

Лабораторная работа №1
по дисциплине
«Технологии машинного обучения»
на тему
«Разведочный анализ данных. Исследование и
визуализация данных»

Выполнила:
студент группы ИУ5-646
Подопригорова С. С.

0.0.1. Задание

- Выбрать набор данных (датасет). Вы можете найти список свободно распространяемых датасетов [здесь](#).
- Для первой лабораторной работы рекомендуется использовать датасет без пропусков в данных, например из Scikit-learn.
- Пример преобразования датасетов Scikit-learn в Pandas Dataframe можно посмотреть [здесь](#). Для лабораторных работ не рекомендуется выбирать датасеты большого размера.
- Создать ноутбук, который содержит следующие разделы:
 1. Текстовое описание выбранного Вами набора данных.
 2. Основные характеристики датасета.
 3. Визуальное исследование датасета.
 4. Информация о корреляции признаков.
- Сформировать отчет и разместить его в своем репозитории на github.

1. Дома в Бостоне

```
[25]: import numpy as np
import pandas as pd
import seaborn as sns

import matplotlib.pyplot as plt
%matplotlib inline
```

```
[26]: from sklearn.datasets import load_boston
```

```
[27]: data = load_boston()
```

```
[5]: for x in data:
    print(x)
```

```
data
target
feature_names
DESCR
filename
```

```
[6]: X, y = data['data'], data['target']
```

```
[7]: print(data['DESCR'])
```

```
.. _boston_dataset:
```

```
Boston house prices dataset
```

```
-----
```

****Data Set Characteristics:****

:Number of Instances: 506

:Number of Attributes: 13 numeric/categorical predictive. Median Value (attribute 14) is usually the target.

:Attribute Information (in order):

- CRIM per capita crime rate by town
- ZN proportion of residential land zoned for lots over 25,000 sq.ft.
- INDUS proportion of non-retail business acres per town
- CHAS Charles River dummy variable (= 1 if tract bounds river; 0 otherwise)
- NOX nitric oxides concentration (parts per 10 million)
- RM average number of rooms per dwelling
- AGE proportion of owner-occupied units built prior to 1940
- DIS weighted distances to five Boston employment centres
- RAD index of accessibility to radial highways
- TAX full-value property-tax rate per \$10,000
- PTRATIO pupil-teacher ratio by town
- B $1000(B_k - 0.63)^2$ where B_k is the proportion of blacks by town
- LSTAT % lower status of the population
- MEDV Median value of owner-occupied homes in \$1000's

:Missing Attribute Values: None

:Creator: Harrison, D. and Rubinfeld, D.L.

This is a copy of UCI ML housing dataset.

<https://archive.ics.uci.edu/ml/machine-learning-databases/housing/>

This dataset was taken from the StatLib library which is maintained at [Carnegie Mellon University](#).

The Boston house-price data of Harrison, D. and Rubinfeld, D.L. 'Hedonic prices and the demand for clean air', J. Environ. Economics & Management, vol.5, 81-102, 1978. Used in Belsley, Kuh & Welsch, 'Regression diagnostics ...', Wiley, 1980. N.B. Various transformations are used in the table on pages 244-261 of the latter.

The Boston house-price data has been used in many machine learning papers that address regression problems.

.. topic:: References

- Belsley, Kuh & Welsch, 'Regression diagnostics: Identifying Influential Data and Sources of Collinearity', Wiley, 1980. 244-261.
- Quinlan, R. (1993). Combining Instance-Based and Model-Based Learning. In Proceedings on the Tenth International Conference of Machine Learning, [University of Massachusetts, Amherst. Morgan Kaufmann.](#) 236-243,

```

-
-
- INDUS
- CHAS Charles River (= 1,
- NOX ( 10 )
-
-
- DIS
-
- PTRATIO
- B 1000 (Bk - 0,63) ^ 2, Bk -
- LSTAT%
- MEDV 1000

```

Предсказываем среднюю стоимость частных домов в тысячах долларов.

```
[9]: data.feature_names
```

```
[9]: array(['CRIM', 'ZN', 'INDUS', 'CHAS', 'NOX', 'RM', 'AGE', 'DIS', 'RAD',
        'TAX', 'PTRATIO', 'B', 'LSTAT'], dtype='<U7')
```

```
[10]: X_df = pd.DataFrame(data = X, columns=data.feature_names)
      y_df = pd.Series(y)
```

```
[31]: data = X_df.merge(y_df.rename('MEDV'), left_index=True, right_index=True)
```

```
[32]: data.head()
```

```
[32]:
```

