

Машинное обучение, ФКН ВШЭ

Практическое задание 2

Общая информация

Дата выдачи: 17.09.2018

Мягкий дедлайн: 2:59MSK 24.09.2018

Жесткий дедлайн: 23:59MSK 26.09.2018

О задании

Задание состоит из двух частей:

Первая часть задания посвящена получению студентами навыков работы с библиотекой векторного вычисления numpy и библиотеками для построения графиков matplotlib/plotly. Это задание сдается частично в Яндекс.Контест (см. информацию ниже). В случае проблем с доступом к Яндекс.Контесту обращайтесь к своему семинаристу(ке).

Вторая часть задания посвящена обучению линейной регрессии, подбору гиперпараметров и работе с данными.

Оценивание и штрафы

Каждая из задач имеет определенную «стоимость» (указана в скобках около задачи). Максимально допустимая оценка за работу — 10 баллов.

Сдавать задание после указанного срока сдачи нельзя. При выставлении неполного балла за задание в связи с наличием ошибок на усмотрение проверяющего предусмотрена возможность исправить работу на указанных в ответном письме условиях.

Задание выполняется самостоятельно. «Похожие» решения считаются плагиатом и все задействованные студенты (в том числе те, у кого списали) не могут получить за него больше 0 баллов (подробнее о плагиате см. на странице курса). Если вы нашли решение какого-то из заданий (или его часть) в открытом источнике, необходимо указать ссылку на этот источник в отдельном блоке в конце Вашей работы (скорее всего вы будете не единственным, кто это нашел, поэтому чтобы исключить подозрение в плагиате, необходима ссылка на источник).

Неэффективная реализация кода может негативно отразиться на оценке.

Формат сдачи

Ipython Notebook с выполненным заданием необходимо сдать в Anytask (lab_02). Также обратите внимание, что в первой части про Numpy также необходимо сделать две посылки в Яндекс.Контест.

Часть 1. Numpy

Задачи 1-6

(2.4 балла)

Ниже приведены задачи на работу с numpy-массивами. Для каждой из задач нужно привести 2 реализации: одна без использования numpy (считайте, что там, где на входе или выходе должны быть numpy array, будут просто списки), а вторая полностью векторизованная (без использования питоновских циклов/map/list comprehension). Невекторизованная реализация каждой из задач оценивается в **0.15 балла**, векторизованная – в **0.25 балла**.

Реализации без использования векторизации нужно записать в файл functions.py, а векторизованные — в файл functions_vectorized.py (см. шаблоны). Далее эти файлы необходимо сдать в Яндекс.Контест:

<https://official.contest.yandex.ru/contest/9148/enter/>

(<https://official.contest.yandex.ru/contest/9148/enter/>) в соответствующие задачи. По техническим причинам тестирование проводится на этапе компиляции, поэтому в случае любой ошибки вы будете получать вердикт СЕ, и в логе компиляции можно будет посмотреть, в чем проблема. Частичное выполнение задания (не все задачи) будет оцениваться, хотя и будет получать вердикт СЕ. Для удобства проверки приложите в ячейке ниже ссылки на самые успешные посылки.

- **Задача 1:** Подсчитать произведение ненулевых элементов на диагонали прямоугольной матрицы.
Например, для `X = np.array([[1, 0, 1], [2, 0, 2], [3, 0, 3], [4, 4, 4]])` ответ – 3.
- **Задача 2:** Даны два вектора `x` и `y`. Проверить, задают ли они одно и то же мультимножество.
Например, для `x = np.array([1, 2, 2, 4])`, `y = np.array([4, 2, 1, 2])` ответ – True.
- **Задача 3:** Найти максимальный элемент в векторе `x` среди элементов, перед которыми стоит нулевой.
Например, для `x = np.array([6, 2, 0, 3, 0, 0, 5, 7, 0])` ответ – 5.
- **__ Задача 4 __:** Дан трёхмерный массив, содержащий изображение, размера (height, width, numChannels), а также вектор длины numChannels. Сложить каналы изображения с указанными весами, и вернуть результат в виде матрицы размера (height, width). В ноутбуке приведите пример работы функции – преобразуйте цветное изображение в оттенки серого, используя коэффициенты `np.array([0.299, 0.587, 0.114])`. Считать реальное изображение можно при помощи функции `scipy.misc.imread` (если изображение не в формате png, установите пакет pillow). **Обратите внимание, что в изображении может быть не три канала! За решения, явно или неявно использующие трёхканальность изображения будет присуждена половина баллов.**
- **Задача 5:** Реализовать кодирование длин серий (Run-length encoding). Для некоторого вектора `x` необходимо вернуть кортеж из двух векторов одинаковой длины. Первый содержит числа, а второй - сколько раз их нужно повторить.
Например, для `x = np.array([2, 2, 2, 3, 3, 3, 5])` ответ – `(np.array([2, 3, 5]), np.array([3, 3, 1]))`.
- **Задача 6:** Даны две выборки объектов - `X` и `Y`. Вычислить матрицу евклидовых

расстояний между объектами. Сравните с функцией `scipy.spatial.distance.cdist` по скорости работы (сравнения приведите ниже в ноутбуке).

Замечание. Можно считать, что все указанные объекты непустые (к примеру, в **задаче 1** на диагонали матрицы есть ненулевые элементы) и корректные.

Укажите ссылки на посылки

Посылка по не векторизованным функциям: <https://contest.yandex.ru/contest/9148/run-report/12055442/> (<https://contest.yandex.ru/contest/9148/run-report/12055442/>).

Посылка по векторизованным функциям: <https://contest.yandex.ru/contest/9148/run-report/12055411/> (<https://contest.yandex.ru/contest/9148/run-report/12055411/>).

