

# Opt Lec 8

## • Stochastic Gradient Descent II

### 1. Stochastic Mirror Descent (S-MD)

$$\Rightarrow x_{t+1} = \operatorname{argmin}_{x \in \mathcal{X} \cap D} \gamma_t \tilde{g}(x_t)^T x + D_\phi(x, x_t) \quad \text{where } \mathbb{E}[\tilde{g}(x_t)] \in \partial f(x_t)$$

★ Assumptions for simple analysis

①  $\phi$  is a mirror map that is 1-strongly convex on  $\mathcal{X} \cap D$  w.r.t.  $\|\cdot\|$

$$② R^* = \sup_{x \in \mathcal{X} \cap D} \phi(x) - \phi(x_*)$$

### 2. Mini-batch SGD $\rightarrow$ reduce variance

$$\Rightarrow x_{t+1} = x_t - \gamma \cdot \frac{1}{m} \sum_{i=1}^m g_i(x_t) \quad \text{where } g_i(x_t) \text{'s are i.i.d. stochastic gradients}$$

(need more computation than SGD)

$$\Rightarrow \mathbb{E}\left[\left\|\frac{1}{m} \sum_{i=1}^m g_i(x) - \nabla f(x)\right\|^2\right] = \mathbb{E}\left[\|g_i(x) - \nabla f(x)\|^2\right]$$

### 3. Epoch-based algorithm (SVRG) $\rightarrow$ reduce variance

For  $s=1, 2, \dots, \infty$   $\rightarrow$  outer loop

$$x_i^{(s)} = y^{(s)}$$

For  $t=1, 2, \dots, k$   $\rightarrow$  inner loop

Draw  $i_t^{(s)}$  uniformly at random

$$x_{t+1}^{(s)} = x_t^{(s)} - \gamma \left( \nabla f_{i_t^{(s)}}(x_t^{(s)}) - \nabla f_{i_t^{(s)}}(y^{(s)}) + \frac{1}{n} \sum_{i=1}^n \nabla f_i(y^{(s)}) \right)$$

$$y^{(s+1)} = \frac{1}{K} \sum_{t=1}^K x_t^{(s)}$$

$\Rightarrow$   $K \uparrow \rightarrow$  computation time  $\downarrow$ , variance  $\uparrow$   
 $K \downarrow \rightarrow$  computation time  $\uparrow$ , variance  $\downarrow$

★ Convergence rate for stochastic mirror descent

\*  $f$  is convex and  $\mathbb{E}[\|\tilde{g}(x)\|_*^2] \leq B^2$  ( $\gamma = \frac{R}{B}\sqrt{\frac{2}{T}}$ )

- $D_\phi(x_t, y_{t+1}) - D_\phi(x_{t+1}, y_{t+1})$

$$\begin{aligned}
 &= \cancel{\phi(x_t)} - \cancel{\phi(y_{t+1})} - \nabla \phi(y_{t+1})^\top (x_t - y_{t+1}) - \cancel{\phi(x_{t+1})} + \cancel{\phi(y_{t+1})} + \nabla \phi(y_{t+1})^\top (x_{t+1} - y_{t+1}) \\
 &= \phi(x_t) - \phi(x_{t+1}) - \nabla \phi(y_{t+1})^\top (x_t - x_{t+1}) \\
 &\leq (\nabla \phi(x_t) - \nabla \phi(y_{t+1}))^\top (x_t - x_{t+1}) - \frac{1}{2} \|x_t - x_{t+1}\|^2 \\
 &= \gamma \cdot g_t^\top (x_t - x_{t+1}) - \frac{1}{2} \|x_t - x_{t+1}\|^2 \\
 &\leq \gamma \cdot \|g_t\|_* \|x_t - x_{t+1}\| - \frac{1}{2} \|x_t - x_{t+1}\|^2 \\
 &\leq \frac{\gamma^2 \|g_t\|_*^2}{2}
 \end{aligned}$$

- $g_t^\top (x_t - x) = \frac{1}{\gamma} (\nabla \phi(x_t) - \nabla \phi(y_{t+1}))^\top (x_t - x)$

$$\begin{aligned}
 &= \frac{1}{\gamma} (D_\phi(x_t, y_{t+1}) + D_\phi(x, x_t) - D_\phi(x, y_{t+1})) \\
 &\leq \frac{1}{\gamma} (D_\phi(x_t, y_{t+1}) + D_\phi(x, x_t) - D_\phi(x, x_{t+1}) - D_\phi(x_{t+1}, y_{t+1})) \\
 \Rightarrow \sum_{t=1}^T g_t^\top (x_t - x) &\leq \frac{1}{\gamma} \cdot (D_\phi(x, x_1) - D_\phi(x, x_{T+1}) + \sum_{t=1}^T \frac{\gamma^2 \|g_t\|_*^2}{2}) \\
 &\leq \frac{R^2}{\gamma} + \frac{\gamma}{2} \sum_{t=1}^T \|g_t\|_*^2
 \end{aligned}$$

- $\mathbb{E} \left[ f \left( \frac{1}{T} \sum_{t=1}^T x_t \right) - f(x^*) \right] \leq \frac{1}{T} \mathbb{E} \left[ \sum_{t=1}^T (f(x_t) - f(x^*)) \right]$

$$\begin{aligned}
 &\leq \frac{1}{T} \mathbb{E} \left[ \sum_{t=1}^T \mathbb{E} [g_t^\top] (x_t - x^*) \right] \\
 &\leq \frac{1}{T} \left( \frac{R^2}{\gamma} + \frac{\gamma}{2} \sum_{t=1}^T \mathbb{E} [\|g_t\|_*^2] \right) \\
 &\leq \frac{1}{T} \left( \frac{R^2}{\gamma} + \frac{\gamma}{2} T B^2 \right) \\
 &= \frac{1}{T} \left( R B \sqrt{\frac{T}{2}} + R B \sqrt{\frac{T}{2}} \right) = R B \sqrt{\frac{2}{T}} \quad (\gamma = \frac{R}{B} \sqrt{\frac{2}{T}})
 \end{aligned}$$

(same as mirror-descent)

\*  $f$  is convex,  $\beta$ -smooth and  $\mathbb{E}[\|\nabla f(x) - \tilde{g}(x)\|_*^2] \leq \sigma^2$  ( $\gamma = \frac{1}{\beta + \frac{\sigma}{R}\sqrt{\frac{1}{T}}}$ )

- $\hat{f}(x_{t+1}) - \hat{f}(x_t)$

$$\leq \nabla \hat{f}(x_t)^T (x_{t+1} - x_t) + \frac{\beta}{2} \|x_{t+1} - x_t\|^2$$

$$= g_t^T (x_{t+1} - x_t) + (\nabla \hat{f}(x_t) - g_t)^T (x_{t+1} - x_t) + \frac{\beta}{2} \|x_{t+1} - x_t\|^2$$

$$\leq g_t^T (x_{t+1} - x_t) + \frac{\eta}{2} \|\nabla \hat{f}(x_t) - g_t\|_*^2 + \frac{1}{2} \left(\beta + \frac{1}{\eta}\right) \|x_{t+1} - x_t\|^2$$

( $\because$  Cauchy-Schwarz :  $2ab \leq x a^2 + b^2/x$ )

$$\leq g_t^T (x_{t+1} - x_t) + \frac{\eta}{2} \|\nabla \hat{f}(x_t) - g_t\|_*^2 + \left(\beta + \frac{1}{\eta}\right) D_\phi(x_{t+1}, x_t)$$

