

ECE 8101: Nonconvex Optimization for Machine Learning

Lecture Note 2-4: Stochastic Gradient Descent

Jia (Kevin) Liu

Associate Professor
Department of Electrical and Computer Engineering
The Ohio State University, Columbus, OH, USA

Autumn 2024

Outline

In this lecture:

- Noisy unbiased gradient
- Stochastic gradient method
- Convergence results

Unbiased Stochastic Gradient

- Random vector $\tilde{\mathbf{g}} \in \mathbb{R}^n$ is a **unbiased stochastic gradient** if it can be written as $\tilde{\mathbf{g}} = \mathbf{g} + \mathbf{n}$ where \mathbf{g} is the true gradient and $\mathbb{E}[\mathbf{n}] = \mathbf{0}$
- \mathbf{n} can be interpreted as error in computing \mathbf{g} , measurement noise, Monte Carlo sampling errors, etc.
- If $f(\cdot)$ is non-smooth, $\tilde{\mathbf{g}}$ is a noisy unbiased subgradient at \mathbf{x} if

$$f(\mathbf{z}) \geq f(\mathbf{x}) + (\mathbb{E}[\tilde{\mathbf{g}}|\mathbf{x}])^\top (\mathbf{z} - \mathbf{x}), \quad \forall \mathbf{z}$$

holds almost surely.

\uparrow
convex $f(\cdot)$

\Downarrow
convex

Stochastic Gradient Descent Method

- Consider $\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x})$. Following standard GD, we should do:

$$\mathbf{x}_{k+1} = \mathbf{x}_k - s_k \mathbb{E}[\tilde{\mathbf{g}}_k | \mathbf{x}_k]$$

(Handwritten red note: $= \nabla f(\mathbf{x}_k)$)

- However, $\mathbb{E}[\tilde{\mathbf{g}}_k | \mathbf{x}_k]$ is **difficult** to compute: Unknown distribution, too costly to sample at each iteration k , etc.
- Idea:** Simply use a noisy unbiased subgradient to replace $\mathbb{E}[\tilde{\mathbf{g}}_k | \mathbf{x}_k]$
- The **stochastic subgradient** method works as follows:

$$\mathbf{x}_{k+1} = \mathbf{x}_k - s_k \tilde{\mathbf{g}}_k$$

- ▶ \mathbf{x}_k is the k -th iterate
- ▶ $\tilde{\mathbf{g}}_k$ is any noisy gradient of at \mathbf{x}_k , i.e., $\mathbb{E}[\tilde{\mathbf{g}}_k | \mathbf{x}_k] = \nabla f(\mathbf{x}_k)$
- ▶ s_k is the step size
- ▶ Let $f_{\text{best}}^{(k)} \triangleq \min_{i=1, \dots, k} \{f(\mathbf{x}_i)\}$ and $\|\nabla f_{\text{best}}^{(k)}\| \triangleq \min_{i=1, \dots, k} \{\|\nabla f(\mathbf{x}_i)\|\}$

Historical Perspective

- Also referred to as **stochastic approximation** in the literature, first introduced by [Robbins, Monro '51] and [Keifer, Wolfowitz '52]
- The original work [Robbins, Monro '51] is motivated by finding a root of a continuous function:

vector-valued

$$f(\mathbf{x}) = \mathbb{E}[F(\mathbf{x}, \theta)] = 0,$$

where $F(\cdot, \cdot)$ is **unknown** and depends on a random variable θ . But the experimenter can take random samples (noisy measurements) of $F(\mathbf{x}, \theta)$



Herbert Robbins



Sutton Monro

Historical Perspective

- **Robbins-Monro:** $\mathbf{x}_{k+1} = \mathbf{x}_k + s_k Y(\mathbf{x}_k, \theta)$, where:
 - ▶ $\mathbb{E}[Y(\mathbf{x}, \theta) | \mathbf{x} = \mathbf{x}_k] = f(\mathbf{x}_k)$ is an unbiased estimator of $f(\mathbf{x}_k)$
 - ▶ Robbins-Monro originally showed convergence in L^2 and in probability
 - ▶ Blum later prove convergence is actually w.p.1. (almost surely)
 - ▶ **Key idea:** Diminishing step-size provides **implicit averaging** of the observations
- Robbins-Monro's scheme can also be used in **stochastic optimization** of the form $f(\mathbf{x}^*) = \min_{\mathbf{x}} \mathbb{E}[F(\mathbf{x}, \theta)]$ (equivalent to solving $\nabla f(\mathbf{x}^*) = 0$)
- Stochastic approximation, or more generally, stochastic gradient has found applications in many areas
 - ▶ Adaptive signal processing
 - ▶ Dynamic network control and optimization
 - ▶ Statistical machine learning
 - ▶ Workhorse algorithm for training **deep neural networks**

Convergence of R.V.

1. Convergence in Distr. (weak convergence)

A seq. of (real-valued) r.v. $\{X_n\}$ converges in distr. to X if $\lim_{n \rightarrow \infty} F_n(x) = F(x)$, where F_n and F are cdf of X_n and X , resp. Denote as $X_n \xrightarrow{D} X$.

