

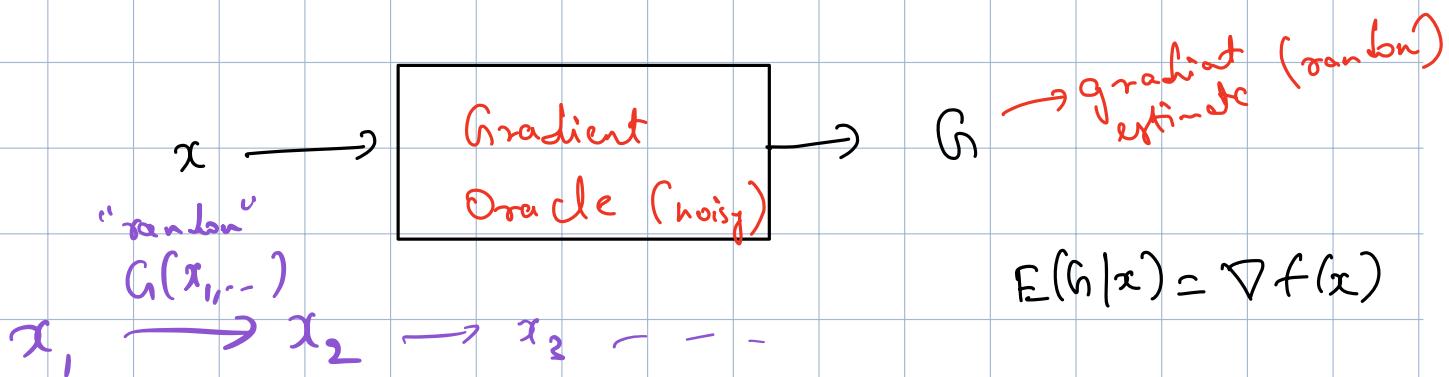
Lecture - 15 (contd)

Non-asymptotic analysis of stochastic gradient algorithms

Aim: Solve $\min_x f(x)$

\curvearrowright smooth

Setting: Unbiased gradient information



SG algorithm: $x_{k+1} = x_k - \alpha(k) G(x_k, \xi_k)$ ①

\downarrow
step-size or learning rate

\uparrow noise

Assumption:

$$(A) \quad E_{\xi_k} (G(x_k, \xi_k)) = \nabla f(x_k), \forall k \geq 1$$

$$E_{\xi_k} \|G(x_k, \xi_k) - \nabla f(x_k)\|^2 \leq \sigma^2$$

for some $\sigma > 0$.

(A0) f is L -smooth

$$\|\nabla f(x_1) - \nabla f(x_2)\| \leq L \|x_1 - x_2\|, \forall x_1, x_2 \in \mathbb{R}^d$$

Asymptotic Convergence:

(A2)

$$\sum_k \alpha(k) = \infty, \quad \sum_k \alpha(k)^2 < \infty$$

e.g. $\alpha(k) = \frac{c}{k} \rightarrow$ standard stochastic approximation conditions.

Under (A0) - (A2), the parameter x_k governed

by ① converges almost surely to the set $\{x^* \mid \nabla f(x^*) = 0\}$

Asymptotic guarantee: $x_k \rightarrow \{x^* \mid \nabla f(x^*) = 0\}$ as $k \rightarrow \infty$ w.p. 1.

Tool used for proving this claim: Kushner-Clark Lemma

Eq ① (\Rightarrow)

$$x_{k+1} = x_k - \alpha(k) (\nabla f(x_k) + M_{k+1})$$

$$M_{k+1} = G(x_k, \xi_k) - \nabla f(x_k)$$

("MDS")

We want to understand the non-asymptotic performance of the Sh algorithm above.

Stationary point, say x^* , has $\nabla f(x^*) = 0$

ϵ -stationary point: \bar{x} is an ϵ -stationary point if $E \|\nabla f(\bar{x})\| \leq \epsilon$, where the expectation is over the randomness in the SG algorithm.

[Ghadimi-Lan, 2014, SIAM J. Opt, "Stochastic first/zeroth order"]

"Randomized Stochastic gradient" (RSG)

Suppose we run ① for N iterations.

$\{x_1, x_2, x_3, \dots, x_N\} \xrightarrow{\text{(random)}} \text{Set of}$
iterations visited by SG algorithm.

RSG will pick a random iterate as follows:

r is picked uniformly at random from $\{x_1, \dots, x_N\}$

i.e., $P(x_R = x_i) = \frac{1}{N}$ $\forall i$

$$E x_R = \frac{1}{N} \sum_{i=1}^N x_i \quad \begin{cases} \text{average} \\ \text{iterate.} \end{cases}$$

The non-asymptotic bound for RSG is of the form

$$E \|\nabla f(x_R)\|^2 \leq \frac{\text{const}}{\sqrt{N}}$$

under suitable choice of step-size.

FL connection: $x \rightarrow$ policy parameter

" " Objective $f \rightarrow J_\pi(s^0)$ or $J(x)$

"Policy gradient"
 Value function
 with state s^0
 in a discounted/SSP

Proof of the non-asymptotic bound for RSG:

First \rightarrow derive a bound for a general step-size

Second \rightarrow Specialize the bound above with a particular stepsize.

First step: Recall $x_{k+1} = x_k - \alpha(k) \mathcal{G}(x_k, \xi_k)$

f is L -smooth:

$$f(x_{k+1}) \leq f(x_k) + \nabla f(x_k)^T (x_{k+1} - x_k) + \frac{L}{2} \|x_{k+1} - x_k\|^2$$

$$= f(x_k) - \alpha(k) \nabla f(x_k)^T \mathcal{G}(x_k, \xi_k) + \frac{L}{2} \alpha(k)^2 \|\mathcal{G}(x_k, \xi_k)\|^2 \quad \text{--- (1)}$$

Taking expectation wrt \mathcal{F}_k , denote this by E_k $\xrightarrow{\sigma\text{-field with info up to } k}$

$$E_k f(x_{k+1}) \leq f(x_k) - \alpha(k) \|\nabla f(x_k)\|^2 + \frac{L}{2} \alpha(k)^2 (\|\nabla f(x_k)\|^2 + \sigma^2) \quad \text{--- (2)}$$

(A)

$$= f(x_k) - \left(\alpha(k) - \frac{L}{2} \alpha(k)^2 \right) \|\nabla f(x_k)\|^2 + \frac{L}{2} \alpha(k)^2 \sigma^2$$

$$\left(\alpha(k) - \frac{L}{2} \alpha(k)^2 \right) \|\nabla f(x_k)\|^2 \leq f(x_k) - E_k(f(x_{k+1})) + \frac{L}{2} \alpha(k)^2 \sigma^2$$

(X)

$$\alpha(k) \|\nabla f(x_k)\|^2 \leq \frac{2(f(x_k) - E_k(f(x_{k+1})))}{2 - L\alpha(k)} + \frac{L\alpha(k)^2 \sigma^2}{(2 - L\alpha(k))}$$

