

Градиентный спуск. Теоремы сходимости в
гладком случае (выпуклые, сильно
выпуклые, PL). Верхние и нижние оценки
сходимости.

Даня Меркулов, Петр Остроухов

Оптимизация для всех! ЦУ

Когда остановить?

$$\|\nabla f(x_k)\| \leq \epsilon$$

$$\min_{x \in \mathbb{R}^n} f(x)$$

$$f(x) = x_1 + x_2 - x_3$$

$$x \in \mathbb{R}^3$$

$$|f(x_k) - f(x^*)| < \epsilon$$

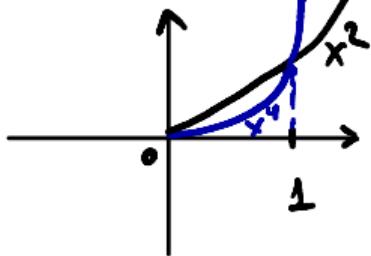
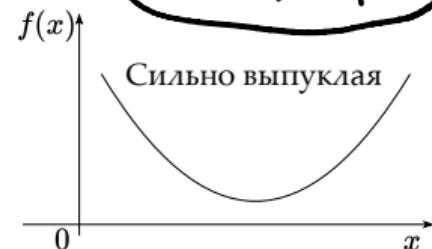
Повторение

$$X_{k+1} = X_k - d \frac{\nabla f(x_k)}{1 \times 1}$$

$$n \times 1 \quad n \times 1 \quad 1 \times 1 \quad n \times 1$$

Виды выпуклости

$\forall \text{лок. мин.} \Rightarrow \text{ГНБ.}$

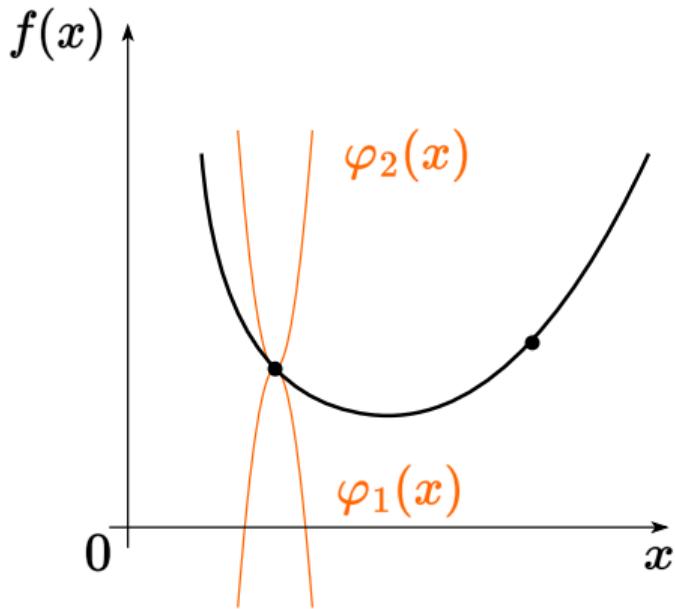


$$\nabla^2 f(x_k) \sum \mu_i I$$

$\forall x \in \mathbb{R} \quad f''(x) \geq \mu > 0$

Рис. 1: Примеры выпуклых функций

Гладкость



Определение: Будем говорить, что функция $f : \mathbb{R}^n \rightarrow \mathbb{R}$ является L -гладкой, если $\forall x, y \in \mathbb{R}^n$ выполнено:

$$\|\nabla f(y) - \nabla f(x)\| \leq L \|y - x\|$$

Липшицевость $g(x)$:

$$\forall x, y: |g(x) - g(y)| \leq L \cdot \|x - y\|$$

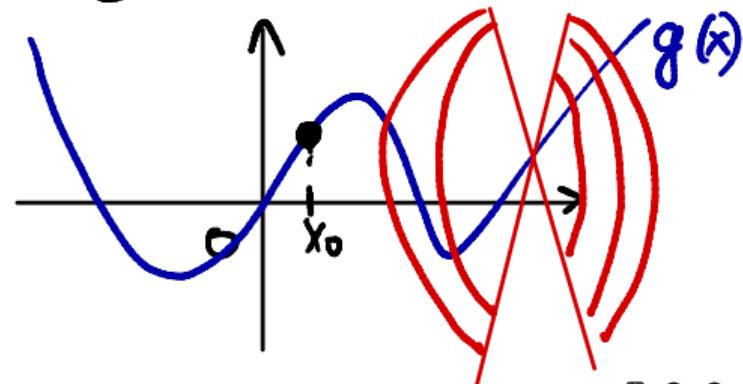
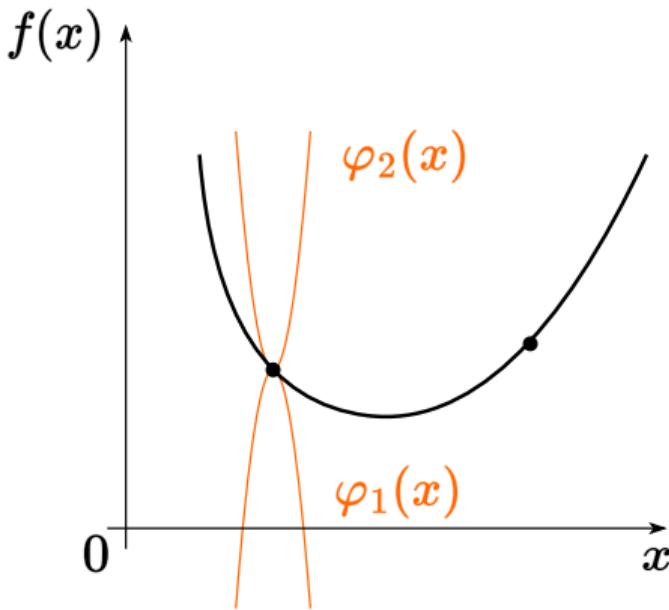


Рис. 2: Иллюстрация Липшицевых парабол, между которыми зажата гладкая функция. Чаще нас интересует мажорирующая из них.

