



for SVMstruct, can show that  $C_f^{(i)} \leq \frac{4R^2}{\lambda}$

$$\Rightarrow C_f^{(i)} \leq \frac{4R^2}{\lambda n} \approx \frac{C_f}{n}$$

$$C_f \leq \frac{4R^2}{\lambda}$$

uses affine invariance  
critically

$$C_f \leq L_{\|\cdot\|} \cdot \text{diam}_{\|\cdot\|}(M)^2$$

if use  $l_2$ -norm, get a bad bound!  $\text{diam}(M)^2 = 2n$

Lipschitz constant in  $l_2$ -norm = largest e-value of Hessian

$$\lambda A^T A = \frac{1}{\lambda n^2} \left( \langle \psi_i(y), \psi_j(y) \rangle \right)_{(i,y), (j,y)}$$

say e.g.  $\langle \psi_i(y), \psi_j(y) \rangle \approx 1$  for lots of outputs

$$(\mathbf{1}\mathbf{1}^T) \mathbf{1} = \mathbf{1} \cdot \text{dim}$$

→ get largest e-value can scale with dim. of matrix

(here, dim is exponential  
⇒ useless bound)

ie. BCFW is "n times" faster  
than batch FW for SVMstruct?

• instead, want to use  $l_2$ -norm on  $\Delta_{n \times 1}$   
 i.e.  $l_p$ -norm for Lipschitz constant...

$$l_p(v) \geq l_q(v) \text{ when } p \leq q$$

then get  $L(\text{diam}(M))^2 \approx \frac{4B^2}{\lambda}$

on  $M$ , use  $l_1$ - $l_2$  norm ( $l_1$ - $l_2$ ?)

i.e.  $\|x\| = \max_i \|x_i\|_2$

### Variance reduced SGD

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

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

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

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

$\nabla f(x_t)$   
 $O(n)$  to compute

linear rate  
 $\downarrow$   
 $\gamma = \frac{1}{L}$   
 $f(x_t) - f^* \leq (1-\rho)^t (f(x_0) - f^*)$   
 where  $\rho \approx \frac{\mu}{L} = \frac{1}{\kappa}$   
 $\kappa \triangleq \frac{L}{\mu}$  "condition #"

Stochastic gradient method  
 [Robbins & Monro 1951]

aka. incremental gradient method

$\text{Var}(f(x_t) - f^*) \ll$  batch gradient method

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

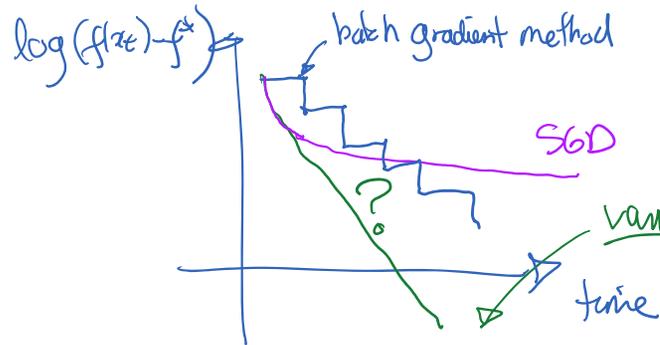
where  $i_t \sim \text{unif}(\{1, \dots, n\})$   
 $O(1)$  to compute

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

var. minimization of convex function

$O(n)$  to compute

$\mu$   $\rightarrow$   $\frac{1}{n} \sum_{i=1}^n \mu_i$  /  $\mu$  case  
ie. sablenian



variance reduced SGD methods?

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

SAG (stochastic average gradient) [Le Roux, Schmitt & Bach NIPS 2012]

(15h30)

(Lagrange opt. prize in 2018)

- store past gradient for each  $i$  ( $g_i$ )
- update one at step  $t$

SAG { pick  $i_t \in N$  unif; update  $g_{i_t}^{(t+1)} \triangleq \nabla f_{i_t}(\alpha_t)$   
 $g_j^{(t+1)} = g_j^{(t)} \quad \forall j \neq i_t$

$$\alpha_{t+1} = \alpha_t - \gamma \frac{1}{n} \sum_{i=1}^n g_i^{(t+1)}$$

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

$\rightarrow g_{i_t}^{(t+1)}$  stable gradient

$\frac{1}{n} \sum_i g_i \approx$  an approximation of  $\nabla f(\alpha_t)$

$O(1)$  cost per iteration  
[but  $O(n)$  storage cost]

big surprise: converge linearly and fast

$\rightarrow$  vs,

"increment aggregated gradient" (IAG) [Blatt et al. 2007]  
 where you cycle deterministically  
 through  $\Sigma_1, \dots, \Sigma_n$

→ linear rate for quadratic function

"big step size"

but  $\delta_{\max} \approx O\left(\frac{1}{n}\right)$

(step size)

SAG convergence rate: thm: with  $\delta_t = \frac{1}{16L}$  where  $L = \max_i \text{Lipschitz}(Df_i)$

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

constant

$$D_{SAG} = \min\left\{\frac{1}{16K_{SAG}}, \frac{1}{8n}\right\} \text{ vs } f_{\text{grad}} \approx \frac{1}{K_{\text{grad}}}$$

example: <sup>b-reg.</sup> log regression on RCVR

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

rate comparison:

gradient method  $\left(\frac{L-\mu}{L+\mu}\right)^2 = 0.99998$

accelerated grad. (Nesterov)  $(1 - \sqrt{\frac{\mu}{L}}) = 0.99761$

Nesterov lower bound:  $\left(\frac{\sqrt{L} - \sqrt{\mu}}{\sqrt{L} + \sqrt{\mu}}\right) = 0.99048$

-vLTVM)

SAG  
(n iterations)

$$(1 - \rho_{SAG})^n = 0.88250$$

practical aspects

(see Schmidt Math. prog & paper) <sup>2016 big prize!</sup>

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

input data  
↓  
 $x_i$

instead of  $O(n \cdot d)$  storage  $\rightsquigarrow O(n)$  storage

b) initialization?  
of  $\alpha_i$ 's memory

best: run SGD for one pass, then start SAG/SAGA

c) step size?

•  $\frac{1}{4L}$

• cheap line search heuristic  
(comes from FISTA)

adaptive step size  
↓

$$\left\{ \begin{array}{l} \text{while } f_i(w_t - \frac{1}{L_i} \nabla f_i(w_t)) \geq f_i(w_t) - \frac{1}{2L_i} \|\nabla f_i(w_t)\|^2 \\ \text{set } \hat{L}_{i,\text{new}} = 2\hat{L}_{i,\text{old}} \\ \text{else } \hat{L}_{i,\text{new}} = \left(\frac{1}{2}\right)^{1/n} \hat{L}_{i,\text{old}} \end{array} \right.$$

use  $\frac{1}{\max_i L_i}$  as step size (see paper)

d) non-uniform sampling?

sample  $i \sim \frac{L_i}{\sum_j L_j}$

e) Stopping criteria!

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

f) sparse features?

$$x_{t+1} = x_t - \delta \left[ \underbrace{\nabla f_{it}(x_t)}_{\text{sparse}} - g_i^t + \underbrace{P_{\lambda} \left[ \frac{1}{n} \sum_j g_j^t \right]}_{\text{dense?}} \right] \quad (\text{sparse SAGA})$$

[Leibin & al.]  
2017

weighted projection on support of  $z_i$

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