

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
 $\alpha \in [0, 1]$

$\Theta_\alpha \triangleq \alpha(X - Y) + \mathbb{E}Y$ to approximate $\mathbb{E}X$

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

SGD setting:

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

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

$\mathbb{E}Y = \frac{1}{n} \sum_i g_i$

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

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

SAG: $x_{t+1} = x_t - \gamma \left[\frac{1}{n} [\nabla f_i(x_t) - g_{i_t}] + \frac{1}{n} \sum_j g_j \right]$ (biased)

SAGA: $x_{t+1} = x_t - \gamma \left[\nabla f_i(x_t) - g_{i_t} + \frac{1}{n} \sum_j g_j \right]$ (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)

SVRG: $x_{t+1} = x_t - \gamma \left[\nabla f_{i_t}(x_t) - \nabla f_{i_t}(x_{old}) + \frac{1}{n} \sum_{j=1}^n \nabla f_j(x_{old}) \right]$ (unbiased)
 (stochastic variance reduced gradient)

↘ x_{old} is updated from outer loop

SVRG algorithm

```

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

Questions:

- what is T_{max} ?
- what is γ ?

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

$T_{max} \geq \frac{L}{\mu} = K \rightarrow$ to run alg., need to know K ;
 \Rightarrow not adaptive to local strong convexity

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

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

then, get same convergence result as SAGA
 \hookrightarrow size of inner loop $\mathbb{E}[T_{max}] = n$; overall cost of SVR $\approx 3 \cdot (\text{SGD for } n \text{ updates})$

SAG/SAGA, loopless SVRG, convergence for convex set ($\mu=0$) get $\min \{ \mathbb{E}[f(x_t)] - f^* \} = O\left(\frac{1}{t}\right)$
 [contrast with $1/\sqrt{t}$ for SGD]

15h34

CRF optimization

$\max_{\theta} \ell(\theta) - w^T \psi(\theta)$

CRF optimization

CRF objective :	primal	dual
SVM struct	$\min_w \frac{\lambda \ w\ ^2}{2} + \frac{1}{n} \sum_{i=1}^n H_i(w)$	$\max_{\alpha_i \in \Delta_{ S_i }} -\lambda \ w(\alpha)\ ^2 + \frac{1}{n} \sum_{i=1}^n p_i^T \alpha_i$
CRF	$\min_w \frac{\lambda \ w\ ^2}{2} + \frac{1}{n} \sum_{i=1}^n -\log p(y^{(i)} x^{(i)}; w)$ <p style="text-align: center;"> $\log \left(\sum_{\tilde{y}} \exp(-w^T \psi_i(\tilde{y})) \right)$ </p>	$\max_{\alpha_i \in \Delta_{ S_i }} -\lambda \ w(\alpha)\ ^2 + \frac{1}{n} \sum_{i=1}^n H_i(\alpha_i)$ <p style="text-align: center;"> $\hat{=} \sum_{\tilde{y}} \alpha_i(\tilde{y}) \log \alpha_i(\tilde{y})$ </p>

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

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

\downarrow
 from MRF

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

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

$\alpha_i^* \in \text{interior of } \Delta_{|S_i|}$

unlike sparse solution in structured SVM

CRF optimization:

- primal is smooth & strongly convex [vs. non-smooth for SVM struct]
- for a while, batch L-BFGS was method of choice [batch \Rightarrow slow for large n]
- [Collins & d. JMLR 2005]: online exponentiated gradient (OEG)

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

$\alpha_i(\tilde{y})^{(t+1)} \propto \alpha_i(\tilde{y})^{(t)} \exp(-\gamma_t \nabla_{\alpha_i} \nabla^2 \mathcal{L}^{(t)})$

EG alg \rightarrow proximal gradient step using $KL(\alpha/\alpha_t)$ as Bregman divergence

\rightarrow get linear conv. rate with cheap $O(1)$ updates (like SGD)

[vs. $O(n)$ for batch method]

[can think of it as variance reduced method as well.]

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

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

$E[0,1]$ (stabilize)
 \downarrow
 $\sim t^{-1}$

2015]

- SDCA
stochastic dual
coord. ascent

$$\alpha_{i_t}^{(t)}(\tilde{y}) = (1-\alpha_t)\alpha_{i_t}^{(t-1)}(\tilde{y}) + \alpha_t \underbrace{\tilde{S}_{i_t}^{(t)}(\tilde{y})}_{p(\tilde{y} | x_{i_t}; w(\alpha^{(t)}))} \in [0,1] \text{ (stabilize)}$$

$\Rightarrow \alpha(w)$

as a related fixed pt. update
 $\alpha_i^* = p(y | x_i; w(\alpha^*)) \forall i$

related to subgradient in primal

[note: BCFW is a special case of SDCA in S/M shift alg. J.]

→ SOTA for CRF [LePrid et al. IJCV 2018] (thanks to LSJ)

proximal gradient method

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

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

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

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

proximal gradient update:

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

• if $\alpha_t \leq \frac{1}{L}$ then $f(w) \leq B_{\alpha_t}(w) \forall w$

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

\Rightarrow if $\Omega(w) = S_M(w)$; we get the projected gradient alg.

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

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

↳ 'proximal operator'

$$\text{prox}_{\delta}^{\Omega}(z) \triangleq \arg \min_w \left\{ \Omega(w) + \frac{1}{2\delta} \|w - z\|_2^2 \right\}$$

could replace with
Bregman divergence to get other
gener. (e.g. OEG)

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

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

⇒ convergence rate of prox. grad. method on $F = f + \Omega$
one same as unconst. gradient descent on f