

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

• SGD

$$x^{k+1} = x^k - \gamma \nabla f_{i_k}(x^k)$$

↑
случайно
выб. и раб.

Проблема: сходимость к оптимальности

$$x = x^* - \gamma \nabla f_i(x^*)$$

↑
выб. и раб.

$\neq 0$

$\nabla f(x^*) = 0$

$\nabla f_i(x^*) \neq 0$

$$x^{k+1} = x^k - \gamma \nabla f_{i_k}(x^k)$$

$\nabla f_{i_k}(x^*) \neq 0$

$x^k \rightarrow x^*$

замечание на $g^k: g^k \rightarrow 0 \quad x^k \rightarrow x^*$

• SAGA

y_i^k — переменные "памяти"

$\nabla f_{i_0}(x^k), \nabla f_{i_{20}}(x^{k+1})$

$x^k \approx x^{k+1}$

запускаем 2 зап. блока 1

$$y_i^k = \begin{cases} \nabla f_{i_k}(x^k), & i = i_k \\ y_i^{k-1}, & \text{иначе} \end{cases}$$

1) $x^{k+1} = x^k - \gamma \cdot \frac{1}{n} \sum_{i=1}^n y_i^k$

$\approx \nabla f(\text{выбранные моменты})$

SAGA

$$= x^k - \gamma \cdot \frac{1}{n} \sum_{i=1}^n (y_i^{k-1} + \nabla f_{i_k}(x^k) - y_{i_k}^{k-1})$$

$$2) \quad \begin{cases} g^k = \frac{1}{n} \sum_{i=1}^n y_i^{k-1} + \nabla f_{i_k}(x^k) - y_{i_k}^{k-1} \\ x^{k+1} = x^k - \gamma g^k \end{cases}$$

Доказ-во сходимости:

- f_i — L -выпукле
- f — μ -сильно выпукле

гол-во сходимости SGD $x^{k+1} = x^k - \gamma g^k$

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

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

$$\mathbb{E}[\mathbb{E}[\langle g^k; x^k - x^* \rangle | x^k]]$$

$$\mathbb{E}[g^k | x^k] = \nabla f(x^k) \quad \text{гол SGD (выпукле сильное)}$$

$$\mathbb{E}[g^k | x^k] = \mathbb{E}\left[\frac{1}{n} \sum_{i=1}^n y_i^k | x^k\right]$$

$$= \frac{1}{n} \sum_{i=1}^n \mathbb{E}[y_i^k | x^k] =$$

$$= \frac{1}{n} \sum_{i=1}^n \frac{1}{n} \nabla f_i(x^k) + \left(1 - \frac{1}{n}\right) y_i^{k-1}$$

$$= \frac{1}{n} \cdot \underbrace{\frac{1}{n} \sum_{i=1}^n \nabla f_i(x^k)}_{\nabla f(x^k)} + \left(1 - \frac{1}{n}\right) \cdot \frac{1}{n} \sum_{i=1}^n y_i^{k-1}$$

$$= \frac{1}{n} \nabla f(x^k) + \left(1 - \frac{1}{n}\right) \cdot \frac{1}{n} \sum_{i=1}^n g_i^{k-1} \neq \nabla f(x^k)$$

$$\mathbb{E}[g^k | x^k] = \frac{1}{n} \sum_{i=1}^n g_i^{k-1} + \frac{1}{n} \mathbb{E}[\nabla f_{i_k}(x^k) g_{i_k}^{k-1} | x^k]$$

$$= \frac{1}{n} \nabla f(x^k) + \left(1 - \frac{1}{n}\right) \cdot \frac{1}{n} \sum_{i=1}^n g_i^{k-1}$$

$$= \alpha \nabla f(x^k) + (1 - \alpha) \cdot \frac{1}{n} \sum_{i=1}^n g_i^{k-1}$$

$$\alpha = 1$$

$$= \nabla f(x^k)$$

гипотеза 2) неслучайно выбрана!