Гладкость

Определение: Будем говорить, что функция $f : \mathbb{R}^n \rightarrow \mathbb{R}$ является L -гладкой, если $\forall x, y \in \mathbb{R}^n$ выполнено:

$$\|\nabla f(y) - \nabla f(x)\| \leq L\|y - x\|$$



Не пытайся

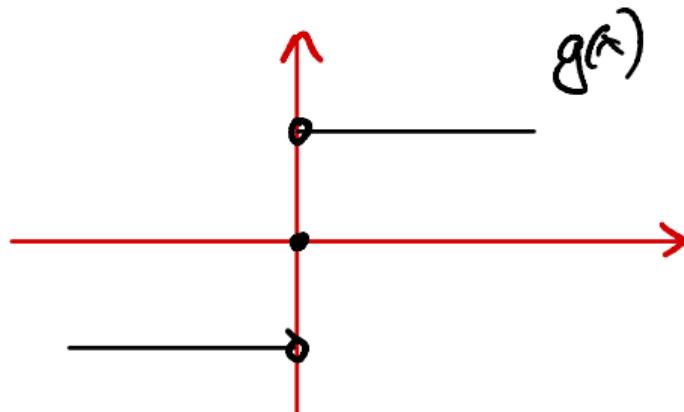
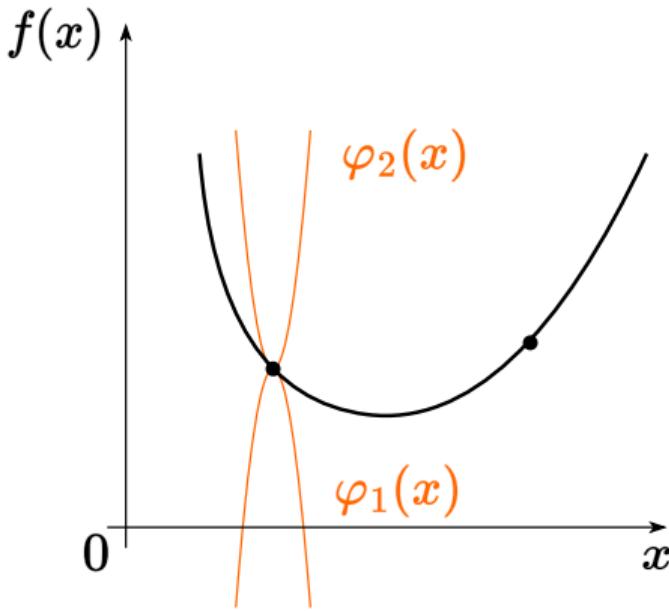


Рис. 2: Иллюстрация Липшицевых парабол, между которыми зажата гладкая функция. Чаще нас интересует мажорирующая из них.

Гладкость

Определение: Будем говорить, что функция $f : \mathbb{R}^n \rightarrow \mathbb{R}$ является L -гладкой, если $\forall x, y \in \mathbb{R}^n$ выполнено:

$$\|\nabla f(y) - \nabla f(x)\| \leq L\|y - x\|$$



Пример не гладкой непрерывной

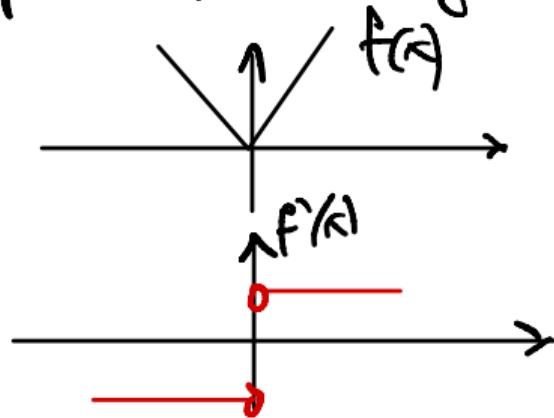
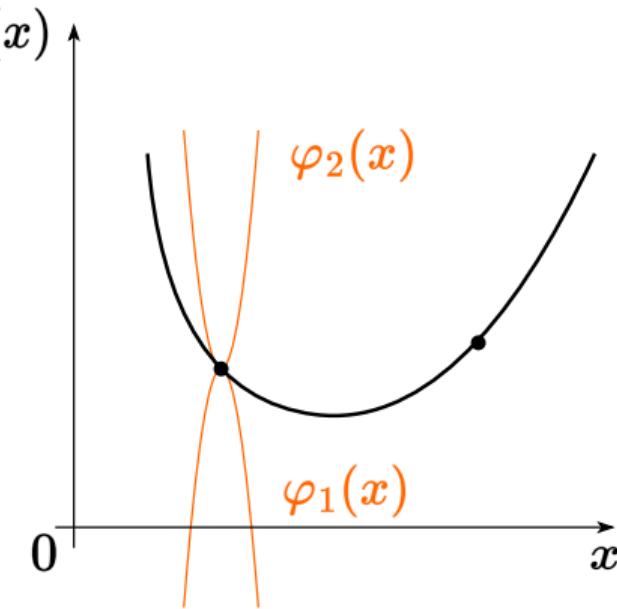


Рис. 2: Иллюстрация Липшицевых парабол, между которыми зажата гладкая функция. Чаще нас интересует мажорирующая из них.

Гладкость

Определение: Будем говорить, что функция $f : \mathbb{R}^n \rightarrow \mathbb{R}$ является L -гладкой, если $\forall x, y \in \mathbb{R}^n$ выполнено:

$$\|\nabla f(y) - \nabla f(x)\| \leq L\|y - x\|$$



Обратим внимание, что значение константы гладкости (Липшицевости градиента) зависит от выбора нормы. Если $f : \mathbb{R}^n \rightarrow \mathbb{R}$ - непрерывно дифференцируема и градиент Липшицев с константой L , то $\forall x, y \in \mathbb{R}^n$:

$$\|f(y) - f(x) - \langle \nabla f(x), y - x \rangle\| \leq \frac{L}{2}\|y - x\|^2$$

Рис. 2: Иллюстрация Липшицевых парабол, между которыми зажата гладкая функция. Чаще нас интересует мажорирующая из них.

