

1. БАЗА

2. Теор.

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

learning rate

Gradient Descent

Daniil Merkulov

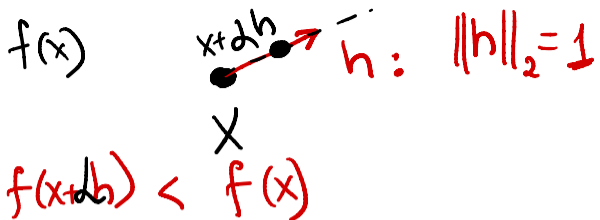
Optimization methods. MIPT

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

$n \times 1$ $n \times 1$ 1×1 $n \times 1$

Direction of local steepest descent

Let's consider a linear approximation of the differentiable function f along some direction h , $\|h\|_2 = 1$:



$f(x)$

$x+h$

$h: \|h\|_2 = 1$

$f(x+h) < f(x)$

Direction of local steepest descent

Let's consider a linear approximation of the differentiable function f along some direction

h , $\|h\|_2 = 1$:

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

Direction of local steepest descent

Let's consider a linear approximation of the differentiable function f along some direction h , $\|h\|_2 = 1$:

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

We want h to be a decreasing direction:

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

Треб. убыв.

$$\underline{f(x) + \alpha \langle f'(x), h \rangle + o(\alpha) < f(x)}$$

Direction of local steepest descent

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

Let's consider a linear approximation of the differentiable function f along some direction h , $\|h\|_2 = 1$:

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

We want h to be a decreasing direction:

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

$$\cancel{f(x)} + \alpha \langle f'(x), h \rangle + o(\alpha) < \cancel{f(x)}$$

and going to the limit at $\alpha \rightarrow 0$:

$$\langle f'(x), h \rangle \leq 0$$

Direction of local steepest descent

Let's consider a linear approximation of the differentiable function f along some direction h , $\|h\|_2 = 1$:

$$\boxed{f(x + \alpha h)} = f(x) + \underbrace{\alpha \langle f'(x), h \rangle}_{\leq 0} + o(\alpha)$$

We want h to be a decreasing direction:

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

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

and going to the limit at $\alpha \rightarrow 0$:

$$\langle f'(x), h \rangle \leq 0$$

$$h = -\frac{\nabla f}{\|\nabla f\|}$$

Also from Cauchy-Bunyakovsky-Schwarz inequality:

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

$$\langle f'(x), h \rangle \geq -\|f'(x)\|_2 \|h\|_2 = -\|f'(x)\|_2$$

$$\max -\langle \nabla f, h \rangle \Leftrightarrow \min \langle \nabla f, h \rangle$$

$\|h\|_2 = 1$

$$L = \nabla f^T h + \lambda (h^T h - 1)$$

$$\frac{\partial L}{\partial h} = \nabla f + 2\lambda h \Rightarrow h = -\frac{1}{2\lambda} \nabla f$$

$$\frac{\partial L}{\partial \lambda} \Rightarrow \|h\|_2^2 = 1 \quad \frac{1}{4\lambda^2} \|\nabla f\|^2 = 1$$

$$\Rightarrow 2\lambda = \|\nabla f\|$$

Direction of local steepest descent

Let's consider a linear approximation of the differentiable function f along some direction h , $\|h\|_2 = 1$:

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

We want h to be a decreasing direction:

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

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

and going to the limit at $\alpha \rightarrow 0$:

$$\langle f'(x), h \rangle \leq 0$$

Also from Cauchy–Bunyakovsky–Schwarz inequality:

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

$$\langle f'(x), h \rangle \geq -\|f'(x)\|_2 \|h\|_2 = -\|f'(x)\|_2$$

Thus, the direction of the antigradient

$$h = -\frac{f'(x)}{\|f'(x)\|_2}$$

gives the direction of the **steepest local** decreasing of the function f .

$$\begin{aligned} \max_{\|h\|_2=1} \langle \nabla f(x), h \rangle &\Leftrightarrow \min_{\|h\|_2=1} \langle -\nabla f(x), h \rangle \\ L &= -\nabla f(x)^T h + \lambda (\|h\|_2^2 - 1) \\ \frac{\partial L}{\partial h} &= -\nabla f + 2\lambda h \\ h &= \nabla f \frac{1}{2\lambda} \end{aligned}$$

