

REAL STATE PRICE PRDICTER

```
In [55]: import pandas as pd

In [56]: housing = pd.read_csv('data.csv')

In [57]: housing.head()
```

	CRIM	ZN	INDUS	CHAS	NOX	RM	AGE	DIS	RAD	TAX	PTRATIO	B	LSTAT	MEDV
0	0.0632	18.0	2.31	0	0.538	6.575	65.2	4.0900	1	296	15.3	396.90	4.98	24.0
1	0.0273	0.0	7.07	0	0.469	6.421	78.9	4.9671	2	242	17.8	396.90	9.14	21.6
2	0.0278	0.0	7.07	0	0.469	7.185	61.1	4.9671	2	242	17.8	392.83	4.03	34.7
3	0.0323	0.0	2.18	0	0.458	6.998	45.8	6.0622	3	222	18.7	394.63	2.94	33.4
4	0.0695	0.0	2.18	0	0.458	7.147	54.2	6.0622	3	222	18.7	396.90	5.33	36.2

```
In [58]: housing.info()
```

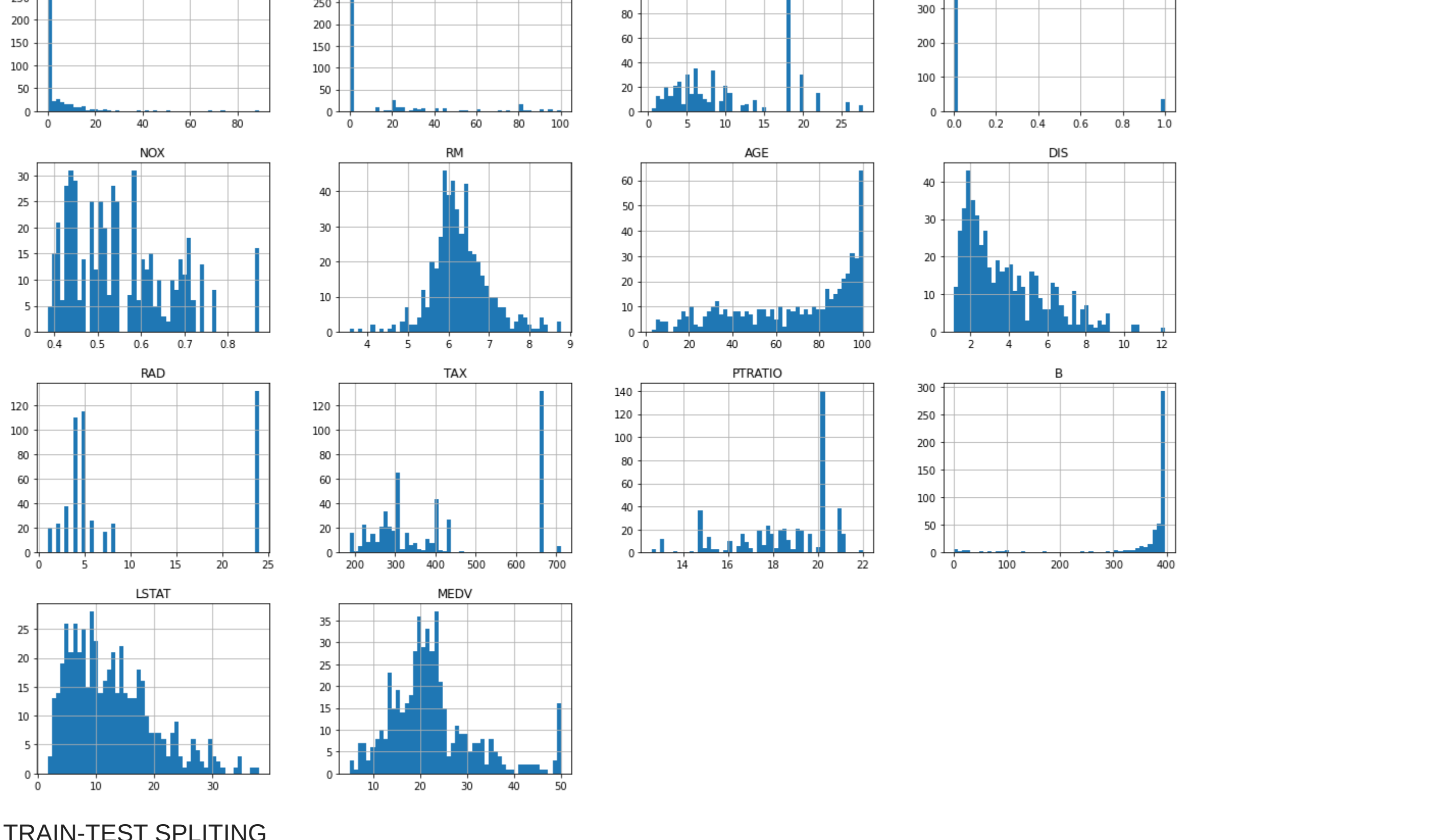
```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 506 entries, 0 to 505
Data columns (total 14 columns):
 #   Column  Non-Null Count  Dtype
---  --
 0   CRIM    506 non-null    float64
 1   ZN      506 non-null    float64
 2   INDUS   506 non-null    float64
 3   CHAS     506 non-null    float64
 4   NOX      506 non-null    float64
 5   RM       506 non-null    float64
 6   AGE      506 non-null    float64
 7   DIS      506 non-null    float64
 8   RAD      506 non-null    float64
 9   TAX      506 non-null    float64
10  PTRATIO  506 non-null    float64
11  B        506 non-null    float64
12  LSTAT    506 non-null    float64
13  MEDV     506 non-null    float64
dtypes: float64(11), int64(3)
memory usage: 51.1 KB
```

```
In [59]: housing.describe()
```

