

# Задача оптимизации

$$\min_{\substack{f(x) \leq 0, \\ i=1, \dots, m, \\ x \in Q}} f(x^*) \quad (1)$$

- $Q \subseteq \mathbb{R}^d$  — подмножество  $d$ -мерного пространства
- $f : Q \rightarrow \mathbb{R}$  — некоторая функция, заданная на множестве  $Q$
- В качестве  $\&$  берётся  $\leq, \geq$  либо  $=$
- $g_i(x) : Q \rightarrow \mathbb{R}, i = 1, \dots, m$  — функции, задающие ограничения

## Матрица Гессе [\[ править | править код \]](#)

Матрица этой квадратичной формы образована вторыми частными производными функции. Если все производные существуют, то

$$H(f) = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} & \cdots & \frac{\partial^2 f}{\partial x_2 \partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \frac{\partial^2 f}{\partial x_n \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_n^2} \end{bmatrix}$$

[Определитель](#) этой матрицы называется **определителем Гессе**, или просто **гессианом** <sup>[[источник не указан 3595 дней](#)]</sup>.

Матрицы Гессе используются в задачах [оптимизации методом Ньютона](#). Полное вычисление матрицы Гессе может быть затруднительно, поэтому были разработаны [квазиньютоновские](#) алгоритмы, основанные на приближённых выражениях для матрицы Гессе. Наиболее известный из них — [алгоритм Бroyдена](#) — [Флетчера](#) — [Гольдфарба](#) — [Шанно](#).

$$\min_{x \in \mathbb{R}^d} \left\{ f(x) = \frac{1}{n} \sum_{i=1}^n \log(1 + \exp(-y_i \cdot (\mathbf{A}x)_i)) \right\}, \quad (2)$$

- $\mathbf{A} \in \mathbb{R}^{n \times d}$  — матрица признаков,  $n$  — количество объектов
- $x \in \mathbb{R}^d$  — вектор параметров,  $d$  — количество параметров
- $y \in \{-1, 1\}^n$  — вектор ответов
- $(\mathbf{A}x)_i, y_i$  —  $i$ -е компоненты векторов  $\mathbf{A}x$  и  $y$  соответственно

Данная задача возникает как один из подходов к решению задачи бинарной классификации.

## Задачи оптимизации. Первые наблюдения.

- ❶ В общем случае задачи оптимизации могут не иметь решения. Например, задача  $\min_{x \in \mathbb{R}} x$  не имеет решения.
- ❷ Задачи оптимизации часто нельзя решить аналитически.
- ❸ Их сложность зависит от вида целевой функции  $f$ , множества  $Q$  и может зависеть от размерности  $x$ .

Если же задача оптимизации имеет решение, то на практике её обычно решают, вообще говоря, приближённо. Для этого применяются специальные алгоритмы, которые и называют методами оптимизации.

# Методы оптимизации I

- Нет смысла искать лучший метод для решения конкретной задачи. Например, лучший метод для решений задачи  $\min_{x \in \mathbb{R}^d} \|x\|^2$  сходится за 1 итерацию: этот метод просто всегда выдаёт ответ  $x^* = 0$ . Очевидно, что для других задач такой метод не пригоден.
- Эффективность метода определяется для класса задач, т.к. обычно численные методы разрабатываются для *приближённого* решения множества однотипных задач.
- Метод разрабатывается для класса задач  $\implies$  метод не может иметь с самого начала полной информации о задаче. Вместо этого метод использует модель задачи, например, формулировку задачи, описание функциональных компонент, множества, на котором происходит оптимизация и т.д.

# Методы оптимизации II

- Предполагается, что численный метод может накапливать специфическую информацию о задаче при помощи некоторого *оракула*. Под оракулом можно понимать некоторое устройство, которое отвечает на последовательные вопросы численного метода.

## Примеры оракулов

- **Оракул нулевого порядка** в запрашиваемой точке  $x$  возвращает значение целевой функции  $f(x)$ .
- **Оракул первого порядка** в запрашиваемой точке возвращает значение функции  $f(x)$  и её градиент в данной точке  $\nabla f(x) = \left( \frac{\partial f(x)}{\partial x_1}, \dots, \frac{\partial f(x)}{\partial x_n} \right)$ .

