



## Смешанный градиентный спуск

$$x^{k+1} = x^k - \gamma \nabla f(x^k, \xi^k)$$

смешан. шаг

$\xi^k$  - независимо и равномерно

по 1-ой форме  
сд. формулы

- Независимость

$$E_{\xi} [\nabla f(x, \xi)]$$

- Равномерность

$$E_{\xi} [\nabla f(x, \xi)] = (\text{описание равномерн.})$$

$$= \sum_{i=1}^n P\{\xi = \xi_i\} \nabla f(x, \xi_i) = (\text{равномерность})$$

$$= \sum_{i=1}^n \frac{1}{n} \nabla f(x, \xi_i) = \frac{1}{n} \sum_{i=1}^n \nabla f(x, \xi_i) = \nabla f(x)$$

## Уровневое многомерное суммирование

$$E[\cdot | x^k] = E[\cdot | \mathcal{F}_k]$$

$\mathcal{F}_k$  -  $\sigma$ -алгебра, порожд.  $x^0, \xi^0, \xi^{k-1}$

гипотеза, что процесс  $x^k$  (бессмешанный)

tower property:

$$E[E[X|Y]] = E[X]$$

Доказ-во сходимости:

Предположения:

- $f$  —  $\mu$ -сильно выпуклая
- $f(\cdot, g)$  —  $L$ -выпуклая (макс  $L = L_{\max}$  по  $g$  по выпуклости)
- $\mathbb{E}_g [\nabla f(x, g)] = \nabla f(x)$      $\mathbb{E}_g [f(x, g)] = f(x)$

Доказ-во:

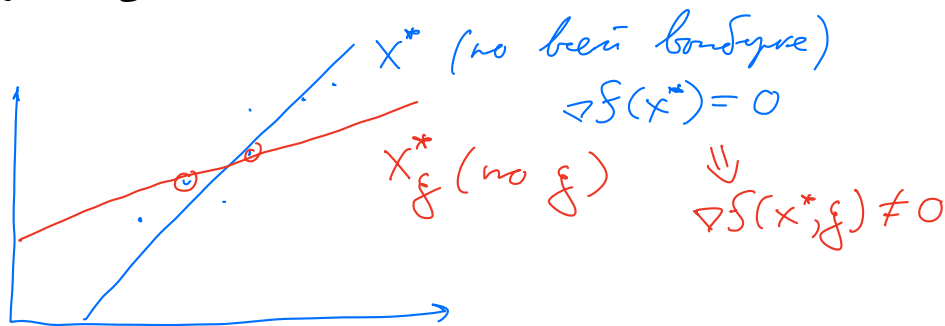
$$\begin{aligned}\|x^{k+1} - x^*\|_2^2 &= \|x^k - \gamma \nabla f(x^k, g^k) - x^*\|_2^2 \\ &= \|x^k - x^*\|_2^2 - 2\gamma \langle \nabla f(x^k, g^k); x^k - x^* \rangle \\ &\quad + \gamma^2 \|\nabla f(x^k, g^k)\|_2^2\end{aligned}$$

Получим н.о. от обеих частей:

$$\begin{aligned}\mathbb{E}[\|x^{k+1} - x^*\|_2^2] &= \mathbb{E}[\|x^k - x^*\|_2^2] - 2\gamma \mathbb{E}[\langle \nabla f(x^k, g^k); x^k - x^* \rangle] \\ &\quad + \gamma^2 \mathbb{E}[\|\nabla f(x^k, g^k)\|_2^2]\end{aligned} \quad (*)$$

$$\mathbb{E}[\|\nabla f(x^k, g^k)\|_2^2]:$$

$$\mathbb{E}[\|\nabla f(x^k, g^k)\|_2^2] = \mathbb{E}[\|\nabla f(x^k, g^k) - \nabla f(x^*, g^k) + \nabla f(x^*, g^k)\|_2^2]$$



$$\|a+b\|_2^2 \leq 2\|a\|_2^2 + 2\|b\|_2^2$$

$$\leq 2\mathbb{E}[\|\nabla f(x^k, g^k) - \nabla f(x^*, g^k)\|_2^2] + 2\mathbb{E}[\|\nabla f(x^*, g^k)\|_2^2]$$

$L$ -majorant  $f(\cdot, \xi)$

$$\leq 4L \mathbb{E} \left[ f(x^t, \xi^t) - f(x^*, \xi^t) + \langle \nabla f(x^*, \xi^t), x^* - x^t \rangle \right] + 2 \mathbb{E} \left[ \|\nabla f(x^*, \xi^t)\|_2^2 \right] \ominus$$

