

# Лекция 2

## Методы оптимизации

# Градиентный спуск

# Структура задач машинного обучения

$D = \{x, y\}_{i=1}^N$  — обучающая выборка.

$y^* = A(x; \mu)$  — метод предсказания  
 $\mu$  — параметры

$L(D, \mu) = E(D, \mu) + R(\mu)$  — функция потерь  
 $E(D, \mu)$  — функция ошибки  
 $R(\mu)$  — функция регуляризации

$L(D, \mu) \rightarrow \min_{\mu}$  — процедура обучения

Ключевая задача в машинном обучении: **оптимизация**.

# Напоминание: градиентный спуск

$$f(x) \rightarrow \min_x$$

$\eta$  — величина шага (гиперпараметр)

- 1 Инициализация

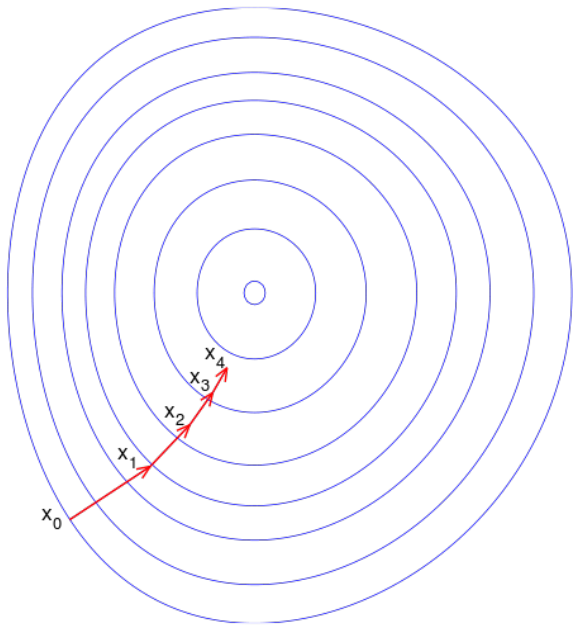
$$k = 0, \quad x_k = \text{начальное приближение}$$

- 2 Шаг в сторону сильнейшего убывания

$$x_{k+1} = x_k - \eta \nabla f(x_k)$$

- 3 Повторение до сходимости

$$k := k + 1, \quad \text{перейти к 2}$$



# Особенности градиентного спуска

## Плюсы

- Достаточно универсальный метод
- Легко реализуем

## Минусы

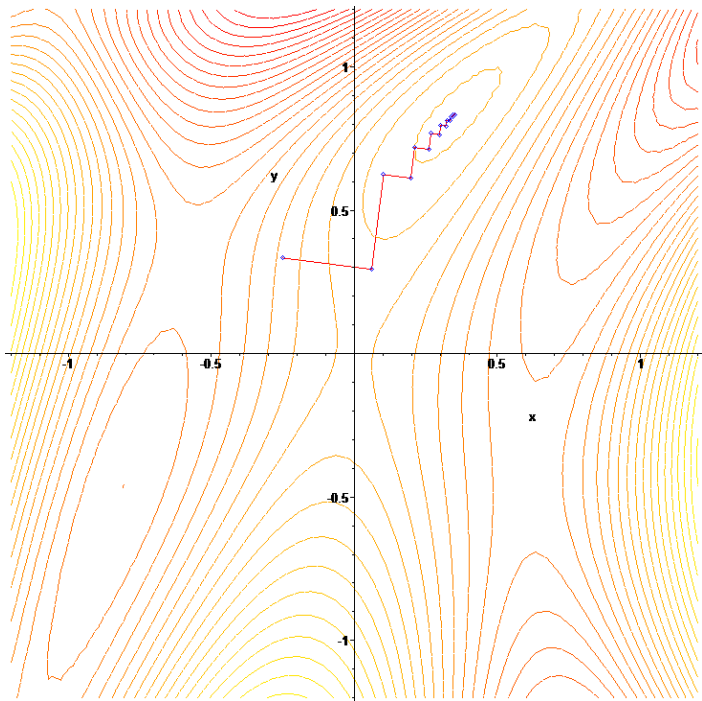
# Особенности градиентного спуска

## Плюсы

- Достаточно универсальный метод
- Легко реализуем

## Минусы

- Попадает в локальные минимумы
- Шаги могут быть медленными
- Неэффективен для больших выборок
- Неприменим для недифференцируемых функций





# Стохастические методы

# Проблема эффективности

Пусть  $D$  содержит очень много элементов.