\hookrightarrow assuming $a(k) \leq \frac{1}{L} \forall k$

Summing \otimes over $k=1 \text{ to } N$,

$$\sum_{k=1}^N a(k) \|\nabla f(x_k)\|^2 \leq 2 \sum_{k=1}^N \frac{f(x_k) - E_k f(x_{k+1})}{2 - L a(k)} + L \sigma^2 \sum_{k=1}^N \frac{a(k)^2}{(2 - L a(k))}$$

(3)

Take expectation wrt. $\{x_1, \dots, x_N\}$, i.e., $E_N(\cdot)$

$$\sum_{k=1}^N a(k) E_N \|\nabla f(x_k)\|^2 \leq 2 \sum_{k=1}^N \frac{E_N f(x_k) - E_N f(x_{k+1})}{2 - L a(k)} + L \sigma^2 \sum_{k=1}^N \frac{a(k)^2}{(2 - L a(k))}$$

"End of lecture 5"

Set $a(k) = a \quad \forall k=1 \dots N$

$$a \sum_{k=1}^N E_N \|\nabla f(x_k)\|^2 \leq \frac{2}{2 - La} (f(x_1) - f(x^*)) + \frac{L \sigma^2 N a^2}{2 - La}$$

↑ global optima of f

Choose $a \leq \frac{1}{L}$

$\hookrightarrow E_N f(x_{N+1}) \geq f(x^*)$

$$E_{x_k} \|\nabla f(x_k)\|^2 = \frac{1}{N} \sum_{k=1}^N E_x \|\nabla f(x_k)\|^2$$

So,

$$\begin{aligned}
 E_{x_k} \|\nabla f(x_k)\|^2 &\leq \frac{1}{Na} \left[\frac{2}{2-L\alpha} (f(x_1) - f(x^*)) + \frac{L\sigma^2 N \alpha^2}{2-L\alpha} \right] \\
 &\leq \frac{2}{Na} (f(x_1) - f(x^*)) \\
 &\quad + \frac{L\sigma^2 N \alpha^2}{Na} \\
 &= \frac{2}{Na} (f(x_1) - f(x^*)) + \underbrace{L\sigma^2 \alpha}_{\text{Variance due to gradient estimation}}
 \end{aligned}$$

Initial error (bias)
 Sampling error (variance)

Initial error

We have,

$$E \|\nabla f(x_k)\|^2 \leq \frac{\text{const}}{Na} + \overline{\text{const}} \alpha$$

Setting $\alpha = \min \left\{ \frac{1}{L}, \frac{1}{\sqrt{N}} \right\}$

$x_{k+1} = x_k - \alpha_k g(x_{k+1})$
 $\alpha_k = \frac{1}{k} \rightarrow \text{diminishing}$
 $\alpha_k = \alpha \rightarrow \text{constant}$
 Iteration, $\alpha = \frac{1}{N}$

$$E \|\nabla f(x_k)\|^2 \leq \frac{2}{N} (f(x_1) - f(x^*)) \max(L, \sqrt{N})$$

$$+ \frac{L\sigma^2}{\sqrt{N}}$$

$$= \frac{2L(f(x_1) - f(x^*))}{N}$$

$$+ \frac{1}{\sqrt{N}} \left(2(f(x_1) - f(x^*)) + L\sigma^2 \right)$$

$$E \|\nabla f(x_p)\|^2 \leq \frac{\text{const}}{\sqrt{N}}$$

→ Bound for Smooth + Unbiased gradient oracle.

(\Leftarrow) For RSG to converge to an ϵ -stationary

point, i.e., $E \|\nabla f(x_p)\|^2 \leq \epsilon$,

in order $\frac{1}{\epsilon^2}$ number of iterations

of RSG algorithm are sufficient.

Stochastic optimization with a biased gradient oracle

(Regular) stochastic optimization

$$\min_{x \in S} \{f(x) = E_{\xi} (F(x, \xi))\}$$

Zeroth-order
input



"gradient" not directly available

Stochastic gradient algorithm:

$$x_{k+1} = \Pi(x_k - \gamma_k g_k)$$

↗
 projection operator ↗ stepsize ↗ gradient estimate

With zeroth-order input, can form gradient estimates that satisfy

$$\|E_{\xi} [g(x, \xi)] - \nabla f(x)\| \leq c_1 \delta^2$$

Bias

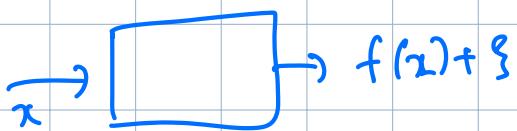
$$E_{\xi} \|g(x, \xi) - E_{\xi} (g(x, \xi))\|^2 \leq \frac{c_2}{\delta^2}$$

Variance

" δ " → bias-variance tradeoff.

How to form such an estimate? "Simultaneous perturbation" trick.

Here is a sample for $f \in \mathcal{L}^3$



Take two measurements : ① $f(x + \delta \Delta) + \xi^+$
② $f(x - \delta \Delta) - \xi^-$

$$\Delta = (\Delta_1, \dots, \Delta_d)$$

$$\Delta_i \text{ Rademacher} = \begin{cases} +1 & \text{w.p. } \frac{1}{2} \\ -1 & \text{w.p. } \frac{1}{2} \end{cases} \text{ iid}$$

Gradient estimate $g^i = \frac{\textcircled{1} - \textcircled{2}}{2 \delta \Delta^i}$

Use Taylor expansion

$$f(x \pm \delta \Delta) = f(x) \pm \delta \Delta^T \nabla f(x) + \frac{\delta^2}{2} \Delta^T \nabla^2 f(x) \Delta + O(\delta^3)$$

$$\frac{f(x + \delta \Delta) - f(x - \delta \Delta)}{2 \delta \Delta^i} = \nabla_i f(x) + \sum_{j \neq i} \frac{\Delta_j}{\Delta_i} \nabla_j f(x) + O(\delta^2)$$

↑
zero-mean

$$E(g^i) = \nabla_i f(x) + O(\delta^2)$$

Can show $E(\|g - E g\|^2) \leq \frac{C \delta^2}{2}$

D=Rademacher vector \Rightarrow SPSA "Spall 1992"

Can we other distributions for random perturbations.