2. Convergence in prob. to a r.v. ("stronger")

$\{X_n\}$ converges in prob. to a r.v. X if $\forall \epsilon > 0$,

$$\lim_{n \rightarrow \infty} \Pr\{|X_n - X| > \epsilon\} = 0. \quad \text{Denote as: } X_n \xrightarrow{P} X.$$

3. Almost sure convergence (pt.-wise convergence in Real Analysis)

$\{X_n\}$ converges a.s. (a.e. or w.p.1, or strongly) to X .

$$\Pr\left\{\lim_{n \rightarrow \infty} X_n = X\right\} = 1. \quad \text{Denoted as } X_n \xrightarrow{\text{a.s.}} X.$$

4. Convergence in expectation = Given $r \geq 1$, $\{X_n\}$ converges

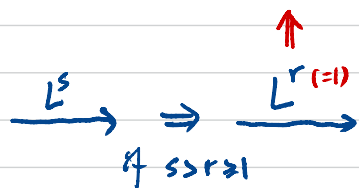
in r -th mean to r.v. X if the r -th abs. moments

$$\mathbb{E}[|X_n|^r] \text{ and } \mathbb{E}[|X|^r] \text{ exist, and}$$

$$\lim_{n \rightarrow \infty} \mathbb{E}[|X_n - X|^r] = 0. \quad \text{Denote as } X_n \xrightarrow{L^r} X$$

* $r=1$: X_n converges in mean to X .

* $r=2$: ————— mean square to X .



Markov Ineq: X : non-neg. r.v. For some $a > 0$

$$\Pr\{X \geq a\} \leq \frac{E[X]}{a}$$

* For r.v. Z_1, \dots, Z_n that are indep. with mean 0.

$$E[\|Z_1 + \dots + Z_n\|^2] \leq E[\|Z_1\|^2 + \dots + \|Z_n\|^2] \quad \left(\begin{array}{l} \text{impdres} \\ E[\|X_i\|^2] < \infty \end{array} \right)$$

* For r.v. Z_1, \dots, Z_n that are not nec. indep.

$$E[\|Z_1 + \dots + Z_n\|^2] \leq n E[\|Z_1\|^2 + \dots + \|Z_n\|^2]$$

Assumptions and Step Size Rules

- $f^* = \inf_x f(\mathbf{x}_k) > -\infty$, with $f(\mathbf{x}^*) = f^*$
- $\mathbb{E}[\|\tilde{\mathbf{g}}_k\|_2^2] \leq G^2$, for all k
- $\mathbb{E}[\|\mathbf{x}_0 - \mathbf{x}^*\|_2^2] \leq R^2$

Commonly used step-size strategies:

- Constant step-size: $s_k = s, \forall k$
- Step-size is square summable, but not summable

$$s_k > 0, \forall k, \quad \sum_{k=1}^{\infty} s_k^2 < \infty, \quad \sum_{k=1}^{\infty} s_k = \infty$$

not needed

Note: This is stronger than needed, but just to simplify proof

Convergence of SGD (Convex)

Asymptotic

- Convergence in expectation:

$$\lim_{k \rightarrow \infty} \mathbb{E}[f_{\text{best}}^{(k)}] = f^*$$

- Convergence in probability: for any $\epsilon > 0$,

$$\lim_{k \rightarrow \infty} \Pr\{|f_{\text{best}}^{(k)} - f^*| > \epsilon\} = 0$$

- Almost sure convergence

$$\Pr\left\{\lim_{k \rightarrow \infty} f_{\text{best}}^{(k)} = f^*\right\} = 1$$

- See [Kushner, Yin '97] for a complete treatment on convergence analysis

(Convex)

Then: If $\mathbb{E}[\|\tilde{g}_k\|_2] \leq G$, $\forall k$, $\mathbb{E}\{\|z_k - z^*\|\} \leq R$ and step-sizes $\{s_k\}_{k=1}^{\infty}$ satisfy: $s_k > 0$, $\forall k$, $\sum_{k=1}^{\infty} s_k^2 < \infty$, $\sum_{k=1}^{\infty} s_k \rightarrow \infty$.

then: $\lim_{k \rightarrow \infty} \mathbb{E}\{f_{\text{best}}^{(k)}\} = f^*$ and $\lim_{k \rightarrow \infty} \mathbb{P}\{|f_{\text{best}}^{(k)} - f^*| > \epsilon\} = 0, \forall \epsilon > 0$.

Proof. Consider cond. expectation square Euclidean dist:

$$\begin{aligned} \mathbb{E}[\|z_{k+1} - z^*\|^2 | z_k] &= \mathbb{E}[\|z_k - s_k \tilde{g}_k - z^*\|^2 | z_k] \\ &= \mathbb{E}[\|z_k - z^*\|^2 + s_k^2 \|\tilde{g}_k\|^2 - 2s_k \tilde{g}_k^T (z_k - z^*) | z_k] \\ &= \|z_k - z^*\|^2 + s_k^2 \mathbb{E}[\|\tilde{g}_k\|^2 | z_k] - 2s_k \underbrace{\mathbb{E}[\tilde{g}_k | z_k]^T}_{= \nabla f(z_k)} (z_k - z^*) \quad (*) \end{aligned}$$

Note: $f(z^*) \geq f(z_k) + \mathbb{E}\{\tilde{g}_k | z_k\}^T (z^* - z_k)$

$$\Rightarrow -\mathbb{E}[\tilde{g}_k | z_k]^T (z_k - z^*) \leq -(f(z_k) - f^*)$$

$$(*) \leq \|z_k - z^*\|^2 + s_k^2 \mathbb{E}[\|\tilde{g}_k\|^2 | z_k] - 2s_k (f(z_k) - f^*)$$

Note: from SGP dynamic, z_{k+1} only dep. z_k . ($z_{k+1} = z_k - s_k \tilde{g}_k$)
and indep. of z_{k-1}, \dots, z_1 .

$$\mathbb{E}[\|z_{k+1} - z^*\|^2 | z_k] = \mathbb{E}[\|z_{k+1} - z^*\|^2 | z_k, \dots, z_1]$$

