

# Машинное обучение, ФКН ВШЭ

## Практическое задание 3

### Общая информация

Дата выдачи: 22.09.2018

Мягкий дедлайн: 7:59MSK 02.10.2018

Жесткий дедлайн: 23:59MSK 03.10.2018

## О задании

Задание состоит из двух частей: в **первой** части вы научитесь применять готовые модели из `sklearn` на данных и исследуете особенности разных видов регуляризации; во **второй** части вы реализуете собственный класс линейной регрессии для нестандартной функции потерь (включая процесс обучения) и исследуете скорость сходимости различных градиентных методов для этой модели.

## 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

Задания сдаются через систему `anytask`. Присылать необходимо ноутбук с выполненным заданием. Часть задания сдаётся в Яндекс.Контест.

Для удобства проверки самостоятельно посчитайте свою максимальную оценку (исходя из набора решенных задач) и укажите ниже.

Оценка: 8.5

In [2]:

```
1 import pandas
2 import matplotlib.pyplot as plt
3 from sklearn.model_selection import train_test_split
4 from sklearn.preprocessing import MinMaxScaler
5 from sklearn.linear_model import LinearRegression, Lasso
6 from sklearn.metrics import r2_score
```

```
In [163]: 1 def scale_data(x_data):
2         scaler = MinMaxScaler()
3         return scaler.fit_transform(x_data)
4
5 def pair_plot(df, target):
6     ncol, nrow = 7, df.shape[1] // 7 + (df.shape[1] % 7 > 0)
7     plt.figure(figsize=(ncol * 4, nrow * 4))
8
9     for i, feature in enumerate(df.columns):
10         plt.subplot(nrow, ncol, i + 1)
11         plt.scatter(df[feature], target, s=10, marker='o', alpha=.5)
12         plt.xlabel(feature)
13         if i % ncol == 0:
14             plt.ylabel('target')
15
16 def graph_by_points(x, y, x_name, y_name):
17     f, (ax2) = plt.subplots(1, 1, figsize=(13, 5))
18
19     ax2.plot(x, y)
20     ax2.set_xlabel(x_name)
21     ax2.set_ylabel(y_name)
22     ax2.grid()
23     ax2.legend()
24
25     f.show()
26
```

```
In [4]: 1 practice_data = pandas.read_csv('./homework-practice-03-data.csv')
2       x_practice_data = practice_data.drop(columns=['f1'])
3       y_practice_data = practice_data['f1']
```

## Линейная регрессия из коробки

1. [0.5 балла] Разбейте выборку, загруженную в ячейке выше в переменную `practice_data`, на обучающую и тестовую части в соотношении 8:2.

```
In [5]: 1 x_train, x_test, y_train, y_test = train_test_split(
2         x_practice_data, y_practice_data, train_size=0.8, random_s
3 x_train.head()
```

```
/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)
```

Out[5]:

	f0	f2	f3	f4	f5	f6
<b>228</b>	13.28	8.947412	1.854125	0.304199	0.637632	724.955750
<b>208</b>	24.27	0.731584	4.289619	0.991890	0.873039	1514.920900
<b>96</b>	27.28	1.730259	1.431718	0.859581	0.513208	58.940865
<b>167</b>	31.71	-2.800439	-2.605946	0.878373	2.113398	7614.162600
<b>84</b>	15.98	1.332029	9.779562	0.983426	1.473478	1886.239400

**2. [0.5 балла]** В качестве целевой переменной для задачи регрессии будем использовать значения признака f1. Обучите, а затем провалидируйте на тестовых данных следующие модели, используя в качестве метрики качества  $R^2$ :

- `LinearRegression` ([http://scikit-learn.org/stable/modules/generated/sklearn.linear\\_model.LinearRegression.html](http://scikit-learn.org/stable/modules/generated/sklearn.linear_model.LinearRegression.html));
- `Lasso` ([http://scikit-learn.org/stable/modules/generated/sklearn.linear\\_model.Lasso.html](http://scikit-learn.org/stable/modules/generated/sklearn.linear_model.Lasso.html)) (линейная регрессия с  $L1$ -регуляризатором) с коэффициентом регуляризации, равным 0.01.

