

1 БАЗА

2. Теор.

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

Gradient Descent

learning rate

Daniil Merkulov

Optimization methods. MIPT

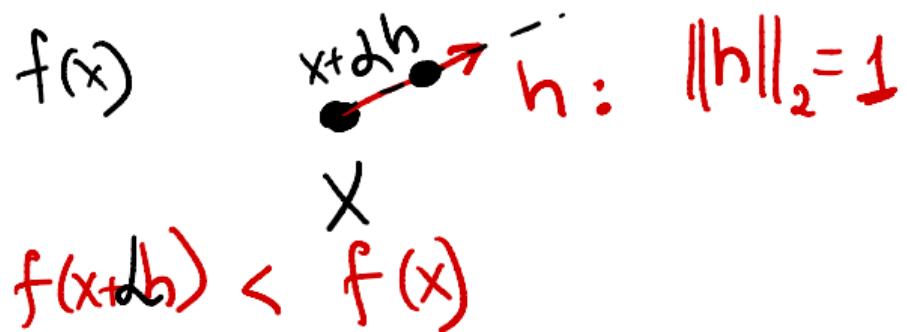
$$X_{k+1} = X_k - d_k \nabla f(X_k)$$

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



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:

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

Требуем.

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

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:

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

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

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

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

Direction of local steepest descent

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

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

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:

$$|\langle f'(x), h \rangle| \leq \|f'(x)\|_2 \|h\|_2$$
$$\langle f'(x), h \rangle \geq -\|f'(x)\|_2 \|h\|_2 = -\|f'(x)\|_2$$
$$\|h\|_2 = 1$$

$$\max -\langle \nabla f, h \rangle \Leftrightarrow \min_{\|h\|^2=1} \langle \nabla f, h \rangle$$

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

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

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

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 .

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

Direction of local steepest descent

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

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

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.

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

Gradient flow ODE

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

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

guckpetuzayus
(GF)

and discretize it on a uniform grid with α step:

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

↖ Эйлер

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 ♣

guckpetenzus
zunepca

ODE
GF

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 ♣

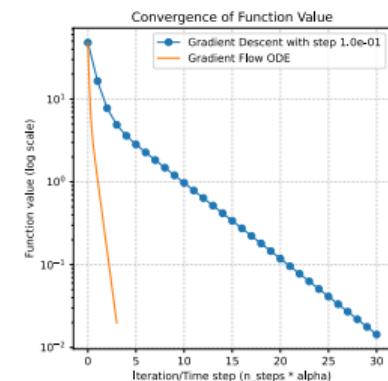
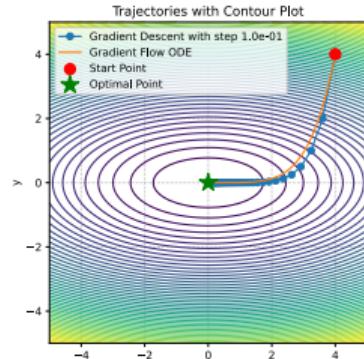


Figure 1: Gradient flow trajectory

Necessary local minimum condition

$$\begin{array}{c}
 f'(x) = 0 \\
 \hline
 -\eta f'(x) = 0 \quad +X \\
 \hline
 x - \eta f'(x) = x \\
 \hline
 \overbrace{x_k - \eta f'(x_k)}^{<} = x_{k+1} \quad \overbrace{\quad}^{K+1}
 \end{array}$$

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 $\forall x, y \in \mathbb{R}^n$:

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

∇f - L ununyel

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



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:

$$\begin{aligned}\phi_1(x) &= f(x_0) + \langle \nabla f(x_0), x - x_0 \rangle - \frac{L}{2} \|x - x_0\|^2, \\ \phi_2(x) &= f(x_0) + \langle \nabla f(x_0), x - x_0 \rangle + \frac{L}{2} \|x - x_0\|^2.\end{aligned}$$

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

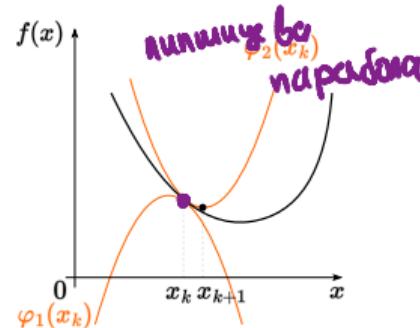


Figure 2: Illustration

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

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 $\sqrt{L^2} \cdot 0 + \nabla f(x_0) + L \cdot (x - x_0)$
 $\phi_1(x) \leq f(x) \leq \phi_2(x) \quad \forall x \in \mathbb{R}^n.$

