

535520: Optimization Algorithms

Lecture 7 – Stochastic Gradient Descent and Variance Reduction

Ping-Chun Hsieh (謝秉均)

October 21, 2024

Announcement

- ▶ Final project: Please select a paper and fill out the Google form: <https://forms.gle/XAAoeWFxi5ZUDhpG7>

This Lecture

1. Stochastic Gradient Descent

2. Accelerating SGD: Variance Reduction

- Reading Material:
 - Chapter 4 of the textbook “Optimization Methods for Large-Scale Machine Learning” by Leon Bottou, Frank Curtis, and Jorge Nocedal.
 - Available at <https://arxiv.org/abs/1606.04838>
 - R. Johnson and T. Zhang, “Accelerating stochastic gradient descent using predictive variance reduction,” NIPS 2013

Today's Lecture: Convergence Results of SGD

	Fixed step size $(\eta_t \equiv \eta)$	Diminishing step sizes $(\eta_t = \Theta(1/t))$
μ -strongly convex functions	$\mathbb{E}[F(x_t) - F_*] \leq \frac{\eta L \sigma^2}{2\mu} + (1 - \eta\mu)^t (F(x_0) - F_*)$	$\mathbb{E}[F(x_t) - F_*] = O\left(\frac{1}{t}\right)$
General functions	$\mathbb{E}\left[\frac{1}{t} \sum_{i=1}^t \ \nabla F(x_i)\ ^2\right] \leq \frac{\eta L \sigma^2 + \frac{2(F(x_0) - F_*)}{\eta t}}{t}$	$\lim_{t \rightarrow \infty} \mathbb{E}\left[\frac{\sum_{i=1}^t \eta_i \ \nabla F(x_i)\ ^2}{\sum_{i=1}^t \eta_i}\right] = 0$

- For general non-convex functions, SGD converges to a (nearly-)stationary point

Review: Stochastic Optimization in Machine Learning

- Optimization problems in ML arise mainly in two forms:

1. Expected risk/loss minimization

$$F^* := \min_{x \in X} F(x), \text{ where } F(x) := \mathbb{E}_{\varepsilon \sim D}[f(x; \varepsilon)]$$

where ε is the randomness (possibly unknown) in our problem

2. Empirical risk minimization (aka Finite-sum problem)

$$F^* := \min_{x \in X} F(x), \text{ where } F(x) := \frac{1}{n} \sum_{i=1}^n f(x; d_i)$$

where $\{d_i\}_{i=1}^n$ are n random data samples

$$\begin{aligned}\nabla F(x) &= \nabla \mathbb{E}[f(x; \varepsilon)] \\ &= \mathbb{E}[\nabla f(x; \varepsilon)]\end{aligned}$$

Stochastic GD for Empirical Risk Minimization

- Idea: Use sampling to estimate expectation [Robbins & Monro, 1951]

At each iteration, we randomly pick an integer $i \in \{1, 2, \dots, n\}$

$$x_{k+1} = x_k - \eta_k \cdot \frac{1}{n} \sum_{i=1}^n \nabla f(x_k; d_i) \quad \xrightarrow{\text{(GD)}} \quad x_{k+1} = x_k - \eta_k \cdot \nabla f(x_k; d_i) \quad \text{(SGD)}$$

- The update requires only gradient of one data sample d_i
- $\nabla f(x_k; d_i)$ is an unbiased estimate of $\nabla f(x_k)$, i.e., $\mathbb{E}[\nabla f(x_k; d_i)] = \nabla f(x_k)$

Review: Stochastic GD for Expected Risk Minimization

- Idea: Use **sampling** to estimate expectation [Robbins & Monro, 1951]

$$x_{k+1} = x_k - \eta_k \cdot \mathbb{E}[\nabla f(x_k; \varepsilon_k)] \quad \rightarrow \quad x_{k+1} = x_k - \eta_k \cdot g(x_k; \varepsilon_k)$$

(GD) → (SGD)

- ▶ $g(x_k; \varepsilon_k)$ is an unbiased estimate of $\mathbb{E}[\nabla f(x; \varepsilon_k)]$
 - ▶ Equivalently:
$$g(x_k; \varepsilon_k) = \mathbb{E}[\nabla f(x_k; \varepsilon_k)] + \underbrace{\varepsilon}_{\text{zero mean noise}}$$
 - ▶ $g(x_k; \varepsilon_k)$ is constructed using one or multiple samples (= mini-batch)

Some Mild Technical Assumptions

► Quick recap:

$$\text{(Objective)} \quad F^* := \min_{x \in X} F(x) \quad (\text{where } F(x) = \mathbb{E}[f(x; \varepsilon)] \text{ or } F(x) := \frac{1}{n} \sum_{i=1}^n f(x; d_i))$$

$$\text{(SGD)} \quad x_{k+1} = x_k - \underbrace{\eta_k \cdot g(x_k; \varepsilon_k)}_{\text{unbiased estimate}}$$

► To prove convergence of SGD, consider the following mild technical assumptions:

(A1) $F(x)$ is bounded below, i.e. $F^* > -\infty$

(A2) $F(x)$ is L -smooth *stochastic gradient*

(A3) $g(x_k; \varepsilon_k)$ has bounded variance, i.e.

$$\mathbb{V}[g(x_k; \varepsilon_k) | x_k] := \mathbb{E}[\|g(x_k; \varepsilon_k)\|^2 | x_k] - \|\mathbb{E}[g(x_k; \varepsilon_k) | x_k]\|^2 \leq M + M_V \|\nabla F(x_k)\|^2$$

Review: Three Useful Lemmas

Lemma 1: PL conditions for μ -strongly convex functions

$$F(x) - F^* \leq \frac{1}{2\mu} \|\nabla F(x)\|^2$$

If $g(x_k; \varepsilon_k)$ is biased, then
this term becomes $-\eta_k \|\nabla F(x_k)\|^2$

Lemma 2: First stochastic descent lemma of SGD

$$\mathbb{E}[F(x_{k+1})|x_k] - F(x_k) \leq -\eta_k \nabla F(x_k)^\top \mathbb{E}[g(x_k; \varepsilon_k)|x_k] + \frac{1}{2} \eta_k^2 L \cdot \mathbb{E}[\|g(x_k; \varepsilon_k)\|^2 | x_k]$$

expected improvement
attained improvement by the 1st-order term
effect of 2nd-order term.

Lemma 3: Second stochastic descent lemma of SGD

$$\mathbb{E}[F(x_{k+1})|x_k] - F(x_k) \leq -\left(1 - \frac{1}{2}\eta_k L(M_\nu + 1)\right) \eta_k \|\nabla F(x_k)\|^2 + \frac{1}{2} \eta_k^2 LM$$

Choose η_k such that RHS is negative

Convergence of SGD: Strong Convexity + Fixed Step Size (SC+F)

We proved the following result in Lecture 6:

Theorem (Convergence of SGD under SC+F):

Suppose the following conditions hold.

- (1) Assumptions (A1)-(A3)
- (2) $F(x)$ is L -smooth and μ -strongly convex

- (3) Step sizes $\eta_k \equiv \eta$ satisfy that $\eta < \frac{\mu}{L(M_\nu + 1)}$

Then, SGD achieves

$$\mathbb{E}[F(x_k)] - F^* \leq \frac{\eta LM}{2\mu} + (1 - \eta\mu)^k(F(x_0) - F^* - \frac{\eta LM}{2\mu})$$

A Closer Look at Theorem (SC+F)

$$\mathbb{E}[F(x_k) - F_*] \leq \frac{\eta LM}{2\mu} + (1 - \eta\mu)^k (F(x_0) - F_*)$$

- What if there is no noise in the estimated gradients?

$M=0$, we have linear convergence.

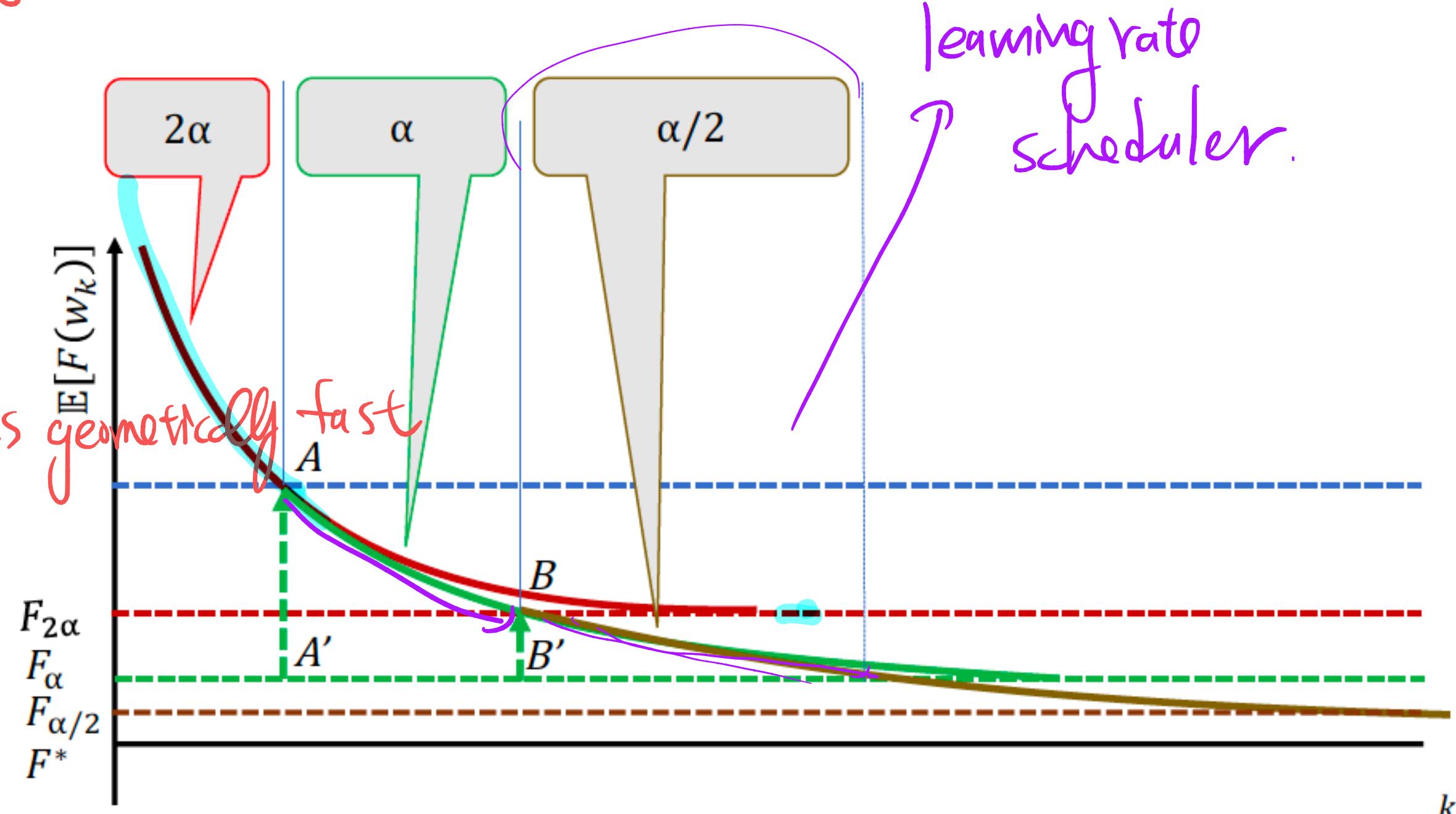
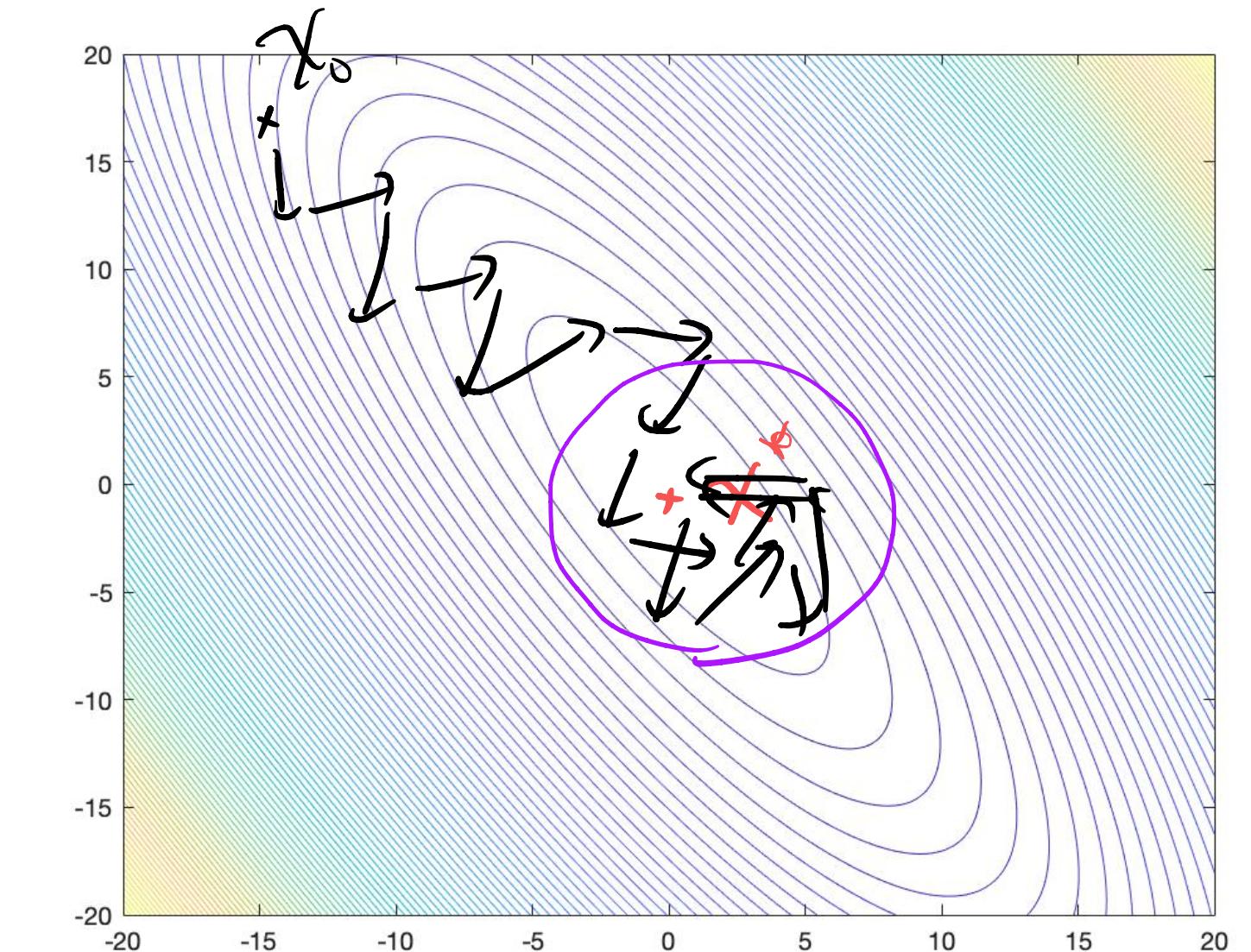
- Why do we need “strong convexity”?