Alternate framework: "Biased function measurement".

$$f(x) = \min_y E_{\xi} [H_x(y, \xi)] \quad \text{--- (x)}$$

↑
The objective at a certain input is the solution to an optimization problem

Suppose a batch-size parameter "m" decides how well (x) is solved, i.e.,

$$F(x, m) = \min_y E_{\xi} [H_x(y, \xi)] + \underbrace{\epsilon(m)}_{\text{error term}}$$

that has positive mean

"increasing m lowers $\epsilon(m)$ ".

e.g. Risk-sensitive optimization ϵ : risk-measure $C_x = \epsilon(Y_x)$

$\hat{\epsilon}_{m,x}$ → estimate of $\epsilon(Y_x)$ using m samples

$E(\hat{\epsilon}_{m,x}) \neq \epsilon(Y_x)$, but

$$E|\hat{\epsilon}_{m,x} - \epsilon(Y_x)|^2 \leq \frac{\text{const.}}{\sqrt{m}} \quad \forall x$$

With biased function measurements, w.r.t. the simultaneous perturbation trick lead to



Fact: $\| E_{\xi} [g(x, \xi_m)] - \nabla f(x) \| \leq C_1 \delta^2 + \frac{C_2}{\delta \sqrt{m}}$

Variance: $E_{\xi} \| g - E_g \|^2 \leq \frac{C_3}{\delta^2}$

Algorithm: Randomized Stochastic gradient
[Hadrini-Lan, SIOpt 2015]

Perform $x_{k+1} = T(x_k - \alpha_k g(x_k, \xi_{k,m_k}))$

for N iterations

& return a random iterate R picked uniformly at random from $\{x_1, \dots, x_N\}$

Assumption (A1) $\|\nabla f(x) - \nabla f(y)\| \leq L \|x-y\| \quad \forall x, y \in X$.

Assumption (A2) $\|\nabla f(x)\| \leq B < \infty, \forall x \in X$

Main claim: With $\alpha_k = \min \left\{ \frac{1}{L}, \frac{1}{N^2 \delta_k^2} \right\}$, $\delta_k = \frac{1}{N^{1/6}}$,

$$m_k = N$$

$$\forall N \geq 1, \quad E \|\nabla f(x_R)\|^2 \leq \frac{2L(f(x_1) - f(x^*))}{N} + \frac{\text{const.}}{N^{4/3}}$$

A biased gradient oracle variant:

Def:

$$\|E_{\xi} [g(x, \xi_m)] - \nabla f(x)\|_{\infty} \leq C_1 \delta^2 + \frac{C_2}{\delta \sqrt{m}}$$

Variance:

$$E_{\xi} \|g - E_g\|^2 \leq C_3 \delta + \tilde{C}_3$$

In this oracle model,
we can obtain
the rate
 $O(\sqrt{\delta})$

Motivation! Biased function measurements in a "noise-less" setting

Recall

$$g^i = \frac{(f(x + \delta \Delta) + \xi^+) - f(x - \delta \Delta) + \xi^-)}{2 \delta^i \delta}$$

$$E \|g - E_g\|^2 \leq \frac{\text{Const}}{\delta^2} \quad \text{because of noise elements } \xi^+, \xi^-$$

$$\text{If no noise, then } E \|g - E_g\|^2 \leq C_3 \delta^2 + \tilde{C}_3$$

Now: Use a Gaussian-smoothing estimator.

[Nesterov - Spokoiny, 2017]

$$g = \Delta \left(\frac{f(x + \delta \Delta) - f(x - \delta \Delta)}{2 \delta} \right)$$

$$\Delta = (\Delta^1, \dots, \Delta^d) \quad \Delta^i \sim N(0, 1)$$

$$x_{k+1} = x_k - \gamma_k g(x_k, \xi_k)$$

Main claim:

With $\gamma_k = \min \left\{ \frac{1}{L}, \frac{1}{\sqrt{N}} \right\}$, $\delta_k = \frac{1}{\sqrt{N}}$,

$$m_p = N$$

$$\forall N \geq 1, \quad E \|\nabla f(x_p)\|^2 \leq \frac{L(f(x_0) - f(x^*))}{N} + \frac{\text{const.}}{\sqrt{N}}$$

For an oracle without

bias-variance tradeoff

\rightarrow End of Lecture-16 \times

Lecture 17 Non-asymptotic bounds for zeroth-order optimization

(0) Input: x

Output: $g(x, \xi)$

$$(i) E_\xi(g(x, \xi)) = \nabla f(x) + c_1 \delta^2 \mathbf{1}_{d \times 1}$$

$$(ii) E_\xi[\|g(x, \xi) - E_\xi(g(x, \xi))\|^2] \leq \frac{c_2}{\delta^2}$$

Oracle with
variance
trick

Algorithm: Randomized Stochastic gradient

[Ghadimi-Lan, SiOpt 2015]

Perform $x_{k+1} = x_k - \gamma_k g(x_k, \xi_k)$

for N iterations

& return a random iterate R picked uniformly
at random from $\{x_1, \dots, x_N\}$

Assumption (A1) $\|\nabla f(x) - \nabla f(y)\| \leq L \|x-y\| \quad \forall x, y \in \mathbb{R}^d.$

Main claim: With $\gamma_k = \min \left\{ \frac{1}{L}, \frac{1}{(2^k N)^{1/3}} \right\}$, $\delta_k = \frac{1}{(2^k N)^{1/6}}$,

$$\forall N \geq 1, \quad E \|\nabla f(x_R)\|^2 \leq \frac{2L(f(x_1) - f(x^*))}{N} + \frac{\text{const.}}{N^{1/3}}$$

Pf: We prove a result for general γ_k & then specialize

With $P_R(k) = \text{Prob}(R=k) = \frac{\gamma_k}{\sum_{i=1}^N \gamma_i}$,

$$E \|\nabla f(x_R)\|^2 \leq \frac{1}{\sum_{i=1}^N \gamma_i} \left[\frac{2(f(x_1) - f(x^*))}{(2-L\gamma_1)} + 2B \sum_{k=1}^N c_1 \delta_k^2 \left(\frac{\gamma_k + L\gamma_1^2}{2-L\gamma_k} \right) \right]$$

$$+ L \sum_{k=1}^N \frac{\gamma_k^2}{(2-L\gamma_k)} \left[2c_1^2 \delta_k^4 + \frac{c_2}{\delta_k^2} \right]$$

Pf: Since f is L -smooth,

$$\begin{aligned} f(x_{k+1}) &\leq f(x_k) + \langle \nabla f(x_k), x_{k+1} - x_k \rangle + \frac{L}{2} \|x_{k+1} - x_k\|^2 \\ &= f(x_k) - \gamma_k \langle \nabla f(x_k), g(x_k, \xi_k) \rangle + \frac{L\gamma_k^2}{2} \|g(x_k, \xi_k)\|^2 \end{aligned}$$

$$E_k(f(x_{k+1}))$$

Conditional expectation
i.e., w.r.t info upto time k

$$\begin{aligned} &\leq f(x_k) - \gamma_k \langle \nabla f(x_k), \nabla f(x_k) + c_1 \delta_k^2 I_{d \times d} \rangle \\ &\quad + \frac{L}{2} \gamma_k^2 \left[\|E_k(g(x_k, \xi_k))\|^2 + \frac{c_2}{\delta_k^2} \right] \end{aligned}$$

$$\text{and } E_k g = \nabla f + \delta^2 I_d$$

$$\leq f(x_k) - \gamma_k \|\nabla f(x_k)\|^2 + \gamma_k c_1 \delta_k^2 \|\nabla f(x_k)\|,$$

$$+ \frac{L}{2} \gamma_k^2 \left[\|\nabla f(x_k)\|^2 + 2c_1 \delta_k^2 \|\nabla f(x_k)\| + 2c_1^2 \delta_k^4 + \frac{c_2}{\delta_k^2} \right]$$