Direction of local steepest descent

Let's consider a linear approximation of the differentiable function f along some direction h , $\|h\|_2 = 1$:

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

We want h to be a decreasing direction:

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

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

and going to the limit at $\alpha \rightarrow 0$:

$$\langle f'(x), h \rangle \leq 0$$

Also from Cauchy–Bunyakovsky–Schwarz inequality:

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

$$\langle f'(x), h \rangle \geq -\|f'(x)\|_2 \|h\|_2 = -\|f'(x)\|_2$$

Thus, the direction of the antigradient

$$h = -\frac{f'(x)}{\|f'(x)\|_2}$$

gives the direction of the **steepest local** decreasing of the function f .
The result of this method is

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

Gradient flow ODE

Let's consider the following ODE, which is referred as Gradient Flow equation.

$$\frac{dx}{dt} = -f'(x(t))$$

(GF)

Gradient flow ODE

Let's consider the following ODE, which is referred as Gradient Flow equation.

$$\frac{dx}{dt} = -f'(x(t))$$

and discretize it on a uniform grid with α step:

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

дискретизация (GF)
Эйлера

Gradient flow ODE

Let's consider the following ODE, which is referred as Gradient Flow equation.

$$\frac{dx}{dt} = -f'(x(t)) \quad (\text{GF})$$

and discretize it on a uniform grid with α step:

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

where $x_k \equiv x(t_k)$ and $\alpha = t_{k+1} - t_k$ - is the grid step.

From here we get the expression for x_{k+1}

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

which is exactly gradient descent.

Open In Colab 

гукретауауауа
түнера

ODE
GF

Gradient flow ODE

Let's consider the following ODE, which is referred as Gradient Flow equation.

$$\frac{dx}{dt} = -f'(x(t))$$

and discretize it on a uniform grid with α step:

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

where $x_k \equiv x(t_k)$ and $\alpha = t_{k+1} - t_k$ - is the grid step.

From here we get the expression for x_{k+1}

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

which is exactly gradient descent.

Open In Colab 

(GF)

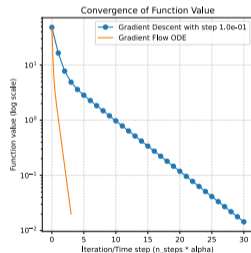
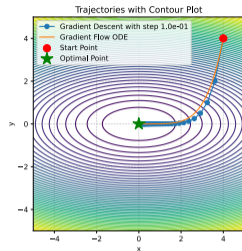



Figure 1: Gradient flow trajectory

Necessary local minimum condition

$$\begin{array}{l} \boxed{f'(x) = 0} \\ \rightarrow -\eta f'(x) = 0 \quad + \times \\ x - \eta f'(x) = x \\ \hline x_k - \eta f'(x_k) = x_{k+1} \\ \hline \quad \leftarrow \quad \quad \quad \leftarrow k+1 \end{array}$$

Minimizer of Lipschitz parabola

If a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is continuously differentiable and its gradient satisfies Lipschitz conditions with constant L , then $\forall x, y \in \mathbb{R}^n$:

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


$\nabla f - L$ normunges

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

Minimizer of Lipschitz parabola

If a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is continuously differentiable and its gradient satisfies Lipschitz conditions with constant L , then $\forall x, y \in \mathbb{R}^n$:

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

which geometrically means, that if we'll fix some point $x_0 \in \mathbb{R}^n$ and define two parabolas:

$$\phi_1(x) = f(x_0) + \langle \nabla f(x_0), x - x_0 \rangle - \frac{L}{2} \|x - x_0\|^2,$$

$$\phi_2(x) = f(x_0) + \langle \nabla f(x_0), x - x_0 \rangle + \frac{L}{2} \|x - x_0\|^2.$$

Minimizer of Lipschitz parabola

If a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is continuously differentiable and its gradient satisfies Lipschitz conditions with constant L , then $\forall x, y \in \mathbb{R}^n$:

