

## Project housing real estate code

1)import pandas as pd

Now the Pandas package can be referred to as pd instead of pandas

Pandas DataFrame is a 2 dimensional data structure, like a 2 dimensional array, or a table with rows and columns

df = pd.DataFrame(data)

2) housing =pd.read\_csv("data.csv")

here housing is a dataframe

3)housing.head()

now dataframe named housing has data of data.csv and .head will show a quick glance at first 5 rows of this data

4)housing.info()

shows all the information about data like how many entries are there this data set we are using has 506 entries

5) housing['CHAS'].value\_counts()

tells in 471 rows data is 0

and in 35 rows data is 1

note chas is in capital

6)housing.describe()

told us min mean max std counts 25% 50% means 25% values are lower than given data in table and all in data.csv

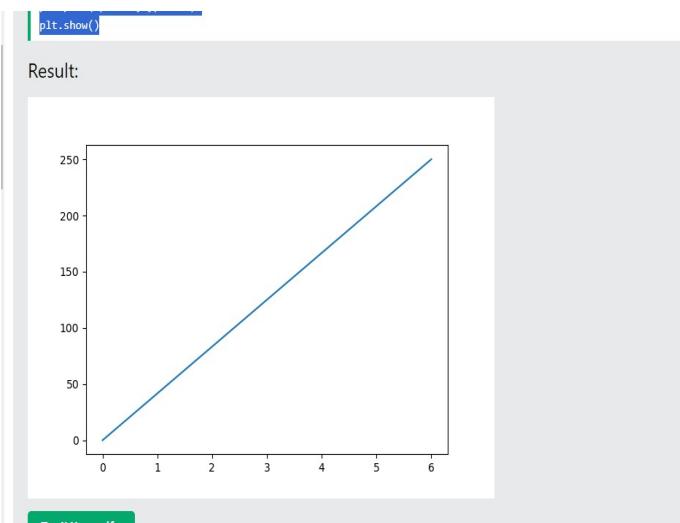
7)import matplotlib.pyplot as plt

imports pyplot

```
import matplotlib.pyplot as plt
import numpy as np

xpoints = np.array([0, 6]) //shows x axis 0 to 6
ypoints = np.array([0, 250]) //shows y axis 0 to 250

plt.plot(xpoints, ypoints)
plt.show()
```



```
6)%matplotlib inline
```

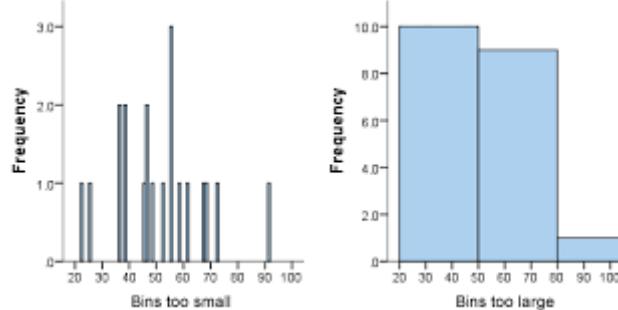
This command is used to enable the inline rendering of matplotlib plots directly in the notebook interface.

```
9)housing.hist(bins=50,figsize(20,15))
```

get plots of every term like age crm b lstat and all

`housing.hist`: This part of the code suggests that you have a variable or DataFrame named `housing`

`bins=50`: This argument specifies the number of bins (or bars) that the histogram will be divided into. Bins represent intervals or ranges of data values. In this case, you're specifying 50 bins, meaning the data will be divided into 50 intervals for the histogram.



`figsize=(20, 15)`: This argument sets the size of the figure (the entire plot) in which the histogram will be displayed. It's a tuple representing the width and height of the figure in inches. In this case, the figure will be 20 inches wide and 15 inches tall.

Now we will do train test splitting

evaluate the performance of a model on unseen data. The process involves dividing a dataset into two distinct subsets: a training set and a testing (or validation) set. The main purpose of this splitting is to assess how well a trained model generalizes to new, unseen data

```
1)import numpy as np
```

```
2 def split_train_test(data,test_ratio):
```

created a function to split test and train the data test ratio will be like 20% to 50%

because we want to test data so we wont give majority of data so its value will be like 0.2 or 0.3 something

2.5) `np.random.seed(42)`

explained neeche ignore for now

```
3)shuffled =np.random.permutation(len(data))
```

`np.random.permutation(len(6))` lets say it was this output will be anything from 0 to 5 random and permutation are functions of numpy library and shuffled is a variable we defined here

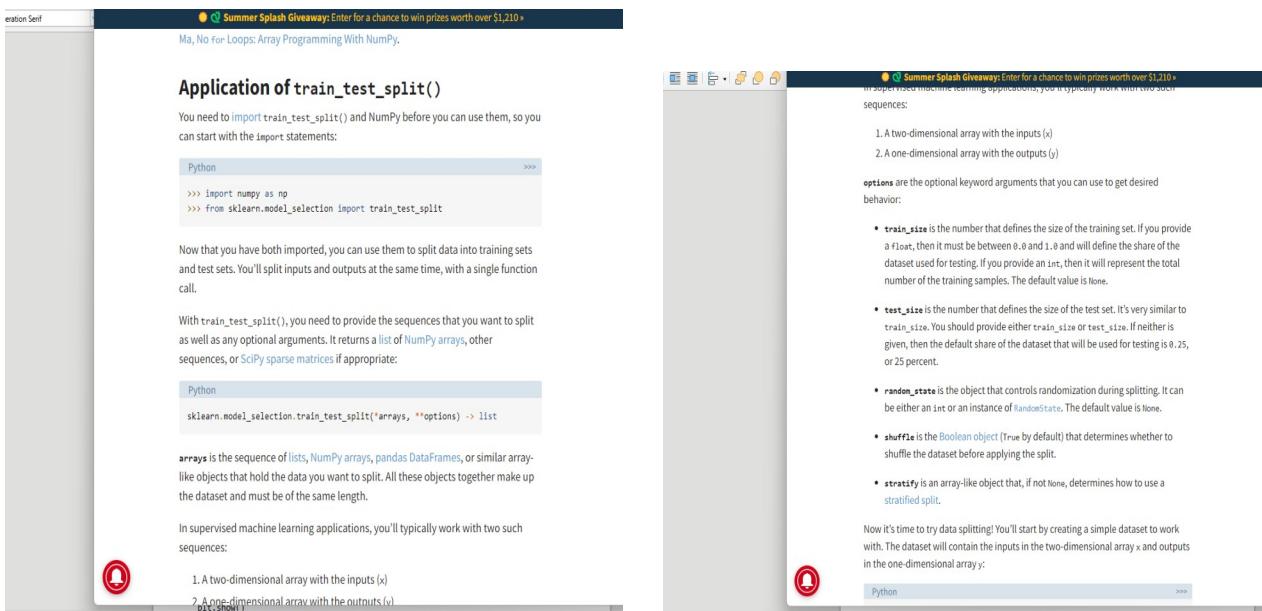
```
4)test_set_size=int(len(data) * test_ratio)
```

