$$\min_{x \in \mathbb{R}^d} f(x) \quad \xcancel{\longrightarrow} \quad \nabla f(x)$$

$$\Vert$$

$$\mathbb{E}\left[ f(x, \xi) \right] \longleftarrow \mathbb{E}_D\left[ l\left( g(x, \xi), b \right) \right] \longrightarrow \min$$

$$\xi \sim D$$

<span style="color:red">↑ веса  ↑ ответ</span>

<span style="color:red">$\xi \in D$ ← не знаем</span>

<span style="color:red">заменяем стохастич.</span>

<span style="color:red">$\boxed{\xi_1 \dots \xi_n \in D}$</span>

$$\frac{1}{n} \sum_{i=1}^{n} l\left( g(x, \xi_i); b_i \right)$$

---

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

- GD: $\quad x^{k+1} = x^k - \gamma \nabla f(x^k)$  <span style="color:red">← дело $n$-белное</span>

- SGD: $\quad x^{k+1} = x^k - \gamma \nabla f_{i_k}(x^k)$  <span style="color:red">← каждую итерацию считаем $\nabla f_{i_k}$ случайного индекса $i_k$</span>

$$\mathbb{E}\left[ \nabla f_{i_k}(x) \right] = \frac{1}{n} \sum_{i=1}^{n} \nabla f_i(x) = \nabla f(x)$$

---

$$\| x^{k+1} - x^* \|^2 = \| x^k - x^* \|^2 - 2\gamma_k \langle \nabla f_{i_k}(x^k); x^k - x^* \rangle + \gamma_k^2 \| \nabla f_{i_k}(x^k) \|^2$$

$$\mathbb{E} \text{ по } i_k \qquad \mathbb{E}\left[ \; \big| x^k \right]$$

<span style="color:red">↑ фикс всего пред.</span>

$$\mathbb{E}\left[ \| x^{k+1} - x^* \|^2 \,\big|\, x^k \right] = \| x^k - x^* \|^2 - 2\gamma_k \mathbb{E}\left[ \langle \nabla f_{i_k}(x^k); \boxed{x^k - x^*} \big| x^k \rangle \right]$$

$$+ \gamma_k^2 \mathbb{E}\left[ \| \nabla f_{i_k}(x^k) \|^2 \,\big|\, x^k \right]$$

$$= \|X^k - X^*\|^2 - 2\gamma_k \langle \mathbb{E}[\nabla f_{i_k}(X^k)|X^k]; X^k - X^* \rangle$$
$$+ \gamma_k^2 \mathbb{E}[\|\nabla f_{i_k}(X^k)\|^2|X^k]$$

$$= \|X^k - X^*\|^2 - 2\gamma_k \langle \nabla f(X^k); X^k - X^* \rangle \quad \text{\color{red}{смотр вот}}$$
$$+ \gamma_k^2 \mathbb{E}[\|\nabla f_{i_k}(X^k)\|^2|X^k]$$

$$\leq \|X^k - X^*\|^2 - 2\gamma_k \left( f(X^k) - f(X^*) + \frac{\mu}{2}\|X^k - X^*\|^2 \right)$$
$$+ \gamma^2 \mathbb{E}[\|\nabla f_{i_k}(X^k)\|^2|X^k]$$

$$= \underbrace{(1 - \gamma_k\mu)\|X^k - X^*\|^2 - 2\gamma_k\left( f(X^k) - f(X^*) \right)}_{\text{\color{yellow}{highlight}}} \quad \text{\color{red}{было GD}}$$
$$\boxed{+ \gamma^2 \mathbb{E}[\|\nabla f_{i_k}(X^k)\|^2|X^k]}$$

$$\color{red}{=}$$

$$\mathbb{E}[\|\nabla f_{i_k}(X^k)\|^2|X^k] =$$
$$\mathbb{E}[\|\nabla f_{i_k}(X^k) - \nabla f(X^k) + \nabla f(X^k)\|^2|X^k]$$
$$= \mathbb{E}\Big[ \underbrace{\|\nabla f_{i_k}(X^k) - \nabla f(X^k)\|^2}_{\color{red}{c}} + \underbrace{\|\nabla f(X^k)\|^2}_{\color{red}{h}} +$$
$$+ 2 \underbrace{\langle \nabla f_{i_k}(X^k) - \nabla f(X^k);}_{\color{red}{c}} \underbrace{\nabla f(X^k)}_{\color{red}{h}} \rangle |X^k \Big]$$

