

today: • finish BCFW  
• SAG

application of BCFW to SUM struct:

• getting  $s_i$  is one loss-augmented decoding call for example  $i$

$$\alpha_i^{(t+1)} = \alpha_i^{(t)} + \delta_t (s_i^{(t)} - \alpha_i^{(t)})$$

$\downarrow$   
LDA

$$w = A\alpha = \sum_i A_i \alpha_i$$

$\downarrow$   
 $\hat{=} w_i$

$\Rightarrow$  you update

$$w_i^{(t+1)} = w_i^{(t)} + \delta_t (w_s^{(t)} - w_i^{(t)})$$

$$\downarrow \psi_i(\hat{y}_i^{(t)})$$

$$w_j^{(t+1)} = w_j^{(t)} \quad \forall j \neq i$$

(worst case)  
 $\rightarrow O(nd)$   
storage

need to store these in memory

or  
store  $\alpha_i$ 's

(memory / computation tradeoff)

note: for line search, also need to store  $b_i^T \alpha_i^{(t)} \hat{=} l_{i,s}^{(t)}$

convergence constants

for SUM struct, can show that  $C_S^{(i)} \leq \frac{4R_i^2}{\lambda n^2}$

$$R_i \hat{=} \max_{\xi \in \mathcal{S}_i} \|\psi_i(\xi)\|_2$$

$$R = \max_i R_i$$

Hessian:  $(\lambda A^T A)_{(i,y), (i,y')}$

$$= \frac{\lambda}{\lambda^2 n^2} \langle \psi_i(y), \psi_i(y') \rangle = \frac{1}{\lambda n^2} \langle \dots \rangle$$

$$(\lambda n^2) d_i^T H_i d_i = \sum_{y, y'} d_{i,y} d_{i,y'} \langle \psi_i(y), \psi_i(y') \rangle$$

$$\leq \sum_{y, y'} |d_{i,y}| |d_{i,y'}| \underbrace{\langle \psi_i(y), \psi_i(y') \rangle}_{\leq R_i^2}$$

$$\leq R_i^2 \left( \sum_y |d_{i,y}| \right) \left( \sum_{y'} |d_{i,y'}| \right)$$

$$\leq 4R_i^2 \|d_i\|_2 \|d_i\|_2$$

$$\max_{s \in \Delta} \|s - x\|_2 = 2$$

$$\text{compare with } C_S \leq \frac{4R^2}{\lambda n^2}$$

$$d_i^T H_i d_i \leq \frac{4R_i^2}{\lambda n^2}$$

$$\Rightarrow C_S^{(i)} = \sum_{i=1}^n C_S^{(i)} \leq \frac{4R^2}{\lambda n^2} \approx \frac{C_S}{n}$$

ie. BCFW is "n times" faster than batch FW for SUM struct

skipped:

importance of affine invariance:

importance of affine invariance:  
 if you use  $\ell_2$ -norm, we get a bad bound?  $\text{diam}(M) = 2n$   
 $L \leq L_{\ell_1}$   $\text{diam}(M)^2$   
 Lipschitz constant in  $\ell_2$ -norm = largest e-value of Hessian

recall Hessian  $\lambda A^T A = \mathbb{1} \left( \langle \psi_i(y), \psi_j(y) \rangle \right)_{(i,y), (j,y)}$

say e.g.  $\langle \psi_i(y), \psi_i(y) \rangle \approx 1$  for lots of output  
 $\lambda n^2 (\mathbb{1} \mathbb{1}^T) \mathbb{1} = (\text{dim}) \cdot \mathbb{1}$

→ get largest e-value can scale with dim of a matrix  
 ⇒ really bad here because dim is exponentially

• instead, want to use  $\ell_1$ -norm of  $\Delta_{RS_i}$

i.e.  $\ell_1$ -norm for Lipschitz constant

then get  $L_i (\text{diam}(\Delta_{RS_i}))^2 \approx \frac{4B^2}{\lambda}$

Variance reduced SGD

setup:  $\min_x \frac{1}{n} \sum_{i=1}^n f_i(x)$   
 $\approx f(x)$

where  $f$  is  $\mu$ -strongly convex  
 $L$ -smooth

$f(x_t) - f(x^*) \leq (1-\rho)^t (f(x_0) - f^*)$

batch gradient method  
 [Cauchy 18<sup>th</sup> century]

$x_{t+1} = x_t - \gamma \frac{1}{n} \sum_{i=1}^n \nabla f_i(x_t)$   
 $O(n)$  to compute

$\gamma = \frac{1}{L}$  linear rate  
 where  $\rho \approx \frac{\mu}{L} = \frac{1}{K}$   
 $K \leq \frac{L}{\mu}$  "condition #"

stochastic gradient method

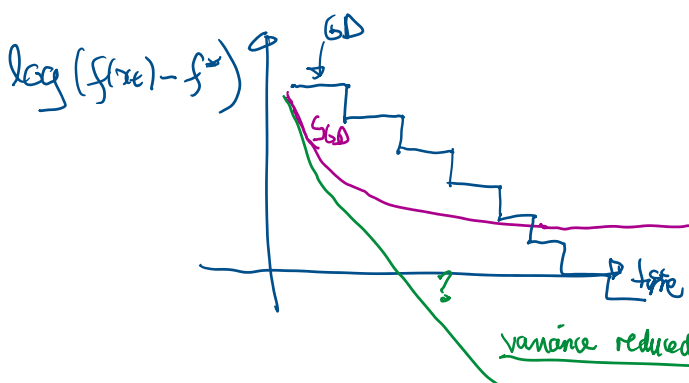
aka. incremental gradient method

[Robbins & Monro 1951]

$x_{t+1} = x_t - \gamma_t \nabla f_{i_t}(x_t)$

where  $i_t \sim \text{unif}(\{1, \dots, n\})$   
 $O(1)$  to compute

$\gamma_t = \text{const} \cdot \gamma$   
 ⇒ linear rate up to a ball of radius  $\gamma$   
 $\gamma_t = \frac{1}{\mu t}$   
 ⇒  $\tilde{O}\left(\frac{1}{\mu t}\right)$  rate  
 i.e. sublinear



variance reduced "SGD" methods?