( $\because$  1-strongly convex :  $\phi(x_{t+1}) - \phi(x_t) - \nabla \phi(x_t)^T (x_{t+1} - x_t) \geq \frac{1}{2} \|x_{t+1} - x_t\|^2$ )

$$= g_t^T (x^* - x_t) + g_t^T (x_{t+1} - x^*) + \left(\beta + \frac{1}{\eta}\right) D_\phi(x_{t+1}, x_t) + \frac{\eta}{2} \|\nabla \hat{f}(x_t) - g_t\|_*^2$$

$$\leq g_t^T (x^* - x_t) + \left(\beta + \frac{1}{\eta}\right) (\nabla \phi(x_t) - \nabla \phi(y_{t+1}))^T (x_{t+1} - x^*) + \left(\beta + \frac{1}{\eta}\right) D_\phi(x_{t+1}, x_t) + \frac{\eta}{2} \|\nabla \hat{f}(x_t) - g_t\|_*^2$$

$$\leq g_t^T (x^* - x_t) + \left(\beta + \frac{1}{\eta}\right) (\nabla \phi(x_t) - \nabla \phi(y_{t+1}))^T (x_{t+1} - x^*) + \left(\beta + \frac{1}{\eta}\right) D_\phi(x_{t+1}, x_t) + \frac{\eta}{2} \|\nabla \hat{f}(x_t) - g_t\|_*^2$$

( $\because \nabla_{x_{t+1}} D(x_{t+1}, y_{t+1})^T (x_{t+1} - x^*) = (\nabla \phi(x_{t+1}) - \nabla \phi(y_{t+1}))^T (x_{t+1} - x^*) \leq 0$ )

$$\leq g_t^T (x^* - x_t) + \left(\beta + \frac{1}{\eta}\right) (D_\phi(x^*, x_t) - D_\phi(x^*, x_{t+1})) + \frac{\eta}{2} \|\nabla \hat{f}(x_t) - g_t\|_*^2$$

( $\because (\nabla \phi(x_t) - \nabla \phi(x_{t+1}))^T (x_{t+1} - x^*) = D_\phi(x^*, x_{t+1}) - D_\phi(x^*, x_t) - D_\phi(x_{t+1}, x_t)$ )

$$= \nabla \hat{f}(x_t)^T (x^* - x_t) + (g_t - \nabla \hat{f}(x_t))^T (x^* - x_t) + \left(\beta + \frac{1}{\eta}\right) (D_\phi(x^*, x_t) - D_\phi(x^*, x_{t+1}))$$

$$+ \frac{\eta}{2} \|\nabla \hat{f}(x_t) - g_t\|_*^2$$

$$\leq \hat{f}(x^*) - \hat{f}(x_t) + (g_t - \nabla \hat{f}(x_t))^T (x^* - x_t) + \left(\beta + \frac{1}{\eta}\right) (D_\phi(x^*, x_t) - D_\phi(x^*, x_{t+1}))$$

$$+ \frac{\eta}{2} \|\nabla \hat{f}(x_t) - g_t\|_*^2$$

$$\therefore \hat{f}(x_{t+1}) - \hat{f}(x^*) \leq (g_t - \nabla \hat{f}(x_t))^T (x^* - x_t) + \left(\beta + \frac{1}{\eta}\right) (D_\phi(x^*, x_t) - D_\phi(x^*, x_{t+1}))$$

$$+ \frac{\eta}{2} \|\nabla \hat{f}(x_t) - g_t\|_*^2$$

$$\Rightarrow \mathbb{E}[\hat{f}(x_{t+1}) - \hat{f}(x^*)] \leq \left(\beta + \frac{1}{\eta}\right) \mathbb{E}[D_\phi(x^*, x_t) - D_\phi(x^*, x_{t+1})] + \frac{\eta}{2} \sigma^2$$

$$\Rightarrow \mathbb{E}[\hat{f}(\frac{1}{T} \sum_{t=1}^T x_{t+1}) - \hat{f}(x^*)] \leq \frac{1}{T} \left(\beta + \frac{1}{\eta}\right) \mathbb{E}[D_\phi(x^*, x_t)] + \frac{\eta}{2} \sigma^2$$

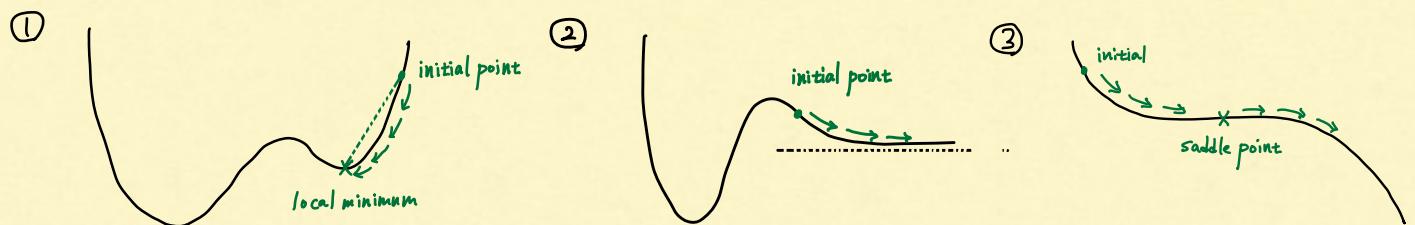
$$= \frac{1}{T} \left(\beta + \frac{\sigma}{R} \sqrt{\frac{1}{T}}\right) R^2 + \frac{1}{2} \cdot \frac{R}{\sigma} \sqrt{\frac{2}{T}} \sigma^2 \quad (\eta = \frac{R}{\sigma} \sqrt{\frac{2}{T}})$$

$$= \frac{\beta R^2}{T} + R \sigma \sqrt{\frac{2}{T}}$$

# Opt Lec 9

- Non-convex optimization

## 1. Gradient descent on non-convex functions

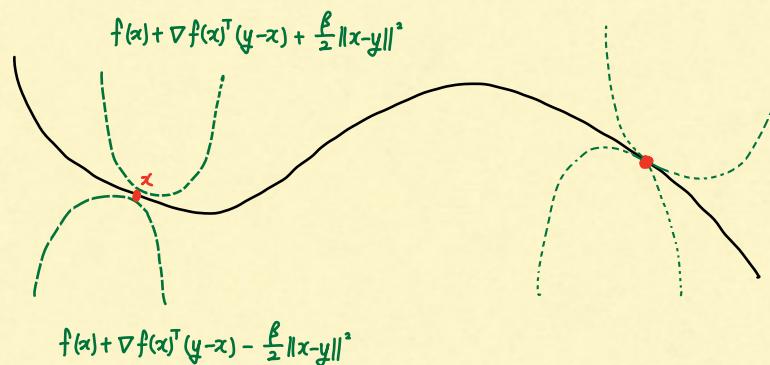


## 2. Smooth with parameter $\beta$

: a differentiable  $f$  satisfies follows over a convex set  $X \subset \text{dom}(f)$

$$\textcircled{1} \quad f(y) \leq f(x) + \nabla f(x)^T (y-x) + \frac{\beta}{2} \|x-y\|^2 \quad \forall x, y \in X$$

$$\textcircled{2} \quad f(y) \geq f(x) + \nabla f(x)^T (y-x) - \frac{\beta}{2} \|x-y\|^2 \quad \forall x, y \in X$$

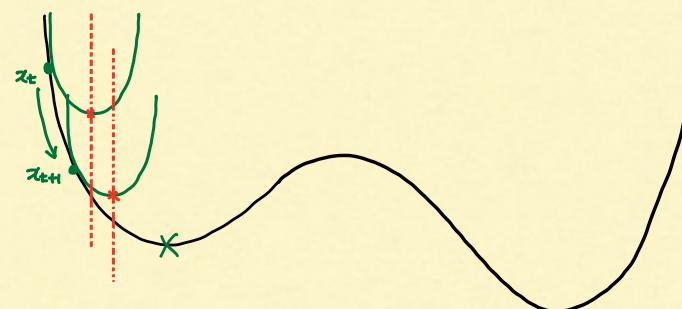


$\Rightarrow$  if  $f$  is twice differentiable, bounded Hessian ( $\|\nabla^2 f(x)\| \leq \beta, \forall x \in X$ ) = smooth

