

# Семинар: Градиентный спуск.

```

В [6]: 1 from typing import Iterable, List
        2
        3 import matplotlib.pyplot as plt
        4 import numpy as np
        5
        6 %matplotlib inline

```

## Часть 1. Градиентный спуск

Функционал ошибки, который мы применяем в задаче регр 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 нашей линейной модели.

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

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

или

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

где  $X$  — это матрица объекты-признаки, а  $y$  — вектор прав

Для того чтобы воспользоваться методом градиентного спуска, нам нужен градиент нашего функционала. Для MSE он будет выглядеть

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

Ниже приведён базовый класс `BaseLoss`, который мы будем использовать для реализации всех наших функционалов ошибки (= функций потерь). У него есть два абстрактных метода:

1. Метод `calc_loss`, который будет принимать на вход ответы  $y$  и веса  $w$  и вычислять значения функционала ошибки.
2. Метод `calc_grad`, который будет принимать на вход ответы  $y$  и веса  $w$  и вычислять значения градиента функционала.

```

B [7]: 1 import abc
        2
        3
        4 class BaseLoss(abc.ABC):
        5
        6
        7     @abc.abstractmethod
        8     def calc_loss(self, X: np.ndarray, y: np.n
        9
        10         raise NotImplementedError
        11
        12     @abc.abstractmethod
        13     def calc_grad(self, X: np.ndarray, y: np.n
        14         raise NotImplementedError

```

Реализация этого абстрактного класса: Mean Squared Error

Данный код определяет BaseLoss класс с использованием модель линейной регрессии, обучает ее на масштабированном датасете, предсказывает значения  $y_{train}$  и  $y_{test}$  с помощью выводит значение средней квадратичной ошибки (MSE) и коэффициент детерминации ( $R^2$ ) для train и test датасетов, чтобы проверить. Метод `calc_loss()` вычисляет значение потерь, а `calc_grad()` функции потерь

### Задание 1.1: Реализуйте класс MSELoss

Он должен вычислять значение функционала ошибки (потери)  $\nabla_w Q(w, X, y)$  по формулам (выше)

```

B [8]: 1 class MSELoss(BaseLoss):
        2     def calc_loss(self, X: np.ndarray, y: np.n
        3         return (np.linalg.norm(np.matmul(X,w)
        4
        5     def calc_grad(self, X: np.ndarray, y: np.n
        6         return 2/len(y)*np.matmul(np.transpose

```

Данный определяет класс с именем MSELoss, который реализует MSE. Он состоит из двух методов: `calc_loss` и `calc_grad`.

Метод `calc_loss` вычисляет потери MSE с учетом входных данных  $X$  и  $y$  и весов модели  $w$ . Он использует формулу для вычисления разницы между прогнозируемыми и фактическими значениями:  $\sum (y_i - \hat{y}_i)^2$ , где  $\hat{y}_i$  — предсказанные значения. Он делит сумму квадратов на количество точек данных, чтобы получить среднюю квадратичную ошибку.

Метод `calc_grad` вычисляет градиент потерь MSE по отношению к весам  $w$ . Он принимает те же входные параметры, что и `calc_loss`, и возвращает вектор градиента. Он представляет направление и величину изменения весов для минимизации потерь.

Теперь мы можем создать объект `MSELoss` и при помощи нашего функционала ошибки и градиента:

```
B [9]: 1 loss = MSELoss()
2
3 X = np.arange(200).reshape(20, 10)
4 y = np.arange(20)
5
6 w = np.arange(10)
7
8 print('X:\n', X)
9 print('y:\n', y)
10 print('w:\n', w, '\n')
11
12
13 print(loss.calc_loss(X, y, w))
14 print(loss.calc_grad(X, y, w))
15
16
17 assert loss.calc_loss(X, y, w) == 27410283.5,
18 assert np.allclose(
19     loss.calc_grad(X, y, w),
20     np.array(
21         [
22             1163180.0,
23             1172281.0,
24             1181382.0,
25             1190483.0,
26             1199584.0,
27             1208685.0,
28             1217786.0,
29             1226887.0,
30             1235988.0,
31             1245089.0,
32         ]
33     ),
34     ), "Метод calc_grad реализован неверно"
35 print("Всё верно!")
```

```

X:
[[ 0  1  2  3  4  5  6  7  8  9]
 [10 11 12 13 14 15 16 17 18 19]
 [20 21 22 23 24 25 26 27 28 29]
 [30 31 32 33 34 35 36 37 38 39]
 [40 41 42 43 44 45 46 47 48 49]
 [50 51 52 53 54 55 56 57 58 59]
 [60 61 62 63 64 65 66 67 68 69]
 [70 71 72 73 74 75 76 77 78 79]
 [80 81 82 83 84 85 86 87 88 89]
 [90 91 92 93 94 95 96 97 98 99]
[100 101 102 103 104 105 106 107 108 109]
[110 111 112 113 114 115 116 117 118 119]
[120 121 122 123 124 125 126 127 128 129]
[130 131 132 133 134 135 136 137 138 139]
[140 141 142 143 144 145 146 147 148 149]
[150 151 152 153 154 155 156 157 158 159]
[160 161 162 163 164 165 166 167 168 169]
[170 171 172 173 174 175 176 177 178 179]
[180 181 182 183 184 185 186 187 188 189]
[190 191 192 193 194 195 196 197 198 199]]

y:
[ 0  1  2  3  4  5  6  7  8  9 10 11 12 13 14 15]

w:
[0 1 2 3 4 5 6 7 8 9]

```

```

27410283.5
[1163180. 1172281. 1181382. 1190483. 1199584. 1208
 1235988. 1245089.]
Всё верно!

```

Сначала я импортирую необходимые библиотеки для числ  
определяю класс с именем `MSELoss`, который содержит де  
`calc_grad`, описанные выше.

Метод `calc_loss` принимает на вход три параметра: `X`, `y` и `w`

- `X` - входные данные;
- `y` - фактические выходные данные;
- `w` - веса.

Внутри метода я вычисляю прогнозируемый результат (`y_h`  
входных данных (`X`) на веса (`w`). Затем вычисляю потери `M`  
между фактическим результатом (`y`) и прогнозируемым рез  
его по количеству выборок (`n`).

Метод `calc_grad` также принимает те же параметры, что и с  
градиент потерь `MSE` по отношению к весам (`w`). Градиент  
использованием производной формулы потерь `MSE`. Отри  
гарантирует, что градиент указывает в направлении уменьш

Наконец, мы создаем экземпляр `MSELoss` класса и демонс  
передавая образцы данных (`X`, `y` и `w`) методам `calc_loss` an  
выводятся на консоль

Для вычисления градиента все готово, реализуем градиент  
формула для одной итерации градиентного спуска выгляди

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

Где  $w^t$  — значение вектора весов на  $t$ -ой итерации, а  $\eta$  — отвечающий за размер шага.

### Задание 1.2: Реализуйте функцию `gradient_descent`

Функция должна принимать на вход начальное значение весов `w_init`, матрицу объектов-признаков `X`, вектор правил функции потерь `loss`, размер шага `lr` и количество итераций `n_iterations`.

Функция должна реализовывать цикл, в котором осуществлен спуск (градиенты берутся из `loss` посредством вызова метода `calc_grad` формуле выше и возвращать траекторию спуска (список значений параметров на каждом шаге).

```

В [10]: 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 = 100000,
        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
       14

```

Функция `gradient_descent` принимает следующие параметры:

- `w_init` - начальные значения параметров.
- `X` - входные данные.
- `y` - целевые значения.
- `loss` - экземпляр класса функции потерь, который вычисляет значения потерь.
- `lr` - скорость обучения.
- `n_iterations` - количество итераций (по умолчанию установлено 100000).

Функция инициализирует пустой список `grad_list` для сохранения значений параметров на каждой итерации. Затем она выполняет итерации `n_iterations`, обновляя значения параметров с помощью правила обновления:

`w_init = w_init - lr * loss.calc_grad(X, y, w_init)` На каждой итерации значения параметров `w_init` добавляются к `grad_list`. Наконец, функция возвращает список значений параметров.

Мы определяем простую задачу линейной регрессии с входными данными `X` и целевыми значениями `y`. Мы инициализируем значения параметров `w_init` нулями, создаем экземпляр `BaseLoss` класса, который вычисляет значения потерь, устанавливаем скорость обучения на `MeanSquaredError` и `lr=0.01`, количество итераций `n_iterations=1000`.

Затем мы вызываем `gradient_descent` функцию, передавая параметров, входные данные, целевые значения, функций и количество итераций. Функция возвращает список, содер

Теперь создадим синтетический датасет и функцию, которая траекторию градиентного спуска по истории.

(Если не оговорено иное, то в задачах используются указанные `n_features`, `n_objects`, `batch_size`, `num_steps`).

```
In [11]: 1 np.random.seed(1337)
          2
          3 n_features = 2
          4 n_objects = 300
          5 batch_size = 10
          6 num_steps = 43
          7
          8 w_true = np.random.normal(size=(n_features,))
          9
         10 X = np.random.uniform(-5, 5, (n_objects, n_fea
         11 X *= (np.arange(n_features) * 2 + 1)[np.newaxis]
         12 y = X.dot(w_true) + np.random.normal(0, 1, (n_
         13 w_init = np.random.uniform(-2, 2, (n_features)
         14
         15 print(X.shape)
         16 print(y.shape)
```

```
(300, 2)
```

```
(300,)
```

Предоставленный код реализует алгоритм градиентного спуска для регрессии. Давайте разберем структуру кода:

Функция `gradient_descent` принимает следующие параметры:

- `w_init`: Начальные веса для модели линейной регрессии.
- `X`: Матрица входных данных (независимые переменные).
- `y`: Целевая переменная (зависимая переменная).
- `loss`: Объект функции потерь, который вычисляет градиент.
- `lr`: Скорость обучения, которая определяет размер шага.
- `n_iterations`: Количество итераций для алгоритма (по умолчанию значение 100 000). Программа инициализирует необходимые массивы для алгоритма. Она генерирует случайные данные и устанавливает начальные веса `w_init` и печатает форму

Цикл градиентного спуска: основная часть кода - это `for` цикл, который выполняется `n_iterations` несколько раз. На каждой итерации он вычисляет потерю относительно весов и обновляет веса, используя скалярное произведение градиента и скорости обучения.