- Долго считать градиент
- Можно сделать мало итераций за разумное время
- Плохое решение

Идея: приближенно оценивать градиент.

# Стохастический градиентный спуск (SGD)

Разложим функцию потерь:

$$L(D, \mu) = \frac{1}{N} \sum_{i=1}^N l(x_i, \mu).$$

# Стохастический градиентный спуск (SGD)

Разложим функцию потерь:

$$L(D, \mu) = \frac{1}{N} \sum_{i=1}^N l(x_i, \mu).$$

Разложим градиент:

$$\nabla_{\mu} L(D, \mu) = \frac{1}{N} \sum_{i=1}^N \nabla_{\mu} l(x_i, \mu).$$

# Стохастический градиентный спуск (SGD)

Разложим функцию потерь:

$$L(D, \mu) = \frac{1}{N} \sum_{i=1}^N l(x_i, \mu).$$

Разложим градиент:

$$\nabla_{\mu} L(D, \mu) = \frac{1}{N} \sum_{i=1}^N \nabla_{\mu} l(x_i, \mu).$$

Пусть  $I = (i_1, i_2, \dots, i_m)$  — небольшая подвыборка  $D$

Приблизим градиент:

$$\nabla_{\mu} L(D, \mu) \approx \frac{1}{m} \sum_I \nabla_{\mu} l(x_i, \mu).$$

# Стохастический градиентный спуск (SGD)

$$\sum_{i=1}^N f_i(x) \rightarrow \min_x$$

$\eta$  — величина шага,  $m$  — размер подвыборки

❶ Инициализация

$k = 0$ ,  $x_k$  = начальное приближение

❷ Шаг *примерно* в сторону сильнейшего убывания

$I :=$  случайная подвыборка размера  $m$

$$x_{k+1} = x_k - \eta \sum_I \nabla f_i(x_k)$$

❸ Повторение до сходимости

$k := k + 1$ , перейти к 2

# SGD + Mini-batches

На каждой итерации:

- Перемешать выборку
- Разбить выборку на равные части размера  $m$  (мини-батчи)

$$I_1 + I_2 + \dots + I_{N/m} = \{1, \dots, N\}$$

- Для  $j = 1, \dots, N/m$

$$x := x - \eta \sum_{I_j} \nabla f_i(x)$$

Возможный вариант:  $m = 1$ .

## Методы 2 порядка



# Ряд Тейлора

$f(x)$  — одномерная бесконечно дифференцируемая функция

$$\begin{aligned} f(x) &= f(x_0) + \frac{f'(x_0)}{1!}(x - x_0) + \frac{f''(x_0)}{2!}(x - x_0)^2 + \dots = \\ &= \sum_{n=0}^{\infty} \frac{f^{(n)}(x_0)}{n!}(x - x_0)^n. \end{aligned}$$

(не всегда сходится)

