

$$f(x) = \frac{1}{n} \sum_{i=1}^n f_i(x) \quad - \quad \text{Empirical Risk Minimization Problem}$$

Problem 1

$$\boxed{x_{k+1} = x_k - \eta \nabla f_i(x_k)} \quad \text{SGD}$$

$i$  is sampled uniformly at random from  $\{1, \dots, n\}$

The cost of computing full gradient is  $\underline{\underline{O(n \cdot d)}}$   
 stoch gradient is  $\underline{\underline{O(d)}}$   
 $n \gg 1$

1) How many samples are left unseen in expectation after one epoch ( $n$  iterations)  $\epsilon_c$

$$X_i = \begin{cases} 1, & \text{if index } i \text{ is never sampled in one epoch} \\ 0, & \text{otherwise} \end{cases}$$

$$\mathbb{E} \left[ \sum_{i=1}^n X_i \right] = ?$$

$$\sum_{i=1}^n \mathbb{E}[X_i] = n \cdot \mathbb{E}[X_1]$$

$$\begin{aligned} \mathbb{E}[X_1] &= P(X_1=1) \cdot 1 + P(X_1=0) \cdot 0 \\ &= P(X_1=1) \end{aligned}$$

$$P(\text{index } i \text{ is sampled at current iteration}) = \frac{1}{n}$$

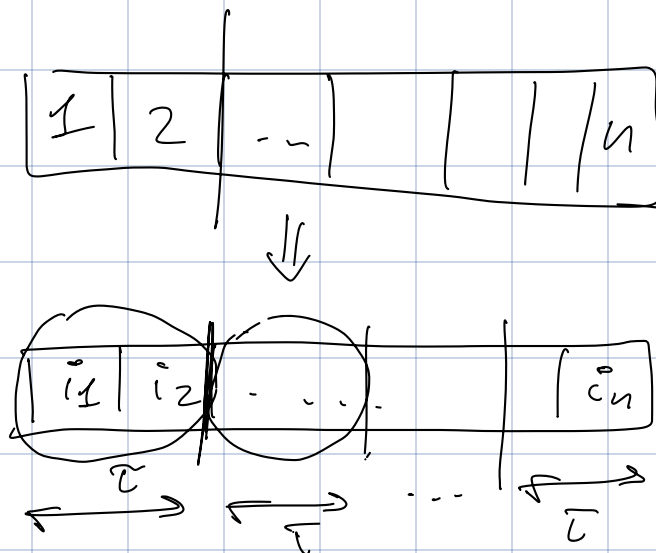
$$P(\text{index } i \text{ is not sampled at current iteration}) = 1 - \frac{1}{n}$$

$$\underline{P(X_1=1) = \left(1 - \frac{1}{n}\right)^n = \mathbb{E}[X_1]}$$

$$\mathbb{E}\left[\sum_{i=1}^n X_i\right] = n \left(1 - \frac{1}{n}\right)^n \xrightarrow{n \rightarrow \infty} \frac{n}{e} \approx 0.368 n$$

$$\lim_{n \rightarrow \infty} \left(1 - \frac{1}{n}\right)^n = e^{-1}$$

Random Reshuffling:



$$2) \mathbb{E}[\|x_{k+1} - x^*\|^2] \geq \left(\frac{1}{2L}\right)^2 \mathbb{E}[\|f(x_k)\|^2] \quad \left(\eta = \frac{1}{2L}\right)$$

$$\mathbb{E}[\|x_{u+1} - x^*\|^2] = \mathbb{E}[\|x_u - \zeta \nabla f(x_u) - x^*\|^2] \quad (\text{step of SGD})$$

$$= \mathbb{E}[\|x_u - x^*\|^2] - 2\zeta \mathbb{E}[\langle x_u - x^*, \nabla f(x_u) \rangle] + \zeta^2 \mathbb{E}[\|\nabla f(x_u)\|^2] \quad (\text{simple linear algebra})$$

$$= \mathbb{E}[\|x_u - x^*\|^2] - 2\zeta \mathbb{E}[\langle x_u - x^*, \nabla f(x_u) \rangle] + \zeta^2 \mathbb{E}[\|\nabla f(x_u)\|^2] \quad (\text{unbiasedness of } \nabla f(x_u))$$