Гладкость

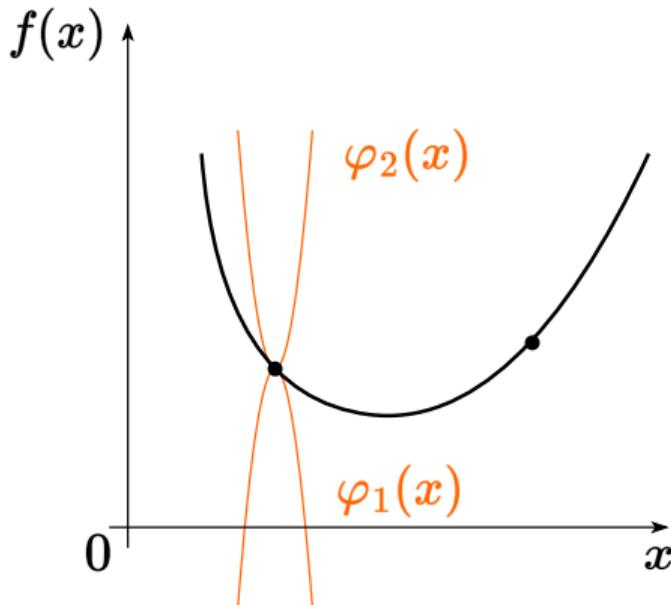


Рис. 2: Иллюстрация Липшицевых парабол, между которыми зажата гладкая функция. Чаще нас интересует мажорирующая из них.

Определение: Будем говорить, что функция $f : \mathbb{R}^n \rightarrow \mathbb{R}$ является L -гладкой, если $\forall x, y \in \mathbb{R}^n$ выполнено:

$$\|\nabla f(y) - \nabla f(x)\| \leq L\|y - x\|$$

Обратим внимание, что значение константы гладкости (Липшицевости градиента) зависит от выбора нормы. Если $f : \mathbb{R}^n \rightarrow \mathbb{R}$ - непрерывно дифференцируема и градиент Липшицев с константой L , то $\forall x, y \in \mathbb{R}^n$:

$$\|f(y) - f(x) - \langle \nabla f(x), y - x \rangle\| \leq \frac{L}{2}\|y - x\|^2$$

Если зафиксируем $x_0 \in \mathbb{R}^n$, то:

$$\varphi_1(x) = f(x_0) + \langle \nabla f(x_0), x - x_0 \rangle - \frac{L}{2}\|x - x_0\|^2$$

$$\varphi_2(x) = f(x_0) + \langle \nabla f(x_0), x - x_0 \rangle + \frac{L}{2}\|x - x_0\|^2$$

Гладкость



Определение: Будем говорить, что функция $f : \mathbb{R}^n \rightarrow \mathbb{R}$ является L -гладкой, если $\forall x, y \in \mathbb{R}^n$ выполнено:

$$\|\nabla f(y) - \nabla f(x)\| \leq L\|y - x\|$$

Обратим внимание, что значение константы гладкости (Липшицевости градиента) зависит от выбора нормы. Если $f : \mathbb{R}^n \rightarrow \mathbb{R}$ - непрерывно дифференцируема и градиент Липшицев с константой L , то $\forall x, y \in \mathbb{R}^n$:

$$\|f(y) - f(x) - \langle \nabla f(x), y - x \rangle\| \leq \frac{L}{2}\|y - x\|^2$$

Если зафиксируем $x_0 \in \mathbb{R}^n$, то:

$$\varphi_1(x) = f(x_0) + \cancel{\langle \nabla f(x_0), x - x_0 \rangle} - \frac{L}{2}\|x - x_0\|^2$$

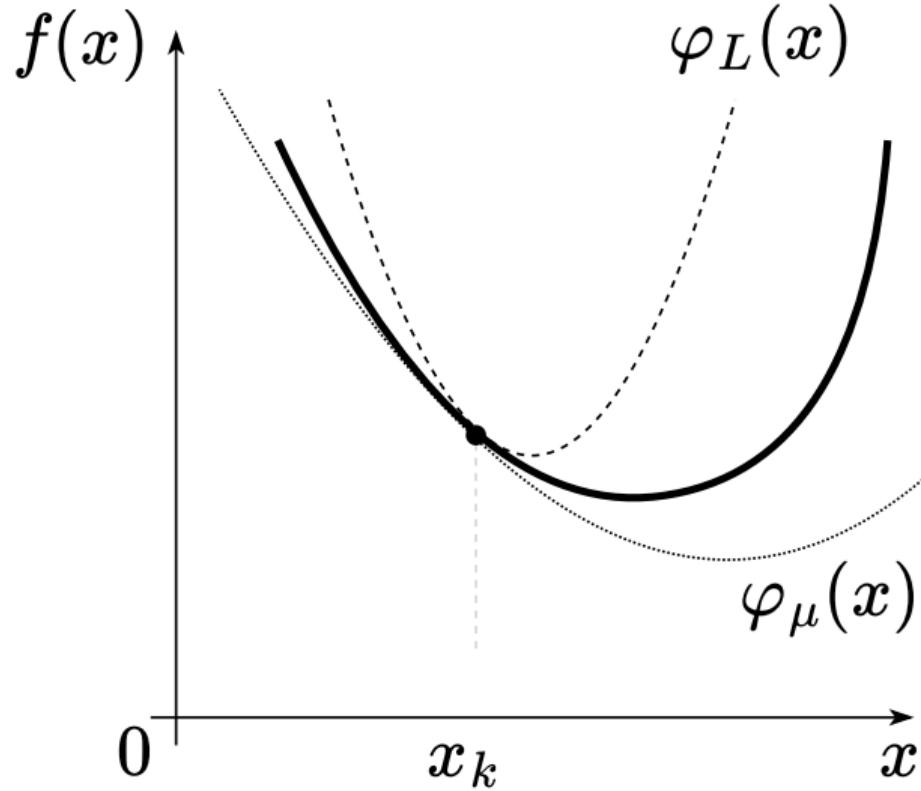
$$\varphi_2(x) = f(x_0) + \cancel{\langle \nabla f(x_0), x - x_0 \rangle} + \frac{L}{2}\|x - x_0\|^2$$

Рис. 2: Иллюстрация Липшицевых парабол, между которыми зажата гладкая функция. Чаще нас интересует мажорирующая из них.

