

## Lecture 20 - variance reduction

Thursday, March 30, 2023 1:29 PM

- today:
  - variance reduction perspective
  - application CRF

Variance reduction idea:

$X \{ Y$  are A.U.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:

$$\alpha \in [0,1] \quad \hat{O}_\alpha \stackrel{\Delta}{=} \alpha(X-Y) + \mathbb{E}Y \text{ to approximate } \mathbb{E}X$$

1 convex comb. coeff. between  $\mathbb{E}X \{ \mathbb{E}Y$

properties:  $\mathbb{E}\hat{O}_\alpha = \alpha \mathbb{E}X + (1-\alpha) \mathbb{E}Y \rightarrow \text{unbiased} \text{ [i.e. } \mathbb{E}\hat{O}_\alpha = \mathbb{E}X]$

$$\begin{cases} \text{if } \alpha=1 \quad \mathbb{E}X=\mathbb{E}Y \text{ [not interesting]} \\ \alpha=0 \end{cases}$$

$$\text{Variance: } \text{Var}(\hat{O}_\alpha) = \alpha^2 \left[ \text{Var}X + \text{Var}Y - 2\text{cov}(X,Y) \right]$$

Variance reduction?

$$\text{for } \alpha=1 \text{ (unbiased setting)} \quad \hat{O}_\alpha \approx X + \underbrace{(\mathbb{E}Y - Y)}_{\text{correction}}$$

SGD setting:

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

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

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

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

SAGA alg.:  $\alpha=1 \Rightarrow \text{unbiased}$

$$\hat{O}_\alpha = \underbrace{\alpha(X-Y) + \mathbb{E}Y}_{\text{SGD}} + \underbrace{\alpha(Y-g_i) + \mathbb{E}g_i}_{\text{SAGA}} + \underbrace{\alpha(g_i - \mathbb{E}g_i)}_{\text{SAGA correction}}$$

$$\text{SAG: } x_{t+1} = x_t - \gamma \left[ \frac{1}{n} \left[ \nabla_{i_t} f(x_t) - g_{i_t}^{(t)} \right] + \frac{1}{n} \sum_{i=1}^n g_i^{(t)} \right] \quad (\text{biased})$$

$$\text{SAGA: } x_{t+1} = x_t - \gamma \left[ \frac{1}{n} \left[ \cdot - \cdot \right] + \cdot \right] \quad (\text{unbiased})$$

$$\text{SVRG: } x_{t+1} = x_t - \gamma \left[ \frac{1}{n} \left[ \nabla_{i_t} f(x_t) - \nabla_{i_t} f(x_{\text{all}}) \right] + \frac{1}{n} \sum_{i=1}^n \nabla_i f(x_{\text{all}}) \right] \quad (\text{unbiased})$$

(stochastic variance reduced gradient)

$\hookrightarrow x_{\text{all}}$  is updated from outer loop

SVRG algorithm:

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

$$\text{compute } g_{\text{ref}} \triangleq \frac{1}{n} \sum_{i=1}^n \nabla_i f(x^{(k)})$$

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

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

end

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

end

questions:

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

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

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

$\Rightarrow$  not adaptive to local strong convexity

fixes of SVRG (now called "loopless")

[Hoffmann & al. NeurIPS 2015]  $T_{\max} \sim \text{Geom}(-)$

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

then get same convergence as SAGA

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

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

SAG / SAGA / logless SVAG, convergence  
for convex f. ( $\mu=0$ ) get  $\min_{t \in T} \{ \mathbb{E} f(x_t) - f^* \} = O\left(\frac{1}{T}\right)$

[contrast this with  $O\left(\frac{1}{T}\right)$   
for SGD]

⊕ note: interpretation requires:  $\| \nabla f_i(x^*) \| = 0 \quad \forall i$

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

↳ get similar rate as SVRG / SAGA with SGD

### CRF optimization

$\text{SVM struct}$	<u>primal</u> $\min_w \frac{\lambda \ w\ ^2}{2} + \frac{1}{n} \sum_{i=1}^n H_i(w)$	$\max_{\vec{y}} L_i(\vec{y}) - w^T \vec{\psi}_i(\vec{y})$	<u>dual</u> $\max_{\alpha_i \in \Delta_{\mathcal{Y}_i}} -\frac{\lambda}{2} \ w(\alpha)\ ^2 + \frac{1}{n} \sum_{i=1}^n \alpha_i^T \vec{\psi}_i$
$\text{CRF}$	$\min_w \frac{\lambda \ w\ ^2}{2} + \frac{1}{n} \sum_{i=1}^n -\log p(y^{(i)}   x^{(i)}, w)$	$\max_{\alpha_i \in \Delta_{\mathcal{Y}_i}} -\frac{\lambda}{2} \ w(\alpha)\ ^2 + \frac{1}{n} \sum_{i=1}^n H_i(\alpha_i)$	$\stackrel{\text{entropy}}{=} - \sum_{\vec{y}} \exp(-w^T \vec{\psi}(\vec{y})) \log \exp(-w^T \vec{\psi}(\vec{y}))$