- What’s the role of initial condition?

x_0 only takes part in $(1 - \eta\mu)^k (F(x_0) - F^*)$, which decreases geometrically

- A practical strategy of choosing η in SGD:

- Run SGD with fixed step sizes
- Whenever the progress stalls, we reduce the step sizes and proceed with SGD



(Figure Source: [Bottou, Curtis, Nocedal, 2018])

Convergence of SGD: Strong Convexity + Diminishing Step Size (SC+D)

Theorem (Convergence of SGD under SC+D):

Suppose the following conditions hold.

(1) Assumptions (A1)-(A3)

(2) $F(x)$ is L -smooth and μ -strongly convex

(3) Step sizes $\eta_k = \frac{1}{k + \gamma}$, where $\beta > \frac{1}{\mu}$ and $\gamma > 0$ such that $\eta_1 \leq \frac{1}{L(M_\nu + 1)}$

Then, SGD achieves

$$\mathbb{E}[F(x_k)] - F^* \leq \frac{\nu}{k + \gamma} = O\left(\frac{1}{k}\right)$$

$$\text{where } \nu := \max \left\{ \frac{\beta^2 LM}{2(\beta\mu - 1)}, (\gamma + 1)(F(x_0) - F^*) \right\}$$

$$\begin{aligned} \eta_1 \cdot L(M_\nu + 1) &\leq 1 \\ \Rightarrow \eta_k \cdot L(M_\nu + 1) &\leq 1 \text{ for all } k \end{aligned}$$

$$\frac{1}{k + \gamma} \leq \frac{1}{k}$$

$$\begin{aligned} \frac{\nu}{k + \gamma} &= \varepsilon \\ \Rightarrow k &= \frac{\nu - \varepsilon\gamma}{\varepsilon} \end{aligned}$$

Proof of Theorem (SC+D)

$$E[E[X|Y]] = E[X]$$

Step 1: By the 2nd stochastic descent lemma,

$$\begin{aligned} E\left[E[F(x_{k+1})|x_k] - F(x_k)\right] &\leq -\left(1 - \frac{1}{2}\eta_k \cdot L \cdot (M_\nu + 1)\right) \cdot \eta_k \cdot \|\nabla F(x_k)\|^2 + \frac{1}{2} \eta_k^2 \cdot L \cdot M \\ &\stackrel{\text{Law of iterated expectation:}}{\leq} \left[-\frac{1}{2} \eta_k \cdot \|\nabla F(x_k)\|^2 + \frac{1}{2} \eta_k^2 \cdot L \cdot M\right] \\ &\leq ((E[F(x_k)] - F^*) \cdot 2\mu) - \frac{1}{2} \eta_k \end{aligned}$$

Step 2: By taking total expectation and subtracting F^* ,

$$E[F(x_{k+1}) - F^*] \leq \left(1 - \eta_k \cdot \mu\right) \cdot E[F(x_k) - F^*] + \frac{1}{2} \eta_k^2 \cdot L \cdot M \quad (*)$$

Step 3: We prove

$$E[F(x_k) - F^*] \leq \frac{\gamma}{K+\gamma} \quad \text{by induction}$$

(Cont.). Want: $E[F(x_k) - F^*] \leq \frac{\gamma}{\kappa + \gamma}$, where $\gamma := \max \left\{ \frac{\beta^2 L M}{2(\beta \mu - 1)}, (\gamma+1) \cdot (F(x_0) - F^*) \right\}$

For $k=1$: By (*),

$$E[F(x_1) - F^*] \leq (1 - \gamma_1 \cdot \mu) \cdot (F(x_0) - F^*) + \frac{1}{2} \gamma_1^2 L \cdot M$$

Suppose (***) holds under K and consider $K+1$:

$$E[F(x_{K+1}) - F^*] \leq \left(1 - \frac{\beta}{\gamma+K} \mu\right) \cdot \frac{\gamma}{\gamma+K} + \frac{\beta^2 L \cdot M}{2(\gamma+K)^2}$$

$$= \left(\frac{(\gamma+K)-1}{(\gamma+K)^2}\right) \gamma - \underbrace{\left(\frac{\beta \cdot \mu - 1}{(\gamma+K)^2}\right)}_{\text{non-positive}} \gamma + \frac{\beta^2 L \cdot M}{2(\gamma+K)^2}$$

$\leq \frac{\gamma}{(\gamma+K)+1}$

since $(\gamma+K)^2 \geq (\gamma+K+1)(\gamma+K-1)$

A Closer Look at Theorem (SC+D)

$$\mathbb{E}[F(x_k) - F_*] \leq \frac{\nu}{k + \gamma}$$

$$\nu = \max \left\{ \frac{\beta^2 LM}{2(\beta\mu - 1)}, (\gamma + 1)(F(x_0) - F(x^*)) \right\}$$

- ▶ Why do we need “strong convexity”? *Connecting $\|\nabla F(x_k)\|^2$ with sub-optimality gap $(F(x_k) - F^*)$*
- ▶ Why do we need $\beta > \frac{1}{\mu}$? *For sufficiently fast convergence!*
- ▶ A toy example in [Nemirovski et al., 2009]

A Toy Example: Slow Convergence Due to Improper β

► Consider $f(x) = x^2/10$ with $x_0 = 1$

► What is μ in this example?

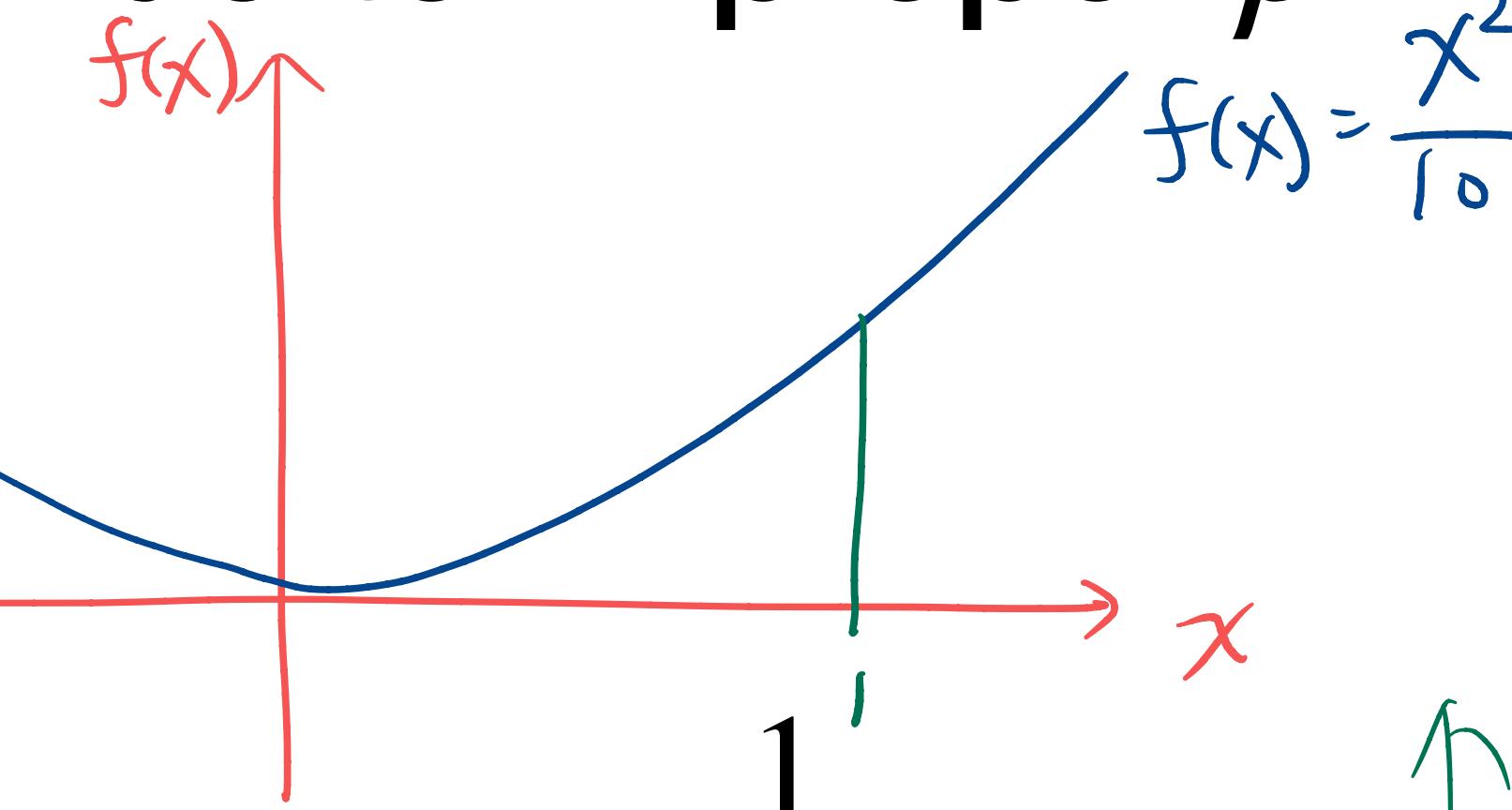
$f(x)$ is $\frac{1}{5}$ -strongly convex