(Cauchy-Schwartz inequality says that  $\forall a, b \in \mathbb{R}^d$

$$|\langle a, b \rangle| \leq \|a\| \cdot \|b\| \rightarrow -\langle a, b \rangle \geq -\|a\| \cdot \|b\|$$

$$\Downarrow$$

$$-\|a\| \cdot \|b\| \leq \langle a, b \rangle \leq \|a\| \cdot \|b\|$$

$$\geq \mathbb{E}[\|x_u - x^*\|^2] - 2\zeta \mathbb{E}[\|x_u - x^*\| \cdot \|\nabla f(x_u)\|] + \zeta^2 \mathbb{E}[\|\nabla f(x_u)\|^2] \quad (\text{Cauchy-Schwartz})$$

(We assume that  $f$  is  $L$ -smooth)

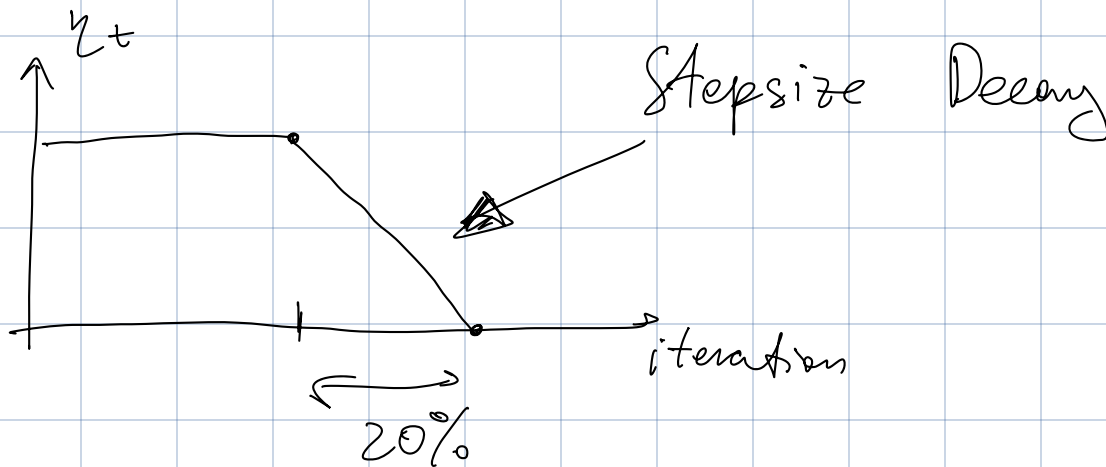
$$\|\nabla f(x_u)\| = \|\nabla f(x_u) - \nabla f(x^*)\| \leq L \|x_u - x^*\|$$

$$\geq \mathbb{E}[\|x_u - x^*\|^2] - \underbrace{(2\zeta L)}_{\ll 1} \mathbb{E}[\|x_u - x^*\|^2]$$

$$+ \zeta^2 \mathbb{E} [\|\nabla f_i(x_k)\|^2]$$

$$\zeta = \frac{1}{2L} \Rightarrow \left(\frac{1}{2L}\right)^2 \mathbb{E} [\|\nabla f_i(x_k)\|^2]$$

$$\mathbb{E} [\|x_{k+1} - x^*\|^2] \geq \left(\frac{1}{2L}\right)^2 \mathbb{E} [\|\nabla f_i(x_k)\|^2] >$$



2) Variance Reduction  
i is sampled

SVRG:  $g_k = \nabla f_i(x_k) - \nabla f_i(w_k) + \nabla f(w_k)$

SARAH:  $g_k = \nabla f_{i_k}(x_k) - \nabla f_{i_k}(w_k) + \nabla f_{i_{k-1}}(w_{k-1})$

