

today:

- variance reduced SGD
- prox method, catalyst, etc ..

### Variance reduction idea :

$$\text{here } f(x) = \frac{1}{n} \sum_i g_i(x)$$

$$\nabla f = \frac{1}{n} \sum_i \nabla g_i \xrightarrow{\text{r.v.}} \mathbb{E}[X]$$

$$= \mathbb{E}[\nabla g_i] \approx \mathbb{E}[X]$$

$\mathbb{E}[\nabla g_i] \approx \mathbb{E}[X]$

SG.D. approximate  $\mathbb{E}[X]$  with just  $X$  (i.e.  $\nabla g_i$ )

general idea:

goal: estimate  $\mathbb{E}X$  using M.-C. samples

suppose:  $\mathbb{E}Y$  is cheap to compute and  $Y$  is correlated with  $X$

$\alpha \in [0, 1]$   
consider estimator  $\hat{G}_\alpha \stackrel{?}{=} \alpha(X - Y) + \mathbb{E}Y$  to approximate  $\mathbb{E}X$   
(here  $X, Y$  are r.v.)

$$\mathbb{E}\hat{G}_\alpha = \alpha \mathbb{E}X + (1-\alpha) \mathbb{E}Y \quad \rightsquigarrow \text{unbiased if } \mathbb{E}Y = \mathbb{E}X \quad [\text{never happens; otherwise use } \mathbb{E}Y]$$

$$\text{Variance } \text{Var}(\hat{G}_\alpha) = \underbrace{\alpha^2}_{\text{Variance reduction aspect}} [\text{Var}(X) + \text{Var}(Y) - 2\text{Cov}(X, Y)]$$

$$\text{for } \alpha=1 \text{ (unbiased)}: \quad \hat{G}_\alpha = X + \underbrace{[\mathbb{E}Y - Y]}_{\text{correction}}$$

for SGD context ..

~~your writing~~

$\nabla X$  is batch gradient ;  $X$  is  $\nabla f_i$

SAG/SAGA algorithm  $Y$  is  $g_i$  [past stored gradient]  
 $\bar{Y} = \frac{1}{n} \sum g_i$

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

standard SAG  $w_{t+1} = w_t - \gamma \frac{1}{n} \sum g_j^{t+1}$

---

$$\sum g_i^t + \nabla f_i(w_t) - g_i^t$$

SAG:

$$w_{t+1} = w_t - \gamma \left[ \nabla f_i(w_t) - g_i^t + \frac{1}{n} \sum_j g_j^t \right] \quad (\text{biased})$$

SAGA:

$$w_{t+1} = w_t - \gamma \left[ \nabla f_i(w_t) - g_i^t + \frac{1}{n} \sum_j g_j^t \right] \quad (\text{unbiased})$$

i.e.  $\mathbb{E}[ \cdot ]|w_t = \nabla f(w_t)$

SVRG  
(stochastic  
variance reduced  
gradient)

$$w_{t+1} = w_t - \gamma \left[ \nabla f_i(w_t) - \nabla f_i(w_{\text{old}}) + \frac{1}{n} \sum_j \nabla f_j(w_{\text{old}}) \right] \quad (\text{for fixed wold  
from outer loop})$$

SVRG for

$K=0, \dots$ (outer loop) compute $g_{\text{ref}} = \frac{1}{n} \sum_i \nabla f_i(w^{(k)})$ $w_0^{(k)} = w^{(k)}$ for $t=0, \dots, t_{\max}$ sample $i_t$ $w_{t+1}^{(k)} = w_t^{(k)} - \gamma \left[ \nabla f_{i_t}(w_t^{(k)}) - \nabla f_{i_t}(w_0^{(k)}) + g_{\text{ref}} \right]$ end $w^{(k+1)} = w_{t_{\max}}^{(k)}$
---

Questions:

- what is  $t_{\max}$ ?
- what is  $\gamma$ ?

SAG:

- need to store  $g_i^t$
- $O(nd)$

- no  $t_{\max}$  to tune
- no wend 2 loops

SVRG:

- fixed  $t_{\max}$
- wend 2 loops
- 2 gradients per iteration

- only need to store  $w^{(k)}$

SVRG convergence result:

$$\text{need } \gamma \leq \frac{1}{L}$$

$$t_{\max} \geq \frac{L}{\mu} = K$$

$\Rightarrow$  in practice,  $t_{\max} = \max \{ n, K \}$

important remark: "adversarial, for strong convexity"

important concept: "adaptivity for strong convexity"

SAG result  $\gamma \approx \frac{1}{L} \Rightarrow$  algorithmic parameters do not depend on  $\mu$  (strong convexity)