Suppose we take the step sizes $\eta_k = \frac{1}{k}$ (i.e., $\beta = 1 < \frac{1}{\mu} = \frac{10}{5} = 2$)

$$x_{k+1} = x_k - \frac{1}{k} f'(x_k) = \left(1 - \frac{1}{5k}\right) x_k$$

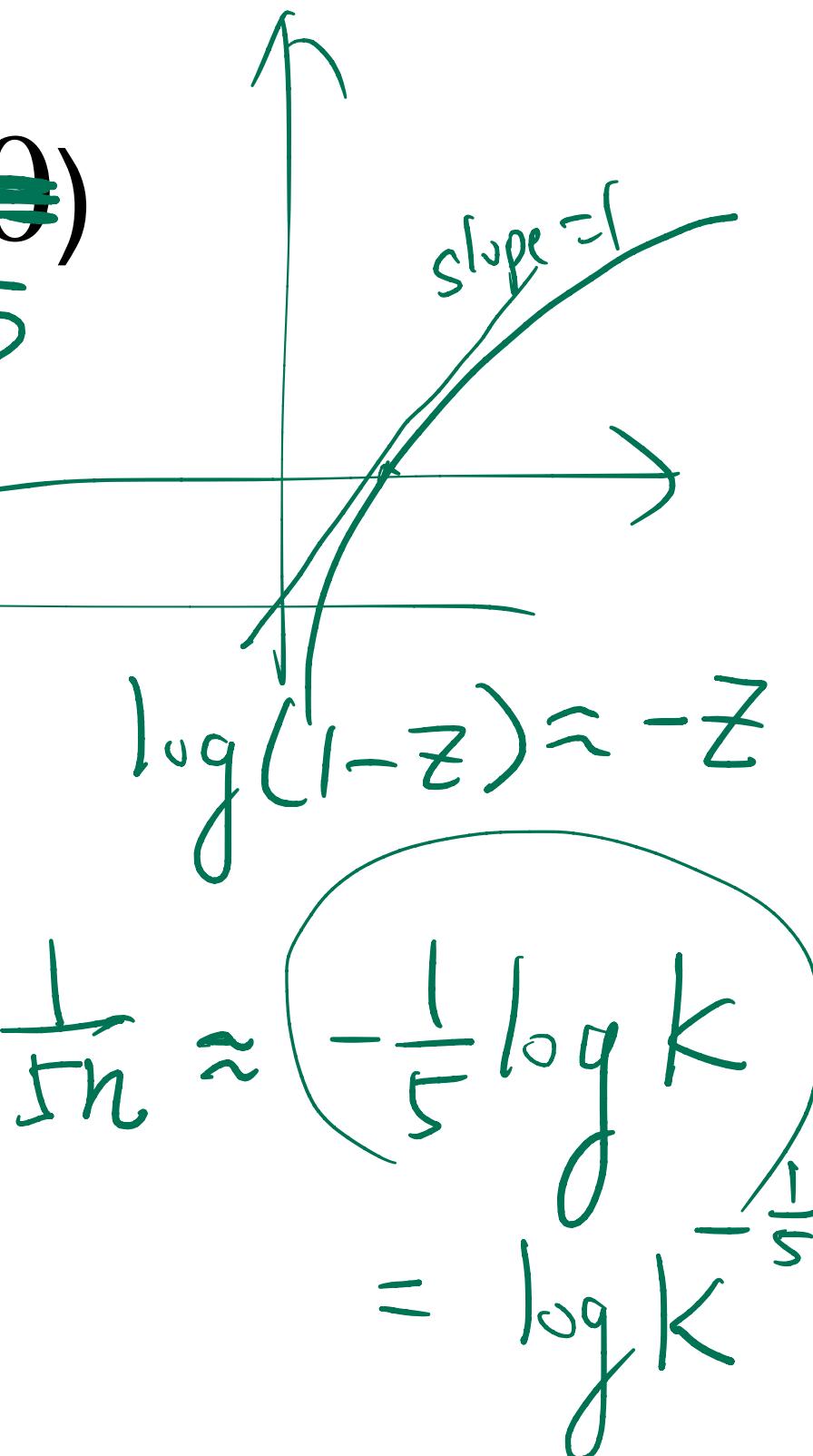
$$x_k = \prod_{n=0}^{k-1} \left(1 - \frac{1}{5n}\right) x_0 > 0.8 \cdot k^{-\frac{1}{5}}$$

$$f(x_k) > \frac{(0.8 \cdot k^{-\frac{1}{5}})^2}{10} \left(1 - \frac{1}{5}\right) \cdot \left(1 - \frac{1}{10}\right) \cdot \left(1 - \frac{1}{15}\right) \cdots$$



$$f'(x_k) = \frac{x_k}{5}$$

$$\begin{aligned} & \log \left(\prod_{n=1}^k \left(1 - \frac{1}{5n}\right) \right) \\ &= \sum_{n=1}^k \log \left(1 - \frac{1}{5n}\right) \approx \sum_{n=1}^k -\frac{1}{5n} \approx -\frac{1}{5} \log k \\ &\Rightarrow \prod_{n=1}^k \left(1 - \frac{1}{5n}\right) \approx k^{-\frac{1}{5}} \end{aligned}$$



Next Question: Is $O\left(\frac{1}{k}\right)$ optimal?

- ▶ (Informal) For the minimization of strongly convex functions, there is **no** algorithm that can achieve an accuracy better than $O(1/k)$ on performing k queries

- ▶ Therefore, SGD with step sizes $\Theta(1/k)$ is optimal

Information-theoretic lower bounds on the oracle complexity of stochastic convex optimization

Alekh Agarwal¹
alekh@cs.berkeley.edu

Peter L. Bartlett^{1,2,3}
peter@berkeley.edu

Pradeep Ravikumar⁴
pradeep@cs.utexas.edu

Martin J. Wainwright^{1,2}
wainwrig@stat.berkeley.edu

Department of Electrical Engineering and Computer Sciences¹
Department of Statistics²
UC Berkeley, Berkeley, CA

Mathematical Sciences³
QUT, Brisbane, Australia

Department of Computer Sciences⁴
UT Austin, Austin, TX

November 22, 2011

Abstract

Relative to the large literature on upper bounds on complexity of convex optimization, lesser attention has been paid to the fundamental hardness of these problems. Given the extensive use of convex optimization in machine learning and statistics, gaining an understanding of these complexity-theoretic issues is important. In this paper, we study the complexity of stochastic convex optimization in an oracle model of computation. We improve upon known results and obtain tight minimax complexity estimates for various function classes.



A. Nemirovski, D. Yudin, “Problem complexity and method efficiency in optimization,” Wiley, 1983

Alekh Agarwal
@ Google

A. Agarwal, P. Barlett, P. Ravikumar, M. Wainright, “Information-theoretic lower bounds on the oracle complexity of stochastic convex optimization,” IEEE Transactions on Information Theory, 2011

Comparison: SGD and Batch GD

- Consider strongly convex “empirical risk minimization” with n samples
- To achieve an accuracy of ϵ , we need: $(\kappa := L/\mu \text{ is the condition number})$

	Iteration Complexity	Per-iteration Cost	Total Computation Cost
Batch GD	$\kappa \log \frac{1}{\epsilon}$	n	$n\kappa \log \frac{1}{\epsilon}$
SGD	$\kappa^2 \frac{1}{\epsilon}$	1	$\kappa^2 \frac{1}{\epsilon}$

SGD has an advantage for large n and moderate ϵ (“big data” regime!)

That is, $\frac{1}{\epsilon} < n \log \frac{1}{\epsilon}$

SGD for Non-Convex and Smooth Functions

Convergence of SGD: Non-Convex Functions + Fixed Step Size (NC+F)

Theorem (Convergence of SGD under NC+F):

Suppose the following conditions hold.

(1) Assumptions (A1)-(A3)

(2) $F(x)$ is L -smooth (but not necessarily convex)

(3) Step sizes $\eta_k \equiv \eta$ satisfy that $\eta < \frac{\mu}{L(M_\nu + 1)}$

Then, SGD achieves

$$\frac{1}{K} \mathbb{E} \left[\sum_{k=0}^{K-1} \|\nabla F(x_k)\|^2 \right] \leq \eta L M + \frac{2(F(x_0) - F^*)}{K\eta}$$

Proof :

Step 1: By the 2nd stochastic descent lemma and taking "total expectation"

$$\begin{aligned} & \frac{E[F(x_{k+1})] - E[F(x_k)]}{\eta} \\ & \leq -(1 - \frac{1}{2}\eta \cdot L \cdot (1 + M_v)) \cdot \eta \cdot E[\|\nabla F(x_k)\|^2] + \frac{1}{2}\eta^2 \cdot L \cdot M \\ & \leq -\frac{1}{2}\eta \cdot E[\|\nabla F(x_k)\|^2] + \frac{1}{2}\eta^2 \cdot L \cdot M \end{aligned}$$

Step 2: By taking the "telescoping sum", we have

$$E[F(x_{K+1})] - F(x_0) \leq -\frac{1}{2}\eta \cdot \sum_{k=0}^{K-1} E[\|\nabla F(x_k)\|^2] + \frac{1}{2}K\eta^2 \cdot L \cdot M$$

$$\begin{aligned} & E[F(x_{k+1})] - E[F(x_k)] \leq \dots \\ & E[F(x_k)] - E[F(x_{k-1})] \leq \dots \\ & \vdots \\ & +) E[F(x_1)] - E[F(x_0)] \leq \dots \end{aligned}$$

Convergence of SGD: Non-Convex Functions + Diminishing Step Size (NC+D)

Theorem (Convergence of SGD under NC+D):

Suppose the following conditions hold.

(1) Assumptions (A1)-(A3)

(2) $F(x)$ is L -smooth (but not necessarily convex)

(3) Step sizes η_k satisfy

Then, SGD achieves

and therefore
non-summable

total distance

Ensure that SGD can reach at least a stationary point.

2nd-order term does NOT explode

"square-summable"

$$\sum_{k=1}^{\infty} \eta_k = \infty \text{ and } \sum_{k=1}^{\infty} \eta_k^2 < \infty$$
$$\lim_{K \rightarrow \infty} \mathbb{E} \left[\sum_{k=1}^K \eta_k \cdot \|\nabla F(x_k)\|^2 \right] < \infty$$
$$\lim_{K \rightarrow \infty} \mathbb{E} \left[\frac{\sum_{k=1}^K \eta_k \|\nabla F(x_k)\|^2}{\sum_{k=1}^K \eta_k} \right] = 0$$

Proof: Recall from the 2nd descent lemma

Step 1: $E[F(x_{k+1})] - E[F(x_k)]$

$$\leq -\left(1 - \frac{1}{2}R_k L \cdot (M_{k+1})\right) R_k \cdot E[\|\nabla F(x_k)\|^2] + \frac{1}{2}R_k^2 LM$$

$$\leq -\frac{1}{2} R_k \cdot E[\|\nabla F(x_k)\|^2] + \frac{1}{2}R_k^2 LM \quad \dots ($$

Step 2: By summing the above over $k=1, \dots, K$,

$$E[F(x_{K+1})] - E[F(x_1)] \leq -\frac{1}{2} \sum_{k=1}^K R_k \cdot E[\|\nabla F(x_k)\|^2] + \frac{1}{2}LM \cdot \sum_{k=1}^K R_k^2$$

" $F^* - E[F(x_1)]$

Recall: Convergence Rates of SGD

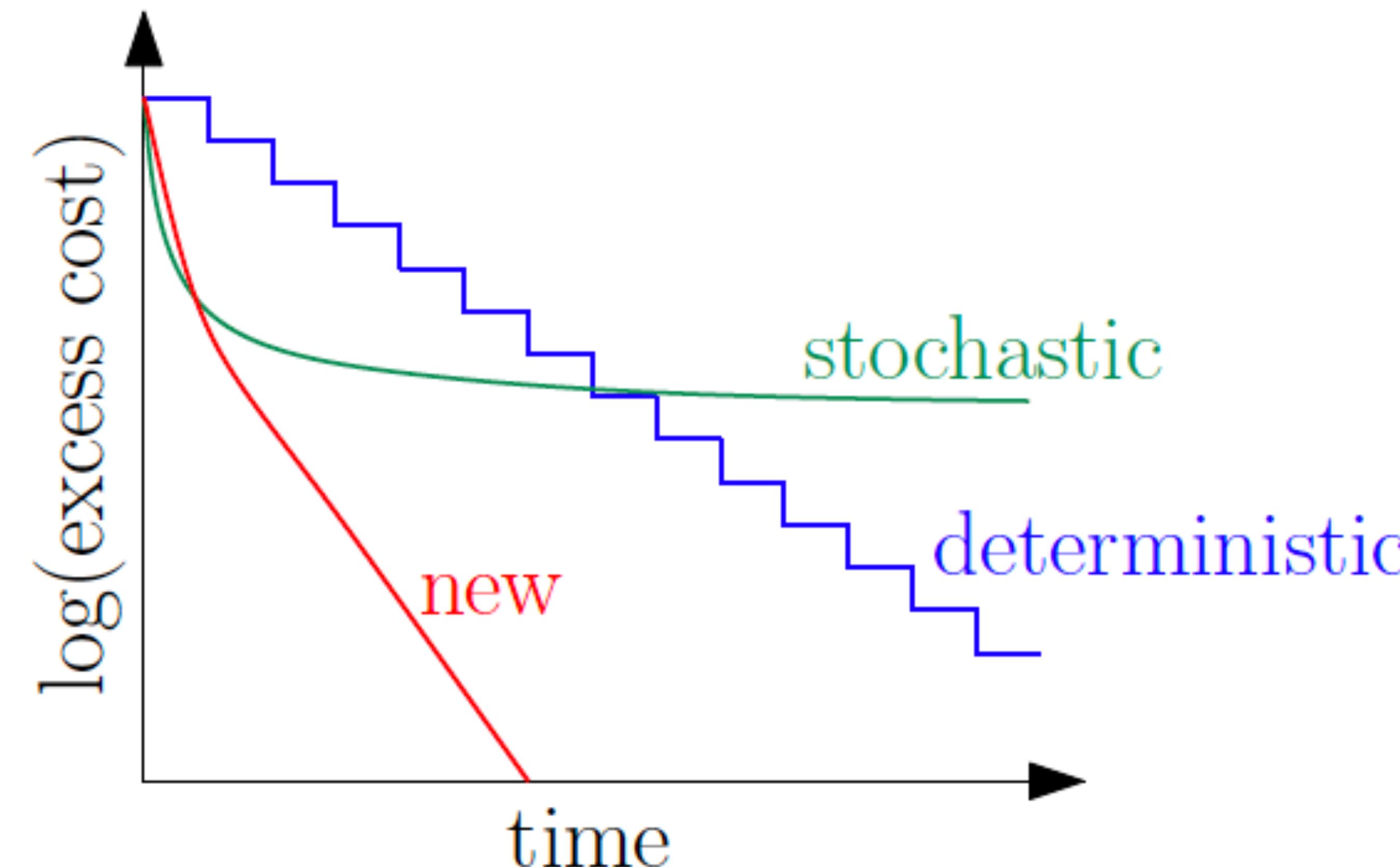
- ▶ We know the convergence rate of SGD for smooth strongly-convex functions

	Fixed step size $(\eta_t \equiv \eta)$	Diminishing step sizes $(\eta_t = \Theta(1/t))$
μ -strongly convex functions	$\mathbb{E}[F(x_t) - F_*] \leq \frac{\eta L \sigma^2}{2\mu} + (1 - \eta\mu)^t (F(x_0) - F_*)$	$\mathbb{E}[F(x_t) - F_*] = O\left(\frac{1}{t}\right)$

- ▶ Any issue with the use of a fixed step size? **Wandering around a small neighborhood**
- ▶ Any issue with the use of $\eta_t = \Theta(1/t)$? **Slow convergence**

Best of Both World?

