

Variance reduction for stochastic gradient methods



Yuxin Chen

Princeton University, Fall 2019

Outline

- Stochastic variance reduced gradient (SVRG)
 - Convergence analysis for strongly convex problems
- Stochastic recursive gradient algorithm (SARAH)
 - Convergence analysis for nonconvex problems
- Other variance reduced stochastic methods
 - Stochastic dual coordinate ascent (SDCA)

Finite-sum optimization

$$\text{minimize}_{\mathbf{x} \in \mathbb{R}^d} \quad F(\mathbf{x}) := \frac{1}{n} \sum_{i=1}^n \underbrace{f_i(\mathbf{x})}_{\substack{\text{loss for } i\text{th sample} \\ (\mathbf{a}_i, y_i)}} + \underbrace{\psi(\mathbf{x})}_{\text{regularizer}}$$

common task in machine learning

- linear regression: $f_i(\mathbf{x}) = \frac{1}{2}(\mathbf{a}_i^\top \mathbf{x} - y_i)^2$, $\psi(\mathbf{x}) = 0$
- logistic regression: $f_i(\mathbf{x}) = \log(1 + e^{-y_i \mathbf{a}_i^\top \mathbf{x}})$, $\psi(\mathbf{x}) = 0$
- Lasso: $f_i(\mathbf{x}) = \frac{1}{2}(\mathbf{a}_i^\top \mathbf{x} - y_i)^2$, $\psi(\mathbf{x}) = \lambda \|\mathbf{x}\|_1$
- SVM: $f_i(\mathbf{x}) = \max\{0, 1 - y_i \mathbf{a}_i^\top \mathbf{x}\}$, $\psi(\mathbf{x}) = \frac{\lambda}{2} \|\mathbf{x}\|_2^2$
- ...

Stochastic gradient descent (SGD)

Algorithm 12.1 Stochastic gradient descent (SGD)

- 1: **for** $t = 1, 2, \dots$ **do**
 - 2: pick $i_t \sim \text{Unif}(1, \dots, n)$
 - 3: $\mathbf{x}^{t+1} = \mathbf{x}^t - \eta_t \nabla f_{i_t}(\mathbf{x}^t)$
-

As we have shown in the last lecture

- large stepsizes poorly suppress variability of stochastic gradients
 \implies SGD with $\eta_t \asymp 1$ tends to oscillate around global mins
- choosing $\eta_t \asymp 1/t$ mitigates oscillation, but is too conservative

Recall: SGD theory with fixed stepsizes

$$\mathbf{x}^{t+1} = \mathbf{x}^t - \eta_t \mathbf{g}^t$$

- \mathbf{g}^t : an unbiased estimate of $F(\mathbf{x}^t)$
- $\mathbb{E}[\|\mathbf{g}^t\|_2^2] \leq \sigma_g^2 + c_g \|\nabla F(\mathbf{x}^t)\|_2^2$
- $F(\cdot)$: μ -strongly convex; L -smooth

From the last lecture, we know

$$\mathbb{E}[F(\mathbf{x}^t) - F(\mathbf{x}^*)] \leq \frac{\eta L \sigma_g^2}{2\mu} + (1 - \eta\mu)^t (F(\mathbf{x}^0) - F(\mathbf{x}^*))$$

Recall: SGD theory with fixed stepsizes

$$\mathbb{E}[F(\mathbf{x}^t) - F(\mathbf{x}^*)] \leq \frac{\eta L \sigma_g^2}{2\mu} + (1 - \eta\mu)^t (F(\mathbf{x}^0) - F(\mathbf{x}^*))$$

- vanilla SGD: $\mathbf{g}^t = \nabla f_{i_t}(\mathbf{x}^t)$
 - **issue:** σ_g^2 is non-negligible even when $\mathbf{x}^t = \mathbf{x}^*$
- **question:** it is possible to design \mathbf{g}^t with reduced variability σ_g^2 ?

A simple idea

Imagine we take some \mathbf{v}^t with $\mathbb{E}[\mathbf{v}^t] = \mathbf{0}$ and set

$$\mathbf{g}^t = \nabla f_{i_t}(\mathbf{x}^t) - \mathbf{v}^t$$

— so \mathbf{g}^t is still an unbiased estimate of $\nabla F(\mathbf{x}^t)$

question: how to reduce variability (i.e. $\mathbb{E}[\|\mathbf{g}^t\|_2^2] < \mathbb{E}[\|\nabla f_{i_t}(\mathbf{x}^t)\|_2^2]$)?

answer: find some zero-mean \mathbf{v}^t that is positively correlated with $\nabla f_{i_t}(\mathbf{x}^t)$ (i.e. $\langle \mathbf{v}^t, \nabla f_{i_t}(\mathbf{x}^t) \rangle > 0$) ([why?](#))

Reducing variance via gradient aggregation

If the current iterate is not too far away from previous iterates, then historical gradient info might be useful in producing such a v^t to reduce variance

main idea of this lecture: aggregate previous gradient info to help improve the convergence rate

Stochastic variance reduced gradient (SVRG)

Strongly convex and smooth problems (no regularization)

$$\text{minimize}_{\mathbf{x} \in \mathbb{R}^d} \quad F(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{x})$$

- f_i : convex and L -smooth
- F : μ -strongly convex
- $\kappa := L/\mu$: condition number

Stochastic variance reduced gradient (SVRG)

— Johnson, Zhang '13

key idea: if we have access to a history point \mathbf{x}^{old} and $\nabla F(\mathbf{x}^{\text{old}})$, then

$$\underbrace{\nabla f_{i_t}(\mathbf{x}^t) - \nabla f_{i_t}(\mathbf{x}^{\text{old}})}_{\rightarrow \mathbf{0} \text{ if } \mathbf{x}^t \approx \mathbf{x}^{\text{old}}} + \underbrace{\nabla F(\mathbf{x}^{\text{old}})}_{\rightarrow \mathbf{0} \text{ if } \mathbf{x}^{\text{old}} \approx \mathbf{x}^*} \quad \text{with } i_t \sim \text{Unif}(1, \dots, n)$$

- is an unbiased estimate of $\nabla F(\mathbf{x}^t)$
- converges to $\mathbf{0}$ if $\mathbf{x}^t \approx \mathbf{x}^{\text{old}} \approx \mathbf{x}^*$
variability is reduced!