	CRIM	ZN	INDUS	CHAS	NOX	RM	AGE	DIS	RAD	TAX	\
0	0.00632	18.0	2.31	0.0	0.538	6.575	65.2	4.0900	1.0	296.0	
1	0.02731	0.0	7.07	0.0	0.469	6.421	78.9	4.9671	2.0	242.0	
2	0.02729	0.0	7.07	0.0	0.469	7.185	61.1	4.9671	2.0	242.0	
3	0.03237	0.0	2.18	0.0	0.458	6.998	45.8	6.0622	3.0	222.0	
4	0.06905	0.0	2.18	0.0	0.458	7.147	54.2	6.0622	3.0	222.0	

	PTRATIO	B	LSTAT	MEDV
0	15.3	396.90	4.98	24.0
1	17.8	396.90	9.14	21.6
2	17.8	392.83	4.03	34.7
3	18.7	394.63	2.94	33.4
4	18.7	396.90	5.33	36.2

```
[33]: data.shape
```

```
[33]: (506, 14)
```

```
[34]: data.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(14)
memory usage: 55.5 KB
```

Все столбцы ненулевые.

```
[35]: data.describe()
```

```
[35]:
```

	CRIM	ZN	INDUS	CHAS	NOX	
↪RM \						
count	506.000000	506.000000	506.000000	506.000000	506.000000	506.000000
mean	3.613524	11.363636	11.136779	0.069170	0.554695	6.284634
std	8.601545	23.322453	6.860353	0.253994	0.115878	0.702617
min	0.006320	0.000000	0.460000	0.000000	0.385000	3.561000
25%	0.082045	0.000000	5.190000	0.000000	0.449000	5.885500
50%	0.256510	0.000000	9.690000	0.000000	0.538000	6.208500
75%	3.677083	12.500000	18.100000	0.000000	0.624000	6.623500
max	88.976200	100.000000	27.740000	1.000000	0.871000	8.780000

	AGE	DIS	RAD	TAX	PTRATIO	
↪B \						
count	506.000000	506.000000	506.000000	506.000000	506.000000	506.000000
mean	68.574901	3.795043	9.549407	408.237154	18.455534	356.674032
std	28.148861	2.105710	8.707259	168.537116	2.164946	91.294864
min	2.900000	1.129600	1.000000	187.000000	12.600000	0.320000
25%	45.025000	2.100175	4.000000	279.000000	17.400000	375.377500
50%	77.500000	3.207450	5.000000	330.000000	19.050000	391.440000

