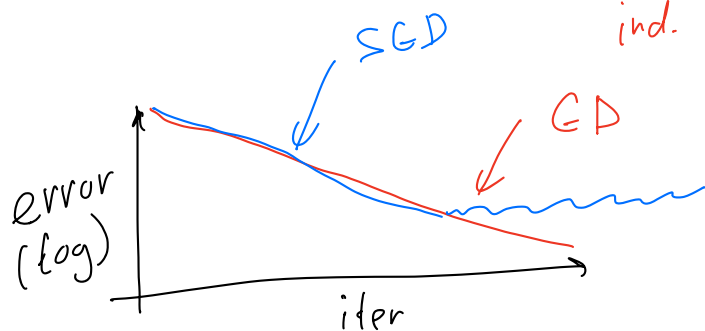


SGD:

$$\min_{x \in \mathbb{R}^d} f(x) = \frac{1}{n} \sum_{i=1}^n f_i(x) = \frac{1}{n} \sum_{i=1}^n \underline{L(\underline{g(x, a_i)}; \underline{b_i})}$$

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

↑
ind. unif.



GD: $\nabla f(x) \rightarrow \nabla f(x^*) \rightarrow 0$

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

$$\boxed{\nabla f(x) = 0}$$

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

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

what to use?
 $g^k \rightarrow 0$

SAGA:

$$y_i^0 = \nabla f_i(x^0)$$

$k \neq 0$

$$y_i^k = \begin{cases} y_i^{k-1} \\ \nabla f_i(x^k) \end{cases}$$

$i \neq i_k$
 $i = i_k$

$|S_k| = b$

SGD

1 vector

SAGA

n vectors
memory

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

with mod.
↑
 $\frac{1}{n}$

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

indeed unif [1...n]

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

$$\mathbb{E} [\] = \mathbb{E} [\underbrace{\mathbb{E}_k [\]}_{\text{cond. (take into account only } k \text{ iteration)}}]$$

$$\begin{aligned} \mathbb{E} [\langle g^k; x^k - x^* \rangle] &= \mathbb{E} [\underbrace{\mathbb{E}_k [\langle g^k; x^k - x^* \rangle]}_{\text{is not stoch}}] \\ &= \mathbb{E} [\langle \underbrace{\mathbb{E}_k [g^k]}_{\text{SGD: } g^k = \nabla f_{i_k}(x^k)}; x^k - x^* \rangle] \end{aligned}$$

$$\begin{aligned} \mathbb{E}_k [g^k] &= \sum_{i=1}^n \mathbb{P}\{i_k = i\} \cdot \nabla f_i(x^k) \\ &= \frac{1}{n} \sum_{i=1}^n \nabla f_i(x^k) = \nabla f(x^k) \end{aligned}$$

SAGA: $\mathbb{E}_k [g^k] =$

$$\begin{aligned} &= \mathbb{E}_k \left[\frac{1}{n} \sum y_i^k \right] = \\ &= \frac{1}{n} \sum \mathbb{E}_k [y_i^k] = \end{aligned}$$

$$\begin{aligned} \mathbb{E}[a+b] &= \mathbb{E}[a] \\ &\quad + \mathbb{E}[b] \end{aligned}$$

correct

$$\begin{aligned} &= \frac{1}{n} \sum \left(\nabla f_i(x^k) \cdot \frac{1}{n} + \left(1 - \frac{1}{n}\right) y_i^{k-1} \right) \\ &= \frac{1}{n^2} \sum \nabla f_i(x^k) + \frac{1}{n} \left(1 - \frac{1}{n}\right) \sum y_i^{k-1} \end{aligned}$$

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

want? $\Rightarrow \nabla f(x^k)$

after modif

$$\mathbb{E}_k[g^k] = \mathbb{E}_k\left[\frac{1}{n} \sum_i y_i^{k-1} + \nabla f_{i_k}(x^k) - y_{i_k}^{k-1}\right]$$

not rand

$$= \frac{1}{n} \sum_i y_i^{k-1} + \mathbb{E}\left[\nabla f_{i_k}(x^k) - \underline{y_{i_k}^{k-1}}\right]$$

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

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

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

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

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

rand.