Upper bound
on $\|\nabla f(x_k)\|$

$$\leq f(x_k) - \left(\gamma_k - \frac{L}{2} \gamma_k^2 \right) \|\nabla f(x_k)\|^2 + c_1 \delta_k^2 (\gamma_k + L\gamma_k^2) \beta$$

$$+ \frac{L}{2} \gamma_k^2 \left[2c_1^2 \delta_k^4 + \frac{c_2}{\delta_k^2} \right]$$

Re-arranging,

$$\left(\gamma_k - \frac{L}{2} \gamma_k^2 \right) \|\nabla f(x_k)\|^2 \leq f(x_k) - E_k f(x_{k+1}) + c_1 \delta_k^2 (\gamma_k + L\gamma_k^2) \beta$$

$$+ \frac{L}{2} \gamma_k^2 \left[2c_1^2 \delta_k^4 + \frac{c_2}{\delta_k^2} \right]$$

Assuming $\gamma_k \leq \frac{1}{L}$

$$\gamma_k \|\nabla f(x_k)\|^2 \leq \frac{2}{(2-L\gamma_k)} \left[f(x_k) - E_k(f(x_{k+1})) + C_1 \delta_k^2 (\gamma_k + L\gamma_k^2) B \right] + \frac{L\gamma_k^2}{2-L\gamma_k} \left[2C_1^2 \delta_k^4 + \frac{C_2}{\delta_k^2} \right]$$

Summing & taking expectation,

$$\sum_{k=1}^N \gamma_k E_N \|\nabla f(x_k)\|^2 \leq 2 \sum_{k=1}^N \frac{E_N(f(x_k)) - E_N(f(x_{k+1}))}{(2-L\gamma_k)} + 2 \sum_{k=1}^N C_1 \delta_k^2 \left(\frac{\gamma_k + L\gamma_k^2}{2-L\gamma_k} \right) B + L \sum_{k=1}^N \frac{\gamma_k^2}{(2-L\gamma_k)} \left(2C_1^2 \delta_k^4 + \frac{C_2}{\delta_k^2} \right)$$

$$\leq 2 \left[\frac{f(x_1)}{2-L\gamma_1} - \frac{E_N(f(x_{N+1}))}{2-L\gamma_1} \right] + 2 \sum_{k=1}^N C_1 \delta_k^2 \left(\frac{\gamma_k + L\gamma_k^2}{2-L\gamma_k} \right) B + L \sum_{k=1}^N \frac{\gamma_k^2}{(2-L\gamma_k)} \left(2C_1^2 \delta_k^4 + \frac{C_2}{\delta_k^2} \right)$$

$$\leq 2 \left[\frac{f(x_1) - f(x^*)}{(2-L\gamma_1)} \right] + 2 \sum_{k=1}^N C_1 \delta_k^2 \left(\frac{\gamma_k + L\gamma_k^2}{2-L\gamma_k} \right) B + L \sum_{k=1}^N \frac{\gamma_k^2}{(2-L\gamma_k)} \left(2C_1^2 \delta_k^4 + \frac{C_2}{\delta_k^2} \right)$$

$$E \|\nabla f(x_p)\|^2 \leq \frac{1}{\sum_{k=1}^N \gamma_k} \left(\dots \right)$$

Specialize the result above for

$$\gamma_k = \left(\gamma = \min \left\{ \frac{1}{L}, \frac{1}{(\Delta^2 N)^{4/3}} \right\} \right), \quad \delta_k = \left(\delta = \frac{1}{(\Delta^5 N)^{1/3}} \right)$$

$$E \| \nabla f(x_k) \|^2 = \frac{1}{N} \sum_{k=1}^N E \| \nabla f(x_k) \|^2$$

$$\leq \frac{1}{N\gamma} \left[\frac{2(f(x_1) - f(x^*))}{2-L\gamma} + \frac{2BN C_1 \delta^2 (\gamma + L\gamma^2)}{(2-L\gamma)} \right. \\ \left. + \frac{LN\gamma^2}{2-L\gamma} \left(2C_1^2 \delta^4 + \frac{C_2}{\delta^2} \right) \right]$$

$$\text{using } \gamma < \frac{1}{L}$$

$$\leq \boxed{\frac{2(f(x_1) - f(x^*))}{N\gamma} + 4\delta C_1 \delta^2 LN\gamma^2 \left(2C_1^2 \delta^4 + \frac{C_2}{\delta^2} \right)}$$

$$\leq \frac{2(f(x_1) - f(x^*))}{N} \max \left(L, (\Delta^2 N)^{2/3} \right) + \frac{LN C_1}{(\Delta^5 N)^{1/3}} \\ + L \left(\frac{d C_1^2}{(\Delta^5 N)^{2/3}} + \frac{2 d \delta^2}{N^{-1/3}} \right) \frac{1}{(\Delta^2 N)^{2/3}}$$

$$= \frac{2L(f(x_1) - f(x^*))}{N} + \frac{2d^{4/3}}{N^{1/3}} (f(x_1) - f(x^*)) + \frac{4BC_1}{d^{4/3} N^{1/3}}$$

$$+ L \left(\frac{C_1^2}{d^{7/3} N^{2/3}} + \frac{C_2 d^{5/3}}{N^{-1/3}} \right) \frac{1}{(\Delta^2 N)^{2/3}}$$

$$= \frac{2L(f(x_1) - f(x^*))}{N} + \frac{1}{N^{1/3}} \left[2d^{4/3} (f(x_1) - f(x^*)) + \frac{4BC_1}{d^{4/3}} + \frac{LC_1^2}{d^{11/3} N} + \frac{C_2 d^{11/3}}{N} \right]$$

Sampling error

Sampling error

$$E \|\nabla f(x_k)\|^2 \leq \frac{\text{const } d^{4/3}}{N^{1/3}} \leq \epsilon$$

$$\text{if } N \geq \frac{d^4}{\epsilon^3}$$

Extend to biased gradient oracle without
bias-variance tradeoff, i.e., a setting like

$$\|\nabla g - \nabla f\| \leq C_1 \delta$$

$$E \|g - \nabla g\|^2 \leq C_2 \delta + C_3$$

End of Lecture-17

Analysis of stochastic gradient algorithms: Strongly convex case

7.1 SG algorithm and a useful lemma

$$x_{k+1} = x_k - \alpha_k g(x_k, \xi_k). \quad (\text{Stochastic gradient algo}) \quad (7.1)$$

Assumption 7.1. $\mathbb{E}_{\xi_k}[g(x_k, \xi_k)] = \nabla f(x_k)$ and $\mathbb{E}_{\xi_k}[\|g(x_k, \xi_k)\|_2^2] - \|\mathbb{E}_{\xi_k}[g(x_k, \xi_k)]\|_2^2 \leq \sigma^2$.