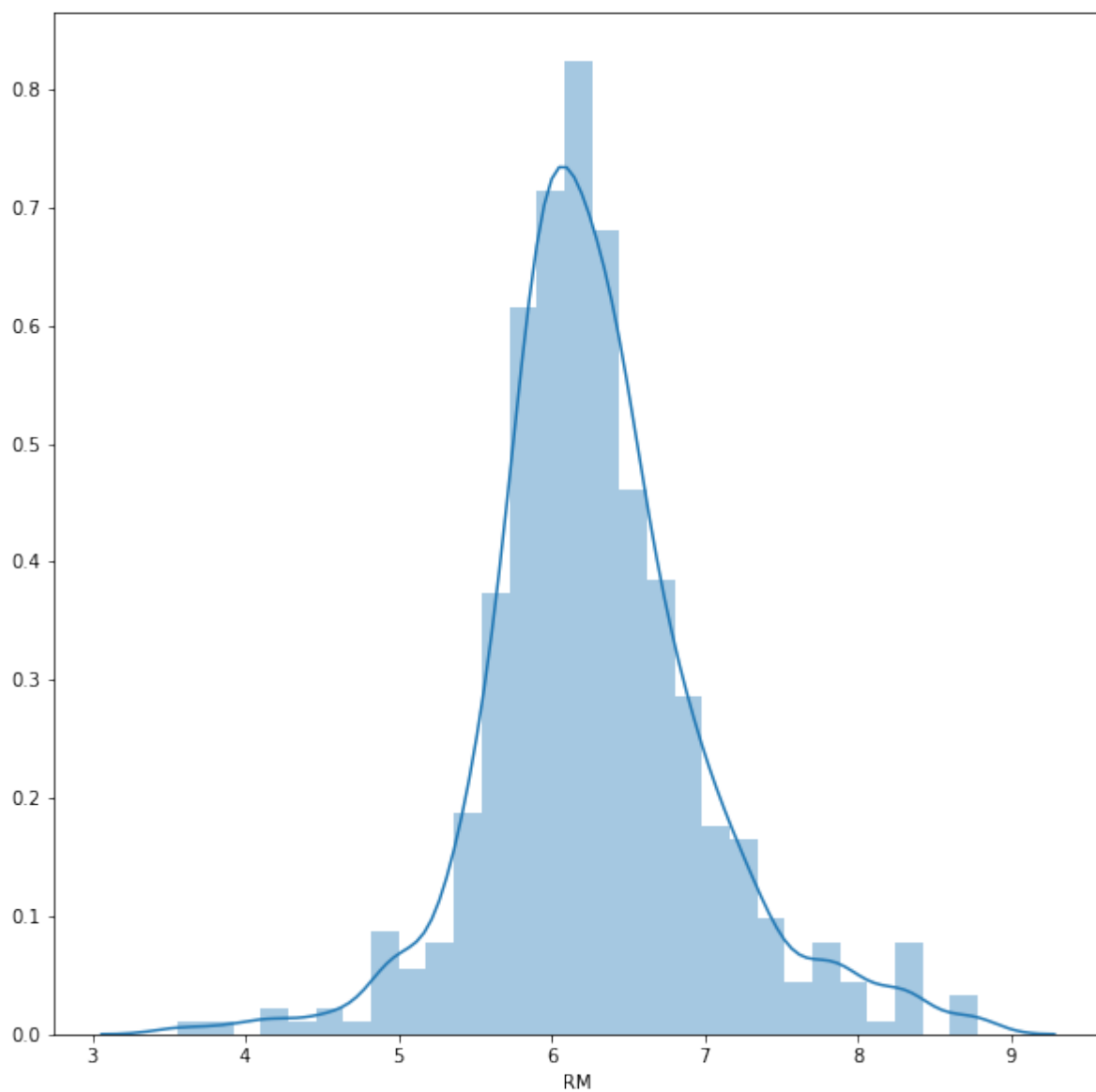
75%	94.075000	5.188425	24.000000	666.000000	20.200000	396.225000
max	100.000000	12.126500	24.000000	711.000000	22.000000	396.900000

	LSTAT	MEDV
count	506.000000	506.000000
mean	12.653063	22.532806
std	7.141062	9.197104
min	1.730000	5.000000
25%	6.950000	17.025000
50%	11.360000	21.200000
75%	16.955000	25.000000
max	37.970000	50.000000

2. Визуальное исследование датасета

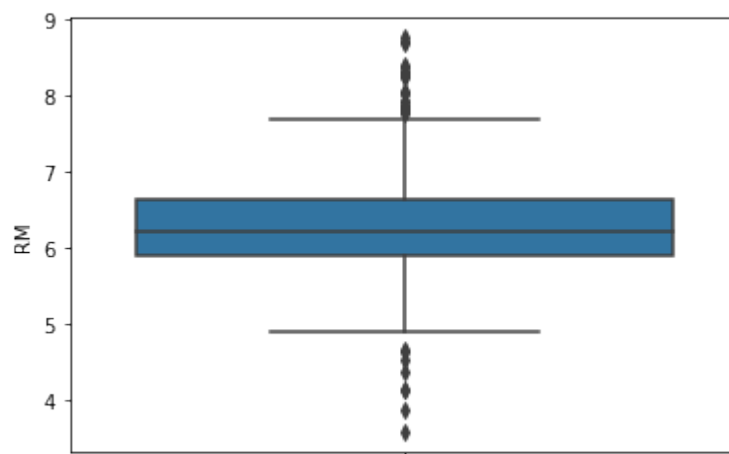
```
[46]: fig, ax = plt.subplots(figsize=(10,10))
      sns.distplot(data['RM'])
```

```
[46]: <matplotlib.axes._subplots.AxesSubplot at 0x7ff5730d15b0>
```