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

which geometrically means, that if we'll fix some point $x_0 \in \mathbb{R}^n$ and define two parabolas:

$$\phi_1(x) = f(x_0) + \langle \nabla f(x_0), x - x_0 \rangle - \frac{L}{2} \|x - x_0\|^2,$$

$$\phi_2(x) = f(x_0) + \langle \nabla f(x_0), x - x_0 \rangle + \frac{L}{2} \|x - x_0\|^2.$$

Then

$$\phi_1(x) \leq f(x) \leq \phi_2(x) \quad \forall x \in \mathbb{R}^n.$$

Minimizer of Lipschitz parabola

If a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is continuously differentiable and its gradient satisfies Lipschitz conditions with constant L , then $\forall x, y \in \mathbb{R}^n$:

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

which geometrically means, that if we'll fix some point $x_0 \in \mathbb{R}^n$ and define two parabolas:

$$\phi_1(x) = f(x_0) + \langle \nabla f(x_0), x - x_0 \rangle - \frac{L}{2} \|x - x_0\|^2,$$

$$\phi_2(x) = f(x_0) + \langle \nabla f(x_0), x - x_0 \rangle + \frac{L}{2} \|x - x_0\|^2.$$

Then

$$\phi_1(x) \leq f(x) \leq \phi_2(x) \quad \forall x \in \mathbb{R}^n.$$

Now, if we have global upper bound on the function, in a form of parabola, we can try to go directly to its minimum.

Minimizer of Lipschitz parabola

If a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is continuously differentiable and its gradient satisfies Lipschitz conditions with constant L , then $\forall x, y \in \mathbb{R}^n$:

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

which geometrically means, that if we'll fix some point $x_0 \in \mathbb{R}^n$ and define two parabolas:

$$\phi_1(x) = f(x_0) + \langle \nabla f(x_0), x - x_0 \rangle - \frac{L}{2} \|x - x_0\|^2,$$

$$\phi_2(x) = f(x_0) + \langle \nabla f(x_0), x - x_0 \rangle + \frac{L}{2} \|x - x_0\|^2.$$

Then

$$\phi_1(x) \leq f(x) \leq \phi_2(x) \quad \forall x \in \mathbb{R}^n.$$

Now, if we have global upper bound on the function, in a form of parabola, we can try to go directly to its minimum.

L-smooth function

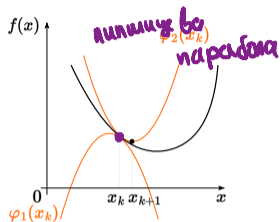


Figure 2: Illustration

$$\nabla \phi_2(x_{k+1}) = 0$$

Minimizer of Lipschitz parabola

If a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is continuously differentiable and its gradient satisfies Lipschitz conditions with constant L , then $\forall x, y \in \mathbb{R}^n$:

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

which geometrically means, that if we'll fix some point $x_0 \in \mathbb{R}^n$ and define two parabolas:

$$\phi_1(x) = f(x_0) + \langle \nabla f(x_0), x - x_0 \rangle - \frac{L}{2} \|x - x_0\|^2,$$

$$\phi_2(x) = f(x_0) + \langle \nabla f(x_0), x - x_0 \rangle + \frac{L}{2} \|x - x_0\|^2.$$

Then

$$\underbrace{\nabla \phi_2}_{0} + \nabla f(x_0) + L \cdot (x - x_0)$$

$$\phi_1(x) \leq f(x) \leq \phi_2(x) \quad \forall x \in \mathbb{R}^n.$$

Now, if we have global upper bound on the function, in a form of parabola, we can try to go directly to its minimum.

