

Алёшин Александр Денисович ИУ5-63Б. РК1. Технологии разведочного анализа и обработки данных.

```
In [55]: import numpy as np
import matplotlib.pyplot as plt

import pandas as pd
import seaborn as sns

%matplotlib inline

In [27]: from sklearn.datasets import load_boston
boston_dataset = load_boston()
#boston_dataset = pd.read_csv('housing_data1.txt', sep="\s+|\t+|\s+\t+|\t+\s+", engine = 'python')

In [29]: print(boston_dataset.keys())

dict_keys(['data', 'target', 'feature_names', 'DESCR', 'filename'])

In [30]: boston_dataset.DESCR

"..._boston_dataset:\n\nBoston house prices dataset\n-----\n\n**Data Set Characteristics:** \n\n :Number of Instances: 506 \n\n :Numb
er of Attributes: 13 numeric/categorical predictive. Median Value (attribute 14) is usually the target.\n\n :Attribute Information (in order):\n      - CRIM
per capita crime rate by town\n      - ZN      proportion of residential land zoned for lots over 25,000 sq.ft.\n      - INDUS      proportion of non-retail bu
siness acres per town\n      - CHAS      Charles River dummy variable (= 1 if tract bounds river; 0 otherwise)\n      - NOX      nitric oxides concentration (p