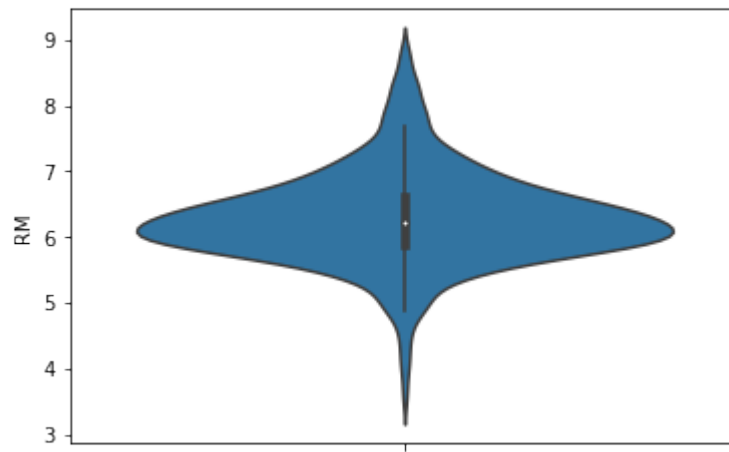
```
[47]: sns.boxplot(y=data['RM'])
```

```
[47]: <matplotlib.axes._subplots.AxesSubplot at 0x7ff5726400d0>
```



```
[48]: sns.violinplot(y=data['RM'])
```

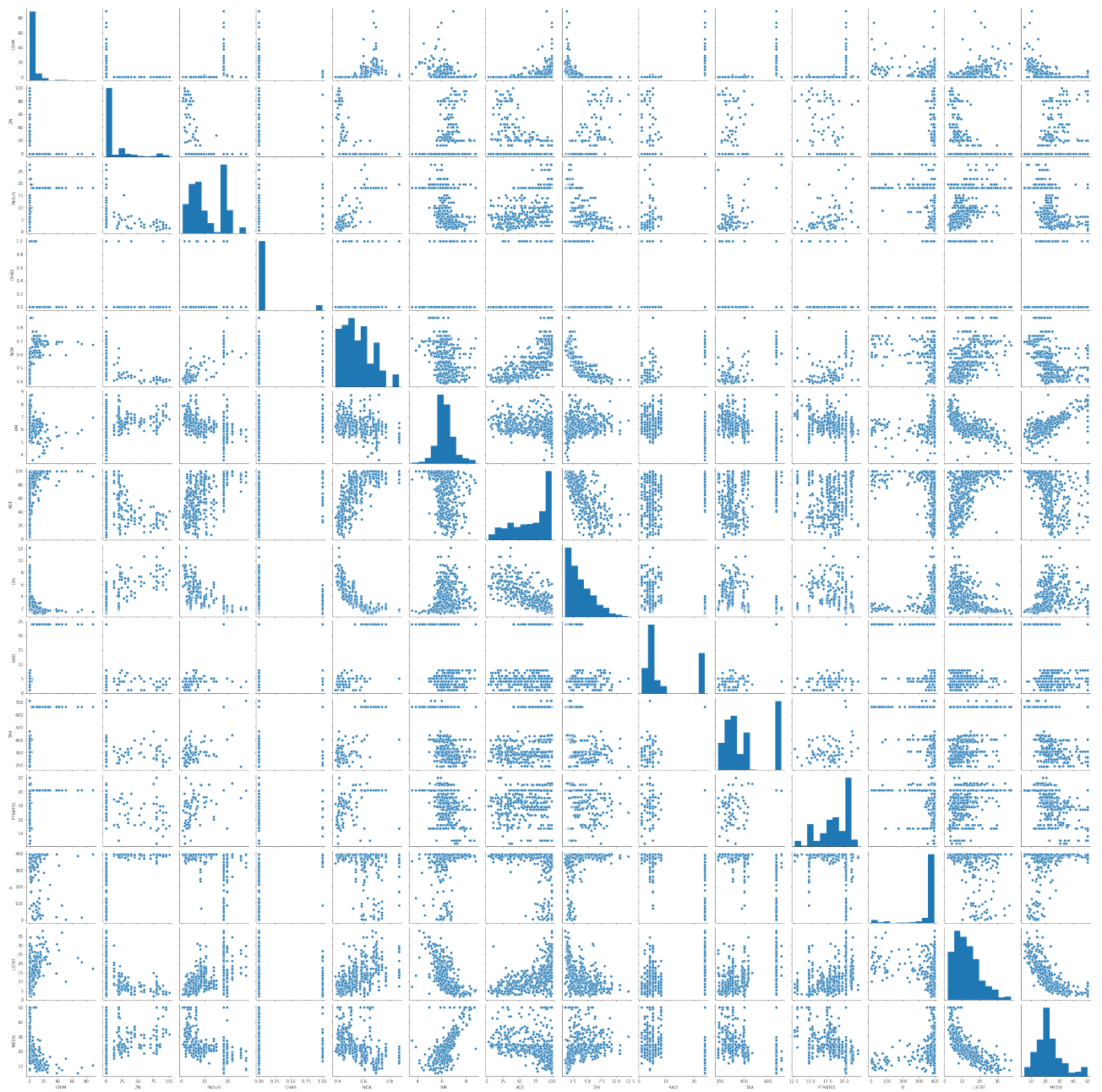
```
[48]: <matplotlib.axes._subplots.AxesSubplot at 0x7ff5727529d0>
```



