

Семинар: Градиентный спуск. Задачи

```

В [37]: 1 from typing import Iterable, List
        2
        3 import matplotlib.pyplot as plt
        4 import numpy as np
        5 import pandas as pd
        6 import abc
        7 from sklearn.impute import SimpleImputer
        8 from sklearn.model_selection import train_test_split
        9 from sklearn.preprocessing import StandardScaler
       10
       11 %matplotlib inline

```

Часть 1. Градиентный спуск (вспомним формулы)

Функционал ошибки, который мы применяем в задаче регрессии — Mean Squared Error:

$$Q(w, X, y) = \frac{1}{\ell} \sum_{i=1}^{\ell} (\langle x_i, w \rangle - y_i)^2$$

где x_i — это i -ый объект датасета, y_i — правильный ответ для i -го объекта, а w — веса нашей линейной модели.

Можно показать, что для линейной модели, функционал ошибки можно записать в матричном виде следующим образом:

$$Q(w, X, y) = \frac{1}{l} (y - Xw)^T (y - Xw)$$

или

$$Q(w, X, y) = \frac{1}{l} \|Xw - y\|^2$$

где X — это матрица объекты-признаки, а y — вектор правильных ответов.

Для того чтобы воспользоваться методом градиентного спуска, нам нужно посчитать градиент нашего функционала. Он будет выглядеть так:

$$\nabla_w Q(w, X, y) = \frac{2}{l} X^T (Xw - y)$$

Формула для одной итерации градиентного спуска выглядит следующим образом:

$$w^t = w^{t-1} - \eta \nabla_w Q(w^{t-1}, X, y)$$

Где w^t — значение вектора весов на t -ой итерации, а η — параметр learning rate, отвечающий за размер шага.

Часть 2. Линейная регрессия (продолжение части 1 - используем часть 1).

Создадим класс для линейной регрессии. Он будет использовать интерфейс, знакомый нам из библиотеки sklearn.

В методе `fit` мы будем подбирать веса w при помощи градиентного спуска нашим методом `gradient_descent`.

В методе `predict` мы будем применять нашу регрессию к датасету.

Задание 2.1: Допишите код в методах `fit` и `predict` класса `LinearRegression_1`.

В методе `fit` вам нужно инициализировать веса w (например, из случайного распределения), применить градиентный спуск и сохранить последние веса w из траектории.

В методе `predict` вам нужно применить линейную регрессию и вернуть вектор ответов.

Обратите внимание, что объект лосса (функционала ошибки) передаётся в момент инициализации и хранится в атрибуте `loss`. Используйте его в методе `fit` для градиентного спуска (например, `c_n_iterations: int = 1000`).

```

B [38]: 1 class BaseLoss(abc.ABC):
2         @abc.abstractmethod
3         def calc_loss(self, X: np.ndarray, y: np.ndarray, w: np.ndarray) -> float:
4             raise NotImplementedError
5
6         @abc.abstractmethod
7         def calc_grad(self, X: np.ndarray, y: np.ndarray, w: np.ndarray) -> np.ndarray:
8             raise NotImplementedError
9
10        class MSELoss(BaseLoss):
11            def calc_loss(self, X: np.ndarray, y: np.ndarray, w: np.ndarray) -> float:
12                return (np.linalg.norm(np.matmul(X,w) - y)**2)/len(y)
13
14            def calc_grad(self, X: np.ndarray, y: np.ndarray, w: np.ndarray) -> np.ndarray:
15                return 2/len(y)*np.matmul(np.transpose(X), (np.matmul(X, w) - y))
16
B [39]: 1 def gradient_descent(
2         w_init: np.ndarray,
3         X: np.ndarray,
4         y: np.ndarray,
5         loss: BaseLoss,
6         lr: float,
7         n_iterations: int = 1000,
8     ) -> List[np.ndarray]:
9         grad_list = []
10        for i in range(n_iterations):
11            grad_list.append(w_init)
12            w_init = w_init - lr * loss.calc_grad(X, y, w_init)
13        return grad_list
B [40]: 1 class LinearRegression_1:
2         def __init__(self, loss: BaseLoss, lr: float = 0.01) -> None:
3             self.loss = loss #функционал ошибки
4             self.lr = lr #градиентный шаг
5
6         def fit(self, X: np.ndarray, y: np.ndarray) -> "LinearRegression_1":
7             X = np.asarray(X)
8             y = np.asarray(y)
9             X = np.hstack([X, np.ones([X.shape[0], 1])])
10            w_init = np.random.normal(size=(X.shape[1],))
11
12            w_arr = gradient_descent(w_init, X, y, self.loss, self.lr, 1000)
13            self.w = w_arr[-1]
14
15            return self
16
17
18        def predict(self, X: np.ndarray) -> np.ndarray:
19            assert hasattr(self, "w"), "Linear regression must be fitted first" #проверяем, о
20            X = np.asarray(X)
21            X = np.hstack([X, np.ones([X.shape[0], 1])])
22
23            y = []
24            for i in range(X.shape[0]):
25                y1=np.sum(X[i]*self.w)
26                y.append(y1)
27
28            return y