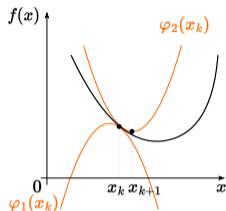


Figure 2: Illustration

$$\nabla \phi_2(x) = 0$$

$$\nabla f(x_0) + L(x^* - x_0) = 0$$

$$x^* = x_0 - \frac{1}{L} \nabla f(x_0)$$

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

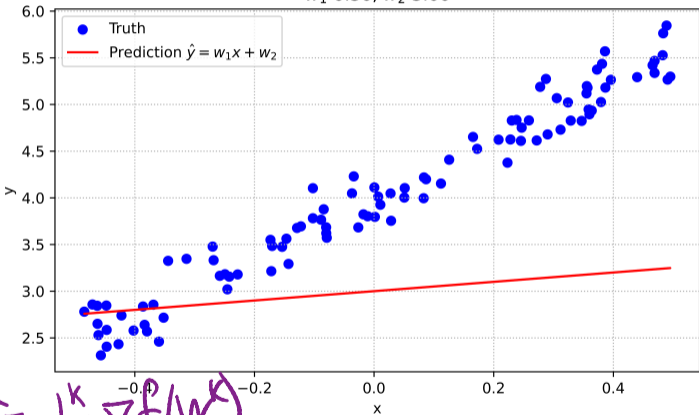
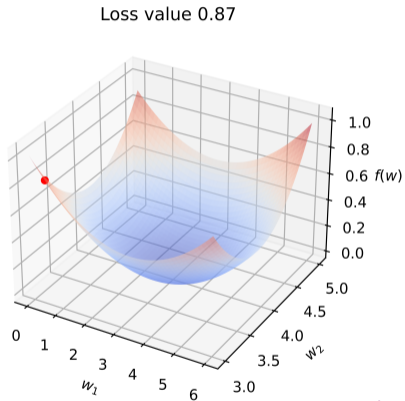
This way leads to the $\frac{1}{L}$ stepsize choosing. However, often the L constant is not known.

Convergence of Gradient Descent algorithm

Heavily depends on the choice of the learning rate α :

$$y = w_1 x + w_2$$

$w_1 = 0.50, w_2 = 3.00$



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

Exact line search aka steepest descent

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

More theoretical than practical approach. It also allows you to analyze the convergence, but often exact line search can be difficult if the function calculation takes too long or costs a lot. Interesting theoretical property of this method is that each following iteration is orthogonal to the previous one:

Метод наискорейшего спуска

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

Optimality conditions:

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

Пример задачи минимизации
 $\frac{1}{2} \|Ax - b\|_2^2 \rightarrow \min_{x \in \mathbb{R}^n}$

$$\nabla f = A^\top (Ax - b)$$

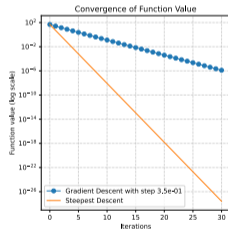
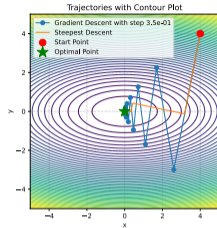


Figure 3: Steepest Descent

Open In Colab

Exact line search aka steepest descent

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

More theoretical than practical approach. It also allows you to analyze the convergence, but often exact line search can be difficult if the function calculation takes too long or costs a lot. Interesting theoretical property of this method is that each following iteration is orthogonal to the previous one:

$$\frac{\partial f}{\partial \alpha} = \frac{\partial f}{\partial x_{k+1}} \cdot \frac{\partial x_{k+1}}{\partial \alpha} = 0$$

$$\alpha_k = \arg \min_{\alpha \in \mathbb{R}^+} f(x_k - \alpha \nabla f(x_k)) \quad \left[A^T (A x_{k+1} - b) \right]^T \left[A (x_k - b) \right] = 0$$

Optimality conditions:

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

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

$$\alpha_k = \arg \min_{\alpha} f(x_{k+1}) = \arg \min_{\alpha} f(x_k - \alpha \cdot A^T (A x_k - b))$$

$$\left[A^T (A (x_k - \alpha g_k) - b) \right]^T g_k = 0$$

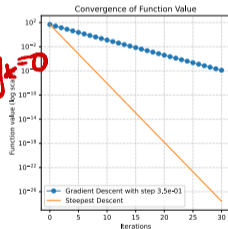
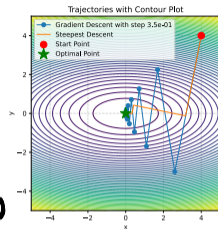