L-SVRG

SAGA

## Problem 2

Assumptions:  $f$  is  $\mu$ -strongly convex;  $L$ -smooth  
H-Lipschitz Hessian  
B - Bounded stoch gradients.

$$\|\nabla^2 f(x) - \nabla^2 f(y)\| \leq H \cdot \|x - y\| \quad (\text{H-Lipschitz of } H_{\text{H}})$$
$$\|\nabla f_i(x_k)\|^2 \leq B^2 \quad \forall i \in \{1, \dots, h\}$$

$$1) g(x) = \|\nabla f(x)\|^2$$

$$\nabla g(x) = \begin{pmatrix} \frac{\partial}{\partial x_1} g(x) \\ \vdots \\ \frac{\partial}{\partial x_d} g(x) \end{pmatrix} - \text{gradient of } g(x)$$

$$\begin{aligned} \frac{\partial}{\partial x_u} g(x) &= \frac{\partial}{\partial x_u} \|\nabla f(x)\|^2 = \frac{\partial}{\partial x_u} \left( \sum_{\ell=1}^d \left( \frac{\partial f}{\partial x_\ell} \right)^2 \right) \\ &= \sum_{\ell=1}^d \frac{\partial}{\partial x_u} \left( \frac{\partial f}{\partial x_\ell} \right)^2 = \sum_{\ell=1}^d 2 \cdot \frac{\partial f}{\partial x_\ell} \frac{\partial}{\partial x_u} \left( \frac{\partial f}{\partial x_\ell} \right) \\ &= \sum_{\ell=1}^d 2 \cdot \underbrace{\frac{\partial^2 f}{\partial x_u \partial x_\ell}}_{\rightarrow} \cdot \underbrace{\frac{\partial f}{\partial x_\ell}}_{\rightarrow} \quad (=) \end{aligned}$$

$$\nabla^2 f(x) = \begin{pmatrix} \frac{\partial^2 f}{\partial x_1 \partial x_1} & \dots & \frac{\partial^2 f}{\partial x_1 \partial x_d} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_d \partial x_1} & \dots & \frac{\partial^2 f}{\partial x_d \partial x_d} \end{pmatrix} \quad \text{— Hessian}$$

$$\Leftrightarrow 2 [\nabla^2 f(x) \cdot \nabla f(x)]_k$$

$$\boxed{\nabla g(x) = 2 \nabla^2 f(x) \cdot \nabla f(x)}$$

Let's show that  $g$  is  $\tilde{L}$ -smooth

$$\|\nabla g(x) - \nabla g(y)\| = \left\| \frac{2 \nabla^2 f(x) \nabla f(x)}{= \nabla g(x)} - \frac{2 \nabla^2 f(y) \nabla f(y)}{= \nabla g(y)} \right\|$$

$$= 2 \left\| \underbrace{\nabla^2 f(x) \nabla f(x) - \nabla^2 f(y) \nabla f(x)}_A + \underbrace{\nabla^2 f(y) \nabla f(x) - \nabla^2 f(y) \nabla f(y)}_B \right\|$$

$$= 2 \left\| \underbrace{(\nabla^2 f(x) - \nabla^2 f(y)) \nabla f(x)}_A + \underbrace{\nabla^2 f(y) (\nabla f(x) - \nabla f(y))}_B \right\|$$

(Triangle inequality:  $\|a+b\| \leq \|a\| + \|b\| \quad \forall a, b \in \mathbb{R}^d$ )

$$\leq 2 \left\| \underbrace{(\nabla^2 f(x) - \nabla^2 f(y))}_A \underbrace{\nabla f(x)}_x \right\| + 2 \left\| \underbrace{\nabla^2 f(y)}_A \underbrace{(\nabla f(x) - \nabla f(y))}_x \right\|$$