Question: Can we have an algorithm that achieves "best of both world"
That is, "linear convergence rate" + "low per-iteration cost"



Next topic: Variance reduction

(Figure Credit: Suvrit Sra)

Variance Reduction for SGD: Finite-Sum Problems

Intuition Behind Variance Reduction

$$I_k \sim \text{Unif}(1, n)$$

- In vanilla SGD, we use an unbiased estimator $g(x_k, \varepsilon_k)$

$$\nabla f(x_k; d_{I_k}).$$

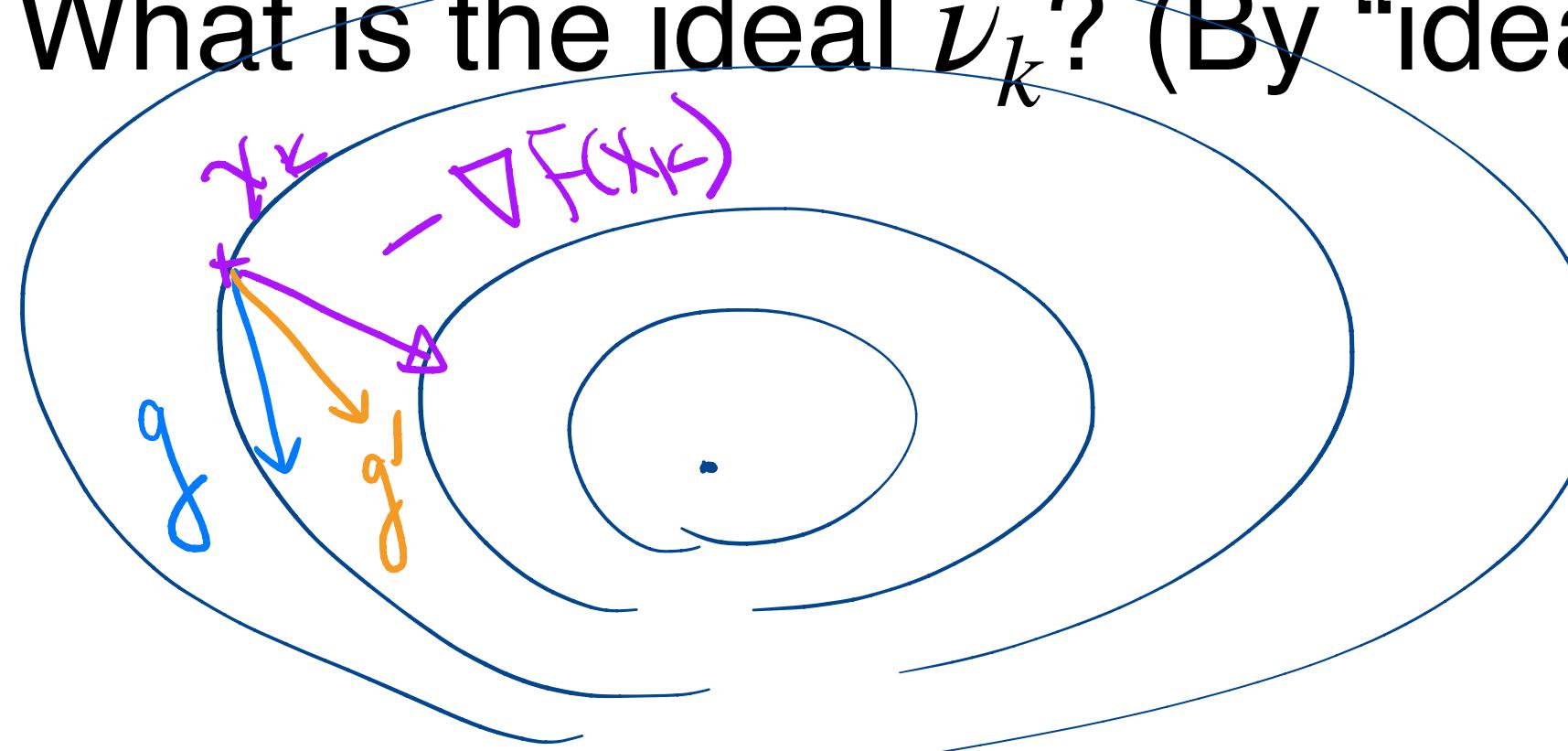
- Question: Can we use $\tilde{g} := \underline{g(x_k, \varepsilon_k)} + \nu_k$ with $\mathbb{E}[\nu_k] = 0$?

$$\mathbb{E}[\tilde{g}] = \mathbb{E}[\underline{g(x_k; \varepsilon_k)} + \nu_k] = \underline{\nabla \bar{F}(x_k)} \Rightarrow \tilde{g} \text{ is an unbiased estimator}$$

- Question: What kind of condition do we want about ν_k ?

To use \tilde{g} in SGD, we need ν_k to be of bounded variance.

- Question: What is the ideal ν_k ? (By “ideal”, we mean zero variance)



Choose ν_k such that $\tilde{g} = \nabla \bar{F}(x_k)$

$$\text{That is, } \nu_k = \tilde{g} - g(x_k; \varepsilon_k) = \nabla \bar{F}(x_k) - g(x_k; \varepsilon_k)$$

A Intuitive Idea of Variance Reduction: Gradient Aggregation

- **Example:** SGD for empirical risk minimization

$$\min_{x \in X} F(x) := -\frac{1}{n} \sum_{i=1}^n f(x; d_i) \quad (\{d_i\}_{i=1}^n \text{ are data samples})$$

Vanilla SGD:

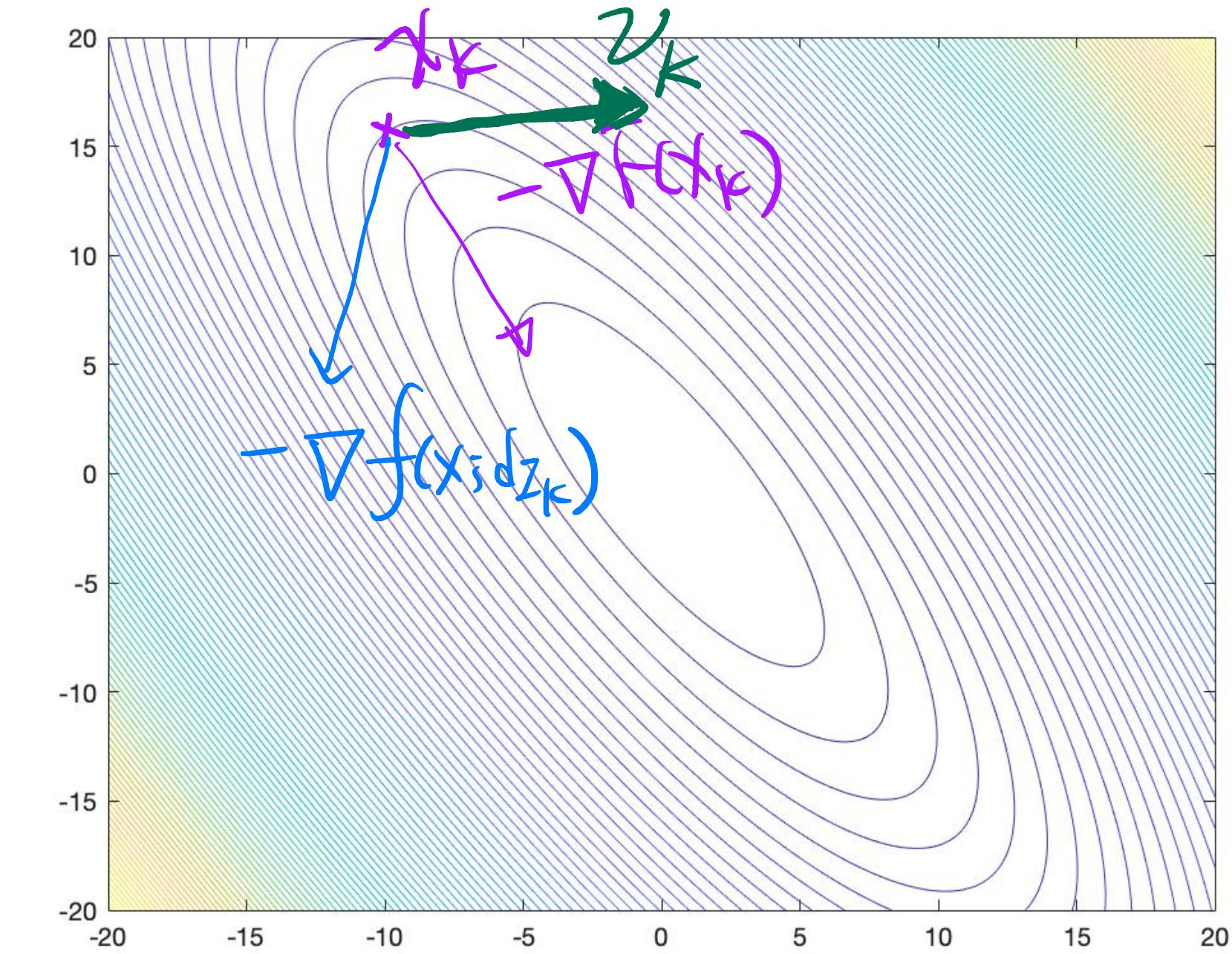
$$x_{k+1} = x_k - \eta \nabla f(x; d_{I_k}), \quad I_k \sim \text{Unif}(1, n)$$

Stochastic gradient

Gradient Aggregation:

$$x_{k+1} = x_k - \eta \left(\nabla f(x; d_{I_k}) - \nu_k \right), \quad I_k \sim \text{Unif}(1, n)$$

where ν_k and $\nabla f(x; d_{I_k})$ are positively correlated (why?)



Stochastic Variance-Reduced Gradient (SVRG)

Accelerating Stochastic Gradient Descent using Predictive Variance Reduction

Rie Johnson
RJ Research Consulting
Tarrytown NY, USA

Tong Zhang
Baidu Inc., Beijing, China
Rutgers University, New Jersey, USA

Abstract

Stochastic gradient descent is popular for large scale optimization but has slow convergence asymptotically due to the inherent variance. To remedy this problem, we introduce an explicit variance reduction method for stochastic gradient descent which we call stochastic variance reduced gradient (SVRG). For smooth and strongly convex functions, we prove that this method enjoys the same fast convergence rate as those of stochastic dual coordinate ascent (SDCA) and Stochastic Average Gradient (SAG). However, our analysis is significantly simpler and more intuitive. Moreover, unlike SDCA or SAG, our method does not require the storage of gradients, and thus is more easily applicable to complex problems such as some structured prediction problems and neural network learning.

- ▶ SVRG serves as a simple and intuitive way to achieve gradient aggregation
- ▶ This idea has been applied to various other problem settings, including constrained optimization problems and RL

Variance-Reduced and Projection-Free Stochastic Optimization

Elad Hazan

Princeton University, Princeton, NJ 08540, USA

EHAZAN@CS.PRINCETON.EDU

Haipeng Luo

Princeton University, Princeton, NJ 08540, USA

HAIPENGL@CS.PRINCETON.EDU

Abstract

The Frank-Wolfe optimization algorithm has recently regained popularity for machine learning applications due to its projection-free property and its ability to handle structured constraints. However, in the stochastic learning setting, it is still relatively understudied compared to the gradient descent counterpart. In this work, leveraging a recent variance reduction technique, we propose two stochastic Frank-Wolfe variants which substantially improve previous results in terms of the number of stochastic gradient evaluations needed to achieve $1 - \epsilon$ accuracy. For example, we improve from $\mathcal{O}(\frac{1}{\epsilon})$ to $\mathcal{O}(\ln \frac{1}{\epsilon})$ if the objective function is smooth and strongly convex, and from $\mathcal{O}(\frac{1}{\epsilon^2})$ to $\mathcal{O}(\frac{1}{\epsilon^{1.5}})$ if the objective function is smooth and Lipschitz. The theoretical improvement is also observed in experiments on real-world datasets for a multiclass classification application.

more (see for example (Hazan & Kale, 2012; Hazan et al., 2012; Jaggi, 2013; Dudik et al., 2012; Zhang et al., 2012; Harchaoui et al., 2015)).