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

(=)

$$g^k = \frac{1}{n} \sum_{i=1}^n y_i^k$$

ist

$$= \left[\frac{1}{n} \sum_{i \neq i_k} y_i^{k-1} \right] + \frac{1}{n} (\nabla f_{i_k}(x^k))$$

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

in new version to get
 $(F_k[g^k] = \nabla f(x^k))$

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

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

$$= \frac{1}{n} \sum_{j=1}^n \left\| \underbrace{\nabla f_j(x^k) - \nabla f_j(x^*)}_{\text{add}} + \underbrace{\nabla f_j(x^*)}_{\text{remove}} \right\|_2^2$$

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

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

$$\leq \frac{2}{n} \sum_{j=1}^n \left\| \nabla f_j(x^k) - \nabla f_j(x^*) \right\|_2^2$$

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

L -smoothness of f_j convex $\|\nabla f(x) - \nabla f(y)\|_2 \leq 2L(f(x) - f(y) - \langle \nabla f(y), x - y \rangle)$

$$\leq \frac{4L}{n} \sum_{j=1}^n \left(\underline{f_j(x^k)} - \underline{f_j(x^*)} - \langle \underline{\nabla f_j(x^*)}, \underline{x^k - x^*} \rangle \right)$$

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

$$= 4L(f(x^k) - f(x^*))$$

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

$$\mathbb{E} \left[\left\| \underset{\substack{\uparrow \\ \text{r.v.}}}{\xi} - \underset{\substack{\uparrow \\ \text{exp.}}}{\mathbb{E} \xi} \right\|_2^2 \right] \leq \mathbb{E} \left[\|\xi\|_2^2 \right]$$

$$2 \mathbb{E}_k \left[\left\| \underbrace{\frac{1}{n} \sum_i y_i^{k-1} - \frac{1}{n} \sum_i \nabla f_i(x^*)}_{\mathbb{E}_k} - \underbrace{y_{i_k}^{k-1} + \nabla f_{i_k}(x^*)}_{\mathbb{E}_k} \right\|_2^2 \right]$$

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

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

$$\begin{aligned} \mathbb{E} \|x^{k+1} - x^*\|_2^2 &= \mathbb{E} [\|x^k - x^*\|_2^2] - 2\gamma \mathbb{E} [\langle \nabla f(x^k), x^k - x^* \rangle] \\ &\quad + \gamma^2 \cdot 4L (f(x^k) - f(x^*)) \\ &\quad + \gamma^2 \cdot 2 \cdot \frac{1}{n} \sum_{i=1}^n \|\nabla f_i(x^*) - g_i^{k+1}\|_2^2 \end{aligned}$$

if $g_i^k \rightarrow \nabla f_i(x^*)$

$$\begin{aligned} \mathbb{E}_k [\|\nabla f_i(x^*) - g_i^k\|_2^2] &= \\ &= \frac{1}{n} \cdot \|\nabla f_i(x^*) - \nabla f_i(x^k)\|_2^2 + \left(1 - \frac{1}{n}\right) \|\nabla f_i(x^*) - g_i^{k-1}\|_2^2 \\ \mathbb{E}_k \left[\frac{1}{n} \sum_{i=1}^n \|\nabla f_i(x^*) - g_i^k\|_2^2 \right] &= \frac{1}{n} \sum_{i=1}^n \|\nabla f_i(x^*) - \nabla f_i(x^k)\|_2^2 \\ &\quad + \left(1 - \frac{1}{n}\right) \frac{1}{n} \sum_{i=1}^n \|\nabla f_i(x^*) - g_i^{k-1}\|_2^2 \\ &\stackrel{\text{L-smoothness and con.}}{\leq} 2L \cdot \frac{1}{n} \sum_{i=1}^n \underbrace{(f_i(x^k) - f_i(x^*))}_{-\langle \nabla f_i(x^*), x^k - x^* \rangle} \\ &= 2L (f(x^k) - f(x^*)) \end{aligned}$$

$$\mathbb{E} \|x^{k+1} - x^*\|_2^2 = \mathbb{E} \|x^k - x^*\|_2^2 - \underbrace{2\gamma \mathbb{E} \langle \nabla f(x^k), x^k - x^* \rangle}_{\text{sfr conv.}}$$

$$+ \gamma^2 \cdot 4L (f(x^k) - f(x^*))$$

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

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

$$\mathbb{E}_k \left[\|x^{k+1} - x^*\|_2^2 + \underbrace{\frac{M\sigma_k^2}{>0}} \right] \leq \underbrace{(1 - \mu\gamma) \|x^k - x^*\|_2^2}_{\text{need to remove}} + (\gamma^2 \cdot 4L - \gamma) (f(x^k) - f(x^*)) + \underbrace{2\sigma_{k-1}^2}_{\text{need to remove}} + 2L M (f(x^k) - f(x^*)) + \underbrace{(1 - \frac{1}{n}) M \sigma_{k-1}^2}_{\text{need to remove}}$$

$$M = 3n \quad 2 + 3n \cdot (1 - \frac{1}{n}) =$$

$$\begin{aligned} &= 3n - 3 + 2 = 3n - 1 = \\ &= 3n \left(1 - \frac{1}{3n}\right) \\ &= M \left(1 - \frac{1}{3n}\right) \end{aligned}$$

$$\|x^{k+1} - x^*\|_2^2 + M \Theta_k^2 \leq (1 - \mu\gamma) \|x^k - x^*\|_2^2 + \left(1 - \frac{1}{3n}\right) M \Theta_{k-1}^2$$