Не забудьте отмасштабировать данные перед обучением моделей!

```
In [118]: 1 x_train_scaled = scale_data(x_train)
2 x_test_scaled = scale_data(x_test)
3
4 def train_and_print_r2_score(model, train_data, test_data):
5     model.fit(*train_data)
6     print(" model has been trained")
7     prediction = model.predict(test_data[0])
8     print(" model r2_score:", r2_score(prediction, test_data[1]))
9
10 print("simple_linear_model:")
11 simple_linear_model = LinearRegression()
12 train_and_print_r2_score(simple_linear_model, (x_train_scaled, y_train_scaled))
13
14 print("lasso_linear_model:")
15 lasso_linear_model = Lasso(alpha=0.01)
16 train_and_print_r2_score(lasso_linear_model, (x_train_scaled, y_train_scaled))
17
```

```
simple_linear_model:
  model has been trained
  model r2_score: -0.19011579351920727
lasso_linear_model:
  model has been trained
  model r2_score: -0.10458100714571472
```

**3. [1 балл] Изучите значения параметров получившихся моделей и сравните количество строго нулевых весов в них.**

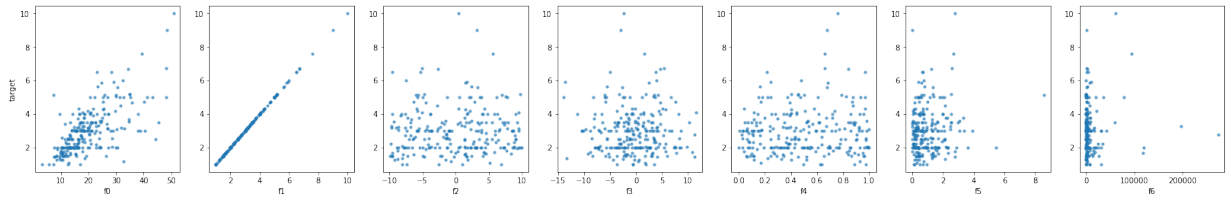
```
In [119]: 1 print("simple_linear_model coefficients:", simple_linear_model.coef_)
2 print("lasso_linear_model coefficients:", lasso_linear_model.coef_)
```

```
simple_linear_model coefficients: [ 4.76700358 -0.08580226 -0.288497
96 0.10520103 1.38159276 0.01781414]
lasso_linear_model coefficients: [ 4.52409084 -0.          -0.
0.03807787 0.72527362 0.          ]
```

Можно видеть, что L1-регуляризация действительно вывела 3 нулевых веса для входных аргументов.

**4. [1 балл] Нарисуйте попарные scatter plot для всех признаков в исходных данных (включая f1). Можно ли что-то сказать о связи признаков, веса которых были занулены методом Lasso, с целевой переменной?**

```
In [121]: 1 pair_plot(practice_data, practice_data['f1'])
```



Можно видеть, что L1-регуляризация выявила самый полезный (конечно, не рассматривая f1) с точки зрения линейной зависимости признак f0. Кроме того, с незначительными весами были взяты признаки f4, f5.

## Реализация градиентного спуска

5. [4 балла] Реализуйте модель линейной регрессии для функции потерь Huber loss, обучаемую градиентным спуском:

$$L_{\delta}(a, y) = \begin{cases} \frac{1}{2}(y - a)^2, & |y - a| \leq \delta, \\ \delta |y - a| - \frac{1}{2}\delta^2 & \text{иначе.} \end{cases}$$

Все вычисления должны быть векторизованы, циклы средствами python допускается использовать только для итераций градиентного спуска. В качестве критерия останова необходимо использовать (одновременно):

- проверку на евклидову норму разности весов на двух соседних итерациях (например, меньше некоторого малого числа порядка  $10^{-6}$ , задаваемого параметром `tolerance`);
- достижение максимального числа итераций (например, 10000, задаваемого параметром `max_iter`).

Необходимо реализовать метод полного и стохастического градиентных спусков, а также поддержать метод `momentum` при помощи параметра `alpha` (способ оценивания градиента должен задаваться при помощи параметра `gd_type`).

Чтобы проследить, что оптимизационный процесс действительно сходится, будем использовать атрибут класса `loss_history` — в нём после вызова метода `fit` должны содержаться значения функции потерь для всех итераций, начиная с первой (до совершения первого шага по антиградиенту).

Инициализировать веса можно случайным образом или нулевым вектором. Ниже приведён шаблон класса, который должен содержать код реализации модели.

