

today: • variance reduction perspective  
• application to CRF

Variance reduction idea

$X$  &  $Y$  are R.V.'s

goal: estimate  $\mathbb{E}X$  using M.C. samples

Suppose:  $\mathbb{E}Y$  is cheap to compute and  $Y$  is correlated with  $X$

consider estimator "control variate" technique

$\alpha \in [0, 1]$   $\Theta_\alpha \triangleq \alpha(X - Y) + \mathbb{E}Y$  to approximate  $\mathbb{E}X$   
 ↑  
 convex combo. coeff. between  $\mathbb{E}X$  &  $\mathbb{E}Y$

properties:  $\mathbb{E}\Theta_\alpha = \alpha \mathbb{E}X + (1-\alpha) \mathbb{E}Y \rightarrow$  unbiased (i.e.  $\mathbb{E}\Theta_\alpha = \mathbb{E}X$ )

if  $\alpha = 1$   
 $\mathbb{E}X = \mathbb{E}Y$  [not interesting]

Variance:  $\text{Var}(\Theta_\alpha) = \alpha^2 [\text{Var}X + \text{Var}Y - 2\text{Cov}(X, Y)]$   
 ↓  
Variance reduction?

for  $\alpha = 1$  (unbiased setting)  $\Theta_\alpha = X + (\mathbb{E}Y - Y)$   
 correction

SAG setting

$X$  is  $\nabla f_i(x_t)$ ;  $\mathbb{E}X =$  batch gradient

SAG/SAGA algorithm:  $Y$  is  $g_i$  [past stored gradient]

$\mathbb{E}Y = \frac{1}{n} \sum_{i=1}^n g_i$

SAG alg.:  $\alpha = \frac{1}{n}$  (biased)

SAGA alg.:  $\alpha = 1 \Rightarrow$  (unbiased)

SAG:	$x_{t+1} = x_t - \gamma \left[ \frac{1}{n} (\nabla f_t(x_t) - g_t^{(t)}) + \frac{1}{n} \sum_j g_j^{(t)} \right]$	(biased)
SAGA:	$x_{t+1} = x_t - \gamma \left[ \frac{1}{n} (\nabla f_{i_t}(x_t) - g_{i_t}^{(t)}) + \frac{1}{n} \sum_j g_j^{(t)} \right]$	(unbiased)
SVRG:	$x_{t+1} = x_t - \gamma \left[ \nabla f_{i_t}(x_t) - \nabla f_{i_t}(x_{\text{old}}) + \frac{1}{n} \sum_j \nabla f_j(x_{\text{old}}) \right]$	(unbiased)

SAG:  $x_{t+1} = x_t - \gamma \left[ \frac{1}{n} \sum_j g_j^i(x_t) - y_i \right] + \frac{1}{n} \sum_j g_j^i(x_t)$  (unbiased)

SVRG:  $x_{t+1} = x_t - \gamma \left[ (\nabla f_i(x_t) - \nabla f_i(x_{old})) + \frac{1}{n} \sum_j \nabla f_j(x_{old}) \right]$  (unbiased)

(stochastic variance reduced gradient)

$\rightarrow x_{old}$  is updated from outer loop (infrequently)

SVRG alg.:

for  $k=0, \dots$ , (outer loop)

compute  $g_{ref} \triangleq \frac{1}{n} \sum_j \nabla f_j(x^{(k)})$

for  $t=0, \dots, T_{max}$

sample  $i_t$

$x_{t+1} = x_t - \gamma \left[ \nabla f_{i_t}(x_t) - \nabla f_{i_t}(x^{(k)}) + g_{ref} \right]$

end

$x^{(k+1)} = x_{T_{max}}^{(k)}$

- question is:
- what is  $T_{max}$ ?
  - what is  $\gamma$ ?

original SVRG convergence result need  $\gamma \leq \frac{1}{L}$

$T_{max} \geq \approx \frac{L}{\mu} = k \rightarrow$  to run alg., need to know  $k$

$\Rightarrow$  not adaptive to local strong convexity

break of SVRG (now called "loopless")