10h39

SAG, SAGA, SVRG, SDCA, OEG, etc.

# SAG (stochastic average gradient)

[Le Roux, Schmidt & Bach NeurIPS 2012]

(Lagrange opt. prize 2018)

- SAG: • store past gradients for each  $i$  ( $g_i$ )
- update one at step  $t$

SAG { pick  $i_t$  unif.; update  $g_{i_t}^{(t+1)} = \nabla f_{i_t}(x_t)$   
 $g_j^{(t+1)} = g_j^{(t)}$  for  $j \neq i_t$

$$x_{t+1} = x_t - \gamma \left[ \frac{1}{n} \sum_{i=1}^n g_i^{(t+1)} \right] \leftarrow \text{store state gradients}$$

$$\frac{1}{n} \left[ \sum_{i=1}^n g_i^{(t)} + g_{i_t}^{(t+1)} - g_{i_t}^{(t)} \right]$$

$\frac{1}{n} \sum_i g_i \approx$  approx of  $\nabla f(x_t)$

O(1) cost per iteration      big surprise: converge linearly and fast

[but O(nd) storage cost in worst case]

"independent aggregated gradient" (IAG) [Blatt & d. 2007]  
where you cycle deterministically through  $\{1, \dots, n\}$

→ linear rate for quadratic functions

but  $\gamma_{max} \approx O\left(\frac{1}{n}\right)$  (if  $\mu$  tiny rate)

⇒ SAG Convergence rate

thm. with  $\gamma_t = \frac{1}{16L}$  where  $L = \max_i (\text{Lipschitz } \nabla f_i)$

$$\mathbb{E} f(x_t) - f^* \leq \left(1 - \min\left\{\frac{\mu}{16L}, \frac{1}{8n}\right\}\right)^t C_0$$

C<sub>0</sub> constant

$D_{SAG} = \min\left\{\frac{1}{16k_{SAG}}, \frac{1}{8n}\right\}$  vs.  $\rho_{grad} \approx \frac{1}{k_{grad}}$

example:  $l_2$ -reg. regression on RCVL

$n = 700k$      $L = 0.25$      $\mu = \frac{1}{n}$     ( $k = \frac{n}{4}$ )

rate comparison      gradient method     $\left(\frac{L-\mu}{L+\mu}\right)^2 = 0.9998$

$$\alpha(\text{accelerated gradient (Nesterov)}) (1 - \sqrt{\frac{\mu}{L}}) = 0.9976$$

$$\text{Nesterov lower bound } \left( \frac{\sqrt{L} - \sqrt{\mu}}{\sqrt{L} + \sqrt{\mu}} \right) = 0.9905$$

$$\text{SAG (n iterations)} (1 - \rho_{\text{SAG}})^n \approx 0.8825$$

practical aspects (see Schmidt & al. Math. Prog. 2016 paper)

a) storage: if  $f_i(w) = h(x_i^T w)$   $\Rightarrow \nabla f_i(w) = \underbrace{h'(x_i^T w)}_{\text{scalar}} \cdot \underbrace{x_i}_{\text{input data}}$

instead of  $O(n)$  storage  $\hookrightarrow O(1)$  storage

b) initialization of  $g_i$ 's? best: run SGD for one pass then switch SAG/SAGA

c) step size?  $\bullet \frac{1}{(4)L}$

• cheap line search heuristic (comes from FISTA)

$$\text{while } \left\{ \begin{array}{l} f_i(w_k - \frac{1}{\tilde{L}_i} \nabla f_i(w_k)) \geq f_i(w_k) - \frac{1}{2\tilde{L}_i} \|\nabla f_i(w_k)\|^2 \\ \text{set } \tilde{L}_{i,\text{new}} = 2 \tilde{L}_{i,\text{old}} \end{array} \right.$$

$$\text{else } \tilde{L}_{i,\text{new}} = \left(\frac{1}{2}\right)^n \tilde{L}_{i,\text{old}}$$

d) non-uniform sampling?

sample  $i \sim \frac{\tilde{L}_i}{\sum_j \tilde{L}_j}$

e) stopping criterion?

you can use  $\frac{1}{n} \sum_j g_j^{(t)}$  as approx  $\nabla f(x_t) = \frac{1}{n} \sum_i \nabla f_i(x_t)$

f) sparse features?

$$x_{t+1} = x_t - \gamma \left[ \underbrace{\nabla f_i(x)}_{\text{sparse}} - g_i^t + \underbrace{P_{\frac{1}{S} \sum_j g_j^t}}_{\text{dense}} \right] \quad (\text{sparse SAGA})$$

weighted projection on support of  $x_i$

$$S_i = \{u : (x_i)_u \neq 0\}$$