```

Класс линейной регрессии создан. Более того, мы можем управлять тем, какой функционал ошибки мы оптимизируем, передавая разные классы в параметр loss при инициализации.

Будем применять нашу регрессию на реальном датафрейме. Загрузим датафрейм с машинами (см. семинар)

```

B [41]: 1 X_raw = pd.read_csv(
2         "http://archive.ics.uci.edu/ml/machine-learning-databases/autos/imports-85.data",
3         header=None,
4         na_values=["?"],
5     )
6
7     X_raw.head()
8     X_raw = X_raw[~X_raw[25].isna()].reset_index(drop=True)
B [42]: 1 y = X_raw[25]
2     X_raw = X_raw.drop(25, axis=1)

```

Задание 2.2: Обработайте датасет нужными методами, чтобы на нём можно было обучать линейную регрессию (используя sklearn-knn.ipynb):

- Заполните пропуски средними (библиотека SimpleImputer)
- Переведите категориальные признаки в числовые (в методе get_dummies использовать drop_first=True.)

- Разделите датасет на обучающую и тестовую выборку (задать: доля тестовой выборки равна 0.3, `random shuffle=True`)
- Нормализуйте числовые признаки (при помощи библиотеки `StandardScaler`)

B [43]: 1 `X_raw.isna().any()` *# сначала проверим есть ли пропуски вообще*

```
Out[43]: 0    False
         1     True
         2    False
         3    False
         4    False
         5     True
         6    False
         7    False
         8    False
         9    False
        10    False
        11    False
        12    False
        13    False
        14    False
        15    False
        16    False
        17    False
        18     True
        19     True
        20    False
        21     True
        22     True
        23    False
        24    False
dtype: bool
```

B [44]: 1 `cat_features_mask = (X_raw.dtypes == "object").values` *# делаем маску для определения с*
 2 `X_real = X_raw[X_raw.columns[~cat_features_mask]]` *# отбираем столбцы с числовыми д*

B [45]: 1 `mis_replacer = SimpleImputer(strategy="mean")` *# заводим переменную реплэйсера*
 2 `data = mis_replacer.fit_transform(X_real)` *# заполняем пропуски*
 3 `X_no_mis_real = pd.DataFrame(data, columns=X_real.columns)` *# возвращаем все к виду дат*

B [46]: 1 `X_cat = X_raw[X_raw.columns[cat_features_mask]].fillna("")` *# в остальных столбцах в*
 2 `X_no_mis = pd.concat([X_no_mis_real, X_cat],axis=1)` *# соединяем таблички*
 3 `X_no_mis.head(10)`

```
Out[46]:
```

	0	1	9	10	11	12	13	16	18	19	...	2	3	4	5	6	7	
0	3.0	122.0	88.6	168.8	64.1	48.8	2548.0	130.0	3.47	2.68	...	alfa-romero	gas	std	two	convertible	rwd	f
1	3.0	122.0	88.6	168.8	64.1	48.8	2548.0	130.0	3.47	2.68	...	alfa-romero	gas	std	two	convertible	rwd	f
2	1.0	122.0	94.5	171.2	65.5	52.4	2823.0	152.0	2.68	3.47	...	alfa-romero	gas	std	two	hatchback	rwd	f
3	2.0	164.0	99.8	176.6	66.2	54.3	2337.0	109.0	3.19	3.40	...	audi	gas	std	four	sedan	fwd	f
4	2.0	164.0	99.4	176.6	66.4	54.3	2824.0	136.0	3.19	3.40	...	audi	gas	std	four	sedan	4wd	f
5	2.0	122.0	99.8	177.3	66.3	53.1	2507.0	136.0	3.19	3.40	...	audi	gas	std	two	sedan	fwd	f
6	1.0	158.0	105.8	192.7	71.4	55.7	2844.0	136.0	3.19	3.40	...	audi	gas	std	four	sedan	fwd	f
7	1.0	122.0	105.8	192.7	71.4	55.7	2954.0	136.0	3.19	3.40	...	audi	gas	std	four	wagon	fwd	f
8	1.0	158.0	105.8	192.7	71.4	55.9	3086.0	131.0	3.13	3.40	...	audi	gas	turbo	four	sedan	fwd	f
9	2.0	192.0	101.2	176.8	64.8	54.3	2395.0	108.0	3.50	2.80	...	bmw	gas	std	two	sedan	rwd	f