Figure 3: Steepest Descent

Open In Colab

$$[A^T(Ax_{k+1}-b)]^T [Ax_k - b] = 0$$

$$[A^T(A(x_k - \alpha g_k) - b)]^T g_k = 0$$

$$g_k^T \cdot A^T (Ax_k - \alpha Ag_k - b) = 0$$

$$g_k^T A^T (Ax_k - b - \alpha Ag_k) = 0$$

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

$$\alpha = \frac{g_k^T g_k}{g_k^T A^T A g_k} = \frac{g_k^T g_k}{g_k^T \nabla^2 f(x_k) g_k}$$

Exact line search aka steepest descent

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

More theoretical than practical approach. It also allows you to analyze the convergence, but often exact line search can be difficult if the function calculation takes too long or costs a lot. Interesting theoretical property of this method is that each following iteration is orthogonal to the previous one:

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

$$\nabla f = Ax_k$$

Optimality conditions:

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

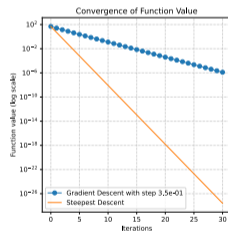
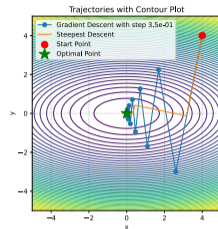


Figure 3: Steepest Descent

Open In Colab

Convergence rates

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

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

smooth

convex

smooth & convex

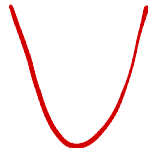
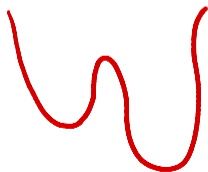
smooth & strongly convex (or PL)

$$\|\nabla f(x_k)\|^2 \approx \mathcal{O}\left(\frac{1}{k}\right)$$

$$f(x_k) - f^* \approx \mathcal{O}\left(\frac{1}{\sqrt{k}}\right)$$

$$f(x_k) - f^* \approx \mathcal{O}\left(\frac{1}{k}\right)$$

$$\|x_k - x^*\|^2 \approx \mathcal{O}\left(\left(1 - \frac{\mu}{L}\right)^k\right)$$



Gradient Descent convergence. Smooth convex case

ФАКТЫ 1 Пусть $f(x): \mathbb{R}^n \rightarrow \mathbb{R}^m$

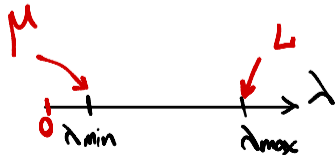
Тогда f - L -Липшицева iff f

$$\forall x \in \mathbb{R}^n : \|Df(x)\| \leq L$$

эквивалентно
если ∇f - L -Липшицев

$$\|\nabla^2 f(x)\| \leq L$$

$$\lambda_{\max}(\nabla^2 f(x)) \leq L$$



2 Если $f: \mathbb{R}^n \rightarrow \mathbb{R}$ - L -выпуклая
(Липшицев $\nabla f \in \text{кон. } L$)

то $\forall x, y \in \mathbb{R}^n : f(y) \leq f(x) + \langle \nabla f(x), y-x \rangle + \frac{L}{2} \|y-x\|^2$

3 Сходимость GD в выпуклом выпуклом случае

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

т.к. f - L -выпуклая

$$f(x^{k+1}) \leq f(x^k) + \langle \nabla f(x^k), -d \nabla f(x^k) \rangle + \frac{L}{2} d^2 \|\nabla f(x^k)\|^2$$

$$f(x^{k+1}) \leq f(x^k) - d \|\nabla f(x^k)\|^2 + \frac{L}{2} d^2 \|\nabla f(x^k)\|^2$$

$$f(x^{k+1}) \leq f(x^k) + \left(\frac{L}{2} d^2 - d\right) \|\nabla f(x^k)\|^2$$