[Hoffman & d. NIPS 2015]  $T_{max} \sim \text{Geom}(\dots)$

[at inner loop, do a batch gradient computation with prob.  $\frac{1}{n}$ ]

then, get same convergence result as SGD

comp. cost.: size of inner loop  $E[T_{max}] = n$

overall cost of SVRG  $\approx 3$  (SGD cost) for  $n$  updates

note: interpolation regime:  $\|\nabla f_i(x^*)\| = 0 \forall i$

vs. just  $\|\frac{1}{n} \sum_i \nabla f_i(x^*)\| = 0$

$\hookrightarrow$  get similar rate as SVRG/SGD with SGD with constant size

not  $\dots$  - max  $\dots$

↳ get similar rate as SVM/Ada with SGD with constant size

CRF optimization

SVM struct

$$\min_w \frac{\lambda \|w\|^2 + 1}{2} + \frac{1}{n} \sum_{i=1}^n H_i(w)$$

(note)  $\rightarrow \max_{\tilde{y}} \ell_i(\tilde{y}) - w^T \psi_i(\tilde{y})$

$$\max_{\alpha_i \in \Delta_{|\mathcal{Y}|}} \frac{\lambda \|w(\alpha)\|^2 + 1}{2} + \frac{1}{n} \sum_{i=1}^n p_i^T \alpha_i$$

CRF

$$\min_w \frac{\lambda \|w\|^2 + 1}{2} + \frac{1}{n} \sum_{i=1}^n -\log p(y^{(i)} | x^{(i)}, w)$$

$\downarrow$   
 $\log(\sum_{\tilde{y}} \exp(-w^T \psi_i(\tilde{y})))$

$$\max_{\alpha_i \in \Delta_{|\mathcal{Y}|}} \frac{\lambda \|w(\alpha)\|^2 + 1}{2} + \frac{1}{n} \sum_{i=1}^n H_i(\alpha_i)$$

(convex)  
 $\triangleq \sum_{\tilde{y}} \alpha_i(\tilde{y}) \log \alpha_i(\tilde{y})$

KKT  $\rightarrow w(\alpha) = \frac{1}{\lambda n} \sum_i \sum_{\tilde{y}} \alpha_i(\tilde{y}) \psi_i(\tilde{y})$

$$= \frac{1}{\lambda n} \sum_i \sum_{c \in \mathcal{Y}} M_{i,c}(\tilde{y}_c) \psi_{i,c}(\tilde{y}_c)$$

$\leftarrow$  from MRF

$$p(y|x; w) \propto \exp(\langle w, \phi(x,y) \rangle)$$

at optimality

$$\alpha_i^*(y) = p(y | x^{(i)}; w(\alpha^*))$$

$\alpha_i^* \in \text{interior of } \Delta_{|\mathcal{Y}|}$  (fully dense)  
unlike sparse soln for SVM struct

10:20

CRF optimization

• primal is smooth [vs. non-smooth for SVM struct]

• for a while, batch L-BFGS was method of choice [batch  $\Rightarrow$  slow for large  $n$ ]

• [Collins & d. JMLA 2002]: online exponential gradient (OEG)

block-coordinate method on dual; exponentiated gradient step on block

$$\alpha_i(\tilde{y})^{(t+1)} \propto \alpha_i(\tilde{y})^{(t)} \exp(-\gamma_t \nabla_{\alpha_i} D(\alpha^{(t)}))$$

EG alg  $\rightarrow$  proximal gradient step using  $KL(\alpha || \alpha^{(t)})$  as Bregman divergence for prox term.

$\rightarrow$  get linear convergence rate with cheap  $O(1)$  updates (like SGD) [vs.  $O(n)$  for batch methods]

[can think of it as a variance reduced method as well?]