10 rows × 25 columns

```
B [47]: 1 X_dum = pd.get_dummies(X_no_mis, drop_first=True) # кодируем категориальные признаки one
        2 X_dum.head(10)
```

```
Out[47]:
```

	0	1	9	10	11	12	13	16	18	19	...	15_three	15_twelve	15_two	17_2bbl	17_4bbl
0	3.0	122.0	88.6	168.8	64.1	48.8	2548.0	130.0	3.47	2.68	...	0	0	0	0	0
1	3.0	122.0	88.6	168.8	64.1	48.8	2548.0	130.0	3.47	2.68	...	0	0	0	0	0
2	1.0	122.0	94.5	171.2	65.5	52.4	2823.0	152.0	2.68	3.47	...	0	0	0	0	0
3	2.0	164.0	99.8	176.6	66.2	54.3	2337.0	109.0	3.19	3.40	...	0	0	0	0	0
4	2.0	164.0	99.4	176.6	66.4	54.3	2824.0	136.0	3.19	3.40	...	0	0	0	0	0
5	2.0	122.0	99.8	177.3	66.3	53.1	2507.0	136.0	3.19	3.40	...	0	0	0	0	0
6	1.0	158.0	105.8	192.7	71.4	55.7	2844.0	136.0	3.19	3.40	...	0	0	0	0	0
7	1.0	122.0	105.8	192.7	71.4	55.7	2954.0	136.0	3.19	3.40	...	0	0	0	0	0
8	1.0	158.0	105.8	192.7	71.4	55.9	3086.0	131.0	3.13	3.40	...	0	0	0	0	0
9	2.0	192.0	101.2	176.8	64.8	54.3	2395.0	108.0	3.50	2.80	...	0	0	0	0	0

10 rows × 65 columns

```
B [48]: 1 # при нормализации ошибку выдало, что в столбцах есть числовые названия, так что нужно пе
        2 X_dum.columns = [str(i) for i in X_dum.columns]
```

```
B [49]: 1 X_train, X_test, y_train, y_test = train_test_split(X_dum, y, random_state=42, test_size=
```

```
B [50]: 1 scaler = StandardScaler() # создаем scaler
        2 scaler.fit(X_train) # обучаем его
        3 scaled_X_train = scaler.transform(X_train) # нормализуем X_train
        4 scaled_X_test = scaler.transform(X_test) # нормализуем X_test
```

Задание 2.3: Обучите написанную вами линейную регрессию на обучающей выборке

Создаем объект линейной регрессии для MSELoss :

```
B [51]: 1 lr_1 = LinearRegression_1(MSELoss()) #создаем регрессор
```

```
B [52]: 1 lr_1.fit(scaled_X_train, y_train) # обучаем модель
        2 y_pred = lr_1.predict(scaled_X_test) # делаем предикты
```

Задание 2.4: Посчитайте ошибку обученной регрессии на обучающей и тестовой выборке при помощи метода r2_score из sklearn.metrics .

```
B [53]: 1 from sklearn.metrics import mean_squared_error
        2 from sklearn.metrics import r2_score
```

```
B [54]: 1 print(f"MSE: {mean_squared_error(y_test, y_pred)*0.5}") #смотрим среднеквадратичн
        2
        3 print('Коэффициент детерминации: ', r2_score(y_test, y_pred))
```

MSE: 2517.9510100373222

Коэффициент детерминации: 0.9331156622103828

Добавим к модели L2 регуляризацию (для борьбы с переобучением). Для этого нам нужно создать новый класс ошибки и его градиента.

Формула функционала ошибки для MSE с L2 регуляризацией выглядит так:

$$Q(w, X, y) = \frac{1}{\ell} \sum_{i=1}^{\ell} (\langle x_i, w \rangle - y_i)^2 + \lambda \|w\|^2$$

Или в матричном виде:

$$Q(w, X, y) = \frac{1}{\ell} \|Xw - y\|^2 + \lambda \|w\|^2,$$

где λ — коэффициент регуляризации

Заметим, что (удобно для вычислений):

$$Q(w, X, y) = \frac{1}{\ell} \|Xw - y\|^2 + \lambda \|w\|^2 = \frac{1}{\ell} (y - Xw)^T (y - Xw) + \lambda w^T w.$$