setting test size like `506 * test ratio` which will be anything from 0 to 1

```
5)test_indices= shuffled [:test_set_size]  
6)train_indices= shuffled [test_set_size:0]
```

lets say test\_set\_size comes to be 208 out of 506  
so test indices will shuffled data from 0 to 208  
while model will train data from 208 to remaining 506  
indices in short term is like index

7)return data.iloc[train\_indices], data.iloc[test\_indices]  
when you combine `data.iloc`, you are using integer-based indexing to select specific rows and columns from the DataFrame or Series. Here's how you might use it:



8)train\_set,test\_set= split\_train\_test(housing,0.2)  
calling out split train test function with housing variable which we happened to define as dataframe for this project containing data from data.csv and 0.2 test ratio

```
9)print(len(train_set),len(test_set))  
we get values 405 101 as output  
if we go back to this code  
shuffled =np.random.permutation(len(data))
```

and print (shuffled) it will give us random values everytime  
and data in train can end up in test set or vice versa and model can overfit patterns  
and that would be bad practice  
to fix this we add np.random.seed(42)

`p.random.seed(42)` is used in the NumPy library to set the random seed for generating random numbers. This allows you to reproduce the same sequence of random numbers every time you run your code with the same seed value. It's often used in situations where you want to ensure reproducibility of your results, especially in scientific experiments or data analysis where randomness is involved.

All we did we could have directly done with `sklearn.model_selection` library and import `train_test_split` as a function

1)from `sklearn.model_selection` import `train_test_split`

2)`train_set,test_set = train_test_split(housing,test_size=0.2,random_state=42)`

42 for random seed you can pick any number 42 is like industry fav number

`test_ratio 0.25%` by default although can pick anything

housing is data frame we used

3)`print(len(train_set))`

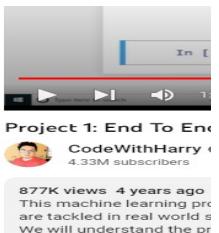
4)`print(len(test_set))`

now lets say we earlier coded this

The screenshot shows a Jupyter Notebook interface with several code cells. The first cell displays the first few rows of the 'housing' dataset. The second cell shows the value counts for the 'CHAS' variable, which has two categories: 0 and 1. The third cell uses the `describe()` method to get summary statistics for all numerical columns. The fourth cell contains a blank command line.

```
505 396.90 7.88 11.9      NaN      NaN      NaN      NaN
[...]
[6]: housing['CHAS'].value_counts()
[6]: CHAS
[6]: 0    471
[6]: 1    35
[6]: Name: count, dtype: int64
[7]: housing.describe()
[7]:          CRIM      ZN      INDUS      CHAS      NOX      RM      AGE      DIS      RAD      TAX

```



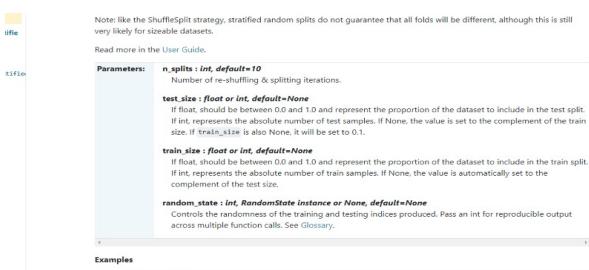
lets say if chas was a very important variable that company needed and 65% of times is 0 how would model distribute test and train data the way where model has somewhat fair 0 and 1 's in both test and train sets so we will use stratified sampling to eliminate this problem

1)from `sklearn.model_selection` import `StratifiedShuffleSplit`

importing `stratified split` function from `sklearn` library

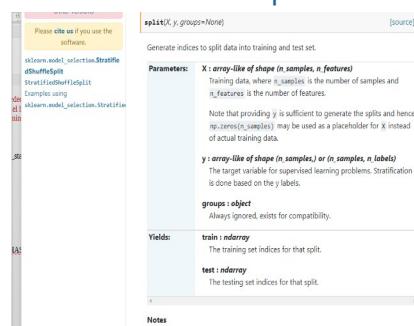
2)`split=StratifiedShuffleSplit(n_splits=1,test_size=0.2,random_state=42)`

default parameters



3)for `train_index, test_index in split.split(housing,housing['CHAS']):`

Generate indices to split data into training and test set.



4)`strat_train_set=housing.loc[train_index]`

he `.loc` accessor is used in pandas, a popular data manipulation library in Python, to select rows and columns from a DataFrame using labels. It provides a way to access data

```

based on the labels of rows
5)strat_train_set=housing.loc[test_index]
6) strat_test_set
shows test set
7)strat_test_set.info()
describes test set
8)strat_test_set['CHAS'].value_counts()
output:
CHAS
0  95
1   7
Name: count, dtype: int64
so 95 entries has 0 as CHAS
7 entries have 1 as CHAS
9)strat_train_set['CHAS'].value_counts()
output:CHAS
0  376
1   28
Name: count, dtype: int64
same for train set

```



same ratio for 0 and 1 's in test set and training set

looking for correlations now

data set we have is very small so we will find out about correlation coefficient

1)corr\_matrix = housing.corr()

In the code `corr_matrix = housing.corr()`, you are calculating the correlation matrix for a DataFrame named `housing`. The correlation matrix is a crucial statistical tool used to understand how variables in a dataset are related to each other in terms of their linear relationships.