$(\|\cdot\| : \text{spectral norm} \rightarrow \text{maximum singular value of a matrix}$   
 $\rightarrow \text{maximum scale of vectors})$

\* No overshooting with  $\gamma = \frac{1}{\beta}$

$\Rightarrow$  Never go across the local minimum to reach global minimum



④ Convergence rate for gradient descent

\*  $f$  is twice differentiable and smooth with parameter  $\beta$  ( $R = \frac{1}{\beta}$ )

$$\bullet f(x_{t+1}) - f(x_t) \leq \nabla f(x_t)^T (x_{t+1} - x_t) + \frac{\beta}{2} \|x_{t+1} - x_t\|^2$$

$$= -\frac{1}{\beta} \|\nabla f(x_t)\|^2 + \frac{1}{2\beta} \|\nabla f(x_t)\|^2$$

$$= -\frac{1}{2\beta} \|\nabla f(x_t)\|^2 \leq 0 \quad \rightarrow \text{non-increasing}$$

$$\Rightarrow \sum_{t=0}^{T-1} \|\nabla f(x_t)\|^2 \leq 2\beta (f(x_0) - f(x_T)) \leq 2\beta (f(x_0) - f(x^*))$$

$$\Rightarrow \sum_{t=0}^{\infty} \|\nabla f(x_t)\|^2 \leq 2\beta (f(x_0) - f(x^*)) < \infty$$

$$\therefore \lim_{t \rightarrow \infty} \|\nabla f(x_t)\|^2 = 0 \quad \Rightarrow \quad \|\nabla f(x_t)\| \rightarrow 0$$

(but,  $f(x_t) - f(x^*)$  may not converge to 0)

## \* Linear model

### ① Linear regression ( $y_i = w x_i + \varepsilon_i$ )

$$\Rightarrow w^* = \underset{w}{\operatorname{argmin}} \sum_{i=1}^n (w x_i - y_i)^2 \quad (= f(w)) \rightarrow \text{convex}$$

$$\rightarrow \nabla^2 f(w) = 2 \cdot \sum_{i=1}^n x_i \cdot x_i^\top$$

$$\rightarrow \|\nabla^2 f(w)\| \leq \rho (\geq 0) \quad (\text{largest eigenvalue})$$

→ gradient descent with  $\gamma = \frac{1}{\rho}$

### ② Linear Neural Net ( $y_i = w_L \cdots w_2 w_1 x_i + \varepsilon_i$ )

$$\Rightarrow (w_1^*, \dots, w_L^*) = \underset{w_1, \dots, w_L}{\operatorname{argmin}} \sum_{i=1}^n (w_L \cdots w_2 w_1 x_i - y_i)^2$$

$\Downarrow x_i=1, y_i=1, n=1$  (relaxation)

$$\Rightarrow (w_1^*, \dots, w_L^*) = \underset{w_1, \dots, w_L}{\operatorname{argmin}} (w_L \cdots w_2 w_1 - 1)^2 \quad (= f(w)) \rightarrow \text{non-convex}$$

$$\rightarrow \nabla f(w) = 2(w_L \cdots w_2 w_1 - 1) \begin{bmatrix} w_L \cdots w_3 w_2 \\ w_L \cdots w_3 w_1 \\ \vdots \\ w_{L-1} \cdots w_2 w_1 \end{bmatrix} = 0$$

$$\rightarrow \begin{cases} w_L \cdots w_2 w_1 = 1 \rightarrow f(w) = 0 \quad (\text{global minimum}) \\ \text{At least two of } w_i \text{'s are 0} \rightarrow f(w) = 1 \quad (\text{local minimum}) \end{cases}$$

⇒ assume  $w^{(0)}$  satisfies C-balanced conditions and  $w_1^{(0)} w_2^{(0)} \cdots w_L^{(0)} < 1$

→ C-balanced :  $w_i \leq C w_j \quad \forall i \neq j, \quad w_i \geq 0 \quad \forall i$

→  $w^{(1)} = w^{(0)} - \gamma \nabla f(w^{(0)})$  is again C-balanced.

$$\begin{aligned} \text{pf)} \quad \nabla f(w^{(0)})_i &= 2(w_L \cdots w_2 w_1 - 1) \cdot \frac{w_L \cdots w_1}{w_i} \\ (w^{(1)})_i &= (w^{(0)})_i - 2\gamma (w_L \cdots w_2 w_1 - 1) \cdot \frac{w_L \cdots w_1}{w_i} \\ &\leq C \left\{ (w^{(0)})_j - 2\gamma (w_L \cdots w_2 w_1 - 1) \cdot \frac{w_L \cdots w_1}{C w_i} \right\} \\ &\leq C \left\{ (w^{(0)})_j - 2\gamma (w_L \cdots w_2 w_1 - 1) \cdot \frac{w_L \cdots w_1}{w_j} \right\} \\ &= C (w^{(1)})_j \end{aligned}$$

∴  $w^{(1)}$  is C-balanced

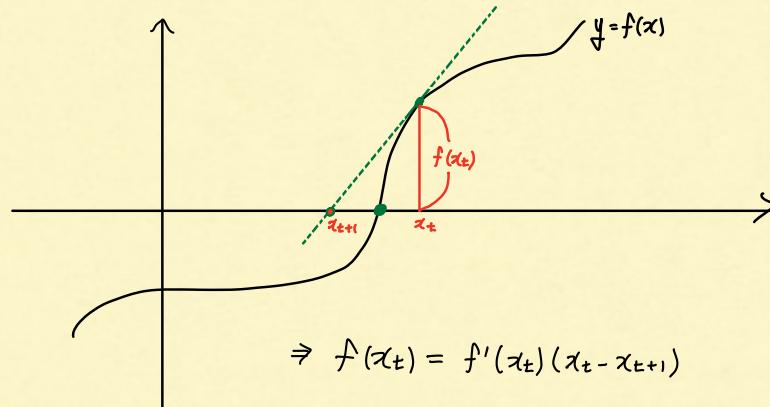
# Opt Lec 10

## • Newton's method

### 1. Newton-Raphson method

: find a zero of differentiable function  $f$

$$\Rightarrow x_{t+1} = x_t - \frac{f(x_t)}{f'(x_t)}$$



$\Rightarrow$  Babylonian method ( $f(x) = x^2 - R$ )

: find a square root of  $R$ . assuming  $0 < x_t - \sqrt{R} < 0.5$  and  $R \geq \frac{1}{4}$

$$\rightarrow x_{t+1} = x_t - \frac{x_t^2 - R}{2x_t} = \frac{1}{2} \left( x_t + \frac{R}{x_t} \right)$$

$$\rightarrow x_{t+1} - \sqrt{R} = \frac{1}{2} \cdot \frac{1}{x_t} (x_t^2 + R - \sqrt{R}x_t) = \frac{1}{2x_t} (x_t - \sqrt{R})^2$$

$$\rightarrow x_t - \sqrt{R} \leq (x_{t+1} - \sqrt{R})^2 \leq \dots \leq (x_0 - \sqrt{R})^{2^T} \leq 2^{-2^T} \rightarrow 0$$

### 2. Newton's method

: find a critical point  $\rightarrow$  global minimum of a differentiable convex function  $f$

$\Rightarrow$  control the update scale and direction according to the local geometry at  $x_t$

$\Rightarrow$  equivalent to gradient descent when the second derivative is constant