Задача 7

(1.6 балла)

Для каждой задачи сравните скорость работы не векторизованной и векторизованной реализации. С помощью пакета `matplotlib` или `plotly` постройте графики времени работы в зависимости от размера данных. **Графики должны выглядеть опрятно!** То есть должны быть подписаны оси, названия графиков, и т.д. Например, ниже представлены хороший и плохой графики:

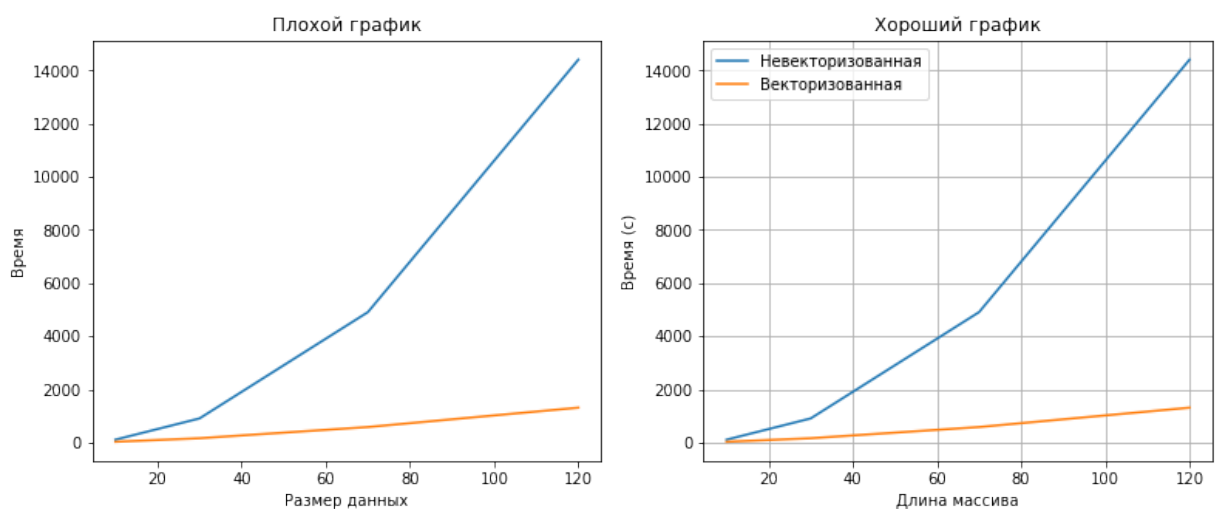
```

In [1]: 1 import matplotlib.pyplot as plt
        2 import numpy as np
        3
        4 %matplotlib inline
        5
        6 data_size = np.array([10, 30, 70, 120])
        7 time_non_vectorized = data_size ** 2 + 10
        8 time_vectorized = data_size ** 1.5
        9
       10 f, (ax1, ax2) = plt.subplots(1, 2, figsize=(13, 5))
       11
       12 ax1.plot(data_size, time_non_vectorized)
       13 ax1.plot(data_size, time_vectorized)
       14 ax1.set_title(u"Плохой график")
       15 ax1.set_xlabel(u"Размер данных")
       16 ax1.set_ylabel(u"Время")
       17
       18 ax2.plot(data_size, time_non_vectorized, label=u"Невекторизованная")
       19 ax2.plot(data_size, time_vectorized, label=u"Векторизованная")
       20 ax2.set_title(u"Хороший график")
       21 ax2.set_xlabel(u"Длина массива")
       22 ax2.set_ylabel(u"Время (с)")
       23 ax2.grid()
       24 ax2.legend()
       25
       26 f.show()

```

/usr/local/lib/python3.7/site-packages/matplotlib/figure.py:448: UserWarning: Matplotlib is currently using module://ipykernel.pylab.backend_inline, which is a non-GUI backend, so cannot show the figure.

```
% get_backend())
```



```
In [2]: 1 def build_execution_time_graph(time_non_vectorized, time_vectorized,
2      f, (ax2) = plt.subplots(1, 1, figsize=(13, 5))
3
4      ax2.plot(data_sizes, time_non_vectorized, label=u"non_vectorized")
5      ax2.plot(data_sizes, time_vectorized, label=u"vectorized")
6      ax2.set_title(function_name + " execution time")
7      ax2.set_xlabel(u"data size")
8      ax2.set_ylabel(u"execution time, seconds")
9      ax2.grid()
10     ax2.legend()
11
12     f.show()
```

```
In [3]: 1 import timeit
2
3 def calculate_execution_time(function, parameters):
4     start_timer = timeit.default_timer()
5     function(*parameters)
6     return (timeit.default_timer() - start_timer)
```

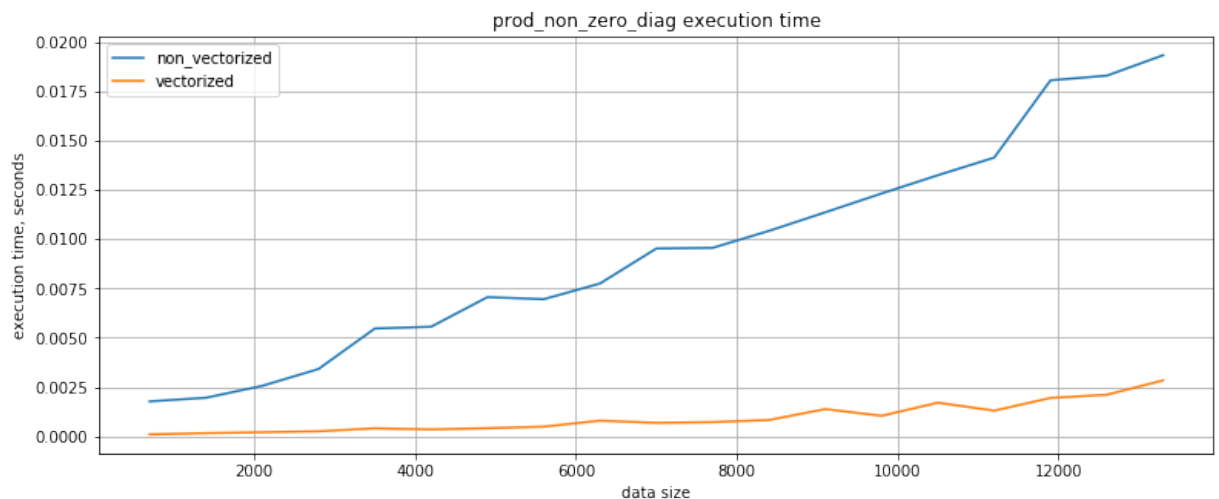
```
In [4]: 1 import functions as funct
2 import functions_vectorized as vfunct
```

```

In [7]: 1 def compare_prod_non_zero_diag():
2         data_sizes = [700 * step for step in range(1, 20)]
3         time_non_vectorized = [
4             calculate_execution_time(func.prod_non_zero_diag, [np.random.randn(
5                 data_size, data_size)
6                 for data_size in data_sizes
7             ])
8         ]
9         time_vectorized = [
10             calculate_execution_time(vfunc.prod_non_zero_diag, [np.random.randn(
11                 data_size, data_size)
12                 for data_size in data_sizes
13             ])
14         ]
15         build_execution_time_graph(
16             time_non_vectorized,
17             time_vectorized,
18             data_sizes,
19             u"prod_non_zero_diag"
20         )
21     compare_prod_non_zero_diag()