# Общая итеративная схема метода оптимизации $\mathcal{M}$

**Входные данные:** начальная точка  $x^0$  (0 – верхний индекс), требуемая точность решения задачи  $\varepsilon > 0$ .

**Настройка.** Задать  $k = 0$  (счётчик итераций) и  $I_{-1} = \emptyset$  (накапливаемая информационная модель решаемой задачи).

**Основной цикл**

- ① Задать вопрос к оракулу  $\mathcal{O}$  в точке  $x^k$ .
- ② Пересчитать информационную модель:  $I_k = I_{k-1} \cup (x_k, \mathcal{O}(x^k))$ .
- ③ Применить правило метода  $\mathcal{M}$  для получения новой точки  $x^{k+1}$  по модели  $I_k$ .
- ④ Проверить критерий остановки  $\mathcal{T}_\varepsilon$ . Если критерий выполнен, то выдать ответ  $\bar{x}$ , иначе положить  $k := k + 1$  и вернуться на шаг 1.

## Сложность методов оптимизации

- **Аналитическая сложность** — число обращений к оракулу, необходимое для решения задачи с точностью  $\varepsilon$ .
- **Арифметическая сложность** — общее число вычислений (включая работу оракула), необходимых для решения задачи с точностью  $\varepsilon$ .



# Примеры итерационных методов. Градиентный спуск

Рассмотрим задачу оптимизации

$$\min_{x \in \mathbb{R}^d} f(x), \quad (5)$$

где функция  $f(x)$  дифференцируема. Предположим, что в любой точке мы можем посчитать её градиент.

---

## Алгоритм 1 Градиентный спуск с постоянным размером шага

---

**Вход:** размер шага  $\gamma > 0$ , стартовая точка  $x^0 \in \mathbb{R}^d$ , количество итераций  $N$

- 1: **for**  $k = 0, 1, \dots, N - 1$  **do**
- 2:     Вычислить  $\nabla f(x^k)$
- 3:      $x^{k+1} = x^k - \gamma \nabla f(x^k)$
- 4: **end for**

**Выход:**  $x^N$

---

# Примеры итерационных методов. Метод Ньютона

Рассмотрим задачу оптимизации

$$\min_{x \in \mathbb{R}^d} f(x), \quad (6)$$

где функция  $f(x)$  дважды непрерывно дифференцируема.

Предположим, что в любой точке мы можем посчитать её градиент и матрицу вторых производных  $\nabla^2 f(x)$ .

---

## Алгоритм 2 Метод Ньютона

---

**Вход:** стартовая точка  $x^0 \in \mathbb{R}^d$ , количество итераций  $N$

- 1: **for**  $k = 0, 1, \dots, N - 1$  **do**
- 2:     Вычислить  $\nabla f(x^k)$  и  $\nabla^2 f(x^k)$
- 3:      $x^{k+1} = x^k - (\nabla^2 f(x^k))^{-1} \nabla f(x^k)$
- 4: **end for**

**Выход:**  $x^N$

---

# Сложность задач оптимизации. Класс задач минимизации липшицевых функций

$$\min_{x \in B_d} f(x) \quad (7)$$

- $B_d = \{x \in \mathbb{R}^d \mid 0 \leq x_i \leq 1, \quad i = 1, \dots, d\}$
- Функция  $f(x)$  является  $M$ -липшицевой на  $B_d$  относительно  $\ell_\infty$ -нормы:

$$\forall x, y \quad |f(x) - f(y)| \leq M \|x - y\|_\infty = M \max_{i=1, \dots, d} |x_i - y_i|. \quad (8)$$

Отображение  $f$  метрического пространства  $(X, \rho_X)$  в метрическое пространство  $(Y, \rho_Y)$  называется липшицевым, если найдётся такая константа  $L$  (константа Липшица этого отображения), что  $\rho_Y(f(x), f(y)) \leq L \cdot \rho_X(x, y)$  при любых  $x, y \in X$ . Это условие называют **условием Липшица**. Отображение с  $L = 1$  (1-липшицево отображение) называют также **коротким отображением**.

Липшицево отображение  $f: X \rightarrow Y$  называется **билипшицевым**, если у него существует обратное  $f^{-1}: Y \rightarrow X$ , которое также является липшицевым.