Stochastic variance reduced gradient (SVRG)

- operate in epochs
- in the s^{th} epoch
 - **very beginning:** take a snapshot $\mathbf{x}_s^{\text{old}}$ of the current iterate, and compute the *batch gradient* $\nabla F(\mathbf{x}_s^{\text{old}})$
 - **inner loop:** use the snapshot point to help reduce variance

$$\mathbf{x}_s^{t+1} = \mathbf{x}_s^t - \eta \{ \nabla f_{i_t}(\mathbf{x}_s^t) - \nabla f_{i_t}(\mathbf{x}_s^{\text{old}}) + \nabla F(\mathbf{x}_s^{\text{old}}) \}$$

a hybrid approach: the batch gradient is computed only once per epoch

SVRG algorithm (Johnson, Zhang '13)

Algorithm 12.2 SVRG for finite-sum optimization

1: **for** $s = 1, 2, \dots$ **do**

2: $\mathbf{x}_s^{\text{old}} \leftarrow \mathbf{x}_{s-1}^m$, and compute $\underbrace{\nabla F(\mathbf{x}_s^{\text{old}})}_{\text{batch gradient}}$ // update snapshot

3: initialize $\mathbf{x}_s^0 \leftarrow \mathbf{x}_s^{\text{old}}$

4: **for** $\underbrace{t = 0, \dots, m-1}_{\text{each epoch contains } m \text{ iterations}}$ **do**

5: choose i_t uniformly from $\{1, \dots, n\}$, and

$$\mathbf{x}_s^{t+1} = \mathbf{x}_s^t - \eta \left\{ \underbrace{\nabla f_{i_t}(\mathbf{x}_s^t) - \nabla f_{i_t}(\mathbf{x}_s^{\text{old}})}_{\text{stochastic gradient}} + \nabla F(\mathbf{x}_s^{\text{old}}) \right\}$$

Remark

- constant stepsize η
- each epoch contains $2m + n$ gradient computations
 - the batch gradient is computed only once every m iterations
 - the average per-iteration cost of SVRG is comparable to that of SGD if $m \gtrsim n$

Convergence analysis of SVRG

Theorem 12.1

Assume each f_i is convex and L -smooth, and F is μ -strongly convex. Choose m large enough s.t. $\rho = \frac{1}{\mu\eta(1-2L\eta)^m} + \frac{2L\eta}{1-2L\eta} < 1$, then

$$\mathbb{E}[F(\mathbf{x}_s^{\text{old}}) - F(\mathbf{x}^*)] \leq \rho^s [F(\mathbf{x}_0^{\text{old}}) - F(\mathbf{x}^*)]$$

- **linear convergence:** choosing $m \gtrsim L/\mu = \kappa$ and constant stepsizes $\eta \asymp 1/L$ yields $0 < \rho < 1/2$
 $\implies O(\log \frac{1}{\varepsilon})$ epochs to attain ε accuracy

Convergence analysis of SVRG

Theorem 12.1

Assume each f_i is convex and L -smooth, and F is μ -strongly convex. Choose m large enough s.t. $\rho = \frac{1}{\mu\eta(1-2L\eta)^m} + \frac{2L\eta}{1-2L\eta} < 1$, then

$$\mathbb{E}[F(\mathbf{x}_s^{\text{old}}) - F(\mathbf{x}^*)] \leq \rho^s [F(\mathbf{x}_0^{\text{old}}) - F(\mathbf{x}^*)]$$

- **total computational cost:**

$$\underbrace{(m+n)}_{\text{\# grad computation per epoch}} \log \frac{1}{\varepsilon} \asymp \underbrace{(n+\kappa)}_{\text{if } m \asymp \max\{n, \kappa\}} \log \frac{1}{\varepsilon}$$

Proof of Theorem 12.1

Here, we provide the proof for an alternative version, where in each epoch,

$$\mathbf{x}_{s+1}^{\text{old}} = \mathbf{x}_s^j \quad \underbrace{\text{with } j \sim \text{Unif}(0, \dots, m-1)}_{\text{rather than } j=m} \quad (12.1)$$

The interested reader is referred to Tan et al. '16 for the proof of the original version

Proof of Theorem 12.1

Let $\mathbf{g}_s^t := \nabla f_{it}(\mathbf{x}_s^t) - \nabla f_{it}(\mathbf{x}_s^{\text{old}}) + \nabla F(\mathbf{x}_s^{\text{old}})$ for simplicity. As usual, conditional on everything prior to \mathbf{x}_s^{t+1} , one has

$$\begin{aligned}\mathbb{E}[\|\mathbf{x}_s^{t+1} - \mathbf{x}^*\|_2^2] &= \mathbb{E}[\|\mathbf{x}_s^t - \eta \mathbf{g}_s^t - \mathbf{x}^*\|_2^2] \\ &= \|\mathbf{x}_s^t - \mathbf{x}^*\|_2^2 - 2\eta(\mathbf{x}_s^t - \mathbf{x}^*)^\top \mathbb{E}[\mathbf{g}_s^t] + \eta^2 \mathbb{E}[\|\mathbf{g}_s^t\|_2^2] \\ &\leq \|\mathbf{x}_s^t - \mathbf{x}^*\|_2^2 - 2\eta(\mathbf{x}_s^t - \mathbf{x}^*)^\top \underbrace{\nabla F(\mathbf{x}_s^t)}_{\text{since } \mathbf{g}_s^t \text{ is an unbiased estimate of } \nabla F(\mathbf{x}_s^t)} + \eta^2 \mathbb{E}[\|\mathbf{g}_s^t\|_2^2] \\ &\leq \|\mathbf{x}_s^t - \mathbf{x}^*\|_2^2 - \underbrace{2\eta(F(\mathbf{x}_s^t) - F(\mathbf{x}^*))}_{\text{by convexity}} + \eta^2 \mathbb{E}[\|\mathbf{g}_s^t\|_2^2]\end{aligned}\quad (12.2)$$

- **key step:** control $\mathbb{E}[\|\mathbf{g}_s^t\|_2^2]$
 - we'd like to upper bound it via the (relative) objective value

Proof of Theorem 12.1

main pillar: control $\mathbb{E}[\|\mathbf{g}_s^t\|_2^2]$ via ...