```

/usr/local/lib/python3.7/site-packages/matplotlib/figure.py:448: UserWarning: Matplotlib is currently using module://ipykernel.pylab.backend_inline, which is a non-GUI backend, so cannot show the figure.
% get_backend())



```

In [8]: 1 def compare_are_multisets_equal():
2       data_sizes = [100000 * step for step in range(1, 50)]
3       time_non_vectorized = [
4           calculate_execution_time(
5               funct.are_multisets_equal, [np.random.rand(data_size),
6               ])
7       for data_size in data_sizes
8   ]
9       time_vectorized = [
10          calculate_execution_time(
11              vfunct.are_multisets_equal, [np.random.rand(data_size)
12          ])
13      for data_size in data_sizes
14  ]
15      build_execution_time_graph(
16          time_non_vectorized,
17          time_vectorized,
18          data_sizes,
19          u"are_multisets_equal"
20      )
21      compare_are_multisets_equal()

```

/usr/local/lib/python3.7/site-packages/matplotlib/figure.py:448: Use
rWarning: Matplotlib is currently using module://ipykernel.pylab.bac
kend_inline, which is a non-GUI backend, so cannot show the figure.
% get_backend())




```

In [9]: 1 def compare_max_after_zero():
2         data_sizes = [100000 * step for step in range(1, 30)]
3         time_non_vectorized = [
4             calculate_execution_time(
5                 funct.max_after_zero, [np.random.randint(10, size = da
6             )
7         for data_size in data_sizes
8     ]
9     time_vectorized = [
10         calculate_execution_time(
11             vfunct.max_after_zero, [np.random.randint(10, size = c
12         )
13     for data_size in data_sizes
14 ]
15 build_execution_time_graph(
16     time_non_vectorized,
17     time_vectorized,
18     data_sizes,
19     u"max_after_zero"
20 )
21 compare_max_after_zero()

```

/usr/local/lib/python3.7/site-packages/matplotlib/figure.py:448: Use
rWarning: Matplotlib is currently using module://ipykernel.pylab.bac
kend_inline, which is a non-GUI backend, so cannot show the figure.
% get_backend())

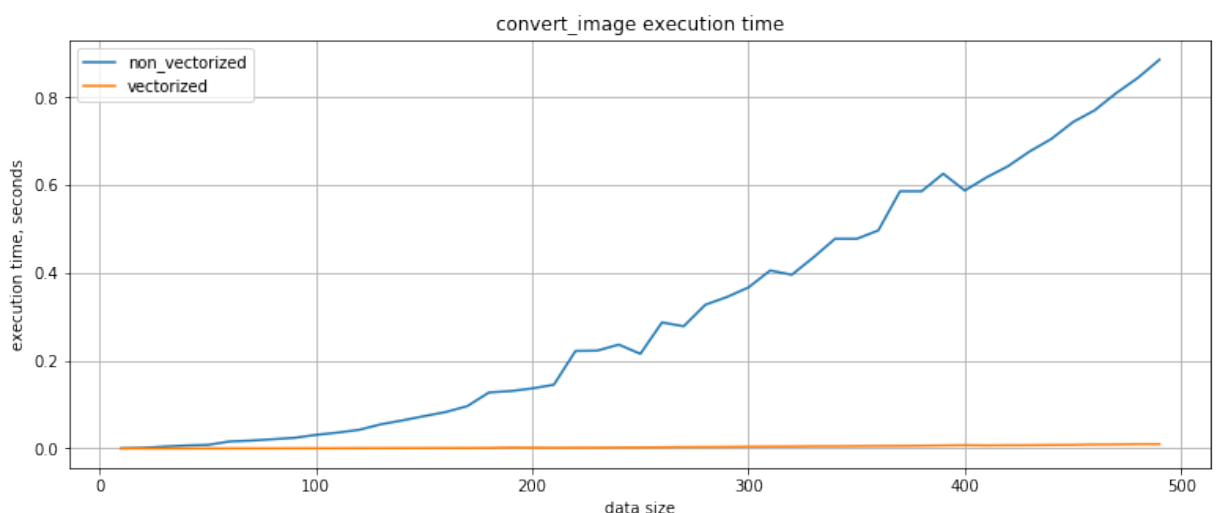


```

In [10]: 1 def compare_convert_image():
2         data_sizes = [10 * step for step in range(1, 50)]
3         data_size = 3
4         import math
5         time_non_vectorized = [
6             calculate_execution_time(
7                 funct.convert_image,
8                 [
9                     np.random.rand(data_size, data_size, int(math.log(
10                      np.random.rand(int(math.log(data_size, 2)))
11                      ]
12                  )
13         for data_size in data_sizes
14     ]
15     time_vectorized = [
16         calculate_execution_time(
17             vfunct.convert_image,
18             [
19                 np.random.rand(data_size, data_size, int(math.log(
20                 np.random.rand(int(math.log(data_size, 2)))
21                 ]
22             )
23         for data_size in data_sizes
24     ]
25     build_execution_time_graph(
26         time_non_vectorized,
27         time_vectorized,
28         data_sizes,
29         u"convert_image"
30     )
31     compare_convert_image()