#### ① 1-dimensional case

$\Rightarrow$  apply Newton-Raphson method to  $f'$

$$\Rightarrow x_{t+1} = x_t - \frac{f'(x_t)}{f''(x_t)}$$

#### ② d-dimensional case

$$\Rightarrow x_{t+1} = x_t - \nabla^2 f(x_t)^{-1} \nabla f(x_t)$$

## \* Affine invariant

: independent to affine change of coordinates

$$\begin{bmatrix} x_t = g(y_t) = Ay_t + b \\ x_{t+1} = x_t - \nabla_{x_t}^2 f(x_t)^{-1} \nabla_{x_t} f(x_t) \\ y_{t+1} = g^{-1}(x_{t+1}) = A^{-1}(x_{t+1} - b) \end{bmatrix} \Leftrightarrow y_{t+1} = y_t - \nabla_{y_t}^2 f \circ g(y_t)^{-1} \nabla_{y_t} f \circ g(y_t)$$

$$pf) \quad \nabla_{y_t}^2 f \circ g(y_t) = A^T \nabla_{x_t}^2 f(Ay_t + b) = A^T \nabla_{x_t}^2 f(x_t)$$

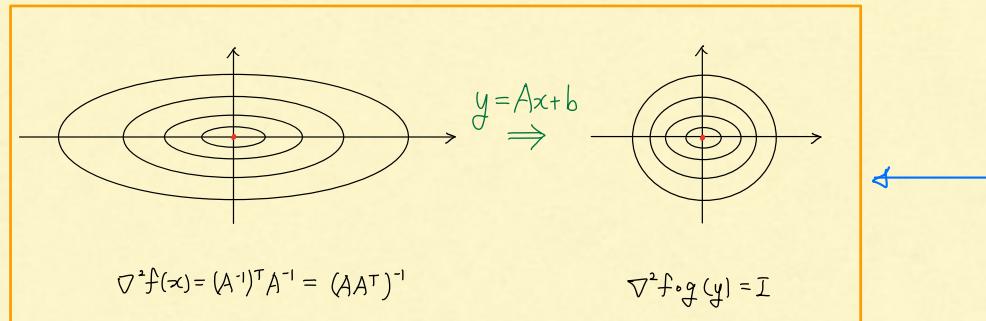
$$\nabla_{y_t}^2 f \circ g(y_t) = A^T \nabla_{x_t}^2 f(Ay_t + b) A = A^T \nabla_{x_t}^2 f(x_t) A$$

$$y_{t+1} = y_t - \nabla_{y_t}^2 f \circ g(y_t)^{-1} \nabla_{y_t} f \circ g(y_t)$$

$$\Rightarrow y_{t+1} = y_t - A^{-1} \nabla_{x_t}^2 f(x_t)^{-1} A^T \cdot A^T \nabla_{x_t} f(x_t)$$

$$\Rightarrow Ay_{t+1} + b = Ay_t + b - \nabla_{x_t}^2 f(x_t)^{-1} \nabla_{x_t} f(x_t)$$

$$\Rightarrow x_{t+1} = x_t - \nabla_{x_t} f(x_t)^{-1} \nabla f(x_t)$$



## \* Another interpretation

$$\Rightarrow x_{t+1} = \underset{x}{\operatorname{argmin}} \left( f(x_t) + \nabla f(x_t)^T (x - x_t) + \frac{1}{2} (x - x_t)^T \nabla^2 f(x_t) (x - x_t) \right)$$

$$\rightarrow \nabla f(x_t) + \nabla^2 f(x_t) (x_{t+1} - x_t) = 0$$

$$\rightarrow x_{t+1} = x_t - \nabla^2 f(x_t)^{-1} \nabla f(x_t)$$

$\Rightarrow$  One-shot algorithm when  $\nabla^2 f(x) = H^{-1} x$

$\Rightarrow$  Super fast when  $\nabla^2 f(x)$  is small ( $\approx$  smooth)

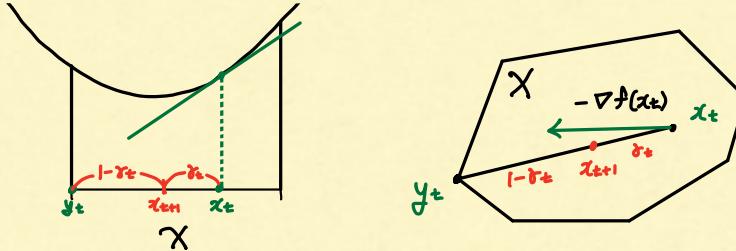
# Opt Lec 11

- Frank-Wolfe

## 1. Frank-Wolfe Algorithm

$$\textcircled{1} \quad y_t \in \arg \min_{y \in X} \nabla f(x_t)^T y \quad \rightarrow \text{boundary point}$$

$$\textcircled{2} \quad x_{t+1} = (1 - \gamma_t) x_t + \gamma_t y_t$$



## 2. Lasso regression

$$\Rightarrow \min_x \|Ax - b\|^2 \quad \text{s.t. } \|x\|_1 \leq 1$$

$$\Rightarrow y_t = -\text{sign}(\nabla f(x_t)_i) e_i \quad \text{where } i = \arg \max_i |\nabla f(x)_i|$$

$$\Rightarrow x_{t+1} = (1 - \gamma_t) x_t + \gamma_t y_t$$

(simpler than projection on L1-ball.)

## 3. Duality gap ( $g(x)$ )

$$\Rightarrow g(x_t) := \nabla f(x_t)^T (x_t - y_t) = \max_{y \in X} \nabla f(x_t)^T (x_t - y)$$

$$\geq \nabla f(x_t)^T (x_t - x^*)$$

$$\geq f(x_t) - f(x^*) \geq 0 \quad (\because \text{convex})$$

$\Rightarrow$  duality gap is non-negative

## 4. Step size

$$\textcircled{1} \quad \gamma_t = \frac{2}{t+1}$$

$$\textcircled{2} \quad \text{line search : } \gamma_t = \arg \min_{\gamma \in [0,1]} f((1-\gamma_t)x_t + \gamma_t y_t)$$

$$\textcircled{3} \quad \text{Gap-based : } \gamma_t = \min \left\{ \frac{g(x_t)}{\beta \|y_t - x_t\|^2}, 1 \right\} \quad \text{when } f \text{ is } \beta\text{-smooth}$$

★ Convergence rate for Frank-Wolfe

\*  $f$  is convex,  $\beta$ -smooth and  $R = \sup_{x,y \in X} \|x-y\|$  ( $\gamma_t = \frac{2}{t+1}$  for  $t \geq 1$ )

$$\rightarrow \text{For } T \geq 2, \quad f(x_T) - f(x^*) \leq \frac{2\beta R^2}{T+1}$$

$$pf) \quad f(x_{t+1}) - f(x_t)$$

$$\leq \nabla f(x_t)^T (x_{t+1} - x_t) + \frac{\beta}{2} \|x_{t+1} - x_t\|^2 \quad (\because \beta\text{-smooth})$$

$$= \nabla f(x_t)^T (-\gamma_t x_t + \gamma_t y_t) + \frac{\beta}{2} \| -\gamma_t x_t + \gamma_t y_t \|^2 \quad (\because x_{t+1} = (1-\gamma_t)x_t + \gamma_t y_t)$$

$$\leq \gamma_t \nabla f(x_t)^T (y_t - x_t) + \frac{\beta}{2} \gamma_t^2 R^2 \quad (\because R = \sup_{x,y \in X} \|x-y\|)$$

$$\leq \gamma_t \nabla f(x_t)^T (x^* - x_t) + \frac{\beta}{2} \gamma_t^2 R^2 \quad (\because y_t = \arg \min_y \nabla f(x_t)^T y)$$

$$\leq \gamma_t (f(x^*) - f(x_t)) + \frac{\beta}{2} \gamma_t^2 R^2 \quad (\because \text{convex})$$

$$\therefore f(x_{t+1}) - f(x^*) \leq (1-\gamma_t) (f(x_t) - f(x^*)) + \frac{\beta}{2} \gamma_t^2 R^2$$