Tower property  $\mathbb{E}[\ ] = \mathbb{E}[\mathbb{E}[\ ] | x^k]$

$$\mathbb{E} \left[ \mathbb{E} \left[ f(x^t, \xi^t) | x^k \right] \right] = \mathbb{E} \left[ f(x^t) \right]$$

$$\mathbb{E} \left[ f(x^*, \xi^t) \right] = f(x^*)$$

$$\mathbb{E} \left[ \mathbb{E} \left[ \langle \nabla f(x^*, \xi^t), x^* - x^t \rangle | x^k \right] \right] =$$

$$= \mathbb{E} \left[ \langle \mathbb{E} \left[ \nabla f(x^*, \xi^t) | x^k \right], x^* - x^t \rangle \right]$$

$$= \mathbb{E} \left[ \langle \underbrace{\nabla f(x^*)}_{=0}, x^* - x^t \rangle \right] = 0$$

$$\mathbb{E} \left[ \mathbb{E} \left[ \|\nabla f(x^*, \xi^t)\|_2^2 | x^k \right] \right]$$

variance =

variance =

$$\overbrace{\frac{1}{n} \sum_{i=1}^n \|\nabla f(x^*, \xi_i)\|_2^2}^{\sigma_*^2}$$

$$\underbrace{\mathbb{E}_{\xi \sim D} \|\nabla f(x^*, \xi)\|_2^2}_{\sigma_*^2}$$

Соединяем

$$\Leftrightarrow 4L \mathbb{E} [f(x^k) - f(x^*)] + 2\sigma_*^2 \quad (**)$$

Теорема (\*\*) & (\*)

$$\mathbb{E} [\|x^{k+1} - x^*\|_2^2] \leq \mathbb{E} [\|x^k - x^*\|_2^2] - 2\gamma \mathbb{E} [\langle \nabla f(x^k, g^k); x^k - x^* \rangle] + \gamma^2 (4L \mathbb{E} [f(x^k) - f(x^*)] + 2\sigma_*^2) \quad (***)$$

$\mathbb{E} [\langle \nabla f(x^k, g^k); x^k - x^* \rangle]$ , tower property

$$\begin{aligned} & \mathbb{E} [\mathbb{E} [\langle \nabla f(x^k, g^k); x^k - x^* \rangle | x^k]] \\ &= \mathbb{E} [\langle \mathbb{E} [\nabla f(x^k, g^k) | x^k]; x^k - x^* \rangle] \\ &= \mathbb{E} [\langle \nabla f(x^k), x^k - x^* \rangle] \quad (****) \end{aligned}$$

Теорема (\*\*\*\*) & (\*\*\*\*)

$$\mathbb{E} [\|x^{k+1} - x^*\|_2^2] \leq \mathbb{E} [\|x^k - x^*\|_2^2] - 2\gamma \mathbb{E} [\langle \nabla f(x^k); x^k - x^* \rangle] + \gamma^2 (4L \mathbb{E} [f(x^k) - f(x^*)] + 2\sigma_*^2)$$

$\mu$ -сильно выпуклый

$$\begin{aligned} \mathbb{E} [\|x^{k+1} - x^*\|_2^2] &\leq \mathbb{E} [\|x^k - x^*\|_2^2] - 2\gamma \mathbb{E} [f(x^k) - f(x^*) + \frac{\mu}{2} \|x^k - x^*\|_2^2] \\ &\quad + \gamma^2 (4L \mathbb{E} [f(x^k) - f(x^*)] + 2\sigma_*^2) \\ &= (1 - \gamma\mu) \mathbb{E} [\|x^k - x^*\|_2^2] \end{aligned}$$

$$-2\gamma(1-2\gamma L) \underbrace{E[f(x^t) - f(x^*)]}_{\geq 0} + 2\gamma^2 \sigma_*^2$$

$$\gamma \leq \frac{1}{2L}$$

$$\leq (1-\gamma\mu) E[\|x^t - x^*\|_2^2] + 2\gamma^2 \sigma_*^2$$

$$\boxed{E[\|x^{k+1} - x^*\|_2^2] \leq \underbrace{(1-\gamma\mu) E[\|x^k - x^*\|_2^2]}_{\text{no net cancel, no "good" GD}} + \underbrace{2\gamma^2 \sigma_*^2}_{\text{the same} + E}}$$

$$R_k^2 = E[\|x^k - x^*\|_2^2]$$

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

3 arguments per step:

$$\begin{aligned} &\leq (1-\gamma\mu) \left( (1-\gamma\mu) R_{k-1}^2 + 2\gamma^2 \sigma_*^2 \right) + 2\gamma^2 \sigma_*^2 \\ &= (1-\gamma\mu)^2 R_{k-1}^2 + 2\gamma^2 \sigma_*^2 [1 + (1-\gamma\mu)] \end{aligned}$$

...