Assumption 7.2. $f(x) \geq \bar{f}$ for all x

$$\mathbb{E}_{\xi_k} g^2 \leq \sigma^2 + \|\mathbb{E}_g\|^2 \\ = \sigma^2 + \|\nabla f\|^2$$

Lemma 7.1 (Improvement in gradient step). *Under Assumption 7.1, we have*

$$\mathbb{E}_{\xi_k}[f(x_{k+1})] - f(x_k) \leq -\alpha_k \left(1 - \frac{1}{2}\alpha_k L\right) \|\nabla f(x_k)\|_2^2 + \frac{1}{2}\alpha_k^2 L\sigma^2. \quad (7.2)$$

Proof. Using L -smoothness of f and the update iteration (7.1), we obtain

$$f(x_{k+1}) - f(x_k) \leq \nabla f(x_k)^T (x_{k+1} - x_k) + \frac{1}{2}L\|x_{k+1} - x_k\|_2^2 \quad (7.3)$$

$$\leq -\alpha_k \nabla f(x_k)^T g(x_k, \xi_k) + \frac{1}{2}\alpha_k^2 L\|g(x_k, \xi_k)\|_2^2. \quad (7.4)$$

Taking expectation wrt ξ_k , we obtain

$$\mathbb{E}_{\xi_k}[f(x_{k+1})] - f(x_k) \leq -\alpha_k \|\nabla f(x_k)\|_2^2 + \frac{1}{2}\alpha_k^2 L\mathbb{E}_{\xi_k}[\|g(x_k, \xi_k)\|_2^2] \quad (7.5)$$

$$\leq -\alpha_k \left(1 - \frac{1}{2}\alpha_k L\right) \|\nabla f(x_k)\|_2^2 + \frac{1}{2}\alpha_k^2 L\sigma^2. \quad (7.6)$$

The claim follows. \square

7.2 Strongly convex case

Theorem 7.2. Let f be a m -strongly convex function. Then, the SG algorithm governed by (7.1) and with $\alpha_k = \alpha$ s.t. $0 < \alpha L < 1$, satisfies

$$\mathbb{E}[f(x_n) - f(x^*)] \leq \frac{\alpha L \sigma^2}{2m} + (1 - \alpha m)^{n-1} \left(f(x_1) - f(x^*) - \frac{\alpha L \sigma^2}{2m} \right). \quad (7.7)$$

Proof. From Lemma 7.1, we have

$$\mathbb{E}_{\xi_k}[f(x_{k+1})] - f(x_k) \leq -\alpha_k(1 - \frac{1}{2}\alpha_k L)\|\nabla f(x_k)\|_2^2 + \frac{1}{2}\alpha_k^2 L \sigma^2.$$

Since $\alpha_k = \alpha$ and $0 < \alpha L < 1$, we have

$$\mathbb{E}_{\xi_k}[f(x_{k+1})] - f(x_k) \leq -\frac{1}{2}\alpha\|\nabla f(x_k)\|_2^2 + \frac{1}{2}\alpha^2 L \sigma^2. \quad (7.8)$$

Since f is m -strongly convex, the PL-condition holds, i.e.,

$$\boxed{f(x) - f(x^*) \leq \frac{1}{2m}\|\nabla f(x)\|_2^2, \forall x.} \rightarrow \text{PL-convex}$$

Using PL-condition in (7.8), we obtain

$$\mathbb{E}_{\xi_k}[f(x_{k+1})] - f(x_k) \leq -m\alpha(f(x_k) - f(x^*)) + \frac{1}{2}\alpha^2 L \sigma^2. \quad (7.9)$$

Subtracting $f(x^*)$ on both sides and re-arranging, we obtain

$$\mathbb{E}_{\xi_k}[f(x_{k+1}) - f(x^*)] \leq (1 - \alpha m)[f(x_k) - f(x^*)] + \frac{1}{2}\alpha^2 L \sigma^2. \quad (7.10)$$

Taking expectations followed by straightforward simplifications, we obtain

$$\begin{aligned} \mathbb{E}[f(x_{k+1}) - f(x^*)] - \frac{\alpha L \sigma^2}{2m} &\leq (1 - \alpha m)\mathbb{E}[f(x_k) - f(x^*)] + \frac{\alpha^2 L \sigma^2}{2} - \frac{\alpha L \sigma^2}{2m} \\ &= (1 - \alpha m) \left(\mathbb{E}[f(x_k) - f(x^*)] - \frac{\alpha L \sigma^2}{2m} \right). \end{aligned} \quad (7.11)$$

Since $\alpha m < \frac{m}{L} < 1$, a repeated application of the above inequality leads to the following bound:

$$\mathbb{E}[f(x_n) - f(x^*)] \leq \frac{\alpha L \sigma^2}{2m} + (1 - \alpha m)^{n-1} \left(f(x_1) - f(x^*) - \frac{\alpha L \sigma^2}{2m} \right). \quad (7.12)$$

The claim follows. \square

Remark 7.1. Taking limits as $k \rightarrow \infty$ in (7.7), we obtain

$$\mathbb{E}[f(x_k) - f(x^*)] \rightarrow \frac{\alpha L \sigma^2}{2m} \text{ as } k \rightarrow \infty. \quad (7.13)$$

The result above implies that a constant stepsize SG algorithm does not converge to the optima, and instead gets to within a ball around the optima.

Diminishing Stepsize:

From improvement lemma,

$$E_k(f(x_{k+1})) - f(x_k) \leq -\alpha_k \left(1 - \frac{\alpha_k L}{2}\right) \|\nabla f(x_k)\|^2 + \frac{\alpha_k^2 L^2}{2}$$

Suppose $\alpha_k L \leq 1$

$$\leq -\frac{\alpha_k}{2} \|\nabla f(x_k)\|^2 + \frac{\alpha_k^2 L \sigma^2}{2}$$

PL-condition \hookrightarrow

$$\leq -\alpha_k m (f(x_k) - f(x^*)) + \frac{\alpha_k^2 L \sigma^2}{2}$$

$$E_k(f(x_{k+1})) - f(x^*) \leq (1 - \alpha_k m) (f(x_k) - f(x^*)) + \frac{\alpha_k^2 L \sigma^2}{2}$$

Set $\alpha_k = \frac{\beta}{k+1}$, $\beta > \frac{1}{m} \rightarrow$ strong convexity forward

Claim: $E(f(x_n) - f(x^*)) \leq \frac{C}{n+1}$,

where $C = \max \left\{ \frac{\beta^2 L \sigma^2}{2(\beta m - 1)}, 2(f(x_0) - f(x^*)) \right\}$

Proof: by induction

Base Case: holds trivially

Assume claim holds for "n"

Rcf:

Boyd, Nocedal, Wright
"Opt for Large-scale learning,"
Survey, 2018

$$E(f(x_{n+1}) - f(\bar{x})) \leq (1-\alpha_n m) E(f(r_n) - f(\bar{x})) + \frac{\alpha_n^2 L \sigma^2}{2}$$