\* Mathematical induction

$$\textcircled{1} \quad T=2$$

$$\rightarrow f(x_2) - f(x^*) \leq \frac{\beta}{2} R^2 \leq \frac{2}{3} R^2$$

$$\textcircled{2} \quad T=k$$

$$\rightarrow f(x_k) - f(x^*) \leq \frac{2\beta R^2}{k+1}$$

$$\textcircled{3} \quad T=k+1$$

$$\begin{aligned} \rightarrow f(x_{k+1}) - f(x^*) &\leq \frac{k-1}{k+1} (f(x_k) - f(x^*)) + \frac{2\beta}{(k+1)^2} R^2 \\ &\leq \frac{k}{k+1} \cdot \frac{2\beta R^2}{k+1} \leq \frac{2\beta R^2}{k+2} \end{aligned}$$

# Opt Lec 12

- Coordinate descent

## 1. Coordinate descent algorithm

: modify only one coordinate at a time

① Select  $i \in [d]$  scalar

②  $x_{t+1} = x_t - \gamma \boxed{\nabla_{i_t} f(x_t)} e_{i_t}$  where  $\nabla_j f(x)$  indicates  $\frac{\partial f(x)}{\partial x_j}$

$\Rightarrow$  random coordinate descent

: coordinate index is selected via uniform distribution  $P(i_t = i) = \frac{1}{d}$

$$\rightarrow \mathbb{E}[\nabla_{i_t} f(x_t) e_{i_t}] = \frac{1}{d} \nabla f(x_t) \quad (\because \sum_i \nabla_i f(x_t) e_i = \nabla f(x_t))$$

## 2. Coordinate-wise $\beta$ -smooth

$$: |\nabla_i f(x + r e_i) - \nabla_i f(x)| \leq \beta |r| \quad \forall x \in \mathbb{R}^d, \forall r \in \mathbb{R}, \forall i \in [d]$$

$$\Rightarrow f(x + r e_i) \leq f(x) + r \nabla_i f(x) + \frac{\beta}{2} r^2 \quad \forall x \in \mathbb{R}^d, \forall r \in \mathbb{R}, \forall i \in [d]$$

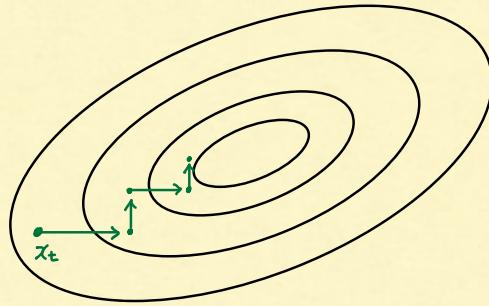
$$\Rightarrow |(\nabla^2 f)_{ii}| \leq \beta \quad (\text{necessary condition for } \beta\text{-smooth})$$

## 3. Coordinate-wise strongly convex with $\alpha$

$$: |\nabla_i f(x + r e_i) - \nabla_i f(x)| \geq \alpha |r| \quad \forall x \in \mathbb{R}^d, \forall r \in \mathbb{R}, \forall i \in [d]$$

$$\Rightarrow f(x + r e_i) \geq f(x) + r \cdot \nabla_i f(x) + \frac{\alpha}{2} r^2 \quad \forall x \in \mathbb{R}^d, \forall r \in \mathbb{R}, \forall i \in [d]$$

$$\Rightarrow |(\nabla^2 f)_{ii}| \geq \alpha \quad (\text{necessary condition for } \alpha\text{-strongly convex})$$



★ Convergence rate for coordinate descent

1)  $P(i_t = i) = \frac{1}{d} \Rightarrow \text{uniform sampling}$

\*  $f$  is  $\alpha$ -strongly convex and  $\beta$ -smooth with  $\gamma = \frac{1}{\beta}$

$$\begin{aligned}
 \bullet f(x_{t+1}) - f(x_t) &\leq \nabla f(x_t)^T (x_{t+1} - x_t) + \frac{\beta}{2} \|x_{t+1} - x_t\|^2 \\
 &\leq -\frac{1}{\beta} \nabla f(x_t)^T \nabla_{i_t} f(x_t) e_{i_t} + \frac{1}{2\beta} \|\nabla_{i_t} f(x_t) e_{i_t}\|^2 \\
 &\leq -\frac{1}{\beta} |\nabla_{i_t} f(x_t)|^2 + \frac{1}{2\beta} |\nabla_{i_t} f(x_t)|^2 \quad (\because \nabla_i f(x) = \nabla f(x)^T e_i) \\
 &= -\frac{1}{2\beta} |\nabla_{i_t} f(x_t)|^2 \\
 \bullet \mathbb{E}[f(x_{t+1}) - f(x_t)] &\leq -\frac{1}{2\beta d} \|\nabla f(x_t)\|^2 \\
 &\leq -\frac{1}{d} \frac{\alpha}{\beta} (f(x_t) - f(x^*)) \quad (\because \text{Polyak-Lojasiewicz condition}) \\
 \bullet \mathbb{E}[f(x_T) - f(x^*)] &\leq \left(1 - \frac{1}{d} \cdot \frac{\alpha}{\beta}\right) (f(x_{T-1}) - f(x^*)) \\
 &\leq \left(1 - \frac{1}{d} \cdot \frac{\alpha}{\beta}\right)^T (f(x_0) - f(x^*))
 \end{aligned}$$

(SGD:  $\mathbb{E}[f(x_T) - f(x^*)] \leq \left(1 - \frac{\alpha}{\beta}\right)^T (f(x_0) - f(x^*)) + \frac{\sigma^2}{2\alpha}$ )

\*  $f$  is coordinate-wise  $\beta$ -smooth and  $\alpha$ -strongly convex

$$\textcircled{1} \quad r = \frac{1}{\beta}$$

- $$f(x_{t+1}) - f(x_t) \leq \nabla_{i_t} f(x_t) e_{i_t}^T (x_{t+1} - x_t) + \frac{\beta}{2} \|x_{t+1} - x_t\|^2$$

$$= -\frac{1}{\beta} \|\nabla_{i_t} f(x_t) e_{i_t}\|^2 + \frac{1}{2\beta} \|\nabla_{i_t} f(x_t) e_{i_t}\|^2$$

$$\leq -\frac{1}{2\beta} \|\nabla_{i_t} f(x_t)\|^2$$

- $$\mathbb{E}[f(x_{t+1}) - f(x_t)] \leq -\frac{1}{2\beta d} \|\nabla f(x_t)\|^2$$

$$\leq -\frac{1}{d} \frac{\alpha}{\beta} (f(x_t) - f(x^*)) \quad (\because \text{Polyak-Lojasiewicz condition})$$

- $$\mathbb{E}[f(x_T) - f(x^*)] \leq \left(1 - \frac{1}{d} \cdot \frac{\alpha}{\beta}\right) (f(x_0) - f(x^*))$$

$$\leq \left(1 - \frac{1}{d} \cdot \frac{\alpha}{\beta}\right)^T (f(x_0) - f(x^*))$$

$$\textcircled{2} \quad r = \frac{1}{\beta_{\max}}$$

- $$\mathbb{E}[f(x_T) - f(x^*)] \leq \left(1 - \frac{1}{d} \cdot \frac{\alpha}{\beta_{\max}}\right)^T (f(x_0) - f(x^*))$$

$$(\because \mathbb{E}[r \nabla_{i_t} f(x_t) e_{i_t}] = \frac{1}{d \beta_{\max}} \nabla f(x_t))$$

2)  $P(i_t = i) = \frac{\beta_i}{\sum_j \beta_j} \Rightarrow \text{importance sampling}$

\*  $f$  is coordinate-wise  $\beta_i$ -smooth and  $\alpha$ -strongly convex with  $r_i = \frac{1}{\beta_i}$