Количество комнат в домах в Бостоне распределено нормально, в среднем комнат 6

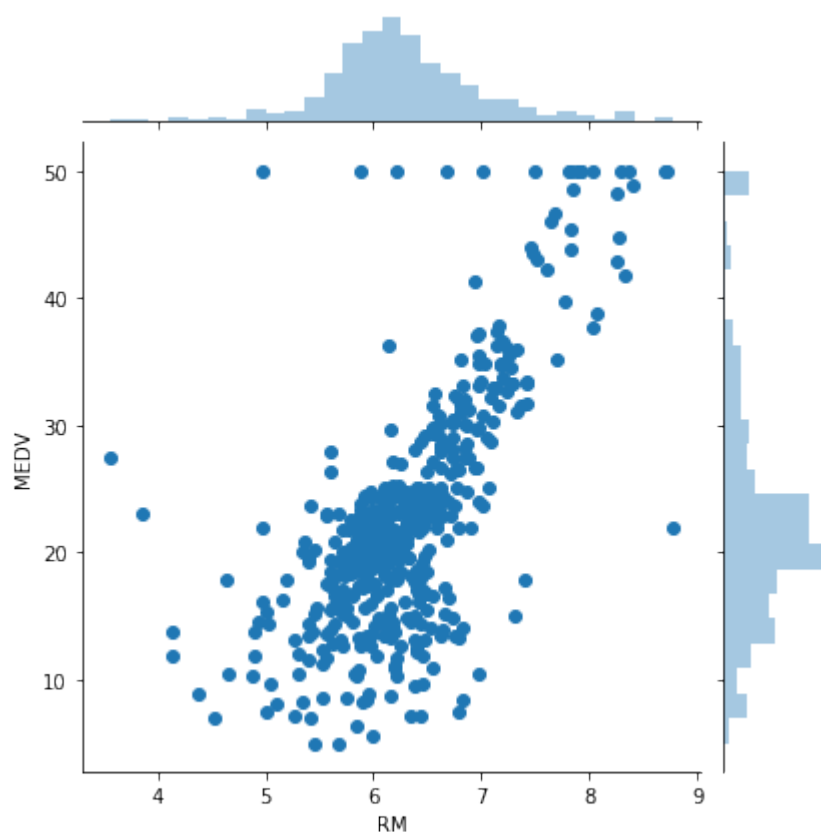
```
[39]: sns.pairplot(data)
```

```
[39]: <seaborn.axisgrid.PairGrid at 0x7ff582c10250>
```

```
[50]: sns.jointplot(x='RM', y='MEDV', data=data)
```

```
[50]: <seaborn.axisgrid.JointGrid at 0x7ff573b0ef40>
```



Прослеживается линейная зависимость между количеством комнат в доме и его ценой

3. Информация о корреляции признаков

[22]: `data.corr()`

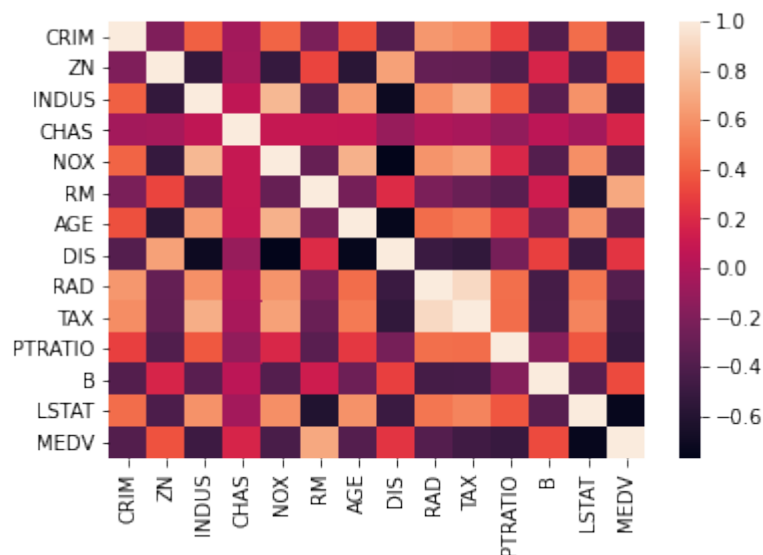
[22]:

