

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

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

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

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

If $\mathbb{E}X = \mathbb{E}Y$ [not interesting]

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

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

SGD setting:

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

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)

	$\alpha(X - Y)$	$+ \mathbb{E}Y$
SAG:	$x_{t+1} = x_t - \gamma \left[\frac{1}{n} [Df_i(x_t) - g_i^{(t)}] + \frac{1}{n} \sum_j g_j^{(t)} \right]$	(biased)
SAGA:	$x_{t+1} = x_t - \gamma \left[\frac{1}{n} [Df_i(x_t) - g_i^{(t)}] + \frac{1}{n} \sum_j g_j^{(t)} \right]$ <small>$(\alpha=1)$</small>	(unbiased)
SVRG, (stochastic variance reduced)	$x_{t+1} = x_t - \gamma \left[Df_i(x_t) - Df_i(x_{\text{ref}}) + \frac{1}{n} \sum_j Df_j(x_{\text{ref}}) \right]$	(unbiased)

\triangleright SVRG, $x_{t+1} = x_t - \gamma [Vf_t(x_t) - Vf_i(x_{\text{old}}) + \frac{1}{n} \sum_j Vf_j(x_{\text{old}})]$ (unbiased)
 (stochastic
variance reduced
gradient)

x_{old} is updated from outer loop

SVRG alg.:

```

for k=0, ... (outer loop)
  compute  $g_{\text{ref}} \triangleq \frac{1}{n} \sum_j \nabla f_j(x^{(k)})$ 
  for t=0, ..., Tmax
    sample  $i_t$ 
     $x_{t+1}^{(k)} = x_t^{(k)} - \gamma [\nabla f_{i_t}(x_t^{(k)}) - \nabla f_{i_t}(x^{(k)}) + g_{\text{ref}}]$ 
  end
   $x^{(k+1)} = x_{T_{\max}}^{(k)}$ 

```

- questions:
- what is T_{\max} ?
 - what is γ ?

original SVRG convergence result: need $\gamma \leq \frac{1}{L}$
 $T_{\max} \geq \pi \frac{L}{\mu} = K \rightarrow$ to run alg., need to know K
 \Rightarrow not adaptive to local strong convexity

tweak of SVRG (now called "Dopplergang")
 [Hoffmann & al. NeurIPS 2015] $T_{\max} \sim \text{Geom}(\dots)$

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

then, get same convergence result as SAGA

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

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

SAG / SAGA / Dopplergang, convergence
 for convex fct. ($\mu=0$)
 \triangleright get $\min_t \{f(x_t) - f^*\} = O\left(\frac{1}{T}\right)$
 [contrast this with $O\left(\frac{1}{T^2}\right)$ for SGD]

* Note: interpolation regime: $\|\nabla f_i(x^*)\| = O(\sqrt{\epsilon})$

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

\hookrightarrow get similar rate as SVRG / SAGA with SGD

CRF optimization

SVM struct $\min_w \frac{1}{2} \ w\ ^2 + \frac{1}{n} \sum_{i=1}^n H_i(w)$	$\max_{\tilde{y}} \tilde{y}^\top H(\tilde{y}) - w^\top H(\tilde{y})$	$\max_{\alpha} \frac{-\lambda \ w(\alpha)\ ^2}{2} + \frac{1}{n} \sum_{i=1}^n \alpha_i^\top \alpha_i$
<u>CRF</u> $\min_w \frac{1}{2} \ w\ ^2 + \frac{1}{n} \sum_{i=1}^n -\log p(y^{(i)} x^{(i)}, w)$	$\log \left(\sum_{\tilde{y}} \exp(-w^\top H(\tilde{y})) \right)$	$\max_{\alpha \in \Delta_{\tilde{Y}}} -\frac{\lambda}{2} \ w(\alpha)\ ^2 + \frac{1}{n} \sum_{i=1}^n H_i(\alpha_i)$ $\equiv \sum_{\tilde{y}} \alpha_i^\top y_i / \log \alpha_i^\top y_i$

$$KKT \rightarrow w(\alpha) = \left(\sum_{\tilde{y}} \sum_{i \in \tilde{y}} \alpha_i (\tilde{y}) \psi_i(\tilde{y}) \right)$$

$$= \frac{1}{n} \sum_i \sum_{y \in C(x_i)} \alpha_i (y) \psi_i (y)$$

↳ from MRF

$$p(y|x; w) \propto \exp(-w^\top \psi(x, y))$$

at optimality

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

$\alpha_i^* \in \text{interior of } \Delta_{\tilde{Y}}$
unlike sparse soln in SVM struct

