

Семинар: Линейная регрессия. Задачи

Задание: Линейная регрессия.

В этом задании мы рассмотрим метод линейной регрессии (метод наименьших квадратов - МНК). Мы будем работать с данными, содержащим информацию о бриллиантах. Описание можно посмотреть [здесь](https://www.kaggle.com/shivam2503/d) (<https://www.kaggle.com/shivam2503/d>)

Content

price - цена в долларах США (326 - 18,823)

carat - вес бриллианта (0.2 - 5.01)

cut - качество среза (Fair, Good, Very Good, Premium, Ideal)

color - цвет бриллианта (от J (худший) до D (лучший))

clarity - степень чистоты бриллианта (от I1 (худшая), SI2, SI1, VS2, VS1, VVS2, VVS1, IF (лучшая))

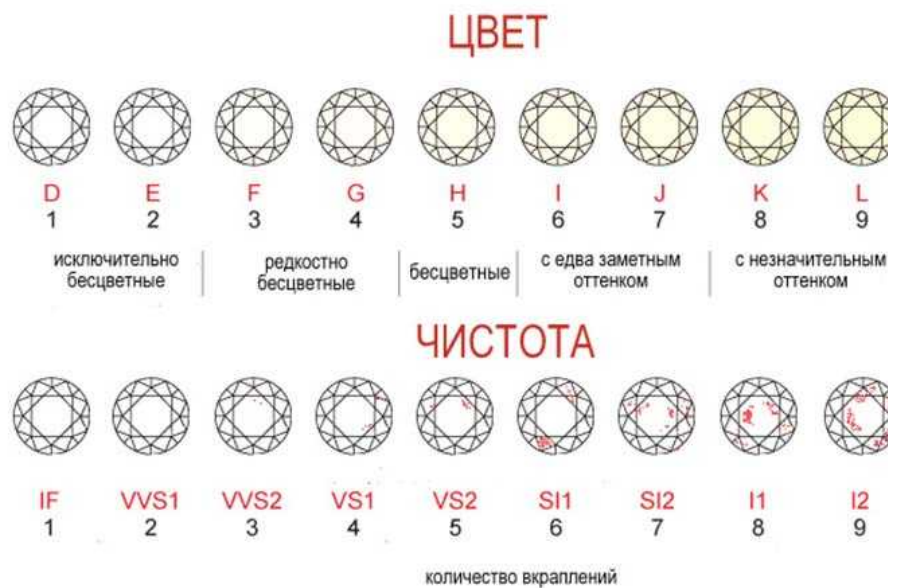
x - длина в миллиметрах (0 - 10.74)

y - ширина в миллиметрах (0 - 58.9)

z - глубина в миллиметрах (0 - 31.8)

depth - общий процент глубины = $z / \text{mean}(x, y) = 2 * z / (x + y)$ (43 - 79)

table - ширина верхней плоскости бриллианта в самом широком месте (43 - 95)



Наглядная иллюстрация цвета и чистоты бриллианта с чуть более расширенными показателями (до M по цвету и до I2 по чистоте)

```

В [1]: 1 import numpy as np
        2 import pandas as pd
        3 import matplotlib.pyplot as plt
        4 import seaborn as sns
        5 from sklearn.metrics import r2_score
        6
        7 import warnings
        8 warnings.filterwarnings("ignore")

В [5]: 1 data = pd.read_csv('diamonds.csv') # загружаем датафрейм
        2 data.head() # выводим первые 5 строк датафрейма

Out[5]:

```

| | Unnamed: 0 | carat | cut | color | clarity | depth | table | price | x | y | z |
|---|------------|-------|---------|-------|---------|-------|-------|-------|------|------|------|
| 0 | 1 | 0.23 | Ideal | E | SI2 | 61.5 | 55.0 | 326 | 3.95 | 3.98 | 2.43 |
| 1 | 2 | 0.21 | Premium | E | SI1 | 59.8 | 61.0 | 326 | 3.89 | 3.84 | 2.31 |
| 2 | 3 | 0.23 | Good | E | VS1 | 56.9 | 65.0 | 327 | 4.05 | 4.07 | 2.31 |
| 3 | 4 | 0.29 | Premium | I | VS2 | 62.4 | 58.0 | 334 | 4.20 | 4.23 | 2.63 |
| 4 | 5 | 0.31 | Good | J | SI2 | 63.3 | 58.0 | 335 | 4.34 | 4.35 | 2.75 |

```
In [7]: 1 data.tail() # выводим последние 5 строк датафрейма
```

```
Out[7]:
```