Python-файл с реализованным классом необходимо сдать на проверку в [Яндекс.Контекст \(https://contest.yandex.ru/contest/9247/\)](https://contest.yandex.ru/contest/9247/)

Укажите ссылку на посылку (run-report): <https://contest.yandex.ru/contest/9247/run-report/12232591/> (<https://contest.yandex.ru/contest/9247/run-report/12232591/>)

In [156]:

```
1 import numpy as np
2 from sklearn.base import BaseEstimator
3 import random
4
5 class HuberReg(BaseEstimator):
6     def __init__(
7         self,
8         delta = 1.0,
9         gd_type = 'stochastic',
10        tolerance = 1e-4,
11        max_iter = 1000,
12        w0 = None,
```

```

13         alpha = 1e-3,
14         eta = 1e-2
15     ):
16         self.delta = delta
17         self.gradient_type = gd_type
18         self.tolerance = tolerance
19         self.max_iterations_count = max_iter
20         self.alpha = alpha
21         self.eta = eta
22         self.loss_history = []
23         self.w = None
24         self.initial_weights = w0
25
26     def fit(self, x_data, y_data):
27         self.init_weights(x_data)
28         self.update_loss_history(x_data, y_data)
29
30         previous, current = 0, 0
31         for iteration in range(self.max_iterations_count):
32             if self.gradient_type == 'full':
33                 current = self.alpha * previous - self.eta * self.
34             elif self.gradient_type == 'stochastic':
35                 row_index = random.randint(0, np.shape(x_data)[0])
36                 current = self.alpha * previous - self.eta * self.
37                     x_data[row_index:row_index + 1],
38                     y_data[row_index:row_index + 1]
39             )
40             else:
41                 raise Exception('unknown gradient type')
42             previous = current
43             self.w += previous
44             self.update_loss_history(x_data, y_data)
45             if abs(np.linalg.norm(current)) < self.tolerance:
46                 break
47         return self
48
49     def init_weights(self, x_data):
50         if self.initial_weights is None:
51             self.w = np.random.rand(x_data.shape[1])
52         else:
53             self.w = self.initial_weights
54
55     def predict(self, x_data):
56         if self.w is None:
57             raise Exception('Not trained yet')
58         else:
59             return x_data.dot(self.w)
60
61     def calc_gradient(self, x_data, y_data):
62         distation = self.calc_distation(x_data, y_data)
63         A = self.delta * x_data.T.dot((distation > self.delta).ast
64         B = self.delta * x_data.T.dot((distation < -self.delta).as
65         C = np.dot(
66             x_data.T, (-distation * (np.absolute(-distation) <= se

```



```

67         return (A + B + C) / x_data.shape[0]
68
69     def update_loss_history(self, x_data, y_data):
70         self.loss_history.append(self.calc_loss(x_data, y_data))
71
72     def calc_loss(self, x_data, y_data):
73         distation = self.calc_distation(x_data, y_data)
74         A = 0.5 * (distation) ** 2
75         B = self.delta * np.absolute(distation) - 0.5 * self.delta
76         return np.mean(np.where(np.absolute(distation) <= self.del
77
78     def calc_distation(self, x_data, y_data):
79         return y_data - x_data.dot(self.w)

```

**6. [1.5 балла] Обучите и провалидируйте модель на тех же данных, сравните качество с предыдущими методами. Исследуйте влияние параметров `max_iter` и `alpha` на процесс оптимизации. Согласуется ли оно с вашими ожиданиями?**

```

In [161]: 1 model = HuberReg()
          2 model.fit(x_train, y_train)
          3 model.calc_gradient(x_train, y_train)
          4 r2_score(model.predict(np.array(x_test)), np.array(y_test))