(for some matrix  $A$  we have  $\|Ax\| \leq \|A\| \cdot \|x\|$   
 some vector  $x$ )

$$\leq 2 \underbrace{\|\nabla^2 f(x) - \nabla^2 f(y)\|} \cdot \underbrace{\|\nabla f(x)\|} + 2 \underbrace{\|\nabla^2 f(y)\|} \cdot \underbrace{\|\nabla f(x) - \nabla f(y)\|}$$

$$\leq 2 \underbrace{H \|x-y\|} \cdot \underbrace{\left\| \frac{1}{n} \sum_{i=1}^n \nabla f_i(x) \right\|} + 2 \underbrace{\|\nabla^2 f(y)\|} \cdot \underbrace{L \|x-y\|}$$

Triangle  
inequality

$$\leq 2H \|x-y\| \cdot \left( \frac{1}{n} \sum_{i=1}^n \|\nabla f_i(x)\| \right) + 2L \cdot L \|x-y\|$$

$$\left( \nabla f \text{ is } L\text{-Lipschitz} \Leftrightarrow \nabla^2 f(x) \preceq L \cdot \text{Id} \right)$$

$$(\|\nabla^2 f(x)\| \leq L)$$

$$\leq 2H \cdot \|x-y\| \underbrace{\left( \frac{1}{n} \sum_{i=1}^n B \right)}_B + 2L^2 \|x-y\|$$

$$= \underbrace{2(BH + L^2)}_{\tilde{L}} \|x-y\|$$

$$g(x) - \tilde{L}\text{-smooth with } \tilde{L} = 2(BH + L^2)$$

$$g(x_{n+1}) \leq g(x_n) + \langle \nabla g(x_n), x_{n+1} - x_n \rangle + \frac{\tilde{L}}{2} \|x_{n+1} - x_n\|^2$$

(smoothness inequality)

$$= g(x_n) + \langle 2 \nabla^2 f(x_n) \nabla f(x_n), -\zeta \nabla f(x_n) \rangle + \frac{\tilde{L}}{2} \|- \zeta \nabla f(x_n)\|^2$$

$$= g(x_n) - 2\zeta \langle \nabla^2 f(x_n) \nabla f(x_n), \nabla f(x_n) \rangle + \frac{\tilde{L}\zeta^2}{2} \|\nabla f(x_n)\|^2$$

Let's use unbiasedness of  $\nabla f(x_n)$ :

$$\mathbb{E}[g(x_{n+1})] \leq \mathbb{E}[g(x_n)] - 2\zeta \mathbb{E}[\langle \nabla^2 f(x_n) \nabla f(x_n), \nabla f(x_n) \rangle] + \frac{\tilde{L}\zeta^2}{2} \mathbb{E}[\|\nabla f(x_n)\|^2]$$

$$\mathbb{E}[g(x_{n+1})] \leq \mathbb{E}[g(x_n)] - \underbrace{2\zeta \mathbb{E}[\langle \nabla^2 f(x_n) \nabla f(x_n), \nabla f(x_n) \rangle]}_{\geq \mu \zeta \mathbb{E}[\|\nabla f(x_n)\|^2]} + \frac{\tilde{L}\zeta^2}{2} \mathbb{E}[\|\nabla f(x_n)\|^2]$$

$$\nabla^2 f(x_n) \geq \mu \cdot \text{Id} \Rightarrow$$



$$\Rightarrow \nabla f(x_u)^T \nabla^2 f(x_u) \nabla f(x_u) \geq \underbrace{\nabla f(x_u)^T \mu \cdot Id \cdot \nabla f(x_u)}_{= \mu \|\nabla f(x_u)\|^2}$$

$$\Rightarrow -\nabla f(x_u)^T \nabla^2 f(x_u) \nabla f(x_u) \leq -\mu \|\nabla f(x_u)\|^2$$

$$\mathbb{E}[g(x_{u+1})] \leq \mathbb{E}[g(x_u)] - 2\gamma\mu \mathbb{E}[\|\nabla f(x_u)\|^2] + \frac{\tilde{L}^2\gamma^2}{2} \mathbb{E}[\|\nabla f(x_u)\|^2]$$

$$= \mathbb{E}[g(x_u)] - 2\gamma\mu \cdot \mathbb{E}[g(x_u)] + \frac{\tilde{L}^2\gamma^2}{2} \underbrace{\mathbb{E}[\|\nabla f(x_u)\|^2]}_{\leq B^2}$$