```

/usr/local/lib/python3.7/site-packages/matplotlib/figure.py:448: Use
rWarning: Matplotlib is currently using module://ipykernel.pylab.bac
kend_inline, which is a non-GUI backend, so cannot show the figure.
% get_backend())

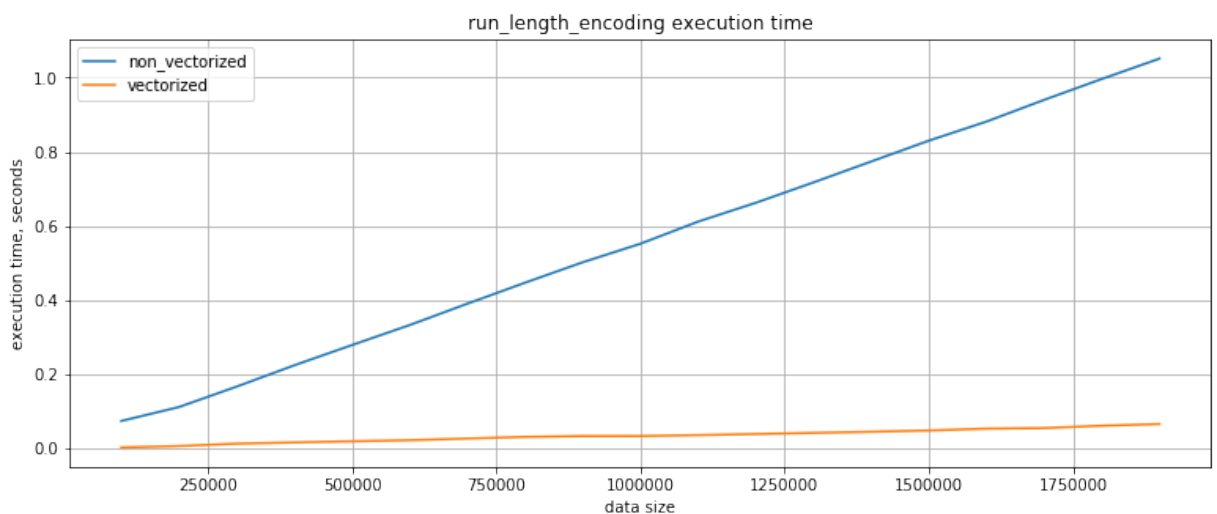


```

In [11]: 1 def compare_run_length_encoding():
2         data_sizes = [100000 * step for step in range(1, 20)]
3         time_non_vectorized = [
4             calculate_execution_time(
5                 funct.run_length_encoding, [np.random.randint(10, size
6             )
7             for data_size in data_sizes
8         ]
9         time_vectorized = [
10            calculate_execution_time(
11                vfunct.run_length_encoding, [np.random.randint(10, siz
12            )
13            for data_size in data_sizes
14        ]
15        build_execution_time_graph(
16            time_non_vectorized,
17            time_vectorized,
18            data_sizes,
19            u"run_length_encoding"
20        )
21        compare_run_length_encoding()

```

/usr/local/lib/python3.7/site-packages/matplotlib/figure.py:448: Use
rWarning: Matplotlib is currently using module://ipykernel.pylab.bac
kend_inline, which is a non-GUI backend, so cannot show the figure.
% get_backend())



In [13]:

```

1  import math
2  def compare_pairwise_distance():
3      data_sizes = [10 * step for step in range(1, 50)]
4      time_non_vectorized = [
5          calculate_execution_time(
6              funct.pairwise_distance,
7              [
8                  np.random.rand(data_size, int(math.log(data_size,
9                  np.random.rand(data_size, int(math.log(data_size,
10             ]
11         )
12     for data_size in data_sizes
13 ]
14 time_vectorized = [
15     calculate_execution_time(
16         vfunct.pairwise_distance,
17         [
18             np.random.rand(data_size, int(math.log(data_size,
19             np.random.rand(data_size, int(math.log(data_size,
20         ]
21     )
22     for data_size in data_sizes
23 ]
24 build_execution_time_graph(
25     time_non_vectorized,
26     time_vectorized,
27     data_sizes,
28     u"pairwise_distance"
29 )
30 compare_pairwise_distance()

```

/usr/local/lib/python3.7/site-packages/matplotlib/figure.py:448: Use rWarning: Matplotlib is currently using module://ipykernel.pylab.backend_inline, which is a non-GUI backend, so cannot show the figure.
% get_backend())



Часть 2. Линейная регрессия

```
In [14]: 1 %pylab inline
          2
          3 import seaborn as sns
          4 sns.set(style="whitegrid")
          5
          6 from sklearn import datasets
          7 from sklearn.model_selection import train_test_split, GridSearchCV
          8 from sklearn.pipeline import Pipeline
          9 from sklearn.linear_model import LinearRegression, Ridge, Lasso
         10 from sklearn.preprocessing import StandardScaler, MinMaxScaler
         11 import pandas as pd
```

Populating the interactive namespace from numpy and matplotlib

```
/usr/local/lib/python3.7/site-packages/IPython/core/magics/pylab.py:
160: UserWarning: pylab import has clobbered these variables: ['f']
`%matplotlib` prevents importing * from pylab and numpy
"\n`%matplotlib` prevents importing * from pylab and numpy"
```

```
In [15]: 1 def pairplot(df, target):
          2     ncol, nrow = 7, df.shape[1] // 7 + (df.shape[1] % 7 > 0)
          3     plt.figure(figsize=(ncol * 4, nrow * 4))
          4
          5     for i, feature in enumerate(df.columns):
          6         plt.subplot(nrow, ncol, i + 1)
          7         plt.scatter(df[feature], target, s=10, marker='o', alpha=.5)
          8         plt.xlabel(feature)
          9         if i % ncol == 0:
         10             plt.ylabel('target')
```

В данном задании мы рассмотрим стандартный датасет для задачи регрессии Boston Housing, в котором необходимо предсказать стоимость недвижимости по 13 признакам.

Датасет достаточно просто загрузить из библиотеки scikit-learn:

```
In [16]: 1 data = datasets.load_boston()
          2 attributes = pd.DataFrame(data.data)
          3 attributes.columns = data.feature_names
          4 target = data.target
```

```
In [17]: 1 print(data['DESCR'])
```

Boston House Prices dataset

=====

Notes

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.