| | Unnamed: 0 | carat | cut | color | clarity | depth | table | price | x | y | z |
|-------|------------|-------|-----------|-------|---------|-------|-------|-------|------|------|------|
| 53935 | 53936 | 0.72 | Ideal | D | SI1 | 60.8 | 57.0 | 2757 | 5.75 | 5.76 | 3.50 |
| 53936 | 53937 | 0.72 | Good | D | SI1 | 63.1 | 55.0 | 2757 | 5.69 | 5.75 | 3.61 |
| 53937 | 53938 | 0.70 | Very Good | D | SI1 | 62.8 | 60.0 | 2757 | 5.66 | 5.68 | 3.56 |
| 53938 | 53939 | 0.86 | Premium | H | SI2 | 61.0 | 58.0 | 2757 | 6.15 | 6.12 | 3.74 |
| 53939 | 53940 | 0.75 | Ideal | D | SI2 | 62.2 | 55.0 | 2757 | 5.83 | 5.87 | 3.64 |

Мы будем решать задачу предсказания цены бриллианта price в зависимости от его характеристик.

Задача 1.1. Есть ли в наборе данных пропущенные значения? Если да, удалите их.

```
In [3]: 1 data.isnull().any() # пропущенных значения в датафрейме отсутствуют
```

```
Out[3]: Unnamed: 0    False
carat           False
cut             False
color           False
clarity         False
depth           False
table           False
price           False
x              False
y              False
z              False
dtype: bool
```

Задача 1.2. Есть ли в наборе данных бессмысленные столбцы (признаки, не несущие дополнительной информации

```
In [4]: 1 data = data.drop('Unnamed: 0', axis=1) # удаляем столбец "Unnamed: 0" с помощью функции drop
2 data.head()
```

```
Out[4]:
```

| | carat | cut | color | clarity | depth | table | price | x | y | z |
|---|-------|---------|-------|---------|-------|-------|-------|------|------|------|
| 0 | 0.23 | Ideal | E | SI2 | 61.5 | 55.0 | 326 | 3.95 | 3.98 | 2.43 |
| 1 | 0.21 | Premium | E | SI1 | 59.8 | 61.0 | 326 | 3.89 | 3.84 | 2.31 |
| 2 | 0.23 | Good | E | VS1 | 56.9 | 65.0 | 327 | 4.05 | 4.07 | 2.31 |
| 3 | 0.29 | Premium | I | VS2 | 62.4 | 58.0 | 334 | 4.20 | 4.23 | 2.63 |
| 4 | 0.31 | Good | J | SI2 | 63.3 | 58.0 | 335 | 4.34 | 4.35 | 2.75 |

Задача 1.3. Линейная регрессия основана на предположении о линейной связи между признаками и целевой переменной. Выбором переменных для включения в модель имеет смысл проверить, насколько эта связь выполняется. Для этого потребуются выборочные корреляции между признаками. Выведите матрицу выборочных корреляций между всеми признаками и целевой переменной (то есть в этой матрице будет $k + 1$ строка, где k – количество вещественных признаков).

Какие вещественные признаки коррелируют с целевой переменной больше всего?

```
In [5]: 1 (data.dtypes == "float").sum() # выводим количество вещественных признаков в датафрейме
```

```
Out[5]: 6
```

```

В [6]: 1 plt.figure(figsize = (12,6)) # создаем фигуру, в которой будет построен график и указываем её
2 sns.heatmap(data.corr(), annot=True, linewidths= 1)
3 # визуализируем матрицу корреляции с помощью тепловой карты из библиотеки seaborn,
4 # параметром annot=True включаем аннотацию к матрице. в каждой клетке - коэф. кор.

```

Out[6]: <Axes: >



Вещественных признаков 6, а в матрице 7 строк (значит все правильно). Видно, что сильнее всего с ценой и carat, дальше идут, с небольшой разницей, параметры длины, ширины и глубины бриллианта. Признак depth переменной очень слабо.