SVRG  $\gamma \approx \frac{1}{L}$  but  $t_{\max} = \max \{n, \frac{L}{\mu}\}$  depends on  $\mu$

for just convex fct. ( $\mu=0$ ), get  $\min_{\epsilon} [\mathbb{E} f(x) - f^*] = O\left(\frac{1}{\epsilon}\right)$

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

SAGA "simple" convergence result:

If you use  $\gamma = \frac{\alpha}{L}$  for SAGA for  $\alpha \leq 1$  ie.  $\mathbb{E} f(x) - f^* \leq (1-p)^t c_0$   
then rate  $p \geq \frac{1}{K} \min \left\{ \frac{1}{n}, \frac{\alpha}{K} \right\}$

two effects:  
• bigger step-size  $\Rightarrow$  bigger variance  
• smaller  $\|\cdot\|$   $\Rightarrow$  slower gradient descent rate (determined by  $\frac{1}{K}$ )

if  $K < n$ ; then as long as  $\frac{\alpha}{K} \geq \frac{1}{n}$  ie.  $\boxed{1 \geq \alpha \geq \frac{K}{n}}$ , we get same rate (roughly)  
 $\rightarrow$  "indifference" to step size when  $K < n$

SAGA complexity to reach  $\epsilon$  error  $O((n+K) \log \frac{1}{\epsilon})$  "SGD steps"

Practical aspects of SAG/SAGA:

a) storage: if  $f_i(w) = h(x_i^\top w)$   $Df_i(w) = \underbrace{h'(x_i^\top w)}_{\text{scalar}} x_i$

i.e. instead of  $O(nd)$  storage, only need  $O(n)$

b) initialization? hack which is to use  $\sum_{i \in S_t} g_i^t$  where  $S_t = \{i : i \text{ has been visited before } t\}$

c) step-size?

$$\cdot \frac{1}{L}$$

• cheap line search heuristic  $\left\{ \begin{array}{l} \text{while } f(w_t - \tilde{\gamma} \nabla f(w_t)) \geq f(w_t) - \frac{1}{2L} \|\nabla f(w_t)\|^2 \\ \text{set } \tilde{\gamma}^{new} = 2\tilde{\gamma}^{old} \end{array} \right.$

(comes from FISTA)

$\tilde{\gamma}$  is surrogate for  $L$

use stepsize  $\frac{1}{\tilde{\gamma}}$

d) non-uniform sampling?  $\rightarrow$  sample  $i \sim \frac{L_i}{\sum_j L_j}$

then get rate with  $\frac{\sum_i L_i}{\mu}$  instead of  $\max_i \frac{L_i}{\mu}$

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

f) sparse features? two tricks 1) "logged update"  $\rightarrow w_t^d = w_{t-1}^d - (\Delta \gamma) g_{d, t}$  [Schmidt et al.]

sparse modification  $w_{t+1} = w_t - \gamma (\nabla f(w_t) - g_t + P_{S_t}(\frac{1}{n} \sum_j g_j^t))$

[Lektor et al. 2017]

weights  
projection on support  
of  $x_t$

$$2) f_i(w) = \frac{\alpha \|w\|^2}{2} + h(x_i^T w)$$

$$w_{t+1} = (1 - \Delta \gamma) w_t - \gamma \left[ \underbrace{\quad}_{\text{store } w \text{ as } \{ \text{direction}, \text{scale} \}} \right]$$

$S_t = \{j : (x_j)_i \neq 0\}$

## Pointers

- variance reduction perspective on SAG / SAGA / SVRG -- in SAGA paper:
  - A. Defazio, F. Bach and S. Lacoste-Julien  
SAGA: A Fast Incremental Gradient Method With Support for Non-Strongly Convex Composite Objectives  
[NIPS 2014](#)
  - other pointers:
    - SVRG paper:  
Rie Johnson and Tong Zhang. Accelerating stochastic gradient descent using predictive variance reduction. [NIPS 2013](#)
    - note that a variant of SVRG that is \*adaptive to local strong convexity\* is given in the following paper (where the end of the inner loop is decided randomly: at every inner loop iteration, with probability 1/n, you end the inner loop):  
T. Hofmann, A. Lucchi, S. Lacoste-Julien, and Brian McWilliams  
Variance Reduced Stochastic Gradient Descent with Neighbors,  
[NIPS 2015](#)
- the practical aspects of SAG are described in the massive journal version:
  - M. Schmidt, N. Le Roux, F. Bach.  
Minimizing Finite Sums with the Stochastic Average Gradient  
[Mathematical Programming](#), 162:83-162, 2017 | [arxiv](#)
  - an alternative to the complicated "lagged updates" when you have sparse features is the Sparse SAGA algorithm; see Section 2 of:  
ASAGA: Asynchronous Parallel SAGA,  
R. Leblond, F. Pedregosa and S. Lacoste-Julien  
[AISTATS 2017](#)