$$\mathbb{E}[g^k | x^k] = \nabla f(x^k)$$

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

μ -свойство берётся

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

$$\text{Тогда} \quad \mathbb{E}[\|g^k\|_2^2]$$

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

$$g^k = \mathbb{E}\left[\left\|\frac{1}{n} \sum_{i=1}^n y_i^{k-1} + \nabla f_{i_k}(x^k) - y_{i_k}^{k-1} - \nabla f(x^*)\right\|_2^2 | x^k\right]$$

$$\pm \nabla f_{i_k}(x^*) = \mathbb{E}\left[\left\|\nabla f_{i_k}(x^k) - \nabla f_{i_k}(x^*) + \nabla f_{i_k}(x^*) - y_{i_k}^{k-1} + \frac{1}{n} \sum_{i=1}^n y_i^{k-1} - \nabla f(x^*)\right\|_2^2 | x^k\right]$$

$$\text{K5W} \quad \|a+b\|_2^2 \leq 2\|a\|_2^2 + 2\|b\|_2^2$$

$$\leq 2\mathbb{E}\left[\left\|\nabla f_{i_k}(x^k) - \nabla f_{i_k}(x^*)\right\|_2^2 | x^k\right]$$

$$+ 2\mathbb{E}\left[\left\|y_{i_k}^{k-1} - \nabla f_{i_k}(x^*) + \frac{1}{n} \sum_{i=1}^n y_i^{k-1} - \nabla f(x^*)\right\|_2^2 | x^k\right]$$

$$\mathbb{E}\left[y_{i_k}^{k-1} - \nabla f_{i_k}(x^*) | x^k\right] = \frac{1}{n} \sum_{i=1}^n y_i^{k-1} - \nabla f(x^*)$$

$$\mathbb{E}\left[\|g - \mathbb{E}[g | x^k]\|_2^2 | x^k\right] = \text{ID}[g | x^k] \leq \mathbb{E}[\|g\|_2^2 | x^k]$$

$$\leq 2\mathbb{E}\left[\left\|\nabla f_{i_k}(x^k) - \nabla f_{i_k}(x^*)\right\|_2^2 | x^k\right]$$

$$+ 2\mathbb{E}\left[\left\|y_{i_k}^{k-1} - \nabla f_{i_k}(x^*)\right\|_2^2 | x^k\right]$$

$$= 2 \cdot \frac{1}{n} \sum_{i=1}^n \left\|\nabla f_i(x^k) - \nabla f_i(x^*)\right\|_2^2$$

$$+ 2 \cdot \frac{1}{n} \sum_{i=1}^n \left\|y_i^{k-1} - \nabla f_i(x^*)\right\|_2^2$$

L - константа

$$\leq 4L \cdot \frac{1}{n} \sum_{i=1}^n (f_i(x^k) - f_i(x^*) - \langle \nabla f_i(x^*); x^k - x^* \rangle) + 2 \cdot \frac{1}{n} \sum_{i=1}^n \|y_i^{k-1} - \nabla f_i(x^*)\|_2^2$$

$$= 4L (f(x^k) - f(x^*) - \langle \nabla f(x^*); x^k - x^* \rangle) + 2 \cdot \frac{1}{n} \sum_{i=1}^n \|y_i^{k-1} - \nabla f_i(x^*)\|_2^2$$

Вернемся к результату:

$$\mathbb{E}[\|x^{k+1} - x^*\|_2^2] \leq (1 - \gamma\mu) \mathbb{E}[\|x^k - x^*\|_2^2] - 2\gamma \mathbb{E}[f(x^k) - f(x^*)] + 4L\gamma^2 \mathbb{E}[f(x^k) - f(x^*)] + 2\gamma^2 \cdot \mathbb{E}\left[\frac{1}{n} \sum_{i=1}^n \|y_i^{k-1} - \nabla f_i(x^*)\|_2^2\right]$$

$$y_i^k \approx \nabla f_i(\tilde{x}^k)$$

\uparrow
ошибка

если $x^k \rightarrow x^*$, то $\tilde{x}^k \rightarrow x^*$
 $y_i^k \rightarrow \nabla f_i(x^*)$