arts per 10 million)\n      - RM      average number of rooms per dwelling\n      - AGE      proportion of owner-occupied units built prior to 1940\n
- DIS      weighted distances to five Boston employment centres\n      - RAD      index of accessibility to radial highways\n      - TAX      full-value propo
rty-tax rate per $10,000\n      - PTRATIO      pupil-teacher ratio by town\n      - B      1000(Bk - 0.63)^2 where Bk is the proportion of black people by town
\n      - LSTAT      % lower status of the population\n      - MEDV      Median value of owner-occupied homes in $1000's\n\n :Missing Attribute Values: None\n
\n :Creator: Harrison, D. and Rubinfeld, D.L.\n\nThis is a copy of UCI ML housing dataset.\nhttps://archive.ics.uci.edu/ml/machine-learning-databases/housing/
\n\nThis dataset was taken from the StatLib library which is maintained at Carnegie Mellon University.\n\nThe Boston house-price data of Harrison, D. and Rubinf
eld, D.L. 'Hedonic\prices and the demand for clean air', J. Environ. Economics & Management, \nvol.5, 81-102, 1978. Used in Belsley, Kuh & Welsch, 'Regression d
iagnostics\n...', Wiley, 1980. N.B. Various transformations are used in the table on\npages 244-261 of the latter.\n\nThe Boston house-price data has been used
in many machine learning papers that address regression\nproblems. \n\n... topic:: References\n\n - Belsley, Kuh & Welsch, 'Regression diagnostics: Identi
fying Influential Data and Sources of Collinearity', Wiley, 1980. 244-261.\n - Quinlan,R. (1993). Combining Instance-Based and Model-Based Learning. In Proceedi
ngs on the Tenth International Conference of Machine Learning, 236-243, University of Massachusetts, Amherst. Morgan Kaufmann.\n"
```

```
In [31]: boston = pd.DataFrame(boston_dataset.data, columns=boston_dataset.feature_names)
boston.head()

Out[31]:
```

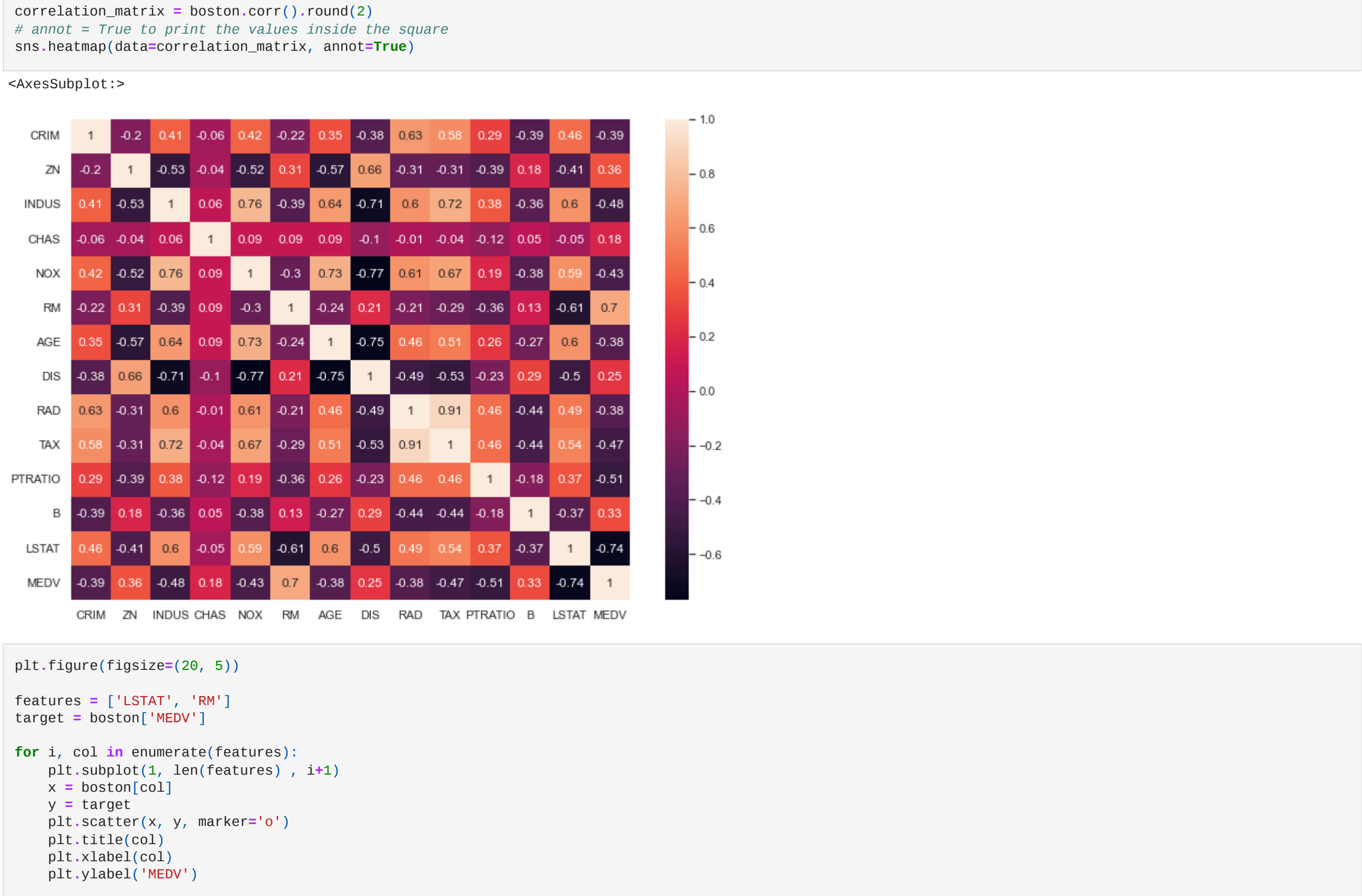
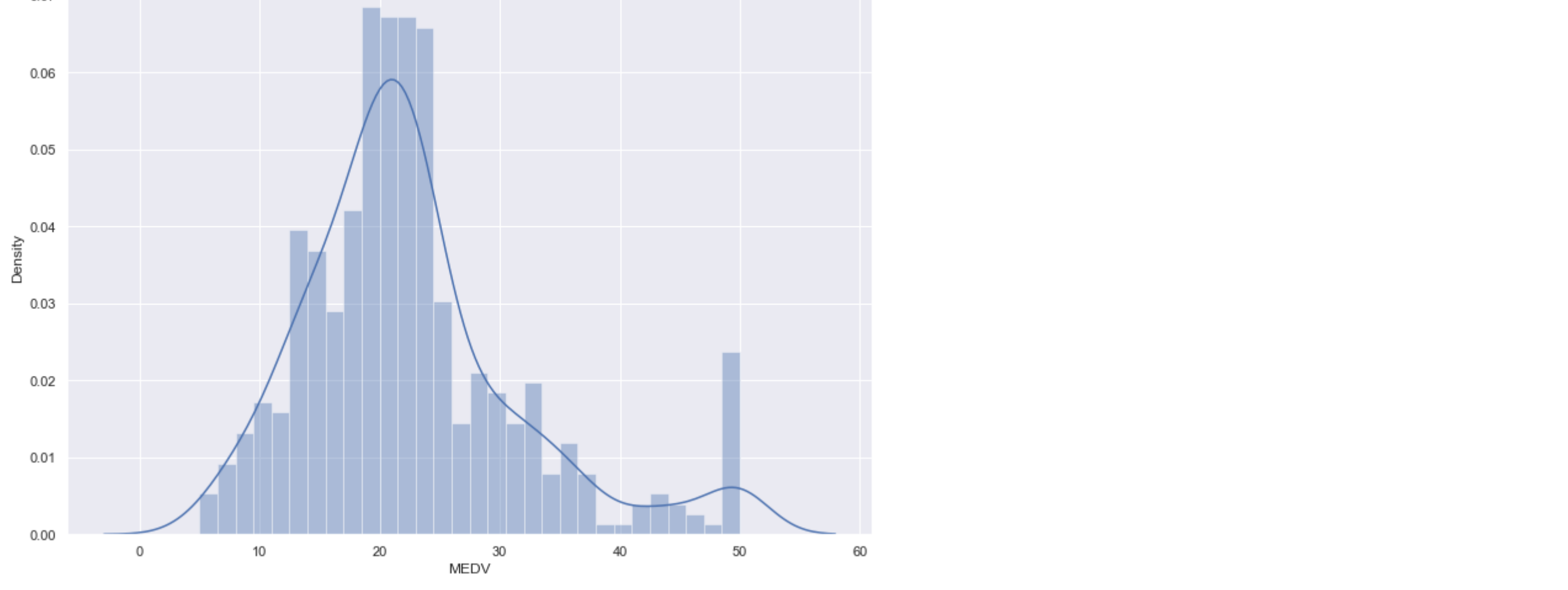
	CRIM	ZN	INDUS	CHAS	NOX	RM	AGE	DIS	RAD	TAX	PTRATIO	B	LSTAT
0	0.00632	18.0	2.31	0.0	0.538	6.575	65.2	4.0900	1.0	296.0	15.3	396.90	4.98
1	0.02731	0.0	7.07	0.0	0.469	6.421	78.9	4.9671	2.0	242.0	17.8	396.90	9.14