The Frank-Wolfe algorithm (Frank & Wolfe, 1956) (also known as *conditional gradient*) and its variants are natural candidates for solving these problems, due to its projection-free property and its ability to handle structured constraints. However, despite gaining more popularity recently, its applicability and efficiency in the stochastic learning setting, where computing stochastic gradients is much faster than computing exact gradients, is still relatively understudied compared to variants of projected gradient descent methods.

In this work, we thus try to answer the following question: *what running time can a projection-free algorithm achieve in terms of the number of stochastic gradient evaluations and the number of linear optimizations needed to achieve a certain accuracy?* Utilizing Nesterov's acceleration technique (Nesterov, 1983) and the recent variance reduction idea (Johnson & Zhang, 2013; Mahdavi et al., 2013), we

[Hazan and Luo, ICML 2016]

Stochastic Variance-Reduced Policy Gradient

Matteo Papini *¹ Damiano Binaghi *¹ Giuseppe Canonaco *¹ Matteo Pirotta² Marcello Restelli¹

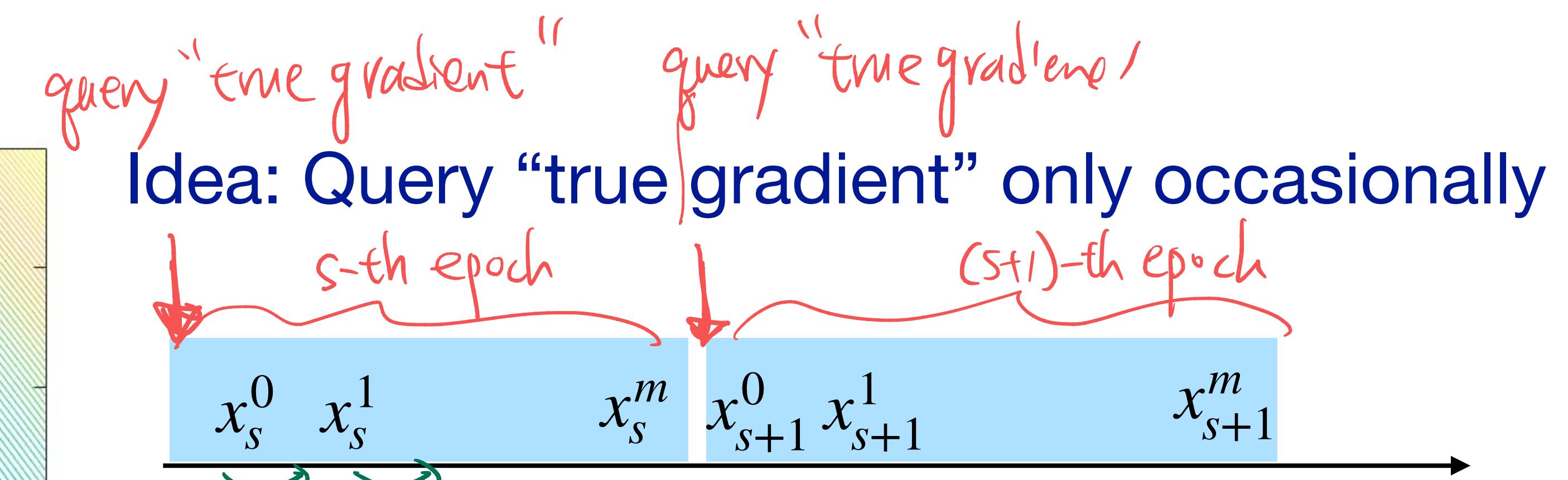
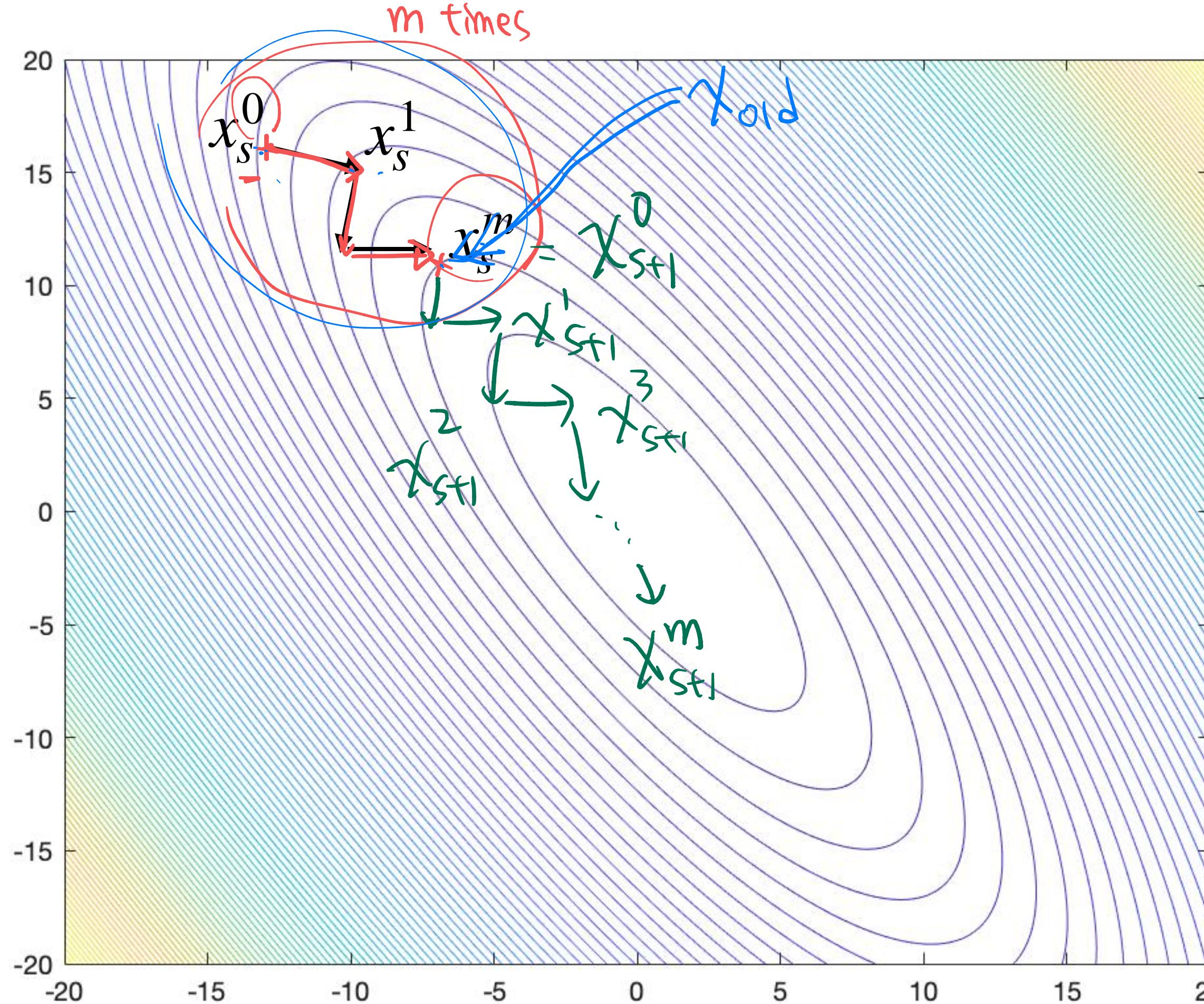
Abstract

In this paper, we propose a novel reinforcement-learning algorithm consisting in a stochastic variance-reduced version of policy gradient for solving Markov Decision Processes (MDPs). Stochastic variance-reduced gradient (SVRG) methods have proven to be very successful in supervised learning. However, their adaptation to policy gradient is not straightforward and needs to account for I) a non-concave objective function; II) approximations in the full gradient computation; and III) a non-stationary sampling process. The result is SVRPG, a stochastic variance-reduced policy gradient algorithm that leverages on importance weights to preserve the unbiasedness of the gradient estimate. Under standard assumptions on the MDP, we provide convergence guarantees for SVRPG with a convergence rate that is linear under increasing batch sizes. Finally, we suggest practical variants of SVRPG, and we empirically evaluate them on continuous MDPs.

a value function, or directly a policy defining the agent's behaviour. Furthermore, when the tasks are characterized by large or continuous state-action spaces, RL needs the powerful function approximators (e.g., neural networks) that are the main subject of study of SL. In a typical SL setting, a performance function $J(\theta)$ has to be optimized w.r.t. to model parameters θ . The set of data that are available for training is often a subset of all the cases of interest, which may even be infinite, leading to optimization of finite sums that approximate the expected performance over an unknown data distribution. When generalization to the complete dataset is not taken into consideration, we talk about Empirical Risk Minimization (ERM). Even in this case, stochastic optimization is often used for reasons of efficiency. The idea of stochastic gradient (SG) ascent (Nesterov, 2013) is to iteratively focus on a random subset of the available data to obtain an approximate improvement direction. At the level of the single iteration, this can be much less expensive than taking into account all the data. However, the sub-sampling of data is a source of variance that can potentially compromise convergence, so that per-iteration efficiency and convergence rate must be traded off

[Papini et al., ICML 2018]

Intuition Behind SVRG



In each stochastic update, use

$$x_s^{k+1} = x_s^k - \eta (\nabla f(x_s^k; d_{I_k}) - \nu_k), \quad I_k \sim \text{Unif}(1, n)$$

$$\nu_k = \nabla f(x_{old}; d_{I_k}) - \nabla F(x_{old})$$

Pseudo Code of SVRG

Notation: $\nabla f(x; d_i) = \nabla f_i(x)$

Outer loop for snapshots

1: **for** $s = 1, 2, \dots$ **do**
2: $x_s^{\text{old}} \leftarrow x_s^j$ and compute $\underbrace{\nabla F(x_s^{\text{old}})}_{(j \sim \text{Unif}(0, \dots, m-1))}$ **batch gradient** // update snapshot

3: initialize $x_s^0 \leftarrow x_s^{\text{old}}$

4: **for** $t = 0, \dots, m-1$ **do** Inner loop for iterate updates

 each epoch contains m iterations

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

$$x_s^{t+1} = x_s^t - \eta \{ \underbrace{\nabla f_{i_t}(x_s^t) - \nabla f_{i_t}(x_s^{\text{old}})}_{\text{stochastic gradient}} + \nabla F(x_s^{\text{old}}) \}$$

► Question: How many gradient computations are needed under SVRG in one epoch?

How about SGD? Per-outer-loop complexity: $(2m+n)$ gradient computations

(Slide Credit: Yuxin Chen)

Comparison: SGD / Batch GD / SVRG

- Consider strongly convex “empirical risk minimization” with n samples

	Iteration Complexity	Per-iteration Cost	Total Computation Cost
Batch GD	$\kappa \log \frac{1}{\epsilon}$	n	$n\kappa \log \frac{1}{\epsilon}$
SGD	$\kappa^2 \frac{1}{\epsilon}$	1	$\kappa^2 \frac{1}{\epsilon}$
SVRG			$(n + \kappa) \log \frac{1}{\epsilon}$

$(\kappa := L/\mu$ is the condition number)

Convergence Rate of SVRG

Theorem Suppose the following conditions hold:

(1) All $f_i(x)$ are convex and L -smooth

(2) $F(x)$ is μ -strongly convex

(3). The snapshot period m is sufficiently large such that

$$\alpha = \frac{1}{\mu \cdot \eta \cdot (1 - 2L\eta) \cdot m} + \frac{2L\eta}{1 - 2L\eta} < 1$$

Then, we have $E[F(x) - F(x^*)] \leq \alpha^s \cdot (F(x_0) - F(x^*)) = \epsilon$

A Lemma for Quantifying Suboptimality $F(x) - F(x^*)$

Lemma:

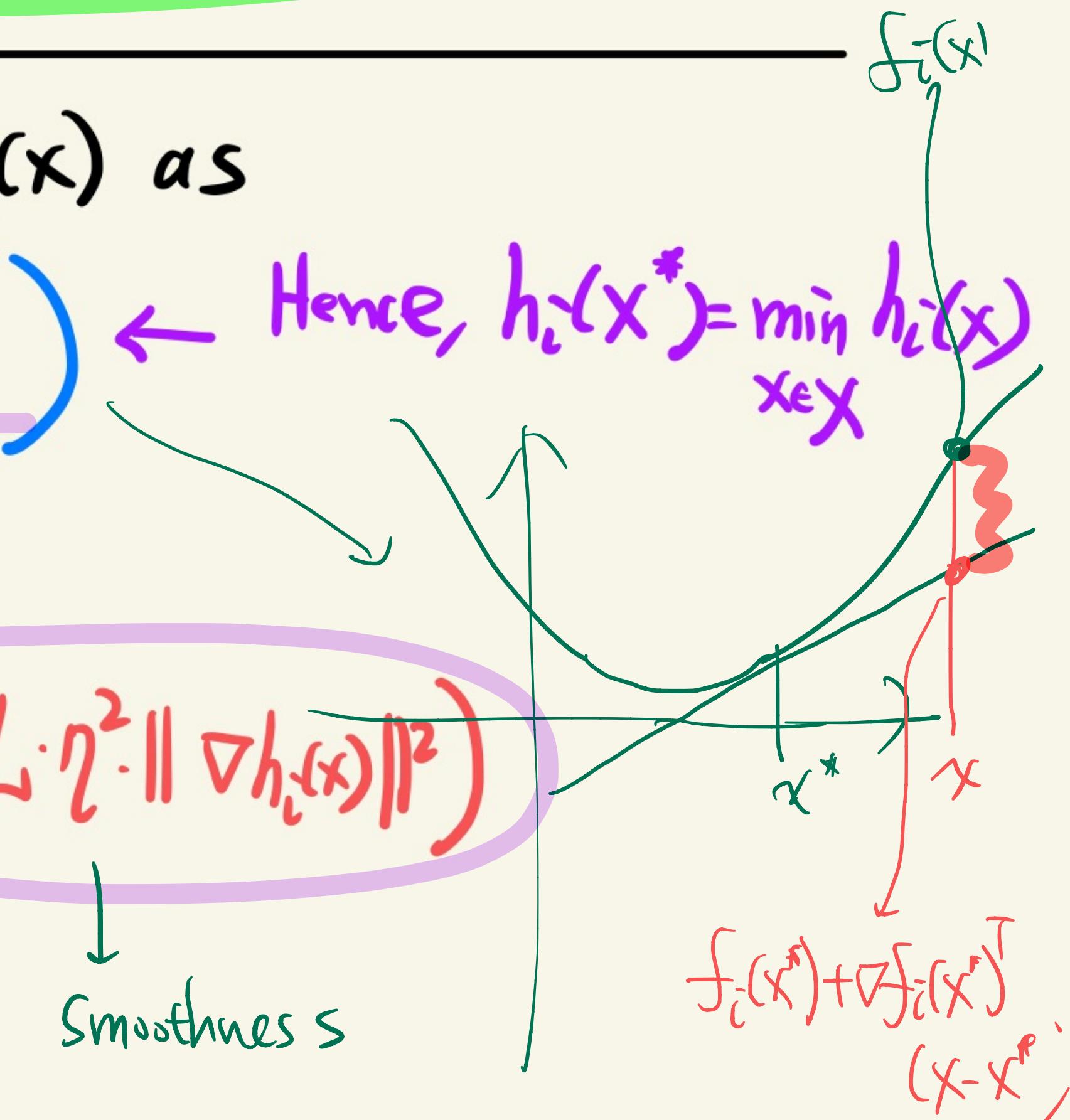
$$\frac{1}{n} \sum_{i=1}^n \| \nabla f_i(x) - \nabla f_i(x^*) \|^2 \leq 2L(F(x) - F(x^*))$$

Proof of Lemma

Let us construct a helper function $h_i(x)$ as

Step 1: $h_i(x) = f_i(x) - \left(f_i(x^*) + \nabla f_i(x^*)^\top (x - x^*) \right)$ ← Hence, $h_i(x^*) = \min_{x \in X} h_i(x)$

Step 2: $0 = h_i(x^*) \leq \min_{\eta} (h_i(x - \eta \cdot \nabla h_i(x)))$
 $\leq \min_{\eta} (h_i(x) - \eta \cdot \|\nabla h_i(x)\|^2 + \frac{1}{2} L \cdot \eta^2 \|\nabla h_i(x)\|^2)$
 $= h_i(x) - \frac{1}{2L} \|\nabla h_i(x)\|^2$



(Cont.).

Step 3: $0 \leq \underline{h_i(x)} - \frac{1}{2L} \|\nabla \underline{h_i(x)}\|^2$ is equivalent to

$$\|\nabla \underline{f_i(x)} - \nabla \underline{f_i(x^*)}\|^2 \leq 2L(f_i(x) - f_i(x^*) - \nabla f_i(x^*)^T(x - x^*)) \quad (*)$$

Step 4: By taking the summation of $(*)$ over $i=1, \dots, n$, we have

$$\frac{1}{n} \sum_{i=1}^n \|\nabla \underline{f_i(x)} - \nabla \underline{f_i(x^*)}\|^2 \leq 2L(F(x) - F(x^*))$$

A Lemma for Quantifying the Variance of SVRG

Lemma: $E[\|g_s^t\|^2] \leq 4 \cdot L \cdot (F(x_s^{t-1}) - F(x^*) + F(x_s^{\text{old}}) - F(x^*))$

Question: If $x_s^t \approx x_s^{\text{old}} \approx x^*$, then what can we say about $E[\|g_s^t\|^2]$?

Proof of Lemma

Recall that $g_s^t = \nabla f_{i_t}(x_s^{t-1}) - \nabla f_{i_t}(x_s^{\text{old}}) + \nabla F(x_s^{\text{old}})$

$$\begin{aligned}
 E[\|g_s^t\|^2] &\leq 2 \cdot E\left[\|\nabla f_{i_t}(x_s^{t-1}) - \nabla f_{i_t}(x^*)\|^2\right] \\
 &\quad + 2 \cdot E\left[\left\|(\nabla f_{i_t}(x_s^{\text{old}}) - \nabla f_{i_t}(x^*)) - \underbrace{\nabla F(x_s^{\text{old}})}_{\text{"}}\right\|^2\right] \\
 &\quad \quad \quad E[\nabla f_{i_t}(x_s^{\text{old}}) - \nabla f_{i_t}(x^*)] \\
 &\leq 2 \cdot E\left[\|\nabla f_{i_t}(x_s^{t-1}) - \nabla f_{i_t}(x^*)\|^2\right] + 2 \cdot E\left[\|\nabla f_{i_t}(x_s^{\text{old}}) - \nabla f_{i_t}(x^*)\|^2\right] \text{ (why?)} \\
 &\leq 4L \cdot \left((F(x_s^{t-1}) - \bar{F}(x^*)) + (F(x_s^{\text{old}}) - \bar{F}(x^*)) \right) \text{ (why?)}
 \end{aligned}$$

A Useful Property

For any random vector \mathbf{z} , we always have

$$E[\|\mathbf{z} - E[\mathbf{z}]\|^2] = E[\|\mathbf{z}\|^2] - \|E[\mathbf{z}]\|^2 \leq E[\|\mathbf{z}\|^2].$$

Let's put everything together and prove the convergence of SVRG

Proof of Convergence of SVRG

Step 1: $E\left[\|x_s^t - x^*\|^2 \mid x_s^{t-1}\right] = E\left[\|x_s^{t-1} - \eta \cdot g_s^t - x^*\|^2 \mid x_s^{t-1}\right]$

$$\begin{aligned} &= \|x_s^{t-1} - x^*\|^2 - 2 \cdot \eta \cdot (x_s^{t-1} - x^*)^\top \cdot E[g_s^t \mid x_s^{t-1}] + \eta^2 \cdot E[\|g_s^t\|^2 \mid x_s^{t-1}] \\ &\leq \|x_s^{t-1} - x^*\|^2 - 2 \cdot \eta \cdot (x_s^{t-1} - x^*)^\top \cdot \underbrace{\nabla F(x_s^{t-1})}_{\geq F(x_s^{t-1}) - F(x^*)} + \eta^2 \cdot 4L(F(x_s^{t-1}) - F(x^*)) \\ &\quad + F(x_s^{\text{old}}) - F(x^*) \\ &\leq \|x_s^{t-1} - x^*\|^2 - 2 \cdot \eta \cdot (F(x_s^{t-1}) - F(x^*)) + \eta^2 \cdot 4L(F(x_s^{t-1}) - F(x^*) + F(x_s^{\text{old}}) \\ &\quad - F(x^*)) \\ &= \|x_s^{t-1} - x^*\|^2 - 2 \cdot \eta \cdot (1 - 2L) \cdot (F(x_s^{t-1}) - F(x^*)) + 4\eta^2 L \cdot (F(x_s^{\text{old}}) - F(x^*)) - (*) \end{aligned}$$

(Cont.)

Step 2: Consider a fixed epoch S

By taking the sum of (*) over $t=1, \dots, m$ and taking the total expectation,

$$E\left[\|x_s^m - x^*\|^2\right] + 2L(1-2\eta)m \cdot E\left[F(x_{S+1}^{old}) - F(x^*)\right]$$

$$\leq E\left[\|x_s^0 - x^*\|^2\right] + 4Lm\eta^2 \cdot E\left[F(x_s^{old}) - F(x^*)\right] \dots (1)$$

$$= E\left[\|x_s^{old} - x^*\|^2\right] + 4Lm\eta^2 \cdot E\left[F(x_s^{old}) - F(x^*)\right] \dots (2)$$

$$\leq \frac{2}{\mu} E\left[F(x_s^{old}) - F(x^*)\right] + 4Lm\eta^2 \cdot E\left[F(x_s^{old}) - F(x^*)\right] \dots (3)$$