`corr()`: This method calculates the Pearson correlation coefficients between all pairs of numeric columns in the DataFrame. The Pearson correlation coefficient measures the linear relationship between two variables. It ranges from -1 to 1, where -1 indicates a perfect negative correlation, 1 indicates a perfect positive correlation, and 0 indicates no linear correlation.

In simple words if one variable increase other might increase or decrease to explain this we use pearson coefficient of correlation

2)corr\_matrix['MEDV'].sort\_values(ascending=False)

medv is median value of occupied homesas you can see in attribute info

```

21    Concerns housing values in suburbs of Boston.
22
23
24 5. Number of Instances: 506
25
26 6. Number of Attributes: 13 continuous attributes (including "class"
27      attribute "MEDV"), 1 binary-valued attribute.
28
29 7. Attribute Information:
30
31 1. CRIM per capita crime rate by town
32 2. ZN proportion of residential land zoned for lots over
33 25,000 sq.ft.
34 3. INDUS proportion of non-retail business acres per town
35 4. CHAS Charles River dummy variable (= 1 if tract bounds
36 river; 0 otherwise)
37 5. NOX nitric oxides concentration (parts per 10 million)
38 6. RM average number of rooms per dwelling
39 7. AGE proportion of owner-occupied units built prior to 1940
40 8. DIS weighted distances to five Boston employment centres
41 9. RAD index of accessibility to radial highways
42 10. TAX full-value property-tax rate per $10,000
43 11. PTRATIO pupil-teacher ratio by town
44 12. B  $1000(Bk - 0.63)^2$  where Bk is the proportion of blacks
45 by town
46 13. LSTAT % lower status of the population
47 14. MEDV Median value of owner-occupied homes in $1000's
48
49 8. Missing Attribute Values: None.
50
51
52

```

ation coefficient

re calculating the

matrix is a crucial

to each other in t

coefficients between

on coefficient mea

, where -1 indicate

tion, and 0 indicate

r decrease to explain

so we passed corr\_matrix via MEDV

`.sort_values(ascending=False)`: This method sorts the correlation values in descending order. The `ascending=False` argument specifies that you want the values in descending order

output:

```

[39]: 13.428571428571429

looking for correlations

[40]: corr_matrix = housing.corr()

[41]: corr_matrix['MEDV'].sort_values(ascending=False)

[42]: MEDV      1.000000
      RM      0.693600
      ZN      0.360445
      B       0.333461
      DIS      0.249929
      CHAS     0.175260
      AGE     -0.376955
      RAD     -0.381626
      crim    -0.383836
      NOX     -0.421221
      TAX      -0.468536
      INDUS    -0.483725
      PTRATIO   -0.507787
      LSTAT    -0.737663
      Unnamed: 14   NaN
      Unnamed: 15   NaN
      Unnamed: 16   NaN
      Unnamed: 17   NaN
      Unnamed: 18   NaN
      Unnamed: 19   NaN
      Name: MEDV, dtype: float64

```

so if we increase medv attribute

rm is in positive correlations

means more the price more the rooms in rough words

AGE -0.376955

here we see age has negative correlation

so homes built prior to 1940 arent that expensive

3)from pandas.plotting import scatter\_matrix

for making graphs using pandas scatter plots

4}attributes=["RM","ZN","MEDV","LSTAT","CHAS"]

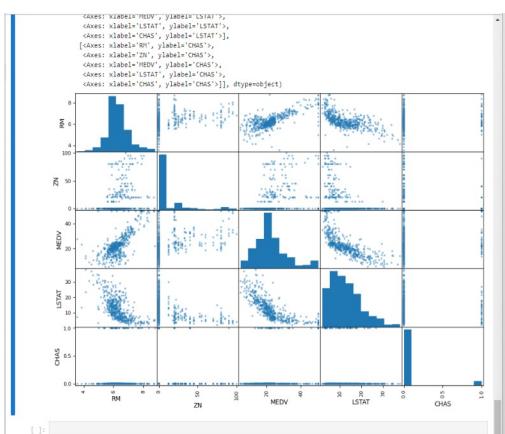
attributes graphs we want to see

```
5)scatter_matrix(housing[attributes],figsize=(12,8))
```

**scatter\_matrix**: This function is used to create a scatter plot matrix. It's part of the pandas plotting library and is used to visualize relationships between multiple variables.

1. **housing[attributes]**: This selects specific columns (attributes) from the DataFrame **housing**. **attributes** is assumed to be a list of column names that you want to include in the scatter plot matrix.
2. **figsize=(12, 8)**: This parameter specifies the size of the figure (plot) in inches. It sets the width to 12 inches and the height to 8 inches.

Output:



notice it doesn't show straight graphs when medv is matched with medv and so on  
it shows histograms to show variations in plotting

to plot one graph

```
6) housing.plot(kind= "scatter",x="RM",y="MEDV",alpha=0.8)
```

where kind = scatter

x and y are mapped to rm and medv

and alpha

seems like you're referring to setting the alpha (transparency) value to 0.8 in a scatter plot. In the context of data visualization, particularly with libraries like Matplotlib in Python, you can use the alpha parameter to control the transparency of data points in a scatter plot. An alpha value of 0.0 makes the points completely transparent, while an alpha value of 1.0 makes them fully opaque.