<http://archive.ics.uci.edu/ml/datasets/Housing>
(<http://archive.ics.uci.edu/ml/datasets/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 ta

ble on
pages 244–261 of the latter.

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

****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, 236–243, University of Massachusetts, Amherst. Morgan Kaufmann.
- many more! (see <http://archive.ics.uci.edu/ml/datasets/Housing>) (<http://archive.ics.uci.edu/ml/datasets/Housing>)

In [18]: 1 attributes.head()

Out[18]:

	CRIM	ZN	INDUS	CHAS	NOX	RM	AGE	DIS	RAD	TAX	PTRATIO	B	L
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	
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	
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	
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	
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	

In [19]: 1 print(target[:4])

[24. 21.6 34.7 33.4]

Разделим выборку на обучающую и тестовую в отношении 8/2:

In [20]: 1 attributes_train, attributes_test, target_train, target_test = tra
/usr/local/lib/python3.7/site-packages/sklearn/model_selection/_split.py:2026: FutureWarning: From version 0.21, test_size will always complement train_size unless both are specified.
FutureWarning)

Задание 1. Обучение линейной регрессии.**(0.5 балл)**

Обучите стандартную линейную регрессию, а также с L_1 и L_2 регуляризаторами (используйте параметры по умолчанию). Посчитайте метрику R^2 для каждого метода (метод score).

```
In [21]: 1 import sklearn.linear_model as lm
2 import sklearn.preprocessing as pr
3 def calculate_standart_regression():
4     regression = lm.LinearRegression()
5     regression.fit(attributes_train, target_train)
6     print("standart regulization:")
7     print("score on train sample:", regression.score(attributes_train, target_train))
8     print("score on test sample:", regression.score(attributes_test, target_test))
9     calculate_standart_regression()
```

```
standart regulization:
score on train sample: 0.7508837786732915
score on test sample: 0.6684825753971597
```

```
/usr/local/lib/python3.7/site-packages/sklearn/linear_model/base.py:
509: RuntimeWarning: internal gelsd driver lwork query error, required iwork dimension not returned. This is likely the result of LAPACK bug 0038, fixed in LAPACK 3.2.2 (released July 21, 2010). Falling back to 'gelss' driver.
      linalg.lstsq(X, y)
```

```
In [22]: 1 def calculate_l1_regression():
2     regression = lm.Lasso()
3     regression.fit(attributes_train, target_train)
4     print("l1 regulization:")
5     print("score on train sample:", regression.score(attributes_train, target_train))
6     print("score on test sample:", regression.score(attributes_test, target_test))
7     calculate_l1_regression()
```

```
l1 regulization:
score on train sample: 0.6958999474620655
score on test sample: 0.6668687223368213
```



```
In [23]: 1 def calculate_l2_regression():
2         regression = lm.Ridge()
3         regression.fit(attributes_train, target_train)
4         print("l2 regulization:")
5         print("score on train sample:", regression.score(attributes_train, target_train))
6         print("score on test sample:", regression.score(attributes_test, target_test))
7 calculate_l2_regression()
```

```
l2 regulization:
score on train sample: 0.7487713711418569
score on test sample: 0.6659608075261689
```

Задание 2. Подбор гиперпараметров.

(1 балл)

Для Lasso- и Ridge-регресий подберите коэффициент регуляризации по обучающей выборке с помощью кросс-валидации. Параметры для перебора возьмите по логарифмической сетке от 10^{-6} до 10^6). Также посчитайте метрику R^2 для тестовой выборки и сравните с предыдущими результатами. Заметно ли изменилось качество?

Useful: GridSearchCV, RidgeCV, LassoCV

```
In [24]: 1 def calculate_ridge_regularization_coefficient():
2         regression = lm.RidgeCV(
3             alphas=np.log2(np.arange(1e-6, 1e6, 100)),
4             store_cv_values=True,
5         )
6         regression.fit(attributes_train, target_train)
7         print("RidgeCV regression:")
8         print("regularization coefficient:", regression.alpha_)
9         print("score on train sample:", regression.score(attributes_train, target_train))
10        print("score on test sample:", regression.score(attributes_test, target_test))
11 calculate_ridge_regularization_coefficient()
```

```
RidgeCV regression:
regularization coefficient: 6.643856204201675
score on train sample: 0.7429713402782756
score on test sample: 0.6627401932541358
```

```
In [25]: 1 def calculate_lasso_regularization_coefficient():
2         regression = lm.LassoCV(
3             alphas=np.log2(np.arange(1e-6, 1e6, 100))
4         )
5         regression.fit(attributes_train, target_train)
6         print("LassoCV regression:")
7         print("regularization coefficient:", regression.alpha_)
8         print("score on train sample:", regression.score(attributes_train, target_train))
9         print("score on test sample:", regression.score(attributes_test, target_test))
10        calculate_lasso_regularization_coefficient()
```

```
LassoCV regression:
regularization coefficient: 6.643856204201675
score on train sample: 0.5510142723812895
score on test sample: 0.549865805812861
```