Возвращаемое значение - список массивов весов на каждой итерации. Это позволяет нам анализировать сходимость алгоритма.

```
B [12]: 1 loss = MSELoss()  
2 w_list = gradient_descent(w_init, X, y, loss,  
3 print(loss.calc_loss(X, y, w_list[0]))  
4 print(loss.calc_loss(X, y, w_list[-1]))
```

425.5891768045026

0.867074968324295

Инициализация потерь: мы создаем экземпляр функции по среднеквадратичной ошибке (MSE), используя MSELoss(). использоваться для оценки производительности нашей мо,

Градиентный спуск: мы вызываем gradient\_descent функции веса (w\_init), входные данные (X), целевые значения (y), скорость обучения (0.01) и количество шагов (num\_steps). : алгоритм градиентного спуска и возвращает список значен

Расчет потерь: Мы вычисляем потери, используя calc\_loss передаем входные данные (X), целевые значения (y) и пер веса из w\_list списка. Это позволяет нам сравнить начальн потерь.

```
B [13]: 1 w_init
```

Out[13]: array([0.62074297, 1.79288146])

```
B [14]: 1 w_list[-1]
```

Out[14]: array([-0.66688623, -0.48984218])

```

B [15]: 1 def plot_gd(w_list: Iterable, X: np.ndarray, y
2
3     w_list = np.array(w_list)
4     meshgrid_space = np.linspace(-2, 2, 100)
5     A, B = np.meshgrid(meshgrid_space, meshgrid
6
7     levels = np.empty_like(A)
8     for i in range(A.shape[0]):
9         for j in range(A.shape[1]):
10             w_tmp = np.array([A[i, j], B[i, j]
11                 levels[i, j] = loss.calc_loss(X, y
12
13     plt.figure(figsize=(15, 6))
14     plt.title("GD trajectory")
15     plt.xlabel(r"$w_1$")
16     plt.ylabel(r"$w_2$")
17     plt.xlim(w_list[:, 0].min() - 0.1, w_list[
18     plt.ylim(w_list[:, 1].min() - 0.1, w_list[
19     plt.gca().set_aspect("equal")
20
21     CS = plt.contour(
22         A, B, levels, levels=np.logspace(0, 1,
23     )
24     CB = plt.colorbar(CS, shrink=0.8, extend="
25
26
27     plt.scatter(w_list[:, 0], w_list[:, 1])
28     plt.plot(w_list[:, 0], w_list[:, 1])
29
30     plt.show()

```

Код состоит из функции с именем `plot_gd`, которая принимает

- `w_list` - параметр, представляющий список весовых векторов, где каждый вектор соответствует определенной итерации алгоритма.
- `x` и `y` - параметры, представляющие входные функции соответственно.
- `loss` - параметр, представляющий функцию потерь, используемую для вычисления потерь для каждого вектора веса.

Затем код переходит к созданию пространства сетки с помощью `np.linspace`. Это пространство сетки используется для создания сетки точек.

Далее используется вложенный цикл для вычисления потерь в сеточном пространстве. Потери вычисляются с использованием объекта `loss`.

После вычисления потерь для каждого балла код создает фигуру `plt` и устанавливает заголовок, метку по оси `x` и метку по оси `y` с помощью `plt.xlabel` и `plt.ylabel` соответственно.

Ограничения по осям `x` и `y` устанавливаются на основе минимальных и максимальных значений весовых векторов в `w_list`. Соотношение сторон устанавливается равным с помощью `plt.gca().set_aspect("equal")`.



Контурный график создается с помощью `plt.contour` визуализирует пространство `meshgrid`. `levels` Параметр определяет уровни контуры. `star` Параметр определяет цветовую карту, используемую для

К графику добавляется цветная полоса, использующая `plt.colorbar` для отображения диапазона значений потерь.

Весовые векторы в `w_list` отображаются в виде точек рассеяния. Связанные весовые векторы отображаются с помощью

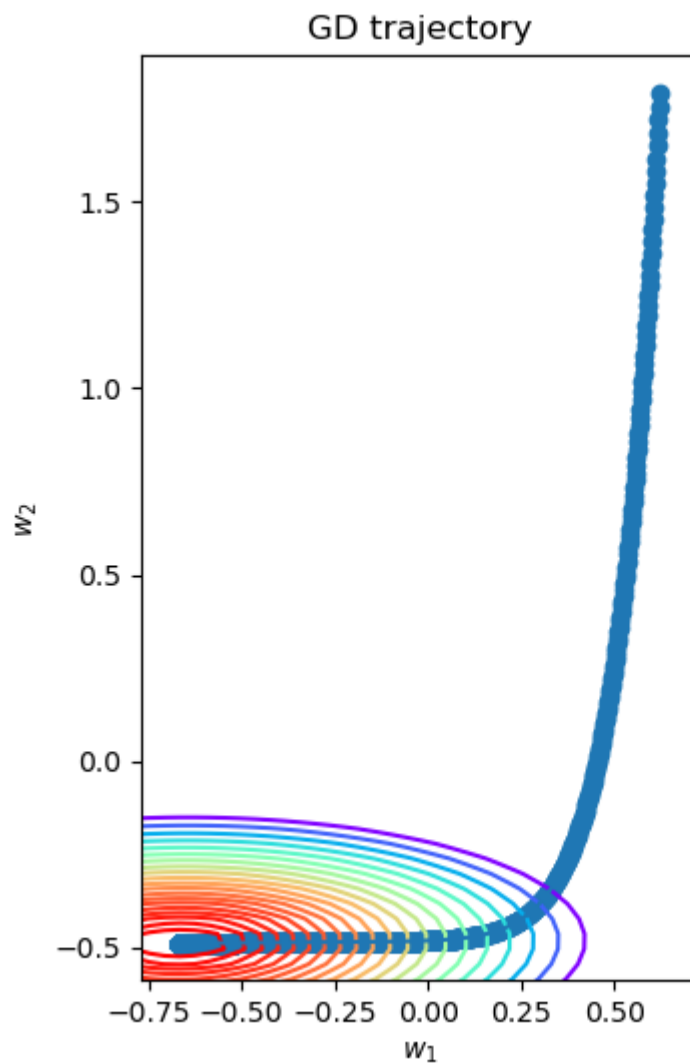
**Задание 1.3: При помощи функций `gradient_descent` и `plot_gd` построить траекторию градиентного спуска для разных значений `lr` ). Используйте не менее четырёх разных значений для `lr`. Вычисляйте значение функционала ошибки при помощи `cost` в первой и последней итерациях градиентного спуска.**

Сделайте и опишите свои выводы о том, как параметр `lr` влияет на траекторию градиентного спуска.

Подсказки:

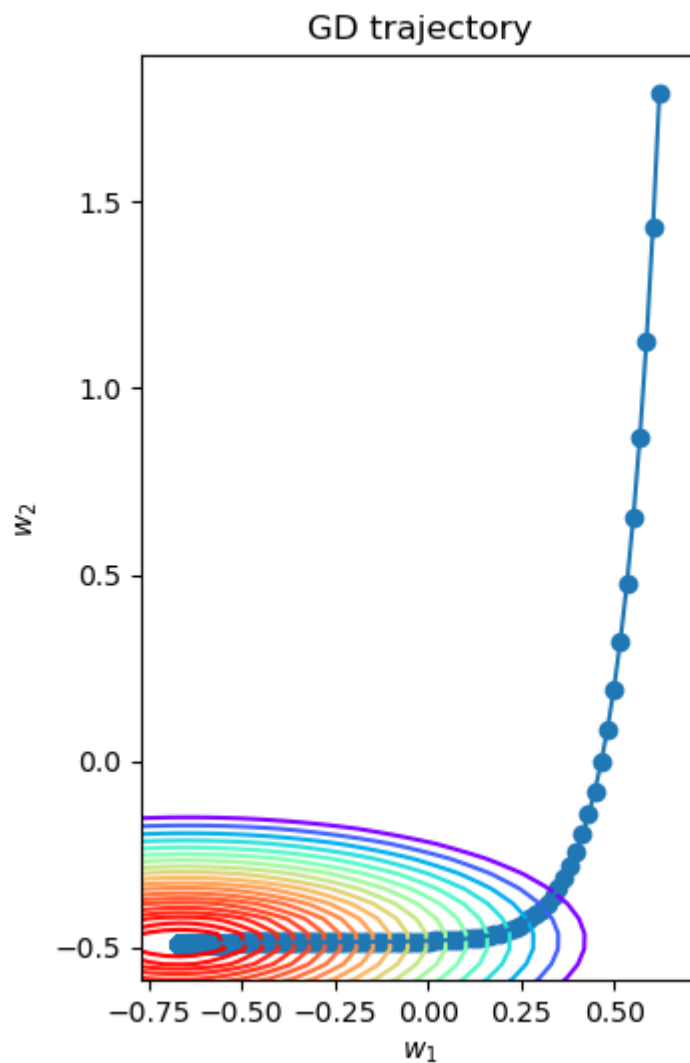
- Функция `gradient_descent` возвращает историю весовых векторов, которую можно передать в функцию `plot_gd`.
- Хорошие значения для `lr` могут лежать в промежутке `[0.001, 0.01]`.

```
B [16]: 1 w_list = gradient_descent(w_init, X, y, loss,
2 plot_gd(w_list, X, y, loss)
3 print('First iteration loss: ',loss.calc_loss(X
4 print('Last iteration loss: ',loss.calc_loss(X
```



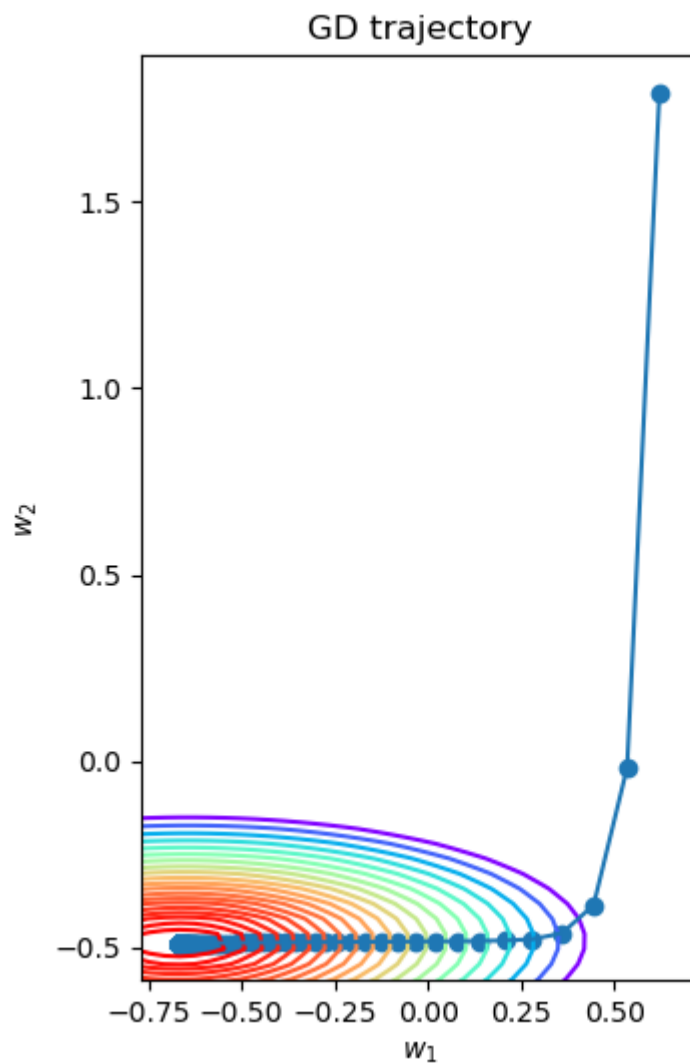
First iteration loss: 412.53549322933384  
Last iteration loss: 0.8670644395649092

```
B [17]: 1 w_list = gradient_descent(w_init, X, y, loss,
2 plot_gd(w_list, X, y, loss)
3 print('First iteration loss: ',loss.calc_loss(X
4 print('Last iteration loss: ',loss.calc_loss(X
```