$$\leq (1-2\gamma\mu) \mathbb{E}[g(x_u)] + \frac{\tilde{L}^2\gamma^2 B^2}{2}$$

Final inequality:

$$\mathbb{E}[g(x_{u+1})] \leq (1-2\gamma\mu) \mathbb{E}[g(x_u)] + \frac{\tilde{L}^2\gamma^2 B^2}{2}$$

$$\begin{aligned}
&\leq (1-2\gamma\mu) \left( (1-2\gamma\mu) \mathbb{E}[g(x_{n-1})] + \frac{\tilde{L}^2 B^2}{2} \right) + \frac{\tilde{L}^2 B^2}{2} \\
&= (1-2\gamma\mu)^2 \mathbb{E}[g(x_{n-1})] + \frac{\tilde{L}^2 B^2}{2} \left( 1 + (1-2\gamma\mu) \right) \\
&\dots \\
&\leq (1-2\gamma\mu)^{k+1} \mathbb{E}[g(x_0)] + \frac{\tilde{L}^2 B^2}{2} \sum_{\ell=0}^k \underbrace{(1-2\gamma\mu)^\ell}_q
\end{aligned}$$

$$\sum_{\ell=0}^{\infty} q^\ell = \frac{1}{1-q} \quad \text{for } q < 1 \Rightarrow \sum_{\ell=0}^k q^\ell \leq \frac{1}{1-q}$$

$$\leq (1-2\gamma\mu)^{k+1} \mathbb{E}[g(x_0)] + \frac{\tilde{L}^2 B^2}{2} \cdot \frac{1}{1-(1-2\gamma\mu)}$$

$$= (1-2\gamma\mu)^{k+1} \mathbb{E}[g(x_0)] + \frac{\tilde{L}^2 B^2}{2} \cdot \frac{1}{2\gamma\mu}$$

$$\mathbb{E}[g(x_n)] \leq (1-2\gamma\mu)^k g(x_0) + \frac{\tilde{L}^2 B^2}{4\mu}$$

### Problem 3

1)  $L$ -smoothness of  $f$ :

$$\begin{aligned} f(x_{u+1}) &\leq f(x_u) + \langle \nabla f(x_u), x_{u+1} - x_u \rangle + \frac{L}{2} \|x_{u+1} - x_u\|^2 \\ &= f(x_u) + \langle \nabla f(x_u), -\zeta \nabla f_i(x_u) \rangle + \frac{L}{2} \|- \zeta \nabla f_i(x_u)\|^2 \\ &= f(x_u) - \zeta \langle \nabla f(x_u), \nabla f_i(x_u) \rangle + \frac{L\zeta^2}{2} \|\nabla f_i(x_u)\|^2 \end{aligned}$$

$$\begin{aligned} \mathbb{E}[f(x_{u+1})] &\leq \mathbb{E}[f(x_u)] - \zeta \mathbb{E}[\langle \nabla f(x_u), \nabla f_i(x_u) \rangle] \\ &\quad + \frac{L\zeta^2}{2} \mathbb{E}[\|\nabla f_i(x_u)\|^2] \end{aligned}$$

Unbiasedness

$$\begin{aligned} &= \mathbb{E}[f(x_u)] - \zeta \mathbb{E}[\langle \nabla f(x_u), \nabla f(x_u) \rangle] \\ &\quad + \frac{L\zeta^2}{2} B^2 \end{aligned}$$

$\|\nabla f_i(x_u)\| \leq B$

$$= \mathbb{E}[f(x_u)] - \zeta \mathbb{E}[\|\nabla f(x_u)\|^2] + \frac{L\zeta^2}{2} B^2$$

$\mu$ -PL inequality:  $\|\nabla f(x_u)\|^2 \geq 2\mu (f(x_u) - f^*)$

$$\Leftrightarrow -\|\nabla f(x_u)\|^2 \leq -2\mu (f(x_u) - f^*)$$

$$\leq \mathbb{E}[f(x_u)] - 2\gamma\mu \mathbb{E}[f(x_u) - f^*] + \frac{L\gamma^2 B^2}{2}$$