Получим с помощью замены:

$$\mathbb{E}\left[\frac{1}{n} \sum_{i=1}^n \|y_i^k - \nabla f_i(x^*)\|_2^2 \mid x^k\right]$$

$$= \frac{1}{n} \sum_{i=1}^n \mathbb{E} [\|y_i^k - \nabla f_i(x^*)\|_2^2 | x^k]$$

$$= \frac{1}{n} \sum_{i=1}^n \left(\frac{1}{n} \|\nabla f_i(x^k) - \nabla f_i(x^*)\|_2^2 + \left(1 - \frac{1}{n}\right) \|y_i^{k-1} - \nabla f_i(x^*)\|_2^2 \right)$$

$$= \left(1 - \frac{1}{n}\right) \cdot \frac{1}{n} \sum_{i=1}^n \|y_i^{k-1} - \nabla f_i(x^*)\|_2^2 + \frac{1}{n} \cdot \frac{1}{n} \sum_{i=1}^n \|\nabla f_i(x^k) - \nabla f_i(x^*)\|_2^2$$

L-magnitude

$$\leq \left(1 - \frac{1}{n}\right) \cdot \frac{1}{n} \sum_{i=1}^n \|y_i^{k-1} - \nabla f_i(x^*)\|_2^2 + \frac{1}{n} \cdot 2L(f(x^k) - f(x^*))$$

Umformung:

$$\mathbb{E} [\|x^{k+1} - x^*\|_2^2] \leq (1 - \gamma\mu) \mathbb{E} [\|x^k - x^*\|_2^2] + 2\gamma \mathbb{E} [f(x^k) - f(x^*)] + 4L\gamma^2 \mathbb{E} [f(x^k) - f(x^*)] + 2\gamma^2 \cdot \mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n \|y_i^{k-1} - \nabla f_i(x^*)\|_2^2 \right]$$

⊕ max. error

$$\mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n \|y_i^k - \nabla f_i(x^*)\|_2^2 \right] \leq \left(1 - \frac{1}{n}\right) \cdot \mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n \|y_i^{k-1} - \nabla f_i(x^*)\|_2^2 \right] + \frac{2L}{n} \cdot \mathbb{E} [f(x^k) - f(x^*)]$$

⊕ max. error

⊗_{1/n}
↓
divergenz gegen 0

⊗_{1/n-1}

$$\mathbb{E} \left[\|x^{k+1} - x^*\|_2^2 + \underbrace{M \cdot G_k^2}_{\text{noisy term}} \right]$$

$$\leq (1-\gamma\mu) \mathbb{E} [\|x^k - x^*\|_2^2]$$

$$- (2\gamma - 4\gamma^2 L) \mathbb{E} [f(x^k) - f(x^*)]$$

$$+ 2\gamma^2 \mathbb{E} [G_{k-1}^2]$$

$$+ \left(1 - \frac{1}{n}\right) M \mathbb{E} [G_{k-1}^2] + \frac{2L}{n} M \cdot \mathbb{E} [f(x^k) - f(x^*)]$$

$$= (1-\gamma\mu) \mathbb{E} [\|x^k - x^*\|_2^2]$$

$$+ \left(1 - \frac{1}{n} + \frac{2\gamma^2}{\mu}\right) M \cdot \mathbb{E} [G_{k-1}^2]$$

$$- \left(2\gamma - 4\gamma^2 L - \frac{2LM}{n}\right) \mathbb{E} [f(x^k) - f(x^*)]$$

$$\left(1 - \frac{1}{n} + \frac{2\gamma^2}{\mu}\right) = 1 - \frac{1}{2n} \Rightarrow \boxed{M = 4\gamma^2 n}$$

$$2\gamma - 4\gamma^2 L - 8\gamma^2 L \geq 0 \Rightarrow \gamma - 6\gamma^2 L \geq 0 \Rightarrow \boxed{\gamma \leq \frac{1}{6L}}$$

$$\boxed{\mathbb{E} \left[\|x^{k+1} - x^*\|_2^2 + M \cdot G_k^2 \right] \leq (1-\gamma\mu) \mathbb{E} [\|x^k - x^*\|_2^2] + \left(1 - \frac{1}{2n}\right) \mathbb{E} [M \cdot G_{k-1}^2]}$$

$$\leq \max \left\{ (1-\gamma\mu), \left(1 - \frac{1}{2n}\right) \right\} \mathbb{E} [\|x^k - x^*\|_2^2 + M \cdot G_{k-1}^2]$$

