

today: variance reduced SGD
• CRF

Variance reduced SGD

setup: $\min_x \underbrace{\frac{1}{n} \sum_i f_i(x)}_{\triangleq f(x)}$ where f is μ -strongly convex
 L -smooth

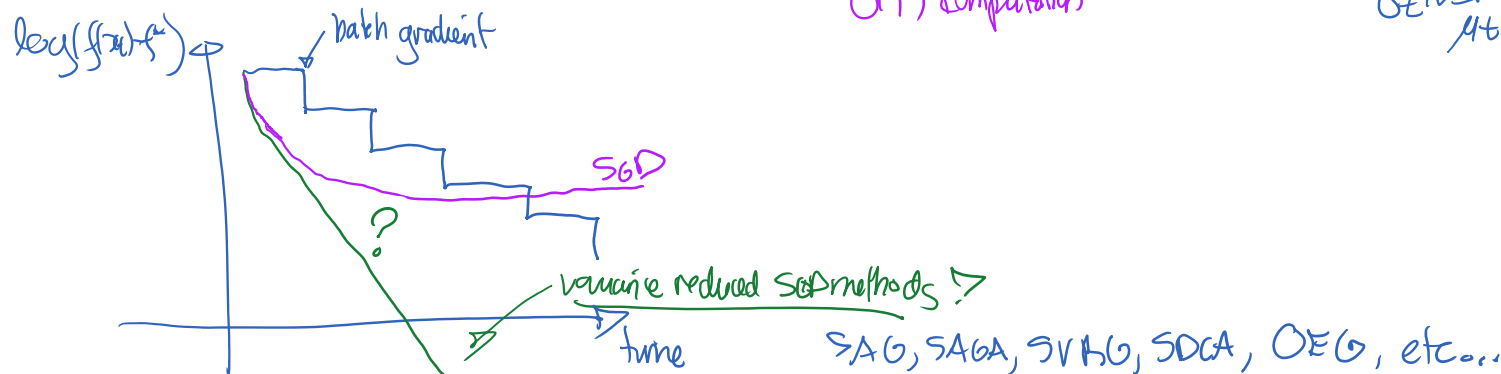
batch gradient method: $x_{t+1} = x_t - \gamma \underbrace{\left(\frac{1}{n} \sum_i \nabla f_i(x_t) \right)}_{O(n) \text{ to compute}}$

$\gamma = \frac{1}{L}$ linear rate
 $f(x_t) - f^* \leq (1-\rho)^t (f(x_0) - f^*)$
where $\rho \triangleq \frac{L}{L+\mu} = \frac{1}{2}$ $K \triangleq \frac{L}{\mu}$
condition #

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

$x_{t+1} = x_t - \gamma_t \nabla f_{i_t}(x_t)$
where $i_t \sim \text{unif} \{1, \dots, n\}$
 $O(1)$ computation

$\gamma_t = \text{const } \gamma \rightarrow$ linear rate up to a ball of radius γ
 $\gamma_t \sim \frac{1}{\sqrt{t}} \rightarrow O\left(\frac{1}{\sqrt{t}}\right)$ rate (sublinear)



SAG (stochastic average gradient) [LeKorun, Schmidt & Bach 2012]

SAG: • store past gradients for each i
 • update one at step t

SAG $\left\{ \begin{array}{l} \text{pick } i_t \sim \text{unif}; \text{ update } g_{i_t}^{(t+1)} \triangleq \nabla f_{i_t}(x_t); \quad g_j^{(t+1)} = g_j^{(t)} \quad \forall j \neq i_t \\ x_{t+1} = x_t - \frac{1}{n} \sum_{i=1}^n g_i^{(t+1)} \end{array} \right.$

$\underbrace{\frac{1}{n} \sum_{i=1}^n g_i^{(t+1)}}_{\text{is an approximation of } \nabla F(x_t)}$ ← stable gradients

$O(1)$ cost per iteration
 (but $O(n)$ storage cost)

big surprise: converge linearly and fast

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

→ linear rate for quadratic function
 but $\delta_{\max} \approx O(\frac{1}{n})$

"big surprise" ←

SAG convergence rate: thm: with $\delta_t = \frac{1}{16L}$ 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

ie.

$$\mathcal{R}_{\text{SAG}} = \min\left\{\frac{1}{16\kappa_{\text{SAG}}}, \frac{1}{8n}\right\} \quad \text{compare with } \mathcal{R}_{\text{grad}} \approx \frac{1}{\kappa_{\text{grad}}}$$

example: log regression on RCV1

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

rate comparison:

gradient method $\left(\frac{L-\mu}{L+\mu}\right)^2 = 0.99998$

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

Nesterov lower bound: has rate $\left(\frac{\sqrt{L}-\sqrt{\mu}}{\sqrt{L}+\sqrt{\mu}}\right)^2 = 0.99048$