Induction
hypothesis

$$\leq \left(1 - \frac{\beta_m}{n+1}\right) \frac{c}{n+1} + \frac{\beta^2 L \sigma^2}{2(n+1)^2}$$

$$= \frac{(n+1 - \beta_m)}{(n+1)^2} c + \frac{\beta^2 L \sigma^2}{2(n+1)^2}$$

$$= \frac{n}{(n+1)^2} c - \frac{(\beta_{m-1})}{(n+1)^2} c + \frac{\beta^2 L \sigma^2}{2(n+1)^2}$$

$$\frac{\beta^2 L \sigma^2}{2(n+1)^2} - c \frac{(\beta_{m-1})}{(n+1)^2} \leq 0 \quad \text{since}$$

$$c = \max \left(\frac{\beta^2 L \sigma^2}{2(\beta_{m-1})}, 2\{f(r_m) - f(\bar{x})\} \right)$$

So,

$$E(f(x_{n+1}) - f(\bar{x})) \leq \frac{nc}{(n+1)^2} \leq \frac{c}{n+2}$$

SG for convex functions: unbiased gradients case

Problem: $\min_x f(x)$

SG: $x_{k+1} = x_k - \alpha_k g(x_k, \xi_k)$

Assume: (1) f is L-smooth

(2) **Assumption 7.1.** $\mathbb{E}_{\xi_k}[g(x_k, \xi_k)] = \nabla f(x_k)$ and $\mathbb{E}_{\xi_k}[\|g(x_k, \xi_k)\|_2^2] - \|\mathbb{E}_{\xi_k}[g(x_k, \xi_k)]\|_2^2 \leq \sigma^2$.

$$\|x_{k+1} - x^*\|^2 = \|x_k - x^*\|^2 - 2\alpha_k \langle g(x_k, \xi_k), x_k - x^* \rangle + \gamma_k^2 \|g(x_k, \xi_k)\|^2$$

$$\begin{aligned} \mathbb{E} \|x_{k+1} - x^*\|^2 &\leq \mathbb{E} \|x_k - x^*\|^2 - 2\alpha_k \langle \nabla f(x_k), x_k - x^* \rangle \\ &\quad + \alpha_k^2 (\sigma^2 + \|\nabla f\|^2) \end{aligned}$$

Convex f: $f(x^*) \geq f(x_k) + \langle \nabla f(x_k), x^* - x_k \rangle$

$\|\nabla f(x_k)\|^2 \leq L \langle \nabla f(x_k), x_k - x^* \rangle \rightarrow \text{How? P.W.}$

$$(2\alpha_k - L\alpha_k^2)(\mathbb{E} f(x_k) - f(x^*)) \leq (\mathbb{E} \|x_k - x^*\|^2 - \mathbb{E} \|x_{k+1} - x^*\|^2) + \alpha_k^2 \sigma^2$$

$$\alpha_k (\mathbb{E} f(x_k) - f(x^*)) \leq \frac{1}{(2-L\alpha_k)} \left(\mathbb{E} \|x_k - x^*\|^2 - \mathbb{E} \|x_{k+1} - x^*\|^2 + \alpha_k^2 \sigma^2 \right)$$

$$\sum_{k=1}^N \alpha_k (E(f(x_k)) - f(x^*)) \leq$$

$$\sum_{k=1}^N \frac{1}{(2-L\alpha_k)} \left(E \|x_k - x^*\|^2 - E \|x_{k+1} - x^*\|^2 \right) + \sum_{k=1}^N \frac{\alpha_k^2 \sigma^2}{(2-L\alpha_k)}$$

Set $L\alpha_k \leq 1$

$$\sum_{k=1}^N \alpha_k (E(f(x_k)) - f(x^*))$$

$$\leq \|x_1 - x^*\|^2 - E \|x_{N+1} - x^*\|^2 - \left(\sum_{k=2}^N \left(\frac{1}{2-L\alpha_{k-1}} - \frac{1}{2-L\alpha_k} \right) E \|x_k - x^*\|^2 \right)$$

Can be ignored

$$+ \sum_{k=1}^N \alpha_k^2 \sigma^2$$

$$E f(x_R) - f(x^*) \leq \frac{1}{\sum \alpha_k} \left(\|x_1 - x^*\|^2 + \sum \alpha_k^2 \sigma^2 \right)$$

$$\alpha_k = \alpha$$

$$E f(x_R) - f(x^*) \leq \frac{1}{N\alpha} \left(\|x_1 - x^*\|^2 + N\alpha^2 \sigma^2 \right)$$

$$= \frac{\|x_1 - x^*\|^2}{N\alpha} + \alpha \sigma^2$$

$$\alpha = \frac{1}{\sqrt{N}} \Rightarrow E f(x_R) - f(x^*) \leq \frac{\|x_1 - x^*\|^2}{\sqrt{N}} + \frac{\sigma^2}{\sqrt{N}}$$

Lecture - 19

SG for 200 case

↳ Convex

↳ Strongly convex

Algo: $x_{k+1} = \Pi(x_k - \alpha g(x_k, \xi_k))$

project onto a convex set C with diameter D
 $(\|x - y\| \leq D)$
 $\forall x, y \in C$

Bound gradients: $\|E_k g - \nabla f\| \leq C_1 \delta^2$

$$E \|g - E g\|^2 \leq C_2 / \delta^2$$

Assume: (i) f is L -smooth & convex

(ii) $\|\nabla f(x)\| \leq B$ ← Bounded gradients

$$\begin{aligned} \|x_{k+1} - x^*\|^2 &= \|x_k - \alpha g(x_k, \xi_k) - x^*\|^2 \\ &= \|x_k - x^*\|^2 - 2\alpha g^T (x_k - x^*) + \alpha^2 \|g\|^2 \end{aligned}$$

Let $\Delta_k = g - \nabla f$

$$\begin{aligned} E_k \|x_{k+1} - x^*\|^2 &\leq \|x_k - x^*\|^2 - 2\alpha \nabla f^T (x_k - x^*) - 2\alpha \Delta_k^T (x_k - x^*) \\ &\quad + \alpha^2 \left(\|E_k g\|^2 + \frac{C_2}{\delta^2} \right) \end{aligned}$$

$$\leq \|x_k - x^*\|^2 - 2\alpha \nabla f^T (x_k - x^*) + 2\alpha C_1 \delta^2 \|x_k - x^*\|,$$

$$+ \alpha^2 \left(\|\nabla f\|^2 + 2\sqrt{2} C_1 \delta^2 \|\nabla f\| + \alpha C_1^2 \delta^4 + \frac{C_2}{\delta^2} \right)$$

— (x)

$$\text{Since } - \sum_{i=1}^d y_i \leq \|y\|_1$$

$$\begin{aligned}\|E_k g\|^2 &= \|\nabla f + c_1 \delta \frac{\mathbf{1}}{2}\|^2 \xrightarrow{\text{vector of ones}} \\ &\leq \|\nabla f\|^2 + c_1^2 \delta^4 \mathbf{1}^\top \mathbf{1} + 2c_1 \delta^2 \nabla f^\top \frac{\mathbf{1}}{2} \\ &\leq \|\nabla f\|^2 + c_1^2 \delta^4 \mathbf{1}^\top \mathbf{1} + 2c_1 \delta^2 \|\nabla f\| \sqrt{2}\end{aligned}$$

$$\text{Letting } z_k = x_k - x^* \text{ and using } \|\nabla f(x)\|^2 \leq L \nabla f(x)^\top (x - x^*)$$