Lemma 12.2

$$\mathbb{E}[\|\mathbf{g}_s^t\|_2^2] \leq 4L[F(\mathbf{x}_s^t) - F(\mathbf{x}^*) + F(\mathbf{x}_s^{\text{old}}) - F(\mathbf{x}^*)]$$

this means if $\mathbf{x}_s^t \approx \mathbf{x}_s^{\text{old}} \approx \mathbf{x}^*$, then $\mathbb{E}[\|\mathbf{g}_s^t\|_2^2] \approx 0$ (reduced variance)

Proof of Theorem 12.1

main pillar: control $\mathbb{E}[\|\mathbf{g}_s^t\|_2^2]$ via ...

Lemma 12.2

$$\mathbb{E}[\|\mathbf{g}_s^t\|_2^2] \leq 4L[F(\mathbf{x}_s^t) - F(\mathbf{x}^*) + F(\mathbf{x}_s^{\text{old}}) - F(\mathbf{x}^*)]$$

this allows one to obtain: conditional on everything prior to \mathbf{x}_s^{t+1} ,

$$\begin{aligned} \mathbb{E}[\|\mathbf{x}_s^{t+1} - \mathbf{x}^*\|_2^2] &\leq (12.2) \\ &\leq \|\mathbf{x}_s^t - \mathbf{x}^*\|_2^2 - 2\eta[F(\mathbf{x}_s^t) - F(\mathbf{x}^*)] \\ &\quad + 4L\eta^2[F(\mathbf{x}_s^t) - F(\mathbf{x}^*) + F(\mathbf{x}_s^{\text{old}}) - F(\mathbf{x}^*)] \\ &= \|\mathbf{x}_s^t - \mathbf{x}^*\|_2^2 - 2\eta(1 - 2L\eta)[F(\mathbf{x}_s^t) - F(\mathbf{x}^*)] \\ &\quad + 4L\eta^2[F(\mathbf{x}_s^{\text{old}}) - F(\mathbf{x}^*)] \end{aligned} \quad (12.3)$$

Proof of Theorem 12.1 (cont.)

Taking expectation w.r.t. all history, we have

$$\begin{aligned} & 2\eta(1 - 2L\eta)m \mathbb{E}[F(\mathbf{x}_{s+1}^{\text{old}}) - F(\mathbf{x}^*)] \\ &= 2\eta(1 - 2L\eta) \sum_{t=0}^{m-1} \mathbb{E}[F(\mathbf{x}_s^t) - F(\mathbf{x}^*)] && \text{by (12.1)} \\ &\leq \underbrace{\mathbb{E}[\|\mathbf{x}_{s+1}^m - \mathbf{x}^*\|_2^2]}_{\geq 0} + 2\eta(1 - 2L\eta) \sum_{t=0}^{m-1} \mathbb{E}[F(\mathbf{x}_s^t) - F(\mathbf{x}^*)] \\ &\leq \mathbb{E}[\|\mathbf{x}_{s+1}^0 - \mathbf{x}^*\|_2^2] + 4Lm\eta^2[F(\mathbf{x}_s^{\text{old}}) - F(\mathbf{x}^*)] && \text{(apply (12.3) recursively)} \\ &= \mathbb{E}[\|\mathbf{x}_s^{\text{old}} - \mathbf{x}^*\|_2^2] + 4Lm\eta^2\mathbb{E}[F(\mathbf{x}_s^{\text{old}}) - F(\mathbf{x}^*)] \\ &\leq \frac{2}{\mu}\mathbb{E}[F(\mathbf{x}_s^{\text{old}}) - F(\mathbf{x}^*)] + 4Lm\eta^2\mathbb{E}[F(\mathbf{x}_s^{\text{old}}) - F(\mathbf{x}^*)] && \text{(strong convexity)} \\ &= \left(\frac{2}{\mu} + 4Lm\eta^2\right) \mathbb{E}[F(\mathbf{x}_s^{\text{old}}) - F(\mathbf{x}^*)] \end{aligned}$$

Proof of Theorem 12.1 (cont.)

Consequently,

$$\begin{aligned} & \mathbb{E}[F(\mathbf{x}_{s+1}^{\text{old}}) - F(\mathbf{x}^*)] \\ & \leq \frac{\frac{2}{\mu} + 4Lm\eta^2}{2\eta(1 - 2L\eta)m} \mathbb{E}[F(\mathbf{x}_s^{\text{old}}) - F(\mathbf{x}^*)] \\ & = \underbrace{\left(\frac{1}{\mu\eta(1 - 2L\eta)m} + \frac{2L\eta}{1 - 2L\eta} \right)}_{=\rho} \mathbb{E}[F(\mathbf{x}_s^{\text{old}}) - F(\mathbf{x}^*)] \end{aligned}$$

Applying this bound recursively establishes the theorem.

Proof of Lemma 12.2

$$\begin{aligned} & \mathbb{E}[\|\nabla f_{i_t}(\mathbf{x}_s^t) - \nabla f_{i_t}(\mathbf{x}_s^{\text{old}}) + \nabla F(\mathbf{x}_s^{\text{old}})\|_2^2] \\ &= \mathbb{E}[\|\nabla f_{i_t}(\mathbf{x}_s^t) - \nabla f_{i_t}(\mathbf{x}^*) - (\nabla f_{i_t}(\mathbf{x}_s^{\text{old}}) - \nabla f_{i_t}(\mathbf{x}^*) - \nabla F(\mathbf{x}_s^{\text{old}}))\|_2^2] \\ &\leq 2\mathbb{E}[\|\nabla f_{i_t}(\mathbf{x}_s^t) - \nabla f_{i_t}(\mathbf{x}^*)\|_2^2] + 2\mathbb{E}[\|\nabla f_{i_t}(\mathbf{x}_s^{\text{old}}) - \nabla f_{i_t}(\mathbf{x}^*) - \nabla F(\mathbf{x}_s^{\text{old}})\|_2^2] \\ &= 2\mathbb{E}[\|\nabla f_{i_t}(\mathbf{x}_s^t) - \nabla f_{i_t}(\mathbf{x}^*)\|_2^2] \\ &\quad + 2\mathbb{E}[\|\nabla f_{i_t}(\mathbf{x}_s^{\text{old}}) - \nabla f_{i_t}(\mathbf{x}^*) - \underbrace{\mathbb{E}[\nabla f_{i_t}(\mathbf{x}_s^{\text{old}}) - \nabla f_{i_t}(\mathbf{x}^*)]}_{\text{since } \mathbb{E}[\nabla f_{i_t}(\mathbf{x}^*)] = \nabla F(\mathbf{x}^*) = \mathbf{0}}\|_2^2] \\ &\leq 2\mathbb{E}[\|\nabla f_{i_t}(\mathbf{x}_s^t) - \nabla f_{i_t}(\mathbf{x}^*)\|_2^2] + 2\mathbb{E}[\|\nabla f_{i_t}(\mathbf{x}_s^{\text{old}}) - \nabla f_{i_t}(\mathbf{x}^*)\|_2^2] \\ &\leq 4L[F(\mathbf{x}_s^t) - F(\mathbf{x}^*) + F(\mathbf{x}_s^{\text{old}}) - F(\mathbf{x}^*)] \end{aligned}$$