SA6: $(1 - \rho_{\text{SA6}})^n = 0.88250$
(n iterations)

practical aspects (see Schmidt journal paper)

a) storage: if $f_i(w) = h(x_i^T w) \Rightarrow \nabla f_i(w) = \underbrace{h'(x_i^T w)}_{\text{scalar}} \overset{\substack{\text{input data} \\ \downarrow}}{x_i}$

instead of $O(nd)$ storage $\leadsto O(n)$

b) initialization? best is run SGD for one pass, then start SA6/SA6A etc...
of g_i 's memory

c) step-size? • $1/L$

• cheap one search heuristic

(comes from FISTA)

$$\left\{ \begin{array}{l} \text{while } f_i(w_t - \frac{1}{\tilde{L}_i} \nabla f_i(w_t)) \geq f_i(w_t) - \frac{1}{2\tilde{L}_i} \|\nabla f_i(w_t)\|^2 \\ \text{set } \tilde{L}_{\text{new}} = 2 \tilde{L}_{\text{old}} \end{array} \right.$$

\tilde{L}_i is surrogate for L_i

use step-size $\frac{1}{\cdot}$

... ..

d) non-uniform sampling?

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

use step-size $\frac{1}{L_i}$

e) stopping criterion?

you can use $\frac{1}{n} \sum_j g_j^t$ as approximate $\nabla f(w_t)$

f) sparse features?

$$w_{t+1} = w_t - \gamma \left(\underbrace{\nabla_{S_i}(w_t)}_{\text{sparse}} - g_i^t + P_{S_i} \left(\underbrace{\frac{1}{n} \sum_j g_j^t}_{\text{dense}} \right) \right) \quad (\text{Sparse SA6A})$$

[Leblond et al. 2017]

weighted projection on support of x_i

$$S_i \triangleq \{u : (x_i)_u \neq 0\}$$

15h47

Variance reduction idea:

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

properties: $\mathbb{E}G_\alpha = \alpha \mathbb{E}X + (1-\alpha) \mathbb{E}Y \rightsquigarrow$ unbiased (ie. $\mathbb{E}G_\alpha = \mathbb{E}X$)
if $\mathbb{E}Y = \mathbb{E}X$ [not interesting]
 $\alpha = 1$

$$\text{variance } \text{Var}(G_\alpha) = \alpha^2 [\text{var}(X) + \text{var}(Y) - 2\text{cov}(X, Y)]$$

$$\text{variance} \quad \text{variance} \rightarrow \underbrace{\frac{1}{n} \sum_{i=1}^n \text{var}(g_i)}_{\text{variance reduction}} + \underbrace{\frac{1}{n} \sum_{i=1}^n \text{var}(g_i)}_{\text{variance reduction}}$$

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

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

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

$$\text{SAG: } w_{t+1} = w_t - \underbrace{\delta}_{\alpha = \frac{1}{n}} \left[\underbrace{\nabla f_i(w_t)}_X - \underbrace{g_{i_t}^t}_Y + \underbrace{\frac{1}{n} \sum_j g_j^t}_{\mathbb{E}Y} \right] \quad (\text{biased})$$

$$\text{SAGA: } w_{t+1} = w_t - \underbrace{\delta}_{\alpha = 1} \left[\nabla f_{i_t}(w_t) - g_{i_t}^t + \frac{1}{n} \sum_j g_j^t \right] \quad (\text{unbiased})$$

SVRG: (stochastic variance reduced gradient)

$$w_{t+1} = w_t - \delta \left[\nabla f_{i_t}(w_t) - \nabla f_{i_t}(w_{\text{old}}) + \frac{1}{n} \sum_{j=1}^n \nabla f_j(w_{\text{old}}) \right] \quad (\text{unbiased})$$

w_{old} is updated from outer loop

CRF objective

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

$$\begin{aligned}
 \text{ker} \rightarrow w(\alpha) &= \frac{1}{\lambda n} \sum_i \sum_{\tilde{y}} \alpha_i(\tilde{y}) \psi_i(\tilde{y}) \\
 &= \frac{1}{\lambda n} \sum_i \sum_{\tilde{y}} M_{i,c}(\tilde{y}_c) \psi_{i,c}(\tilde{y}_c)
 \end{aligned}$$

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

\rightarrow from MRF

\Rightarrow at optimality

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

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

unlike sparse solution in structured SVM

CRF optimization:

SAG for CRF

$$w^{(t+1)} = (1 - \lambda \delta_t) w^{(t)} - \delta_t \left[\underbrace{\nabla_{\tilde{y}_i} \psi_i(w^{(t)})}_{\text{neg. dual objective}} + \frac{1}{n} \sum_j \tilde{y}_j \right]$$

$\odot \text{EG}$
 online exponentiated gradient

$$\alpha_{i,t}^*(\tilde{y})^{(t+1)} \propto \alpha_{i,t}^*(\tilde{y})^{(t)} \exp(-\delta_t \nabla_{\alpha_{i,t}^*(\tilde{y})} D(\alpha^{(t)}))$$

neg. dual objective

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

SDCA
Stochastic dual
coordinate ascent

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

$$Q_{i,t}(\tilde{y})^{(t+1)} = (1-\delta_t) Q_{i,t}(\tilde{y})^{(t)} + \delta_t \underbrace{\tilde{S}_i(\tilde{y})^{(t)}}_{\substack{\text{get this using marginal} \\ \text{inference}}} \quad \text{prox term}$$

$$p(\tilde{y} | x_i; w(\alpha_t))$$