Take expectation over joint distr. of $\{z_k, \dots, z_1\}$ yields

$$\begin{aligned} \mathbb{E}[\|z_{k+1} - z^*\|^2] &\leq \mathbb{E}[\|z_k - z^*\|^2] - 2s_k \mathbb{E}[f(z_k) - f^*] \\ &\quad + s_k^2 \mathbb{E}[\|\tilde{g}_k\|^2] \\ &\quad \leq G \end{aligned}$$

Apply this process recursively:

$$\mathbb{E}[\|z_{k+1} - z^*\|^2] \leq \underbrace{\mathbb{E}[\|z_1 - z^*\|^2]}_{\leq R^2} - 2 \sum_{i=1}^k s_i (\underbrace{\mathbb{E}[f(z_i)] - f^*}_{\geq \min_{i=1 \dots k} f(z_i) - f^*}) + G^2 \underbrace{\sum_{k=1}^k s_k^2}_{\leq B}$$

$$\Rightarrow \min_{i=1 \dots k} \{ \mathbb{E}[f(z_i)] - f^* \} \leq \frac{R^2 + G^2 B}{2 \underbrace{\sum_{i=1}^k s_i}_{\rightarrow \infty}} \rightarrow 0 \text{ as } k \rightarrow \infty$$

Claim: The fn $g(y) \triangleq \min_{i=1 \dots k} \{ y_i \}$ is concave. (HW)

Then Jensen's inequality:

* If f is convex: $f(\mathbb{E}X) \leq \mathbb{E}f(X)$

* --- concave: \geq

$$\mathbb{E}[f_{\text{best}}^{(k)}] = \mathbb{E}[\underbrace{\min_{i=1 \dots k} f(z_i)}_{\text{concave}}] \stackrel{\text{Jensen's}}{\leq} \min_{i=1 \dots k} \mathbb{E}[f(z_i)] \rightarrow f^*$$

\uparrow
 done earlier



Convergence in Expectation and Probability (Convex)

Proof Sketch:

- **Key quantity:** **Expected** squared Euclidean distance to the optimal set. Let \mathbf{x}^* be any minimizer of f . We can show that

$$\mathbb{E}[\|\mathbf{x}_{k+1} - \mathbf{x}^*\|_2^2 | \mathbf{x}_k] \leq \|\mathbf{x}_k - \mathbf{x}^*\|_2^2 - 2s_k(f(\mathbf{x}_k) - f^*) + s_k^2 \mathbb{E}[\|\tilde{\mathbf{g}}_k\|_2^2 | \mathbf{x}_k]$$

- which can further lead to

$$\min_{i=1, \dots, k} \left\{ \mathbb{E}[f(\mathbf{x}_i)] - f^* \right\} \leq \frac{R^2 + G^2 \|s\|^2}{2 \sum_{i=1}^k s_i}$$

- The result $\min_{i=1, \dots, k} \mathbb{E}[f(\mathbf{x}_i)] \rightarrow f^*$ simply follows from the divergent step-size series rule

Convergence in Expectation and Probability (Convex)

- Jensen's inequality and concavity of minimum yields

$$\mathbb{E}[f_{\text{best}}^{(k)}] = \mathbb{E}\left[\min_{i=1,\dots,k} f(\mathbf{x}_i)\right] \leq \min_{i=1,\dots,k} \mathbb{E}[f(\mathbf{x}_i)]$$

Therefore, $\mathbb{E}[f_{\text{best}}^{(k)}] \rightarrow f^*$ (convergence in expectation)

- Convergence in expectation also implies convergence in probability: By Markov's inequality, for any $\epsilon > 0$,

$$\Pr\{f_{\text{best}}^{(k)} - f^* \geq \epsilon\} \leq \frac{\mathbb{E}[f_{\text{best}}^{(k)} - f^*]}{\epsilon},$$

i.e., RHS goes to 0, which proves convergence in probability. □

Convergence Rate (Convex)

$$\lim_{k \rightarrow \infty} \mathbb{E} \left\{ \min_{i=1, \dots, k} f(x_i) - f^* \right\} \leq \frac{R^2 + G \sum_{k=1}^{\infty} s_k^2}{2 \sum_{k=1}^{\infty} s_k}$$

(Nicola Cresima)

$$1 + \frac{1}{2} + \frac{1}{3} + \frac{1}{4} + \frac{1}{5} + \dots$$

$$= 1 + \frac{1}{2} + (\frac{1}{3} + \frac{1}{4}) + (\frac{1}{5} + \frac{1}{6} + \frac{1}{7} + \frac{1}{8}) + \dots$$

$$> 1 + \frac{1}{2} + (\frac{1}{4} + \frac{1}{4}) + (\frac{1}{8} + \frac{1}{8} + \frac{1}{8} + \frac{1}{8}) + \dots$$

$$= 1 + \frac{1}{2} + \frac{1}{2} + \frac{1}{2} + \dots = \infty$$

$\Theta(\log(t))$.

- Classical diminishing step-sizes $s_k = \alpha/k$ for some $\alpha > 0$:

$\sum_k s_k = O(\log(t))$ and $\sum_k s_k^2 = O(1)$. So convergence rate is $O(1/\log(t))$

$\int_1^n f(x) dx \leq \sum_{n=1}^{\infty} f(n) \leq f(1) + \int_1^{\infty} f(x) dx$

$\frac{1}{1^2} + \frac{1}{2^2} + \frac{1}{3^2} + \dots = \frac{\pi^2}{6}$

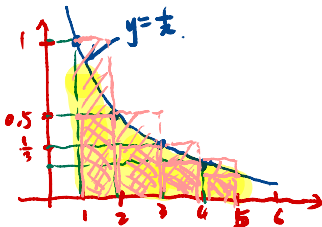
integral test

- Diminishing step-sizes $s_k = \alpha/\sqrt{k}$ for some $\alpha > 0$: $\sum_k s_k = O(\sqrt{t})$ and $\sum_k s_k^2 = O(\log(t))$. So convergence rate is $O(\log(t)/\sqrt{t}) = \tilde{O}(1/\sqrt{t})$