Отображение  $f: X \rightarrow Y$  называется **колипшицевым**, если существует константа  $L$  такая, что для любых  $x \in X$  и  $y \in Y$  найдётся  $x' \in f^{-1}(y)$  такое, что  $\rho_Y(f(x), y) \leq L \cdot \rho_X(x, x')$ .

## Сложность задач оптимизации. Класс задач минимизации липшицевых функций

### Наблюдение

Множество  $B_d$  является ограниченным и замкнутым, т.е. компактом, а из липшицевости функции  $f$  следует и её непрерывность, поэтому задача (7) имеет решение, ибо непрерывная на компакте функция достигает своих минимального и максимального значений. Пусть  $f_* = \min_{x \in B_d} f(x)$ .

- **Класс методов.** Для данной задачи рассмотрим методы нулевого порядка.
- **Цель:** найти  $\bar{x} \in B_d$ :  $f(\bar{x}) - f_* \leq \varepsilon$ .

# Сложность задач оптимизации. Класс задач минимизации липшицевых функций

Рассмотрим один из самых простых способов решения этой задачи — метод равномерного перебора.

---

## Алгоритм 3 Метод равномерного перебора

---

**Вход:** целочисленный параметр перебора  $p \geq 1$

- 1: Сформировать  $(p + 1)^d$  точек вида  $x_{(i_1, \dots, i_d)} = \left(\frac{i_1}{p}, \frac{i_2}{p}, \dots, \frac{i_d}{p}\right)^T$ , где  $(i_1, \dots, i_d) \in \{0, 1, \dots, p\}^n$
- 2: Среди точек  $x_{(i_1, \dots, i_d)}$  найти точку  $\bar{x}$  с наименьшим значением целевой функции  $f$ .

**Выход:**  $\bar{x}, f(\bar{x})$

---

# Сложность задач оптимизации. Класс задач минимизации липшицевых функций

## Теорема 1 (Теорема 1.1.1 из книги Нестерова 2010 года)

Алгоритм 3 с параметром  $p$  возвращает такую точку  $\bar{x}$ , что

$$f(\bar{x}) - f_* \leq \frac{M}{2p}, \quad (9)$$

откуда следует, что методу равномерного перебора нужно в худшем случае

$$\left(\left\lceil \frac{M}{2\varepsilon} \right\rceil + 2\right)^d \quad (10)$$

обращений к оракулу, чтобы гарантировать  $f(\bar{x}) - f_* \leq \varepsilon$ .



# Сложность задач оптимизации. Класс задач минимизации липшицевых функций

- Предположим  $M = 2$ ,  $d = 13$  И  $\varepsilon = 0.01$ , то есть размерность задачи сравнительно небольшая и точность решения задачи не слишком высокая.
- Необходимое число обращений к оракулу:  
 $(\lfloor \frac{M}{2\varepsilon} \rfloor + 2)^d = 102^{13} > 10^{26}$ .
- Сложность оракула при этом составляет не меньше  $d = 13$  арифметических операций.
- Производительность компьютера:  $10^{11}$  арифметических операций в секунду.
- Общее время: хотя бы  $10^{16}$  секунд, что больше 300 миллионов лет.

## Выпуклые и гладкие функции

### Хорошие новости:

- Богатая и интересная теория
- Существуют эффективные алгоритмы приближённого решения

### Плохие новости:

- Класс выпуклых и гладких задач не очень широк
- На практике часто приходится сталкиваться с невыпуклыми задачами

Тем не менее, иногда методы выпуклой гладкой оптимизации хорошо себя проявляют на практике и на невыпуклых или негладких задачах, имеющих локальную выпуклость или гладкость, поэтому имеет смысл для начала разобраться, как они работают для выпуклых и гладких задач.

## Определение 1

Множество  $Q \subseteq \mathbb{R}^d$  называется выпуклым, если для любых двух точек  $x, y \in Q$  и для любого числа  $\alpha \in [0, 1]$  точка  $z = \alpha x + (1 - \alpha)y$  принадлежит множеству  $Q$ .

Это означает, что вместе с любыми двумя точками во множестве содержится и отрезок, их соединяющий.