	CRIM	ZN	INDUS	CHAS	NOX	RM	AGE	DIS	RAD	TAX	PTRATIO	B	LSTAT	MEDV
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	3.61524	11.36336	11.13679	0.069170	0.554695	6.287096	68.574901	3.785043	9.548407	408.237154	18.455534	356.674032	12.653063	22.523066
std	8.601545	23.322453	6.860353	0.253984	0.115878	0.705510	28.148861	2.105710	8.707259	168.537116	2.164946	91.294864	7.141062	9.197104
min	0.006320	0.000000	0.460000	0.000000	0.388000	3.561000	2.900000	1.129600	1.000000	187.000000	12.000000	0.320000	1.730000	5.000000
25%	0.082045	0.000000	5.190000	0.000000	0.448000	5.885000	45.025000	1.210071	4.000000	279.000000	17.400000	375.377500	6.950000	17.025000
50%	0.256510	0.000000	9.690000	0.000000	0.538000	6.209000	77.500000	3.207450	5.000000	330.000000	19.050000	391.440000	11.360000	21.200000
75%	3.677093	12.500000	18.100000	0.000000	0.624000	6.629000	94.075000	5.188425	24.000000	666.000000	20.200000	386.225000	16.950000	25.000000
max	88.976200	100.000000	27.740000	1.000000	0.871000	8.780000	100.000000	12.126500	24.000000	711.000000	22.000000	386.900000	37.970000	50.000000

```
In [70]: matplotlib inline

In [71]: #for plotting histogram
import matplotlib.pyplot as plt
housing.hist(bins=5, figsize=(28,15)) #we have to write "plt.show()" in the end to see the graphs if we are coding any IDE (PYCHARM ETC.)
```



TRAIN-TEST SPLITTING

```
In [72]: #for learning purpose

# import numpy as np

# def split_train_test(data, test_ratio):
#     no_random_seed42if dont use the "seed" then random fun. will give diff. test,train data everytime we run it due to which some test data will be seen by algo before testing
#     shuffled = np.random.permutation(len(data))
#     print(shuffled)
#     test_set_size = int(len(data)*test_ratio)
#     test_indices = shuffled[test_set_size:]
#     train_indices = shuffled[:test_set_size]
#     return data.iloc[train_indices], data.iloc[test_indices]

In [73]: train_set, test_set = split_train_test(housing, 0.2)

In [74]: # print(f'rows in train set: {len(train_set)} \nrows in test set: {len(test_set)}\n')

In [75]: from sklearn.model_selection import train_test_split
train_set, test_set = train_test_split(housing, test_size=0.2, random_state=42)
print(f'rows in train set: {len(train_set)} \nrows in test set: {len(test_set)}\n')
```

rows in train set: 404
rows in test set: 102

```
In [76]: #Stratified Shuffling Split = It insure that model should get every type of data so that it can train accordingly(eg.ex. = CHAS ki value 1 bhi hoasakti he ye to tumhe bataya hi nahit ki 0 hoasakti he)
```

```
In [77]: from sklearn.model_selection import StratifiedShuffleSplit
split = StratifiedShuffleSplit(n_splits=1, test_size=0.2, random_state=42)
for train_index, test_index in split.split(housing, housing['CHAS']):
    strat_train_set = housing.loc[train_index]
    strat_test_set = housing.loc[test_index]
```

```
In [78]: strat_train_set["CHAS"].value_counts()
```