- Constant step-sizes $s_k = \alpha$ for some $\alpha > 0$: $\sum_k s_k = k\alpha$ and $\sum_k s_k^2 = k\alpha^2$. So convergence rate is $O(1/t) + O(\alpha)$

$\int_1^n \frac{1}{x} dx = \log(n)$

$\sum_{k=1}^n \frac{1}{k+1} \leq \int_1^n \frac{1}{x} dx \leq \sum_{k=1}^n \frac{1}{k}$



Convergence Rate (Strongly Convex)

Theorem 1 (Optimality Gap)

If $f(\cdot)$ is μ -strongly convex, then the SGD method with a constant step-size $s_k = s < 2/\mu$ satisfies:

$$\mathbb{E}[\|\mathbf{x}_k - \mathbf{x}^*\|^2] \leq (1 - 2s\mu)^k \|\mathbf{x}_0 - \mathbf{x}^*\|^2 + \frac{s\sigma^2}{2\mu}$$

Remark:

- If $\sigma^2 = 0$ (GD), constant step-size $s \Rightarrow$ linear convergence to \mathbf{x}^* .
- If $\sigma^2 > 0$, SGD with constant step-size $s \Rightarrow$ linear convergence to $\frac{s\sigma^2}{2\mu}$ -neighborhood of \mathbf{x}^*

strongly convex:

$$f(y) \geq f(x) + \nabla f(x)^T (y-x) + \frac{\mu}{2} \|y-x\|^2$$

$$f(x) \geq f(y) + \nabla f(y)^T (x-y) + \frac{\mu}{2} \|y-x\|^2$$

Add together: $[\nabla f(y) - \nabla f(x)]^T (y-x) \geq \mu \|y-x\|^2$

Annotations: $\mathbb{E}[\tilde{g}_k] \rightarrow \nabla f(x_k)$, $\nabla f(x^) = 0$*

Recall: $\mathbb{E}[\|z_{k+1} - x^*\|^2 | z_k] \leq \|z_k - x^*\|^2 + s_k^2 \mathbb{E}[\|\tilde{g}_k\|^2 | z_k] - 2s_k \mathbb{E}[\tilde{g}_k^T | z_k]^T (z_k - x^*)$

Taking full expectation:

$$\mathbb{E}[\|z_{k+1} - x^*\|^2] \leq \mathbb{E}[\|z_k - x^*\|^2] + s_k^2 \mathbb{E}[\|\tilde{g}_k\|^2] - 2\mu s_k \mathbb{E}[\|z_k - x^*\|^2]$$

Annotation: $\leq \sigma^2$

$$\mathbb{E}[\tilde{g}_k^T (z_k - x^*) | z_k] \geq \mu \|z_k - x^*\|^2$$

$$\rightarrow = (1 - 2\mu s_k) \mathbb{E}[\|z_k - x^*\|^2] + s_k^2 \sigma^2 \quad (1) \quad s < \frac{2}{\mu}$$

Applying (1) recursively from $k-1$ down to 1, letting $s_k = s$ with

$$\mathbb{E}[\|z_k - x^*\|^2] \leq (1 - 2\mu s)^k \|z_0 - x^*\|^2 + \frac{s\sigma^2}{2\mu} \quad \square$$

(HW). What about diminishing step-size?

Convergence Rate (Nonconvex) – Finite Sum

- Consider the following finite-sum minimization

$$\min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}) = \min_{\mathbf{x} \in \mathbb{R}^d} \frac{1}{N} \sum_{i=1}^N f_i(\mathbf{x})$$

where N is typically large, e.g., empirical risk minimization (ERM) in ML

- Consider using SGD to solve this problem under the following assumptions:
 - $f(\cdot)$ is nonconvex and bounded from below
 - ∇f is differentiable with L -Lipschitz continuous gradients (L -smooth)
 - $\mathbb{E}[\|\nabla f_i(\mathbf{x})\|^2] \leq \sigma^2$ for some σ^2 and all \mathbf{x} (bounded gradient, can be relaxed)

can be relaxed: $\mathbb{E}[\|\nabla f_i(\mathbf{z}_k) - \nabla f(\mathbf{z}_k)\|^2] \leq \sigma^2$

can be further relaxed: $\mathbb{E}[\|\nabla f_i(\mathbf{z}_k) - \nabla f(\mathbf{z}_k)\|^2] \leq \sigma \|\nabla f(\mathbf{z}_k)\|^2$
"SNR $\geq \frac{1}{\sigma}$ "

Convergence Rate (Nonconvex) – Finite Sum

Theorem 2 (Stationarity Gap)

If the finite-sum problem $f(\cdot)$ is nonconvex, differentiable, and L -smooth, then the SGD method with step-sizes $\{s_k\}$ satisfies

$$\min_{k=0,1,\dots,t-1} \mathbb{E} \{ \|\nabla f(\mathbf{x}_k)\|_2^2 \} \leq \frac{f(\mathbf{x}_0) - f^*}{\sum_{k=0}^{t-1} s_k} + \frac{L\sigma^2}{2} \frac{\sum_{k=0}^{t-1} s_k^2}{\sum_{k=0}^{t-1} s_k}.$$

Remark:

- If $\sigma^2 = 0$, then a constant step-size yields an $O(1/t)$ rate.
- Classical diminishing step-sizes $s_k = \alpha/k$ for some $\alpha > 0$:
 $\sum_k s_k = O(\log(t))$ and $\sum_k s_k^2 = O(1)$. So convergence rate is $O(1/\log(t))$
- Diminishing step-sizes $s_k = \alpha/\sqrt{k}$ for some $\alpha > 0$: $\sum_k s_k = O(\sqrt{t})$ and $\sum_k s_k^2 = O(\log(t))$. So convergence rate is $O(\log(t)/\sqrt{t}) = \tilde{O}(1/\sqrt{t})$
- Constant step-sizes $s_k = \alpha$ for some $\alpha > 0$: $\sum_k s_k = k\alpha$ and $\sum_k s_k^2 = k\alpha^2$. So convergence rate is $O(1/t) + O(\alpha)$

Theorem 2 (Stationarity Gap)

If the finite-sum problem $f(\cdot)$ is nonconvex, differentiable, and L -smooth, then the SGD method with step-sizes $\{s_k\}$ satisfies

$$\min_{k=0,1,\dots,t-1} \mathbb{E} \{\|\nabla f(\mathbf{x}_k)\|_2^2\} \leq \frac{f(\mathbf{x}_0) - f^*}{\sum_{k=0}^{t-1} s_k} + \frac{L\sigma^2}{2} \frac{\sum_{k=0}^{t-1} s_k^2}{\sum_{k=0}^{t-1} s_k}.$$

Proof: Consider uniform sampling $i_k \in \{1, \dots, N\}$ with $P_r(i_k = i) = \frac{1}{N}$.

$$\mathbf{z}_{k+1} = \mathbf{z}_k - s_k \nabla f_{i_k}(\mathbf{z}_k).$$

$$\mathbb{E}[\nabla f_{i_k}(\mathbf{z}_k)] = \sum_{i=1}^N P_r(i_k = i) \nabla f_i(\mathbf{z}_k) = \frac{1}{N} \sum_{i=1}^N \nabla f_i(\mathbf{z}_k) = \nabla f(\mathbf{z}_k).$$

Recall the descent lemma in GD.

$$f(\mathbf{z}_{k+1}) \leq f(\mathbf{z}_k) + \nabla f(\mathbf{z}_k)^T (\mathbf{z}_{k+1} - \mathbf{z}_k) + \frac{L}{2} \|\mathbf{z}_{k+1} - \mathbf{z}_k\|^2$$

Plugging in SGD iteration yields:

$$f(\mathbf{z}_{k+1}) \leq f(\mathbf{z}_k) - s_k \nabla f(\mathbf{z}_k)^T \nabla f_{i_k}(\mathbf{z}_k) + \frac{L s_k^2}{2} \|\nabla f_{i_k}(\mathbf{z}_k)\|^2$$

Take expectation w.r.t. i_k

$$\begin{aligned} \mathbb{E}[f(\mathbf{z}_{k+1})] &\leq \mathbb{E}\left[f(\mathbf{z}_k) - s_k \nabla f(\mathbf{z}_k)^T \nabla f_{i_k}(\mathbf{z}_k) + \frac{L s_k^2}{2} \|\nabla f_{i_k}(\mathbf{z}_k)\|^2\right] \\ &= \mathbb{E}[f(\mathbf{z}_k)] - s_k \|\nabla f(\mathbf{z}_k)\|^2 + \frac{L s_k^2}{2} \underbrace{\mathbb{E}[\|\nabla f_{i_k}(\mathbf{z}_k)\|^2]}_{\leq \sigma^2} \\ &\leq \underbrace{\mathbb{E}[f(\mathbf{z}_k)]}_{\text{good}} - \underbrace{s_k \|\nabla f(\mathbf{z}_k)\|^2}_{\text{bad}} + \frac{L s_k^2}{2} \sigma^2. \end{aligned}$$

As in GD: rearrange to get the grad norm on LHS:

$$s_k \|\nabla f(x_k)\|^2 \leq \mathbb{E}[f(x_k)] - \mathbb{E}[f(x_{k+1})] + \frac{L s_k^2}{2} \sigma^2 \quad (2)$$

Summing (2) from 1 to t & use iterated expectation to get:

$$\sum_{k=1}^{t-1} s_{k-1} \mathbb{E}[\|\nabla f(x_{k-1})\|^2] \leq \underbrace{\sum_{k=1}^t [\mathbb{E}[f(x_{k-1})] - \mathbb{E}[f(x_k)]]}_{\text{telescope.}} + \frac{L \sigma^2}{2} \sum_{k=0}^{t-1} s_k^2$$

$\geq \min_{k=0, \dots, t-1} \{\mathbb{E}[\|\nabla f(x_{k-1})\|^2]\}$

$$\Rightarrow \min_{k=0, \dots, t-1} \{\mathbb{E}[\|\nabla f(x_k)\|^2]\} \leq \frac{f(x_0) - f(x^*)}{\sum_{k=0}^{t-1} s_k} + \frac{L \sigma^2}{2} \frac{\sum_{k=0}^{t-1} s_k^2}{\sum_{k=0}^{t-1} s_k} \quad \square$$

Convergence Rate (Nonconvex) - Finite Sum+Time Oracle

Theorem 3 ([Ghadimi & Lan '13])

Suppose $f(\cdot)$ is L -smooth and has σ -bounded gradients and it is known a priori that the SGD algorithm will be executed for T iterations. Let $s_k = c/\sqrt{T}$, where

$$c = \sqrt{\frac{2(f(\mathbf{x}_0) - f^*)}{L\sigma^2}}.$$

Then, the iterates of SGD satisfy

$$\min_{0 \leq t \leq T-1} \mathbb{E}[\|\nabla f(\mathbf{x}_t)\|^2] \leq \sqrt{\frac{2(f(\mathbf{x}_0) - f^*)L}{T}}\sigma.$$

Theorem 3 ([Ghadimi & Lan '13])