We get

$$\mathbb{E}[f(x_{u+1})] \leq \mathbb{E}[f(x_u)] - 2\gamma\mu \mathbb{E}[f(x_u) - f^*] + \frac{L\gamma^2 B^2}{2}$$

Let's subst.  $-f^*$

$$\mathbb{E}[f(x_{u+1}) - f^*] \leq \mathbb{E}[f(x_u) - f^*] - 2\gamma\mu \mathbb{E}[f(x_u) - f^*] + \frac{L\gamma^2 B^2}{2}$$

$$2) \mathbb{E}[f(x_{u+1}) - f^*] \leq \mathbb{E}[f(x_u) - f^*] - 2\mu \cdot \gamma_u \mathbb{E}[f(x_u) - f^*] + \frac{L\gamma_u^2 B^2}{2}$$

$$\gamma_u = \frac{2^{k+1}}{2\mu(k+1)^2}$$

$$\mathbb{E}[f(x_{u+1}) - f^*] \leq \mathbb{E}[f(x_u) - f^*] - 2\mu \cdot \frac{2^{k+1}}{2\mu(k+1)^2} \mathbb{E}[f(x_u) - f^*]$$

$$+ \frac{L \cdot (2k+1)^2}{2 \cdot 4\mu^2 (k+1)^4} B^2$$

$$\leq \left(1 - \frac{2k+1}{(k+1)^2}\right) \mathbb{E}[f(x_k) - f^*]$$

$$+ \frac{L \cdot (2k+2)^2}{2 \cdot 4\mu^2 (k+1)^4} B^2$$

$$= \left(\frac{(k+1)^2 - (2k+1)}{(k+1)^2}\right) \mathbb{E}[f(x_k) - f^*]$$

$$+ \frac{L \cdot 4 \cdot (k+1)^2}{2 \cdot 4\mu^2 (k+1)^4} B^2$$

$$= \frac{k^2}{(k+1)^2} \mathbb{E}[f(x_k) - f^*]$$

$$+ \frac{L B^2}{2\mu^2 (k+1)^2}$$

We obtain:  $\mathbb{E}[f(x_{k+1}) - f^*] \leq \frac{k^2}{(k+1)^2} \mathbb{E}[f(x_k) - f^*]$

$$+ \frac{LB^2}{2\mu^2(k+1)^2}$$

$$\mathbb{E} \left[ \underbrace{(k+1)^2 (f(x_{k+1}) - f^*)}_{\delta_f(k+1)} \right] \leq \mathbb{E} \left[ \underbrace{k^2 (f(x_k) - f^*)}_{\delta_f(k)} \right] + \frac{LB^2}{2\mu^2}$$

$$\mathbb{E} [\delta_f(k+1)] \leq \mathbb{E} [\delta_f(k)] + \frac{LB^2}{2\mu^2}$$

$$\leq \mathbb{E} [\delta_f(k-1)] + 2 \cdot \frac{LB^2}{2\mu^2}$$

$$\leq \mathbb{E} [\delta_f(k-2)] + 3 \cdot \frac{LB^2}{2\mu^2}$$

.....

$$\leq \mathbb{E} [\delta_f(0)] + (k+1) \cdot \frac{LB^2}{2\mu^2}$$

$$\delta_f(0) = 0^2 \cdot (f(x_0) - f^*) = 0$$

$$\mathbb{E} [\delta_f(k+1)] \leq (k+1) \cdot \frac{LB^2}{2\mu^2}$$

This implies:

$$\mathbb{E} [k^2 (f(x_n) - f^*)] \leq k \cdot \frac{LB^2}{2\mu^2}$$

$$\mathbb{E} [f(x_n) - f^*] \leq \frac{1}{k} \cdot \frac{LB^2}{2\mu^2} \quad \text{— convergence rate}$$

Or:  $\|x_{n+1} - x^*\|^2 = \|x_n - x^*\|^2 - \dots$