First iteration loss: 304.4618116242326  
Last iteration loss: 0.8670644395649092

```
B [18]: 1 w_list = gradient_descent(w_init, X, y, loss,
2 plot_gd(w_list, X, y, loss)
3 print('First iteration loss: ',loss.calc_loss(X
4 print('Last iteration loss: ',loss.calc_loss(X
```

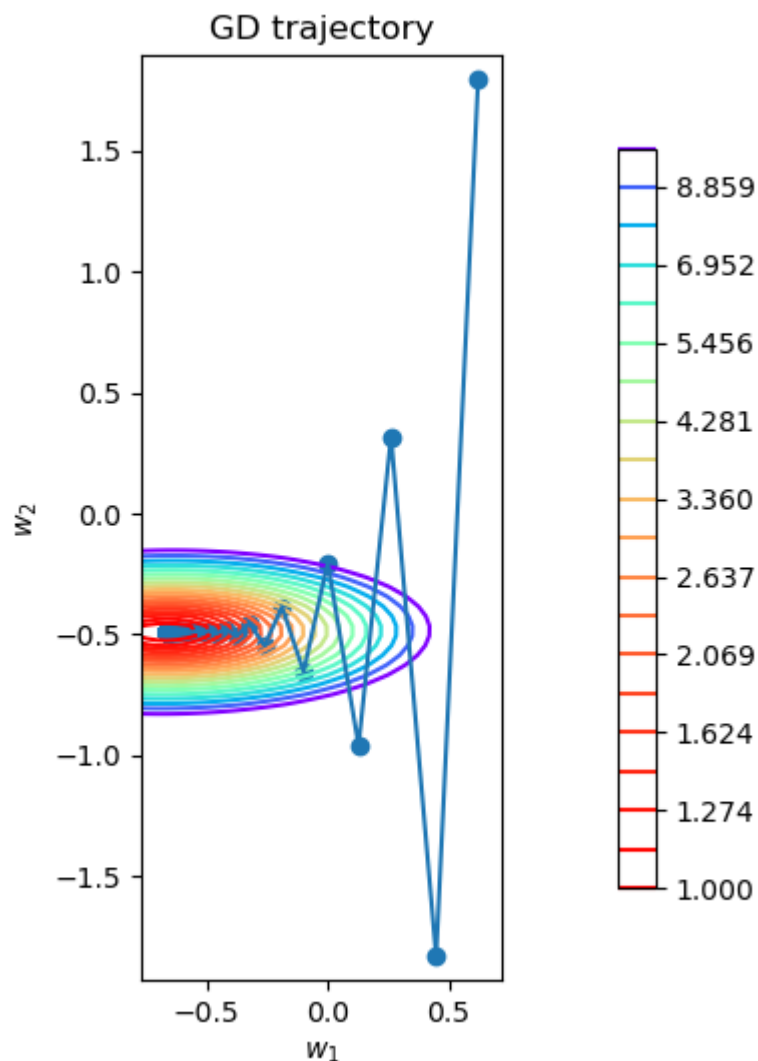


First iteration loss: 29.051696934645758  
Last iteration loss: 0.8670644395649092

```

B [19]: 1 w_list = gradient_descent(w_init, X, y, loss,
2         plot_gd(w_list, X, y, loss)
3         print('First iteration loss: ', loss.calc_loss(
4         print('Last iteration loss: ', loss.calc_loss(X

```



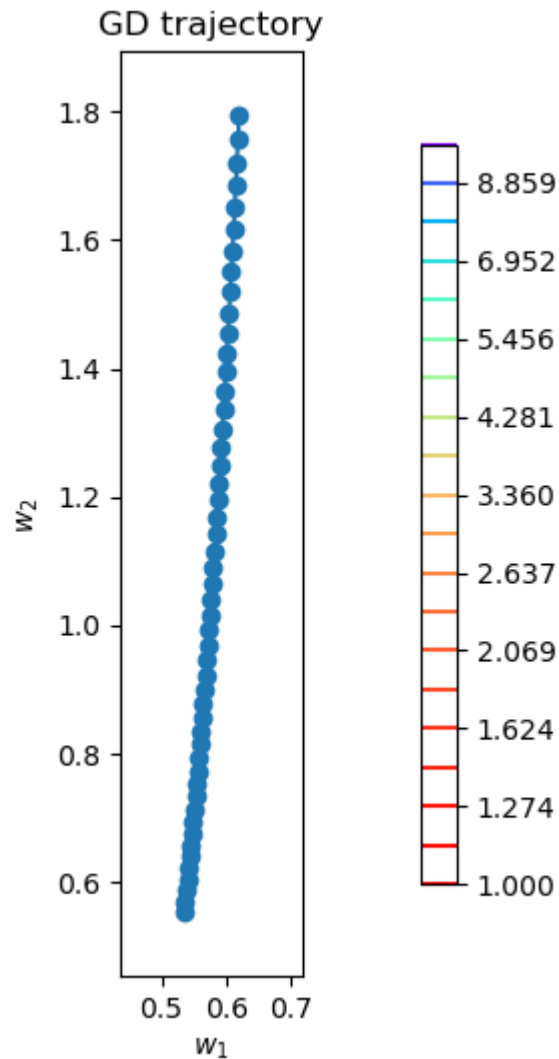
First iteration loss: 155.26258214352075  
 Last iteration loss: 0.8670644395649092

Различия минимальны. Причиной тому, по всей видимости, настройках функции по умолчанию задан слишком большой ( $n_{\text{iterations}} = 100000$ ). Для улучшения наглядности, значит показатель

```

B [22]: 1 w_list = gradient_descent(w_init, X, y, loss,
2         plot_gd(w_list, X, y, loss)
3         print('First iteration loss: ', loss.calc_loss(
4         print('Last iteration loss: ', loss.calc_loss(X

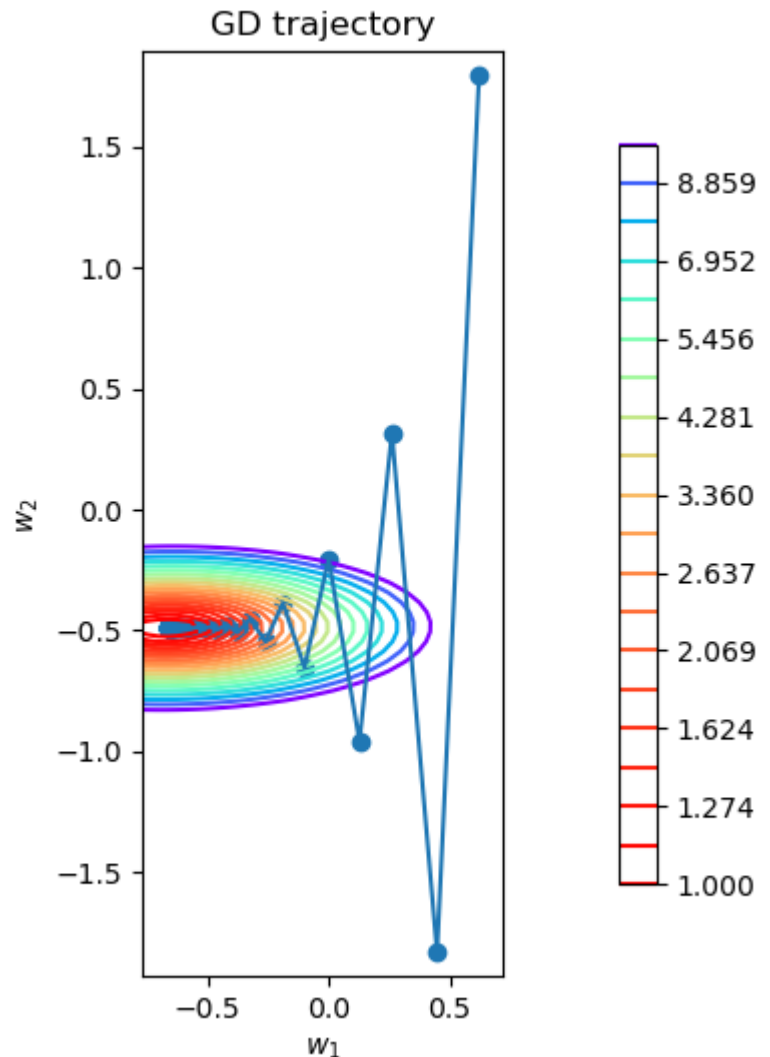
```



First iteration loss: 412.53549322933384  
 Last iteration loss: 97.42290021496474

Очевидно, что мы не достигли минимального значения, и з  
 возросло до 100

```
B [107]: 1 w_list = gradient_descent(w_init, X, y, loss,
2         plot_gd(w_list, X, y, loss)
3         print('First iteration loss: ',loss.calc_loss(
4         print('Last iteration loss: ',loss.calc_loss(X
```



First iteration loss: 155.26258214351958  
 Last iteration loss: 0.8670654544468456

На основании проведенного анализа можно сделать вывод  
 малом количестве шагов градиентный спуск может не дост

С увеличением числа шагов ситуация улучшается, и спуск  
 шаге 0,1 достигается значение потерь, аналогичное значен  
 шагов. На шаге размером 0,3 произошло нечто, приведшее  
 $1,642 \cdot 10^59$ .