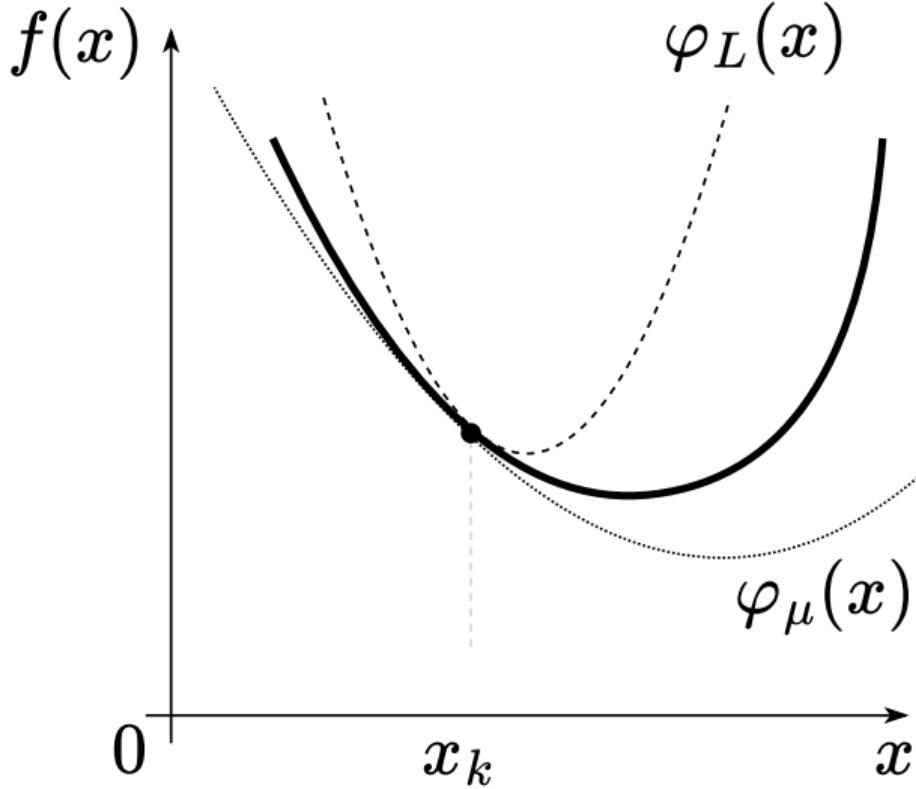
Это две параболы, и для них верно, что

$$\varphi_1(x) \leq f(x) \leq \varphi_2(x) \quad \forall x$$

Гладкость и сильная выпуклость



Гладкость и сильная выпуклость



$$A \in S_{++}^n$$

$$x \in \mathbb{R}^n$$

Пример:

$$f(x) = \frac{1}{2} x^T A x$$

$$\nabla f = Ax$$

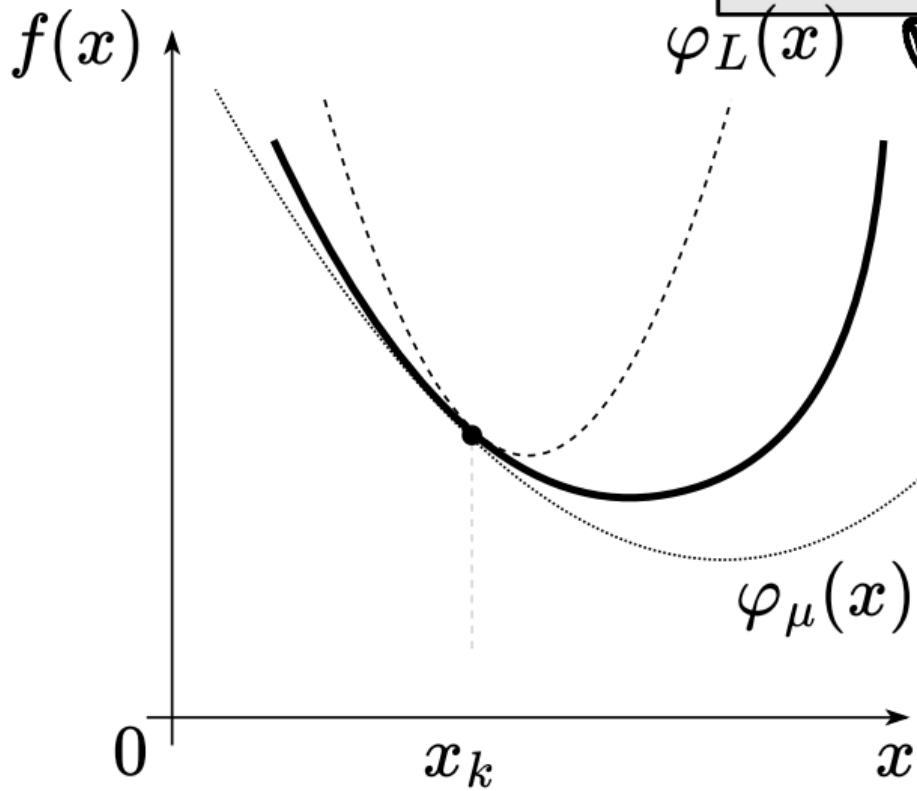
$$\nabla^2 f = A$$

для сильн. вып.

