

today : • finish BCFW  
• SAG

convergence result for BCFW

$$C_f^{(i)} \leq L_0 \text{diam}(M_i)^2 \quad C_f^{(\oplus)} \triangleq \sum_{i=1}^n C_f^{(i)}$$

$$\mathbb{E}[f(x^{(t)}) - f^*] \triangleq \epsilon_t \leq \frac{2n}{t+2n} [C_f^{(\oplus)} + f(x^{(0)}) - f^*] \quad \text{for } \epsilon_t = \frac{2n}{t+2n}$$

if you use line search:  $\epsilon_t \leq \frac{2n C_f^{(\oplus)}}{t-t_0+2n}$  for  $t \geq t_0$   
 $t_0 \leq n \log \frac{2(f(x^{(0)}) - f^*)}{\epsilon}$   
 time to ensure that  $\epsilon_t \leq \epsilon$

batch FW

$$\epsilon_t \leq \frac{2C_f}{t+2}$$

one can show that  $C_f^{(\oplus)} \leq C_f$  for quadratic functions

BCFW

$$\epsilon_t \leq \frac{2n C_f^{(\oplus)}}{t-t_0+2n}$$

$C_f$  vs.  $nC_f^{(\oplus)}$

but BCFW is often n times cheaper than batch FW

Complexity  
# of oracle calls  
FW  $O\left(\frac{nC_f}{\epsilon}\right)$

$\Rightarrow$  BCFW is "never" slower than batch FW (when not using parallelization)

extensions to BCFW:

- non-uniform sampling i.e.  $\left\{ \begin{matrix} C_f^{(i)} \\ \mathcal{O}_i(x) \end{matrix} \right\}$
  - using away-step, etc. to "get" linear convergence
- } ICML 2016

Application of BCFW on SVM duality:

- getting  $S_i$  is one loss-augmented decoding call for example  $\tilde{c}$

$$\alpha_i^{(t+1)} = \alpha_i^{(t)} + \gamma (S_i^{(t)} - \alpha_i^{(t)})$$

$$\rightarrow \text{you update } w_i^{(t+1)} = w_i^{(t)} + \gamma (W_S^{(t)} - w_i^{(t)})$$

$$w_j^{(t+1)} = w_j^{(t)} \quad \frac{1}{\lambda} \sum_i y_i |y_i^{(t)}|$$

need to store these in memory  
 or store  $\alpha_i$  (memory/computation tradeoff)

in memory  
 or store  $\alpha_i$  (memory/computation tradeoff)

note: for line search, also need to store  $b_i^T \alpha_i \triangleq R_{i, \alpha}$

convergence constants

for SVMstruct, can show that  $C_f^{(i)} \leq \frac{4R_i^2}{\lambda n^2}$   $R_i \triangleq \max_{y \in \mathcal{Y}_i} \|\psi_i(y)\|$   
 $R \triangleq \max_i R_i$

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

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

use affine invariance crucially

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

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

if use  $\ell_2$ -norm, we get a bad bound?  $\text{diam}(M)^2 = 2n$

Lipschitz constant in  $\ell_2$ -norm = largest e-value of Hessian

$\lambda \text{ATA} = \frac{1}{\lambda n^2} (\langle \psi_i(y), \psi_j(y) \rangle)_{(i,y), (j,y)}$

say eg.  $\langle \psi_i(y), \psi_j(y) \rangle \approx 1$  for lots of output

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

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

• instead, want to use  $\ell_1$ -norm of  $\Delta_{(i,y)}$

ie.  $\ell_{\infty}$ -norm for Lipschitz constant.

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

on  $M$ , use  $\ell_1$ -norm on  $\Delta_{(i,y)}$  and then max over blocks

ie.  $\|\alpha\| = \max_i \|\alpha_i\|_{\ell_2}$   $\ell_1$ - $\ell_2$  or  $\ell_2$ - $\ell_1$ ?