## Attribute combinations

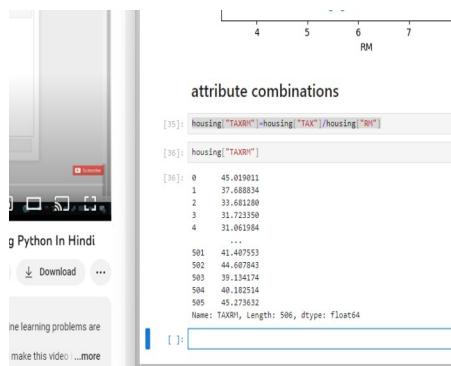
```
1)housing["TAXRM"]=housing["TAX"]/housing["RM"]
```

we created a new attribute taxrm which shows us tax per room

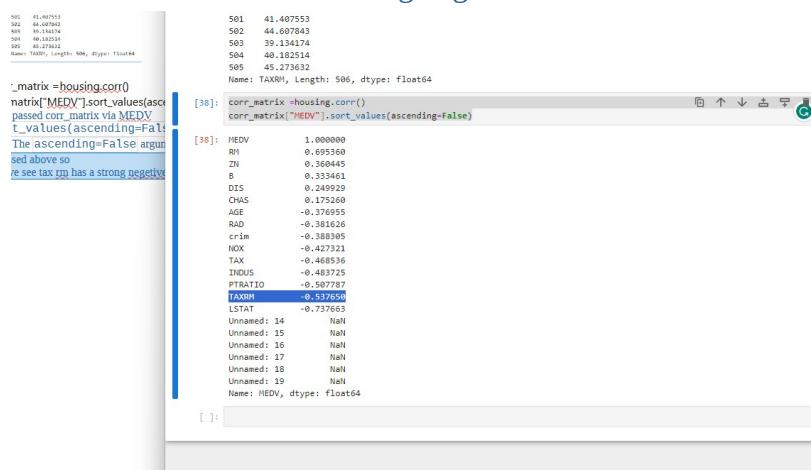
. RM average number of rooms per dwelling

TAX full-value property-tax rate per \$10,000

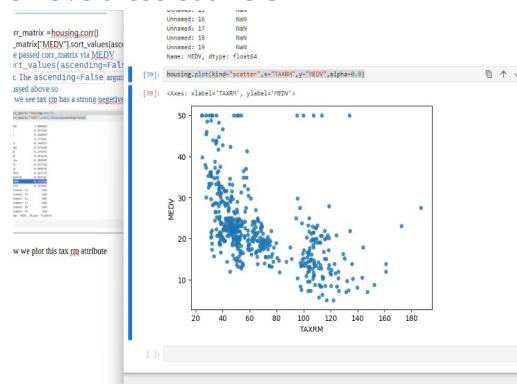
2)housing["TAXRM"]  
to show values of taxrm attribute



3)corr\_matrix =housing.corr()  
corr\_matrix["MEDV"].sort\_values(ascending=False)  
so we passed corr\_matrix via MEDV  
.sort\_values(ascending=False): This method sorts the correlation values in descending order. The ascending=False argument specifies that you want the values in descending order discussed above so here we see tax rm has a strong negative correlation



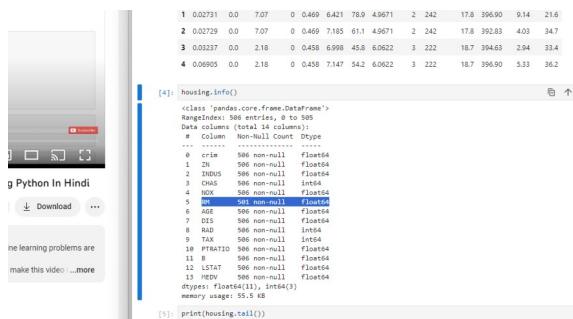
4)now we plot this tax rm attribute  
housing.plot(kind="scatter",x="TAXRM",y="MEDV",alpha=0.8)  
these dots outside of dense cluster are called outliers part of data cleaning is to remove these out liers



if data had missing attributes

for this we deleted some RM entries from data

now we will rerun or restart whole kernal because we changed data it might change past outputs we had



The screenshot shows a Jupyter Notebook cell with the code `[4]: housing.info()`. The output displays information about the 'housing' DataFrame, including column names, data types, and non-null counts. The 'RM' column is listed as having 501 non-null values. The code `[5]: print(housing.tail())` is also visible at the bottom.

```
[4]: housing.info()
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 501 entries, 0 to 500
Data columns (total 14 columns):
 #   Column      Non-Null Count  Dtype  
--- 
 0   crim        500 non-null    float64 
 1   INDUS       500 non-null    float64 
 2   NOBHD      500 non-null    float64 
 3   OHSAS      500 non-null    int64  
 4   NOX         500 non-null    float64 
 5   DIS         500 non-null    float64 
 6   AGE         500 non-null    float64 
 7   DIS         500 non-null    float64 
 8   RAD         500 non-null    float64 
 9   TAX         500 non-null    int64  
 10  PTRATIO     500 non-null    float64 
 11  B           500 non-null    float64 
 12  LSTAT       500 non-null    float64 
 13  MEDV        500 non-null    float64 
dtypes: float64(10), int64(4)
memory usage: 55.5 kB
[5]: print(housing.tail())
```

now see rm has 501 entries now because we deleted 5 entries

To take care of missing Attributes we have 3 options

1) get rid of missing data points

A=housing.dropna(subset=[“RM”])

a.shape

output:501,15

and original housing data frame remains unchanged

2) get rid of whole attribute

housing.drop(“RM”,axis=1)

output :no RM column and original housing data frame remains unchanged

3)set the value to some other value(o ,mean or median)

median=housing[“RM”].median()

housing[“RM”].fillna(median)

using pandas library creates a dataframe median

then fill Rm values having na with median values

can't choose first 2 options because it has positive strong correlation

still original housing data frame is unchanged after 3<sup>rd</sup> option

now to do same stuff easily using library we import simpleimputer from sklearn library

1)from sklearn.impute import SimpleImputer

2)imputer = SimpleImputer(strategy="median")

simple imputer will use median to fill na spots

3)imputer.fit(housing)

to imputer to fit into housing data frame

4)imputer.statistics\_

the `statistics_` attribute is used to access the statistics that were used by the imputer to replace missing values in a dataset.