$$\mu \cdot I \leq \nabla^2 f \leq L \cdot I$$

$$\mu \cdot I \leq A \leq L \cdot I$$

Гладкость и сильная выпуклость



$$\lambda_{\min}(A) \geq \mu$$

$$\mu \cdot I \leq A$$

$$A - \mu \cdot I \succeq 0$$

$$(\dots) - \mu (\dots) \succeq 0$$

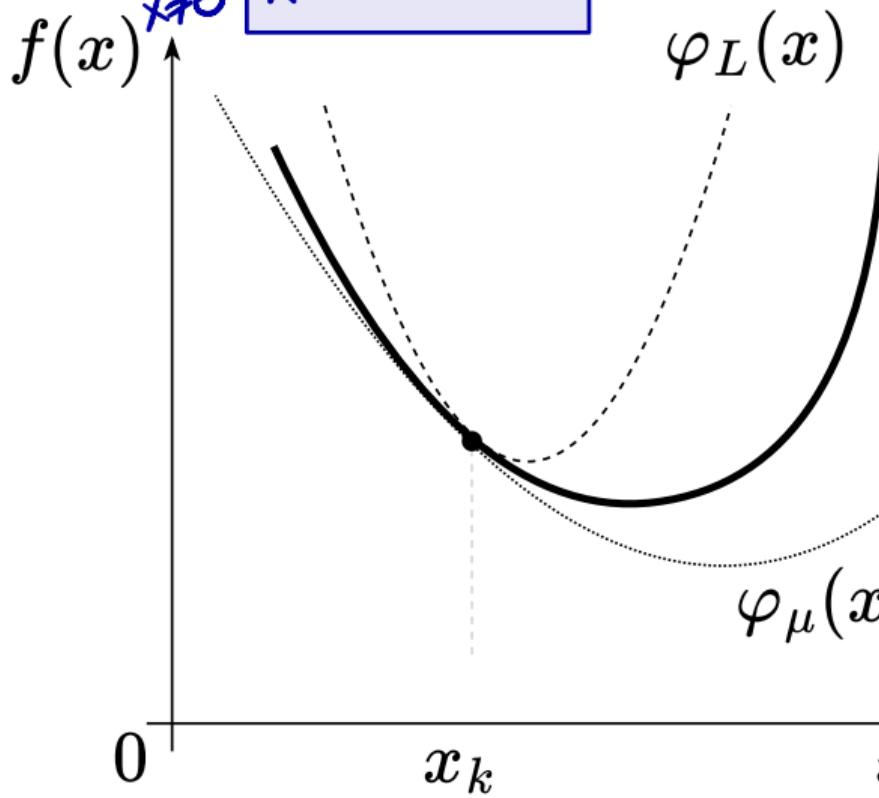
$$\begin{pmatrix} \bullet & -\mu & \circ \\ & \ddots & \\ \circ & & -\mu \end{pmatrix} \succeq 0$$

$$\lambda(\dots) \geq 0$$

$$\lambda_{\min}(\dots) \geq 0$$

Гладкость и сильная выпуклость

$$\forall x \in \mathbb{R}^n_{x \neq 0} \quad x^T A x > 0$$



где гладкость

если

$$\lambda_{\max}(A) \leq L$$

$$\text{TO } f(x) = \frac{1}{2} x^T A x$$

- L-гладкая

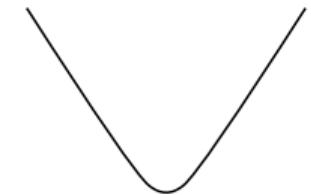
$$A = \begin{pmatrix} 1 & 0 \\ 0 & 2 \end{pmatrix}$$

$$\lambda_{\min} = 1$$
$$\lambda_{\max} = 2$$

$$f(x) = \frac{1}{2} x^T A x$$

$$\mu = 1$$
$$L = 2$$

Гладкость и сильная выпуклость



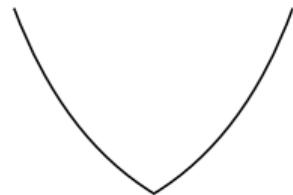
Гладкая
Выпуклая



Гладкая
 μ - сильно выпуклая



Негладкая
Выпуклая



Негладкая
 μ - сильно выпуклая

Градиентный спуск

Направление локального наискорейшего спуска

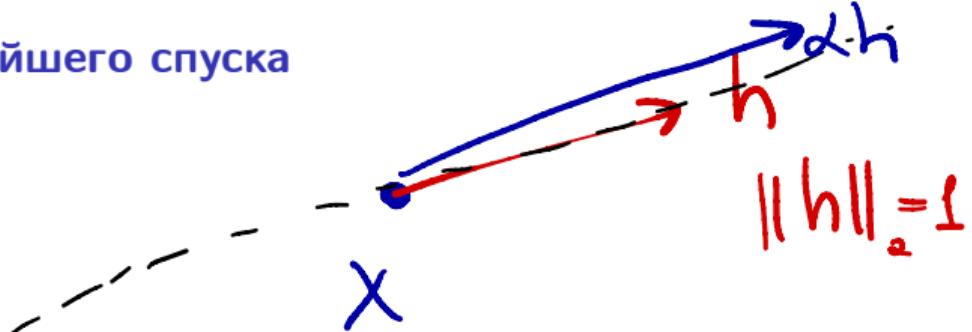
Рассмотрим линейное приближение
дифференцируемой функции f вдоль
направления h , где $\|h\|_2 = 1$:

Направление локального наискорейшего спуска

Рассмотрим линейное приближение дифференцируемой функции f вдоль направления h , где $\|h\|_2 = 1$:

$$f(x + \alpha h) = f(x) + \underbrace{\alpha \langle \nabla f(x), h \rangle}_{f_x^I(x + \alpha h)} + o(\alpha)$$

$$f_x^I(x + \alpha h)$$



Направление локального наискорейшего спуска

Рассмотрим линейное приближение дифференцируемой функции f вдоль направления h , где $\|h\|_2 = 1$:

$$f(x + \alpha h) = f(x) + \underbrace{\alpha \langle \nabla f(x), h \rangle}_{\text{линейное приближение}} + o(\alpha)$$

