

Lecture 20 - variance reduction

Thursday, March 21, 2019 13:27

today:
• variance reduction perspective
• application CRF

Variance reduction idea

$X \& Y$ be R.V.

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 $\Theta_\alpha \triangleq \alpha(X-Y) + \mathbb{E}Y$ to approximate $\mathbb{E}X$
 $\alpha \in [0, 1]$

Properties: $\mathbb{E}\Theta_\alpha = \alpha \mathbb{E}X + (1-\alpha) \mathbb{E}Y \rightarrow$ unbiased (i.e. $\mathbb{E}\Theta_\alpha = \mathbb{E}X$)
 $\downarrow \mathbb{E}Y = \mathbb{E}X$ (not interesting)

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

\downarrow Variance reduction

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

SGD setting:

11

SGD setting :

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

SAG/SAGA algorithms : Y is g_i^t [past stored gradient]

$$\mathbb{E}Y = \frac{1}{n} \sum g_i^t$$

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

SAGA II : $\alpha = 1$ (unbiased)

$$\text{SAG: } x_{t+1} = x_t - \gamma \left[\nabla f_i(x_t) - g_{it}^t \right] + \left[\frac{1}{n} \sum g_j^t \right] \quad (\text{biased})$$

$$\text{SAGA: } x_{t+1} = x_t - \gamma \left[\nabla f_i(x_t) - g_{it}^t + \frac{1}{n} \sum g_j^t \right] \quad (\text{unbiased})$$

$$\text{SVRG: } x_{t+1} = x_t - \gamma \left[\nabla f_i(x_t) - \nabla f_i(x_{\text{old}}) + \sum_{j=1}^n \nabla f_j(x_{\text{old}}) \right] \quad (\text{unbiased})$$

(stochastic variance reduced gradient)

x_{old} is updated from outer loop

SVRG algorithm :

SVRG algorithm :

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

questions:

- what is T_{\max} ?
- what is γ ?

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

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

fixes on SVRG:

[Hoffmann et al. NIPS 2015]

$T_{\max} \sim \text{Geom}(\sim)$

[of inner loop iterations, do batch gradient comp.]

with prob. $\frac{1}{n}$

then, get same convergence result as SAGA

↳ here, size of inner loop $\mathbb{E}[t_{\text{max}}] = n$

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

SAG / SAGA / SURG_{Hoffmann}, rate for convex fct. ($\mu=0$), get $\min_t [\mathbb{E}f(x_t) - f^*] = O\left(\frac{1}{t}\right)$

[$O(n^3)$]

[contrast with $\frac{1}{\sqrt{t}}$ for SGD]

CRF objective:

| | | |
|---|--|--|
| <u>primal</u> SVM struct $\min_w \frac{\ w\ ^2}{2} + \frac{1}{n} \sum_{i=1}^n H_i(w)$ | $\max_{\tilde{y} \gg} l_i(\tilde{y}) - w^T \psi(\tilde{y})$ | <u>dual</u> $\max_{\alpha_i \in \Delta(\mathcal{X}_i)} -\frac{\ w(\alpha)\ ^2}{2} + \frac{1}{n} \sum_{i=1}^n \alpha_i^T \alpha_i$ |
| <u>CRF</u> : $\min_w \frac{\ w\ ^2}{2} + \frac{1}{n} \sum_i -\log p(y^{(i)} x^{(i)}, w)$ | $\log \left(\prod_{\tilde{y}} \exp(-w^T \psi(\tilde{y})) \right)$ | $\max_{\alpha_i \in \Delta(\mathcal{X}_i)} -\frac{\ w(\alpha)\ ^2}{2} + \frac{1}{n} \sum_{i=1}^n H_i(\alpha_i)$ $\triangleq -\sum_{\tilde{y}} \alpha_i(\tilde{y}) \log \alpha_i(\tilde{y})$ |

$$\begin{aligned} \text{KKT} \rightarrow w(\alpha) &= \frac{1}{n} \sum_i \sum_{\tilde{y}} \alpha_i(\tilde{y}) \psi_i(\tilde{y}) \\ &= \frac{1}{n} \sum_i \sum_{c \in \mathcal{C}} \alpha_{i,c}(\tilde{y}_c) \psi_{i,c}(\tilde{y}_c) \end{aligned}$$

* at optimality

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

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

from MRF

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

$$\downarrow$$

$$x_i(y) \rightarrow p(y|x^{(i)}, w^{(i)})$$

$\alpha_i^* \in \text{interior of } \Delta_{\mathcal{C}(i)}$

unlike sparse solution in structured SVM

CRF optimization

- primal objective 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 et al. JMLR 2008] : online exponentiated gradient (OEG)

block-coordinate method on dual ; exponentiated gradient $\alpha_i^{(t+1)}$ stop on block

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

dual fct.

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

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

[can think of it as variance reduced SGD as well]

SAGA for CRF:

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

$\dots \dots \dots \dots$

Schmidt et al.

\rightarrow running LHS
 $w = \text{argmin}_w \text{L}(w)$
 [Schmidt et al.,
 AISTATS 2015]

SDCA
 stochastic dual
 coord. ascent.

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

SOA
 for CRF [Le et al.
 UAI 2015]

$$\alpha_{i,t}^{(t+1)} = (1-\gamma_t)\alpha_i^{(\tilde{y})^{(t)}} + \gamma_t \underbrace{\tilde{S}_i^{(\tilde{y})^{(t)}}}_{p(\tilde{y} | x_i, w(\alpha^{(t)}))}$$

as a relaxed fixed point update for $\alpha_i^* = p(y|x_i, w(\alpha^*))$

[margin inference]

$\in \mathcal{L}(\mathcal{X})$ (stabilize)

Proximal gradient method

↳ generalization of projected gradient method to other non-smooth functions

composite framework: $F(w) \triangleq f(w) + \Omega(w)$ where f is convex & L -smooth

- constrained opt. $\Omega(w) = \delta_M(w) \triangleq \begin{cases} 0 & \text{if } w \in M \\ +\infty & \text{o.w.} \end{cases}$ Ω is convex but not nec. smooth
 "indicator of M "

- L_1 -regularization $\Omega(w) = \|w\|_1$

proximal gradient update:

$$w_{t+1} = \text{prox}_{\Omega} \left(\text{grad} f(w_t) + \nabla F(w_t), w_t \right)$$

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

$$\triangleq B_t(w)$$

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

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

\rightarrow if $Q(w) = S_M(w)$, we get the projected gradient alg.

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

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

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

$$\text{i.e. } \|\operatorname{prox}_\gamma^Q(w) - \operatorname{prox}_\gamma^Q(w')\|_2 \leq \|w - w'\|_2$$

\Rightarrow convergence rate of prox. gradient method
are same as unconstrained gradient descent

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