CHAS	count
0	376
1	28

Name: CHAS, dtype: int64

```
In [79]: strat_test_set["CHAS"].value_counts()
```

CHAS	count
0	95
1	7

Name: CHAS, dtype: int64

```
In [80]: #95/7 = 376/28
```

```
In [81]: housing = strat_train_set.copy()#ye mene imputer k baad kiya
#while working with big data in ML iss step k baad ik copy lele
```

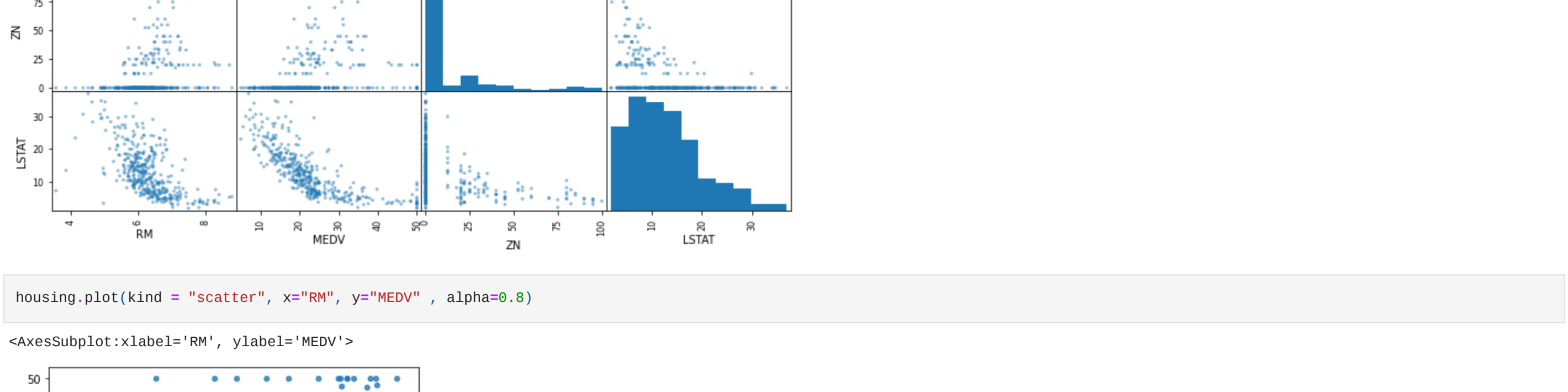
looking for correlation

```
In [82]: corr_matrix = housing.corr()
corr_matrix["MEDV"].sort_values(ascending=False)
```

```
Out[82]: MEDV    1.000000
RM        0.679323
B         0.361761
ZN        0.397471
DIS       0.248451
CHAS      0.255866
AGE       0.364996
RAD       0.374693
CRIM      0.385715
NOX       0.422873
TAX       0.450955
INDUS     0.473510
PTRATIO   0.483534
LSTAT     0.748494
Name: MEDV, dtype: float64
```