Gradient Descent convergence. Smooth convex case

2) $f(x)$ - выпуклая

$$f(y) \geq f(x) + \nabla f(x)^T (y-x)$$

$$\frac{1}{2} \alpha^2 - \alpha \rightarrow \min_{\alpha} \rightarrow \alpha^{\text{opt}} = \frac{1}{L}$$

$$\rightarrow f(x^{k+1}) \leq f(x^k) - \frac{1}{2L} \|\nabla f(x^k)\|^2$$

↑ *уменьшение*

$$f(x^k) - f(x^{k+1}) \geq \frac{1}{2L} \|\nabla f(x^k)\|^2$$

$$\begin{cases} y = x^* \\ x = x^k \end{cases} \rightarrow \begin{aligned} f(x^*) &\geq f(x^k) + \nabla f(x^k)^T (x^* - x^k) \\ f(x^k) &\leq f(x^*) + \nabla f(x^*)^T (x^k - x^*) \end{aligned}$$

$$f(x^k) - f(x^*) \leq \nabla f(x^*)^T (x^k - x^*) \quad \text{conv}$$

$$\begin{aligned} f(x^{k+1}) &\leq f(x^k) - \frac{1}{2L} \|\nabla f(x^k)\|^2 \leq \\ &\leq f(x^*) + \nabla f(x^k)^T (x^k - x^*) - \frac{1}{2L} \|\nabla f(x^k)\|^2 = \\ &= f(x^*) + \frac{1}{2} (\|x^k - x^*\|^2 - \|x^k - x^*\| - \frac{1}{L} \|\nabla f(x^k)\|^2) \end{aligned}$$

$a^T a - b^T b = (a-b)^T (a+b)$

$$\nabla f(x^k)^T \left(x^k - x^* - \frac{1}{2L} \nabla f(x^k) \right)$$

$$\left[x^k - x^* - \frac{1}{L} \nabla f(x^k) \right]^T \left[x^k - x^* + \left(x^k - x^* - \frac{1}{L} \nabla f(x^k) \right) \right]$$

Gradient Descent convergence. Smooth convex case

$$\hookrightarrow \nabla f(x^*)^T [x^k - x^* - \frac{1}{2L} \nabla f(x^*)] =$$

$$\ominus f(x^k) + \frac{L}{2} (\|x^k - x^*\|^2 - \|x^{k+1} - x^*\|^2)$$

$$\sum_{k=0}^{T-1} f(x^{k+1}) - f(x^*) \leq \frac{L}{2} (\|x^0 - x^*\|^2 - \|x^T - x^*\|^2)$$

$$\sum_{k=0}^{T-1} f(x^{k+1}) - f^* \leq \frac{L}{2} (\|x^0 - x^*\|^2 - \|x^T - x^*\|^2) \leq$$

$$\sum_{k=0}^{T-1} (f(x^{k+1}) - f^*) \leq \frac{LR^2}{2} \quad | : T$$

$$\sum_{k=0}^{T-1} \left(\frac{1}{T} f(x^{k+1}) - \frac{1}{T} f^* \right) \leq \frac{LR^2}{2T}$$

$$f\left(\frac{1}{T} \sum_{k=0}^{T-1} x^{k+1}\right) \leq \frac{1}{T} \sum_{k=0}^{T-1} f(x^{k+1})$$

$$\sum_{k=0}^{T-1} \frac{1}{T} f(x^{k+1})$$

$$f(\theta_1 x_1 + \dots + \theta_T x_T) \leq \theta_1 f(x_1) + \dots + \theta_T f(x_T)$$

$$x_1, \dots, x_T$$

$$\theta_1 x_1 + \dots + \theta_T x_T$$

$$\theta_i \geq 0$$

$$\sum \theta_i = 1$$

$$\theta_i = \frac{1}{T}$$

не в
консепт
 $f(\frac{x_1}{T} + \dots + \frac{x_T}{T}) \leq \frac{1}{T} \sum f(x_i)$

$$f\left(\frac{1}{T} x_1 + \dots + \frac{1}{T} x_T\right) \leq \frac{1}{T} \sum_{i=1}^T f(x_i)$$

$$\leq \sum_{k=0}^{T-1} \left(\frac{1}{T} f(x^{k+1}) - \frac{1}{T} f^* \right) \leq \frac{LR^2}{2T}$$

$$T \cdot f(\bar{x}) \leq f(x^0) + f(x^1) + \dots + f(x^T)$$

$$f(\bar{x}) - f^* \leq \frac{LR^2}{2T}$$

сублинейно
O(1/T)