2	0.02729	0.0	7.07	0.0	0.469	7.185	61.1	4.9671	2.0	242.0	17.8	392.83	4.03
3	0.03237	0.0	2.18	0.0	0.458	6.998	45.8	6.0622	3.0	222.0	18.7	394.63	2.94
4	0.06905	0.0	2.18	0.0	0.458	7.147	54.2	6.0622	3.0	222.0	18.7	396.90	5.33

```
In [32]: boston['MEDV'] = boston_dataset.target

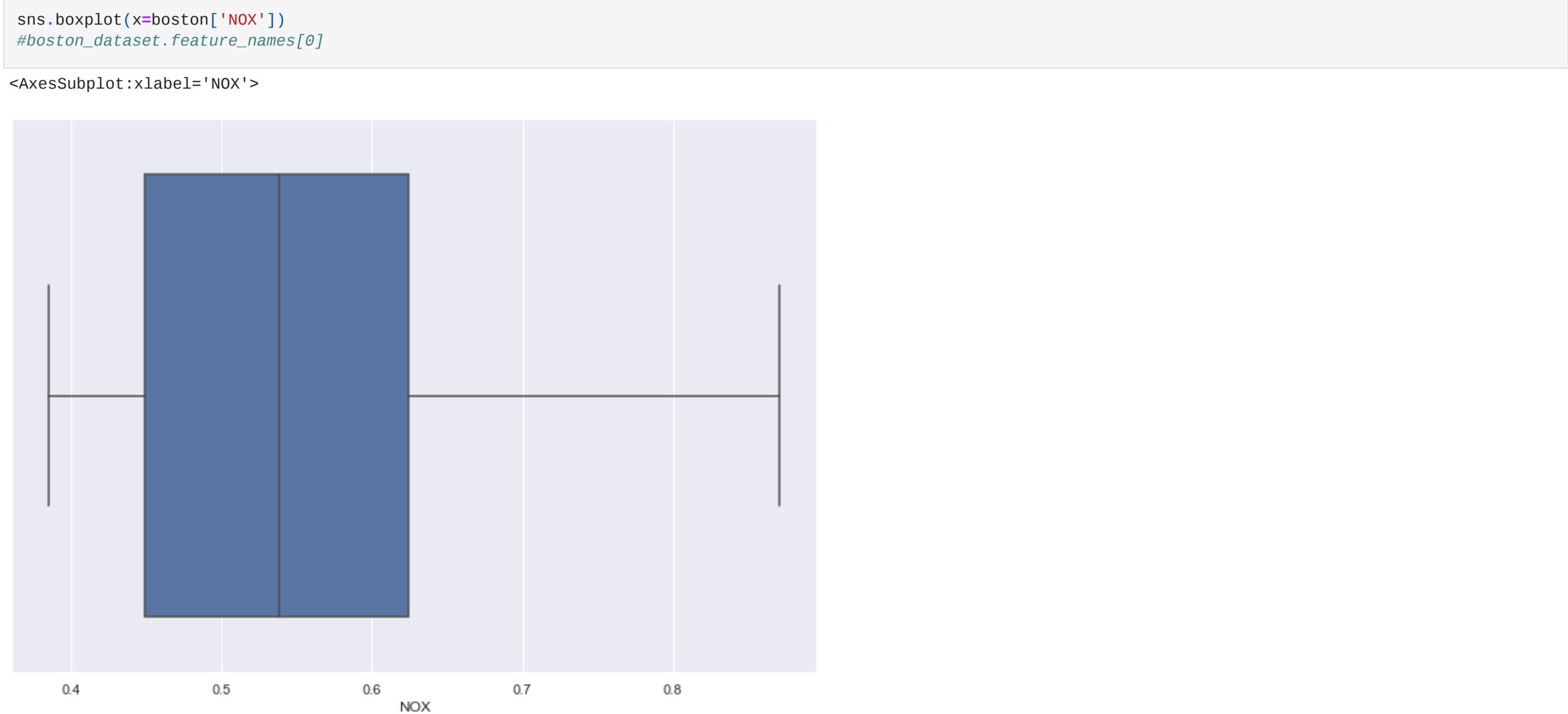
In [33]: boston.isnull().sum()

Out[33]:
```

CRIM	0
ZN	0
INDUS	0
CHAS	0
NOX	0
RM	0
AGE	0
DIS	0
RAD	0
TAX	0
PTRATIO	0
B	0
LSTAT	0
MEDV	0
dtype:	int64



Цены растут по мере линейного увеличения стоимости RM. Существуют выбросы. Данные ограничены 50. Цены имеют тенденцию к снижению с ростом LSTAT.



Чтобы соответствовать модели линейной регрессии, выбираются те особенности, которые имеют высокую корреляцию с целевой переменной MEDV. Глядя на матрицу корреляции видно, что RM имеет сильную положительную корреляцию с MEDV (0.7) где LSTAT имеет высокую отрицательную корреляцию с MEDV (-0.74). RAD,TAX имеют соотношение 0,91. Эти пары признаков сильно связаны друг с другом. Мы не должны выбирать обе эти функции вместе для обучения модели. То же самое касается признаков DIS а также AGE которые имеют корреляцию -0,75.