Задача 1.4. Так как линейная модель складывает значения признаков с некоторыми весами, нам нужно аккуратно с признаками. Закодируйте категориальные переменные при помощи OneHot-кодирования.

```

В [7]: 1 y = data['price'] # выделение зависимой переменной
2 X = data.drop('price', axis=1) # удаление зависимой переменной из матрицы признаков
3 X.head()

```

Out[7]:

| | carat | cut | color | clarity | depth | table | x | y | z |
|---|-------|---------|-------|---------|-------|-------|------|------|------|
| 0 | 0.23 | Ideal | E | SI2 | 61.5 | 55.0 | 3.95 | 3.98 | 2.43 |
| 1 | 0.21 | Premium | E | SI1 | 59.8 | 61.0 | 3.89 | 3.84 | 2.31 |
| 2 | 0.23 | Good | E | VS1 | 56.9 | 65.0 | 4.05 | 4.07 | 2.31 |
| 3 | 0.29 | Premium | I | VS2 | 62.4 | 58.0 | 4.20 | 4.23 | 2.63 |
| 4 | 0.31 | Good | J | SI2 | 63.3 | 58.0 | 4.34 | 4.35 | 2.75 |

```

В [8]: 1 X_dum = pd.get_dummies(X, drop_first=True)
2 # кодирую данные, передаю параметр drop_first=True, чтобы избежать линейной зависимости
3 X_dum # вывод датафрейма

```

Out[8]:

| | carat | depth | table | x | y | z | cut_Good | cut_Ideal | cut_Premium | cut_Very Good | ... | color_H | color_I | color_J | cl |
|-------|-------|-------|-------|------|------|------|----------|-----------|-------------|---------------|-----|---------|---------|---------|----|
| 0 | 0.23 | 61.5 | 55.0 | 3.95 | 3.98 | 2.43 | 0 | 1 | 0 | 0 | ... | 0 | 0 | 0 | |
| 1 | 0.21 | 59.8 | 61.0 | 3.89 | 3.84 | 2.31 | 0 | 0 | 1 | 0 | ... | 0 | 0 | 0 | |
| 2 | 0.23 | 56.9 | 65.0 | 4.05 | 4.07 | 2.31 | 1 | 0 | 0 | 0 | ... | 0 | 0 | 0 | |
| 3 | 0.29 | 62.4 | 58.0 | 4.20 | 4.23 | 2.63 | 0 | 0 | 1 | 0 | ... | 0 | 1 | 0 | |
| 4 | 0.31 | 63.3 | 58.0 | 4.34 | 4.35 | 2.75 | 1 | 0 | 0 | 0 | ... | 0 | 0 | 1 | |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | |
| 53935 | 0.72 | 60.8 | 57.0 | 5.75 | 5.76 | 3.50 | 0 | 1 | 0 | 0 | ... | 0 | 0 | 0 | |
| 53936 | 0.72 | 63.1 | 55.0 | 5.69 | 5.75 | 3.61 | 1 | 0 | 0 | 0 | ... | 0 | 0 | 0 | |
| 53937 | 0.70 | 62.8 | 60.0 | 5.66 | 5.68 | 3.56 | 0 | 0 | 0 | 1 | ... | 0 | 0 | 0 | |
| 53938 | 0.86 | 61.0 | 58.0 | 6.15 | 6.12 | 3.74 | 0 | 0 | 1 | 0 | ... | 1 | 0 | 0 | |
| 53939 | 0.75 | 62.2 | 55.0 | 5.83 | 5.87 | 3.64 | 0 | 1 | 0 | 0 | ... | 0 | 0 | 0 | |

53940 rows × 23 columns



Задача 1.5. Разделите выборку на тренировочную и тестовую. Долю тестовой выборки укажите равной 0.3.

```
B [9]: 1 from sklearn.model_selection import train_test_split # импортирую нужный модуль
2 X_train, X_test, y_train, y_test = train_test_split(X_train, y_train, test_size=0.3, random_state=1)
3 # разделяю выборку. За размер доли тестовой выборки отвечает параметр test_size, его, как указ
4 # кроме того передаю параметр random_state (аналог np.random.seed()), для того, чтобы разбиение
```