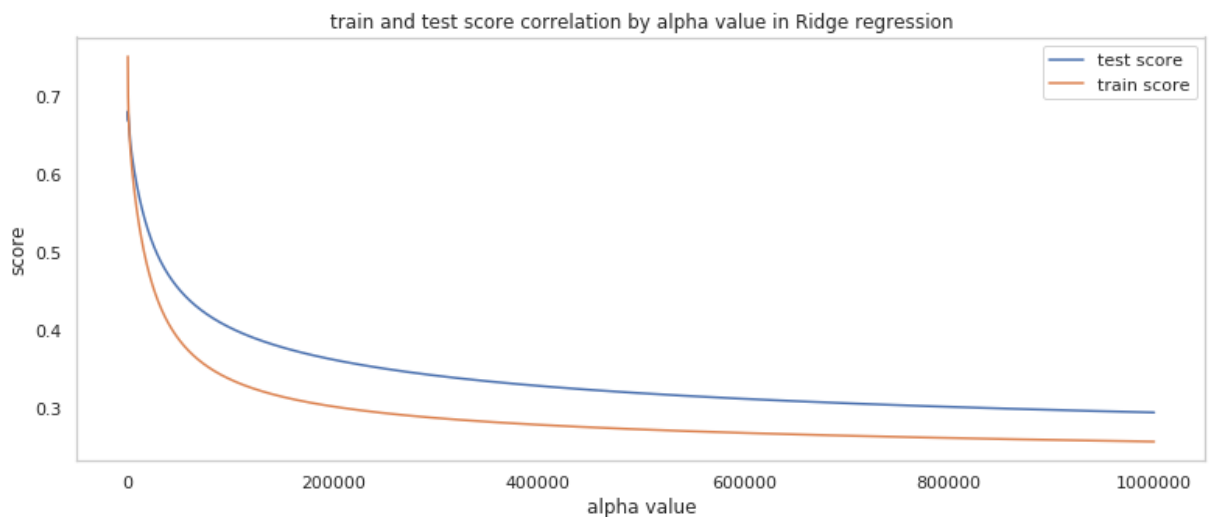
Постройте график зависимости R^2 для обучающей и тестовой (на кросс-валидации) выборок в зависимости от значения гиперпараметра. Для Lasso также постройте график зависимости количества ненулевых весов.

```

In [26]: 1 def draw_ridge_alpha_correlation():
2         scores_train = []
3         scores_test = []
4         alpha_values = np.arange(1e-6, 1e6, 100)
5         for alpha_value in alpha_values:
6             regression = lm.Ridge(alpha=alpha_value)
7             regression.fit(attributes_train, target_train)
8             scores_train.append(regression.score(attributes_train, target_train))
9             scores_test.append(regression.score(attributes_test, target_test))
10
11         f, (ax2) = plt.subplots(1, 1, figsize=(13, 5))
12
13         ax2.plot(alpha_values, scores_test, label=u"test score")
14         ax2.plot(alpha_values, scores_train, label=u"train score")
15         ax2.set_title(u"train and test score correlation by alpha value")
16         ax2.set_xlabel(u"alpha value")
17         ax2.set_ylabel(u"score")
18         ax2.grid()
19         ax2.legend()
20
21         f.show()
22         draw_ridge_alpha_correlation()

```

/usr/local/lib/python3.7/site-packages/matplotlib/figure.py:448: UserWarning: Matplotlib is currently using module://ipykernel.pylab.backend_inline, which is a non-GUI backend, so cannot show the figure.
% get_backend())

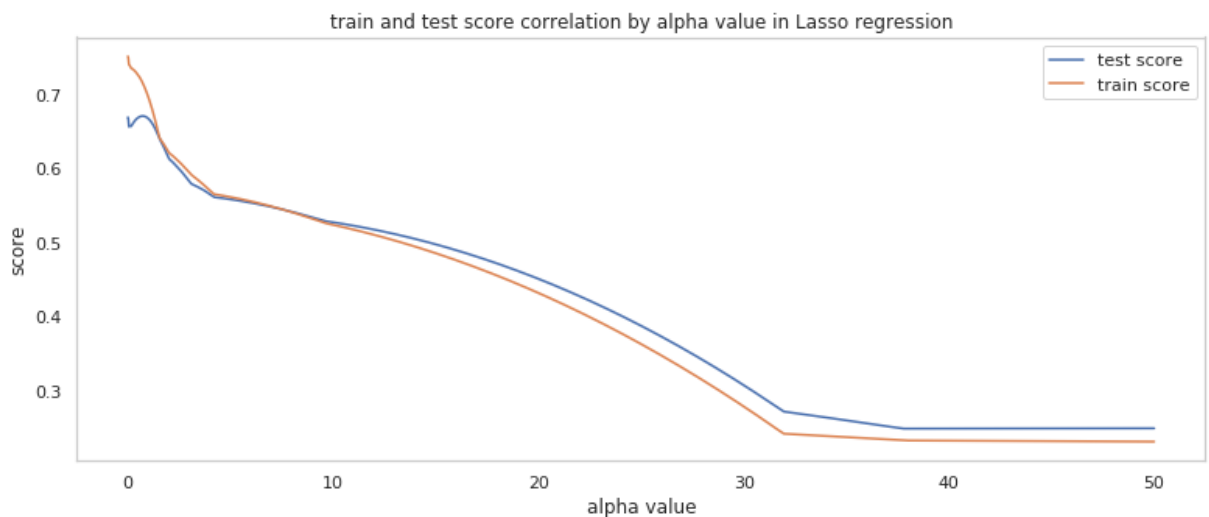


```

In [27]: 1 def draw_lasso_alpha_correlation():
2         scores_train = []
3         scores_test = []
4         alpha_values = np.arange(1e-6, 50, 0.005)
5         for alpha_value in alpha_values:
6             regression = lm.Lasso(alpha=alpha_value)
7             regression.fit(attributes_train, target_train)
8             scores_train.append(regression.score(attributes_train, target_train))
9             scores_test.append(regression.score(attributes_test, target_test))
10
11         f, (ax2) = plt.subplots(1, 1, figsize=(13, 5))
12
13         ax2.plot(alpha_values, scores_test, label=u"test score")
14         ax2.plot(alpha_values, scores_train, label=u"train score")
15         ax2.set_title(u"train and test score correlation by alpha value")
16         ax2.set_xlabel(u"alpha value")
17         ax2.set_ylabel(u"score")
18         ax2.grid()
19         ax2.legend()
20
21         f.show()
22 draw_lasso_alpha_correlation()

```

/usr/local/lib/python3.7/site-packages/matplotlib/figure.py:448: UserWarning: Matplotlib is currently using module://ipykernel.pylab.backend_inline, which is a non-GUI backend, so cannot show the figure.
% get_backend())

