

Lecture 19 - SAG

Thursday, March 25, 2021 13:34

today: • finish BCFW
• SAG

application of BCFW to SVM struct:

• getting s_i is one less-augmented decoding call for example \bar{i}

$$\alpha_i^{(t+1)} = \alpha_i^{(t)} + \gamma_t (S_i^{(t)} - \alpha_i^{(t)})$$

$$\Rightarrow \text{you update } w_i^{(t+1)} = w_i^{(t)} + \gamma_t (w_s^{(t)} - w_i^{(t)})$$

$$w_j^{(t+1)} = w_j^{(t)} \quad t_j \neq i$$

$w = A\alpha = \sum_i A_i \alpha_i \triangleq w_i$

(worst case) \rightarrow dnd) need to store those in memory or store α_i 's

$\downarrow \gamma_i(y_i^{(t)})$
 $\uparrow \gamma_i(\hat{y}_i^{(t)})$

(memory / computation tradeoff)

note: for line search, also need to store $b_i^T \alpha_i \triangleq r_i^{(t)}$

Convergence Constants

for SVM struct, can show that $C_F^{(i)} \leq \frac{4R_i^2}{\lambda n^2}$

$R_i \triangleq \max_{y \in \mathcal{Y}} \|\gamma_i(y)\|_2$
 $R = \max_i R_i$

Hessian: $(\lambda A^T A)_{(i,y), (j,y)}$
 $= \lambda \frac{1}{n^2} \langle \gamma_i(y), \gamma_j(y) \rangle$
 $d_i^T H_i d_i \leq \frac{1}{\lambda n^2} \max_{y \in \mathcal{Y}} \langle \gamma_i(y), \gamma_i(y) \rangle \|d_i\|_2^2$
 R_i^2

use affine invariance crucially
 \downarrow
 compare $C_F \leq \frac{4R^2}{\lambda n}$

$\Rightarrow C_F = \sum_{i=1}^n C_F^{(i)} \leq \frac{4R^2}{\lambda n} \approx \frac{C_F}{n}$

ie. BCFW is "n times" faster than batch FW for SVM struct?

importance of affine invariance:

$C_F \leq L \| \cdot \| \cdot \text{diam}_{p,1}(M)^2$

boo

.2

if you use l_2 -norm, we get a bound? $\text{diam}(M) = 2n$

Lipschitz constant in l_2 -norm = largest e-value of Hessian

recall Hessian $\lambda A^T A = \frac{1}{n} \sum_{i=1}^n \langle \nabla_i f(y), \nabla_i f(y') \rangle_{(y,y), (y',y')}$

say eg. $\langle \nabla_i f(y), \nabla_i f(y') \rangle \approx 1$ for bits of output

$$(11^T) \mathbf{1} = \text{dim} \cdot \mathbf{1}$$

→ get largest e-value can scale with dim of matrix ⇒ really bad here because dim is exponential

• instead, want to use l_1 -norm of $\Delta f_{(S)}$

i.e. l_{∞} -norm for Lipschitz constant

$$\text{then get } L_i (\text{diam}(\Delta f_{(S)})) \approx \frac{4R^2}{\lambda}$$

on M , use l_1 -norm on $\Delta f_{(S)}$ and then max over blocks

$$\text{i.e. } \alpha = \max_i \|g_i\|_1 \quad l_1\text{-}l_{\infty} \text{ or } l_{\infty}\text{-}l_1 ?$$

Variance reduced SGD

setup: $\min_x \frac{1}{n} \sum_{i=1}^n f_i(x)$
 $\stackrel{\Delta}{=} f(x)$

where f is μ -strongly convex

L -smooth

$$f(x_t) - f(x^*) \leq (1-\rho)^t (f(x_0) - f^*)$$

batch gradient method
 [Cauchy 18th century]

$$x_{t+1} = x_t - \gamma \frac{1}{n} \sum_{i=1}^n \nabla f_i(x_t)$$

$O(n)$ to compute

$$\gamma = \frac{1}{L}$$

where $\rho \approx \frac{\mu}{L} = \frac{1}{\kappa}$

$$\kappa \stackrel{\Delta}{=} \frac{L}{\mu} \text{ "condition" \#}$$

Stochastic gradient method
 a.k.a. incremental gradient method
 [Robbins & Monro 1951]

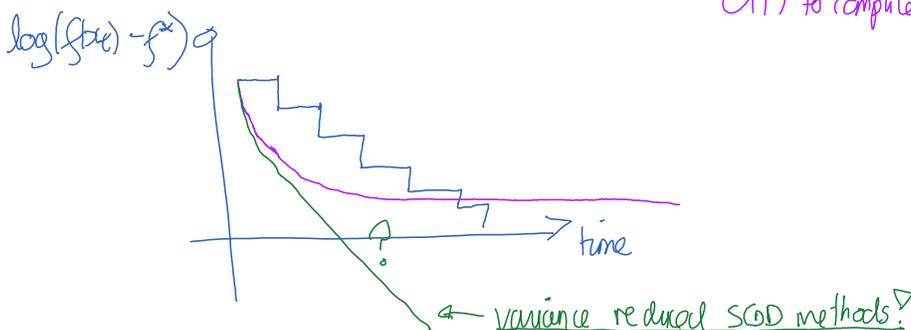
$$x_{t+1} = x_t - \gamma_t \nabla f_{i_t}(x_t)$$

where $i_t \sim \text{Unit}(\{1, \dots, n\})$

$O(1)$ to compute

$\gamma_t = \text{const} \cdot \frac{1}{\sqrt{t}}$
 ⇒ linear rate up to a ball of radius γ

$\gamma_t \approx \frac{1}{\sqrt{t}}$
 ⇒ $\tilde{O}(\frac{1}{\sqrt{t}})$ rate
 i.e. sublinear



time

← variance reduced SGD methods

SAG, SAGA, SVRG, SDCA, OEG, etc...