• SAGA for CRF [Schmitt & d. AISTATS 2015]

$$w^{(t+1)} = (1 - \lambda \gamma_t) w^{(t)} - \gamma_t [ \nabla f_{i_t}(w_t) - g_i^{(t)} + \sum_j g_j^{(t)} ]$$

• SDCA (stochastic dual coordinate ascent)

$$\alpha_i^{(t+1)}(\tilde{y}) = (1 - \gamma_t) \alpha_i^{(t)}(\tilde{y}) + \gamma_t \tilde{g}_i^{(t)}(\tilde{y})$$

$\rightarrow$  related to subgradient

Stochastic away (coordinate ascent)  
 SOTA for CNF  
 [Le Prid & al. IJAI 2016]  
 thanks to efficient care search

$y_t = (1 - \alpha_t) y_t + \alpha_t \dots$   
 as relaxed fixed point update  
 $\alpha_i^* = \rho(y|x_i; w(\alpha^*)) \forall i$   
 $\triangleq \alpha_i(w)$

related to subgradient on primal

[note: BCFW is a special case of SPCA on symmetric obj.]

proximal gradient method

↳ generalization of projected gradient method to other non-smooth fct.

composite framework  $F(w) \triangleq f(w) + \Omega(w)$  where  $f$  is convex &  $L$ -smooth  
 $\Omega$  " " " but not nec. smooth

- constrained opt. :  $\Omega(w) = \delta_M(w) \triangleq \begin{cases} 0 & \text{if } w \in M \\ +\infty & \text{o.w.} \end{cases}$   
 "indicator fct." on  $M$
- $l_1$ -regularization  $\Omega(w) = \|w\|_1$

proximal gradient update:

$$w_{t+1} = \underset{w}{\operatorname{argmin}} \underbrace{f(w_t) + \langle \nabla f(w_t), w - w_t \rangle + \frac{1}{2\delta_t} \|w - w_t\|_2^2}_{B_t(w)} + \Omega(w)$$

- if  $\delta \leq \frac{1}{L}$ , then  $f(w) \leq B_\delta(w) \forall w$
- we can rewrite  $B_\delta(w) = \frac{1}{2\delta_t} \|w - [w_t - \delta_t \nabla f(w_t)]\|_2^2 + \text{const.}$   
 (by completing the square)

$\Rightarrow$  if  $\Omega(w) = \delta_M(w)$ : we get projected grad. alg.

$$w_{t+1} = \operatorname{Prox}_{\delta_t}^{\Omega} (w_t - \delta_t \nabla f(w_t))$$

↳ "proximal operator"

$$\operatorname{prox}_{\delta}^{\Omega}(z) \triangleq \underset{w}{\operatorname{argmin}} \left\{ \Omega(w) + \frac{1}{2\delta} \|w - z\|_2^2 \right\}$$

can be replaced by Fregman divergence generally. (e.g. BEG)

⊗ like projection, prox operator is non-expansive (ie. 1-Lipschitz)  
 ie.  $\|\operatorname{prox}_{\delta}^{\Omega}(w) - \operatorname{prox}_{\delta}^{\Omega}(w')\|_2 \leq \|w - w'\|_2$

(recall lecture 1/ landscape of fields)

ie.  $\| \text{prox}_\gamma^2(w) - \text{prox}_\gamma^2(w') \|_2 \leq \|w - w'\|_2$  (Lipschitz)

$\Rightarrow$  convergence rate for prox grad. method on  $F = f + \lambda Q$  are same as unconstrained gradient descent on  $f$

\* to be useful, need  $\text{prox}_\gamma^2$  to be efficiently computable

$$\text{prox}_\gamma^2(z) = \underset{w}{\text{argmin}} \quad \|w\|_1 + \frac{1}{2\gamma} \|w - z\|_2^2$$

$$\text{"soft-thresholding"} \Rightarrow \begin{cases} \text{sgn}(z_i) [z_i - \gamma] & \text{if } |z_i| \geq \gamma \\ 0 & \text{o.w.} \end{cases}$$

(component-wise)

used eg for Lasso:  $Q_1$ -regularized least-squares

FISTA  $\rightarrow$  accelerated prox gradient method

$\hookrightarrow$  SOTA for batch Lasso

\* split-learning  $\rightarrow$  use (prox) SAGA for Lasso &  $Q_1$ -reg. log regression

prox SAGA  $w_{t+1} = \text{prox}_\gamma^2 \left( w_t - \gamma \left[ \nabla f_t(w_t) - g_{i_t}^{(t)} + \sum_{j \neq i_t} g_j^{(t)} \right] \right)$

could accelerate prox SAGA using "catalyst" (see next class)