Хотим, чтобы h было направлением убывания:

$$f(x + \alpha h) - f(x) < 0$$

убывание

$$\alpha \langle \nabla f(x), h \rangle + o(\alpha) < 0$$

$\downarrow : \mathbb{R} \rightarrow \mathbb{R}$

$$\cdot \langle \nabla f(x), h \rangle$$

Направление локального наискорейшего спуска

Рассмотрим линейное приближение дифференцируемой функции f вдоль направления h , где $\|h\|_2 = 1$:

$$f(x + \alpha h) = f(x) + \alpha \langle \nabla f(x), h \rangle + o(\alpha)$$

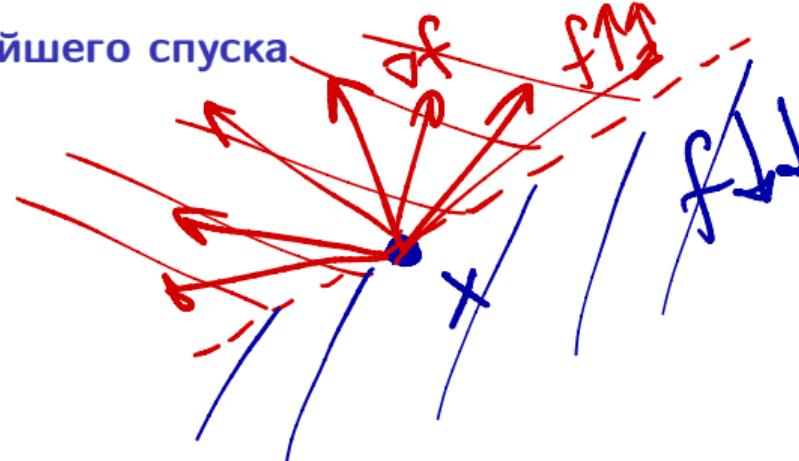
Хотим, чтобы h было направлением убывания:

$$f(x + \alpha h) - f(x) < 0$$

$$\alpha \langle \nabla f(x), h \rangle + o(\alpha) < 0$$

Переходя к пределу при $\alpha \rightarrow 0$:

$$\langle \nabla f(x), h \rangle < 0$$



если h со направлением ∇f $\langle \nabla f, h \rangle \geq 0$ $f \uparrow \uparrow$

если h противоположен ∇f $\langle \nabla f, h \rangle \leq 0$ $f \downarrow \downarrow$

Направление локального наискорейшего спуска

Рассмотрим линейное приближение дифференцируемой функции f вдоль направления h , где $\|h\|_2 = 1$:

$$f(x + \alpha h) = f(x) + \alpha \langle \nabla f(x), h \rangle + o(\alpha)$$

Хотим, чтобы h было направлением убывания:

$$f(x + \alpha h) - f(x) < 0$$

$$\alpha \langle \nabla f(x), h \rangle + o(\alpha) < 0$$

Переходя к пределу при $\alpha \rightarrow 0$:

$$\langle \nabla f(x), h \rangle < 0$$

Более того, мы хотим, чтобы разница $f(x) - f(x + \alpha h)$ была максимальна:

$$h = \arg \max_h (-\langle \nabla f(x), h \rangle) = \arg \min_h \langle \nabla f(x), h \rangle.$$

Также из неравенства Коши–Буняковского получаем:

$$|\langle \nabla f(x), h \rangle| \leq \|\nabla f(x)\|_2 \|h\|_2$$

$$\langle \nabla f(x), h \rangle \geq \underbrace{-\|\nabla f(x)\|_2 \|h\|_2}_{-\|\nabla f(x)\|_2} = -\|\nabla f(x)\|_2$$

3 • 1

$$-3 \leq \langle \nabla f(x), h \rangle \leq 3$$

Направление локального наискорейшего спуска

Рассмотрим линейное приближение дифференцируемой функции f вдоль направления h , где $\|h\|_2 = 1$:

$$f(x + \alpha h) = f(x) + \alpha \langle \nabla f(x), h \rangle + o(\alpha)$$

Хотим, чтобы h было направлением убывания:

$$\begin{cases} f(x + \alpha h) - f(x) < 0 \\ \alpha \langle \nabla f(x), h \rangle + o(\alpha) < 0 \end{cases}$$

Переходя к пределу при $\alpha \rightarrow 0$:

$$\langle \nabla f(x), h \rangle < 0$$

Более того, мы хотим, чтобы разница $f(x) - f(x + \alpha h)$ была максимальна:

$$h = \arg \max_h (-\langle \nabla f(x), h \rangle) = \arg \min_h \langle \nabla f(x), h \rangle.$$

Также из неравенства Коши–Буняковского получаем:

$$\begin{aligned} |\langle \nabla f(x), h \rangle| &\leq \|\nabla f(x)\|_2 \|h\|_2 \\ \langle \nabla f(x), h \rangle &\geq -\|\nabla f(x)\|_2 \|h\|_2 = -\|\nabla f(x)\|_2 \end{aligned}$$

Таким образом, направление антиградиента

$$h = \arg \min_h \langle \nabla f(x), h \rangle = -\frac{\nabla f(x)}{\|\nabla f(x)\|_2}$$

представляет собой направление **наискорейшего локального убывания** функции f .

Направление локального наискорейшего спуска

Рассмотрим линейное приближение дифференцируемой функции f вдоль направления h , где $\|h\|_2 = 1$:

$$f(x + \alpha h) = f(x) + \alpha \langle \nabla f(x), h \rangle + o(\alpha)$$

Хотим, чтобы h было направлением убывания:

$$f(x + \alpha h) - f(x) < 0$$

$$\alpha \langle \nabla f(x), h \rangle + o(\alpha) < 0$$

Переходя к пределу при $\alpha \rightarrow 0$:

$$\langle \nabla f(x), h \rangle < 0$$

Более того, мы хотим, чтобы разница $f(x) - f(x + \alpha h)$ была максимальна:

$$h = \arg \max_h (-\langle \nabla f(x), h \rangle) = \arg \min_h \langle \nabla f(x), h \rangle.$$

Также из неравенства Коши–Буняковского получаем:

$$\begin{aligned} |\langle \nabla f(x), h \rangle| &\leq \|\nabla f(x)\|_2 \|h\|_2 \\ \langle \nabla f(x), h \rangle &\geq -\|\nabla f(x)\|_2 \|h\|_2 = -\|\nabla f(x)\|_2 \end{aligned}$$

Таким образом, направление антиградиента

$$h = \arg \min_h \langle \nabla f(x), h \rangle = -\frac{\nabla f(x)}{\|\nabla f(x)\|_2}$$

представляет собой направление **наискорейшего локального убывания** функции f .