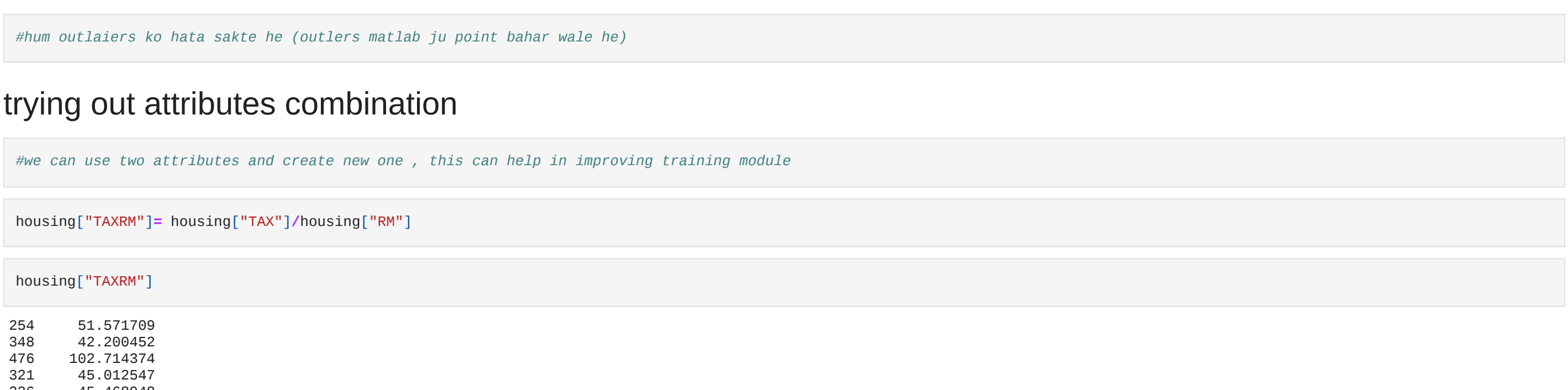
```
In [83]: #pearson correlation coef = dependance of one thing on dec. , inc. of other thing
# (kisi chi ko agar inc. ya dec. kargo to dosri chi inc hogi ya dec.)
# strong + correlation =1 ( ye inc. to wo bhi inc. same rate se)
# strong - correlation =-1 (ye inc. to wo dec. same rate se)
```

```
In [84]: from pandas.plotting import scatter_matrix
attributes = ["RM", "MEDV", "ZN", "LSTAT"]
scatter_matrix(housing[attributes], figsize=(12,8))
```



```
In [85]: housing.plot(kind = "scatter", x="RM", y="MEDV", alpha=0.8)
```

```
Out[85]: <AxesSubplot: xlabel='RM', ylabel='MEDV'>
```



```
In [86]: #hm outliers ko kaha sakte he (outliers matlab ju point bahar wale he)
```

trying out attributes combination

```
In [87]: #we can use two attributes and create new one , this can help in improving training modelle
```

```
In [88]: housing["TAXRM"] = housing["TAX"]/housing["RM"]
```

```
In [89]: housing["TAXRM"]
```

	TAXRM
254	51.571799
348	42.204452
178	182.714574
321	45.612547
326	45.468948
155	65.597152
423	189.126559
98	35.294418
455	182.068956
216	46.970889
Name: TAXRM, Length: 404, dtype: float64	

```
In [90]: housing.head()
```