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

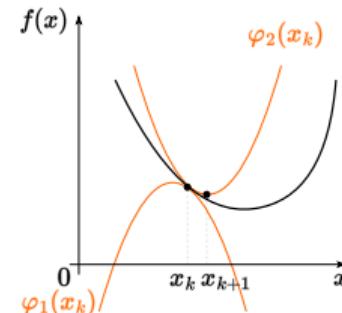


Figure 2: Illustration

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

$$\nabla f(x_0) + L(x^{k+1} - x_0) = 0$$

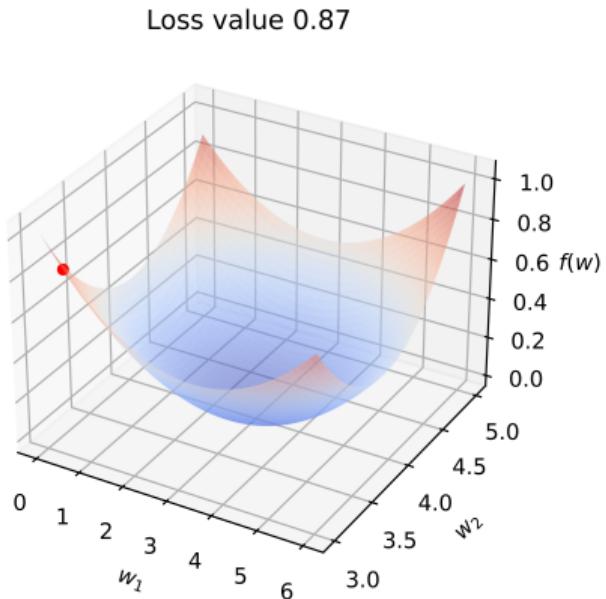
$$x^* = x_0 - \frac{1}{L} \nabla f(x_0)$$

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

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

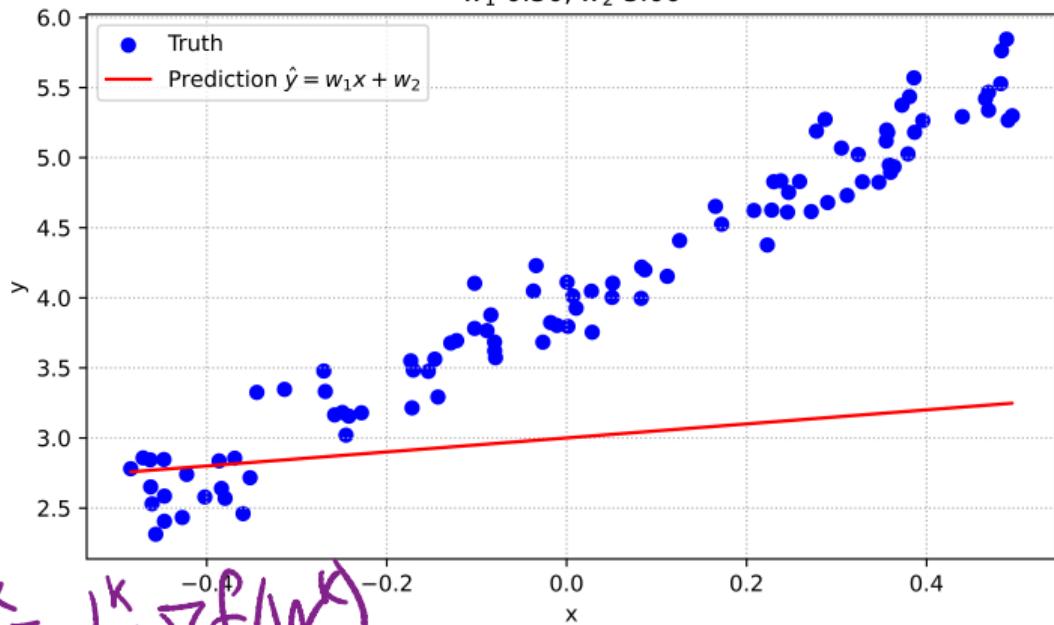
Convergence of Gradient Descent algorithm

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:

Метог наискорейшего спуска

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

Optimality conditions:

Пришлі \exists якщо $f(x)$ має локальні мінімуми.