The screenshot shows a Jupyter Notebook cell with the code for initializing a SimpleImputer with 'median' strategy, fitting it to the 'housing' DataFrame, and then printing its statistics. The statistics are shown as a NumPy array of median values for each feature.

```
[40]: imputer = SimpleImputer(strategy='median')
[41]: imputer.fit(housing)
[42]: SimpleImputer(strategy='median')
[43]: imputer.statistics_
[44]: array([2.56510000e-01, 0.00000000e+00, 9.69000000e+00, 0.00000000e+00,
       5.30000000e-01, 6.20500000e+00, 7.75000000e+01, 3.20745000e+00,
       5.31440000e+00, 1.13000000e+01, 2.12000000e+01, 5.3057273e+01])
[45]: imputer.transform(housing)
[46]: housing_tr=pd.DataFrame(z,columns=housing.columns)
```

5) Q = imputer.transform(housing)

The `transform` method of the `SimpleImputer` instance (`imputer`) is applied to the `housing` dataset, and it returns a new dataset (`Q`) with the missing values replaced using the strategy specified (such as mean, median, etc.).

6) housing\_tr=pd.DataFrame(Q,columns=housing.columns)

It looks like you're creating a new DataFrame called `housing_tr` by using the imputed data `Q` and assigning the columns from the original `housing` DataFrame. This is a common step after performing imputation, as it helps maintain the column names and potentially other metadata. You create the `housing_tr` DataFrame using the imputed data `Q`. By using `columns=housing.columns`, you ensure that the column names from the original `housing` DataFrame are retained in the new `housing_tr` DataFrame.

This is a common practice to keep your data consistent and to have meaningful column names in the imputed dataset.

7) housing\_tr.describe()

```
5.00000000e+00, 3.30000000e+02, 1.90500000e+01, 3.91440000e+02,
1.13600000e+01, 2.12000000e+01, 5.36057273e+01]
[58]: Q = imputer.transform(housing)
[59]: housing_tr=pd.DataFrame(Q,columns=housing.columns)
[60]: housing_tr.describe()
[60]:
```

	crim	ZN	INDUS	CHAS	NOX	RM	AGE	DIS	RAD	TAX
count	506.000000	506.000000	506.000000	506.000000	506.000000	506.000000	506.000000	506.000000	506.000000	506.000000
mean	3.613524	11.363636	11.136779	0.069170	0.554695	6.283500	68.574901	3.795043	9.549407	408.237154
std	8.601545	23.322453	6.860353	0.253994	0.115878	0.702026	28.148861	2.105710	8.707259	168.537116
min	0.006320	0.000000	0.460000	0.000000	0.385000	3.561000	2.900000	1.129600	1.000000	187.000000
25%	0.082045	0.000000	0.519000	0.000000	0.449000	5.885500	45.025000	2.100175	4.000000	279.000000
50%	0.256510	0.000000	9.690000	0.000000	0.538000	6.208000	77.500000	3.207450	5.000000	330.000000
75%	3.677083	12.500000	18.100000	0.000000	0.624000	6.618750	94.075000	5.188425	24.000000	666.000000
max	88.976200	100.000000	27.740000	1.000000	0.871000	8.780000	100.000000	12.126500	24.000000	711.000000

We want to perform this on a train dataset for real projects so before finding correlations we type this

housing=strat\_train\_set.copy()

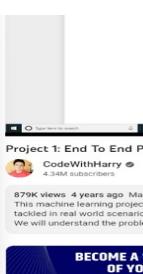
so now its copy of train set

404 entries

so now output will be

```
5.00000000e+00, 3.30000000e+02, 1.90500000e+01, 3.91440000e+02,
1.13600000e+01, 2.12000000e+01, 5.36057273e+01]
[44]: Q = imputer.transform(housing)
[45]: housing_tr=pd.DataFrame(Q,columns=housing.columns)
[46]: housing_tr.describe()
[46]:
```

	crim	ZN	INDUS	CHAS	NOX	RM	AGE	DIS	RAD	TAX
count	404.000000	404.000000	404.000000	404.000000	404.000000	404.000000	404.000000	404.000000	404.000000	404.000000
mean	3.602814	10.836634	11.344950	0.069307	0.558064	6.277856	69.039851	3.746210	9.735149	412.341584
std	8.099383	22.150636	6.877817	0.254290	0.116875	0.712422	28.258248	2.099057	8.731259	168.672623
min	0.006320	0.000000	0.740000	0.000000	0.389000	3.561000	2.900000	1.129600	1.000000	187.000000
25%	0.086962	0.000000	5.190000	0.000000	0.453000	5.878750	44.850000	2.035975	4.000000	284.000000
50%	0.260735	0.000000	9.690000	0.000000	0.538000	6.209000	78.200000	3.122200	5.000000	337.000000
75%	3.731923	12.500000	18.100000	0.000000	0.631000	6.630000	94.100000	5.100400	24.000000	666.000000
max	73.534100	100.000000	27.740000	1.000000	0.871000	8.780000	100.000000	12.126500	24.000000	711.000000



Project 1: End To End P  
CodeWithHarry  
4.3M subscribers

879K Views 4 years ago Mac

This machine learning project

tackled in real world scenarios!

We will understand the problem

BECOME A  
OF YOU

scikit learn has 3 objects  
primarily three types of projects

## 1)ESTIMATORS

estimates some parameters based on datasets eg imputer we used above it has fit method-fits in the datasets and calculates internal parameters

## 2)TRANSFORMERS

transforms method takes input and returns output based on learnings from fit() it also has convenience function called fit\_transform()

## 3)PREDICTORS

linear regression model is an example of predictor

fit() predict() are two common functions it also gives score() function which will evaluate the predictions

feature scaling

two types of feature scaling method

### 1) MIN-MAX SCALING (NORMALIZATION)

(value-min)/(max-min) will come between 0 and 1  
sklearn provides a class called min-max scaler for this

### 2)STANDARDIZATION

(value-mean)/std

sklearn provides a class called standard scaler for this

creating a pipeline

1)from sklearn.pipeline import Pipeline

importing pipeline feature

2)from sklearn.preprocessing import StandardScaler

importing standardisation lib

3)my\_pipeline=Pipeline([  
    ('imputer',SimpleImputer(strategy="median")),  
    ('std\_scaler',StandardScaler()),  
])

**Pipeline:** A **Pipeline** is a way to streamline and automate a sequence of data preprocessing and modeling steps. It ensures that the steps are executed in a specified order, making it easier to manage and deploy machine learning workflows.