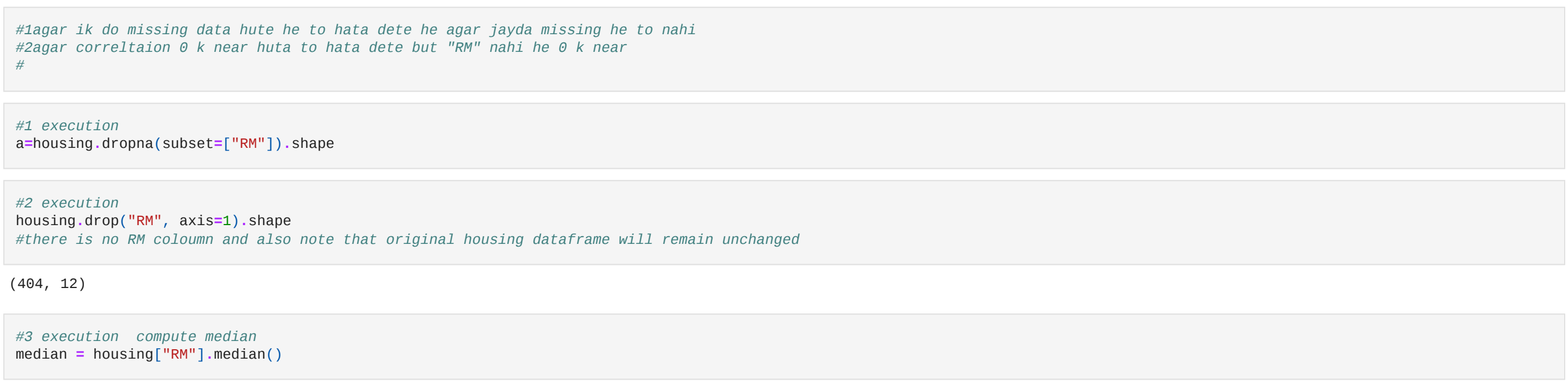
	CRIM	ZN	INDUS	CHAS	NOX	RM	AGE	DIS	RAD	TAX	PTRATIO	B	LSTAT	MEDV	TAXRM
254	0.04819	80.0	3.64	0	0.392	6.108	32.0	9.2203	1	315	16.4	392.89	6.57	21.9	51.571799
348	0.01501	80.0	2.01	0	0.435	6.635	29.7	8.3440	4	280	17.0	390.94	5.99	24.5	42.204452
476	4.87141	0.0	18.10	0	0.614	6.484	93.6	2.3053	24	666	20.2	396.21	18.68	16.7	102.714574
321	0.18159	0.0	7.38	0	0.493	6.376	54.3	4.5404	5	287	19.6	396.90	6.87	23.1	45.612547
326	0.39347	0.0	7.38	0	0.493	6.312	28.9	5.4159	5	287	19.6	396.90	6.15	23.0	45.468948

```
In [91]: corr_matrix = housing.corr()
corr_matrix["MEDV"].sort_values(ascending=False)
```

```
Out[91]: MEDV    1.000000
RM        0.679323
B         0.361761
ZN        0.397471
DIS       0.248451
CHAS      0.255866
AGE       0.364996
RAD       0.374693
CRIM      0.385715
NOX       0.422873
TAX       0.450955
INDUS     0.473510
PTRATIO   0.483534
TAXRM     0.526884
LSTAT     0.748494
Name: MEDV, dtype: float64
```

```
In [92]: housing.plot(kind = "scatter", x="TAXRM", y="MEDV", alpha=0.8)
```

```
Out[92]: <AxesSubplot: xlabel='TAXRM', ylabel='MEDV'>
```



```
In [93]: #ye mene selecting ki k line paar bataya tha
housing = strat_train_set.drop("MEDV", axis=1)
housing_labels = strat_train_set["MEDV"].copy()
```

missing attributes

```
In [94]: #to take care of missing attributes , you have three option:
# 1.get rid of the missing data point
# 2.get rid of the whole attributes
# 3.set the value to some value(0, mean or median)
```

```
In [95]: #agar ik do missing data hote to hata dete he agar jayda missing he to nahi
#agar correlation 0 k near huta to hata dete but "RM" nahi he 0 k near
```

```
In [96]: #1 execution
a=housing.dropna(subset=["RM"]).shape
```

```
In [97]: #2 execution
housing.dropna("RM", axis=1).shape
#there is no RM column and also note that original housing dataframe will remain unchanged
```

```
Out[97]: (404, 12)
```

```
In [98]: #3 execution
median = housing["RM"].median()
median = housing["RM"].median()
```

```
In [99]: median
```

6.2135

```
In [100]: #3
housing["RM"].fillna(median)
#note that original housing dataframe will remain unchanged
```

```
Out[100]: 254    6.198
348    6.635
476    6.484
321    6.376
326    6.312
155    6.152
423    6.183
98     7.829
455    6.525
216    6.885
Name: RM, Length: 404, dtype: float64
```

agar kahi bhi missing value hogi data set me ya fir input me waha ye median

bhar dena (niche)

```
In [101]: housing.describe()#before filling missing attributes
```