14h22

SAG (stochastic average gradient) [Le Roux, Schmidt & Bach NIPS 2012]

(Pontryagin opt. prize 2015)

SAG: • store past gradient for each i (g_i)
• update one at step t

SAG { pick i_t unif. & update $g_{i_t}^{(t+1)} = \nabla f_{i_t}(x_t)$
 $g_j^{(t+1)} = g_j^{(t)}$ for $j \neq i_t$

$$x_{t+1} = x_t - \gamma \left[\sum_{i=1}^n g_i^{(t+1)} \right] \leftarrow \text{store } \underline{\text{stable}} \text{ gradients}$$

$$\left[\sum_{i=1}^n g_i^{(t)} + g_{i_t}^{(t+1)} - g_{i_t}^{(t)} \right]$$

$\frac{1}{n} \sum_i g_i \approx \text{approx of } \nabla f(x^*)$

$O(1)$ cost per iteration

big surprise: converge linearly and fast

[but $O(n)$ storage cost] in worst case

"increment aggregated gradient" (IAG) [Blatt et al. 2007]
where you cycle deterministically through $\{1, \dots, n\}$

→ linear rate for quadratic function
but $\gamma_{\max} \approx O\left(\frac{1}{nL}\right)$ (tiny rate)
→ big step size vs. $\frac{1}{nL}$

SAG convergence rate:

thm: with $\gamma_t = \frac{1}{16L}$

where $L = \max_i (\text{Lipschitz}(\nabla f_i))$

$$\mathbb{E} f(x_t) - f^* \leq \left(1 - \min\left\{\frac{\mu}{16L}, \frac{1}{8n}\right\}\right)^t C_0$$

constant

$$\rho_{\text{SAG}} = \min\left\{\frac{1}{16K_{\text{SAG}}}, \frac{1}{8n}\right\} \text{ vs. } \rho_{\text{grad}} \approx \frac{1}{K_{\text{grad}}}$$

example: l_2 -reg. log regression on RCV1

$$n = 700k \quad L = 0.25 \quad \mu = \frac{1}{n} \quad (K = \frac{n}{4})$$

rate comparison

$$\dots = 1 - \frac{\mu}{16L} \quad (1 - \frac{1}{8n})$$

rate comparison

gradient method $\left(\frac{L-\mu}{L+\mu}\right)^2 = 0.9998$

accelerated gradient (Nesterov) $\left(1-\sqrt{\frac{\mu}{L}}\right) = 0.99761$

Nesterov lower bound $\left(\frac{\sqrt{L}-\sqrt{\mu}}{\sqrt{L}+\sqrt{\mu}}\right) = 0.99048$

SAG $(n \text{ iterations}) \quad (1-\rho_{SAG})^n \approx 0.88250$

practical aspects (see Schmidt & al. Math prog. 2016 paper)

a) storage: if $f_i(w) = h(x_i^T w) \Rightarrow \nabla_{w_i} f_i(w) = \underbrace{h'(x_i^T w)}_{\text{scalar}} x_i^T$
 instead $O(n \cdot d)$ storage $\leadsto O(n)$ storage

input data

b) initialization of g_i ? best: run SGD for one pass then switch SAG/SAGA

c) step-size?

$\frac{1}{(4)L}$

• cheap line search heuristic (comes from FISTA)

while $\left\{ \begin{aligned} f_i(w_t - \frac{1}{L_i} \nabla f_i(w_t)) &\geq f_i(w_t) - \frac{1}{2L_i} \|\nabla f_i(w_t)\|^2 \\ \text{set } \tilde{L}_{i,new} &= 2 \tilde{L}_{i,old} \\ \text{else } \tilde{L}_{i,new} &= \left(\frac{1}{2}\right)^{1/4} \tilde{L}_{i,old} \end{aligned} \right.$

"wrong inequality"

d) non-uniform sampling?

sample $i \sim \frac{\tilde{L}_i}{\sum_j \tilde{L}_j}$

e) stopping criterion?

you can use $\frac{1}{n} \sum_j g_j^{(t)}$ as approx $\nabla f(x_t) = \frac{1}{n} \sum_{i=1}^n \nabla f_i(x_t)$

f) sparse features?

$x_{t+1} = x_t - \gamma \left[\underbrace{\nabla f_i(x_t)}_{\text{sparse}} - g_i^{(t)} + P_{S_i} \left[\frac{1}{n} \sum_{j=1}^n g_j^{(t)} \right] \right]$ (sparse SAGA)

weighted projection on support of x_i
 $S_i = \{u \mid (x_i)_u \neq 0\}$