1. **('imputer', SimpleImputer(strategy="median")):** In this part of the pipeline, the first step is to handle missing values in the dataset using the **SimpleImputer** class. The **SimpleImputer** is used to replace missing values with a specified strategy, and in this case, the strategy is set to "median". This means that missing values in each column will be replaced with the median value of that column.
2. **('std\_scaler', StandardScaler()):** The second step in the pipeline involves standardizing the data using the **StandardScaler** class. Standardization (also called Z-score normalization) scales the features so that they have a mean of 0 and a standard deviation of

1. This preprocessing step is often used to ensure that all features are on a similar scale, which can help improve the performance of some machine learning algorithms.

To summarize, the given pipeline first replaces missing values with the median value for each column and then standardizes the features by scaling them to have a mean of 0 and a standard deviation of 1. This is a common preprocessing pipeline that is used to prepare data for various machine learning models. You can further extend this pipeline by adding more steps or customizing the existing ones based on the specific requirements of your machine learning task.

4)housing\_num\_tr= my\_pipeline.fit\_transform(housing\_tr)  
fitting pipeline into housing\_tr(training set) and transforming it

5)housing\_num\_tr  
outputs a numpy array  
we need numpy as input

```
[49]: my_pipeline=Pipeline([
    ('imputer',SimpleImputer(strategy="median")),
    ('std_scaler',StandardScaler())
])

[50]: housing_num_tr=my_pipeline.fit_transform(housing_tr)

[51]: housing_num_tr
```

selecting a desired model

1)from sklearn.linear\_model import LinearRegression  
model=LinearRegression()

**Import LinearRegression:** You import the `LinearRegression` class from the `sklearn.linear_model` module. This class provides the functionality to create and train a linear regression model. **Create a LinearRegression Model:** You create an instance of the `LinearRegression` class by calling `LinearRegression()`. This initializes the linear regression model with its default settings.

## 2BEFORE THIS WE ADD THIS CODE BEFORE MISSING ATTRIBUTES

```
+ )housing=strat_train_set.drop("MEDV",axis=1)
```

**strat\_train\_set.drop("MEDV", axis=1)**: This line of code is used to create a new DataFrame named `housing` by dropping the column labeled "MEDV" from the `strat_train_set`.

```
+ )housing_labels=strat_train_set["MEDV"].copy()
```

**strat\_train\_set["MEDV"].copy()**: This line of code is creating a new Series named `housing_labels` by extracting the values from the "MEDV" column of the `strat_train_set` DataFrame. The `.copy()` method ensures that you create a copy of the extracted values rather than just referencing them. This is important because modifications to `housing_labels` won't affect the original DataFrame.

## SUPERATING HOUSING AND HOUSING LABEL

resuming to linear regression

2)model.fit(housing\_num\_tr,housing\_labels)

this is house num tr

The screenshot shows a Jupyter Notebook interface with two code cells and their outputs. Cell [49] defines a Pipeline with an imputer (SimpleImputer(strategy='median')) and a standard scaler (StandardScaler()). Cell [50] runs the pipeline on the training data. Cell [51] prints the transformed training data, which is a NumPy array of shape (n\_samples, n\_features). The array contains numerical values representing the scaled features. The right panel shows the notebook's sidebar with various icons and a message from ChatGPT.

```
[49]: my_pipeline=Pipeline([
    ('imputer',SimpleImputer(strategy="median")),
    ('std_scaler',StandardScaler())
])
[50]: housing_num_tr=my_pipeline.fit_transform(housing_tr)
[51]: housing_num_tr
```

Column 1	Column 2	Column 3	Column 4	Column 5	Column 6	Column 7	Column 8	Column 9	Column 10	Column 11	Column 12	Column 13
-0.43942006	3.12628155	-1.12165014	...	-0.86091034	-0.06501087	-0.59572845	...	-0.123392175	-0.12628155	-1.35893781	...	-0.94116739
...	...	...	...	...	...	...	...	...	...	...	...	...
0.15682329	-0.4898311	0.98236806	...	0.81480158	-0.61974213	1.12210986	...	-0.43526567	-0.4898311	-1.23083158	...	-1.27603303
...	...	...	...	...	...	...	...	2.43126115	-1.02383378	...	...	...
...	...	...	...	...	...	...	...	0.14321609	-0.4898311	0.98336806	...	0.73869575
...	...	...	...	...	...	...	...	-0.89710776	1.10156693	...	...	...
...	...	...	...	...	...	...	...	-0.43974024	-0.4898311	0.37049623	...	0.09940681
...	...	...	...	...	...	...	...	0.08433985	-0.6552217	...	...	...

housing labels we created above by extracting values from MEDV

3)some\_data=housing.iloc[:5]

takes 5 enteries from hosuing dataframe

4)housing\_labels=housing\_labels.iloc[:5]

takes five enteries from MEDV from train set

5)prepared\_data=my\_pipeline.transform(some\_data)

5 enteries from housing gets transmormed into pipeline

5)model.predict(prepared\_data)

gives ouput predictions

array([23.94921288, 27.21581435, 20.57575334, 25.06681527, 23.77281415])

now we compare it to our labels

6)some\_labels

```
254  21.9  
348  24.5  
476  16.7  
321  23.1  
326  23.0
```

Name: MEDV, dtype: float64

gives output like this

compare it to predictions they arent that good last one good though

convert it to list

7)list(some\_labels)

we get output like

```
[21.9, 24.5, 16.7, 23.1, 23.0]
```

```
[ ]:
```

evaluating the model

1)from sklearn.metrics import mean\_squared\_error

This statement imports the `mean_squared_error` function from the `sklearn.metrics` module. This function is commonly used to calculate the mean squared error (MSE) between two sets of data, typically predictions and actual target values. The mean squared error is a common metric used to evaluate the performance of regression models.