$$\text{KKT} \rightarrow w(\alpha) = \frac{1}{n} \sum_{i=1}^n \sum_{\vec{y}} \alpha_i \vec{\psi}_i(\vec{y}) \vec{\psi}_i(\vec{y})$$

$$= \frac{1}{n} \sum_{i=1}^n \sum_{c \in \mathcal{C}_i} \sum_{\vec{y}_c} \alpha_i c(\vec{y}_c) \vec{\psi}_{i,c}(\vec{y}_c)$$

from MRF

$$p(\vec{y} | x; w) \propto \exp(w^T \vec{\psi}(x, \vec{y}))$$

⊕ duality

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

$\alpha_i^* \in \text{interior of } \Delta_{\mathcal{Y}_i}$

unlike sparse solution in SVM struct

14h33

### CRF optimization

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

- for a while, batch L-BFGS was method of choice [batch  $\Rightarrow$  slow range n]
- [Fallouts & al. JMLR 2008] = online exponentiated gradient (OEG)

block-coordinate method on dual; exponentiated gradient step on block  
 $\alpha_i^{(t+1)} \text{ or } \alpha_i^{(t+1)} \exp(-\gamma_t D_{\alpha_i} D(\alpha^{(t)})$  dual adj.

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

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

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

- SAGA for CRF  
[Schmidt & al. AISTATS 2015]

$$w^{(t+1)} = (1-\gamma_t) w^{(t)} - \gamma_t [Df_i(w^{(t)}) - g_i^{(t)} + \sum_j g_j^{(t)}]$$

- SDCA (stochastic dual coordinate ascent)  
 $\xrightarrow{\text{SOTA}}$  for CRF  
[Le Pardal & al. UAI 2018]  
thanks line search

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

$\in [0, 1] \rightsquigarrow \text{stabilize}$

as a relaxed fixed pt. iteration

$$\alpha_i(\tilde{y}) = p(y | x_i; w^{(t)}) \alpha_i$$

$p(y | x_i; w^{(t)}) \rightsquigarrow \alpha_i(w)$

related to subgradient form

[note: PFW is a special case of SDCA on SVM obj.]

### proximal gradient method

$\hookrightarrow$  generalization of projected gradient method to other non-smooth ft.

composite framework  $F(w) \triangleq f(w) + \mathcal{L}(w)$  where  $f$  is convex  $L$ -smooth

- constrained opt.  $\mathcal{L}_M(w) = \mathcal{S}_M(w) \triangleq \begin{cases} 0 & \text{if } w \in M \\ +\infty & \text{o.w.} \end{cases}$  but not necessarily smooth  
"Indicator fn." on M

- $\ell_1$ -regularization  $\mathcal{L}(w) = \|w\|_1$

### proximal gradient update:

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

$\stackrel{\triangle}{=} B_t(w)$

- if  $\gamma \leq \frac{1}{L}$ , then  $f(w) \leq B_t(w) \quad \forall w$

• we can rewrite  $B_t(w) = \frac{1}{2\gamma t} \|w - [w_t - \gamma t \nabla f(w_t)]\|_2^2 + c$ ,

$\Rightarrow$  if  $\Omega(w) = \delta_M(w)$ ; we get (by completing square)  
projected gradient alg.

$$w_{t+1} = \operatorname{prox}_{\gamma t}^{-2}(w_t - \gamma t \nabla f(w_t))$$

↳ "proximal operator"

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

could replace by  
proximal divergence to get other  
generalization (e.g. OEG)

⊕ like projection, prox operator is non-expansive (i.e. 1-Lipschitz)

$$\text{i.e. } \|\operatorname{prox}(\omega) - \operatorname{prox}(\omega')\|_2 \leq \|\omega - \omega'\|_2 \quad (\text{recall Lecture 11 landscape of rates})$$

$\Rightarrow$  convergence rate for prox gradient method or  $F = f + \Omega$   
are same as unconstrained gradient descent on  $f$

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

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

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

Used e.g. for lasso:  $l_1$ -reg. least square

FISTA  $\rightarrow$  accelerated prox-gradient method

↳ SOTA for batch lasso

$$\min_w \frac{1}{n} \sum_i \ell_i(w) + \Omega(w)$$

\* scikit-learn  $\rightarrow$  use (prox) SAGA for loss +  $l_1$ -reg. log reg.

$$\operatorname{prox}_{\gamma} \text{SAGA} \quad w_{t+1} = \operatorname{prox}_{\gamma \delta}^{-2} \left[ w_t - \gamma t \left[ \nabla f_t(w_t) - g_t^{(4)} + \left[ \sum_j \xi_j^{(4)} \right] \right] \right]$$

prox SAGA

$$w_{t+1} = \text{prox}_{\frac{\Omega}{n}}^{\Omega} \left[ w_t - \gamma \epsilon \left[ \nabla f_t(w_t) - g_t^{(t)} + \left[ \sum_j \xi_j g_j^{(t)} \right] \right] \right]$$

could accelerate proxSAGA using "catalyst" [see next class]