$$= \mathbb{E}[\|\nabla f_{i_k}(X^k) - \nabla f(X^k)\|^2|X^k] + \|\nabla f(X^k)\|^2$$

$$+2 \langle \mathbb{E}[\nabla f_{i_k}(x^k)|x^k] - \nabla f(x^k) \; ; \; \nabla f(x^k) \rangle$$

$$\underbrace{\phantom{xxx}}_{0}$$

$$= \mathbb{E}\left[\|\nabla f_{i_k}(x^k) - \nabla f(x^k)\|^2 \mid x^k\right] + \|\nabla f(x^k)\|^2 \quad \text{\color{red}{было GD}}$$

<u>Предположение:</u>

$$\mathbb{E}\left[\|\nabla f_i(x) - \nabla f(x)\|^2\right] \leq B\|\nabla f(x)\|^2 + \sigma^2$$

$$\mathbb{E}\left[\|\nabla f_{i_k}(x^k)\|^2 \mid x^k\right] \leq \|\nabla f(x^k)\|^2 + B\|\nabla f(x^k)\|^2 + \sigma^2$$

$$\overset{\color{red}{=}}{\phantom{x}} (1-\gamma_k\mu)\|x^k - x^*\|^2 - 2\gamma_k\left(f(x^k) - f(x^*)\right)$$
$$+ (1+B)\gamma_k^2 \underbrace{\|\nabla f(x^k)\|^2}_{\color{red}{-\nabla f(x^*)}} + \gamma_k^2\sigma^2 \quad \color{red}{\text{типа грд.}}$$
$$\color{red}{\underset{0}{\|}}$$

$$\leq (1-\gamma_k\mu)\|x^k - x^*\|^2 - \underline{2\gamma_k\left(f(x^k) - f(x^*)\right)}$$
$$+ \underline{(1+B)\gamma_k^2\left(2L\left(f(x^k) - f(x^*)\right)\right)} + \color{red}{\gamma_k^2\sigma^2}$$

$$\color{red}{-2 + (1+B)4L\gamma_k = 0}$$

$$\color{red}{\gamma_k \leq \frac{1}{2L(1+B)}}$$

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

$$\mathbb{E}\left[\|X^{k+1}-x^*\|^2 \mid X^k\right] \leq (1-\gamma_k\mu)\|X^k-x^*\|^2 + \gamma_k^2\sigma^2$$

$$\mathbb{E}\left[\mathbb{E}\left[\|X^{k+1}-X^*\|^2 \mid X^k\right]\right] \leq (1-\gamma_k\mu)\mathbb{E}\left[\|X^k-x^*\|^2\right] + \gamma_k^2\sigma^2$$

$$\mathbb{E}\left[\underbrace{\|X^{k+1}-X^*\|^2}_{R_{k+1}^2}\right] \leq (1-\gamma_k\mu)\underbrace{\mathbb{E}\left[\|X^k-x^*\|^2\right]}_{R_k^2} + \gamma_k^2\sigma^2$$

1) $\gamma_k \equiv \gamma$ замечаем $= (1-\gamma\mu)R_{k-1}^2 + \gamma^2\sigma^2$

$$R_{k+1}^2 \leq (1-\gamma\mu)\boxed{R_k^2} + \gamma^2\sigma^2$$

$$\leq (1-\gamma\mu)^2 R_{k-1}^2 + \gamma^2\sigma^2 + (1-\gamma\mu)\gamma^2\sigma^2$$

$$\leq (1-\gamma\mu)^3 R_{k-2}^2 + \gamma^2\sigma^2 + (1-\gamma\mu)\gamma^2\sigma^2 + (1-\gamma\mu)^2\gamma^2\sigma^2$$

$$\leq (1-\gamma\mu)^{k+1}R_0^2 + \gamma^2\sigma^2\sum_{i=0}^{k}(1-\gamma\mu)^i$$