2)housing\_predictions= model.predict(housing\_num\_tr)

where `model= linear regression()`

and `housing_num_tr= my_pipeline.fit_transform(housing_tr)`

where `housing tr` is

The screenshot shows a Jupyter Notebook interface with several code cells and their outputs. The code cells are numbered [43], [44], [45], [46], and [47]. Cell [43] contains `imputer.statistics_`. Cell [44] contains `b = imputer.transform(housing)`. Cell [45] contains `housing_tr=pd.DataFrame(b,columns=housing.columns)`. Cell [46] contains `housing_tr.describe()`, which displays a detailed summary statistics table for the dataset. Cell [47] contains the text "scikit-learn design" and "creating pipeline". The right side of the screen shows a file browser with various files and folders.

```
[43]: imputer.statistics_
[43]: array([ 2.86735e-01,  0.00000e+00,  9.00000e+00,  0.00000e+00,  5.38000e-01,
       6.20900e+00,  7.82000e+01,  3.12200e+00,  5.00000e+00,  3.37000e+02,
      1.90000e+01,  3.30955e+02,  1.15700e+01])

[44]: b = imputer.transform(housing)

[45]: housing_tr=pd.DataFrame(b,columns=housing.columns)

[46]: housing_tr.describe()
    crim      ZN      INDUS      CHAS      NOX      RM      AGE      DIS      RAD      TAX
count  404.000000  404.000000  404.000000  404.000000  404.000000  404.000000  404.000000  404.000000  404.000000
mean   3.602814  10.33634  11.344950  0.69307  0.558064  6.277856  69.039851  3.746210  9.735149  412.341584
std    8.099383  22.150636  6.877917  0.234290  0.116875  0.712422  28.258248  2.099057  8.731259  166.672623
min    0.006320  0.000000  0.740000  0.000000  0.389000  3.561000  2.900000  1.129600  1.000000  187.000000
25%   0.086982  0.000000  5.190000  0.000000  0.453000  5.878750  44.850000  2.035975  4.000000  284.000000
50%   0.286735  0.000000  9.900000  0.000000  0.538000  6.209000  78.200000  3.122300  5.000000  337.000000
75%   3.731923  12.500000  18.100000  0.000000  0.631000  6.630000  94.100000  5.100400  24.000000  666.000000
max   73.534100  100.000000  27.740000  1.000000  0.871000  8.780000  100.000000  12.126500  24.000000  711.000000
```

3)lin\_mse=mean\_squared\_error(housing\_labels,housing\_predictions)

mse is mean squared error function of lib sk learn

where we declared housing predictions above and

housing\_labels we know from this code we wrote before

`housing=strat_train_set.drop("MEDV",axis=1)`

`housing_labels=strat_train_set["MEDV"].copy()`

4)lin\_rmse=np.sqrt(lin\_mse)

we calculate root mean square error

```
[67]: lin_mse  
[67]: 23.352925815303564  
  
[68]: lin_rmse  
[68]: 4.832486504409874
```

```
[ ]:
```



but calculated housing\_label values were  
[21.9, 24.5, 16.7, 23.1, 23.0]

23 very less than these

so we will discard this model

now we will use decision tree regressor

so we will alter above code

we converted code from linear regression to decision tree

commented out linear regression using #

## selecting a desired model

```
*[52]: #from sklearn.linear_model import LinearRegression  
      from sklearn.tree import DecisionTreeRegressor  
  
[55]: #model=LinearRegression()  
      model=DecisionTreeRegressor()  
  
[56]: model.fit(housing_num_tr,housing_labels)  
[56]: ▾ DecisionTreeRegressor  
      DecisionTreeRegressor()
```

```
[ ]: some_data=housing.iloc[:5]
```



```
[ ]: some_labels=housing_labels.iloc[:5]
```

```
[ ]: prepared_data=my_pipeline.transform(some_data)
```

```
[ ]: model.predict(prepared_data)
```

```
[ ]: list(some_labels)
```

but now our root mean square error goes 0 how?

We overfitted our training data and decision tree cleared all noise in training data

Anyways we rename them as rmse and mse because they are not part of lin\_ (linear regression now)

```
[64]: lin_mse
```

```
[64]: 0.0
```

```
[65]: lin_rmse
```

```
[65]: 0.0
```

```
[ ]:
```



using cross validation

1) from sklearn.model\_selection import cross\_val\_score

imports cross val func from sk leanr libraries

2)scores=cross\_val\_score(model,housing\_num\_tr,housing\_labels,scoring="neg\_mean\_squared\_error",cv=10)

**cross\_val\_score**: This function is part of scikit-learn's **model\_selection** module and is used for performing cross-validation. It takes several parameters:

- **model**: This should be your trained machine learning model that you want to evaluate.
- **cv=10**: This parameter specifies the number of folds for cross-validation. In this case, you're using 10-fold cross-validation, meaning the data will be split into 10 parts, and the model will be trained and evaluated 10 times.
- **eg\_mean\_squared\_error**: In scikit-learn's convention, scoring functions typically consider higher values to be better. Since lower values of MSE indicate better performance (as they indicate smaller prediction errors), the **neg\_mean\_squared\_error** scoring function returns the negation of the mean squared error. This is done to align with the convention of higher scores being better.

3)rmse\_scores=np.sqrt(-scores)

scores are already negative so we have to -scores to root square it since

**neg\_mean\_squared\_error** scoring function returns the negation of the mean squared error. This is done to align with the convention of higher scores being better.

