

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

• SGD

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

↑
выбирается шаг и индекс.

Проблема: сходимость к оптимальному

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

↑
градиент не равен 0

≠ 0

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

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

$x^k \rightarrow x^*$

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

• SAGA

y_i^k — накопленный градиент

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

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

новый градиент,
по записанной

SAGA

$$2) \quad x^{k+1} = x^k - \gamma \cdot g^k$$

$$g^k = \frac{1}{n} \sum y_i^k + \left(1 - \frac{1}{n}\right) \nabla f_{i_k}(x^k) - \left(1 - \frac{1}{n}\right) \cdot \frac{1}{n} \sum y_i^{k-1}$$

одно-одно по координатам:

$$g^k = \frac{1}{n} \sum y_i^{k-1} + \nabla f_{i_k}(x^k) - y_{i_k}^{k-1}$$

Докажем корректность:

- f_i — L -гладкие
- f — μ -сильно выпуклая

для SGD basta: $x^{k+1} = x^k - \gamma g^k$

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

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

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

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

для SGD (нужно доказать)
для чего нужно?

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

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

$$\stackrel{\text{нужно } y_i^k}{=} \frac{1}{n} \sum_{i=1}^n \left(\underbrace{\frac{1}{n}}_{p_1 (i=i_k)=\frac{1}{n}} \nabla f_i(x^k) + \underbrace{\left(1 - \frac{1}{n}\right)}_{1-p_1} y_i^{k-1} \right)$$

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

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

$$+ (1 - \frac{1}{n}) \nabla f(x^k) - (1 - \frac{1}{n}) \cdot \frac{1}{n} \sum y_i^{k-1}$$

$$\begin{aligned} 2) \mathbb{E}[g^k | x^k] &= \mathbb{E}\left[\frac{1}{n} \sum y_i^{k-1} + \nabla f_{i_k}(x^k) - y_{i_k}^{k-1} | x^k\right] \\ &= \frac{1}{n} \sum y_i^{k-1} + \mathbb{E}[\nabla f_{i_k}(x^k) | x^k] - \mathbb{E}[y_{i_k}^{k-1} | x^k] \end{aligned}$$

on SGD u borne

$$\begin{aligned} &= \cancel{\frac{1}{n} \sum y_i^{k-1}} + \nabla f(x^k) - \cancel{\frac{1}{n} \sum y_i^{k-1}} \\ &= \nabla f(x^k) \end{aligned}$$

Progreßaussagen $\subset \mathbb{E}[g^k | x^k]$

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

μ -starke Konvergenz

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

Planeps $\|g^k\|_2^2$

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

K5W $\|a+b\|_2^2 \leq 2\|a\|_2^2 + 2\|b\|_2^2$

$$\begin{aligned}
&\leq 2\mathbb{E}\left[\|\nabla f_{i_k}(x^k) - \nabla f_{i_k}(x^*)\|_2^2 | x^k\right] \\
&+ 2\mathbb{E}\left[\left\|y_{i_k}^{k-1} - \nabla f_{i_k}(x^*) - \left[\frac{1}{n} \sum y_i^{k-1} - \nabla f(x^*)\right]\right\|_2^2 | x^k\right]
\end{aligned}$$

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

$$\begin{aligned}
\mathbb{E}[\|\xi - \mathbb{E}[\xi | x^k]\|_2^2 | x^k] &= \mathbb{D}[\|\xi\|_2^2 | x^k] \\
&\leq \mathbb{E}[\|\xi\|_2^2 | x^k]
\end{aligned}$$