#### Задание 1.4: Реализуйте функцию `stochastic_gradien`

Функция должна принимать все те же параметры, что и фу  
 но ещё параметр `batch_size`, отвечающий за размер бат

Функция должна как и раньше реализовывать цикл, в котор  
 градиентного спуска, но на каждом шаге считать градиент  
 только по случайно выбранной части.

Подсказка: для выбора случайной части можно использовать <https://numpy.org/doc/stable/reference/random/generated/numpy.random.choice> с правильным параметром `size`, чтобы выбрать случайные проиндексировать получившимся массивом массив `X`:

```
batch_indices = np.random.choice(X.shape[0], size=batch_size, replace=False)
batch = X[batch_indices]
```

(здесь `np.random.choice` генерирует из `np.arange(X.shape[0])` без повторения)

```
B [39]: 1 X.shape[0]
```

```
Out[39]: 300
```

```
B [40]: 1 def stochastic_gradient_descent(
2     w_init: np.ndarray,
3     X: np.ndarray,
4     y: np.ndarray,
5     loss: BaseLoss,
6     lr: float,
7     batch_size: int,
8     n_iterations: int = 1000,
9 ) -> List[np.ndarray]:
10     grad_list = []
11     batch_indices = np.random.choice(X.shape[0], batch_size, replace=False)
12     batch = X[batch_indices]
13     batch_y = y[batch_indices]
14
15     for i in range(n_iterations):
16         grad_list.append(w_init)
17         w_init = w_init - lr * loss.calc_grad(batch, batch_y)
18         grad_list.append(w_init)
19     return grad_list
20
21
```

Функция `stochastic_gradient_descent` принимает следующие

- `w_init` - начальные значения параметров модели.
- `X` - функции ввода.
- `y` - целевые значения.
- `loss` - экземпляр функции потерь, которая вычисляет градиент.
- `lr` - скорость обучения.
- `batch_size` - размер пакета, используемого для вычисления градиента.
- `n_iterations` - количество итераций (по умолчанию установлено 1000).

Функция инициализирует пустой список `grad_list` для хранения градиентов на каждой итерации. Затем она случайным образом выбирает пакет входных данных, `X` используя `np.random.choice`. Эти индексы используются для создания пакета входных функций `batch` и соответствующего целевого значения `batch_y`.



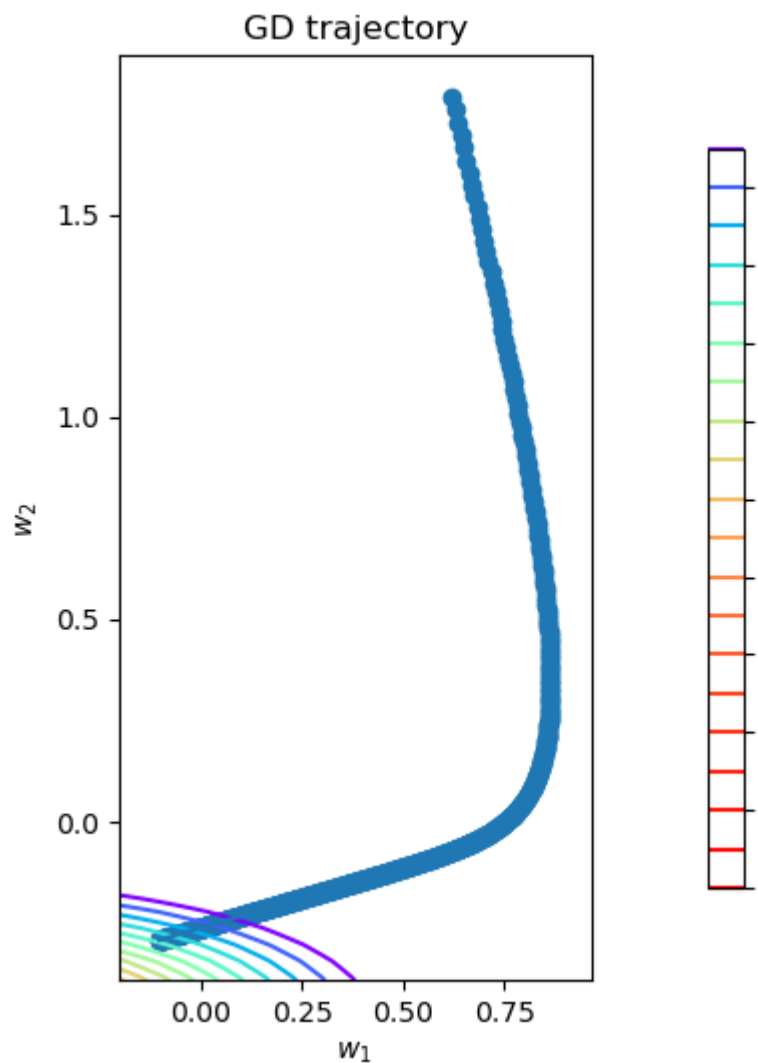
Затем функция входит в цикл, который повторяется  $n\_iterations$  итерации она добавляет текущие значения параметров  $w$  обновляются с помощью формулы,  $w_{init} = w_{init} - lr * loss(w_{init})$  где  $loss.calc\_grad$  вычисляется градиент функции по параметрам.

**Задание 1.5:** При помощи функций `stochastic_gradient_descent` нарисуйте траекторию градиентного спуска для разных значений `lr` и размера подвыборки (параметра `batch_size`). (Для каждого размера подвыборки вычисляйте значение функционала метода `calc_loss` на первой и последней итерациях стохастического спуска).

Сделайте и опишите свои выводы о том, как параметры `lr` и `batch_size` влияют на поведение стохастического градиентного спуска. Как отличается стохастический градиентный спуск от обычного?

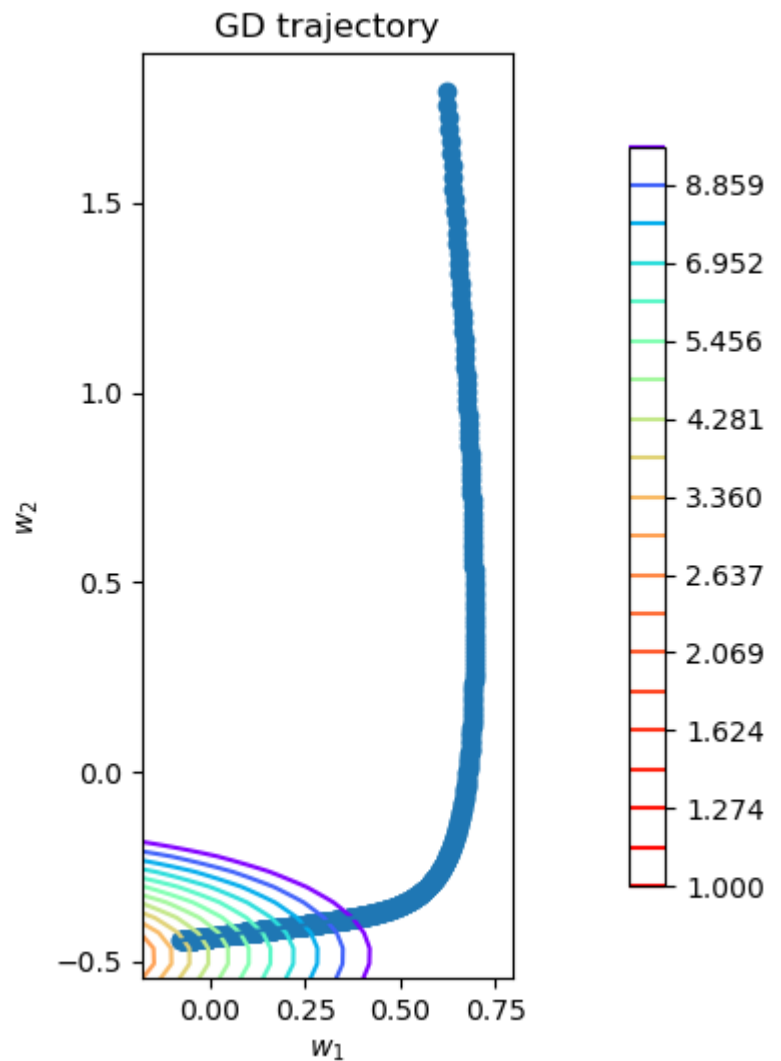
Обратите внимание, что в нашем датасете всего 300 объектов, поэтому больше этого числа не будет иметь смысла.

```
B [64]: 1 w_list = stochastic_gradient_descent(w_init, X
2 plot_gd(w_list, X, y, loss)
3 print('First iteration loss: ', loss.calc_loss(X
4 print('Last iteration loss: ', loss.calc_loss(X
```



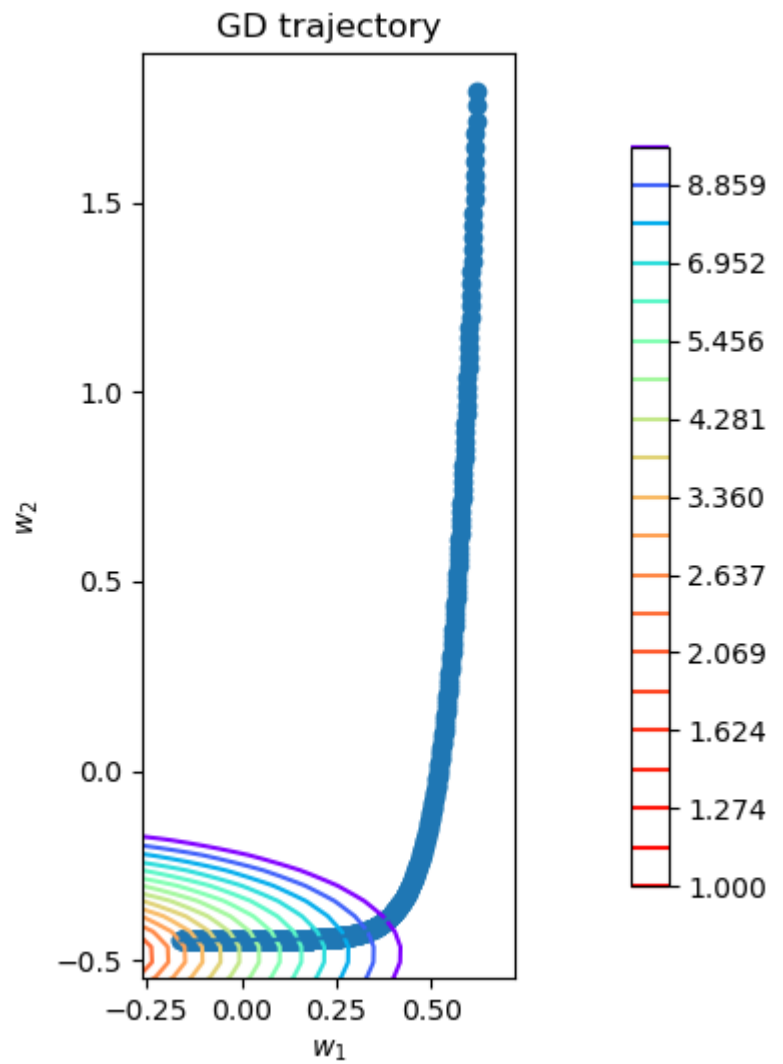
First iteration loss: 413.8806671894657  
Last iteration loss: 6.333267456575259

```
B [65]: 1 w_list = stochastic_gradient_descent(w_init, X
2 plot_gd(w_list, X, y, loss)
3 print('First iteration loss: ',loss.calc_loss(X
4 print('Last iteration loss: ',loss.calc_loss(X
```