Итерация метода имеет вид:

$$x^{k+1} = x^k - \alpha \nabla f(x^k)$$

размер шага

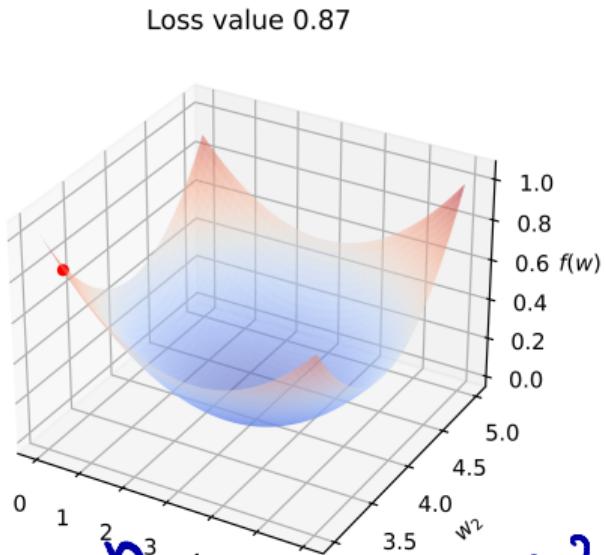
learning
rate

Сходимость алгоритма градиентного спуска

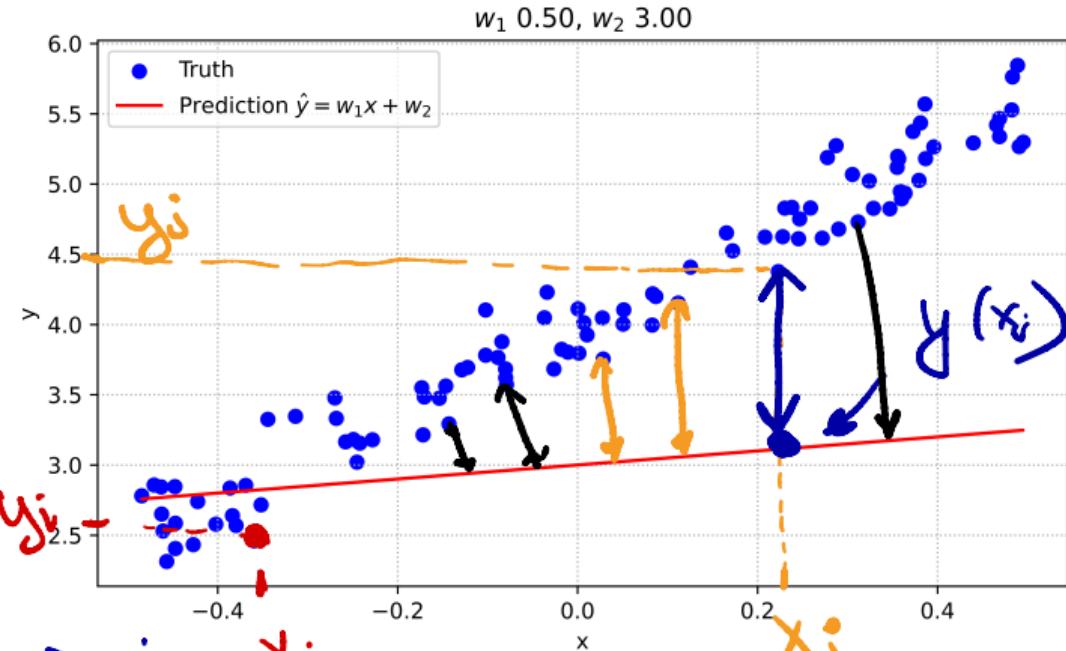
$$y = w_0 + w_1 \cdot x$$

$$y(x_i) = ?$$

Код для построения анимации ниже. Сходимость существенно зависит от выбора шага α :