- $$\mathbb{E}[f(x_T) - f(x^*)] \leq \left(1 - \frac{\alpha}{\sum_j \beta_j}\right)^T (f(x_0) - f(x^*)) \rightarrow \text{faster.}$$

$$(\because \mathbb{E}[r_i \nabla_{i_t} f(x_t) e_{i_t}] = \frac{1}{\beta_i} \frac{\beta_i}{\sum_j \beta_j} \nabla f(x_t) = \frac{1}{\sum_j \beta_i} \nabla f(x_t))$$

3)  $i_t = \operatorname{argmax}_i |\nabla_i f(x_t)| \Rightarrow \text{steepest coordinate descent}$

\*  $f$  is coordinate-wise  $\beta$ -smooth and  $\alpha$ -strongly convex with  $r = \frac{1}{\beta}$

- $$\mathbb{E}[f(x_T) - f(x^*)] \leq \left(1 - \frac{1}{d} \cdot \frac{\alpha}{\beta}\right)^T (f(x_0) - f(x^*)) \rightarrow \text{uniform sampling}$$

$$(\because \max_i |\nabla_i f(x_t)|^2 \geq \frac{1}{d} \sum_j |\nabla_j f(x)|^2)$$

# Opt Lec 13

## Duality

1. Constrained optimization problem

$\Rightarrow$  Primal ( $p^*$ )

$$\min_x f(x) \longrightarrow \text{convex}$$

$$\text{s.t. } h_i(x) = 0 \quad | \leq i \leq m \longrightarrow \text{affine}$$

$$g_j(x) \leq 0 \quad | \leq j \leq r \longrightarrow \text{convex}$$

$\Rightarrow$  Dual ( $d^*$ )

$$\max_{\mu, \lambda} L(\mu, \lambda)$$

$$\text{s.t. } \lambda_j \geq 0 \quad | \leq j \leq r$$

Lagrange function ( $\Delta(x, \mu, \lambda)$ )

$$\rightarrow \Delta(x, \mu, \lambda) = f(x) + \sum_{i=1}^m \mu_i h_i(x) + \sum_{j=1}^r \lambda_j g_j(x)$$

Lagrange dual function ( $L(\mu, \lambda)$ )

$$\rightarrow L(\mu, \lambda) = \min_x \Delta(x, \mu, \lambda)$$

$\Rightarrow$  primal optimal solution is lower bounded by dual feasible solution

$$\rightarrow p^* = \inf_x f(x) \geq \inf_x \Delta(x, \mu, \lambda) = L(\mu, \lambda)$$

$\Rightarrow$  duality gap and stopping criterion can be defined as follow

$$\underbrace{p^* - d^*}_{\text{duality gap}} \leq p^* - L(\mu^{(k)}, \lambda^{(k)}) \leq \underbrace{f(x^{(k)}) - L(\mu^{(k)}, \lambda^{(k)})}_{\text{stopping criterion}} \leq \varepsilon$$

## 2. Conjugate function ( $f^*$ )

$$: f^*(y) = \sup_{x \in \text{dom}(f)} y^T x - f(x)$$

$\Rightarrow f^*$  is always convex as it is a pointwise supremum of affine function of  $y$

$$\begin{aligned} \rightarrow f^*(\alpha y_1 + (1-\alpha)y_2) &= \sup_{x \in \text{dom}(f)} (\alpha y_1 + (1-\alpha)y_2)^T x - f(x) \\ &\leq \alpha \left( \sup_{x \in \text{dom}(f)} y_1^T x - f(x) \right) + (1-\alpha) \left( \sup_{x \in \text{dom}(f)} y_2^T x - f(x) \right) \\ &= \alpha f^*(y_1) + (1-\alpha) f^*(y_2) \end{aligned}$$

\*  $f(x) + f^*(y) \geq x^T y \quad \forall x \in X, \forall y \in X^* \quad (\text{Fenchel's inequality})$

pf)  $f(x) + f^*(y) \geq f(x) + \sup_z z^T y - f(z) \geq f(x) + x^T y - f(x) = x^T y$

\*  $f^{**} \leq f$

( : Fenchel's inequality )

pf)  $f^{**}(x) = \sup_y y^T x - f^*(y) \leq \sup_y f(x) = f(x) \quad \forall x \in X$

(  $f^{**}$  is a closed convex of  $f$ . )

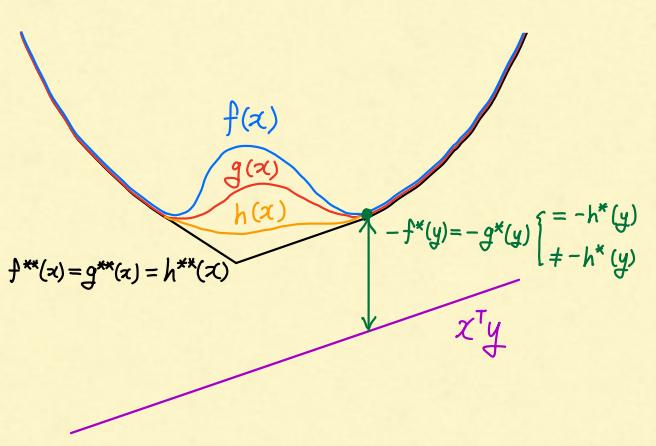
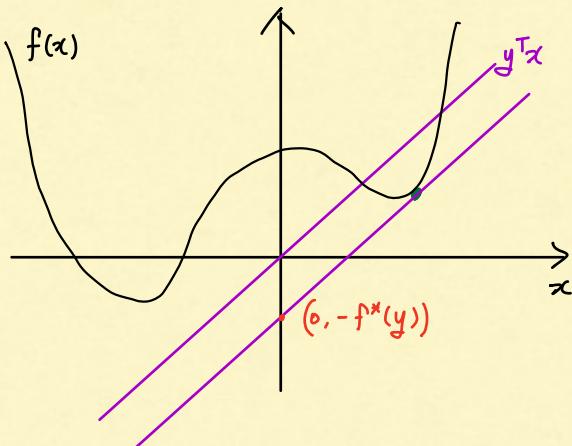
\*  $f$  is closed and convex  $\Rightarrow f^{**} = f^*$

\*  $f$  is closed and convex  $\Rightarrow y \in \partial f(x) \Leftrightarrow x \in \partial f^*(y) \Leftrightarrow x^T y = f(x) + f^*(y)$

pf) ( $\Rightarrow$ )  $y \in \partial f(x) \Rightarrow f^*(y) = \sup_z z^T y - f(z) = x^T y - f(x)$   
 $\Rightarrow f^*(y) = x^T y - f(x) = x^T y - f^{**}(x) \quad (\because f^{**} = f)$

$$\begin{aligned} \Rightarrow f^{**}(x) &= x^T y - f^*(y) = \sup_z z^T x - f^*(z) \\ \Rightarrow x &\in \partial f^*(y) \end{aligned}$$

( $\Leftarrow$ )  $x \in \partial f^*(y) \Rightarrow f^{**}(x) = \sup_z z^T x - f^*(z) = y^T x - f^*(y)$   
 $\Rightarrow f^{**}(x) = f(x) = x^T y - f^*(y) \quad (\because f^{**} = f)$   
 $\Rightarrow f^*(y) = x^T y - f(x) = \sup_z z^T y - f(z)$   
 $\Rightarrow y \in \partial f(x)$



$\Rightarrow$  Maximum gap between  $x^T y$  and  $f(x)$        $\Rightarrow$  minimum gap between  $x^T y$  and  $f(x)$