where the last inequality would hold if we could justify

$$\underbrace{\frac{1}{n} \sum_{i=1}^n \|\nabla f_i(\mathbf{x}) - \nabla f_i(\mathbf{x}^*)\|_2^2}_{\text{relies on both smoothness and convexity of } f_i} \leq 2L[F(\mathbf{x}) - F(\mathbf{x}^*)] \quad (12.4)$$

Proof of Lemma 12.2 (cont.)

To establish (12.4), observe from smoothness and convexity of f_i that

$$\underbrace{\frac{1}{2L} \|\nabla f_i(\mathbf{x}) - \nabla f_i(\mathbf{x}^*)\|_2^2 \leq f_i(\mathbf{x}) - f_i(\mathbf{x}^*) - \nabla f_i(\mathbf{x}^*)^\top (\mathbf{x} - \mathbf{x}^*)}_{\text{an equivalent characterization of } L\text{-smoothness}}$$

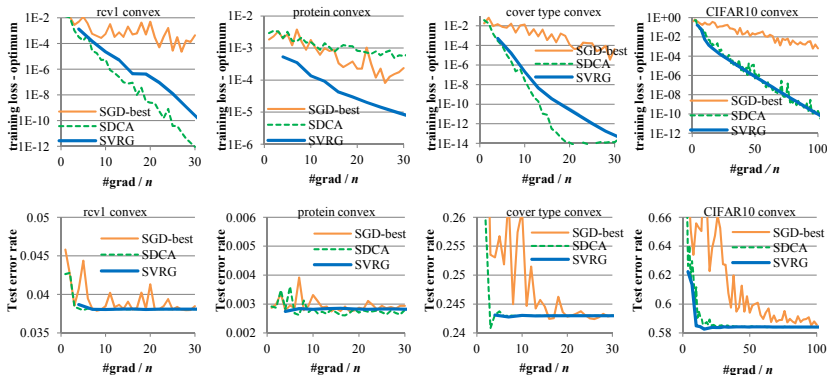
Summing over all i and recognizing that $\nabla F(\mathbf{x}^*) = \mathbf{0}$ yield

$$\begin{aligned} \frac{1}{2L} \sum_{i=1}^n \|\nabla f_i(\mathbf{x}) - \nabla f_i(\mathbf{x}^*)\|_2^2 &\leq nF(\mathbf{x}) - nF(\mathbf{x}^*) - n(\nabla F(\mathbf{x}^*))^\top (\mathbf{x} - \mathbf{x}^*) \\ &= nF(\mathbf{x}) - nF(\mathbf{x}^*) \end{aligned}$$

as claimed

Numerical example: logistic regression

— Johnson, Zhang '13



ℓ_2 -regularized logistic regression on CIFAR-10

Comparisons with GD and SGD

	SVRG	GD	SGD
comp. cost	$(n + \kappa) \log \frac{1}{\epsilon}$	$n\kappa \log \frac{1}{\epsilon}$	$\frac{\kappa^2}{\epsilon}$ (practically often $\frac{\kappa}{\epsilon}$)

Proximal extension

$$\text{minimize}_{\mathbf{x} \in \mathbb{R}^d} \quad F(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{x}) + \psi(\mathbf{x})$$

- f_i : convex and L -smooth
- F : μ -strongly convex
- $\kappa := L/\mu$: condition number
- ψ : potentially non-smooth

Proximal extension (Xiao, Zhang '14)

Algorithm 12.3 Prox-SVRG for finite-sum optimization

1: **for** $s = 1, 2, \dots$ **do**

2: $\mathbf{x}_s^{\text{old}} \leftarrow \mathbf{x}_{s-1}^m$, and compute $\underbrace{\nabla F(\mathbf{x}_s^{\text{old}})}_{\text{batch gradient}}$ // update snapshot

3: initialize $\mathbf{x}_s^0 \leftarrow \mathbf{x}_s^{\text{old}}$

4: **for** $\underbrace{t = 0, \dots, m-1}_{\text{each epoch contains } m \text{ iterations}}$ **do**

5: choose i_t uniformly from $\{1, \dots, n\}$, and

$$\mathbf{x}_s^{t+1} = \text{prox}_{\eta\psi} \left(\mathbf{x}_s^t - \eta \underbrace{\left\{ \nabla f_{i_t}(\mathbf{x}_s^t) - \nabla f_{i_t}(\mathbf{x}_s^{\text{old}}) \right\}}_{\text{stochastic gradient}} + \nabla F(\mathbf{x}_s^{\text{old}}) \right)$$

Stochastic recursive gradient algorithm (SARAH)

Nonconvex and smooth problems

$$\text{minimize}_{\mathbf{x} \in \mathbb{R}^d} \quad F(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{x})$$