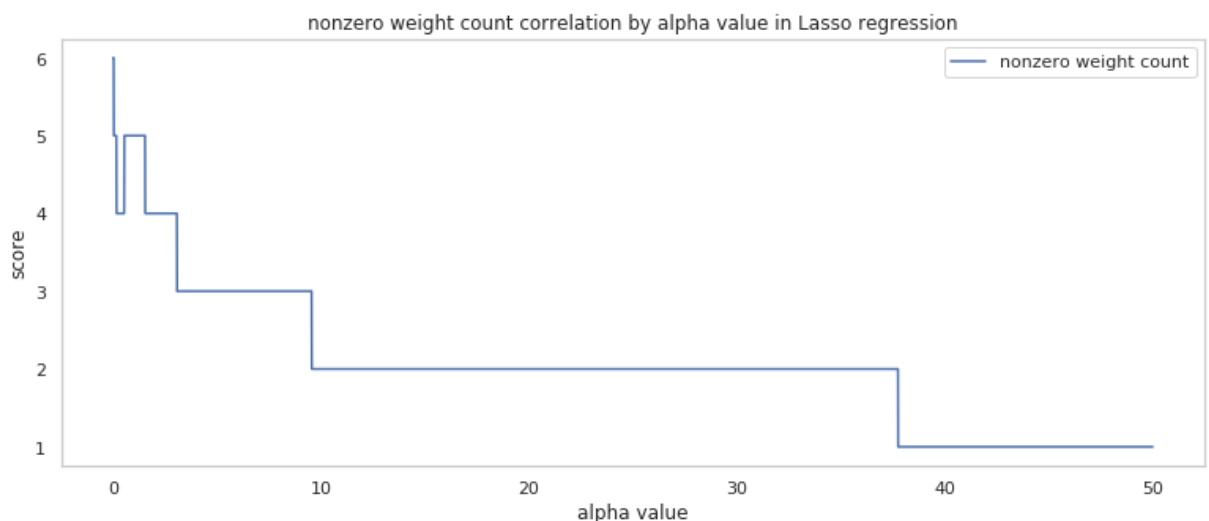


```

In [28]: 1 def draw_lasso_nonzero_weight_count_correlation():
2         nonzero_weight_count = []
3         alpha_values = np.arange(1e-6, 50, 0.005)
4         for alpha_value in alpha_values:
5             regression = lm.Lasso(alpha=alpha_value)
6             regression.fit(attributes_train, target_train)
7             nonzero_weight_count.append(sum(regression.coef_ > 0))
8
9         f, (ax2) = plt.subplots(1, 1, figsize=(13, 5))
10
11         ax2.plot(alpha_values, nonzero_weight_count, label=u"nonzero w
12         ax2.set_title(u"nonzero weight count correlation by alpha valu
13         ax2.set_xlabel(u"alpha value")
14         ax2.set_ylabel(u"score")
15         ax2.grid()
16         ax2.legend()
17
18         f.show()
19         draw_lasso_nonzero_weight_count_correlation()

```

/usr/local/lib/python3.7/site-packages/matplotlib/figure.py:448: Use
rWarning: Matplotlib is currently using module://ipykernel.pylab.bac
kend_inline, which is a non-GUI backend, so cannot show the figure.
% get_backend())



В какой момент происходит недообучение? Почему?

Ответ: По графику видно, что переобучение Ridge происходит при коэффициенте 0.7. Lasso переобучается при коэффициенте 6 -- 10, видимо из-за изменения количества ненулевых векторов (какой-то значимый вектор перестает влиять на результат).

Подготовка данных

Как вы могли заметить, большого прироста качества с помощью подбора коэффициента регуляризации добиться не удалось. Поэтому прежде чем бросаться обучать модели, необходимо изучить и предобработать данные.

Задание 3. Масштабирование.

(1 балл)

Зачастую признаки в сырых данных имеют разный масштаб. Попробуйте применить масштабирование к данным, сравните качество. Заметно ли изменение?

Pipeline, StandardScaler, MinMaxScaler

```
In [29]: 1 def calculate_regressions_with_standart_scaler():
2         scaler = pr.StandardScaler()
3         scaled_attributes_train = scaler.fit_transform(attributes_train)
4         scaled_attributes_test = scaler.fit_transform(attributes_test)
5
6         linear_regression = lm.LinearRegression()
7         linear_regression.fit(scaled_attributes_train, target_train)
8         ridge_regression = lm.Ridge()
9         ridge_regression.fit(scaled_attributes_train, target_train)
10        lasso_regression = lm.Lasso()
11        lasso_regression.fit(scaled_attributes_train, target_train)
12        print("standart scaler:")
13        print("scaled simple linear regression score:", linear_regression.score)
14        print("scaled ridge regression score:", ridge_regression.score)
15        print("scaled lasso regression score:", lasso_regression.score)
16
17 calculate_regressions_with_standart_scaler()
```

```
standart scaler:
scaled simple linear regression score: 0.6260120845037371
scaled ridge regression score: 0.6260848349266783
scaled lasso regression score: 0.6068805552998523
```

```
In [30]: 1 def calculate_regressions_with_minmax_scaler():
2         scaler = pr.MinMaxScaler()
3         scaled_attributes_train = scaler.fit_transform(attributes_train)
4         scaled_attributes_test = scaler.fit_transform(attributes_test)
5
6         linear_regression = lm.LinearRegression()
7         linear_regression.fit(scaled_attributes_train, target_train)
8         ridge_regression = lm.Ridge()
9         ridge_regression.fit(scaled_attributes_train, target_train)
10        lasso_regression = lm.Lasso()
11        lasso_regression.fit(scaled_attributes_train, target_train)
12        print("minmax scaler:")
13        print("scaled simple linear regression score:", linear_regression.score)
14        print("scaled ridge regression score:", ridge_regression.score)
15        print("scaled lasso regression score:", lasso_regression.score)
16
17 calculate_regressions_with_minmax_scaler()
```

```
minmax scaler:
scaled simple linear regression score: 0.569764132570967
scaled ridge regression score: 0.6013620981305376
scaled lasso regression score: 0.26666070646435747
```