# Ряд Тейлора и градиентный спуск

$$f(x) \rightarrow \min_x$$

Разложим в ряд Тейлора до первой производной

$$f(x) \approx f(x_0) + f'(x_0)(x - x_0)$$

# Ряд Тейлора и градиентный спуск

$$f(x) \rightarrow \min_x$$

Разложим в ряд Тейлора до первой производной

$$f(x) \approx f(x_0) + f'(x_0)(x - x_0)$$

Минимизируем

$$f(x_0) + f'(x_0)(x - x_0) \rightarrow \min_x$$

# Ряд Тейлора и градиентный спуск

$$f(x) \rightarrow \min_x$$

Разложим в ряд Тейлора до первой производной

$$f(x) \approx f(x_0) + f'(x_0)(x - x_0)$$

Минимизируем

$$f(x_0) + f'(x_0)(x - x_0) \rightarrow \min_x$$

$f'(x_0) > 0$  — минимум в направлении  $x \rightarrow -\infty$

$f'(x_0) < 0$  — минимум в направлении  $x \rightarrow +\infty$

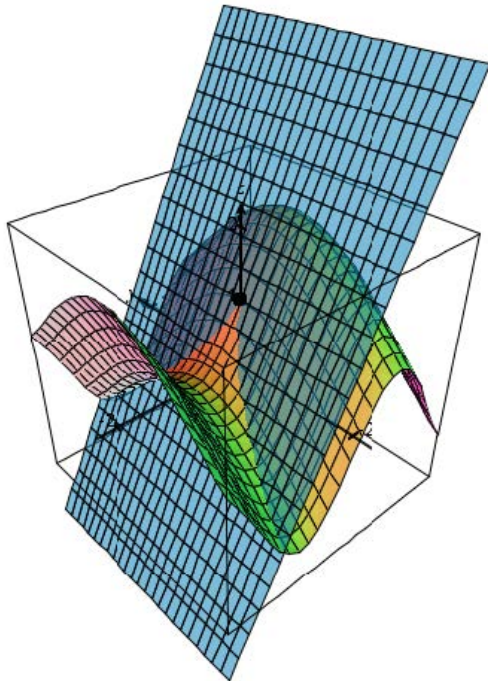
Оптимальный шаг: в направлении  $-f'(x_0)$

# Многомерный градиентный спуск

$$f(x) \rightarrow \min_x$$

Разложим в ряд Тейлора до первой производной

$$f(x) \approx f(x_0) + (x - x_0)^T \nabla f(x_0)$$



# Многомерный градиентный спуск

Задача:

$$f(x) \rightarrow \min_x$$

Разложим в ряд Тейлора до первой производной

$$f(x) \approx f(x_0) + (x - x_0)^T \nabla f(x_0)$$

Когда достигается максимум  $\langle x - x_0, \nabla f(x_0) \rangle$ ?

# Многомерный градиентный спуск

Задача:

$$f(x) \rightarrow \min_x$$

Разложим в ряд Тейлора до первой производной

$$f(x) \approx f(x_0) + (x - x_0)^T \nabla f(x_0)$$

Когда достигается максимум  $\langle x - x_0, \nabla f(x_0) \rangle$ ?

Ответ: когда  $(x - x_0)$  и  $\nabla f(x_0)$  параллельны



# Многомерный градиентный спуск

Задача:

$$f(x) \rightarrow \min_x$$

Разложим в ряд Тейлора до первой производной

$$f(x) \approx f(x_0) + (x - x_0)^T \nabla f(x_0)$$

Когда достигается максимум  $\langle x - x_0, \nabla f(x_0) \rangle$ ?

Ответ: когда  $(x - x_0)$  и  $\nabla f(x_0)$  параллельны

Наибольшее возрастание:  $x - x_0 = \eta \nabla f(x_0)$

Наибольшее убывание:  $x - x_0 = -\eta \nabla f(x_0)$

## Минимизация 2 порядка

$$f(x) \rightarrow \min_x$$

Разложим в ряд Тейлора до второй производной

$$f(x) \approx f(x_0) + f'(x_0)(x - x_0) + \frac{1}{2}f''(x_0)(x - x_0)^2$$

## Минимизация 2 порядка

$$f(x) \rightarrow \min_x$$

Разложим в ряд Тейлора до второй производной

$$f(x) \approx f(x_0) + f'(x_0)(x - x_0) + \frac{1}{2}f''(x_0)(x - x_0)^2$$

Минимизируем

$$f(x_0) + f'(x_0)(x - x_0) + \frac{1}{2}f''(x_0)(x - x_0)^2 \rightarrow \min_x$$

## Минимизация 2 порядка

$$f(x) \rightarrow \min_x$$

Разложим в ряд Тейлора до второй производной

$$f(x) \approx f(x_0) + f'(x_0)(x - x_0) + \frac{1}{2}f''(x_0)(x - x_0)^2$$

Минимизируем

$$f(x_0) + f'(x_0)(x - x_0) + \frac{1}{2}f''(x_0)(x - x_0)^2 \rightarrow \min_x$$