4)rmse\_scores

outputs :array([3.89223743, 5.73748904, 4.98666514, 3.84111264, 3.81562577, 2.93027302, 6.91642971, 3.58514295, 3.40899545, 3.88052831])

[ ]:

less error compared to linear regression

using linear regression

array([4.27001096, 4.26034932, 5.11177484, 3.83151625, 5.3435416 , 4.38995448, 7.46842271, 5.49003038, 4.1392956 , 6.05405933])

so we choose decision tree regressor

5)def print\_scores(scores):

    print("scores:",scores)

    print("Mean:",scores.mean())

    print("standardDeviation:",scores.std())

created a function to print scores its mean and its std deviation

these functions included in numpy mean std mean and all

```
7)print_scores(rmse_scores)
output not we using rmse scores given by decision tree regressor
scores: [4.06276969 5.73236426 5.25725979 4.02922252 3.62287869 2.66622392
 6.82764601 3.82298967 3.51983664 4.29918597]
Mean: 4.384037716238515
standardDeviation: 1.1567149701246144
```

using linear regression we get output:

```
scores: [4.27001096 4.26034932 5.11177484 3.83151625 5.3435416 4.38995448
 7.46842271 5.49003038 4.1392956 6.05405933]
Mean: 5.035895549223018
standardDeviation: 1.0545235927112007
```

using random forest regressor

```
selecting a desired model
[150]: #from sklearn.linear_model import LinearRegression
#from sklearn.tree import DecisionTreeRegressor
from sklearn.ensemble import RandomForestRegressor

[151]: model=LinearRegression()
model=DecisionTreeRegressor()
model=RandomForestRegressor()

[152]: model.fit(housing_num_tr,housing_labels)

[153]: > RandomForestRegressor()
RandomForestRegressor()

[154]: some_data=housing.iloc[:5]
some_labels=housing_labels.iloc[:5]

[155]: prepared_data=my_pipeline.transform(some_data)

[156]: model.predict(prepared_data)
```

we get this output

```
scores: [3.02807391 2.76676348 4.36154577 2.67570716 3.29508102 2.65835356
 4.86175816 3.27670055 3.11379993 3.40758445]
Mean: 3.3445367979314087
standardDeviation: 0.6893969355768271
```

randomforest regressor works better than decision tree regressor and linear regressor

note: Overfitting is a common problem in machine learning where a model learns the training data too well, to the point that it captures the noise and random fluctuations in the data rather than just the underlying patterns. As a result, an overfitted model may perform very well on the training data but fails to generalize well to new, unseen data. In essence, the model becomes too specialized to the training data and loses its ability to make accurate predictions on new data.

That was the thing with decision tree regressor before we cross validated it

saving the model

1) from joblib import dump,load

The `joblib` library in Python provides tools for efficiently saving and loading Python objects, especially large numerical arrays, to and from disk. It is commonly used for saving machine learning models, trained parameters, and other data structures. The `dump` and `load` functions are two important functions from the `joblib` library.

```
2)dump(model, 'my.joblib')
```

**dump** function is used to save Python objects to disk in a binary format that is optimized for speed and efficiency. This function is often used to save trained machine learning models, so you can reuse them later without retraining.

**load Function:** The **load** function is used to load the previously saved objects back into memory. It allows you to retrieve and use the objects without retraining or reprocessing.

Testing model on test data

```
1)x_test=strat_test_set.drop("MEDV",axis=1)
```

```
2)y_test=strat_test_set["MEDV"].copy()
```

already done before attribute thing

```
3)x_test_prepared=my_pipeline.transform(x_test)
```

apply pipeline to transorm x test

my pipeline we definded before

```
4)final_predictions= model.predict(x_test_prepared)
```

predicting from randomforest regressor

```
5)final_mse=mean_squared_error(y_test,final_predictions)
```

applied room=t\_mean\_squared functions to final predictions using randomforestregressor this time on training set

```
6)final_rmse=np.sqrt(final_mse)
```

applied same as before

```
7)final_rmse
```

error reduces output:3.0425337560761325

```
8)print(final_predictions,list(y_test))
```

model working very good

```
[89]: print(final_predictions,list(y_test))

[24, 303 11.48 25.335 22.191 18.377 14.64 19.71 14.24 33.738 42.945
 19.334 11.695 24.725 28.828 19.407 11.084 31.132 14.443 23.865 18.869
 20.181 17.222 17.469 21.737 18.433 31.233 16.895 32.767 9.027 33.388
 24.224 21.409 23.573 10.972 21.263 11.298 43.412 24.339 23.795 41.527
 24.1 29.971 20.561 20.648 18.711 33.537 44.6 20.243 20.681 21.246
 21.189 20.462 20.556 20.556 20.574 40.157 28.085 20.556 23.728
 47.016 18.067 18.805 24.823 14.463 23.1 40.058 30.517 19.465 19.103
 25.886 23.807 21.159 22.669 34.381 13.304 15.467 20.065 20.827 21.348
 22.638 21.184 14.068 23.127 19.974 21.379 14.095 21.136 23.012 23.235
 18.678 26.594 7.311 26.572 19.36 29.363 20.028 31.977 14.595 26.699
 21.107 19.989 14.655 18.2 30.3 23.6 14.4 15.4 19.4 20.4 30.3 35.2 23.1, 13.8, 25.0, 27.9, 19.5, 12.3, 3
 2.1, 12.1, 18.3, 18.4, 17.1, 19.3, 4.4, 21.2, 16.6, 15.2, 18.9, 7.2, 34.6, 20.1, 20.6, 23.6, 13.1, 23.8, 12.
 7, 43.1, 24.7, 22.2, 44.0, 28.1, 31.0, 21.7, 23.4, 19.5, 33.1, 41.7, 18.7, 18.9, 20.6, 21.2, 13.6, 20.3, 17.8, 27.
 1, 11.5, 50.5, 29.1, 18.9, 20.4, 50.0, 7.2, 17.2, 36.2, 14.6, 33.2, 23.8, 19.9, 21.5, 37.3, 27.0, 22.0, 24.3, 19.8,
 33.3, 7.0, 19.4, 20.9, 21.1, 20.4, 22.2, 11.9, 11.7, 21.6, 19.7, 23.0, 16.7, 21.7, 20.6, 23.3, 19.6, 28.0, 5.0, 24.
 4, 20.8, 24.7, 21.8, 23.6, 19.0, 25.0, 20.3, 21.5]
```