### 3. Generalized linear models

$\Rightarrow$  primal ( $p^*$ )

$$\min_{x, w} f(w) + g(x)$$

$$\text{s.t. } w = Ax$$

$\Rightarrow$  dual ( $d^*$ )

$$\max_u \left( \min_{x, w} f(w) + g(x) + u^T (Ax - w) \right)$$

$$= \max_u \left( \min_{x, w} -(u^T w - f(w)) - (-u^T Ax - g(x)) \right)$$

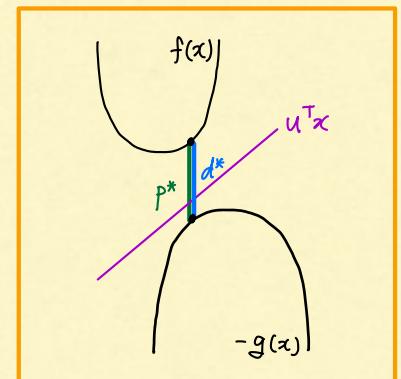
$$\geq \max_u (-f^*(u) - g^*(-u^T A))$$

$$* A = I$$

$$\Rightarrow \min_x f(x) + g(x)$$

$$\Rightarrow \max_u -f^*(u) - g^*(-u)$$

$$= \max_u \left( \min_w (f(w) - u^T w) + \min_w (u^T w - (-g(w))) \right)$$



\* conjugate function of constant function is infinite

$$\Rightarrow g(x) = c^T x \longrightarrow g^*(u) = \infty \quad \forall u$$

# Opt Lec 14

## • Variance reduction

### 1. Variance reduction

: Reducing variance of  $X$  by using variable  $Y$  that is positively correlated to  $X$  with known expectation

$$\Rightarrow Z_\alpha = \alpha(X - Y) + \mathbb{E}[Y]$$

$$\Rightarrow \mathbb{E}[Z_\alpha] = \alpha \mathbb{E}[X] + (1-\alpha) \mathbb{E}[Y] \rightarrow \text{unbiased when } \alpha=1$$

$$\Rightarrow \text{Var}[Z_\alpha] = \alpha^2 (\text{Var}(X) + \text{Var}(Y) - 2\text{Cov}(X, Y)) \rightarrow \text{smaller variance when } \alpha < 1$$

#### ① SVRG

$$\rightarrow X = \nabla f_{i_t}(\theta_t)$$

$$\rightarrow Y = \nabla f_{i_t}(\tilde{\theta}), \quad \mathbb{E}[Y] = \frac{1}{n} \sum_i \nabla f_i(\tilde{\theta}) \rightarrow \text{high computation}$$

$$\rightarrow \alpha = 1 \rightarrow \text{unbiased}$$

$$\therefore x_{t+1} = x_t - \gamma \left( \nabla f_{i_t}(\theta_t) - \nabla f_{i_t}(\tilde{\theta}) + \frac{1}{n} \sum_i \nabla f_i(\tilde{\theta}) \right)$$

#### ② SAG

$$\rightarrow X = \nabla f_{i_t}(\theta_t)$$

$$\rightarrow Y = y_{i_t}^{(t-1)}, \quad \mathbb{E}[Y] = \frac{1}{n} \sum_{i=1}^n y_i^{(t-1)} \rightarrow \text{high memory cost}$$

$$y_i^{(t)} = \begin{cases} \nabla f_i(\theta_t) & \text{if } i = i_t \\ y_i^{(t-1)} & \text{otherwise} \end{cases}$$

$$\rightarrow \alpha = \frac{1}{n} \rightarrow \text{biased}$$

$$\therefore x_{t+1} = x_t - \gamma \left( \frac{1}{n} \left( \nabla f_{i_t}(\theta_t) - y_{i_t}^{(t-1)} \right) + \frac{1}{n} \sum_i y_i^{(t-1)} \right)$$

### ③ SAGA

$$\rightarrow X = \nabla f_{i_t}(\theta_t)$$

$$\rightarrow Y = y_{i_t}^{(t-1)}, \quad \mathbb{E}[Y] = \frac{1}{n} \sum_{i=1}^n y_i^{(t-1)} \rightarrow \text{high memory cost}$$

$$y_i^{(t)} = \begin{cases} \nabla f_i(\theta_t) & \text{if } i = i_t \\ y_i^{(t-1)} & \text{otherwise} \end{cases}$$

$$\rightarrow \alpha = 1 \rightarrow \text{unbiased}$$

$$\therefore x_{t+1} = x_t - \gamma \left( \nabla f_{i_t}(\theta_t) - y_{i_t}^{(t-1)} + \frac{1}{n} \sum_i y_i^{(t-1)} \right)$$

SVRG	SAGA	
2 Grad	1 Grad	Computation
$\tilde{\theta}, \frac{1}{n} \sum_{i=1}^n \nabla f_i(\tilde{\theta})$	$[y_1^t, y_2^t, \dots, y_n^t]$	Memory storage

# Opt Lec 15

- Distributed deep learning

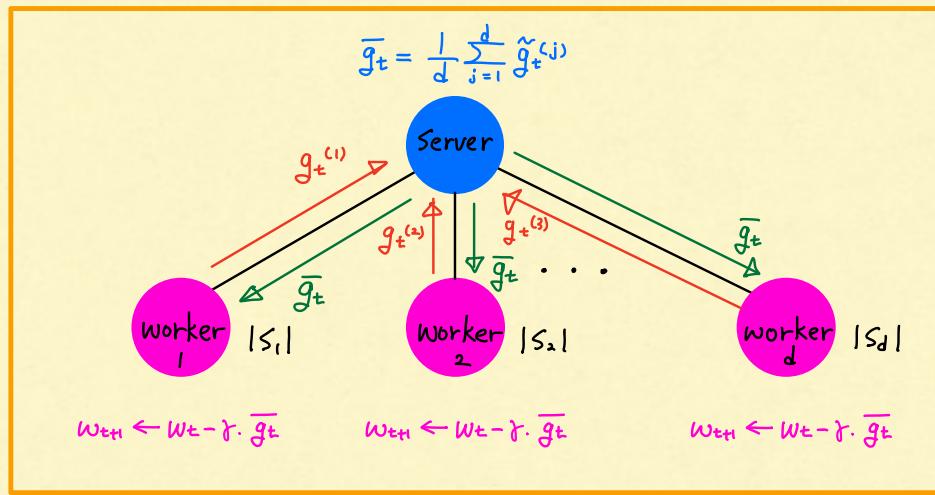
## 1. Distributed system

: data is distributed to separate machines as data can not fit to one device anymore

⇒ data can not be sent to outside → privacy

⇒ communication bottleneck occurs

⇒ distributed SGD with data parallelism



## ① Tern Grad

$$\Rightarrow \tilde{g}_t = \text{ternalize}(g_t) = s_t \cdot \text{sign}(g_t) \circ b_t$$

scalar  
↑  
binary vector  
↓ elementwise product

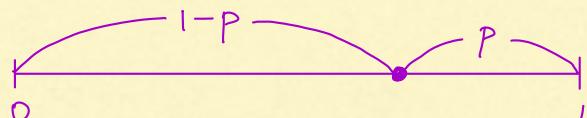
$$\rightarrow s_t = \max(\text{abs}(g_t))$$

$$\rightarrow (b_t)_k = \begin{cases} 1 & \text{with probability } \frac{|(g_t)_k|}{s_t} \\ 0 & \text{with probability } 1 - \frac{|(g_t)_k|}{s_t} \end{cases}$$

$$\Rightarrow \mathbb{E}[(b_t)_k] = \frac{|(g_t)_k|}{s_t}$$

$$\Rightarrow \mathbb{E}[\tilde{g}_t] = g_t \quad (\text{unbiased})$$

$$\Rightarrow \begin{cases} g_t \in \mathbb{R}^P \rightarrow 32P \\ \tilde{g}_t \in \mathbb{R}^P \rightarrow 2P + 32 \end{cases}$$



## ② QSGD

$$\text{(tuning parameter } \geq 1\text{)} \quad \# \text{ of quantization level}$$