	CRIM	ZN	INDUS	CHAS	NOX	RM	AGE	DIS	RAD	TAX	PTRATIO	B	LSTAT	MEDV
count	404.000000	404.000000	404.000000	404.000000	404.000000	404.000000	404.000000	404.000000	404.000000	404.000000	404.000000	404.000000	404.000000	404.000000
mean	3.602814	10.83634	11.34950	0.069397	0.558064	6.285956	69.039851	3.746210	9.735149	412.341584	18.473297	353.392822	12.791609	
std	8.096383	22.150636	6.877817	0.254290	0.116875	0.712516	28.258248	2.090957	8.731259	168.672623	2.129243	96.060235	7.235740	
min	0.006320	0.000000	0.740000	0.000000	0.388000	3.561000	2.900000	1.129600	1.000000	187.000000	13.000000	0.320000	1.730000	
25%	0.086862	0.000000	5.190000	0.000000	0.453000	5.879750	44.850000	1.203975	4.000000	284.000000	17.400000	374.617500	6.847500	
50%	0.286735	0.000000	9.900000	0.000000	0.538000	6.213500	78.200000	3.122200	5.000000	337.000000	19.000000	390.950000	11.570000	
75%	3.731823	12.500000	18.100000	0.000000	0.631000	6.630250	94.100000	5.100400	24.000000	666.000000	20.200000	386.630000	17.102500	
max	73.534100	100.000000	27.740000	1.000000	0.871000	8.780000	100.000000	12.126500	24.000000	711.000000	22.000000	386.900000	36.980000	

```
In [102]: from sklearn.impute import SimpleImputer
imputer = SimpleImputer(strategy="median")
imputer.fit(housing)
```

```
Out[102]: SimpleImputer(strategy='median')
```

```
In [103]: imputer.statistics_
```

array([2.86735e+01, 0.889098e+00, 9.986098e+00, 0.886098e+00, 0.538098e+01, 6.213500e+00, 7.820000e+01, 3.122200e+00, 5.000000e+00, 3.370000e+02, 1.900000e+01, 3.906098e+02, 1.157000e+01])

```
In [104]: x=imputer.transform(housing)
```

```
In [105]: housing_tr = x.DataFrame(x.columns=housing.columns)
```

```
In [106]: housing_tr.describe()
```

	CRIM	ZN	INDUS	CHAS	NOX	RM	AGE	DIS	RAD	TAX	PTRATIO	B	LSTAT
count	404.000000	404.000000	404.000000	404.000000	404.000000	404.000000	404.000000	404.000000	404.000000	404.000000	404.000000	404.000000	404.000000
mean	3.602814	10.83634	11.34950	0.069397	0.558064	6.285956	69.039851	3.746210	9.735149	412.341584	18.473297	353.392822	12.791609
std	8.096383	22.150636	6.877817	0.254290	0.116875	0.712516	28.258248	2.090957	8.731259	168.672623	2.129243	96.060235	7.235740
min	0.006320	0.000000	0.740000	0.000000	0.388000	3.561000	2.900000	1.129600	1.000000	187.000000	13.000000	0.320000	1.730000
25%	0.086862	0.000000	5.190000	0.000000	0.453000	5.879750	44.850000	1.203975	4.000000	284.000000	17.400000	374.617500	6.847500
50%	0.286735	0.000000	9.900000	0.000000	0.538000	6.213500	78.200000	3.122200	5.000000	337.000000	19.000000	390.950000	11.570000
75%	3.731823	12.500000	18.100000	0.000000	0.631000	6.630250	94.100000	5.100400	24.000000	666.000000	20.200000	386.630000	17.102500
max	73.534100	100.000000	27.740000	1.000000	0.871000	8.780000	100.000000	12.126500	24.000000	711.000000	22.000000	386.900000	36.980000

scikit-learn design

primarily, three types of objects

1. estimators - It estimates some parameter based on a dataset. Eg. imputer . It has fit method and transform method. Fit method - fits the datasets and calculates internal parameters. It has some hyperparameters like strategy
2. transformers - transform method takes input an returns output based on the learning from fit(). It has a convenience function called fit_transform() which fits and then transform
3. predictors - Linearregression model is an example of predictor. fit() and predict() are two common functions. It also gives score function which will evaluate the predictions

Feature Scaling

primarily, two types of feature scaling method:

1. Min-Max scaling (Normalisation) value=(min)/(max-min) sklearn provide a class called MinMaxScaler for this
2. Standardization (Value - Mean)/std Sklearn provides a class called StandardScaler for this

creating pipeline

```
In [107]: from sklearn.pipeline import Pipeline
from sklearn.preprocessing import StandardScaler
my_pipeline = Pipeline([
    ("imputer", SimpleImputer(strategy="median")),
    #.....add as many as you want in your pipeline
    ("std_scaler", StandardScaler())
])
```

```
In [108]: housing_num_tr = my_pipeline.fit_transform(housing)
```

```
In [109]: housing_num_tr
```

(404, 13)

```
Out[109]: array([[ 0.43942896,  3.12628155, -1.1216591
```