Задача 1.6. Зачастую при использовании линейных моделей вещественные признаки масштабируются. При этом они теряют прямую статистическую интерпретацию ("при увеличении X_1 на 1, y увеличивается на w_1 "), но приобретаю задач машинного обучения. В этой задаче масштабируйте вещественные признаки тренировочной и тестовой вы StandardScaler .

```
B [10]: 1 from sklearn.preprocessing import StandardScaler # импортирую модуль для нормировки
2
3 scaler = StandardScaler() # создаю "нормировщик"
4 scaler.fit(X_train) # обучаю его на данных из X_train
5 scaled_X_train = scaler.transform(X_train) # с помощью уже обученного "нормир
6 scaled_X_test = scaler.transform(X_test) # датафреймах X_train и X_test
```

Задача 1.7. Оцените линейную регрессию на тренировочной выборке. Выведите среднеквадратичную ошибку на тр выборках.

```
B [11]: 1 from sklearn.linear_model import LinearRegression
2 from sklearn.metrics import mean_squared_error

B [12]: 1 # создаем модель линейной регрессии
2 regr = LinearRegression()
3
4 # обучаем модель на train датасете, используя масштабированные признаки
5 regr.fit(scaled_X_train, y_train)
6
7 # предсказываем значения y_train и y_test с помощью обученной модели
8 y_train_pred = regr.predict(scaled_X_train)
9 y_test_pred = regr.predict(scaled_X_test)
10
11 # выводим значение средней квадратичной ошибки (MSE) для train и test датасетов
12 print("Train MSE: ")
13 print("Test MSE: ")
14
15 # выводим коэффициент квадрата детерминации (R2) для train и test датасетов, чтобы проверить т
16 print(f'\nTrain: {(r2_score(y_train, y_train_pred))}')
17 print(f'Test: {(r2_score(y_test, y_test_pred))}')
```

Train MSE: 1144.60
Test MSE: 1095.78

Train: 0.9189349877290272
Test: 0.9217238681492771

Задача 1.8. Изучите документацию модуля LinearRegression и выведите полученные оценки коэффициентов (вс вещественные переменные, оценки коэффициентов которых по модулю на порядок превышают оценки прочих вещ

```
B [13]: 1 df = pd.DataFrame() # создаю датафрейм
2 for i in range(len(list(X_train.columns.values))): # для каждого признака нахожу соответс
3     df[list(X_train.columns.values)[i]] = pd.Series(abs(regr.coef_[i]))
4
5 cat_features_mask_1 = (X.dtypes == "object").values # создаю маску для отбора только вещес
6 df[list(X[X.columns[~cat_features_mask_1]].columns)].T.sort_values(0, ascending=False)
7 # делаю таблицу и осуществляю коэффициенты по убыванию
```

```
Out[13]:
```

| | 0 |
|-------|-------------|
| carat | 5277.247129 |
| x | 1050.102162 |
| depth | 83.585954 |
| table | 56.869625 |
| z | 12.888928 |
| y | 0.654293 |

Среди вещественных переменных слишком большие веса получились у веса бриллианта и параметра x. Следующи примерно одинаковые по соотношению веса, а вот последняя переменная (параметр y) имеет очень маленький вес

Задача 1.9. Как можно заметить из анализа корреляционной матрицы в задаче 3.3, между некоторыми признаками корреляция, что может быть индикатором проблемы мультиколлинеарности (наличия линейной зависимости между переменными (факторами) модели). Различия в порядке коэффициентов, выявленные в предыдущей задаче также присутствие. Как известно, для решения этой проблемы можно

либо исключить некоторые признаки из модели, либо использовать регуляризацию. Мы воспользуемся вторым в

Вспомним, что смысл регуляризации заключается в том, чтобы изменить функцию потерь так, чтобы устранить проблему мультиколлинеарности. При L1-регуляризации предлагается минимизировать следующую функцию потерь:

$$\|y - X\hat{w}\|^2 + \alpha \sum_{i=1}^k |w_i|$$

Такая модель называется Lasso-регрессией.