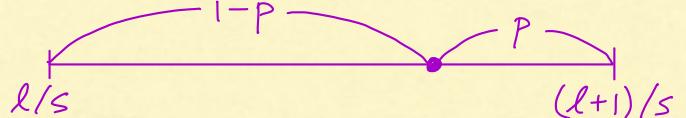
$$\Rightarrow (\tilde{g}_t)_k = Q_s((g_t)_k) = \|g_t\|_2 \cdot \text{sign}((g_t)_k) \cdot \xi_k(g_t, s)$$

$$\rightarrow \xi_k(g_t, s) = \begin{cases} l/s & \text{with probability } 1 - \left( \frac{|(g_t)_k|}{\|g_t\|_2} s - l \right) \\ (l+1)/s & \text{with probability } \frac{|(g_t)_k|}{\|g_t\|_2} s - l \end{cases}$$

$$\rightarrow 0 \leq l < s \quad (\text{integer}) \Rightarrow \frac{|(g_t)_k|}{\|g_t\|_2} \in [l/s, (l+1)/s]$$

$$\Rightarrow \mathbb{E}[\xi_k(g_t, s)] = \frac{|(g_t)_k|}{\|g_t\|_2}$$

$$\Rightarrow \mathbb{E}[(\tilde{g}_t)_k] = (g_t)_k \quad (\text{unbiased})$$



$$\Rightarrow \text{Var}((\tilde{g}_t)_k) = \min\left(\frac{k}{s^2}, \frac{\sqrt{k}}{s}\right) \|g_t\|_2^2$$

$$\Rightarrow \begin{cases} g_t \in \mathbb{R}^P \rightarrow 32P \\ \tilde{g}_t \in \mathbb{R}^P \rightarrow 2P + 32 \end{cases}$$

# Opt Lec 16

- Deep learning optimization

## 1. SGD rules

① SGD

$$\therefore \theta_{t+1} = \theta_t - \eta_t \nabla \ell(\theta_t)$$

② Momentum

$$\therefore v_{t+1} = \gamma v_t + \nabla \ell(\theta_t)$$

③ Nesterov

$$\therefore v_{t+1} = \gamma v_t + \nabla \ell(\theta_t)$$

$$\theta_{t+1} = \theta_t - \eta_t v_{t+1}$$

$$\theta_{t+1} = \theta_t - \eta_t (\gamma v_{t+1} + \nabla \ell(\theta_t))$$

④ Adagrad

$$\therefore v_{t+1} = v_t + \nabla \ell(\theta_t)^2$$

$$\theta_{t+1} = \theta_t - \frac{\eta_t}{\sqrt{v_{t+1} + \epsilon}} \nabla \ell(\theta_t)$$

⑤ RMSProp

$$\therefore v_{t+1} = \rho v_t + (1-\rho) \nabla \ell(\theta_t)^2$$

$$m_{t+1} = \gamma m_t + \frac{\eta_t}{\sqrt{v_{t+1} + \epsilon}} \nabla \ell(\theta_t)$$

$$\theta_{t+1} = \theta_t - m_{t+1}$$

⑥ ADAM

$$\therefore m_{t+1} = \beta_1 m_t + (1-\beta_1) \nabla \ell(\theta_t)$$

$$v_{t+1} = \beta_2 v_t + (1-\beta_2) \nabla \ell(\theta_t)^2$$

$$b_{t+1} = \frac{\sqrt{1-\beta_2^{t+1}}}{1-\beta_1^{t+1}}$$

$$\theta_{t+1} = \theta_t - \alpha_t \frac{m_{t+1}}{\sqrt{v_{t+1} + \epsilon}} b_{t+1}$$

⑦ NADAM

$$\therefore m_{t+1} = \beta_1 m_t + (1-\beta_1) \nabla \ell(\theta_t)$$

$$v_{t+1} = \beta_2 v_t + (1-\beta_2) \nabla \ell(\theta_t)^2$$

$$b_{t+1} = \frac{\sqrt{1-\beta_2^{t+1}}}{1-\beta_1^{t+1}}$$

$$\theta_{t+1} = \theta_t - \alpha_t \frac{\beta_1 m_{t+1} + (1-\beta_1) \nabla \ell(\theta_t)}{\sqrt{v_{t+1} + \epsilon}} b_{t+1}$$

## 2. Inclusion relationship

$\Rightarrow$  SGD  $\subseteq$  Momentum ( $\gamma = 0$ )

$\Rightarrow$  SGD  $\subseteq$  Nesterov ( $\gamma = 0$ )

$\Rightarrow$  Momentum  $\subseteq$  RMSProp ( $\rho = 1, \epsilon = 0$ )

$\Rightarrow$  Momentum  $\subseteq$  ADAM ( $\alpha_t = \epsilon \eta_t (1-\gamma^t), \beta_1 = \gamma, \beta_2 = 0$ )

$\Rightarrow$  Nesterov  $\subseteq$  NADAM ( $\alpha_t = \epsilon \eta_t (1-\gamma^t), \beta_1 = \gamma, \beta_2 = 0$ )

\* Complex algorithm may not perform well with bad hyperparameters

### 3. Learning rate & Batch size

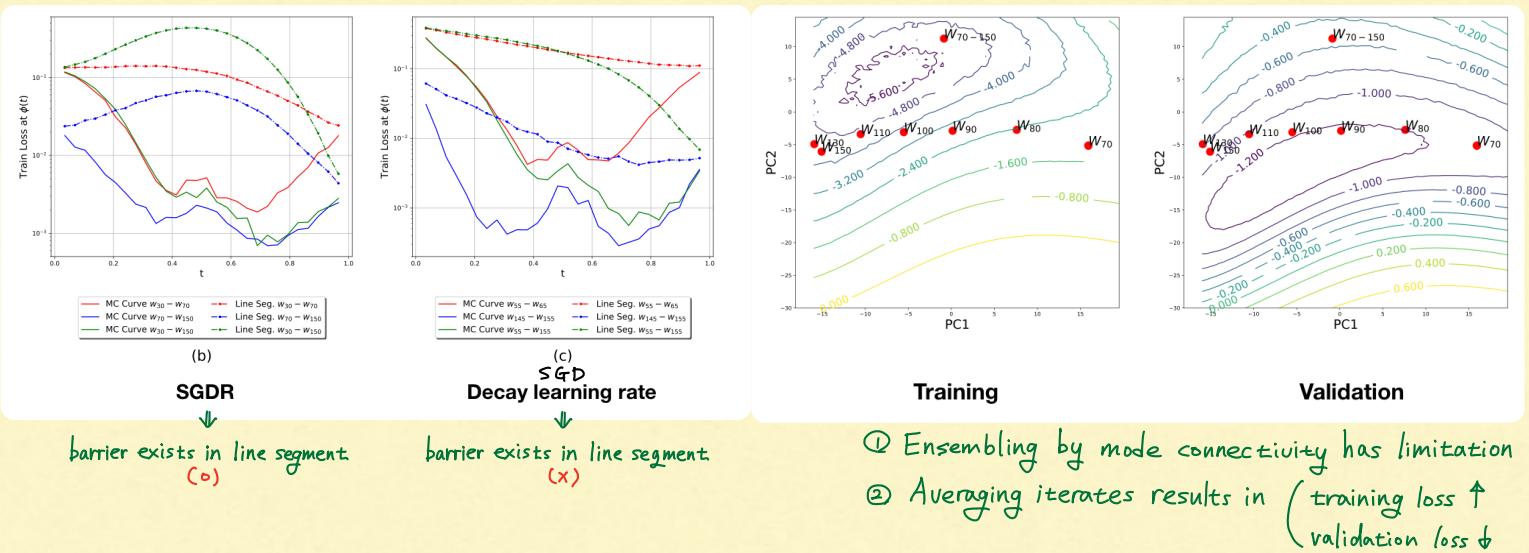
⇒ important for the following 3 cases

① objective function is non-convex

② many local minimas

③ generalization

⇒ Learning rate



⇒ Batch size