- f_i : L -smooth, potentially nonconvex

Recursive stochastic gradient estimates

— Nguyen, Liu, Scheinberg, Takac '17

key idea: recursive / adaptive updates of gradient estimates
stochastic

$$\begin{aligned} \mathbf{g}^t &= \nabla f_{i_t}(\mathbf{x}^t) - \nabla f_{i_t}(\mathbf{x}^{t-1}) + \mathbf{g}^{t-1} \\ \mathbf{x}^{t+1} &= \mathbf{x}^t - \eta \mathbf{g}^t \end{aligned} \quad (12.5)$$

comparison to SVRG (use a **fixed** snapshot point for the entire epoch)

$$(\text{SVRG}) \quad \mathbf{g}^t = \nabla f_{i_t}(\mathbf{x}^t) - \nabla f_{i_t}(\mathbf{x}^{\text{old}}) + \nabla F(\mathbf{x}^{\text{old}})$$

Restarting gradient estimate every epoch

For many (e.g. strongly convex) problems, recursive gradient estimate \mathbf{g}^t may decay fast (variance \downarrow ; bias (relative to $\nabla F(\mathbf{x}^t)$) \uparrow)

- \mathbf{g}^t may quickly deviate from the target gradient $\nabla F(\mathbf{x}^t)$
- progress stalls as \mathbf{g}^t cannot guarantee sufficient descent

solution: reset \mathbf{g}^t every few iterations to calibrate with the true batch gradient

Bias of gradient estimates

Unlike SVRG, \mathbf{g}^t is NOT an unbiased estimate of $\nabla F(\mathbf{x}^t)$

$$\mathbb{E}[\mathbf{g}^t \mid \text{everything prior to } \mathbf{x}_s^t] = \nabla F(\mathbf{x}^t) \underbrace{-\nabla F(\mathbf{x}^{t-1}) + \mathbf{g}^{t-1}}_{\neq 0}$$

But if we average out all randomness, we have (exercise!)

$$\mathbb{E}[\mathbf{g}^t] = \mathbb{E}[\nabla F(\mathbf{x}^t)]$$

Stochastic Recursive gradient algorithm

Algorithm 12.4 SARAH (Nguyen et al. '17)

- 1: **for** $s = 1, 2, \dots, S$ **do**
 - 2: $\mathbf{x}_s^0 \leftarrow \mathbf{x}_{s-1}^{m+1}$, and compute $\underbrace{\mathbf{g}_s^0 = \nabla F(\mathbf{x}_s^0)}_{\text{batch gradient}}$ // restart \mathbf{g} anew
 - 3: $\mathbf{x}_s^1 = \mathbf{x}_s^0 - \eta \mathbf{g}_s^0$
 - 4: **for** $t = 1, \dots, m$ **do**
 - 5: choose i_t uniformly from $\{1, \dots, n\}$
 - 6: $\mathbf{g}_s^t = \underbrace{\nabla f_{i_t}(\mathbf{x}_s^t) - \nabla f_{i_t}(\mathbf{x}_s^{t-1})}_{\text{stochastic gradient}} + \mathbf{g}_s^{t-1}$
 - 7: $\mathbf{x}_s^{t+1} = \mathbf{x}_s^t - \eta \mathbf{g}_s^t$
-

Convergence analysis of SARAH (nonconvex)

Theorem 12.3 (Nguyen et al. '19)

Suppose each f_i is L -smooth. Then SARAH with $\eta \lesssim \frac{1}{L\sqrt{m}}$ obeys

$$\frac{1}{(m+1)S} \sum_{s=1}^S \sum_{t=0}^m \mathbb{E} \left[\|\nabla F(\mathbf{x}_s^t)\|_2^2 \right] \leq \frac{2}{\eta(m+1)S} [F(\mathbf{x}_0^0) - F(\mathbf{x}^*)]$$

- iteration complexity for finding ε -approximate stationary point (i.e. $\|\nabla F(\mathbf{x})\|_2 \leq \varepsilon$):

$$O \left(n + \frac{L\sqrt{n}}{\varepsilon^2} \right) \quad (\text{setting } m \asymp n, \eta \asymp \frac{1}{L\sqrt{m}})$$

Convergence analysis of SARAH (nonconvex)

Theorem 12.3 (Nguyen et al. '19)

Suppose each f_i is L -smooth. Then SARAH with $\eta \lesssim \frac{1}{L\sqrt{m}}$ obeys

$$\frac{1}{(m+1)S} \sum_{s=1}^S \sum_{t=0}^m \mathbb{E} \left[\|\nabla F(\mathbf{x}_s^t)\|_2^2 \right] \leq \frac{2}{\eta(m+1)S} [F(\mathbf{x}_0^0) - F(\mathbf{x}^*)]$$

- also derived by Fang et al. '18 (for a SARAH-like algorithm “Spider”) and improved by Wang et al. '19 (for “SpiderBoost”)

Proof of Theorem 12.3

Theorem 12.3 follows immediately from the following claim on the total objective improvement in one epoch (why?)

$$\mathbb{E}[F(\mathbf{x}_s^{m+1})] \leq \mathbb{E}[F(\mathbf{x}_s^0)] - \frac{\eta}{2} \sum_{t=0}^m \mathbb{E}[\|\nabla F(\mathbf{x}_s^t)\|_2^2] \quad (12.6)$$

We will then focus on establishing (12.6)

Proof of Theorem 12.3 (cont.)

To establish (12.6), recall that the smoothness assumption gives

$$\mathbb{E}[F(\mathbf{x}_s^{t+1})] \leq \mathbb{E}[F(\mathbf{x}_s^t)] - \eta \mathbb{E}[\nabla F(\mathbf{x}_s^t)^\top \mathbf{g}_s^t] + \frac{L\eta^2}{2} \mathbb{E}[\|\mathbf{g}_s^t\|_2^2] \quad (12.7)$$

Since \mathbf{g}_s^t is not an unbiased estimate of $\nabla F(\mathbf{x}_s^t)$, we first decouple

$$2\mathbb{E}[\nabla F(\mathbf{x}_s^t)^\top \mathbf{g}_s^t] = \underbrace{\mathbb{E}[\|\nabla F(\mathbf{x}_s^t)\|_2^2]}_{\text{desired gradient estimate}} + \underbrace{\mathbb{E}[\|\mathbf{g}_s^t\|_2^2]}_{\text{variance}} - \underbrace{\mathbb{E}[\|\nabla F(\mathbf{x}_s^t) - \mathbf{g}_s^t\|_2^2]}_{\text{squared bias of gradient estimate}}$$

Substitution into (12.7) with straightforward algebra gives

$$\begin{aligned} \mathbb{E}[F(\mathbf{x}_s^{t+1})] &\leq \mathbb{E}[F(\mathbf{x}_s^t)] - \frac{\eta}{2} \mathbb{E}[\|\nabla F(\mathbf{x}_s^t)\|_2^2] + \frac{\eta}{2} \mathbb{E}[\|\nabla F(\mathbf{x}_s^t) - \mathbf{g}_s^t\|_2^2] \\ &\quad - \left(\frac{\eta}{2} - \frac{L\eta^2}{2}\right) \mathbb{E}[\|\mathbf{g}_s^t\|_2^2] \end{aligned}$$

Proof of Theorem 12.3 (cont.)