Минимум при

$$x = x_0 - \frac{f'(x_0)}{f''(x_0)}$$

## Минимизация 2 порядка (многомерная)

$$f(x) \rightarrow \min_x$$

Разложим в ряд Тейлора до второй производной

$$f(x) \approx f(x_0) + (x - x_0)^T \nabla f(x_0) + \frac{1}{2} (x - x_0)^T \nabla^2 f(x_0) (x - x_0)$$

## Минимизация 2 порядка (многомерная)

$$f(x) \rightarrow \min_x$$

Разложим в ряд Тейлора до второй производной

$$f(x) \approx f(x_0) + (x - x_0)^T \nabla f(x_0) + \frac{1}{2} (x - x_0)^T \nabla^2 f(x_0) (x - x_0)$$

Минимизируем

$$f(x_0) + (x - x_0)^T \nabla f(x_0) + \frac{1}{2} (x - x_0)^T \nabla^2 f(x_0) (x - x_0) \rightarrow \min_x$$

Минимум при

$$x = x_0 - (\nabla^2 f(x_0))^{-1} \nabla f(x_0)$$

# Метод Ньютона

$$f(x) \rightarrow \min_x$$

- 1 Инициализация

$$k = 0, \quad x_k = \text{начальное приближение}$$

- 2 Шаг в точку примерного минимума

$$x_{k+1} = x_k - (\nabla^2 f(x_k))^{-1} \nabla f(x_k)$$

- 3 Повторение до сходимости

$$k := k + 1, \quad \text{перейти к 2}$$

# Метод Ньютона (модификация)

$$f(x) \rightarrow \min_x$$

$\eta$  — величина шага (гиперпараметр)

❶ Инициализация

$$k = 0, \quad x_k = \text{начальное приближение}$$

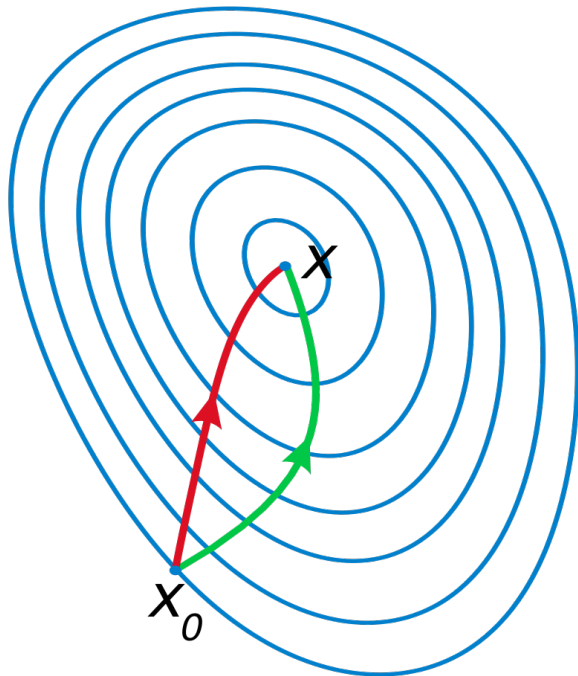
❷ Шаг в сторону примерного минимума

$$x_{k+1} = x_k - \eta \left( \nabla^2 f(x_k) \right)^{-1} \nabla f(x_k)$$

❸ Повторение до сходимости

$$k := k + 1, \quad \text{перейти к 2}$$





# Почему «метод Ньютона»?

Задача:

$$g(x) = 0$$

Итеративное решение:

$$x_{k+1} = x_k - \frac{g(x_k)}{g'(x_k)}$$

# Почему «метод Ньютона»?

Задача:

$$g(x) = 0$$

Итеративное решение:

$$x_{k+1} = x_k - \frac{g(x_k)}{g'(x_k)}$$

Задача:

$$f(x) \rightarrow \min_x \quad \Rightarrow \quad f'(x) = 0$$

Итеративное решение:

$$x_{k+1} = x_k - \frac{f'(x_k)}{f''(x_k)}$$

# Методы 0 порядка

# Покоординатный спуск

(Производную нельзя вычислить)

$$f(x) = f(x_1, x_2, \dots, x_n) \rightarrow \min_x$$

Минимизируем вдоль одной координаты  $i$

$$f(x_1, \dots, x_{i-1}, z, x_{i+1}, \dots, x_n) \rightarrow \min_z$$

Минимизация линейным поиском  
(перебор всех значений в окрестности).