First iteration loss: 413.51183006890705  
Last iteration loss: 3.6912358545991735

```
B [66]: 1 w_list = stochastic_gradient_descent(w_init, X
2 plot_gd(w_list, X, y, loss)
3 print('First iteration loss: ',loss.calc_loss(X
4 print('Last iteration loss: ',loss.calc_loss(X
```



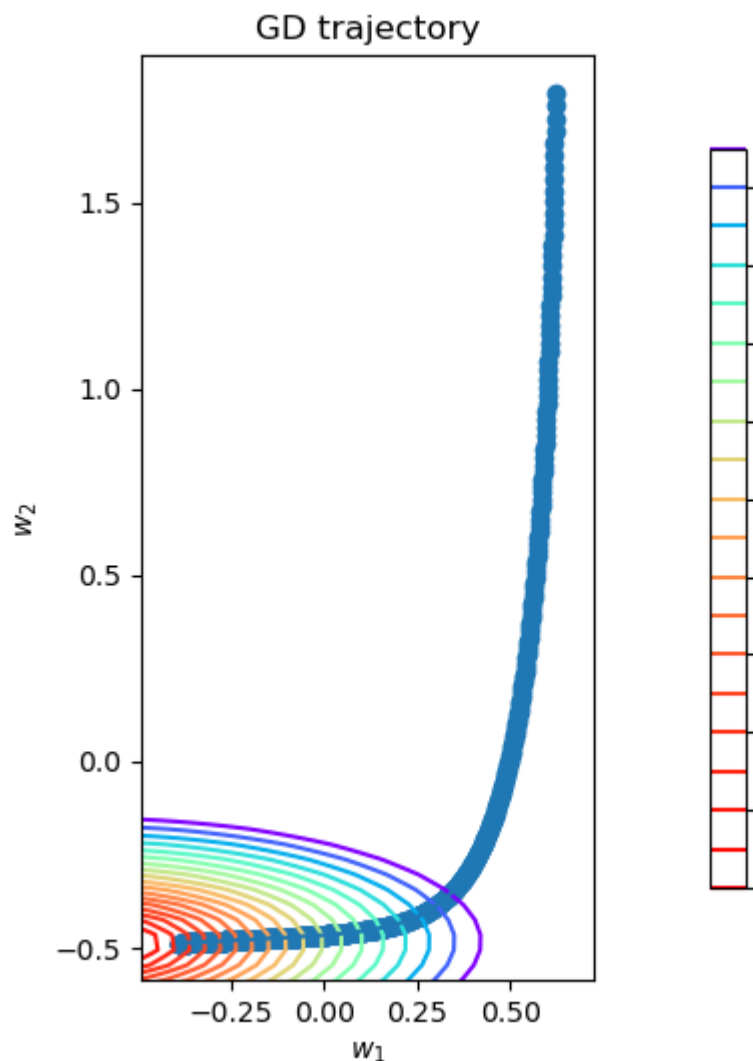
First iteration loss: 411.923969305594

Last iteration loss: 2.9742634047508036

```

B [68]: 1 w_list = stochastic_gradient_descent(w_init, X
2 plot_gd(w_list, X, y, loss)
3 print('First iteration loss: ', loss.calc_loss(
4 print('Last iteration loss: ', loss.calc_loss(X

```



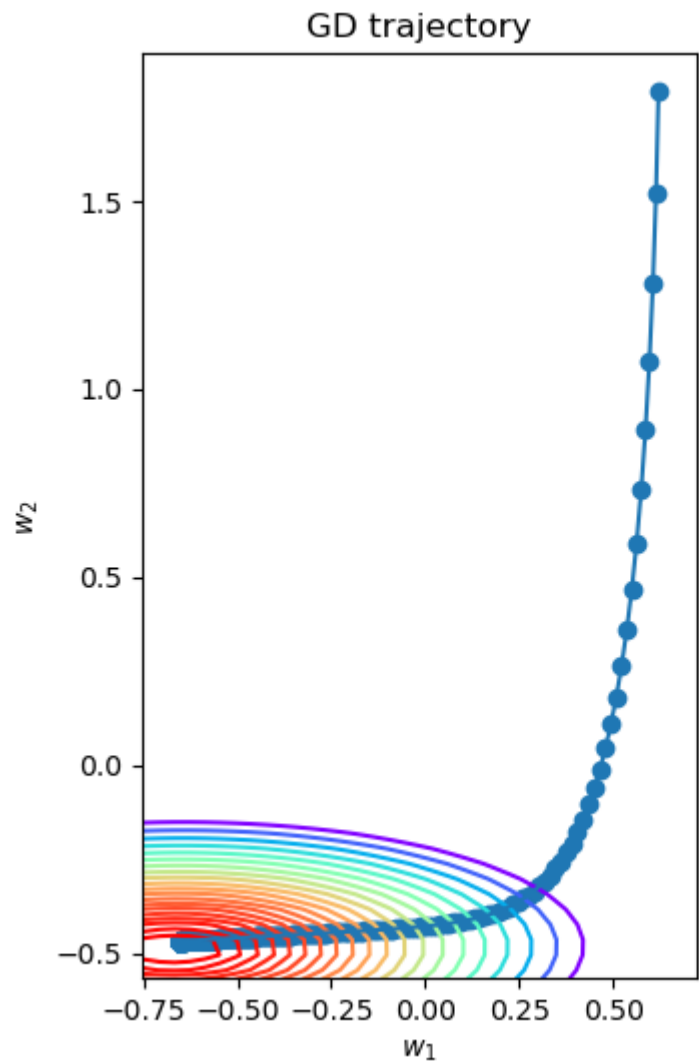
First iteration loss: 413.200675003692

Last iteration loss: 1.4725004316102306

С увеличением количества объектов точность модели пове

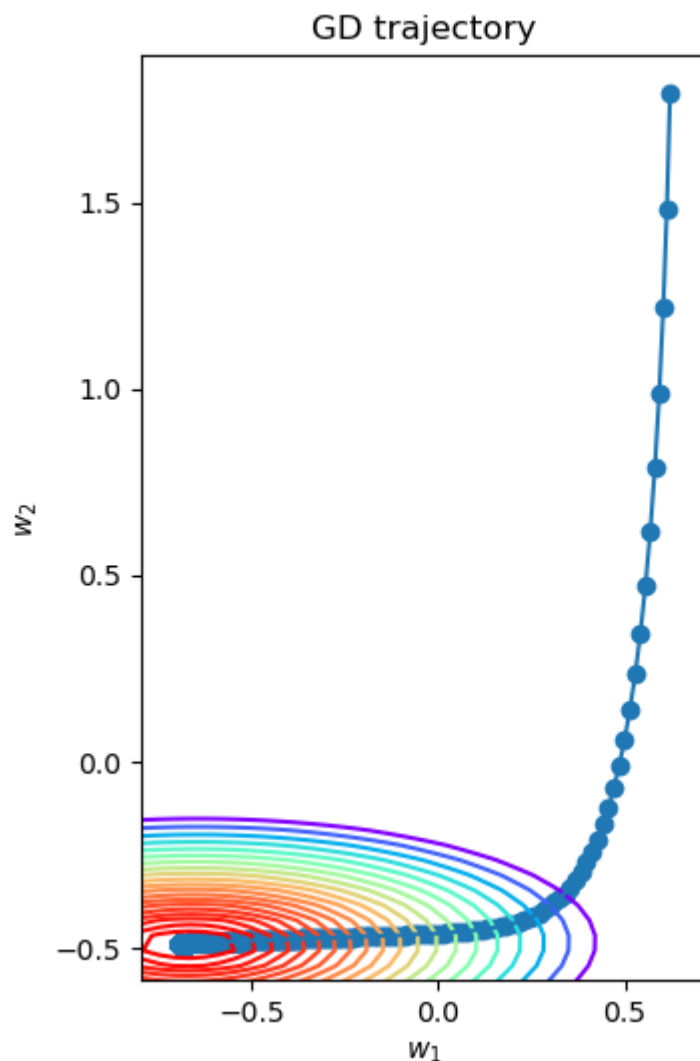
При 100 объектах она уже почти достигает минимума. Далк  
те же значения объектов и увеличить размер шага.

```
B [71]: 1 w_list = stochastic_gradient_descent(w_init, X
2 plot_gd(w_list, X, y, loss)
3 print('First iteration loss: ',loss.calc_loss(
4 print('Last iteration loss: ',loss.calc_loss(X
```



First iteration loss: 333.11166740747166  
Last iteration loss: 0.9072502252882477