При L2-регуляризации предлагается минимизировать следующую функцию потерь:

$$\|y - X\hat{w}\|^2 + \frac{1}{2} \alpha \|w\|^2$$

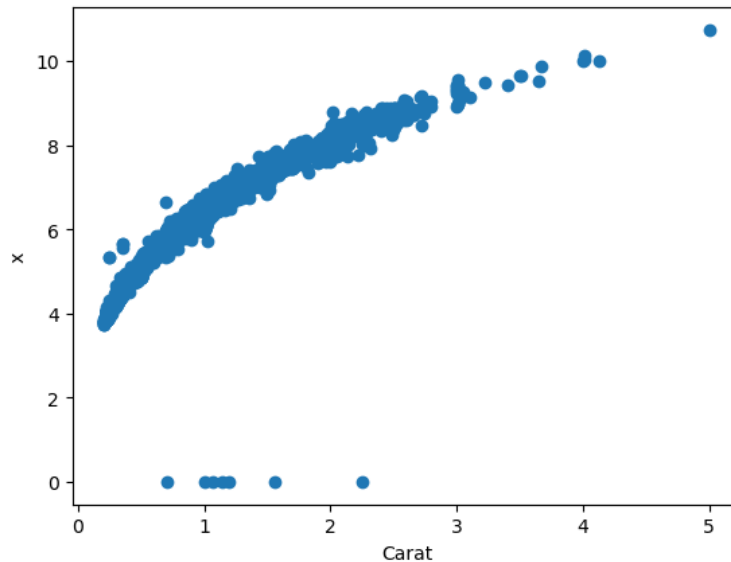
Такая модель называется Ridge-регрессией.

Выберите два признака, сильно коррелирующих между собой, и постройте зависимость (scatter) одного от другого.

Обучите Lasso-регрессию и Ridge-регрессию, установив гиперпараметр регуляризации равным 10. Для этого используйте `sklearn`. Сильно ли уменьшились веса? Сделайте вывод о том, насколько сильно проблема мультиколлинеарности.

```
B [14]: 1 import numpy as np
        2 import matplotlib
        3 from matplotlib import pyplot as plt
```

```
B [15]: 1 plt.scatter(X_train['carat'], X_train['x'])
        2 # на матрице корреляции было видно, что переменные carat и x между собой сильно коррелируют.
        3 # далее вывожу scatter-plot, чтобы проследить их зависимость
        4 plt.xlabel('Carat')
        5 plt.ylabel('x')
        6 plt.show();
```



```
B [16]: 1 from sklearn.linear_model import Lasso, Ridge # импортирую Lasso и Ridge
        2
        3 lasso = Lasso(10.0).fit(scaled_X_train, y_train) # создаю Lasso-модель с λ = 10
        4
        5 ridge = Ridge(10.0).fit(scaled_X_train, y_train) # делаю то же самое с Ridge
```

```
B [17]: 1 print(regr.coef_)
        2 print('\n', lasso.coef_)
        3 print('\n', ridge.coef_)

[ 5.27724713e+03 -8.35859541e+01 -5.68696252e+01 -1.05010216e+03
 -6.54293136e-01 -1.28889279e+01  1.71502354e+02  4.16023988e+02
  3.39486233e+02  3.13132588e+02 -8.45948129e+01 -1.02390339e+02
 -2.00244063e+02 -3.57144807e+02 -4.40474155e+02 -5.25835595e+02
  9.87614017e+02  1.59869330e+03  1.04783493e+03  1.67852566e+03
  1.82339758e+03  1.29267929e+03  1.46168957e+03]

[4783.41094049 -88.18838824 -67.34812205 -605.4571011 -0.
 -9.04245448  45.90699967  207.92331741  150.09822521  139.4370833
 -12.65867611 -35.62012708 -122.72563295 -284.88080225 -369.04859265
 -467.2452535  669.45958233  863.02801363  409.79199968  1053.42448122
 1102.04028249  852.95943422  958.90274725]

[ 5.23973797e+03 -8.18029757e+01 -5.70570935e+01 -1.00923042e+03
 -2.55973659e+00 -1.64756784e+01  1.72041187e+02  4.17140870e+02
  3.40412927e+02  3.14372969e+02 -8.41350398e+01 -1.02177926e+02
 -1.99692501e+02 -3.56285462e+02 -4.39199150e+02 -5.24544818e+02
  9.74870405e+02  1.56782640e+03  1.02157962e+03  1.65263349e+03
  1.79332538e+03  1.27483844e+03  1.44104250e+03]
```

Смотря на коэффициенты кажется, что Lasso справилось с задачей лучше, чем Ridge. В общем веса немног один коэффициент даже занулились. Проблема мультиколлинеарности была, но сказать, что она очень сил нельзя, но нам важно, что мы смогли ее подправить и улучшить результат.

