

# ⊙ Постановка задачи

$$f(\theta) \rightarrow \min_{\theta \in \mathbb{R}^p}$$

$$f(\theta) = \frac{1}{N} \sum_{i=1}^N f_i(\theta) \rightarrow \min_{\theta \in \mathbb{R}^p}$$

Пример: лин. регрессия

есть набор

$$(x_i, y_i) \quad i = \overline{1, N}$$

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

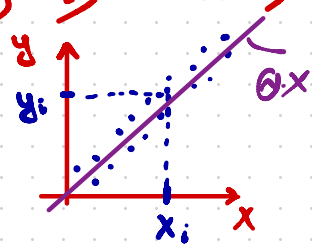
$$y \in \mathbb{R}^N$$

модель

$$y_i \stackrel{\text{хотела}}{=} f(x_i) = f(\theta, x_i) = f_i(\theta)$$

$$f_i(\theta) = \theta^T x_i$$

$$y_i \approx f_i(\theta)$$



$$|y_i - f_i(\theta)| \rightarrow \min_{\theta \in \mathbb{R}^p}$$

$$\sum_{i=1}^N (y_i - f_i(\theta))^2 \rightarrow \min_{\theta \in \mathbb{R}^p}$$

$$f_i(\theta) = \theta^T x_i$$
$$f = X\theta$$

$$\|X\theta - y\|^2 \rightarrow \min_{\theta \in \mathbb{R}^p}$$

$$\langle X\theta - y, X\theta - y \rangle =$$

$$\nabla f = 2X^T(X\theta - y) = \sum_{i=1}^N 2$$

$$f(\theta) = \frac{1}{N} \sum_{i=1}^N (y_i - \theta^T x_i)^2$$

$$\nabla_{\theta} f = \frac{1}{N} \sum_{i=1}^N 2(y_i - \theta^T x_i) \cdot (-x_i) \in \mathbb{R}^p$$

$N$  - размер датасетов

$$10^7 - 10^8 - 10^9$$

$p$  - размер модели

$$p \approx 10^{11} - 10^{12}$$

$$p \cdot 16 \text{ бит} \cdot N = 10^{11} \cdot 10^7 \cdot 10 \approx 10^{19} \text{ бит}$$

↙  
 $10^{11} \cdot 10 \text{ бит}$

$$10^{10} \text{ Гбит}$$

$$10^9 \text{ ГБ}$$

$$10^{12} \text{ бит}$$

$$10^3 \text{ Гбит}$$

$$10^3 \text{ ГБ}$$

Идея SGD: вместо точной суммы  
считать стох. аппроксимацию

$$g^k \approx \nabla f(\theta^k)$$

↑  
стох.  
градиент

SGD (1951)  $g^k = \nabla f_{i_k}(\theta^k)$

выбираем  
случайно одно слагаемое  
из суммы

Замечание: стох. градиент является несмещённой  
оценкой истинного градиента

$$g^k = \sum_{i=1}^N \nabla f_i(\theta^k) \cdot \xi_i$$

$\xi_i$  - сл. вел.

$$\xi_i = \begin{cases} 1, & \text{с вер. } \frac{1}{N} \\ 0, & \text{с вер. } \frac{N-1}{N} \end{cases}$$

$$\mathbb{E} g^k = \sum_{i=1}^N \nabla f_i(\theta^k) \cdot \frac{1}{N} = \nabla f(\theta^k)$$

$\text{Var } g^k \sim \frac{1}{b}$

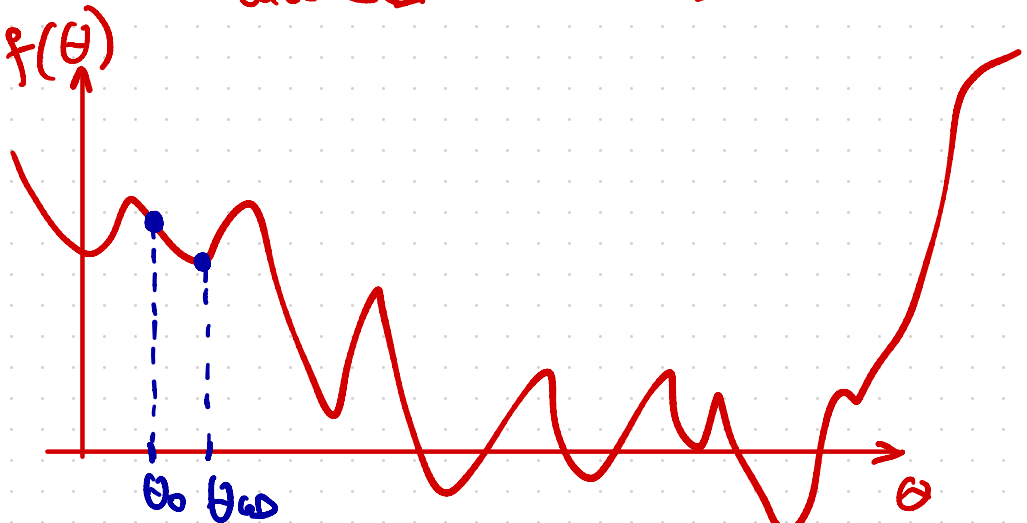
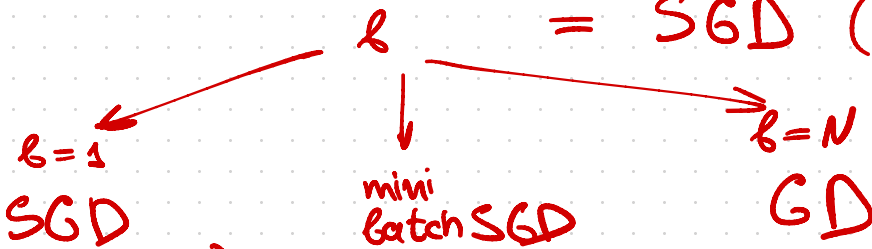
Mini batch SGD:

случаен  $g^k = \frac{1}{b} \sum_{i=1}^b \nabla f_{j_i}(\theta^k)$

случајно

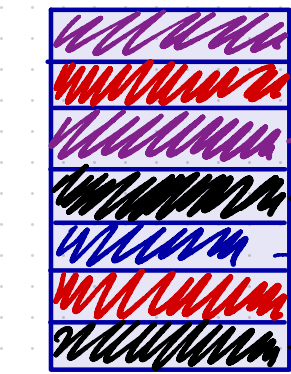
выборку ем б случајно уз N

- Ели пољубави Hardware (GPU), то препорука mini batch SGD = SGD (b=1)



эпоха  $b$  SGD

$b = 2$



$N$  строк

1 шаг SGD  $\theta^1 = \dots$

$$\theta^2 = \theta^1 - \alpha g^1$$

2 шага SGD  $\theta^2 = \dots$

$$\theta^3 = \theta^2 - \alpha g^2$$

3 шага.

эпоха

Как зависит время 1 эпохи от

$N$

$b$  ?

$b$

в одной эпохе

$\frac{N}{b}$  итераций

если позволяет GPU RAM,  
то сокращается

# Summary

Suppose, our target function is the sum of functions.

$$\min_{\theta \in \mathbb{R}^p} g(\theta) := \frac{1}{n} \sum_{i=1}^n f_i(\theta)$$

This problem usually arises in Deep Learning, where the gradient of the loss function is calculating over the huge number of data points, which could be very expensive in terms of the iteration cost (calculation of gradient is linear in  $n$ ).

Thus, we can switch from the full gradient calculation to its unbiased estimator:

$$\theta_{k+1} = \theta_k - \alpha_k \nabla f_{i_k}(\theta),$$

where we randomly choose  $i_k$  index of point at each iteration uniformly:

$$\mathbb{E}[\nabla f_{i_k}(\theta)] = \sum_{i=1}^n p(i_k = i) \nabla f_i(\theta) = \frac{1}{n} \sum_{i=1}^n \nabla f_i(\theta) = \nabla g(\theta)$$

Iterations could be  $n$  times cheaper! But convergence requires  $\alpha_k \rightarrow 0$ .

# Convergence

## General setup

We consider classic finite-sample average minimization:

$$\min_{x \in \mathbb{R}^p} f(x) = \min_{x \in \mathbb{R}^p} \frac{1}{n} \sum_{i=1}^n f_i(x)$$

Let us consider stochastic gradient descent assuming  $\nabla f$  is Lipschitz:

$$x_{k+1} = x_k - \alpha_k \nabla f_{i_k}(x_k) \tag{SGD}$$

Lipschitz continuity implies:

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

using (SGD):

$$f(x_{k+1}) \leq f(x_k) - \alpha_k \langle \nabla f(x_k), \nabla f_{i_k}(x_k) \rangle + \alpha_k^2 \frac{L}{2} \|\nabla f_{i_k}(x_k)\|^2$$

Now let's take expectation with respect to  $i_k$ :

$$\mathbb{E}[f(x_{k+1})] \leq \mathbb{E}[f(x_k) - \alpha_k \langle \nabla f(x_k), \nabla f_{i_k}(x_k) \rangle + \alpha_k^2 \frac{L}{2} \|\nabla f_{i_k}(x_k)\|^2]$$

Using linearity of expectation:

$$\mathbb{E}[f(x_{k+1})] \leq f(x_k) - \alpha_k \langle \nabla f(x_k), \mathbb{E}[\nabla f_{i_k}(x_k)] \rangle + \alpha_k^2 \frac{L}{2} \mathbb{E}[\|\nabla f_{i_k}(x_k)\|^2]$$