Сравните оптимальное значение коэффициента регуляризации для Ridge-регрессии до и после масштабирования. Изменилось ли оно? Предположите, почему так могло произойти.

```
In [31]: 1 def calculate_regression_coeficients_after_standart_scaling():
2         scaler = pr.MinMaxScaler()
3         scaled_attributes_train = scaler.fit_transform(attributes_train)
4
5         ridge_regression = lm.RidgeCV(
6             alphas=np.log2(np.arange(1e-6, 1e6, 100)),
7             store_cv_values=True,
8         )
9         ridge_regression.fit(scaled_attributes_train, target_train)
10
11        print("new ridge coefficient:", ridge_regression.alpha_)
12
13 calculate_regression_coeficients_after_standart_scaling()
```

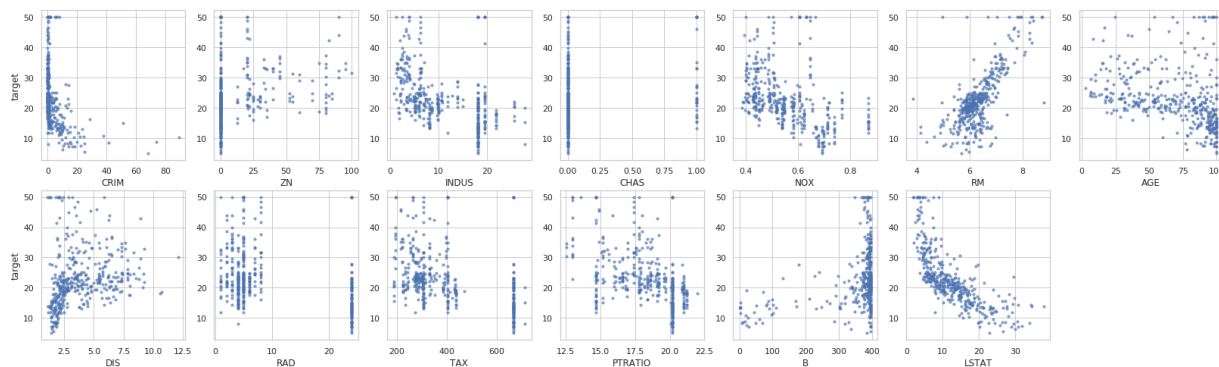
```
new ridge coefficient: 6.643856204201675
```

Ответ: Значение коэффициента не изменилось, поскольку MinMax масштабирование нормирует данные относительно целевого вектора.

Задание 4. Новые признаки.**(1.5 балла)**

Полезным также бывает посмотреть как целевая переменная зависит от каждого признака.

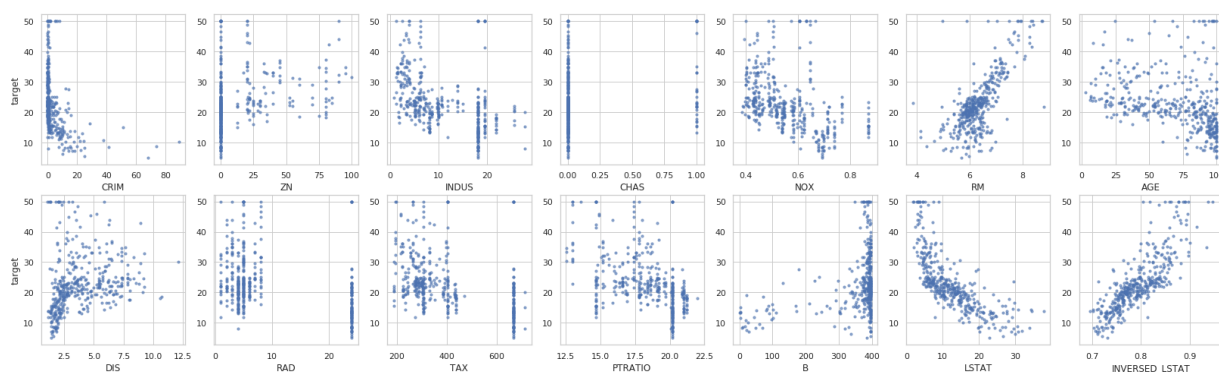
```
In [32]: 1 pairplot(attributes_train, target_train)
```



Обратите внимание на зависимость цены от признака LSTAT. Является ли эта зависимость линейной? А какой? Попробуйте выбрать преобразование для этого признака так, чтобы получившаяся зависимость была более линейной. Добейтесь R^2 на тестовой выборке не меньше 0.71.

Ответ: Прослеживается гиперболическая зависимость. Для линейности можно попробовать взять $1/\sqrt[10]{LSTAT}$.

```
In [43]: 1 extended_attributes_train = attributes_train
2         attributes_train = extended_attributes_train.drop(columns = ["INVE
3         pairplot(extended_attributes_train, target_train)
4
```



Помимо преобразований отдельных признаков полезными бывают их попарные взаимодействия. Воспользуйтесь PolynomialFeatures, чтобы добавить попарные произведения и квадраты всех признаков. Обучите Ridge-регрессию (подберите гиперпараметр!) и посчитайте R^2 на тесте. Сильно ли изменилось качество?


```
In [47]: 1 polynomial_features = pr.PolynomialFeatures(2)
2 ridge_regression = lm.RidgeCV(alphas=np.log2(np.arange(1e-6, 1e6,
3 ridge_regression.fit(polynomial_features.fit_transform(attributes_
4 print(ridge_regression.score(polynomial_features.fit_transform(att

0.799489615663477
```

Да, при попарном взаимодействии качество увеличивается на порядок.

Задание 5. Оптимальная архитектура, анализ.

(2 балла)

При помощи Pipeline и GridSearchSCV выберите оптимальную архитектуру, комбинируя различные методы масштабирования, степень полинома в PolynomialFeatures , а также регуляризаторы. Для Lasso поставьте максимальное количество итераций больше значения по умолчанию, чтобы оптимизация сошлась.

```
In [ ]: 1 # Your code here
2 # ...
```

Для одной комбинаций метода масштабирования и линейной регрессии постройте зависимость метрики R^2 на кросс-валидации (GridSearchCV.cv_results_['mean_test_score']) от значения параметра регуляризации для различных степеней полиномов в PolynomialFeatures .

```
In [ ]: 1 # Your code here
2 # ...
```

При каких значениях происходит переобучение? А недообучение? Почему?

Ответ:

Оценка: 8