Suppose $f(\cdot)$ is L -smooth and has σ -bounded gradients and it is known a priori that the SGD algorithm will be executed for T iterations. Let $s_k = c/\sqrt{T}$, where

$$c = \sqrt{\frac{2(f(x_0) - f^*)}{L\sigma^2}}.$$

Then, the iterates of SGD satisfy

$$\min_{0 \leq t \leq T-1} \mathbb{E}[\|\nabla f(x_t)\|^2] \leq \sqrt{\frac{2(f(x_0) - f^*)L}{T}} \sigma.$$

Proof: We have shown:

$$\min_{k=0, \dots, T-1} \left\{ \mathbb{E}[\|\nabla f(x_k)\|^2] \right\} \leq \frac{f(x_0) - f^*}{\sum_{k=0}^{T-1} s_k} + \frac{L\sigma^2}{2} \frac{\sum_{k=0}^{T-1} s_k^2}{\sum_{k=0}^{T-1} s_k} \quad (c)$$

$$s_k = \frac{c}{\sqrt{T}} \Rightarrow \sum_{k=0}^{T-1} s_k = T \cdot \frac{c}{\sqrt{T}} = c\sqrt{T}$$

$$\sum_{k=0}^{T-1} s_k^2 = T \cdot \frac{c^2}{T} = c^2$$

$$(c) \Rightarrow \min_{k=0, \dots, T-1} \left\{ \mathbb{E}[\|\nabla f(x_k)\|^2] \right\} \leq \frac{f(x_0) - f^*}{c\sqrt{T}} + \frac{L\sigma^2}{2} \cdot \frac{c^2}{c\sqrt{T}}$$

$$= \frac{1}{\sqrt{T}} \left(\frac{f(x_0) - f^*}{c} + \frac{L\sigma^2 c}{2} \right) \quad \left(\frac{2\sigma f + L\sigma^2 c^2}{2c} = O(\sqrt{T}) \right)$$

Young's Ineq. $ab \leq \frac{a^p}{p} + \frac{b^q}{q}$, with $\frac{1}{p} + \frac{1}{q} = 1$.

By picking "a" = \sqrt{a} , "b" = \sqrt{b} , $p=q=2$, $a+b \geq 2\sqrt{ab}$
 \uparrow "=" iff $a=b$.

By forcing $\frac{f(x_0) - f^*}{c} = \frac{L\sigma^2 c}{2} \Rightarrow c = \sqrt{\frac{2(f(x_0) - f^*)}{L\sigma^2}}$

Then, the stated result follows. 

Convergence Rate (Nonconvex) - General Expectation Minimization with Batching

- Consider the following general expectation minimization problem

$$f(\mathbf{x}) = \mathbb{E}_{\xi}[f(\mathbf{x}, \xi)],$$

where ξ is a random variable with distribution \mathcal{D} .

- Consider using SGD to solve this problem under the following assumptions:
 - $f(\cdot)$ is nonconvex and bounded from below
 - ∇f is differentiable with L -Lipschitz continuous gradients (L -smooth)
 - $\mathbb{E}_{\xi}[\nabla f(\mathbf{x}, \xi)] = \nabla f(\mathbf{x})$ and $\mathbb{E}_{\xi}[\|\nabla f(\mathbf{x}, \xi) - \nabla f(\mathbf{x})\|_2^2] \leq \sigma^2$
- A common approach in SGD: Rather than choosing one training sample randomly at a time, use a **larger random mini-batch of samples** \mathcal{B}_k , with $|\mathcal{B}_k| = B_k$. Then, $\mathbf{g}_k = \frac{1}{B_k} \sum_{i=1}^{B_k} \nabla f(\mathbf{x}, \xi_i)$. SGD becomes:

$$\mathbf{x}_{k+1} = \mathbf{x}_k - s_k \mathbf{g}_k = \mathbf{x}_k - \frac{s_k}{B_k} \sum_{i=1}^{B_k} \nabla f(\mathbf{x}, \xi_i),$$

where ξ_1, \dots, ξ_{B_k} are i.i.d. sampled from \mathcal{D}

Convergence Rate (Nonconvex) - General Expectation Minimization with Batching

Theorem 4 (Stationarity Gap)

In the expectation minimization problem, supposed that $f(\cdot)$ is nonconvex, differentiable, and L -smooth. For any given $\epsilon > 0$, then the SGD method with mini-batch size $B_k = B = \max\{1, \frac{2\sigma^2}{\epsilon^2}\}$, $\forall k$, and step-sizes $s_k \leq \frac{1}{2L}$, $\forall k$, satisfies

$$\mathbb{E}[\|\nabla f(\hat{\mathbf{x}}_t)\|_2^2] \leq \frac{4L(f(\mathbf{x}_0) - f^*)}{t} + \frac{\epsilon^2}{2}, \quad (1)$$

where $\hat{\mathbf{x}}_t$ is chosen uniformly at random from $\mathbf{x}_0, \dots, \mathbf{x}_{t-1}$. Thus, Eq. (1) implies that taking $t = \lceil \frac{8L(f(\mathbf{x}_0) - f^*)}{\epsilon^2} \rceil$ yields $\mathbb{E}[\|\nabla f(\hat{\mathbf{x}}_t)\|_2^2] \leq \epsilon^2$.