## Выпуклые функции

### Определение 2

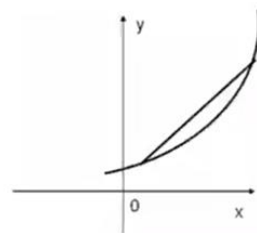
Функция  $f(x)$  заданная на выпуклом множестве  $Q \subseteq \mathbb{R}^d$  называется **выпуклой**, если для любых двух точек  $x, y \in Q$  и для любого числа  $\alpha \in [0, 1]$  выполняется следующее неравенство:

$$f(\alpha x + (1 - \alpha)y) \leq \alpha f(x) + (1 - \alpha)f(y). \quad (12)$$

(если знак  $<$  для всех  $x \neq y$ ,  $\alpha \in (0, 1)$ , то **строго выпуклой**)

В одномерном

случае это означает, что между любыми двумя точками  $x, y \in Q$  график функции  $f$  проходит не выше отрезка, соединяющего  $f(x)$  и  $f(y)$ .

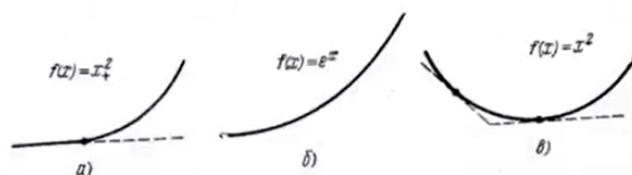


## Определение 3

Функция  $f(x)$  заданная на выпуклом множестве  $Q \subseteq \mathbb{R}^d$  называется  $\mu$ -сильно выпуклой, если для любых двух точек  $x, y \in Q$  и для любого числа  $\alpha \in [0, 1]$  выполняется следующее неравенство:

$$f(\alpha x + (1 - \alpha)y) \leq \alpha f(x) + (1 - \alpha)f(y) - \alpha(1 - \alpha)\frac{\mu}{2}\|x - y\|_2^2. \quad (13)$$

График (строго) выпуклой функции лежит (строго) выше касательной гиперплоскости, а для сильно выпуклой функции график лежит выше некоторого параболоида.

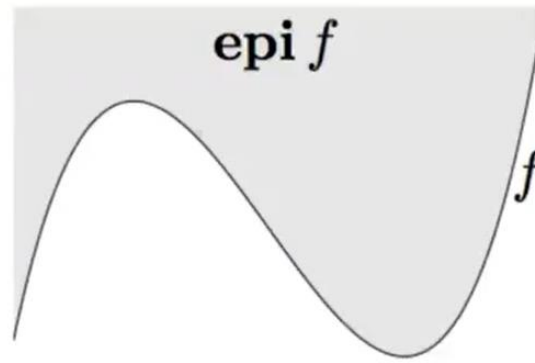


## Свойства выпуклых функций I

- 1) Выпуклые функции непрерывны во всех внутренних точках области определения.
- 2) Сильно выпуклая функция, очевидно, строго выпукла. Обратное неверно.
- 3) Функция выпуклая, тогда и только тогда, когда её надграфик (эпиграф) выпуклое множество, где под надграфиком функции  $f$ , определённой на множестве  $Q \subseteq \mathbb{R}^d$ , понимается следующее множество:

$$\text{epi} f = \{(x, t) \mid x \in Q, t \in \mathbb{R}, t \geq f(x)\} \subseteq \mathbb{R}^{d+1}$$

## Свойства выпуклых функций II



- ④ Если  $f(x)$  выпуклая функция, то множество

$$C_\gamma = \{x \in Q \subseteq \mathbb{R}^d \mid f(x) \leq \gamma, \gamma = \text{const}\}$$

является выпуклым.

## Свойства выпуклых функций III

- ⑤ Для выпуклой функции  $f(x)$  выполнено неравенство Йенсена

$$f\left(\sum_{i=1}^n \lambda_i x_i\right) \leq \sum_{i=1}^n \lambda_i f(x_i)$$

для любых  $\lambda_i \geq 0$  и  $\sum_{i=1}^n \lambda_i = 1$ .

Неравенство Йенсена обобщается и на случай выпуклой комбинации бесконечного (счетного или несчетного) числа точек.