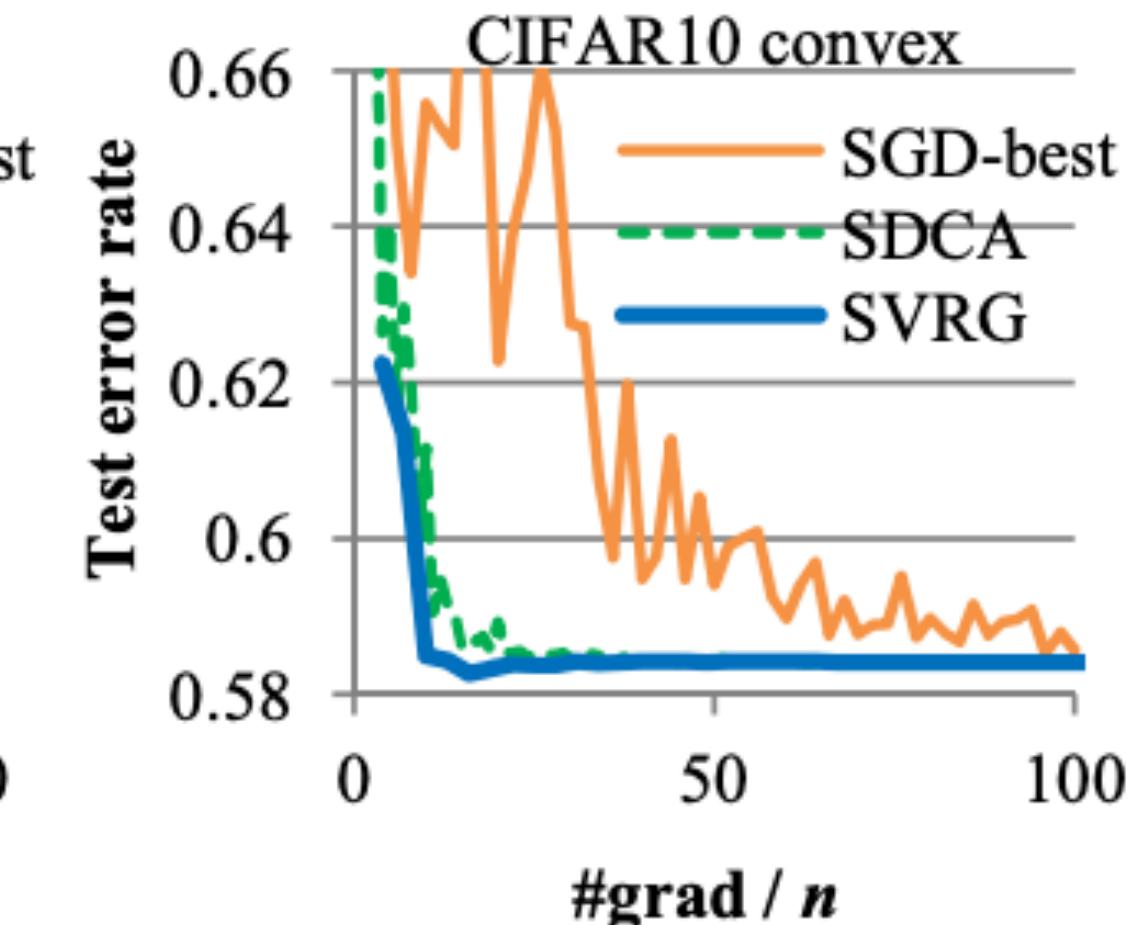
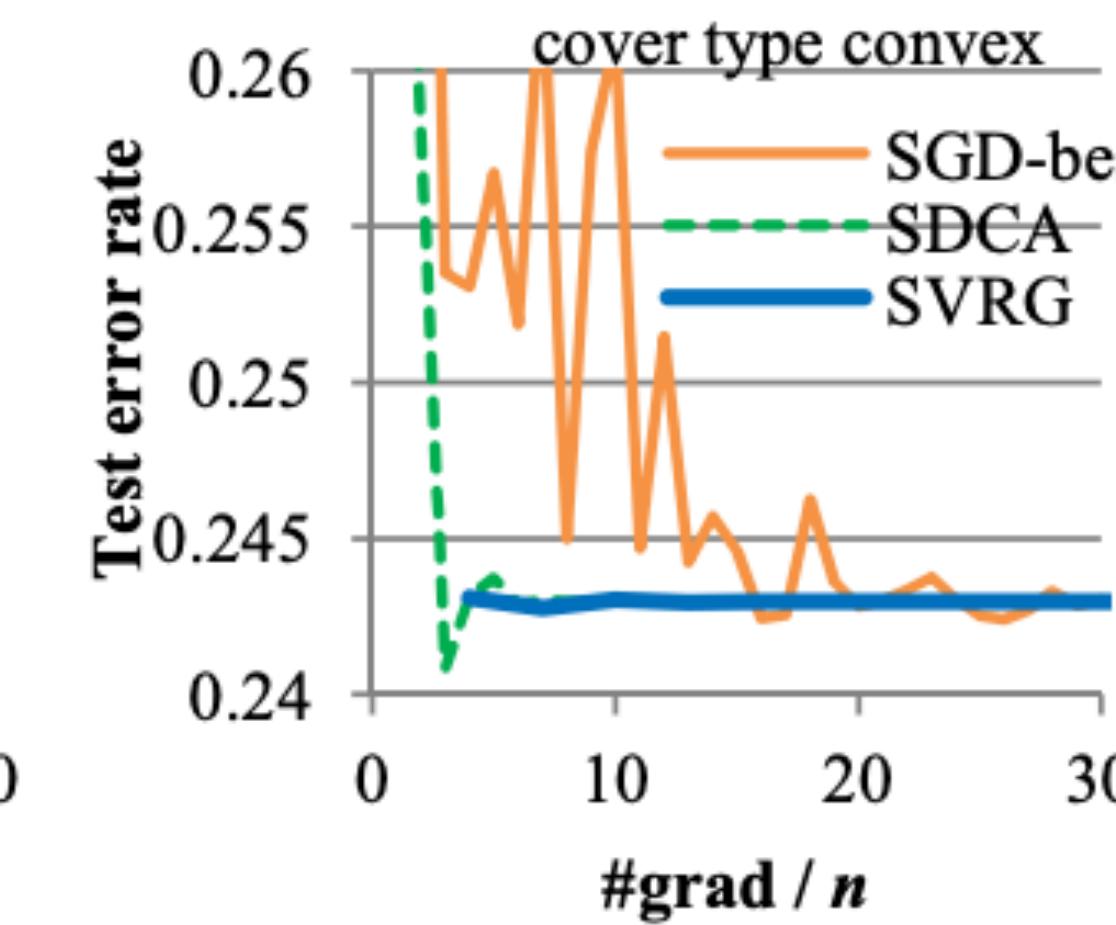
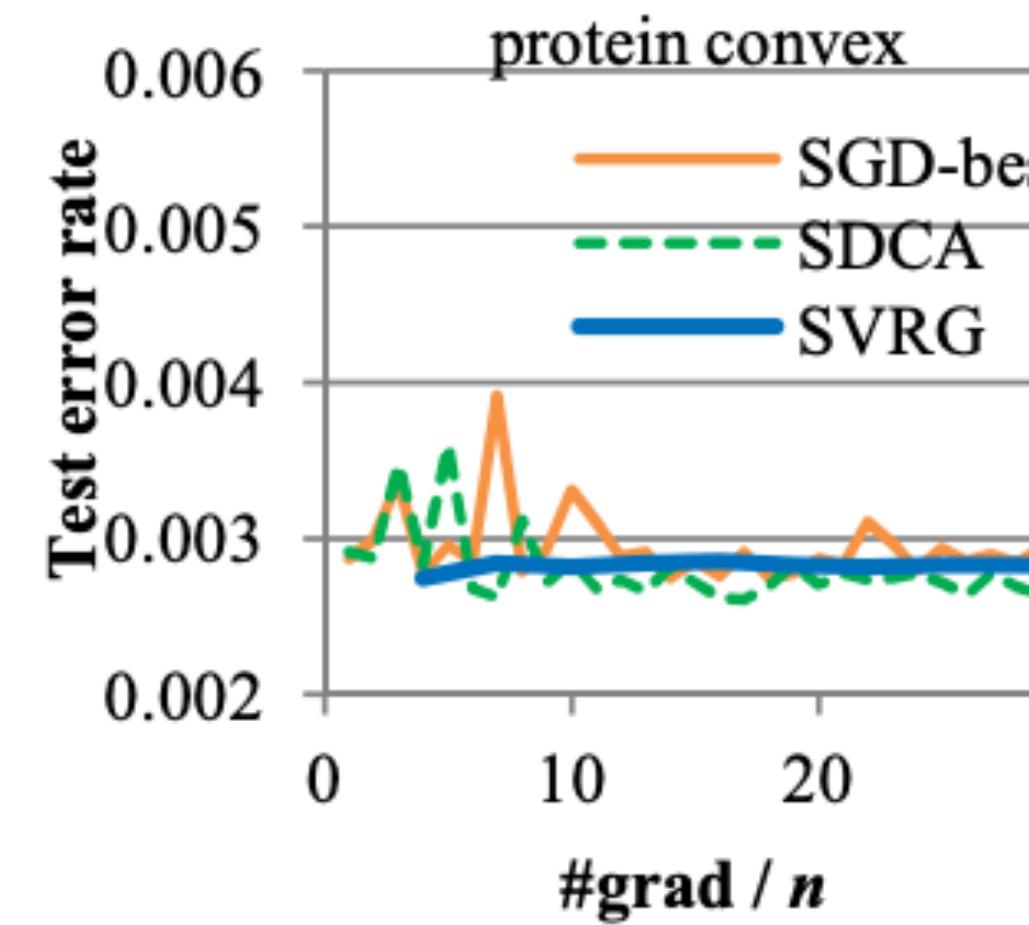
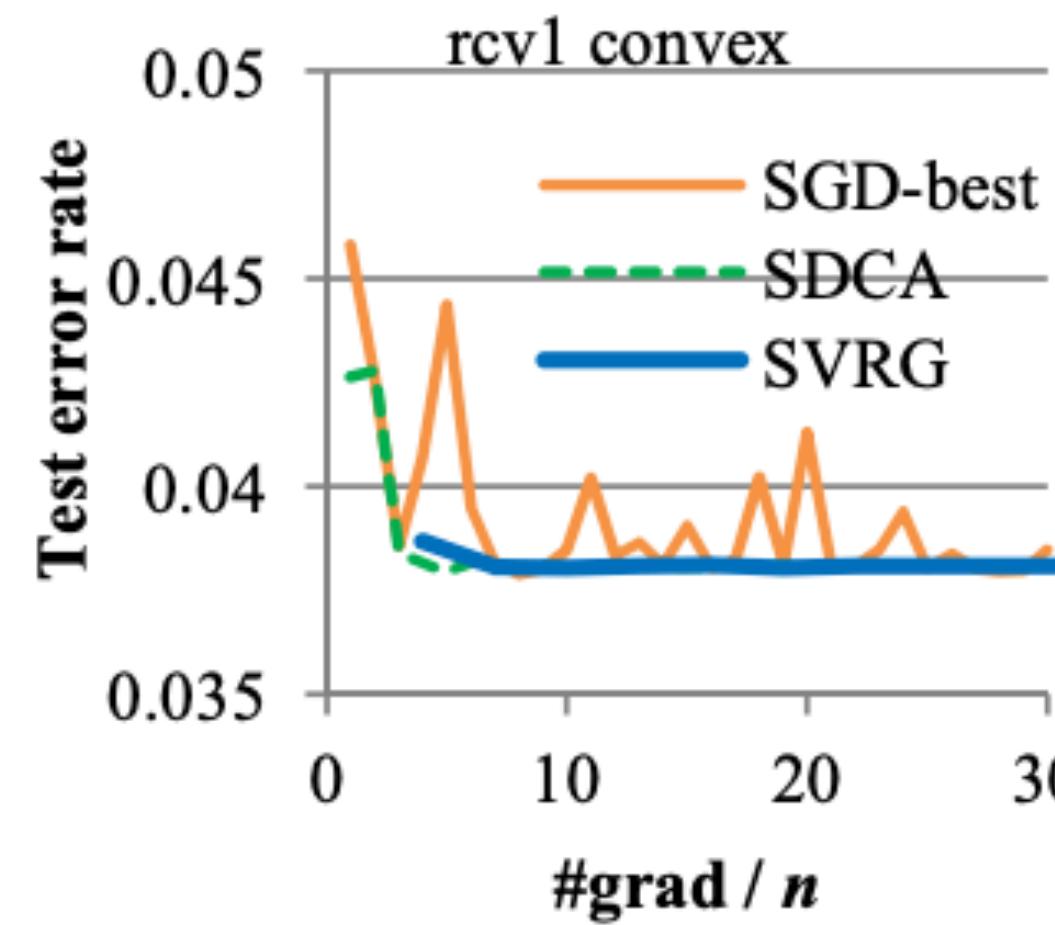
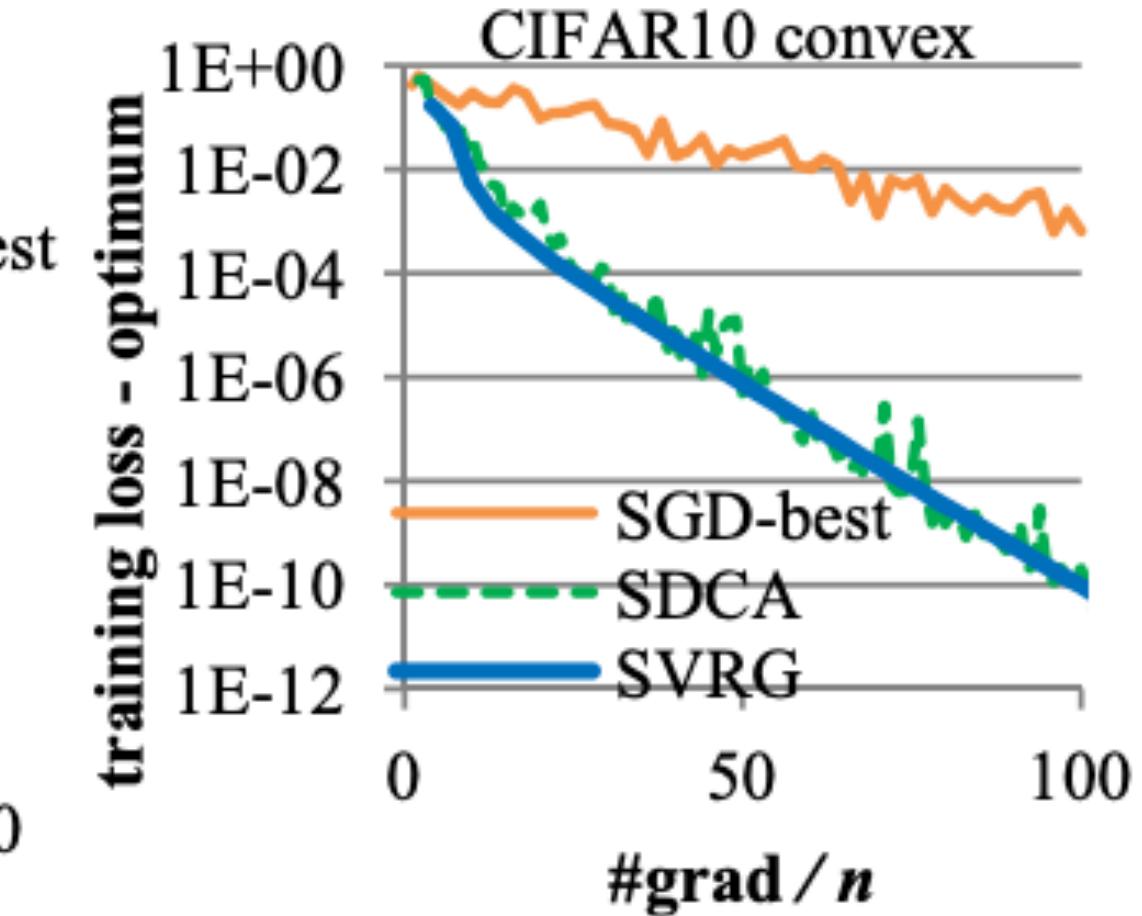
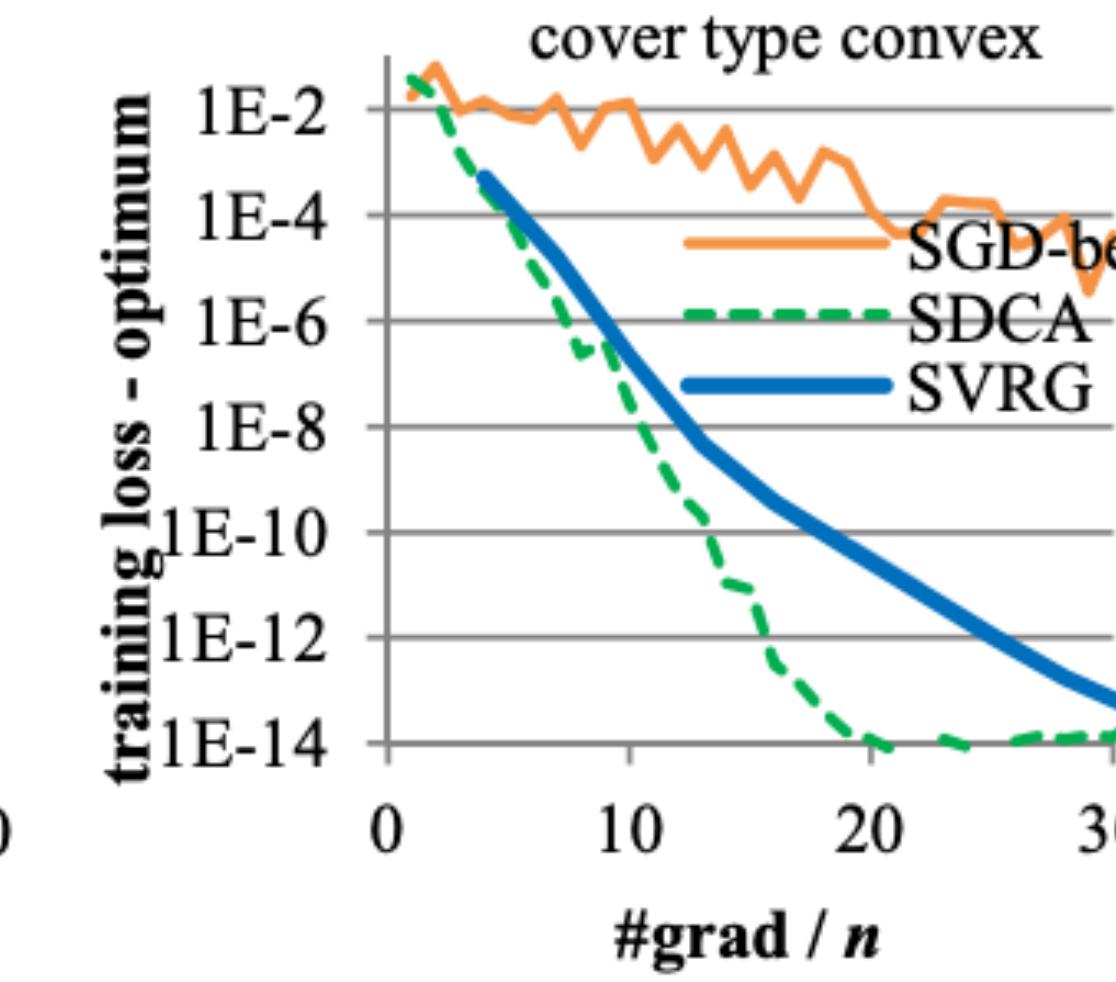
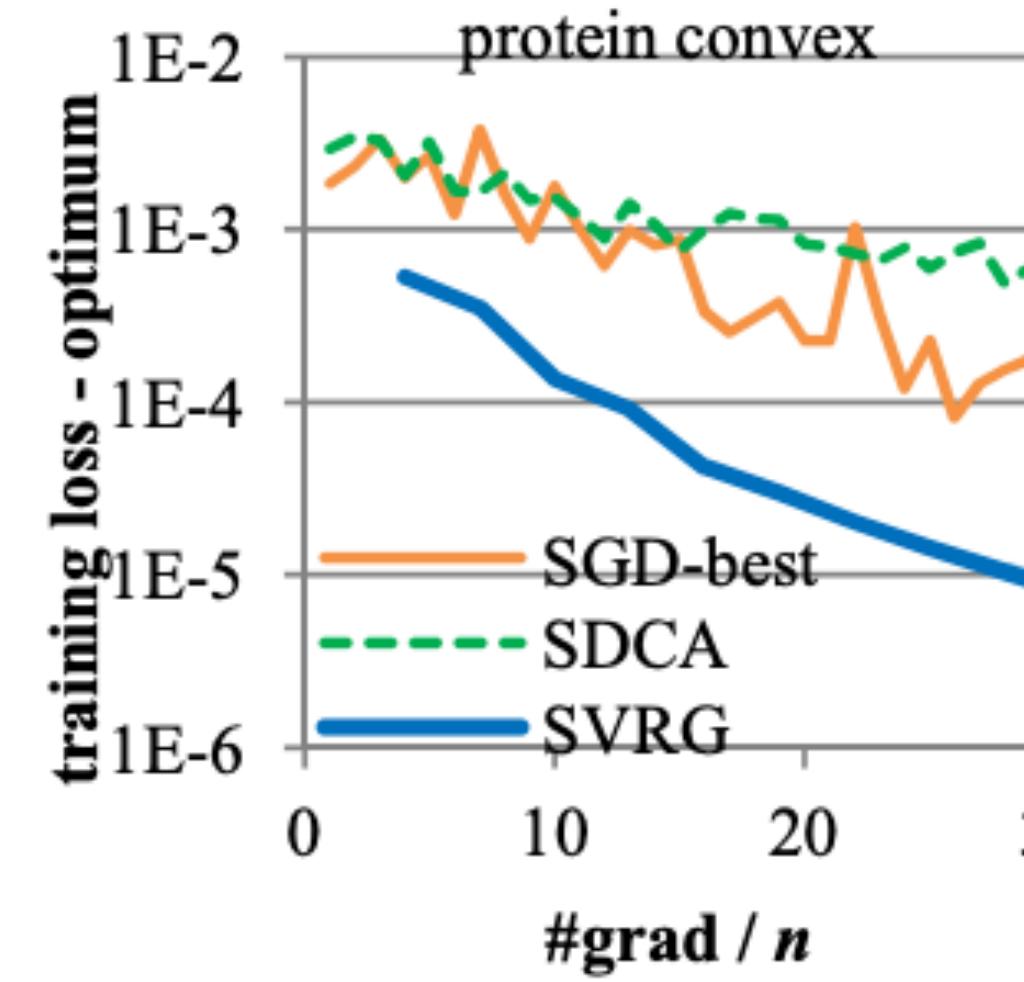
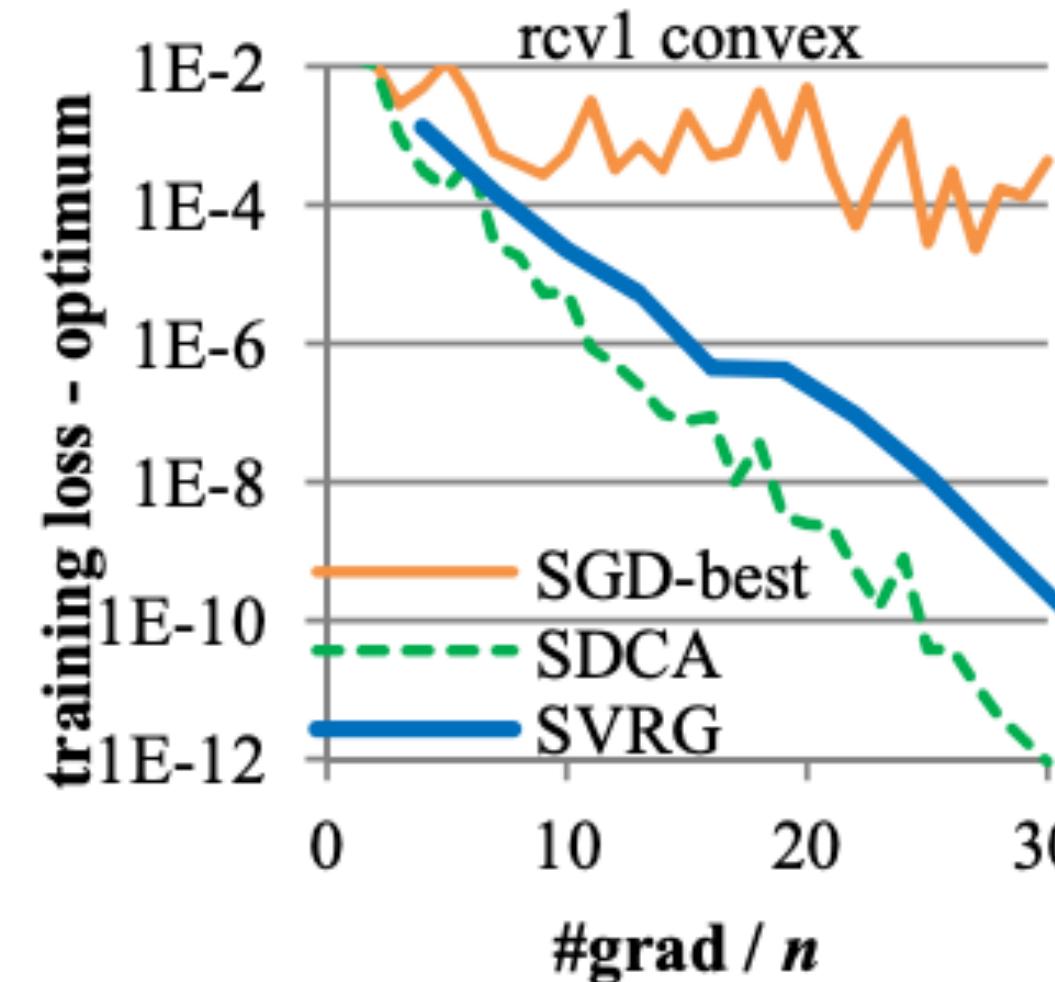
$$= \left(\frac{2}{\mu} + 4Lm\eta^2\right) \cdot E\left[F(x_s^{old}) - F(x^*)\right]$$