(*) =)

$$\begin{aligned}E_k z_{k+1}^2 &\leq z_k^2 - 2\alpha \nabla f^\top z_k + 2c_1 \delta^2 \|z_k\|_1 \\ &\quad + \alpha^2 \left[L \nabla f^\top z_k + 2\sqrt{2} c_1 \delta^2 L \|z_k\|_1 + \alpha c_1^2 \delta^4 + \frac{c_2}{\delta^2} \right]\end{aligned}$$

$$\text{Convex } f: f(x^*) \geq f(x_k) + (\nabla f(x_k), x^* - x_k)$$

Using this for

$$\begin{aligned}E_k z_{k+1}^2 &\leq z_k^2 - (2\alpha - L\alpha^2) (f(x_k) - f(x^*)) \\ &\quad + 2\sqrt{2} \|z_k\|_1 c_1 \delta^2 (\alpha + L\alpha^2) + \alpha^2 \left[\alpha c_1^2 \delta^4 + \frac{c_2}{\delta^2} \right] \\ &\quad \text{Using } \|y\|_1 \leq \sqrt{2} \|y\|_2\end{aligned}$$

Re-arranging,

$$\begin{aligned}&\alpha (f(x_k) - f(x^*)) \\ &\leq \frac{1}{2-L\alpha} \left[z_k^2 - E_k z_{k+1}^2 + 2\sqrt{2} \|z_k\|_1 c_1 \delta^2 (\alpha + L\alpha^2) + \alpha^2 \left(\alpha c_1^2 \delta^4 + \frac{c_2}{\delta^2} \right) \right]\end{aligned}$$

Setting $\lambda \leq \frac{1}{L}$ & summing over $k=1 \dots N$ & taking total expectations,

$$\sum_{k=1}^N \lambda E(f(x_k) - f(x^*))$$

$$\begin{aligned} &\leq \sum_{k=1}^N \left(E z_k^2 - E z_{k+1}^2 \right) + 2\sqrt{\lambda} \sum_{k=1}^N E \|z_k\| C_1 \delta^2 (\lambda + L\lambda^2) \\ &\quad + \sum_{k=1}^N \lambda^2 \left(d C_1 \delta^2 + \frac{C_2}{\delta^2} \right) \\ &\leq z_1^2 + 2\sqrt{\lambda} D N C_1 \delta^2 (\lambda + L\lambda^2) \\ &\quad + N \lambda^2 \left(d C_1^2 \delta^4 + \frac{C_2}{\delta^2} \right) \end{aligned}$$

$$\begin{aligned} &E f(x_R) - f(x^*) \quad x_k \xrightarrow[\text{random}]{} \{x_1, \dots, x_N\} \\ &\leq \frac{1}{Nd} \left[z_1^2 + 2\sqrt{\lambda} D N C_1 \delta^2 (\lambda + L\lambda^2) \right. \\ &\quad \left. + \lambda^2 \left(d C_1^2 \delta^4 + \frac{C_2}{\delta^2} \right) \right] \end{aligned}$$

Setting $\lambda = \min\left(\frac{1}{C}, \frac{1}{N^{1/3}}\right)$, $\delta = \frac{1}{N^{1/6}}$

$$E f(x_R) - f(x^*) \leq \frac{\text{const}}{N^{1/3}}$$

The claim follows. \square

Remark 5.8. In contrast to the constant step size case handled previously, with a diminishing step size, we have a bound that vanishes as $m \rightarrow \infty$. However, the step size choice requires the knowledge of the strong convexity parameter μ , while the constant step size case in Theorem 5.4 did not assume such information. On a related note, it is possible to obtain a bound of $O(1/\sqrt{m})$ with a step size choice that does not require the knowledge of μ , and more importantly, with a bound that does not scale inversely with μ . Such a bound may be preferable for ill-conditioned problems, where μ is very small. The reader is referred to (Nemirovski *et al.*, 2009) for the details.

5.3.2 SG with biased gradient information

As before, we consider the update iteration in (5.21). Unlike the previous section, where we assumed unbiased gradient estimates (i.e., the condition A5.1 holds), here the estimate $\hat{\nabla}f(\theta_k)$ is a biased approximation to the gradient of the objective function f at θ_k .

As in the asymptotic analysis in Section 4.1.2, the biased gradient estimate $\hat{\nabla}f(\theta_k)$ can be decomposed as follows:

$$\begin{aligned}\hat{\nabla}f(\theta_k) &= \nabla f(\theta_k) + \beta_k + \eta_k, \text{ where} \\ \beta_k &= E\left[\hat{\nabla}f(\theta_k) \mid \mathcal{F}_k\right] - \nabla f(\theta_k), \\ \eta_k &= \hat{\nabla}f(\theta_k) - E\left[\hat{\nabla}f(\theta_k) \mid \mathcal{F}_k\right],\end{aligned}\tag{5.29}$$

where \mathcal{F}_k is a σ -field generated by $\{\theta_i, i \leq k\}$. In the above, β_k is the bias in the gradient estimate and η_k , $n \geq 0$, is a martingale difference sequence.

Using a simultaneous perturbation-based gradient estimate implies $\beta_k = O(\delta_k^2)$, where δ_k is the perturbation parameter used in forming the estimate (see Chapter 3 for several examples). While the bias goes down as δ_k^2 , the variance of the gradient estimate scales inversely with δ_k^2 . This has been formalized earlier in assumptions A5.2–A5.3.

We now present a non-asymptotic bound in expectation for the SG algorithm (5.21) with inputs from a biased gradient oracle that satisfies the aforementioned assumptions.

$$\theta_{k+1} = \theta_k - a_k \hat{\nabla} f(\theta_k)$$

136

Non-asymptotic analysis of stochastic gradient algorithms

f is μ -strongly convex

Proposition 5.1. Suppose the objective function f is L -smooth (see Definition 3.1), and assumptions A5.2–A5.3 hold. Then, we have

$$\begin{aligned} \mathbb{E} \|\theta_{m+1} - \theta^*\|^2 &\leq \underbrace{2 \exp(-2\mu\Gamma(m)) \|\theta_0 - \theta^*\|^2}_{\text{initial error}} \\ &+ \underbrace{2 \sum_{k=1}^n a_k^2 \exp(-2\mu(\Gamma(m) - \Gamma_k)) c_1^2 \delta_k^4 +}_{\text{bias error}} \\ &\underbrace{\sum_{k=1}^n a_k^2 \exp(-2\mu(\Gamma(m) - \Gamma_k)) c_2 \delta_k^{-2},}_{\text{sampling error}} \end{aligned} \quad (5.30)$$

where $\Gamma(k) := \sum_{i=1}^k a_i$.