	CRIM	ZN	INDUS	CHAS	NOX	RM	
↪AGE \							
CRIM	1.000000	-0.200469	0.406583	-0.055892	0.420972	-0.219247	0.352734
ZN	-0.200469	1.000000	-0.533828	-0.042697	-0.516604	0.311991	-0.569537
INDUS	0.406583	-0.533828	1.000000	0.062938	0.763651	-0.391676	0.644779
CHAS	-0.055892	-0.042697	0.062938	1.000000	0.091203	0.091251	0.086518
NOX	0.420972	-0.516604	0.763651	0.091203	1.000000	-0.302188	0.731470
RM	-0.219247	0.311991	-0.391676	0.091251	-0.302188	1.000000	-0.240265
AGE	0.352734	-0.569537	0.644779	0.086518	0.731470	-0.240265	1.000000
DIS	-0.379670	0.664408	-0.708027	-0.099176	-0.769230	0.205246	-0.747881
RAD	0.625505	-0.311948	0.595129	-0.007368	0.611441	-0.209847	0.456022
TAX	0.582764	-0.314563	0.720760	-0.035587	0.668023	-0.292048	0.506456
PTRATIO	0.289946	-0.391679	0.383248	-0.121515	0.188933	-0.355501	0.261515
B	-0.385064	0.175520	-0.356977	0.048788	-0.380051	0.128069	-0.273534
LSTAT	0.455621	-0.412995	0.603800	-0.053929	0.590879	-0.613808	0.602339
MEDV	-0.388305	0.360445	-0.483725	0.175260	-0.427321	0.695360	-0.376955
	DIS	RAD	TAX	PTRATIO	B	LSTAT	MEDV

CRIM	-0.379670	0.625505	0.582764	0.289946	-0.385064	0.455621	-0.388305
ZN	0.664408	-0.311948	-0.314563	-0.391679	0.175520	-0.412995	0.360445
INDUS	-0.708027	0.595129	0.720760	0.383248	-0.356977	0.603800	-0.483725
CHAS	-0.099176	-0.007368	-0.035587	-0.121515	0.048788	-0.053929	0.175260
NOX	-0.769230	0.611441	0.668023	0.188933	-0.380051	0.590879	-0.427321
RM	0.205246	-0.209847	-0.292048	-0.355501	0.128069	-0.613808	0.695360
AGE	-0.747881	0.456022	0.506456	0.261515	-0.273534	0.602339	-0.376955
DIS	1.000000	-0.494588	-0.534432	-0.232471	0.291512	-0.496996	0.249929
RAD	-0.494588	1.000000	0.910228	0.464741	-0.444413	0.488676	-0.381626
TAX	-0.534432	0.910228	1.000000	0.460853	-0.441808	0.543993	-0.468536
PTRATIO	-0.232471	0.464741	0.460853	1.000000	-0.177383	0.374044	-0.507787
B	0.291512	-0.444413	-0.441808	-0.177383	1.000000	-0.366087	0.333461
LSTAT	-0.496996	0.488676	0.543993	0.374044	-0.366087	1.000000	-0.737663
MEDV	0.249929	-0.381626	-0.468536	-0.507787	0.333461	-0.737663	1.000000

```
[23]: sns.heatmap(data.corr())
```

```
[23]: <matplotlib.axes._subplots.AxesSubplot at 0x7ff5999d38b0>
```



- Лучше всего цена дома коррелирует со статусом населения LSTAT (-0.737), количеством комнат RM (0.695) и соотношением учеников и учителей PTRATIO (0.508)
- Меньше всего на цену дома влияет наличие водоёмов CHAS (0,175) и расстояние до бостонских центров занятости DIS (0,249). Эти признаки лучше исключить.

```
[ ]:
```