15h35

Variance reduced SGD

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

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

$$x \leftarrow \frac{1}{n} \sum_{i=1}^n \nabla f_i(x)$$
 where  $\nabla$  is  $\mu$ -strongly convex  
 $L$ -smooth

batch gradient method  
 [Candès 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

$$\gamma = \frac{1}{L} \quad f(x_t) - f^* \leq (1-\rho)^t \frac{f(x_0) - f^*}{\rho}$$

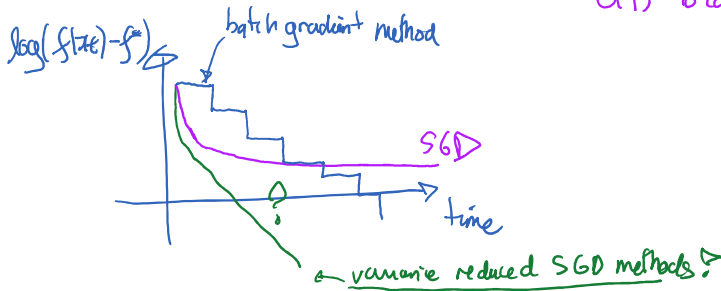
where  $\rho \approx \frac{\mu}{L} = \frac{1}{K}$   
 $K \triangleq \frac{L}{\mu}$  "condition #"  
 i.e. sublinear

stochastic gradient method  
 aka. incremental gradient method  
 [Robbins & Monro 1951]

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

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

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



SAG (stochastic average gradient)

[Le Roux, Schmitt & Bach NIPS 2012]

(Lagrange opt. prize in 2015)

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

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

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

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

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

$\tilde{O}(1)$  cost per iteration  
 [but  $O(n)$  storage cost]

big surprise: converge linearly and fast

"inherent aggregated gradient" (IAG) [Batt & al. 2007]

where you cycle deterministically through  $\{1, \dots, n\}$

→ linear rate for quadratic function

but  $\delta_{\max} \approx O\left(\frac{1}{nL}\right)$  (1/3 tiny rate)

"big step size" vs.

"big step size" vs.  $\delta_{max} \approx O(\frac{1}{\eta L})$  (step size)

SAG convergence rate:  
 f.m. with  $\delta_k = \frac{1}{kL}$  where  $L = \max_i \text{Lipschitz}(\nabla f_i)$   
 $\mathbb{E} \|f(w_t) - f^*\| \leq (1 - \min\{\frac{\mu}{L}, \frac{1}{8n}\})^t C_0$  ← constant

$\rho_{SAG} = \min\{\frac{1}{kK_{SAG}}, \frac{1}{8n}\}$  vs.  $\rho_{grad} \approx \frac{1}{K_{grad}}$

Example:  $l_2$ -reg. regression on RCV1

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

rate comparison:

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

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

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

SAG (n iterations)  $(1 - \rho_{SAG})^n \approx 0.88356$

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

a) storage: if  $f_i(w) = h(x_i^T w)$   $\Rightarrow \nabla_w f_i(w) = h'(x_i^T w) \overset{\text{scalar}}{\downarrow} x_i$   
 instead of  $O(d \cdot n)$  storage  $\rightarrow O(n)$  storage

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

c) step-size?  $\frac{1}{(4)L}$   
 • cheap line search heuristic (comes from FISTA)  $\left\{ \begin{array}{l} \text{while } f_i(w_t - \frac{1}{L} \nabla f_i(w_t)) \geq f_i(w_t) - \frac{1}{2L} \|\nabla f_i(w_t)\|^2 \\ \text{set } \tilde{L}_{i,new} = 2\tilde{L}_{i,old} \\ \text{else } \tilde{L}_{i,new} = (\frac{1}{2})^{k_n} \tilde{L}_{i,old} \end{array} \right.$   
 adaptive step size

d) non-uniform sampling? sample in  $\tilde{L}_i$

e) stopping criterion? you can use  $\frac{\sum_{j=1}^n g_j^{(k)}}{n}$  as approx. of  $\nabla f(x^k)$

f) sparse proximal?  $x_{t+1} = x_t - \gamma \left[ \underbrace{\nabla f_{i_t}(x_t) - g_{i_t}^t}_{\text{sparse}} + P_{S_t} \left[ \underbrace{\frac{1}{n} \sum_{j=1}^n g_j^t}_{\text{dense}} \right] \right]$  (sparse SAGA)

[Lehndal & al. 2017]

weighted projection on support of  $x^*$   
 $S_t = \{u : |x_{i_t}^*| > 0\}$