было GD

$$\leq \gamma^2\sigma^2\sum_{i=0}^{\infty}(1-\gamma\mu)^i$$

$$\leq \frac{\gamma^2\sigma^2}{\gamma\mu} = \frac{\gamma\sigma^2}{\mu}$$

$$R_{k+1}^2 \leq \underbrace{(1-\mu\gamma)^{k+1}R_0^2}_{\text{как у GD}} + \underbrace{\frac{\gamma\sigma^2}{\mu}}_{\text{окр. решения ~ const } \gamma\text{-пост.}}$$

- $\gamma = \dfrac{1}{K+1}$ $\qquad R^2_{K+1} \leq \dots + \dfrac{\sigma^2}{\mu(K+1)}$ <span style="color:red">субълинейная</span>

  <span style="color:red">↑</span>
  <u>общее число итераций</u> <span style="color:red">(теория)</span>
  <span style="color:red">плохо</span>

- $\gamma_k = \dfrac{1}{k+1}$ ( вспоминаем $FW$ ) $\longrightarrow$ <span style="color:red">$F_{k+1} \leq (1-\eta_k)F_k + \eta_k^2 C$</span>

  $$R^2_{k+1} \leq (1 - \mu\gamma_k) R^2_k + \gamma_k^2 \sigma^2$$

Можно доказать

$$\mathbb{E}\left[\|x^{k+1} - x^*\|^2\right] = O\left(\underbrace{\left(1 - \frac{\mu}{2L}\right)^K \|x^0 - x^*\|^2}_{\color{red}{GD}} + \underbrace{\frac{\sigma^2}{\mu^2 K}}_{\substack{\color{red}{стох}\\ \color{red}{сублинейная}}}\right)$$

нижняя $$\Omega\left(\underbrace{\left(1 - \sqrt{\frac{\mu}{L}}\right)^K \|x^0 - x^*\|^2}_{\color{red}{ускоряется}} + \underbrace{\frac{\sigma^2}{\mu^2 K}}_{\color{red}{не\ ускоряется}}\right)$$

1) оптим. метод — Nesterov <span>(стох вариация)</span>

2) стох больше ($\sigma^2 \Rightarrow$ больше) — оптим. метод $SGD$

- брать батч $b$
  $$\nabla f_{i_k}(x^k) \longrightarrow \frac{1}{b} \sum_{i \in S_k} \nabla f_i(x^k)$$
  <span style="color:red">↑<br>вычисл. ресурс<br>$b \ll n$</span>

  $S_k$ набор индексов
  $i_k$
  $i_k$ — независим
  $|S_k| = b$

$$\mathbb{E}\left[\frac{1}{b} \sum_{i \in S_k} \nabla f_i(x^k)\right] = \nabla f(x^k)$$

$$\mathbb{E}\left\| \frac{1}{b} \sum_{i \in S} \left( \nabla f_i(x) - \nabla f(x) \right) \right\|^2 =$$

$$= \mathbb{E}\left[ \frac{1}{b^2} \left( \underbrace{\sum_{i \in S} \| \nabla f_i(x) - \nabla f(x) \|^2}_{\leq B\|\nabla f(x)\|^2 + \sigma^2} + \sum_{\substack{i \neq j \\ i,j \in S}} \langle \underbrace{\nabla f_i(x) - \nabla f(x)}_{\mathbb{E}[\,\cdot\,\mid \text{все } i]} ; \underbrace{\nabla f_j(x) - \nabla f(x)}_{\mathbb{E}[\,\cdot\,\mid \text{все } j]} \rangle \right) \right]$$

(with $=0$ markers over each inner-product factor)

$$\leq \frac{1}{b^2} \cdot b \left( B\|\nabla f(x)\|^2 + \sigma^2 \right)$$

$$\leq \frac{B}{b} \|\nabla f(x)\|^2 + \frac{\sigma^2}{b}$$

- $\delta$ртто $\delta$аsnoz $b_k$          $b_k \uparrow$          $\dfrac{\sigma^2}{\gamma \mu \cdot b_k}$