Sample Complexity Bound:

$$\sum_{k=0}^{t-1} B_k = \frac{2\sigma^2}{\epsilon^2} t = \left\lceil \frac{16L(f(\mathbf{x}_0) - f^*)\sigma^2}{\epsilon^4} \right\rceil = O(\epsilon^{-4})$$

- **Optimal** up to constant factors (see [Arjevani et al. 2019] for lower bound)

Theorem 4 (Stationarity Gap)

In the expectation minimization problem, supposed that $f(\cdot)$ is nonconvex, differentiable, and L -smooth. For any given $\epsilon > 0$, then the SGD method with mini-batch size $B_k = B = \max\{1, \frac{2\sigma^2}{\epsilon^2}\}$, $\forall k$, and step-sizes $s_k \leq \frac{1}{2L}$, $\forall k$, satisfies

$$\mathbb{E}[\|\nabla f(\hat{x}_t)\|_2^2] \leq \frac{4L(f(x_0) - f^*)}{t} + \frac{\epsilon^2}{2}, \quad (1)$$

where \hat{x}_t is chosen uniformly at random from x_0, \dots, x_{t-1} . Thus, Eq. (1) implies that taking $t = \lceil \frac{8L(f(x_0) - f^*)}{\epsilon^2} \rceil$ yields $\mathbb{E}[\|\nabla f(\hat{x}_t)\|_2^2] \leq \epsilon^2$.

Proof: ① WTS: When $B_k = B = \max\{1, \frac{2\sigma^2}{\epsilon^2}\}$, we have

$$\mathbb{E}[\|g(x) - \nabla f(x)\|^2 | x] \leq \frac{\sigma^2}{2}.$$

Note: $g(x) = \frac{1}{B} \sum_{i=1}^B \nabla f_i(x, \xi_i)$, where ξ_1, \dots, ξ_B are i.i.d. sampled from \mathcal{D} .

$$\mathbb{E}[\|g(x) - \nabla f(x)\|^2 | x] = \mathbb{E}\left[\left\|\frac{1}{B} \sum_{i=1}^B \nabla f_i(x, \xi_i) - \nabla f(x)\right\|^2 | x\right]$$

$$\begin{aligned} & \stackrel{\text{expand indep.}}{=} \mathbb{E}\left[\left\|\frac{1}{B} \sum_{i=1}^B [\nabla f_i(x, \xi_i) - \nabla f(x)]\right\|^2 | x\right] \\ & \leq \frac{1}{B} \sum_{i=1}^B \underbrace{\mathbb{E}_{\xi_i}[\|\nabla f_i(x, \xi_i) - \nabla f(x)\|^2 | x]}_{\leq \sigma^2} \leq \frac{\sigma^2}{B} \stackrel{\text{[2\sigma^2]}}{\leq} \frac{\sigma^2}{2}. \end{aligned}$$

Recall "Descent lemma".

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

$$\begin{aligned} & \stackrel{\text{add \& subtract}}{=} f(x_k) + \underbrace{g_k^T}_{s_k g_k} (x_{k+1} - x_k) + (\nabla f(x_k) - g_k)^T (x_{k+1} - x_k) + \frac{L}{2} \|x_{k+1} - x_k\|^2 \end{aligned}$$

$$\leq \|r(\beta_k) - g_k\|^2 + \frac{1}{4s_k} \|z_{k+1} - z_k\|^2 \stackrel{\leq s_k g_k}{\leq} \frac{1}{4s_k} \|z_{k+1} - z_k\|^2 \quad (\text{by prop. } \frac{d}{dt} \frac{1}{2} \text{ in F-Y inequality.})$$

$$= f(z_k) - s_k \|g_k\|^2 + s_k \underbrace{(r(z_k) - g_k)^T (z_{k+1} - z_k)} + \frac{L}{2} s_k^2 \|g_k\|^2 \quad (1)$$

Fenchel-Young's Ineq: $\underline{a}^T \underline{b} \leq \frac{1}{2\alpha} \|\underline{a}\|^2 + \frac{\alpha}{2} \|\underline{b}\|^2$

"Convex Conjugate":

Let X be topo. space, and X^* be dual space

$$\langle \cdot, \cdot \rangle: X^* \times X \rightarrow \mathbb{R}$$

Def (Convex Conjugate): For a fn $f: X \rightarrow \mathbb{R} \cup \{-\infty, +\infty\}$, its convex conjugate is the fn: $f^*: X^* \rightarrow \mathbb{R} \cup \{-\infty, +\infty\}$, where at $z^* \in X^*$, we have:

$$f^*(z^*) \triangleq \sup \{ \langle z^*, z \rangle - f(z) : z \in X \}$$

$$\Rightarrow z \in X, p \in X^*, \underbrace{\langle p, z \rangle}_{\text{inner}} \leq \underbrace{f(z)}_{\|\cdot\|_2} + \underbrace{f^*(p)}_{\|\cdot\|_2}$$

$$(1) \Rightarrow f(z_{k+1}) \leq f(z_k) - s_k \|g_k\|^2 + s_k \|r(z_k) - g_k\|^2 + \left(\frac{1}{4s_k} + \frac{L}{2}\right) s_k^2 \|g_k\|^2$$

$$= f(z_k) - s_k \left[1 - \left(\frac{1}{4} + \frac{Ls_k}{2}\right) \right] \|g_k\|^2 + s_k \|r(z_k) - g_k\|^2 \quad (2)$$

$$\text{Since: } s_k \leq \frac{1}{2L} \Rightarrow Ls_k \leq \frac{1}{2} \Rightarrow \frac{Ls_k}{2} \leq \frac{1}{4} \Rightarrow \frac{1}{4} + \frac{Ls_k}{2} \leq \frac{1}{2}$$

$$\Rightarrow -\left(\frac{1}{4} + \frac{Ls_k}{2}\right) \geq -\frac{1}{2} \Rightarrow -\left[1 - \left(\frac{1}{4} + \frac{Ls_k}{2}\right)\right] \leq -\frac{1}{2}$$

$$\text{Then, (2)} \Rightarrow f(z_{k+1}) \leq f(z_k) - \underbrace{\frac{s_k}{2} \|g_k\|^2}_{\text{good}} + \underbrace{s_k \|r(z_k) - g_k\|^2}_{\text{bad}}$$

