

Lecture 19 - SAG

Tuesday, March 28, 2023 2:32 PM

today: • finish BCFW
• SAG

application of BCFW to SUMstruct:

- getting s_i^* is one loss-augmented decoding call for example i

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

$w = A\alpha = \sum_i A_i \alpha_i$
(worst case) $\rightarrow O(nd)$
need to store these in memory
or
store α_i 's

$$\begin{cases} w_i^{(t+1)} = w_i^{(t)} + \gamma_t (w_s^{(t)} - w_i^{(t)}) \\ w_j^{(t+1)} = w_j^{(t)} \quad \forall j \neq i \end{cases}$$

(memory/computation tradeoff)

note: for line search, also need to store $b^T q_i \triangleq \ell_{i,\alpha}^{(t)}$

convergence constants:

for SUMstruct, can show that $\mathcal{G}^{(i)} \leq \frac{4R_i^2}{\lambda n^2}$ $R_i \triangleq \max_{y \in \mathcal{Y}_i} \|\nabla_i(y)\|_2$

Hessian: $(\nabla A^T A)_{(i,y), (j,y)}$

$$= \lambda \sum_{k=2}^n \langle \nabla_k(y), \nabla_k(y) \rangle$$

$$d_i^T H_i d_i \leq \lambda \max_{(y,y')} \langle \nabla_i(y), \nabla_i(y') \rangle \leq \|d_i\|_2^2 \leq R_i^2$$

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

$$R \triangleq \max_i R_i$$

compare with

$$\mathcal{G} \leq \frac{4R^2}{\lambda 1}$$

i.e. BCFW is "n times" faster than batch FW for SUMstruct?

importance of affine invariance:

$$\mathcal{G} \leq L_{1,1} \text{diag}_{1,1}(M)^2$$

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

Lipschitz constant in ℓ_2 -norm \approx largest eigenvalue of Hessian

recall Hessian $\lambda A^T A = \frac{1}{\lambda n^2} \left(\langle \nabla_i f(y), \nabla_j f(y) \rangle \right)_{(i,y),(j,y)}$

say e.g. $\langle \nabla_i f(y), \nabla_i f(y) \rangle \gg 1$ for lots of output
 $(A^T A)_{ii} = (\text{dim}) \cdot 1$

\rightarrow get largest eigenvalue can scale with dim of a matrix
 \Rightarrow really bad here because dim is exponentially

- instead, want to use ℓ_1 -norm of Δ_{Hess}

i.e. ℓ_1 -norm for Lipschitz constant

$$\text{then get } L_i (\text{diam}(\Delta_{\text{Hess}}))^2 \approx \frac{4R^2}{\lambda}$$

variance reduced SGD

Setup: $\min_x \frac{1}{n} \sum_{i=1}^n f_i(x)$

where f is μ -strongly convex

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

linear rate

batch gradient method
 [Cauchy 18th century]

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

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

$O(n)$ to compute

$$K \triangleq \frac{L}{\mu} \quad \text{"condition #!"}$$

stochastic gradient method
 aka. incremental gradient method
 [Robbins & Monro 1951]

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

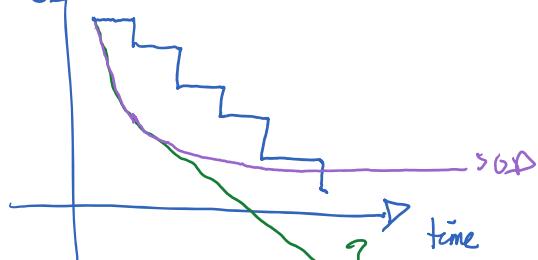
$\gamma_t = \text{const. } \gamma$
 \Rightarrow linear rate up to a ball of radius γ

where $i_t \sim \text{unif}\{1, \dots, n\}$
 $O(1)$ to compute

$$\gamma_t = \frac{1}{\mu t}$$

$\Rightarrow \tilde{O}\left(\frac{1}{\sqrt{t}}\right)$ rate
 i.e. sublinear

$$\log(f(x_t) - f^*)$$



a variance reduced SGD method?

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

15h33

SAG (stochastic average gradient)

[Le Roux, Schmidt & Bach NeurIPS 2012]

(Sagnaneq prize 2018)

- SAG:
 - store past gradient for each i (\bar{g}_i)
 - update one at step t

SAG: pick i_t unif.; update $\bar{g}_{i_t}^{(t+1)} = \nabla f_{i_t}(x_t)$
 $\bar{g}_j^{(t+1)} = \bar{g}_j^{(t)}$ for $j \neq i_t$

$$x_{t+1} = x_t - \gamma \left[\frac{1}{n} \sum_{i=1}^n \bar{g}_i^{(t+1)} \right] \quad \text{& stored stale gradients}$$

$$\left[\frac{1}{n} \sum_{i=1}^n \bar{g}_i^{(t)} + \bar{g}_{i_t}^{(t+1)} - \bar{g}_{i_t}^{(t)} \right]$$

$\frac{1}{n} \bar{g}_i \approx \text{approx of } \nabla f(x_t)$

- O(1) cost per iteration big surprise: converge linearly and fast

[but $O(nd)$ 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 fit.
 big stepsize vs but $\gamma_{\max} \rightarrow O(\frac{1}{n})$ (stingy rule)

SAG convergence rate

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

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

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

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

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

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

$$\text{accelerated gradient (Nesterov)} \quad \left(1 - \frac{\mu}{L}\right) = 0.7761$$

$$\text{SAG (n iterations)} \quad (1 - \rho_{\text{SAG}})^n \approx 0.8825$$

practical aspects [see Schmidt et al. Math Prog. 2016 paper] input data

a) storage: if $f_i(w) = h(x_i^T w) \Rightarrow \nabla_w f_i(w) = h'(x_i^T w) x_i$
instead of $O(n \cdot d)$ storage \rightarrow $O(n)$ storage

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

c) step size? • $\frac{1}{(4)L}$ (wrong inequality!)

• cheap line search heuristic
comes from FISTA while $\begin{cases} f_i\left(w_t - \frac{1}{L} \nabla f_i(w_t)\right) \geq f_i(w_t) - \frac{1}{2L} \|\nabla f_i(w_t)\|^2 \\ \text{set } L_{i,\text{new}} = 2L_{i,\text{old}} \end{cases}$
else $L_{i,\text{new}} = \left(\frac{1}{2}\right)^{1/n} L_{i,\text{old}}$

d) non-uniform sampling?

sample $i \sim \frac{L_i}{\sum_j L_j}$

e) stopping criterion?

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

f) sparse features?

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

weighted projection on support of x_t

$$S_{i,t} \stackrel{\Delta}{=} \{u : (x_t)_u \neq 0\}$$