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

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

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

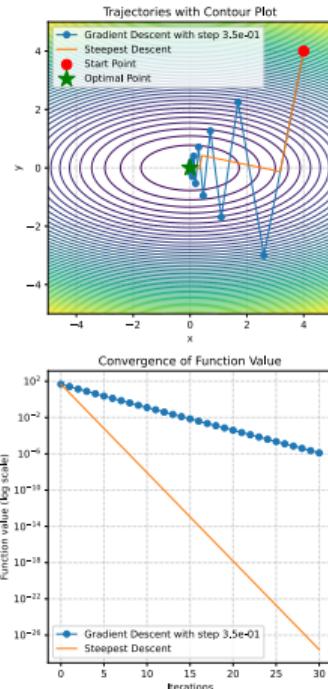


Figure 3: Steepest Descent

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:

$$\frac{\partial f}{\partial d} = \frac{\partial f}{\partial x_{k+1}}^T \frac{\partial x_{k+1}}{\partial d} = 0$$

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

$$[A^T(Ax_k - b)]^T [A(x_k - b)] = 0$$

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

Optimality conditions:

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

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

$$d_k = \arg \min_d f(x_{k+1}) = \arg \min_d f(x_k - d \cdot A^T(Ax_k - b))$$

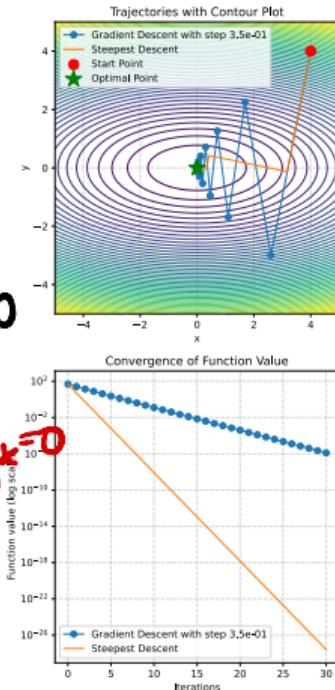


Figure 3: Steepest Descent

Open In Colab ♣

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

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

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

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

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

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

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:

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

$$\nabla f = A x_k$$

Optimality conditions:

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

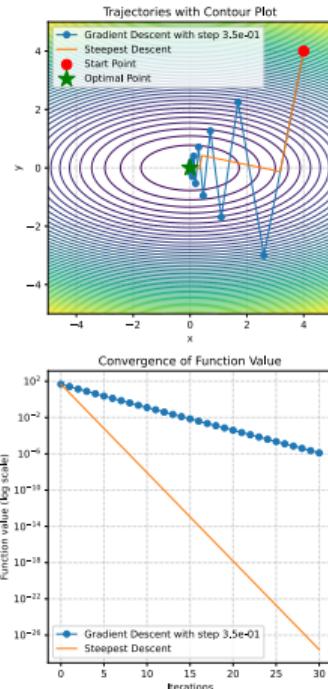


Figure 3: Steepest Descent

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

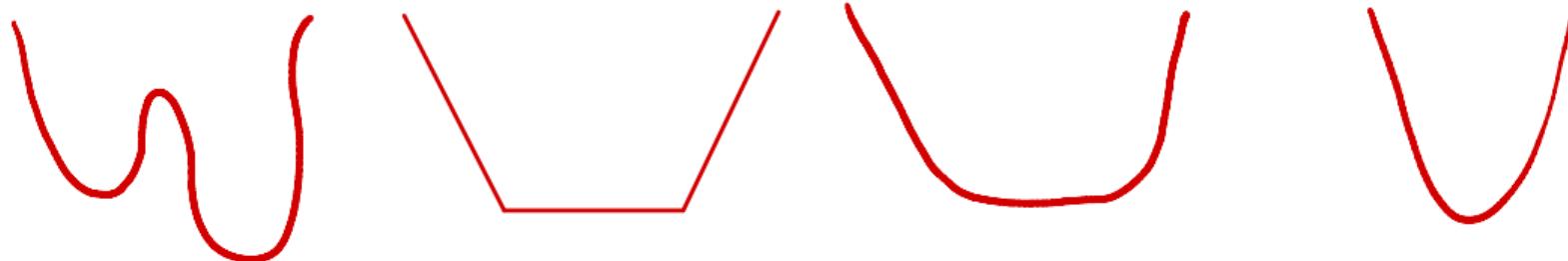
smooth & strongly convex (or PL)