Градиент $\nabla_w Q(w, X, y)$ выглядит так:

$$\nabla_w Q(w, X, y) = 2 \cdot X^T (Xw - y)$$

Задание 2.5: Реализуйте класс `MSEL2Loss`

Он должен вычислять значение функционала ошибки (лосс) $Q(w, X, y)$ и его градиент $\nabla_w Q(w, X, y)$ по фо

Подсказка: обратите внимание, что последний элемент вектора `w` — это bias (сдвиг) (в классе `LinearRegre` добавляется колонка из единиц — константный признак). bias регуляризовать не нужно. Поэтому не забудьте из `w` при подсчёте слагаемого $\lambda \|w\|^2$ в `calc_loss` и занулить его при подсчёте слагаемого $2\lambda w$ в `calc_g`

```
B [78]: 1 class MSEL2Loss(BaseLoss):
2         def __init__(self, coef: float = 1.0):
3             self.coef = coef
4
5         def calc_loss(self, X: np.ndarray, y: np.ndarray, w: np.ndarray) -> float:
6             return (np.linalg.norm(np.matmul(X, w) - y)**2) / len(y) + self.coef * (np.linalg.norm(w
7
8         def calc_grad(self, X: np.ndarray, y: np.ndarray, w: np.ndarray) -> np.ndarray:
9             w_new = np.append(w[:-1], [0])
10            return 2 * np.matmul(np.transpose(X), (np.matmul(X, w) - y)) / len(y) + 2 * self.coef * w_new
```

```
B [79]: 1 # проверка:
2
3 # создадим объект лосса
4 losst = MSEL2Loss()
5
6 # создадим датасет
7 Xt = np.arange(200).reshape(20, 10)
8 yt = np.arange(20)
9
10 # создадим вектор весов
11 wt = np.arange(10)
12
13 #print(Xt)
14 #print(yt)
15 #print(wt)
16
17 # выведем значение лосса и градиента на этом датасете с этим вектором весов
18 print(losst.calc_loss(Xt, yt, wt))
19 print(losst.calc_grad(Xt, yt, wt))
20
21 # проверка, что методы реализованы правильно
22 assert losst.calc_loss(Xt, yt, wt) == 27410487.5, "Метод calc_loss реализован неверно"
23 assert np.allclose(
24     losst.calc_grad(Xt, yt, wt),
25     np.array(
26         [
27             1163180.0,
28             1172283.0,
29             1181386.0,
30             1190489.0,
31             1199592.0,
32             1208695.0,
33             1217798.0,
34             1226901.0,
35             1236004.0,
36             1245089.0,
37         ]
38     ),
39     "Метод calc_grad реализован неверно"
40 )
41 print("Всё верно!")
```

```
27410487.5
[1163180. 1172283. 1181386. 1190489. 1199592. 1208695. 1217798. 1226901.
 1236004. 1245089.]
Всё верно!
```

Теперь мы можем использовать лосс с l2 регуляризацией в нашей регрессии. Пусть:

```
B [57]: 1 lr_1_l2 = LinearRegression_1(MSEL2Loss(0.1))
```

Задание 2.6: Обучите регрессию с лоссом `MSEL2Loss`. Попробуйте использовать другие коэффициенты регуляризации на тестовой выборке? Сравните результат применения регрессии с регуляризацией к тем же выборкам с результатом применения регрессии без регуляризации к тем же выборкам. (Для оценки качества `mean_squared_error` и `r2_score` из `sklearn.metrics`).

```
B [72]: 1 lr_1_l2.fit(scaled_X_train, y_train)
2 y_predict = lr_1_l2.predict(scaled_X_test)
3 y_predict_train = lr_1_l2.predict(scaled_X_train)
4 print(f"MSE on Test: {mean_squared_error(y_test, y_predict)**0.5}")
5 print(f"MSE on Train: {mean_squared_error(y_train, y_predict_train)**0.5}")
6 print('Коэффициент детерминации на тестовой выборке: ', r2_score(y_test, y_predict))
7 print('Коэффициент детерминации на тренировочной выборке', r2_score(y_train, y_predict_train))
```

MSE on Test: 2518.250194286107
MSE on Train: 1393.7382610784412
Коэффициент детерминации на тестовой выборке: 0.9330997668024732
Коэффициент детерминации на тренировочной выборке 0.9596828392692338

```
B [73]: 1 lr_1_l2 = LinearRegression_1(MSEL2Loss(0.01))
2 lr_1_l2.fit(scaled_X_train, y_train)
3 y_predict = lr_1_l2.predict(scaled_X_test)
4 y_predict_train = lr_1_l2.predict(scaled_X_train)
5 print(f"MSE on Test: {mean_squared_error(y_test, y_predict)**0.5}")
6 print(f"MSE on Train: {mean_squared_error(y_train, y_predict_train)**0.5}")
7 print('Коэффициент детерминации на тестовой выборке: ', r2_score(y_test, y_predict))
8 print('Коэффициент детерминации на тренировочной выборке', r2_score(y_train, y_predict_train))
```