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

$$\leq \max\left(1 - \mu\gamma; 1 - \frac{1}{3n}\right) \left(\|x^k - x^*\|_2^2 + M \Theta_{k-1}^2\right)$$

SAGA: $\left(n + \frac{L}{\mu}\right) \log \frac{1}{\epsilon}$ iterations
 $\subseteq \text{ED}$ sublinear

$\subseteq \text{D}$ $\frac{L}{\mu} \log \frac{1}{\epsilon}$ iterations

SAGA in terms of comput n times better than ED

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

SAGA: $g^k \rightarrow 0$ $y_i^k \rightarrow \nabla f_i(x^*) \neq 0$
 $\frac{1}{n} \sum y_i^k \rightarrow 0$

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

$w^k = x^k$ (sometimes e.g. per epoch)

SVRG (not class): $\frac{1}{n} (\nabla f_{i_k}(x^k) - \nabla f_{i_k}(w^k)) + \nabla f(w^k)$
 without

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

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

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

\uparrow
 batch

\Downarrow other stoch.

$$x^{k+1} = x^k - \gamma [\nabla f]_{i_k}(x^k)$$

\uparrow
 random coord.

Example

$$f(x) = Ax - b \quad [\nabla f]_{i_k} = a_{i_k}^T x - b_{i_k}$$

Example $f(x)$ only 2-0 inf

$$\frac{f(x+\tau e) - f(x-\tau e)}{2\tau}$$

$$x^{k+1} = x^k - \gamma \cdot \underbrace{d [\nabla f]_{i_k}(x^k)}_{g^k}$$

uniformly
and indep.

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

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

$$\begin{aligned} \mathbb{E}_k[g^k] &= d \begin{pmatrix} \mathbb{E}_k[g^k]_1 \\ \vdots \end{pmatrix} = \\ &= d \begin{pmatrix} \frac{1}{d} [\nabla f(x^k)]_1 \\ \vdots \end{pmatrix} = \cancel{d} \cdot \frac{1}{\cancel{d}} \nabla f(x^k) \\ &= \nabla f(x^k) \end{aligned}$$

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

$$= \mathbb{E}_k[\|d [\nabla f(x^k)]_{i_k}\|_2^2] =$$

$$= d^2 \mathbb{E}_k[\|[\nabla f(x^k)]_{i_k}\|_2^2]$$

$$= d^2 \mathbb{E}_k[\underbrace{[\nabla f(x^k)]_{i_k}^2}_{\text{the value of } i_k \text{ coord of } \nabla f(x^k)}]$$

the value of i_k coord
of $\nabla f(x^k)$

$$= d^2 \sum_{i=1}^d \frac{1}{d} (\nabla f(x^k))_i^2 =$$

$$= d \sum_{i=1}^d (\nabla f(x^k))_i^2 = d \|\nabla f(x^k)\|_2^2$$

$$\mathbb{E}[\|x^{k+1} - x^*\|_2^2] = \mathbb{E}\|x^k - x^*\|_2^2 - \underbrace{2\gamma \mathbb{E} \langle \nabla f(x^k); x^k - x^* \rangle}_{\text{noise term}} + \gamma^2 d \mathbb{E} \|\nabla f(x^k)\|_2^2$$

μ -strongly conv, L -smoothness

$$\leq \mathbb{E}\|x^k - x^*\|_2^2 - \gamma \mu \|x^k - x^*\|_2^2 - \gamma (f(x^k) - f(x^*)) + \gamma^2 \cdot d \cdot 2L (f(x^k) - f(x^*))$$

$$= \underbrace{(1 - \gamma \mu) \|x^k - x^*\|_2^2}_{\text{linear conv}} - \gamma(1 - \gamma \cdot d \cdot 2L) (f(x^k) - f(x^*))$$

$$1 - \gamma \cdot 2dL \geq 0$$

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

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

$$\frac{1}{\gamma/\mu} \log \frac{1}{\varepsilon} \text{ iterations}$$

$$\frac{dL}{\mu} \log \frac{1}{\varepsilon} \text{ iterations}$$