```

Out[161]: -0.05079108347565242

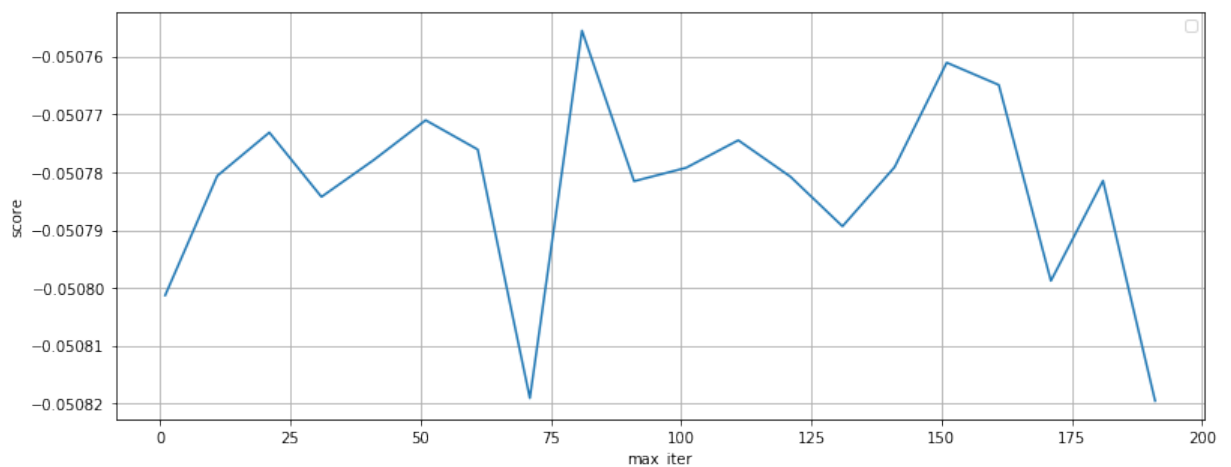
Результат довольно низок, однако перебор параметров может помочь.

```
In [172]: 1 def build_and_print_max_iter_graph():
2         scores = []
3         values = np.arange(1, 200, 10)
4         for it in values:
5             model = HuberReg(max_iter=it)
6             model.fit(x_train, y_train)
7             scores.append(r2_score(model.predict(np.array(x_test)), np
8         graph_by_points(values, scores, "max_iter", "score")
9     build_and_print_max_iter_graph()
```

No handles with labels found to put in legend.

/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 [177]:

```

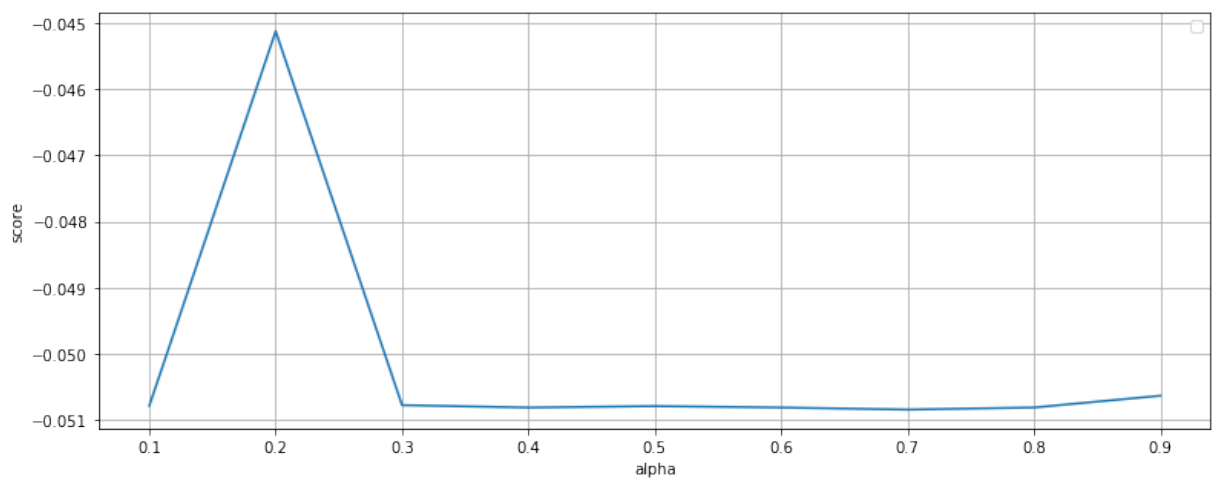
1 def build_and_print_alpha_graph():
2     scores = []
3     values = np.arange(0.1, 1, 0.1)
4     for it in values:
5         model = HuberReg(alpha=it)
6         model.fit(x_train, y_train)
7         scores.append(r2_score(model.predict(np.array(x_test)), np
8     graph_by_points(values, scores, "alpha", "score")
9 build_and_print_alpha_graph()

```

No handles with labels found to put in legend.

/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())
```



Результат целиком и полностью согласуется с нашими ожиданиями.

**7. [1.5 балла] Постройте графики (на одной и той же картинке) зависимости величины функции потерь от номера итерации для полного, стохастического градиентного спуска, а также для полного градиентного спуска с методом инерции. Сделайте выводы о скорости сходимости различных модификаций градиентного спуска.**

Не забывайте о том, что должны из себя представлять *красивые* графики!

In [30]:

1