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

L-marginals f_{i_k}

$$\begin{aligned}
&\leq 4L \mathbb{E} \left[\underbrace{f_{i_k}(x^k) - f_{i_k}(x^*)}_{\text{blue arrow}} - \langle \nabla f_{i_k}(x^*); x^k - x^* \rangle \mid x^k \right] \\
&\quad + 2 \mathbb{E} \left[\|y_{i_k}^{k-1} - \nabla f_{i_k}(x^*)\|_2^2 \mid x^k \right] \\
&= 4L \left[\frac{1}{n} \sum (f_i(x^k) - f_i(x^*)) - \langle \frac{1}{n} \sum \nabla f_i(x^*); x^k - x^* \rangle \right] \\
&\quad + 2 \mathbb{E} \left[\|y_{i_k}^{k-1} - \nabla f_{i_k}(x^*)\|_2^2 \mid x^k \right] \\
&= 4L (f(x^k) - f(x^*) - \langle \nabla f(x^*); x^k - x^* \rangle) \\
&\quad + 2 \mathbb{E} \left[\|y_{i_k}^{k-1} - \nabla f_{i_k}(x^*)\|_2^2 \mid x^k \right]
\end{aligned}$$

Requiesce u. ugenalau:

$$\begin{aligned}
\mathbb{E} \left[\|x^{k+1} - x^*\|_2^2 \right] &\leq (1 - \gamma/n) \mathbb{E} \left[\|x^k - x^*\|_2^2 \right] \\
&\quad - 2\gamma \mathbb{E} \left[f(x^k) - f(x^*) \right] \\
&\quad + \gamma^2 \cdot 4L \mathbb{E} \left[f(x^k) - f(x^*) \right] \\
&\quad + \gamma^2 \cdot 2 \mathbb{E} \left[\|y_{i_k}^{k-1} - \nabla f_{i_k}(x^*)\|_2^2 \right]
\end{aligned}$$

$$\begin{aligned}
&\mathbb{E} \left[\|y_{i_k}^{k-1} - \nabla f_{i_k}(x^*)\|_2^2 \mid x^k \right] = \\
&= \frac{1}{n} \sum \|y_i^{k-1} - \nabla f_i(x^*)\|_2^2
\end{aligned}$$

$$\begin{aligned} \mathbb{E} [\|x^{k+1} - x^*\|_2^2] &\leq (1-\gamma\mu) \mathbb{E} [\|x^k - x^*\|_2^2] \\ &\quad - 2\gamma \mathbb{E} [f(x^k) - f(x^*)] \\ &\quad + \gamma^2 \cdot 4L \mathbb{E} [f(x^k) - f(x^*)] \\ &\quad + \gamma^2 \cdot 2 \mathbb{E} \left[\frac{1}{n} \sum \|y_i^{k-1} - \nabla f_i(x^*)\|_2^2 \right] \end{aligned}$$

$$y_i^k \approx \nabla f_i(\bar{x}^k), \text{ где } x^k \rightarrow x^*, \text{ и } \bar{x}^k \rightarrow x^*$$

↑
сходится

$$y_i^k \rightarrow \nabla f_i(x^*)$$

Тогда можно получить:

$$\begin{aligned} \mathbb{E} [\|y_i^k - \nabla f_i(x^*)\|_2^2] &= \\ &= \mathbb{E} \left[\mathbb{E} [\|y_i^k - \nabla f_i(x^*)\|_2^2 | x^k] \right] = \\ &= \mathbb{E} \left[\frac{1}{n} \|\nabla f_i(x^k) - \nabla f_i(x^*)\|_2^2 + \left(1 - \frac{1}{n}\right) \|y_i^{k-1} - \nabla f_i(x^*)\|_2^2 \right] \end{aligned}$$

L -гладность f_i

$$\begin{aligned} &\leq \mathbb{E} \left[\frac{1}{n} \cdot 2L (f_i(x^k) - f_i(x^*) - \langle \nabla f_i(x^*), x^k - x^* \rangle) \right. \\ &\quad \left. + \left(1 - \frac{1}{n}\right) \|y_i^{k-1} - \nabla f_i(x^*)\|_2^2 \right] \end{aligned}$$

Суммируем по $i = 1 \dots n$

$$\begin{aligned} & \mathbb{E} \left[\underbrace{\frac{1}{n} \sum_{i=1}^n \|y_i^k - \nabla f_i(x^*)\|_2^2}_{\sigma_k^2} \right] = \\ & \leq \mathbb{E} \left[\frac{2L}{n} \cdot \frac{1}{n} \sum (f_i(x^k) - f_i(x^*)) \right. \\ & \quad \left. - \langle \frac{1}{n} \sum \nabla f_i(x^*) ; x^k - x^* \rangle \right] \\ & \quad + \left(1 - \frac{1}{n}\right) \mathbb{E} \left[\underbrace{\frac{1}{n} \sum \|y_i^{k-1} - \nabla f_i(x^*)\|_2^2}_{\sigma_{k-1}^2} \right] \end{aligned}$$