Задача 3.10 Как обсуждалось на семинарах, Lasso-регрессию можно использовать для отбора наиболее информат следующих значений параметра регуляризации α : 0.1, 1, 10, 100, 200 – обучите Lasso- и Ridge-регрессии и построй евклидовой нормы весов (`np.linalg.norm()` от вектора оценок коэффициентов) в зависимости от параметра α . К является численной характеристикой величины вектора, а потому по норме можно судить о том, насколько больши вектор оценок коэффициентов.

Какой метод сильнее уменьшает веса?

In [18]:

```
1 import numpy as np
```

In [19]:

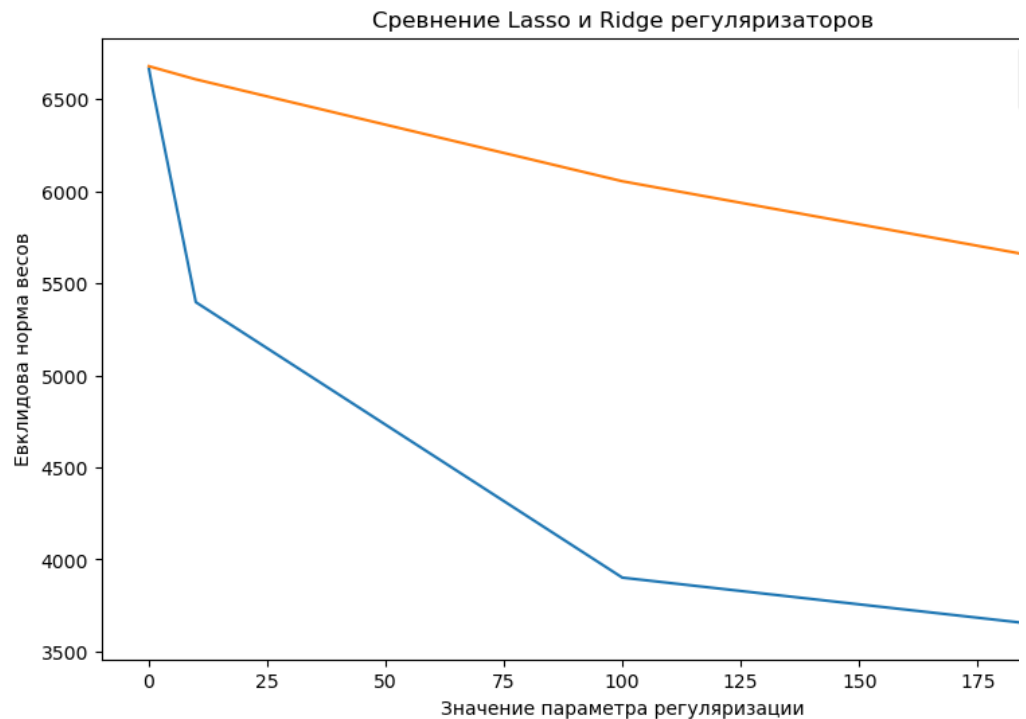
```
1 alpha = [0.1, 1, 10, 100, 200]
2 lasso_evklid_coefs = []
3 ridge_evklid_coefs = []
4 for i in alpha:
5     lasso_evklid_coefs.append(np.linalg.norm(Lasso(i).fit(scaled_X_train, y_train).coef_))
6     ridge_evklid_coefs.append(np.linalg.norm(Ridge(i).fit(scaled_X_train, y_train).coef_))
```

In [20]:

```
1 fig, axes = plt.subplots(figsize=(10,6))
2 axes.plot(alpha, lasso_evklid_coefs, alpha, ridge_evklid_coefs)
3 axes.legend(['lasso', 'ridge'])
4 plt.xlabel('Значение параметра регуляризации')
5 plt.ylabel('Евклидова норма весов')
6 plt.title('Сравнение Lasso и Ridge регуляризаторов')
```

создаем поверхность и задае
рисуем нужные графики
добавляем легенду
подписываем оси графика
даем графику название

Out[20]: Text(0.5, 1.0, 'Сравнение Lasso и Ridge регуляризаторов')



Из графика видно, что Lasso сильнее уменьшает веса как при маленьких значениях параметра регуляризации, так и