Докажем, что GD сходится сублинейно для шагких выпуклых функций

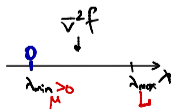
$$d = \frac{1}{2}$$

$$f(x) = \frac{1}{2} \|Ax - b\|^2$$

$$\nabla f = A^T(Ax - b) \rightarrow \nabla^2 f = A^T A$$

$$\|\nabla^2 f(x)\| \leq L$$

$$\lambda_{\max}(A^T A)$$



Gradient Descent convergence. Smooth μ -strongly convex case

Gradient Descent convergence. Polyak-Lojasiewicz case

④ Сильно выпуклый шагкий случай:

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

$$\|x^{k+1} - x^*\|^2 = \|x^k - \alpha \nabla f(x^k) - x^*\|^2 =$$

$$= \|x^k - x^*\|^2 + \alpha^2 \|\nabla f(x^k)\|^2 - 2\alpha \nabla f(x^k)^T (x^k - x^*) \leq$$

$$\leq (1 - 2\mu) \|x^k - x^*\|^2 + \alpha^2 \|\nabla f(x^k)\|^2 - 2\alpha (f(x^k) - f^*) \quad (\leq)$$

сильная выпуклость: $f(x^*) \geq f(x^k) + \nabla f(x^k)^T \cdot (x^* - x^k) + \frac{\mu}{2} \|x^k - x^*\|^2$
 $\rightarrow -2\alpha \cdot \nabla f(x^k)^T (x^k - x^*) \leq 2\alpha (f(x^*) - f(x^k)) - \alpha \mu \|x^k - x^*\|^2$

PL: $f: \mathbb{R}^n \rightarrow \mathbb{R}$ - дифференцируема, $\mu > 0$:

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

$$f(x) - f^* \leq \frac{1}{2\mu} \|\nabla f(x)\|^2$$

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

Лемма: Если $f(x)$ - мещьновыпукла, то она удовл.

PL:

$$\text{Пусть } x^* = \operatorname{argmin} f(x) \quad f^* = f(x^*)$$

Сильная выпуклость:

$$\forall x, y \in \mathbb{R}^n \quad f(y) \geq f(x) + \nabla f(x)^T (y - x) + \frac{\mu}{2} \|y - x\|^2$$

$$y = x^* \quad f(x^*) \geq f(x) + \nabla f(x)^T (x^* - x) + \frac{\mu}{2} \|x^* - x\|^2$$

$$f(x) - f(x^*) \leq \nabla f(x)^T (x - x^*) - \frac{\mu}{2} \|x - x^*\|^2 =$$

$$= \left[\nabla f(x) - \frac{\mu}{2} (x - x^*) \right]^T (x - x^*) =$$

$$= \underbrace{\left[\frac{1}{\sqrt{\mu}} \nabla f(x) - \frac{\sqrt{\mu}}{2} (x - x^*) \right]^T}_{(a-b)^T} \underbrace{\sqrt{\mu} (x - x^*)}_{a+b} \quad (\leq)$$

$(a-b)^T$

$a+b$

$$b = \sqrt{\mu}(x-x^*) - \frac{1}{\sqrt{\mu}} \nabla f(x)$$

$$a = \frac{1}{\sqrt{\mu}} \nabla f(x)$$

$$a+b = \sqrt{\mu}(x-x^*)$$

$$a-b = -\sqrt{\mu}(x-x^*) + \frac{2}{\sqrt{\mu}} \nabla f(x)$$

$$\begin{aligned} & \Leftrightarrow \frac{1}{2} \left[\frac{1}{\mu} \|\nabla f(x)\|^2 - \|\sqrt{\mu}(x-x^*) - \frac{1}{\sqrt{\mu}} \nabla f(x)\|^2 \right] \leq \\ & \leq \frac{1}{2\mu} \|\nabla f(x)\|^2 \quad \text{v.T.g.} \end{aligned}$$