$$\sum_{i=1}^n (y(x_i) - y_i)^2 \rightarrow \min_{w_0, w_1}$$

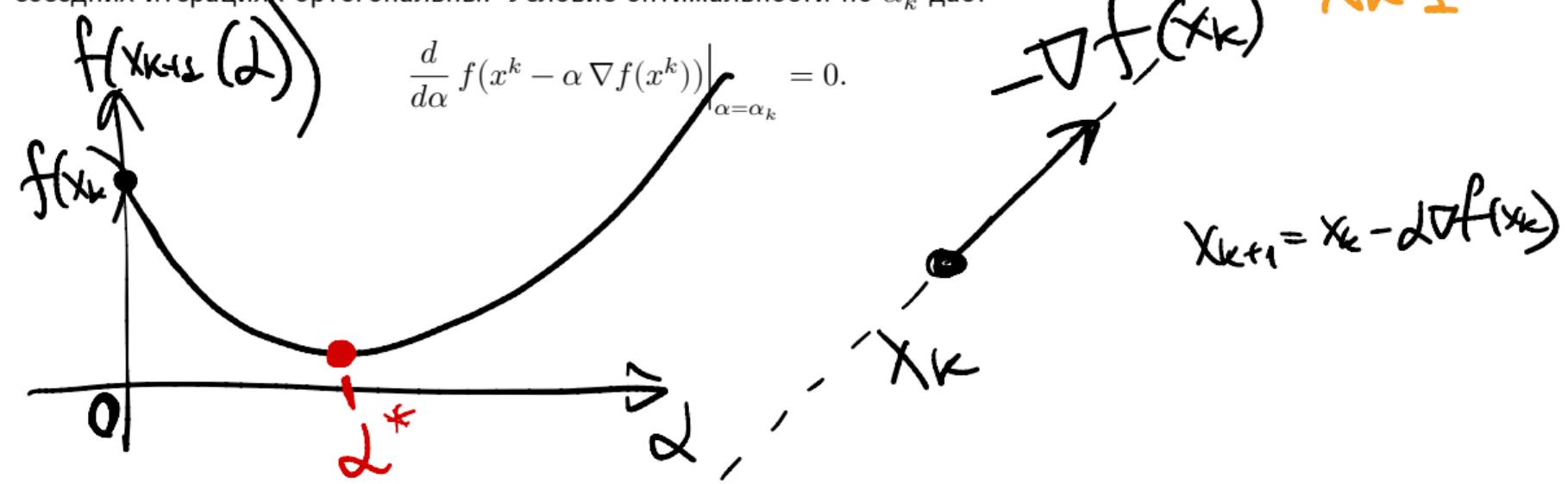


Точный линейный поиск (метод наискорейшего спуска)

$$\alpha_k = \arg \min_{\alpha \in \mathbb{R}^+} f(x^k - \alpha \nabla f(x^k)) = \underset{\alpha \in \mathbb{R}^+}{\operatorname{argmin}} f(x_{k+1})$$

Подход скорее теоретический, чем практический: он удобен для анализа сходимости, но точный линейный поиск часто затруднён, если вычисление функции занимает слишком много времени или стоит слишком дорого.

Интересное теоретическое свойство этого метода заключается в том, что градиенты на соседних итерациях ортогональны. Условие оптимальности по α_k даёт



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

Точный линейный поиск (метод наискорейшего спуска)

$$\alpha_k = \arg \min_{\alpha \in \mathbb{R}^+} f(x^k - \alpha \nabla f(x^k))$$

Подход скорее теоретический, чем практический: он удобен для анализа сходимости, но точный линейный поиск часто затруднён, если вычисление функции занимает слишком много времени или стоит слишком дорого.

Интересное теоретическое свойство этого метода заключается в том, что градиенты на соседних итерациях ортогональны. Условие оптимальности по α_k даёт

$$\frac{d}{d\alpha} f(x^k - \alpha \nabla f(x^k)) \Big|_{\alpha=\alpha_k} = 0.$$

$$\frac{\partial f^T}{\partial x} \cdot \left(\frac{\partial x}{\partial \alpha} \right) = 0$$

$$-\left(\nabla f(x_{k+1}) \right)^T \cdot \nabla f(x_k) = 0$$

$$\frac{\partial X_{k+1}}{\partial \alpha} = \frac{\partial (x_k - \alpha \nabla f(x_k))}{\partial \alpha} = -\nabla f(x_k)$$

Точный линейный поиск (метод наискорейшего спуска)

$$\alpha_k = \arg \min_{\alpha \in \mathbb{R}^+} f(x^k - \alpha \nabla f(x^k))$$

Подход скорее теоретический, чем практический: он удобен для анализа сходимости, но точный линейный поиск часто затруднён, если вычисление функции занимает слишком много времени или стоит слишком дорого.

Интересное теоретическое свойство этого метода заключается в том, что градиенты на соседних итерациях ортогональны. Условие оптимальности по α_k даёт

$$\frac{d}{d\alpha} f(x^k - \alpha \nabla f(x^k)) \Big|_{\alpha=\alpha_k} = 0.$$

Условия оптимальности:

Точный линейный поиск (метод наискорейшего спуска)

$$\alpha_k = \arg \min_{\alpha \in \mathbb{R}^+} f(x^k - \alpha \nabla f(x^k))$$

Подход скорее теоретический, чем практический: он удобен для анализа сходимости, но точный линейный поиск часто затруднён, если вычисление функции занимает слишком много времени или стоит слишком дорого.

Интересное теоретическое свойство этого метода заключается в том, что градиенты на соседних итерациях ортогональны. Условие оптимальности по α_k даёт

$$\frac{d}{d\alpha} f(x^k - \alpha \nabla f(x^k)) \Big|_{\alpha=\alpha_k} = 0.$$

Условия оптимальности:

$$\nabla f(x^{k+1})^\top \nabla f(x^k) = 0$$

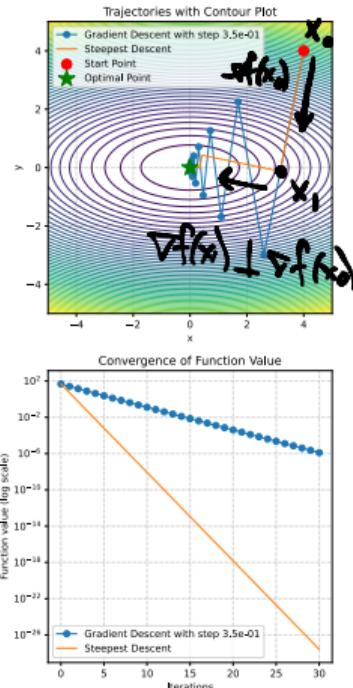


Рис. 3: Наискорейший спуск

Открыть в Colab ♣

$$f(x) = \frac{1}{2} x^T A x$$

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

$$A = A^T$$

$$\nabla f(x) = Ax$$

$$: \nabla f(x_{k+1})^T \nabla f(x_k) = 0$$

$$(Ax_{k+1})^T \cdot Ax_k = 0$$

Сильно выпуклые квадратичные функции

$$(A(x_k - \alpha A x_k))^T A x_k = 0$$

$$(Ax_k - \alpha A A x_k)^T A x_k = 0$$

$$(x_k^T A^T - \alpha x_k^T A^T A^T) A x_k = 0$$

$$\begin{aligned} \nabla f(x_k) &= \\ &= g_k \end{aligned}$$

$$g_k^T g_k - \alpha g_k^T A^T g_k = 0$$

$$\alpha = \frac{g_k^T g_k}{g_k^T A^T g_k}$$