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

$$\leq (1-\gamma\mu)^{k+1} R_0^2 + 2\gamma^2 \sigma_*^2 \underbrace{\sum_{i=0}^{\infty} (1-\gamma\mu)^i}_{\frac{1}{\gamma\mu}}$$

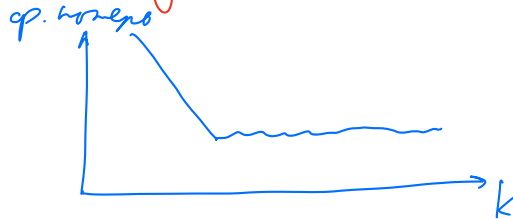
$$\leq (1-\gamma\mu)^{k+1} R_0^2 + \frac{2\gamma \sigma_*^2}{\mu}$$

Сходимость SGD с переменным шагом

$$\boxed{E[\|x^{(k+1)} - x^*\|_2^2] \leq (1-\gamma\mu)^{k+1} E[\|x^0 - x^*\|_2^2] + \frac{2\gamma\sigma_*^2}{\mu}}$$

⊕ сходимость линейная, как у GD

⊖ сходимость до оптимальности:  $\sim \gamma$ ,  $\sim \sigma_*^2$

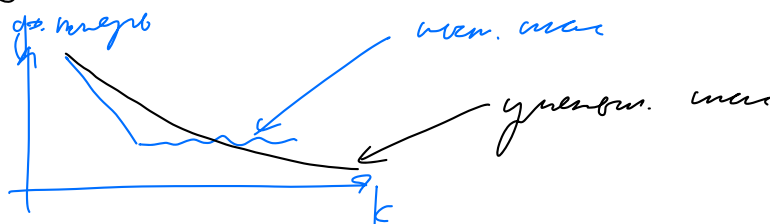


⊕ простота + при мал. шаге лучше качество ML-модели

Как бороться со сходимостью к оптимальности?

1)  $\gamma$  уменьшаем: ⊕ сходимость быстрее  
⊖ медленнее

$\gamma \rightarrow \gamma_k \sim \frac{1}{k}, \frac{1}{\sqrt{k}}$ : ⊕ сходимость к решению  
⊖ требуется линейная св. (линейная зависимость)



2)  $\sigma_*^2$  уменьшить

$$\nabla F(x, f) \rightarrow \frac{1}{b} \sum_{f \in S} \nabla f(x, f)$$

$S$  - batch (набор объектов и меток) размера  $b$   
(все точки имеют равную весовую)

$$\begin{aligned}
 & \bullet E_S \left[ \frac{1}{b} \sum_{f \in S} \nabla f(x, f) \right] = \\
 &= \frac{1}{b} \sum_{f \in S} E_f \left[ \nabla f(x, f) \right] = \text{reg. \& prob.} \\
 &= \frac{1}{b} \sum_{f \in S} \nabla f(x) = \frac{1}{b} \cdot b \nabla f(x) = \nabla f(x)
 \end{aligned}$$

$$\begin{aligned}
 & \bullet E_S \left[ \left\| \frac{1}{b} \sum_{f \in S} \nabla f(x^*, f) \right\|_2^2 \right] \\
 &= \frac{1}{b^2} E_S \left[ \sum_{f \in S} \left\| \nabla f(x^*, f) \right\|_2^2 \right] \\
 &\quad + \frac{1}{b^2} E_S \left[ \sum_{f \neq \eta} \langle \nabla f(x^*, f), \nabla f(x^*, \eta) \rangle \right] \\
 &= \frac{1}{b^2} \sum_{f \in S} E_f \left[ \left\| \nabla f(x^*, f) \right\|_2^2 \right] = \frac{1}{b^2} \cdot b \cdot \sigma_*^2 = \frac{\sigma_*^2}{b} \\
 &\quad + \frac{1}{b^2} \sum_{f \neq \eta} E_{f, \eta} \left[ \langle \nabla f(x^*, f), \nabla f(x^*, \eta) \rangle \right]
 \end{aligned}$$

$f, \eta$  - независимы

$$E_f \left[ \langle \nabla f(x^*, f), \nabla f(x^*, \eta) \rangle \right]$$

$$= \langle E_f [\nabla f(x^*, f)]; \nabla f(x^*, \eta) \rangle =$$

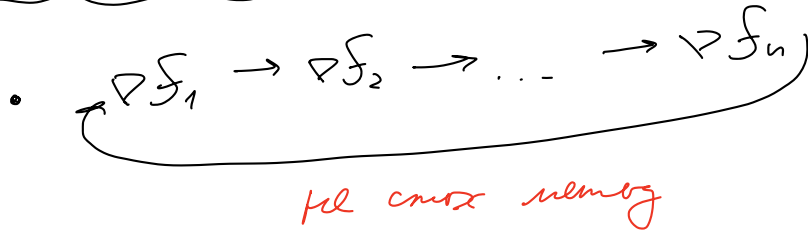
$$= \langle \nabla f(x^*), \nabla f(x^*, \eta) \rangle = 0$$

$$= \boxed{\frac{\sigma_*^2}{b}} \leftarrow \text{экстрем. суммирование}$$



- ⊕ оржемаво гурзаванево б рѣ
  - ⊖ мемораво рачево
- 

на грамме:



- shuffling {
- замуровано гурово 1 рѣ  $\rightarrow f_1 \dots f_n$  доре рачево
  - муровано рачево гурово  $\rightarrow$  ере муровано