$$\|\nabla f(x^k)\|^2 \geq 2\mu \cdot (f(x^k) - f^*)$$

$$\leq (1-2\mu) \|x^k - x^*\|^2 + 2L^2 \|\nabla f(x^k)\|^2 - 2L(f(x^k) - f^*) \leq$$

$$\leq (1-2\mu) \|x^k - x^*\|^2 + 2L^2 \cdot 2L(f(x^k) - f^*) - 2L(f(x^k) - f^*) =$$

$$= (1-2\mu) \|x^k - x^*\|^2 + 2L(f(x^k) - f^*) [2L - 1] =$$

$$= (1-2\mu) \|x^k - x^*\|^2 + \underbrace{2L(f(x^k) - f^*)}_{\leq 0} \underbrace{[2L - 1]}_{\geq 0} \leq$$

$$\boxed{2L \leq 1}$$

$$R^2 = \|x^0 - x^*\|^2$$

$$\leq (1-2\mu) \|x^k - x^*\|^2$$

$$1-2\mu < 1$$

$$\|x^{k+1} - x^*\|^2 \leq (1-2\mu) \|x^k - x^*\|^2$$

$$\|x^k - x^*\|^2 \leq (1-2\mu)^k \cdot R^2$$

$$\begin{aligned}
 d\mu - 1 &< 0 & 1 - d\mu &> 0 \\
 -1 < 1 - d\mu &< 1 & -1 < 1 - d\mu &< 1 \\
 & & & d \leq \frac{1}{L} \\
 & & & d = \frac{1}{L} \\
 & & & \frac{\mu}{L} < 1 \\
 0 < 1 - d\mu &< 1 & & \\
 1 - d\mu > 0 &\Rightarrow d\mu < 1 & & \\
 1 - d\mu < 1 & & & \\
 -d\mu < 0 & & & \\
 d > 0 \quad \mu > 0 & & & \\
 & & & d \leq \frac{1}{L}
 \end{aligned}$$

⑤ Сходимость GD при PL: L - шагкость

$$\begin{aligned}
 f(x^{k+1}) &\leq f(x^k) + \nabla f(x^k)^T (x^{k+1} - x^k) + \frac{L}{2} \|x^{k+1} - x^k\|^2 = \\
 &= f(x^k) + \nabla f(x^k)^T (-d \nabla f(x^k)) + \frac{L}{2} d^2 \|\nabla f(x^k)\|^2 = \\
 &= f(x^k) - d \|\nabla f(x^k)\|^2 + \frac{d^2 L}{2} \|\nabla f(x^k)\|^2 = \\
 &= f(x^k) + \|\nabla f(x^k)\|^2 \left(\frac{d^2 L}{2} - d \right) = \\
 &= f(x^k) + \frac{d}{2} \|\nabla f(x^k)\|^2 (dL - 2) = \\
 &= f(x^k) - \frac{d}{2} \|\nabla f(x^k)\|^2 (2 - dL) \leq
 \end{aligned}$$

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

$$\begin{aligned}
 d &\leq \frac{2}{L} \\
 dL &\leq 2
 \end{aligned}$$

$$\leq f(x^k) - \frac{d}{2} \cdot 2\mu (f(x^k) - f^*) (2 - dL)$$

$$f(x^{k+1}) - f^* \leq f(x^k) - f^* - d\mu (f(x^k) - f^*) (2 - dL)$$

$$f(x^{k+1}) - f(x^k) \leq [1 - d\mu(2 - dL)] \cdot (f(x^k) - f^*)$$

$$f(x^k) - f^* \leq [1 - d\mu(2 - dL)]^k \cdot (f(x^0) - f^*)$$

$$-1 < 1 - d\mu(2 - dL) < 1$$

$$-2 < -d\mu(2 - dL) < 0$$

$$0 < d\mu(2 - dL) < 2$$

$$d\mu(2 - dL) < 2$$

$$\frac{1}{L} \mu (2 - d) < 2$$

$$\frac{\mu}{L} \cdot 1 < 2$$

$$\frac{2}{L} \mu (2 - 2) < 2$$

$$2 - dL > 0$$

$$d < \frac{2}{L}$$

$$d = \frac{1}{L}$$

$$d = \frac{2}{L}$$

$$d = \frac{1}{L}$$