линейная сложность $X^k \rightarrow X^*$
 $\sigma_k^2 \rightarrow 0$

$\gamma = \frac{1}{6L}$ и шаг метода неукручено

$$\mathbb{E} \left[\|X^k - X^*\|_2^2 + M \cdot \sigma_k^2 \right] \leq \max \left\{ \left(1 - \frac{\mu}{8L}\right); \left(1 - \frac{1}{2n}\right) \right\}^k \cdot \mathbb{E} \left[\|X^0 - X^*\|_2^2 + M \sigma_0^2 \right]$$
$$\leq \mathbb{E} \left[\|X^k - X^*\|_2^2 \right]$$

Оценки сложности

$$k = O \left(\left[\frac{L}{\mu} + n \right] \log \frac{1}{\varepsilon} \right) \text{ итераций}$$

для GD:

$$k = O \left(\frac{L}{\mu} \log \frac{1}{\varepsilon} \right) \text{ итераций}$$

Сложности вращательных:

SAGA: $O \left(\left[\frac{L}{\mu} + n \right] \log \frac{1}{\varepsilon} \right)$ вращательных итераций

GD: $O \left(n \frac{L}{\mu} \log \frac{1}{\varepsilon} \right)$ вращательных итераций

- ⊕ сходимость, как у GD (но быстрее) и в 0 точке
- ⊕ шаг ступенчатый $O(1)$ (у GD $O(n)$)
- ⊕ сходимость к минимуму
- ⊖ $O(nd)$ генерации данных
- ⊖ шаг SAGA не более $\frac{1}{L}$ (здесь, а не у GD)

• SVRG

$$\begin{aligned}
 x^{k+1} &= x^k - \gamma g^k \\
 g^k &= \nabla f_{ik}(x^k) - \nabla f_{ik}(\omega^k) + \nabla f(\omega^k) \\
 \omega^k & \text{ — выборка по } \mathcal{D} \\
 \omega^k &= \begin{cases} x^k & \text{если } b \in \text{выборки} \\ \omega^{k-1} & \text{иначе} \end{cases}
 \end{aligned}$$

Результат: $x^k \rightarrow x^*$, где $\omega^k \rightarrow x^*$

$$\begin{aligned}
 g^k &= \cancel{\nabla f_{ik}(x^k)} - \cancel{\nabla f_{ik}(\omega^k)} + \cancel{\nabla f(\omega^k)} \\
 &= \underbrace{\nabla f_{ik}(x^*) - \nabla f_{ik}(x^*)}_{\rightarrow 0} + \nabla f(x^*) = 0
 \end{aligned}$$

$g^k \rightarrow 0$

Сходимость:

$$k = O\left(\left[\frac{L}{\mu} + n\right] \log \frac{1}{\epsilon}\right) \text{ итераций}$$

⊕ меньше SAGA

⊕ меньше $O(d)$

⊖ неограничен размер градиента

⊖ $O(1)$ количество итераций ≈ 3 (y SAGA = 1)

• SARAH (возмущенный ugen SVRG)

$$x^{k+1} = x^k - \gamma g^k$$

$$g^k = \nabla f_{i_k}(x^k) - \nabla f_{i_k}(x^{k-1}) + g^{k-1}$$

$$g^k = \begin{cases} \nabla f(x^k) & \text{ref b t imp.} \\ g^k & \text{иначе} \end{cases} \quad \begin{matrix} \Rightarrow \nabla f(x^{k-1}) \\ \text{смененный} \\ \text{регрессор} \end{matrix}$$

"смененный" вместо SVRG норма

g^k - смененный (y SAGA, SVRG для точечных.)

$$\mathbb{E}[g^k | x^k] \neq \nabla f(x^k)$$

$$= \mathbb{E}[\nabla f_{i_k}(x^k) - \nabla f_{i_k}(x^{k-1}) | x^k] + g^{k-1}$$

$$= \nabla f(x^k) - \underbrace{\nabla f(x^{k-1}) + g^{k-1}}_{\neq 0}$$

Сложность:

$$k = O\left(\left[\frac{L}{\mu} + n\right] \log \frac{1}{\epsilon}\right) \text{ итераций}$$

⊕ *lyone*, „*novel*“, „*crossed*“ re *german*,
re *SURE*

⊖ *remont* *gruguelin*