Since uniform sampling implies unbiased estimate of gradient:  $\mathbb{E}[\nabla f_{i_k}(x_k)] = \nabla f(x_k)$   
:

$$\mathbb{E}[f(x_{k+1})] \leq f(x_k) - \alpha_k \|\nabla f(x_k)\|^2 + \alpha_k^2 \frac{L}{2} \mathbb{E}[\|\nabla f_{i_k}(x_k)\|^2]$$

## Polyak-Lojasiewicz conditions

$$\frac{1}{2} \|\nabla f(x)\|_2^2 \geq \mu(f(x) - f^*), \forall x \in \mathbb{R}^p \quad (\text{PL})$$

This inequality simply requires that the gradient grows faster than a quadratic function as we move away from the optimal function value. Note, that strong convexity implies PL, but not vice versa. Using PL we can write:

$$\mathbb{E}[f(x_{k+1})] - f^* \leq (1 - 2\alpha_k \mu)[f(x_k) - f^*] + \alpha_k^2 \frac{L}{2} \mathbb{E}[\|\nabla f_{i_k}(x_k)\|^2]$$

This bound already indicates, that we have something like linear convergence if far from solution and gradients are similar, but no progress if close to solution or have high variance in gradients at the same time.

## Stochastic subgradient descent

$$x_{k+1} = x_k - \alpha_k g_{i_k} \quad (\text{SSD})$$

for some  $g_{i_k} \in \partial f_{i_k}(x_k)$ .

For convex  $f$  we have

$$\mathbb{E}[\|x_{k+1} - x^*\|^2] = \|x_k - x^*\|^2 - 2\alpha_k \langle g_k, x_k - x^* \rangle + \alpha_k^2 \mathbb{E}[\|g_{i_k}\|^2]$$

Here we can see, that step-size  $\alpha_k$  controls how fast we move towards solution. And squared step-size  $\alpha_k^2$  controls how much variance moves us away. Usually, we bound  $\mathbb{E}[\|g_{i_k}\|^2]$  by some constant  $B^2$ .

$$\mathbb{E}[\|x_{k+1} - x^*\|^2] = \|x_k - x^*\|^2 - 2\alpha_k \langle g_k, x_k - x^* \rangle + \alpha_k^2 B^2$$

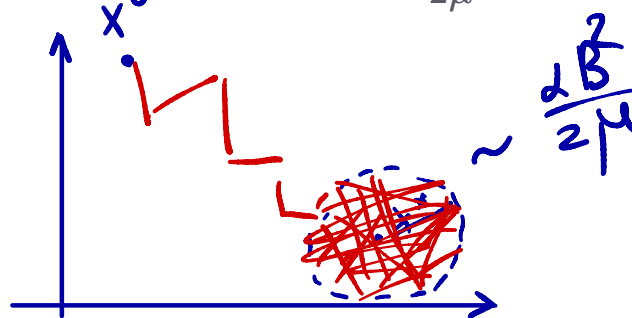
If we also have strong convexity:

$$\mathbb{E}[\|x_k - x^*\|^2] \leq (1 - 2\alpha_k \mu) \|x_{k-1} - x^*\|^2 + \alpha_k^2 B^2$$

And finally, with  $\alpha_k = \alpha < \frac{2}{\mu}$ :

$$\mathbb{E}[\|x_k - x^*\|^2] \leq (1 - 2\alpha_k \mu)^k R^2 + \frac{\alpha B^2}{2\mu},$$

where  $R = \|x_0 - x^*\|$



## Bounds

Conditions	$\ \mathbb{E}[f(x_k)] - f(x^*)\  \leq$	Type of convergence
Convex, Lipschitz-continuous gradient (L)	$\mathcal{O}\left(\frac{1}{\sqrt{k}}\right)$	Sublinear
$\mu$ -Strongly convex, Lipschitz-continuous gradient (L)	$\mathcal{O}\left(\frac{1}{k}\right)$	Sublinear
Convex, non-smooth	$\mathcal{O}\left(\frac{1}{\sqrt{k}}\right)$	Sublinear
$\mu$ -Strongly convex, non-smooth	$\mathcal{O}\left(\frac{1}{k}\right)$	Sublinear



# Code

 Open in Colab

## References

- [Lecture](#) by Mark Schmidt @ University of British Columbia
- [Convergence theorems](#) on major cases of GD, SGD (projected version included)

[https://colab.research.google.com/github/MerkulovDaniil/sber219/blob/main/notebooks/9\\_01.ipynb#scrollTo=UUJA6iqI7MLu](https://colab.research.google.com/github/MerkulovDaniil/sber219/blob/main/notebooks/9_01.ipynb#scrollTo=UUJA6iqI7MLu)