: array([25.17108491, 23.76973578, 30.72188341, 29.02824518, 30.3760959
, 23.85163281, 20.04797089, 21.50391804, 16.58094867, 38.9751735
, 23.36935032, 20.02067188, 19.82781512, 19.4746917 , 20.8123431
, 18.42822969, 19.34729633, 19.84777816, 14.98854566, 17.4545650
, 16.0259169 , 19.62028642, 21.23892795, 18.23713663, 19.2471999
, 16.28980732, 18.23713663, 20.36645933, 24.44311334, 26.0719522
, 17.32716966, 20.53935107, 19.48379137, 17.21797363, 20.8123431
, 19.32909699, 18.49192738, 18.57382441, 19.62938609, 25.3530783
, 29.25683659, 26.9455255 , 21.47661903, 21.85880515, 20.5666520
, 17.0450709 , 17.99144555, 20.21176953, 14.46987339, 16.3711683
, 19.60280708, 20.98523687, 20.58787605, 22.34895752, 18.9196103
, 31.30426226, 23.42394834, 27.3641053 , 21.25822696, 19.2744989
, 17.58196041, 19.62938609, 24.08822422, 26.87272314, 29.9848101
, 22.57767906, 18.00054522, 18.82861516, 16.24430897, 18.8923128
, 23.7333371 , 19.58388774, 20.53825338, 22.16819322, 22.4229846
, 22.54128038, 22.47758269, 21.21272861, 22.04989822, 18.7922164
, 26.5542347 , 25.57147038, 22.68687509, 21.45841969, 23.4785463
, 26.67156674, 20.0752699 , 21.03983328, 29.10214221, 29.7573184
, 23.7333371 , 23.62414107, 23.96828855, 21.85880515, 22.2045926
, 25.62606839, 21.42282101, 38.76599139, 36.56017364, 32.8239071
, 26.5542347 , 27.04561686, 23.62414107, 21.1854296 , 21.4584196
, 18.58292408, 18.44642903, 21.0944329 , 24.25201828, 22.0225992
, 21.71321044, 26.44563866, 19.14710359, 20.77594446, 22.2509099
, 19.28359864, 21.54031672, 20.12986722, 18.77401714, 17.4909637
, 16.7558178 , 19.97517353, 19.58388774, 18.62841483, 18.83
```