15h27

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 & al. JMLR 2008] : online exponentiated gradient (OEG)

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

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

EG alg \rightarrow proximal gradient step using $KL(\alpha || \alpha^{(t)})$ as a 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?]

$$w^{(t+1)} = (1 - \gamma_t) w^{(t)} - \gamma_t \left[\nabla f_{i_t}(w^{(t)}) - g_i^{(t)} + \frac{1}{n} \sum_j g_j^{(t)} \right]$$

$E[\alpha]$ (stabilize)

$$\alpha_{i_t}^{(t+1)}(y) = (1 - \gamma_t) \alpha_{i_t}^{(t)}(y) + \gamma_t \underbrace{s_y^{(t)}(y)}_{p(y|x^{(i)}, w^{(t)})}$$

as related fixed pt.

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

- SDCA
(stochastic dual coordinate ascent)

(coordinate descent)

SOTA for CRF

[LePoutre et al. UAI 2018]

thanks to line search

as related fixed pt.

update

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

$$p(y|x_i; w(\alpha^*))$$

$$\triangleq \alpha_i / w$$

related
to subgradient
point

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

proximal gradient method

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

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

• constrained opt. $\Omega(w) = \delta_M(w) \triangleq \begin{cases} 0 & \text{if } w \in M \\ +\infty & \text{otherwise} \end{cases}$
"indicator fn." on M

• ℓ_1 -regularization $\Omega(w) = \|w\|_1$

proximal gradient update:

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

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

• we can rewrite $B_t(w) = \frac{1}{2\gamma} \|w - [w_t - \gamma \nabla f(w_t)]\|_2^2 + \text{cst.}$
(by completing the square)

\Rightarrow if $\Omega(w) = \delta_M(w)$; we get
projected grad., resp.

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

↳ "proximal operator"

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

could be replaced by
Bregman divergence to get other
generaliz. (e.g. OEG)

⊗ Like projection, prox operator is non-expansive (i.e. 1-Lipschitz)

$$\text{i.e. } \|\operatorname{prox}_\gamma^{\Omega}(w) - \operatorname{prox}_\gamma^{\Omega}(w')\|_1 \leq \|w - w'\|_1$$

(recall Lecture 11
landscape of
rotors)

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

$$\text{i.e. } \|\text{prox}_g^\Omega(w) - \text{prox}_g^\Omega(w')\|_2 \leq \|w - w'\|_2$$

(recall lecture 1
landscape of
rates)

\Rightarrow convergence rate for prox. gradient method on $F = f + \frac{\lambda}{2}\|w\|^2$
are same as unconstrained gradient descent on f

* to be useful, need prox_g^Ω to be efficiently computable

$$\text{prox}_\gamma^{\| \cdot \|_1}(z) = \arg \min_w \|w\|_1 + \frac{1}{2\gamma} \|w - z\|^2$$

$$\begin{aligned} \text{"soft-thresholding"} &\stackrel{\Delta}{=} \begin{cases} \text{sgn}(z)\left[\|z\|_1 - \gamma\right] & \text{if } \|z\|_1 \geq \gamma \\ 0 & \text{o.w.} \end{cases} \\ &\text{(component-wise)} \end{aligned}$$

Used e.g. for Lasso: ℓ_1 -regularized least square

FISTA \rightarrow accelerated prox gradient method

\hookrightarrow SOTA for batch lasso

$$\min_w \sum_i f_i(w) + \Omega(w)$$

* Scikit-learn \rightarrow use tSAGA for lasso ; ℓ_1 -reg., log-reg.

$$\text{prox SAGA} \quad w_{t+1} = \underset{\text{(prox)}}{\text{Prox}}_\gamma^\Omega \left(w_t - \gamma [Df_t(w_t) - g_t^{(t)} + \frac{1}{\eta} \sum_j g_j^{(t)}] \right)$$

Could accelerate proxSAGA using "catalyst"
(see next class)