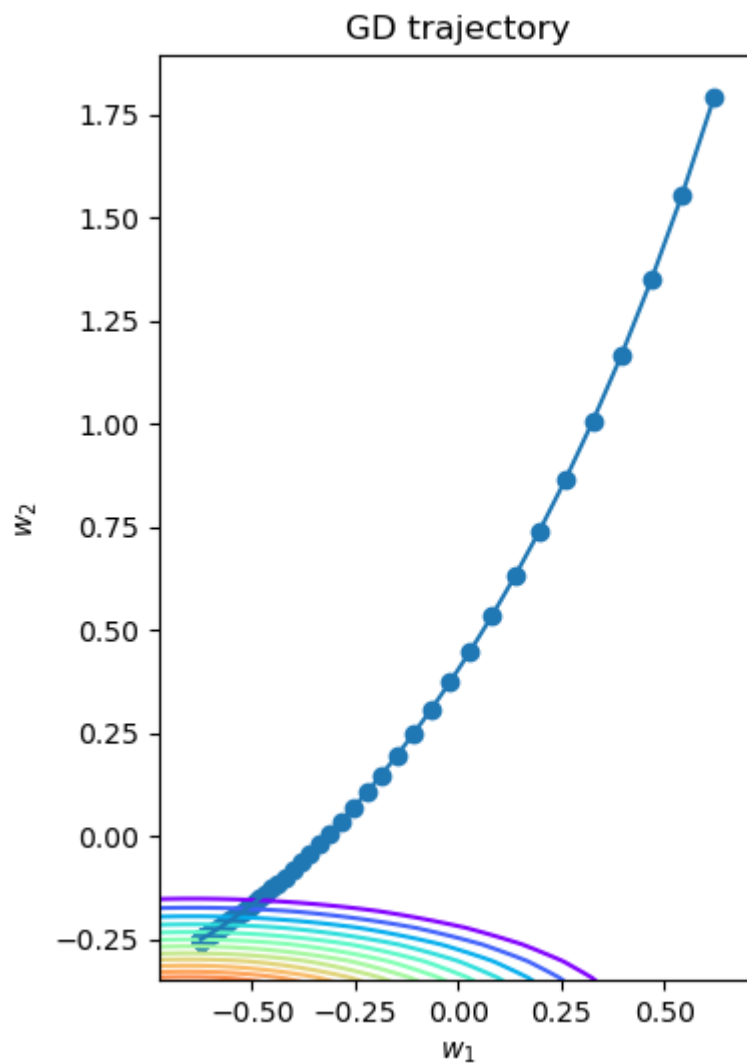
```
B [72]: 1 w_list = stochastic_gradient_descent(w_init, X
2 plot_gd(w_list, X, y, loss)
3 print('First iteration loss: ', loss.calc_loss(X
4 print('Last iteration loss: ', loss.calc_loss(X
```



First iteration loss: 320.97618013755493  
Last iteration loss: 0.8740066409316497

Очень близко к минимуму подошли. Видно, что если увеличивать количество итераций, то можно будет доходить до минимума. Можно еще попробовать взять, например, 100000, как в простом градиентном спуске, мы бы получили эффективный loss.

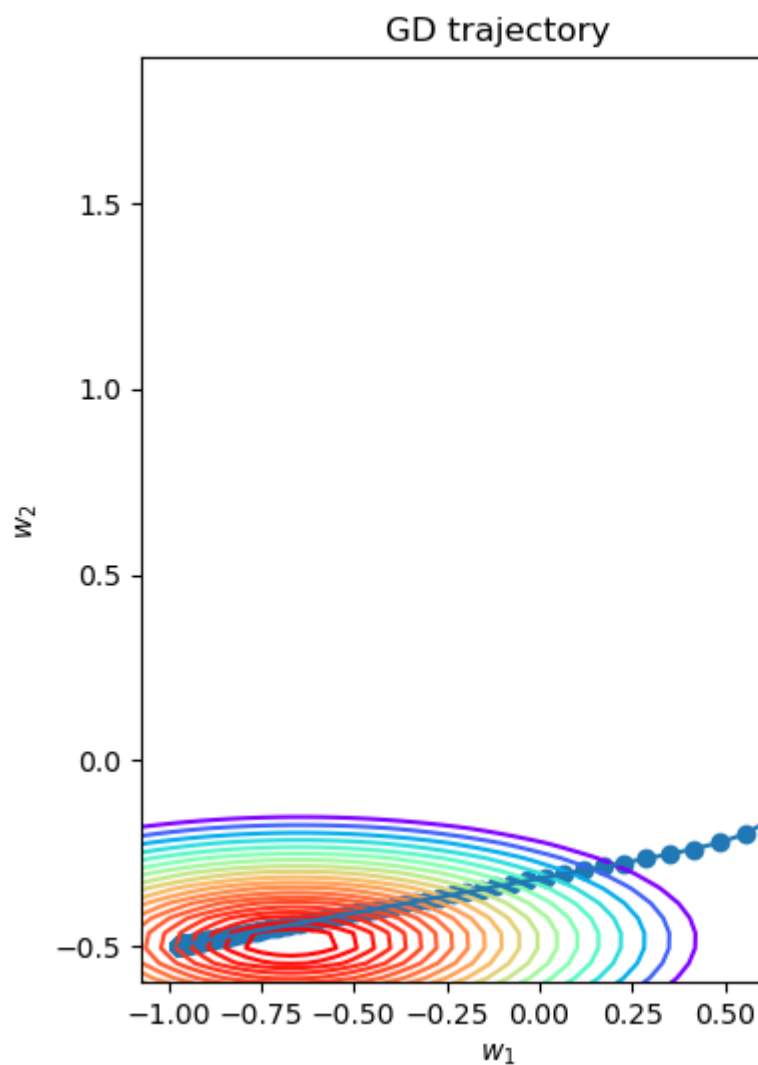
```
B [73]: 1 w_list = stochastic_gradient_descent(w_init, X
2 plot_gd(w_list, X, y, loss)
3 print('First iteration loss: ', loss.calc_loss(X
4 print('Last iteration loss: ', loss.calc_loss(X
```



First iteration loss: 342.9916309778918  
Last iteration loss: 5.447094712720428

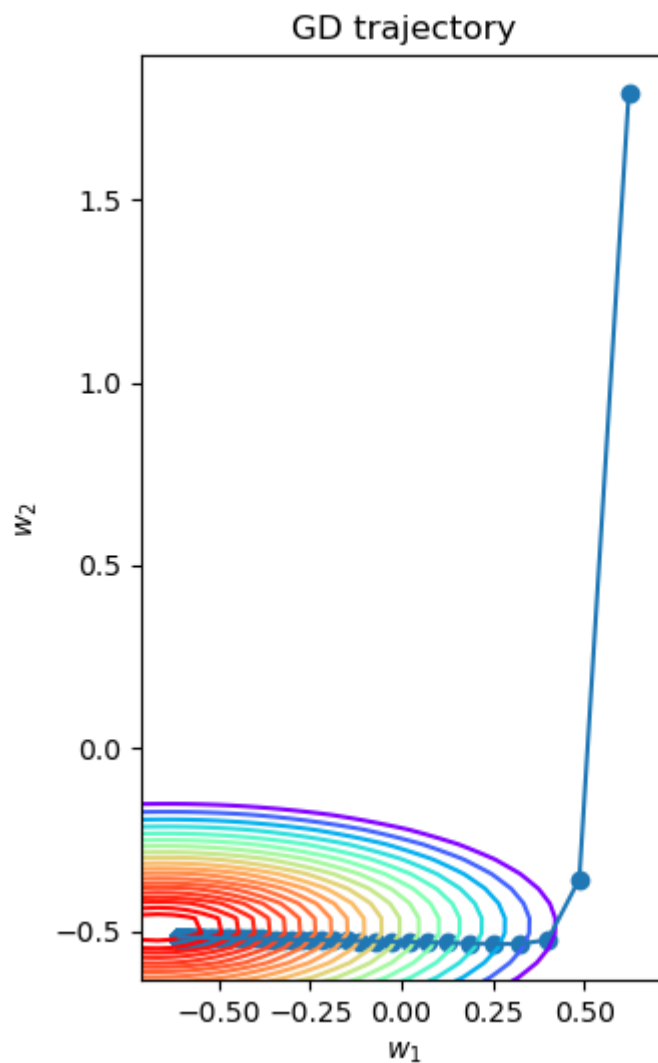


```
B [74]: 1 w_list = stochastic_gradient_descent(w_init, X
2 plot_gd(w_list, X, y, loss)
3 print('First iteration loss: ', loss.calc_loss(X
4 print('Last iteration loss: ', loss.calc_loss(X
```



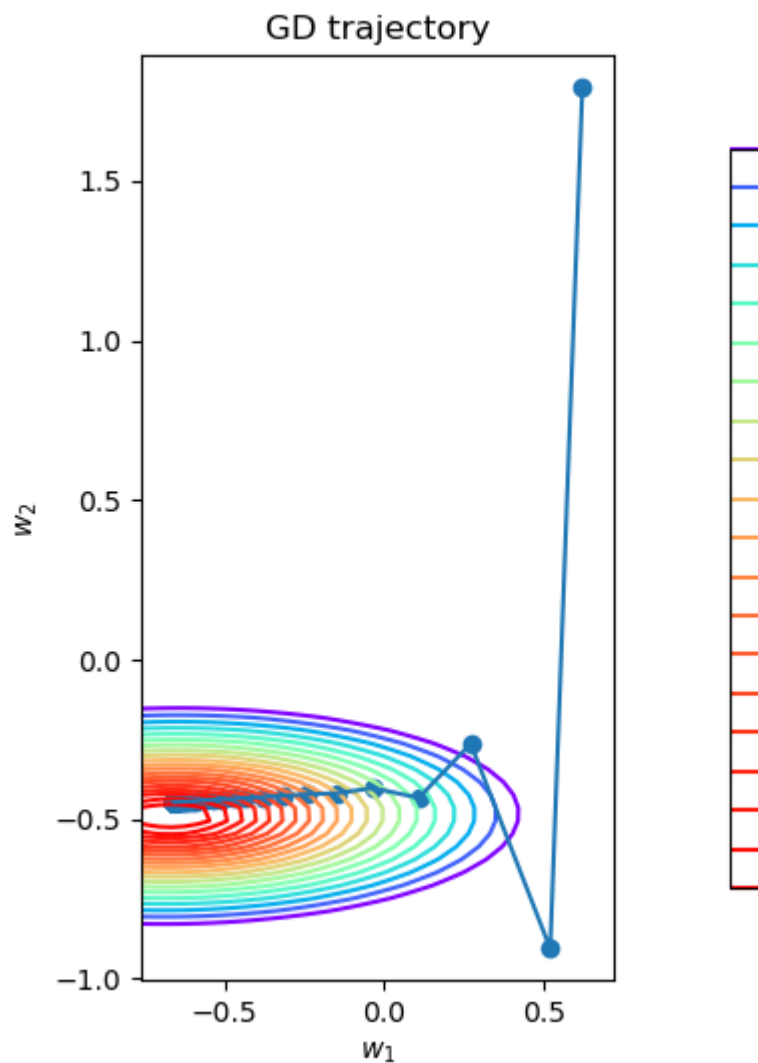
First iteration loss: 108.46361207138666  
Last iteration loss: 1.5739612050157534

```
B [76]: 1 w_list = stochastic_gradient_descent(w_init, X
2 plot_gd(w_list, X, y, loss)
3 print('First iteration loss: ',loss.calc_loss(
4 print('Last iteration loss: ',loss.calc_loss(X
```