$$\int f''(\theta^* + \lambda(\theta_m - \theta^*)) (\theta_m - \theta^*) d\lambda = f'(\theta_m)$$

Proof. Let $z_m = \theta_m - \theta^*$ denote the error at time instant n of the algorithm (5.21). Using $\nabla f(\theta^*) = 0$, we have

$$\left(\int_0^1 \nabla^2 f(\theta^* + \lambda(\theta_m - \theta^*)) d\lambda \right) z_m = \nabla f(\theta_m).$$

Using the fact above, we arrive at a recursion for z_m from (5.29). Letting

$$J_m := \int_0^1 \nabla^2 f(\theta^* + \lambda(\theta_m - \theta^*)) d\lambda,$$

$$z_{m+1} = (I - a(m) J_m) z_m - a(m) (\beta_m + \eta_m)$$

$$= \Pi_m z_0 - \sum_{k=1}^n a(k) \Pi_m \Pi_k^{-1} (\beta_k + \eta_k),$$

$$\text{where } \Pi_m := \prod_{k=1}^n (I - a(k) J_k).$$

By the conditional Jensen's inequality, we obtain

$$(\mathbb{E}_m \|z_{m+1}\|)^2 \leq \mathbb{E}_m (\langle z_m, z_m \rangle)$$

$$\begin{aligned} &= \mathbb{E}_m \left(\|\Pi_m z_0\|^2 + \left\| \sum_{k=1}^n a(k) \Pi_m \Pi_k^{-1} \beta_k \right\|^2 + \left\| \sum_{k=1}^n a(k) \Pi_m \Pi_k^{-1} \eta_k \right\|^2 \right. \\ &\quad \left. - 2 \left\langle \Pi_m z_0, \sum_{k=1}^n a(k) \Pi_m \Pi_k^{-1} \beta_k \right\rangle - 2 \left\langle \Pi_m z_0, \sum_{k=1}^n a(k) \Pi_m \Pi_k^{-1} \eta_k \right\rangle \right) \end{aligned}$$

→ shorthand for
 $\prod_{j=2}^n (I - a(j) J_j)$

$$\textcolor{purple}{\cancel{\left\langle \sum_{k=1}^n a(k) \Pi_m \Pi_k^{-1} \beta_k, \sum_{k=1}^n a(k) \Pi_m \Pi_k^{-1} \eta_k \right\rangle \right)}} \quad (5.31)$$

$$\begin{aligned} & \leq 2 \|\Pi_m z_0\|^2 + 2 \sum_{k=1}^n a(k)^2 \|\Pi_m \Pi_k^{-1}\|^2 c_1^2 \delta_k^4 \\ & \quad + \sum_{k=1}^n a(k)^2 \|\Pi_m \Pi_k^{-1}\|^2 \mathbb{E} \|\eta_k\|^2. \end{aligned} \quad (5.32)$$

For the last inequality, we have used the following facts: (i) η_k is a martingale difference implying the last two cross terms in (5.31) are zero; (ii) $\beta_k \leq c_1 \delta_k^2$ from A5.3; and (iii) Cauchy-Schwarz inequality for the first cross term in (5.31).

Now, we bound each of the square terms in (5.32) separately. Since the objective is strongly convex, we have that $\|I - a(m)J_m\| \leq \exp(-\mu a(m))$. Hence,

$$\begin{aligned} \|\Pi_m \Pi_k^{-1}\|_2 &= \left\| \prod_{j=k+1}^n (I - a_j J_j) \right\|_2 \\ &\leq \prod_{j=k+1}^n \|(1 - a_j \mu)I - a_j (J_j - \mu I)\|_2 \\ &\leq \prod_{j=k+1}^n \|(1 - a_j \mu)I\|_2 \leq \prod_{j=k+1}^n (1 - a_j \mu) \\ &\leq \exp(-\mu(\Gamma(m) - \Gamma(k))). \end{aligned} \quad (5.33)$$

$\Gamma(m) = \sum_{j=1}^m a(j)$

From A5.3, we can infer that the second moment of the martingale difference is bounded above by c_2/δ_k^2 . The main claim now follows by plugging the bound on η_m and (5.33) into (5.32). \square

By specializing the result in the proposition above, we derive a non-asymptotic bound of the order $O(1/\sqrt{m})$.

Theorem 5.6. (Biased gradients and strongly convex objec-

tive) Let $a(k) = c/k$ and $\delta_k = \delta_0/k^\delta$. Then,

$$\begin{aligned}\mathbb{E} \|\theta_m - \theta^*\| &\leq \frac{\sqrt{2} \|\theta_0 - \theta^*\|}{m^{\mu c}} + \frac{\sqrt{2} c c_1 \delta_0^2}{\sqrt{2\mu c - 4\delta - 1}} m^{-\frac{1}{2}-2\delta} \\ &\quad + \frac{\sqrt{c_2} c}{\delta_0 \sqrt{2\mu c + 2\delta - 1}} m^{\delta-\frac{1}{2}}.\end{aligned}$$

Remark 5.9. Choosing $\delta = 0$, one can obtain a bound of the order $O(m^{-1/2})$ for simultaneous perturbation schemes that lead to biased gradient estimates, and this bound matches the corresponding bound with unbiased gradient information up to constant factors. Contrast this with the difference in rates between biased and unbiased gradient information for the non-convex and convex cases in the previous sections.

Remark 5.10. Using L -smoothness of f and $\nabla f(\theta^*) = 0$, we have

$$\mathbb{E}[f(\theta_m)] - f(\theta^*) \leq \frac{L}{2} \mathbb{E} \|\theta_m - \theta^*\|^2 = O\left(\frac{1}{m}\right).$$

Proof. Bounding a sum by an integral, we obtain

$$\exp(-\mu\Gamma(m)) \leq \exp(-\mu c \ln m) = m^{-\mu c}.$$

Plugging $a(k) = c/k$ and $\delta_k = \delta_0/k^\delta$ into the bias error term in (5.30), we obtain

$$\begin{aligned}\sum_{k=1}^m a(k)^2 \exp(-2\mu(\Gamma(m) - \Gamma_k)) c_1^2 \delta_k^4 &\leq \sum_{k=1}^m \frac{c^2}{k^2} n^{-2\mu c} k^{2\mu c} c_1^2 \frac{\delta_0^4}{m^{4\delta}} \\ &\leq c^2 n^{-2\mu c} c_1^2 \delta_0^4 \sum_{k=1}^m k^{2\mu c - 4\delta - 2} \\ &\leq \frac{c^2 c_1^2 \delta_0^4}{(2\mu c - 4\delta - 1)} m^{-1-4\delta}.\end{aligned}$$

Along similar lines, the sampling error term in (5.30) can be upper-bounded as follows:

$$\sum_{k=1}^m a(k)^2 \exp(-2\mu(\Gamma(m) - \Gamma_k)) \frac{c_2}{\delta_k^2} \leq \frac{c^2 c_2}{\delta_0^2 (2\mu c - 4\delta - 1)} m^{-1+2\delta}.$$

□