# Покоординатный спуск

$$f(x) = f(x_1, x_2, \dots, x_n) \rightarrow \min_x$$

- 1 Инициализация

$$k = 0, \quad x^k = \text{начальное приближение}$$

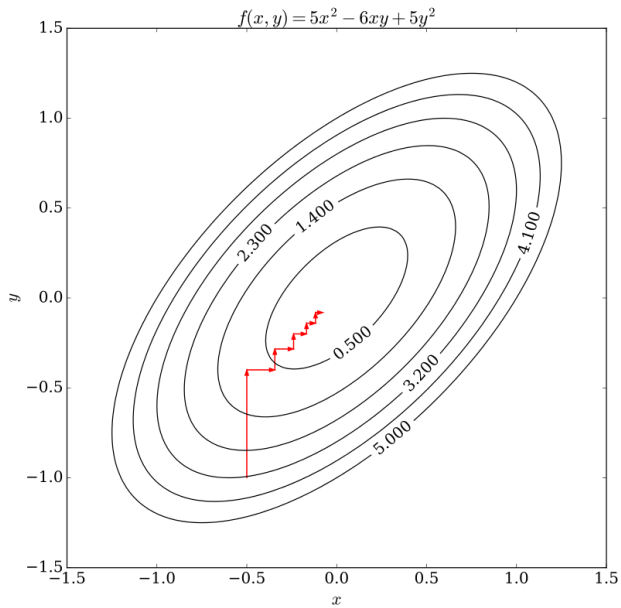
- 2 Минимизируем вдоль координаты  $i$  для  $i = 1, \dots, n$

$$z^* = \arg \min_z f(x_1, \dots, x_{i-1}, z, x_{i+1}, \dots, x_n)$$

$$x^{k+1} = (x_1^k, \dots, x_{i-1}^k, z^*, x_{i+1}^k, \dots, x_n^k)$$

$$k := k + 1$$

- 3 Повторение до сходимости



# Вариации градиентного спуска



# GD + линейный поиск

$$f(x) \rightarrow \min_x$$

## ❶ Инициализация

$$k = 0, \quad x_k = \text{начальное приближение}$$

## ❷ Минимизация в направлении сильнейшего убывания

$$\alpha^* = \arg \min_{\alpha} f(x_k - \alpha \nabla f(x_k))$$

$$x_{k+1} = x_k - \alpha^* \nabla f(x_k)$$

## ❸ Повторение до сходимости

$$k := k + 1, \quad \text{перейти к 2}$$

## GD + переменный шаг

$$f(x) \rightarrow \min_x$$

$\eta$  — базовая величина шага (гиперпараметр)

❶ Инициализация

$$k = 0, \quad x_k = \text{начальное приближение}$$

❷ Шаг в сторону сильнейшего убывания

$$x_{k+1} = x_k - \frac{\eta}{k} \nabla f(x_k)$$

❸ Повторение до сходимости

$$k := k + 1, \quad \text{перейти к 2}$$

## GD + инерция

$$f(x) \rightarrow \min_x$$

$\eta$  — величина шага,  $\alpha$  — влияние инерции

### ❶ Инициализация

$$k = 0, \quad x_k = \text{начальное приближение}, \quad m_k = 0$$

### ❷ Шаг в сторону сильнейшего убывания

$$m_{k+1} := -\eta \nabla f(x_k) + \alpha m_k$$

$$x_{k+1} = x_k + m_{k+1}$$

### ❸ Повторение до сходимости

$$k := k + 1, \quad \text{перейти к 2}$$

## Резюме

# Обзор методов

- Градиентный спуск
- Стохастический градиентный спуск
- Метод Ньютона
- Покоординатный спуск

# Источники

- G. Venter. Review of Optimization Techniques.
- Википедия :)