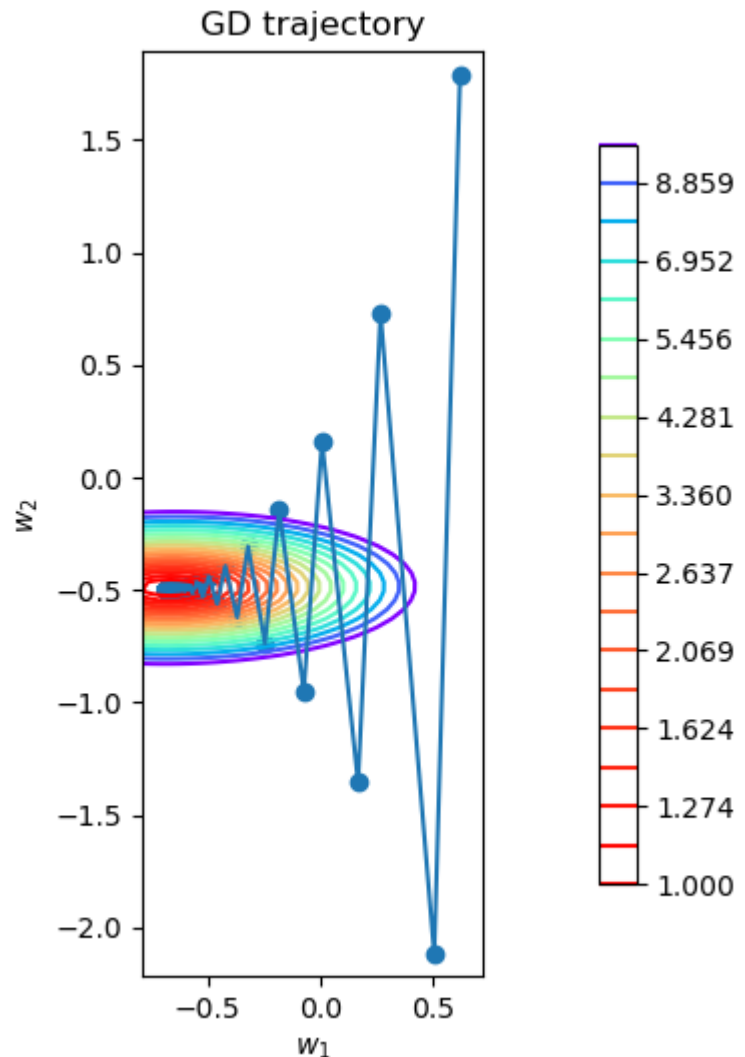
First iteration loss: 12.324328164519052  
Last iteration loss: 0.9494737268620215

```
B [114]: 1 w_list = stochastic_gradient_descent(w_init, X
2 plot_gd(w_list, X, y, loss)
3 print('First iteration loss: ', loss.calc_loss(X
4 print('Last iteration loss: ', loss.calc_loss(X
```



First iteration loss: 26.047064784562675  
Last iteration loss: 1.0157915932020185

```
B [116]: 1 w_list = stochastic_gradient_descent(w_init, X
2 plot_gd(w_list, X, y, loss)
3 print('First iteration loss: ',loss.calc_loss(
4 print('Last iteration loss: ',loss.calc_loss(X
```



First iteration loss: 224.6651069191318  
 Last iteration loss: 0.8728661708031257

Исследование показывает, что длина шага ( $\eta$ ) и batch size  $n$  влияют на процесс градиентного спуска. Малые значения  $n$  означают, что точка минимума не будет достигнута, а большой batch size требует больших вычислительных ресурсов. Необходимо найти оптимальные значения  $\eta$  и  $n$  для достижения минимума. Рекомендуется установить  $\eta$  и  $n$  так, чтобы процесс градиентного спуска был эффективным. Стохастический метод менее требователен и выполняется быстрее, но также демонстрирует хорошие результаты.

Можно заметить, что поведение градиентного спуска, особенно в начале, очень сильно зависит от размера шага  $\eta$ .

Как правило, в начале спуска мы хотим делать большие шаги, чтобы быстро приблизиться к минимуму, а позже мы уже хотим делать маленькие шаги, чтобы более точно достичь минимума и не "перепрыгнуть" его.

Также следует учитывать, что стохастический метод менее выполняется быстрее, но может демонстрировать хорошие результаты, учитывая только часть объектов.

Рекомендуется выбирать оптимальный баланс между скоростью и качеством результата. Например, установка `batch_size` равной половине объектов даст хороший результат и потребует вдвое меньше вычислений.

Чтобы достичь такого поведения мы можем постепенно уменьшать длину шага с увеличением номера итерации. Сделать это можно, например, уменьшая длину шага по следующей формуле:

$$\eta_t = \lambda \left( \frac{s_0}{s_0 + t} \right)^p$$

где  $\eta_t$  — длина шага на итерации  $t$ ,  $\lambda$  — начальная длина шага,  $s_0$  — параметр, влияющий на скорость сходимости,  $p$  — параметр, влияющий на скорость сходимости.

**Задание 1.6: Реализуйте функцию `stochastic_gradient_descent` по формуле выше.** Параметр  $s_0$  возьмите равным количеству объектов в наборе данных. Параметр  $p$  — параметр, влияющий на скорость сходимости.

```
B [118]: 1 def stochastic_gradient_descent(
2         w_init: np.ndarray,
3         X: np.ndarray,
4         y: np.ndarray,
5         loss: BaseLoss,
6         lr: float,
7         batch_size: int,
8         p: float,
9         n_iterations: int = 1000,
10    ) -> List[np.ndarray]:
11
12        grad_list = []
13        batch_indices = np.random.choice(X.shape[0], batch_size)
14        batch = X[batch_indices]
15        batch_y = y[batch_indices]
16        for i in range(n_iterations):
17            grad_list.append(w_init)
18            w_init = w_init - lr * (1 / (1 + i)) ** p * loss(batch, batch_y)
19        return grad_list
20
21
```

Функция `stochastic_gradient_descent` принимает несколько параметров:

- `w_init` - начальные значения параметров модели.
- `X` - входные характеристики набора данных.
- `y` - целевые значения набора данных.
- `loss` - экземпляр класса функции потерь, который вычисляет потерю.
- `lr` - скорость обучения алгоритму SGD.
- `batch_size` - размер пакета, используемого для каждой итерации.
- `p` - параметр мощности, используемый при снижении скорости обучения.
- `n_iterations` - количество итераций для выполнения (по умолчанию 1000).

Функция возвращает список массивов `pumpy`, `grad_list` кот значения параметров на каждой итерации.

Мы сначала определяем класс функции потерь, `MeanSqua` среднеквадратичную ошибку потерь и ее градиент. Затем мы случайные данные  $X$  и  $y$  для демонстрационных целей.

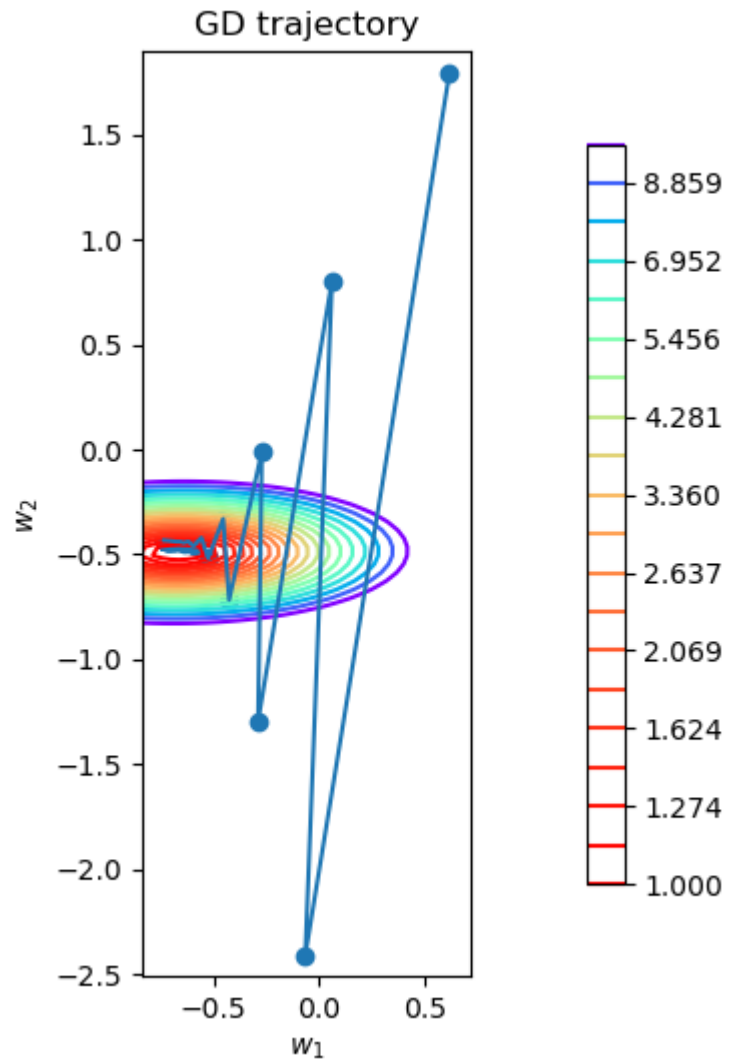
Далее мы инициализируем параметры модели `w_init` в виде создаем экземпляр `MeanSquaredError` функции потерь и за такие как скорость обучения (`lr`), размер пакета (`batch_size`), количество итераций (`n_iterations`).