Sum over $t = 0, \dots, m$ to arrive at

$$\begin{aligned} \mathbb{E}[F(\mathbf{x}_s^{m+1})] &\leq \mathbb{E}[F(\mathbf{x}_s^0)] - \frac{\eta}{2} \sum_{t=0}^m \mathbb{E}[\|\nabla F(\mathbf{x}_s^t)\|_2^2] \\ &\quad + \frac{\eta}{2} \left\{ \sum_{t=0}^m \mathbb{E}[\|\nabla F(\mathbf{x}_s^t) - \mathbf{g}_s^t\|_2^2] - \underbrace{(1 - L\eta)}_{\geq 1/2} \sum_{t=0}^m \mathbb{E}[\|\mathbf{g}_s^t\|_2^2] \right\} \end{aligned}$$

The proof of (12.6) is thus complete if we can justify

Lemma 12.4

If $\eta \leq \frac{1}{L\sqrt{m}}$, then (for fixed η , the epoch length m cannot be too large)

$$\sum_{t=0}^m \underbrace{\mathbb{E}[\|\nabla F(\mathbf{x}_s^t) - \mathbf{g}_s^t\|_2^2]}_{\text{squared bias of gradient estimate}} \leq \frac{1}{2} \sum_{t=0}^m \underbrace{\mathbb{E}[\|\mathbf{g}_s^t\|_2^2]}_{\text{variance}}$$

- informally, this says the accumulated squared bias of gradient estimates (w.r.t. batch gradients) can be controlled by the accumulated variance

Proof of Lemma 12.4

Key step:

Lemma 12.5

$$\mathbb{E} \left[\|\nabla F(\mathbf{x}_s^t) - \mathbf{g}_s^t\|_2^2 \right] \leq \sum_{k=1}^t \mathbb{E} \left[\|\mathbf{g}_s^k - \mathbf{g}_s^{k-1}\|_2^2 \right]$$

- convert the bias of gradient estimates to the differences of consecutive gradient estimates (a consequence of the smoothness and the recursive formula of \mathbf{g}_s^t)

Proof of Lemma 12.4 (cont.)

From Lemma 12.5, it suffices to connect $\{\|\mathbf{g}_s^t - \mathbf{g}_s^{t-1}\|_2\}$ with $\{\|\mathbf{g}_s^t\|_2\}$:

$$\begin{aligned}\|\mathbf{g}_s^t - \mathbf{g}_s^{t-1}\|_2^2 &\stackrel{(12.5)}{=} \|\nabla f_{i_t}(\mathbf{x}_s^t) - \nabla f_{i_t}(\mathbf{x}_s^{t-1})\|_2^2 \stackrel{\text{smoothness}}{\leq} L^2 \|\mathbf{x}_s^t - \mathbf{x}_s^{t-1}\|_2^2 \\ &= \eta^2 L^2 \|\mathbf{g}_s^{t-1}\|_2^2\end{aligned}$$

Invoking Lemma 12.5 then gives

$$\mathbb{E} \left[\|\nabla F(\mathbf{x}_s^t) - \mathbf{g}_s^t\|_2^2 \right] \leq \sum_{k=1}^t \mathbb{E} \left[\|\mathbf{g}_s^k - \mathbf{g}_s^{k-1}\|_2^2 \right] \leq \eta^2 L^2 \sum_{k=1}^t \mathbb{E} \left[\|\mathbf{g}_s^{k-1}\|_2^2 \right]$$

Summing over $t = 0, \dots, m$, we obtain

$$\sum_{t=0}^m \mathbb{E} \left[\|\nabla F(\mathbf{x}_s^t) - \mathbf{g}_s^t\|_2^2 \right] \leq \eta^2 L^2 m \sum_{t=0}^{m-1} \mathbb{E} \left[\|\mathbf{g}_s^t\|_2^2 \right]$$

which establishes Lemma 12.4 if $\eta \lesssim \frac{1}{L\sqrt{m}}$

Proof of Lemma 12.5

Since this lemma only concerns a single epoch, we shall drop the dependency on s for simplicity. Let \mathcal{F}_k contain all info up to \mathbf{x}^k and \mathbf{g}^{k-1} , then

$$\begin{aligned} & \mathbb{E} \left[\left\| \nabla F(\mathbf{x}^k) - \mathbf{g}^k \right\|_2^2 \mid \mathcal{F}_k \right] \\ &= \mathbb{E} \left[\left\| \nabla F(\mathbf{x}^{k-1}) - \mathbf{g}^{k-1} + (\nabla F(\mathbf{x}^k) - \nabla F(\mathbf{x}^{k-1})) - (\mathbf{g}^k - \mathbf{g}^{k-1}) \right\|_2^2 \mid \mathcal{F}_k \right] \\ &= \left\| \nabla F(\mathbf{x}^{k-1}) - \mathbf{g}^{k-1} \right\|_2^2 + \left\| \nabla F(\mathbf{x}^k) - \nabla F(\mathbf{x}^{k-1}) \right\|_2^2 + \mathbb{E} \left[\left\| \mathbf{g}^k - \mathbf{g}^{k-1} \right\|_2^2 \mid \mathcal{F}_k \right] \\ &\quad + 2 \langle \nabla F(\mathbf{x}^{k-1}) - \mathbf{g}^{k-1}, \nabla F(\mathbf{x}^k) - \nabla F(\mathbf{x}^{k-1}) \rangle \\ &\quad - 2 \langle \nabla F(\mathbf{x}^{k-1}) - \mathbf{g}^{k-1}, \mathbb{E}[\mathbf{g}^k - \mathbf{g}^{k-1} \mid \mathcal{F}_k] \rangle \\ &\quad - 2 \langle \nabla F(\mathbf{x}^k) - \nabla F(\mathbf{x}^{k-1}), \mathbb{E}[\mathbf{g}^k - \mathbf{g}^{k-1} \mid \mathcal{F}_k] \rangle \\ &\stackrel{\text{(exercise)}}{=} \left\| \nabla F(\mathbf{x}^{k-1}) - \mathbf{g}^{k-1} \right\|_2^2 - \left\| \nabla F(\mathbf{x}^k) - \nabla F(\mathbf{x}^{k-1}) \right\|_2^2 + \mathbb{E} \left[\left\| \mathbf{g}^k - \mathbf{g}^{k-1} \right\|_2^2 \mid \mathcal{F}_k \right] \end{aligned}$$