MSE on Test: 2525.426536883011
MSE on Train: 1408.5899787098672
Коэффициент детерминации на тестовой выборке: 0.9327179278005473
Коэффициент детерминации на тренировочной выборке 0.9588190194146218

```
B [74]: 1 lr_1_l2 = LinearRegression_1(MSEL2Loss(0.0001))
2 lr_1_l2.fit(scaled_X_train, y_train)
3 y_predict = lr_1_l2.predict(scaled_X_test)
4 y_predict_train = lr_1_l2.predict(scaled_X_train)
5 print(f"MSE on Test: {mean_squared_error(y_test, y_predict)**0.5}")
6 print(f"MSE on Train: {mean_squared_error(y_train, y_predict_train)**0.5}")
7 print('Коэффициент детерминации на тестовой выборке: ', r2_score(y_test, y_predict))
8 print('Коэффициент детерминации на тренировочной выборке', r2_score(y_train, y_predict_train))
```

MSE on Test: 2518.5787904710046
MSE on Train: 1393.8114918752556
Коэффициент детерминации на тестовой выборке: 0.9330823065868887
Коэффициент детерминации на тренировочной выборке 0.959678602411492

```
B [75]: 1 lr_1_l2 = LinearRegression_1(MSEL2Loss(0.5))
2 lr_1_l2.fit(scaled_X_train, y_train)
3 y_predict = lr_1_l2.predict(scaled_X_test)
4 y_predict_train = lr_1_l2.predict(scaled_X_train)
5 print(f"MSE on Test: {mean_squared_error(y_test, y_predict)**0.5}")
6 print(f"MSE on Train: {mean_squared_error(y_train, y_predict_train)**0.5}")
7 print('Коэффициент детерминации на тестовой выборке: ', r2_score(y_test, y_predict))
8 print('Коэффициент детерминации на тренировочной выборке', r2_score(y_train, y_predict_train))
```

MSE on Test: 3197.41511786243
MSE on Train: 1898.5504183773194
Коэффициент детерминации на тестовой выборке: 0.8921480684437594
Коэффициент детерминации на тренировочной выборке 0.9251879139460245

```
B [76]: 1 lr_1_l2 = LinearRegression_1(MSEL2Loss(0.9))
2 lr_1_l2.fit(scaled_X_train, y_train)
3 y_predict = lr_1_l2.predict(scaled_X_test)
4 y_predict_train = lr_1_l2.predict(scaled_X_train)
5 print(f"MSE on Test: {mean_squared_error(y_test, y_predict)**0.5}")
6 print(f"MSE on Train: {mean_squared_error(y_train, y_predict_train)**0.5}")
7 print('Коэффициент детерминации на тестовой выборке: ', r2_score(y_test, y_predict))
8 print('Коэффициент детерминации на тренировочной выборке', r2_score(y_train, y_predict_train))
```

MSE on Test: 3562.7185617962423
MSE on Train: 2117.8816804070384
Коэффициент детерминации на тестовой выборке: 0.8660961967849566
Коэффициент детерминации на тренировочной выборке 0.9069040297555025

```
B [77]: 1 lr_1_l2 = LinearRegression_1(MSEL2Loss(0.00000001))
2 lr_1_l2.fit(scaled_X_train, y_train)
3 y_predict = lr_1_l2.predict(scaled_X_test)
4 y_predict_train = lr_1_l2.predict(scaled_X_train)
5 print(f"MSE on Test: {mean_squared_error(y_test, y_predict)**0.5}")
6 print(f"MSE on Train: {mean_squared_error(y_train, y_predict_train)**0.5}")
7 print('Коэффициент детерминации на тестовой выборке: ', r2_score(y_test, y_predict))
8 print('Коэффициент детерминации на тренировочной выборке', r2_score(y_train, y_predict_train))
```

MSE on Test: 2518.7941514179693
MSE on Train: 1393.8159172538788
Коэффициент детерминации на тестовой выборке: 0.9330708619785003
Коэффициент детерминации на тренировочной выборке 0.9596783463686409

Улучшить результат при изменении коэффициента не получилось. В каких-то случаях может падать значение коэффициента качестве всех значений падало