(Cont.).

Step 3: Consequently, we have

$$\begin{aligned} & E[F(X_{st}^{\text{old}}) - F(x^*)] \\ & \leq \frac{\frac{2}{\mu} + 4L^2m^2}{2\eta(1-2L\eta)m} \cdot E[F(X_s^{\text{old}}) - F(x^*)] - \underbrace{\frac{1}{2L(1-2L\eta)m} E[\|X_s^m - x^*\|^2]}_{\geq 0} \\ & \leq \left(\frac{1}{\mu\eta(1-2L\eta)m} + \frac{2L^2}{1-2L\eta} \right) \cdot E[F(X_s^{\text{old}}) - F(x^*)] \end{aligned}$$

Numerical Comparison: SGD and SVRG



A General Recipe of Variance Reduction

- ▶ Suppose we'd like estimate $E[X]$ (where X is a random variable)
- ▶ **A general VR approach:** If we could compute efficiently $E[Y]$ for another random variable Y that is highly correlated with X , then we could use an estimator as

$$\theta_\alpha = \alpha(X - Y) + E[Y]$$

- ▶ **Question:** Expected value and variance of θ_α ?

$$\mathbb{E}[\theta_\alpha] = \alpha\mathbb{E}[X] + (1 - \alpha)\mathbb{E}[Y]$$

$$\text{Var}(\theta_\alpha) = \alpha^2(\text{Var}(X) + \text{Var}(Y) - 2\text{Cov}(X, Y))$$

- ▶ By varying α from 0 to 1, what do we have regarding **variance and bias**?
- ▶ How to interpret SVRG?

A General Recipe of Variance Reduction: SAGA vs SVRG

SVRG $x_s^{k+1} = x_s^k - \eta \left(\nabla f(x_s^k; d_{I_k}) - \nabla f(x_{old}; d_{I_k}) - \nabla F(x_{old}) \right), \quad I_k \sim \text{Unif}(1,n)$

SAGA $x^{k+1} = x^k - \eta \left(\nabla f(x^k; d_{I_k}) - \nabla f(\phi_{I_k}^k; d_{I_k}) - \frac{1}{n} \sum_{i=1}^n \nabla f(\phi_i^k; d_i) \right), \quad I_k \sim \text{Unif}(1,n)$

where $\nabla f(\phi_i^k; d_i)$ are past stochastic gradients stored in the buffer

General recipe: $\theta_\alpha = \alpha(X - Y) + E[Y]$. What are X, Y, α in. SAGA?

Many More VR Methods

- ▶ **SAG (Stochastic Averaged Gradient)**
- ▶ **SAGA [Defazio, Bach, and Lacoste-Julien, NIPS 2014]**
- ▶ **SARAH (StochAstic Recursive grAdient algoritHm) [Nguyen, Liu, Scheinberg, Takac, ICML 2017]**
- ▶ **SDCA (Stochastic Dual Coordinate Ascent) [Johnson and Zhang, NIPS 2013; Shalev-Shwartz, ICML 2016]**
- ▶ **And more!**

More VR Methods: SARAH

Question: Can we do VR without taking a snapshot?

- ▶ StochAStic RecursiVe grAdient algoritHm (SARAH)

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

$$g^t = (\nabla f_{I_t}(x^t) - \nabla f_{I_t}(x^{t-1}) + g^{t-1}$$

Idea: Recursive updates of gradient estimates

- ▶ Question: Why is this reasonable? Any issue?

Pseudo Code of SARAH

Notation: $\nabla f(x; d_i) = \nabla f_i(x)$

```
1: for  $s = 1, 2, \dots, S$  do
2:    $x_s^0 \leftarrow x_{s-1}^{m+1}$ , 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           Outer loop for periodic resets
5:     choose  $i_t$  uniformly from  $\{1, \dots, n\}$ 
6:      $\mathbf{g}_s^t = \underbrace{\nabla f_{i_t}(\mathbf{x}_s^t) - \nabla f_{i_t}(\mathbf{x}_s^{t-1})}_{\text{stochastic gradient}} + \mathbf{g}_s^{t-1}$  Inner loop for recursive updates
7:      $\mathbf{x}_s^{t+1} = \mathbf{x}_s^t - \eta \mathbf{g}_s^t$ 
```

► **Question:** Why do we need to reset \mathbf{g}^t periodically?

Next Lecture: How to accelerate SGD for expected loss?