Since $\nabla F(\mathbf{x}^0) = \mathbf{g}^0$. Sum over $k = 1, \dots, t$ to obtain

$$\mathbb{E} \left[\left\| \nabla F(\mathbf{x}^k) - \mathbf{g}^k \right\|_2^2 \right] = \sum_{k=1}^t \mathbb{E} \left[\left\| \mathbf{g}^k - \mathbf{g}^{k-1} \right\|_2^2 \right] - \underbrace{\sum_{k=1}^t \left\| \nabla F(\mathbf{x}^k) - \nabla F(\mathbf{x}^{k-1}) \right\|_2^2}_{\leq 0; \text{ done!}}$$

Stochastic dual coordinate ascent (SDCA)

— *a dual perspective*

A class of finite-sum optimization

$$\text{minimize}_{\mathbf{x} \in \mathbb{R}^d} \quad F(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{x}) + \frac{\mu}{2} \|\mathbf{x}\|_2^2 \quad (12.8)$$

- f_i : convex and L -smooth

Dual formulation

The dual problem of (12.8)

$$\text{maximize}_{\boldsymbol{\nu}} \quad D(\boldsymbol{\nu}) = \frac{1}{n} \sum_{i=1}^n -f_i^*(-\boldsymbol{\nu}_i) - \frac{\mu}{2} \left\| \frac{1}{\mu n} \sum_{i=1}^n \boldsymbol{\nu}_i \right\|_2^2 \quad (12.9)$$

- a primal-dual relation

$$\boldsymbol{x}(\boldsymbol{\nu}) = \frac{1}{\mu n} \sum_{i=1}^n \boldsymbol{\nu}_i \quad (12.10)$$

Derivation of the dual formulation

$$\min_{\mathbf{x}} \quad \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{x}) + \frac{\mu}{2} \|\mathbf{x}\|_2^2$$

$$\iff \min_{\mathbf{x}, \{\mathbf{z}_i\}} \quad \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{z}_i) + \frac{\mu}{2} \|\mathbf{x}\|_2^2 \quad \text{s.t. } \mathbf{z}_i = \mathbf{x}$$

$$\iff \max_{\{\boldsymbol{\nu}_i\}} \min_{\mathbf{x}, \{\mathbf{z}_i\}} \quad \underbrace{\frac{1}{n} \sum_{i=1}^n f_i(\mathbf{z}_i) + \frac{\mu}{2} \|\mathbf{x}\|_2^2 + \frac{1}{n} \sum_{i=1}^n \langle \boldsymbol{\nu}_i, \mathbf{z}_i - \mathbf{x} \rangle}_{\text{Lagrangian}}$$

$$\iff \max_{\{\boldsymbol{\nu}_i\}} \min_{\mathbf{x}} \quad \underbrace{\frac{1}{n} \sum_{i=1}^n -f_i^*(-\boldsymbol{\nu}_i)}_{\text{conjugate: } f_i^*(\boldsymbol{\nu}) := \max_{\mathbf{z}} \{\langle \boldsymbol{\nu}, \mathbf{z} \rangle - f_i(\mathbf{z})\}} + \frac{\mu}{2} \|\mathbf{x}\|_2^2 - \frac{1}{n} \sum_{i=1}^n \langle \boldsymbol{\nu}_i, \mathbf{x} \rangle$$

$$\iff \max_{\{\boldsymbol{\nu}_i\}} \quad \frac{1}{n} \sum_{i=1}^n -f_i^*(-\boldsymbol{\nu}_i) - \frac{\mu}{2} \left\| \underbrace{\frac{1}{\mu n} \sum_{i=1}^n \boldsymbol{\nu}_i}_{\text{optimal } \mathbf{x} = \frac{1}{\mu n} \sum_i \boldsymbol{\nu}_i} \right\|_2^2$$

Randomized coordinate ascent on dual problem

— Shalev-Shwartz, Zhang '13

- **randomized coordinate ascent:** at each iteration, randomly pick one **dual** (block) coordinate ν_{i_t} of (12.9) to optimize
- **maintain the primal-dual relation** (12.10)

$$\mathbf{x}^t = \frac{1}{\mu n} \sum_{i=1}^n \nu_i^t \quad (12.11)$$

A variant of SDCA without duality

SDCA might not be applicable if the conjugate functions are difficult to evaluate

This calls for a dual-free version of SDCA

A variant of SDCA without duality

— S. Shalev-Shwartz '16

Algorithm 12.6 SDCA without duality

- 1: initialize $\mathbf{x}^0 = \frac{1}{\mu n} \sum_{i=1}^n \mathbf{v}_i^0$
 - 2: for $t = 0, 1, \dots$ do
 - 3: // choose a random coordinate to optimize
 - 4: choose i_t uniformly from $\{1, \dots, n\}$
 - 5: $\Delta^t \leftarrow -\eta \mu n (\nabla f_{i_t}(\mathbf{x}^t) + \mathbf{v}_{i_t}^t)$
 - 6: $\mathbf{v}_i^{t+1} \leftarrow \underbrace{\mathbf{v}_i^t + \Delta^t \mathbb{1}\{i = i_t\}}_{\text{update only the } i_t^{\text{th}} \text{ coordinate}} \quad (1 \leq i \leq n)$
 - 7: $\mathbf{x}^{t+1} \leftarrow \mathbf{x}^t + \frac{1}{\mu n} \Delta^t$ // based on (12.11)
-