Наконец, мы применяем `stochastic_gradient_descent` функ входными данными и сохраняем обновленные значения па Затем мы печатаем обновленные значения параметров на

**Задание 1.7:** При помощи новой функции `stochastic_` функции `plot_gd` нарисуйте траекторию градиентного значений параметра  $p$ . Используйте не менее четырёх  $p$  Хорошими могут быть значения, лежащие в промежутке от возьмите равным 0.01, а параметр `batch_size` равным 10 параметра  $p$  вычисляйте значение функционала ошибок `calc_loss` на первой и последней итерациях стохастичес

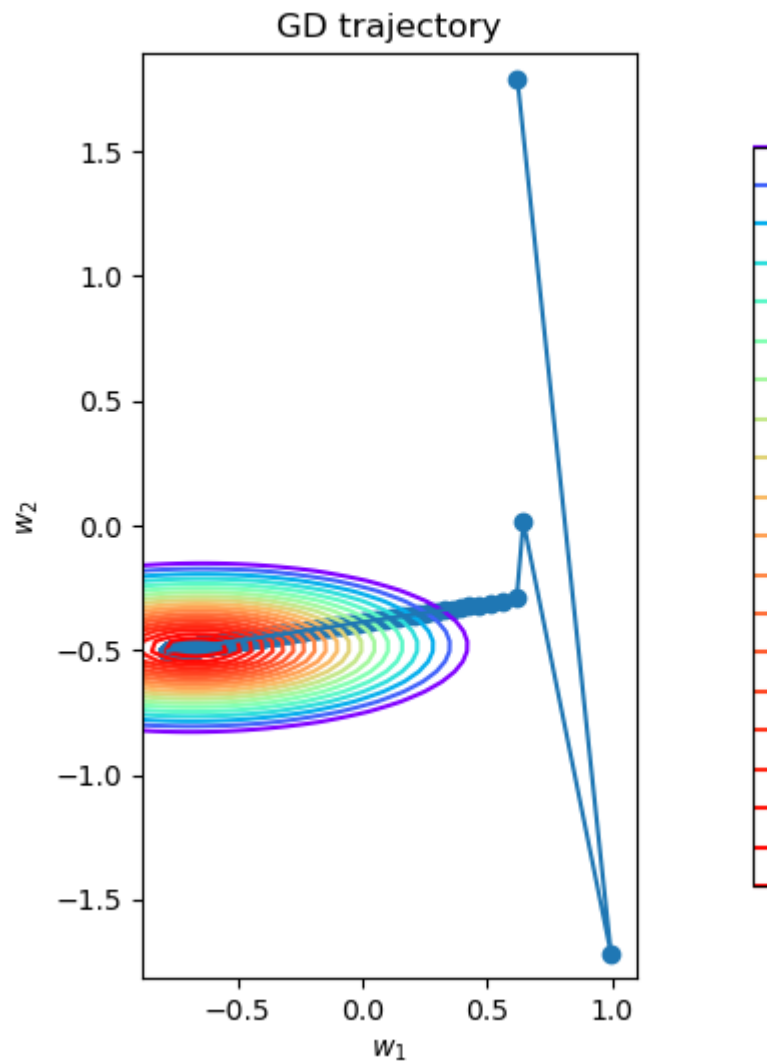
Сделайте и опишите свои выводы о том, как параметр  $p$  в стохастического градиентного спуска

```
B [119]: 1 w_list = stochastic_gradient_descent(w_init, X
2 plot_gd(w_list, X, y, loss)
3 print('First iteration loss: ',loss.calc_loss(
4 print('Last iteration loss: ',loss.calc_loss(X
```



First iteration loss: 298.2310061325282  
Last iteration loss: 1.171475933393098

```
B [121]: 1 w_list = stochastic_gradient_descent(w_init, X
2 plot_gd(w_list, X, y, loss)
3 print('First iteration loss: ',loss.calc_loss(
4 print('Last iteration loss: ',loss.calc_loss(X
```



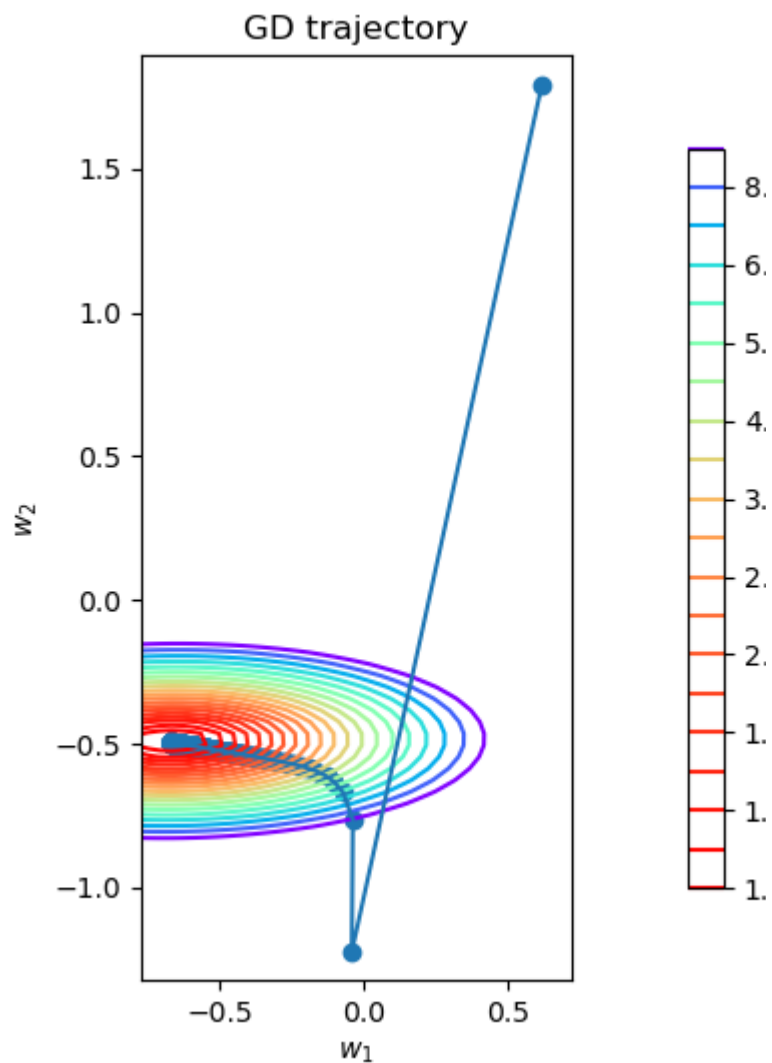
First iteration loss: 144.5542390836199  
Last iteration loss: 0.9786857142942068



```

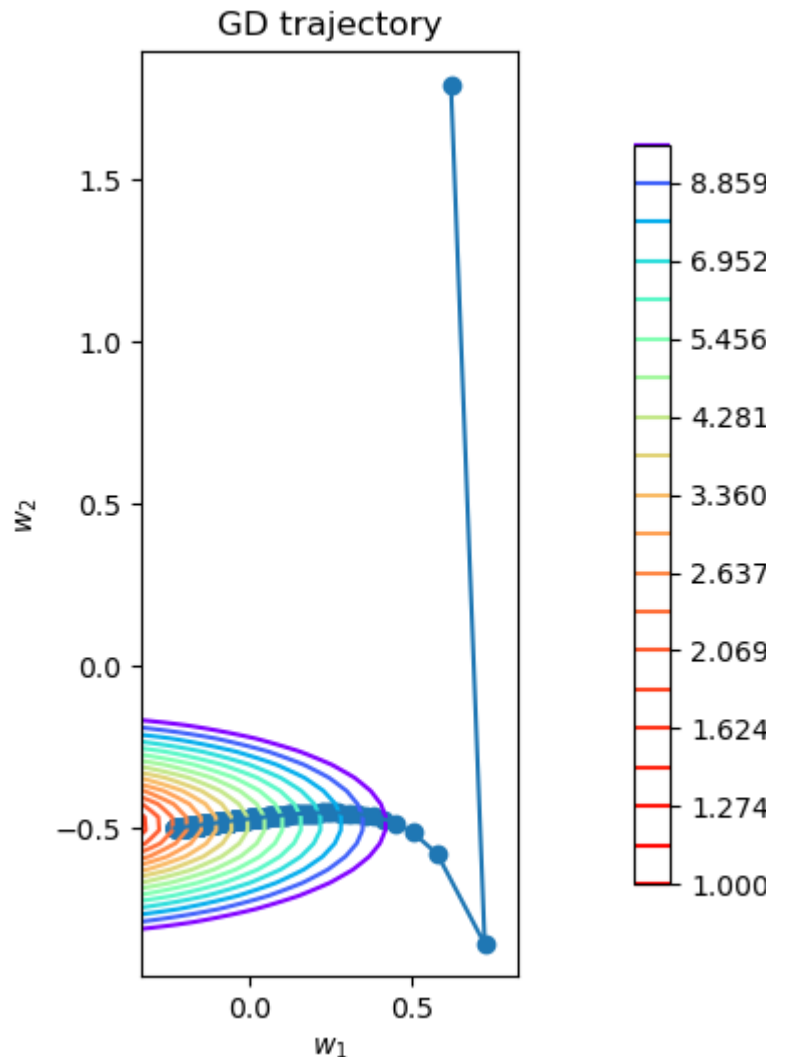
B [122]: 1 w_list = stochastic_gradient_descent(w_init, X
2 plot_gd(w_list, X, y, loss)
3 print('First iteration loss: ',loss.calc_loss(X
4 print('Last iteration loss: ',loss.calc_loss(X

```



First iteration loss: 47.28956578327367  
 Last iteration loss: 0.8697160468376856

```
B [123]: 1 w_list = stochastic_gradient_descent(w_init, X
2 plot_gd(w_list, X, y, loss)
3 print('First iteration loss: ',loss.calc_loss(
4 print('Last iteration loss: ',loss.calc_loss(X
```



First iteration loss: 27.01798314255404  
 Last iteration loss: 2.353688398024119

По результатам исследования видно, что при увеличении  $\alpha$  величина шагов уменьшается. При значениях  $\alpha$  от 0,1 до 0, происходит недостаточно быстро, вследствие чего возможен минимум. При значении 0,75 достигается оптимальное соотношение. Дальнейшее увеличение значения  $\alpha$  приводит к слишком большому результату, вследствие чего алгоритм не доходит до точки минимума (является слишком большим).

**Задание 1.8: Сравните сходимость обычного градиентного стохастической версии:** Нарисуйте график зависимости  $\alpha$  от ошибки (лосса) (его можно посчитать при помощи метода  $\text{loss}$  из датасета и  $w$  с соответствующей итерации) от номера итерации полученных при помощи обычного и стохастического градиентного спуска с одинаковыми параметрами. В SGD параметр `batch_size` равен 1.

```
B [126]: 1 print(n_features)
          2 print(n_objects)
          3 print(batch_size)
          4 print(num_steps)
```

```
2
300
10
43
```

```
B [147]: 1 steps = np.arange(50)
          2
          3 w_list_gd = gradient_descent(w_init, X, y, loss)
          4 grad_list=[(loss.calc_loss(X, y, w_list_gd[i]))
          5
          6 w_list_sgd = stochastic_gradient_descent(w_init, X, y, loss)
          7 st_grad_list=[(loss.calc_loss(X, y, w_list_sgd[i]))
          8
```

```
B [148]: 1 plt.plot(steps, grad_list, label='GD', color='blue')
          2 plt.plot(steps, st_grad_list, label='SGD', color='red')
          3 plt.xlabel('steps')
          4 plt.ylabel('Функционал ошибки')
          5 plt.title('Изменение функционала ошибки от шага')
          6 plt.legend()
          7 plt.show()
```