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

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

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

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



Gradient Descent convergence. Smooth convex case

Факты

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

Tогда f - L -липшицева iff

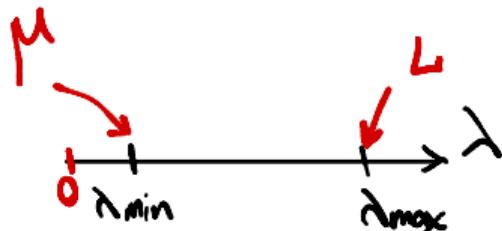
$$\forall x \in \mathbb{R}^n : \|Df(x)\| \leq L$$

специальный случай
если

Df - L -липшицев

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

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



2) $\text{Если } f: \mathbb{R}^n \rightarrow \mathbb{R} \text{ - } L\text{-липшицева}$
(липшицев $\nabla f \in \text{кот. } L$)

, 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 & шагом learning rule

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

т.к. f - L -липшицева

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

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

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

Gradient Descent convergence. Smooth convex case

2) $f(x) - \text{being convex}$

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

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

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

gekennzeichnet

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

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

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

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

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

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

Gradient Descent convergence. Smooth convex case

$$\nabla f(x^*)^\top \left[x^* - x^* - \frac{1}{2L} \nabla f(x^*) \right] =$$

$$\Leftrightarrow f(x^*) + \frac{L}{2} \left(\|x^* - x^*\|^2 - \|x^{*+} - x^*\|^2 \right)$$

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

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

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

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

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

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

$$f(\theta x + (1-\theta)y) \leq \theta f(x) + (1-\theta)f(y)$$

x_1, \dots, x_T

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

$\theta_i \geq 0$

$\sum \theta_i = 1$

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

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

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

$$T \cdot f(x^*) \leq f(x^0) + f(x^1) + \dots + f(x^T)$$

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

сублинейно
O(T)

Доказали, что GD сходится сублинейно
для малых выпуклых функций

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

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

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

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

$$\nabla^2 f$$



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

Gradient Descent convergence. Smooth μ -strongly convex case

Gradient Descent convergence. Polyak-Lojasiewicz case

④ Сильно выпуклый локальный минимум:

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

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

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

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 $f(x)$ — сильно выпуклая, то она убывает.

PL:

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

Сильная выпуклость:

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

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

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

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

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

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

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

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

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

$$\Rightarrow \frac{1}{2} \left[\frac{1}{\mu} \|\nabla f(x)\|^2 - \|\sqrt{\mu}(x - x^*) - \frac{1}{\sqrt{\mu}} \nabla f(x)\|^2 \right] \leq$$

$$\leq \frac{1}{2\mu} \|\nabla f(x)\|^2 \quad \text{R.T.g.}$$

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

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

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

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

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

$$\alpha \leq \frac{1}{L}$$

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

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

$$1-\alpha\mu < 1$$

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

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

$$\alpha\mu - 1 < 0$$

$$1 - \alpha\mu > 0$$

$$-1 < 1 - \alpha\mu < 1$$

$$-1 < 1 - \alpha\mu < 1$$

$$\alpha \leq \frac{1}{L}$$

$$0 < 1 - \alpha\mu < 1$$

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

$$1 - \alpha\mu > 0 \Rightarrow \alpha\mu < 1$$

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

$$1 - \alpha\mu < 1$$

$$-\alpha\mu < 0$$

$$\alpha > 0 \quad \mu > 0$$

$$\alpha \leq \frac{1}{L}$$

⑤ Сходимость GD qua PL:

L - шаг кости

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

$$= f(x^k) + \nabla f(x^k)^T (-\alpha \nabla f(x^k)) + \frac{L}{2} \alpha^2 \|\nabla f(x^k)\|^2 = \quad x^{k+1} - x^k = -\alpha \nabla f(x^k)$$

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

$$= f(x^k) + \|\nabla f(x^k)\|^2 \left(\frac{\alpha^2 L}{2} - \alpha \right) =$$

$$= f(x^k) + \frac{\alpha}{2} \|\nabla f(x^k)\|^2 (\alpha L - 2) = \quad \alpha \leq \frac{2}{L}$$

$$= f(x^k) - \frac{\alpha}{2} \|\nabla f(x^k)\|^2 (2 - \alpha L) \leq \quad \alpha L \leq 2$$

$$\leq f(x^*) - \frac{\alpha}{2} \cdot 2\mu(f(x^*) - f^*)(2 - \alpha L)$$

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

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

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

$$-1 < 1 - \alpha\mu(2 - \alpha L) < 1$$

$$-2 < -\alpha\mu(2 - \alpha L) < 0$$

$$0 < \alpha\mu(2 - \alpha L) < 2$$

$$2 - \alpha L \geq 0$$

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

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

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

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

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