A variant of SDCA without duality

A little intuition

- the optimality condition requires (check!)

$$\boldsymbol{\nu}_i^* = -\nabla f_i(\mathbf{x}^*), \quad \forall i \quad (12.12)$$

- with a modified update rule, one has

$$\boldsymbol{\nu}_{i_t}^{t+1} \leftarrow \underbrace{(1 - \eta\mu n)\boldsymbol{\nu}_{i_t}^t + \eta\mu n(-\nabla f_{i_t}(\mathbf{x}^t))}_{\text{cvx combination of current dual iterate and gradient component}}$$

— when it converges, it will satisfy (12.12)

SDCA as SGD

The SDCA (without duality) update rule reads:

$$\mathbf{x}^{t+1} = \mathbf{x}^t - \eta \underbrace{(\nabla f_{i_t}(\mathbf{x}^t) + \boldsymbol{\nu}_{i_t}^t)}_{:=\mathbf{g}^t}$$

It is straightforward to verify that \mathbf{g}^t is an **unbiased gradient estimate**

$$\mathbb{E}[\mathbf{g}^t] = \mathbb{E}[\nabla f_{i_t}(\mathbf{x}^t)] + \mathbb{E}[\boldsymbol{\nu}_{i_t}^t] = \frac{1}{n} \sum_{i=1}^n \nabla f_i(\mathbf{x}^t) + \underbrace{\frac{1}{n} \sum_{i=1}^n \boldsymbol{\nu}_i^t}_{=\boldsymbol{\mu}\mathbf{x}^t} = \nabla F(\mathbf{x}^t)$$

SDCA as variance-reduced SGD

The SDCA (without duality) update rule reads:

$$\mathbf{x}^{t+1} = \mathbf{x}^t - \eta \underbrace{(\nabla f_{i_t}(\mathbf{x}^t) + \boldsymbol{\nu}_{i_t}^t)}_{:= \mathbf{g}^t}$$

The variance of $\|\mathbf{g}^t\|_2$ goes to 0 as we converge to the optimizer

$$\begin{aligned} \mathbb{E}[\|\mathbf{g}^t\|_2^2] &= \mathbb{E}[\|\boldsymbol{\nu}_{i_t}^t - \boldsymbol{\nu}_{i_t}^* + \boldsymbol{\nu}_{i_t}^* + \nabla f_{i_t}(\mathbf{x}^t)\|_2^2] \\ &\leq 2 \underbrace{\mathbb{E}[\|\boldsymbol{\nu}_{i_t}^t - \boldsymbol{\nu}_{i_t}^*\|_2^2]}_{\rightarrow 0 \text{ as } t \rightarrow \infty} + 2 \underbrace{\mathbb{E}[\|\boldsymbol{\nu}_{i_t}^* + \nabla f_{i_t}(\mathbf{x}^t)\|_2^2]}_{\leq \|\mathbf{w}^t - \mathbf{w}^*\|_2^2 \text{ (Shalev-Shwartz '16)}} \end{aligned}$$

Convergence guarantees of SDCA

Theorem 12.6 (informal, Shalev-Shwartz '16)

Assume each f_i is convex and L -smooth, and set $\eta = \frac{1}{L+\mu n}$. Then it takes SDCA (without duality) $O\left((n + \frac{L}{\mu}) \log \frac{1}{\varepsilon}\right)$ iterations to yield ε accuracy

- the same computational complexity as SVRG
- storage complexity: $O(nd)$ (needs to store $\{\nu_i\}_{1 \leq i \leq n}$)

Reference

- [1] "*Recent advances in stochastic convex and non-convex optimization*," Z. Allen-Zhu, *ICML Tutorial*, 2017.
- [2] "*Accelerating stochastic gradient descent using predictive variance reduction*," R. Johnson, T. Zhang, *NIPS*, 2013.
- [3] "*Barzilai-Borwein step size for stochastic gradient descent*," C. Tan, S. Ma, Y.H. Dai, Y. Qian, *NIPS*, 2016.
- [4] "*A proximal stochastic gradient method with progressive variance reduction*," L. Xiao, T. Zhang, *SIAM Journal on Optimization*, 2014.
- [5] "*Minimizing finite sums with the stochastic average gradient*," M. Schmidt, N. Le Roux, F. Bach, *Mathematical Programming*, 2013.
- [6] "*SAGA: A fast incremental gradient method with support for non-strongly convex composite objectives*," A. Defazio, F. Bach, and S. Lacoste-Julien, *NIPS*, 2014.

Reference

- [7] "Variance reduction for faster non-convex optimization," Z. Allen-Zhu, E. Hazan, *ICML*, 2016.
- [8] "Katyusha: The first direct acceleration of stochastic gradient methods," Z. Allen-Zhu, *STOC*, 2017.
- [9] "SARAH: A novel method for machine learning problems using stochastic recursive gradient," L. Nguyen, J. Liu, K. Scheinberg, M. Takac, *ICML*, 2017.
- [10] "Spider: Near-optimal non-convex optimization via stochastic path-integrated differential estimator," C. Fang, C. Li, Z. Lin, T. Zhang, *NIPS*, 2018.
- [11] "SpiderBoost and momentum: Faster variance reduction algorithms," Z. Wang, K. Ji, Y. Zhou, Y. Liang, V. Tarokh, *NIPS*, 2019.

Reference

- [12] "*Optimal finite-Sum smooth non-convex optimization with SARAH*," L. Nguyen, M. vanDijk, D. Phan, P. Nguyen, T. Weng, J. Kalagnanam, arXiv:1901.07648, 2019.
- [13] "*Stochastic dual coordinate ascent methods for regularized loss minimization*," S. Shalev-Shwartz, T. Zhang, *Journal of Machine Learning Research*, 2013.
- [14] "*SDCA without duality, regularization, and individual convexity*," S. Shalev-Shwartz, ICML, 2016.
- [15] "*Optimization methods for large-scale machine learning*," L. Bottou, F. Curtis, J. Nocedal, 2016.