Умножив:

$$\begin{aligned} & \mathbb{E}[\sigma_k^2] \leq \left(1 - \frac{1}{n}\right) \mathbb{E}[\sigma_{k-1}^2] + \frac{2L}{n} (f(x^k) - f(x^*)) \\ & + \quad \sigma_k^2 \rightarrow 0 \\ & \mathbb{E}[\|x^{k+1} - x^*\|_2^2] \leq \left(1 - \frac{\gamma}{n}\right) \mathbb{E}[\|x^k - x^*\|_2^2] \\ & \quad - 2\gamma \mathbb{E}[f(x^k) - f(x^*)] \\ & \quad + \gamma^2 \cdot 4L \mathbb{E}[f(x^k) - f(x^*)] \\ & \quad + \gamma^2 \cdot 2\mathbb{E}[\sigma_{k-1}^2] \end{aligned}$$

$$\|x^k - x^*\|_2^2 \rightarrow 0$$

Создаем удобный промежуточный результат:

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

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

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

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

$$\mathbb{E}[\|x^{k+1} - x^*\|_2^2 + M \cdot \sigma_k^2] \leq (1 - \mu\gamma) \mathbb{E}[\|x^k - x^*\|_2^2] + (1 - \frac{1}{2n}) \mathbb{E}[M \cdot \sigma_{k-1}^2]$$

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

линейная сходимость по $x^k \rightarrow x^*$
 $\sigma_k^2 \rightarrow 0$

$\gamma = \frac{1}{8L}$ и заданные пред.

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

$$\mathbb{E}[\|x^k - x^*\|_2^2] \leq$$

Оценки сходимости:

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

где GD (наименьшее значение)

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

- ⊕ сходимость, как у GD сгужена (с норм. до $n \log \frac{1}{\epsilon}$)
 - ⊕ шаг сгужен $O(1)$ у SAGA
 - у GD сгужен $O(n)$
 - ⊕ сходимость к моменту времени
 - ⊖ $O(nd)$ нужно хранить
-

• SVRG

$$\begin{aligned} x^{k+1} &= x^k - \gamma g^k \\ g^k &= \nabla f_{i_k}(x^k) - \nabla f_{i_k}(w^k) + \nabla f(w^k) \end{aligned}$$

w^k — случайное пересечение

$$w^k = \begin{cases} x^k & \text{раз } b \text{ итераций} \\ w^{k-1} & \text{иначе} \end{cases}$$

Сходимость:

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

Доказательство: $x^k \rightarrow x^*$ ($w^k \rightarrow x^*$)

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

$$\underbrace{\nabla f_{i_k}(x^*) - \nabla f_{i_k}(x^*)}_{\rightarrow 0}$$

$$\Rightarrow g^k \rightarrow 0$$

⊕ nice SAGA

⊕ nice $O(d)$

⊖ higher variance

• SARAH (improves upon SVRG)

$$\begin{aligned}x^{k+1} &= x^k - \eta g^k \\ g^k &= \nabla f_{i_k}(x^k) - \nabla f_{i_k}(x^{k-1}) + g^{k-1} \\ g^k & \text{ — zero drift, как } \nabla f(x^k)\end{aligned}$$

"variance" better SVRG
 g^k — unbiased (y SAGA, SVRG — stochastic unbiased)

$$\begin{aligned}\mathbb{E}[g^k | x^k] &\neq \nabla f(x^k) \\ &= \mathbb{E}[\nabla f_{i_k}(x^k) - \nabla f_{i_k}(x^{k-1}) | x^k] \\ &\quad + g^{k-1} \\ &= \nabla f(x^k) - \nabla f(x^{k-1}) + g^{k-1}\end{aligned}$$

Сложность:

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

⊕ nice for variance, like SVRG

⊖ higher variance