$$\mathbb{E}[f(\mathbf{x}_{k+1}) | \mathbf{x}_k] \leq f(\mathbf{x}_k) - \frac{s_k}{2} \underbrace{\left[\|\mathbf{g}_k\|^2 \right]}_{\pm \nabla f(\mathbf{x}_k)} + s_k \mathbb{E}[\|\mathbf{g}_k - \nabla f(\mathbf{x}_k)\|^2 | \mathbf{x}_k]$$

add
subtract

$$= f(\mathbf{x}_k) - \frac{s_k}{2} \left[\|\nabla f(\mathbf{x}_k)\|^2 + \mathbb{E}[\|\mathbf{g}_k - \nabla f(\mathbf{x}_k)\|^2 | \mathbf{x}_k] \right]$$

$$+ s_k \mathbb{E}[\|\mathbf{g}_k - \nabla f(\mathbf{x}_k)\|^2 | \mathbf{x}_k]$$

$$= f(\mathbf{x}_k) - \frac{s_k}{2} \|\nabla f(\mathbf{x}_k)\|^2 + \frac{s_k}{2} \mathbb{E}[\|\mathbf{g}_k - \nabla f(\mathbf{x}_k)\|^2 | \mathbf{x}_k]$$

$\leq \frac{\epsilon^2}{2}$

$\leq \frac{s_k \epsilon^2}{4}$

(3)

Taking full expectation on both sides, choosing $s_k \leq \frac{1}{2t}$,
 summary (3) for $k=0, \dots, t-1$, we have:

$$\frac{1}{t} \sum_{k=0}^{t-1} \mathbb{E}[\|\nabla f(\mathbf{x}_k)\|^2] \leq \frac{4L}{t} \sum_{k=0}^{t-1} (\mathbb{E}[f(\mathbf{x}_k)] - \mathbb{E}[f(\mathbf{x}_{k+1})]) + \frac{\epsilon}{2}$$

$$= \frac{4L}{t} (f(\mathbf{x}_0) - \underbrace{\mathbb{E}[f(\mathbf{x}_t)]}_{\leq -f^*}) + \frac{\epsilon}{2}$$

$$\leq \frac{4L}{t} (f(\mathbf{x}_0) - f^*) + \frac{\epsilon}{2}$$

Finally, choosing output $\hat{\mathbf{x}}$ uniformly at random from $\{\mathbf{x}_0, \dots, \mathbf{x}_{t-1}\}$

we have $\mathbb{E}[\|\nabla f(\hat{\mathbf{x}}_t)\|^2] \leq \frac{4L(f(\mathbf{x}_0) - f^*)}{t} + \frac{\epsilon}{2}$ ▣

Mini-Batching SGD as Gradient Descent with Error

- SGD with mini-batch:

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \frac{s_k}{B_k} \sum_{i=1}^{B_k} \nabla f(\mathbf{x}, \xi_i)$$

- This can be viewed as a “gradient descent with error”

$$\mathbf{x}_{k+1} = \mathbf{x}_k - s_k(\nabla f(\mathbf{x}_k) + \mathbf{e}_k)$$

, where \mathbf{e}_k is the difference between approximation and true gradient

- By setting $s_k = 1/L$, it follows from descent lemma that

$$f(\mathbf{x}_{k+1}) \leq f(\mathbf{x}_k) - \underbrace{\frac{1}{2L} \|\nabla f(\mathbf{x}_k)\|^2}_{\text{good}} + \underbrace{\frac{1}{2L} \|\mathbf{e}_k\|^2}_{\text{bad}}$$

Mini-Batching SGD as Gradient Descent with Error

- SGD progress bound with $s_k = 1/L$ and error is:

$$f(\mathbf{x}_{k+1}) \leq f(\mathbf{x}_k) - \underbrace{\frac{1}{2L} \|\nabla f(\mathbf{x}_k)\|^2}_{\text{good}} + \underbrace{\frac{1}{2L} \|\mathbf{e}_k\|^2}_{\text{bad}}$$

- Relationship between “error-free” rate and “with error” rate:
 - ▶ If “error-free” rate is $O(1/k)$, you maintain this rate if $\|\mathbf{e}_k\|^2 = O(1/k)$
 - ▶ If “error-free” rate is $O(\rho^k)$, you maintain this rate if $\|\mathbf{e}_k\|^2 = O(\rho^k)$
 - ▶ If error goes to zero more slowly, error vanishing rate is the “bottleneck”
- So, need to know how batch-size B_k affects $\|\mathbf{e}_k\|^2$

Mini-Batching SGD as Gradient Descent with Error

- Sample with replacement:

$$\mathbb{E}[\|\mathbf{e}_k\|^2] = \frac{1}{B_k} \sigma^2,$$

where σ^2 is the variance of the stochastic gradient norm (i.e., doubling the batch-size cuts the error in half)

- Sample without replacement (from a dataset of size N):

$$\mathbb{E}[\|\mathbf{e}_k\|^2] = \frac{N - B_k}{N - 1} \frac{1}{B_k} \sigma^2,$$

i.e., driving error to zero as batch size approaches N

- Growing batch-size:

- ▶ For $O(\rho^k)$ linear convergence: need $B_{k+1} = B_k/\rho$
- ▶ For $O(1/k)$ sublinear convergence: need $B_{k+1} = B_k + \text{const.}$

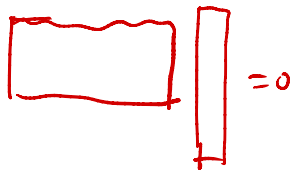
Mini-Batching SGD as Gradient Descent with Error

- SGD with mini-batch:

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \frac{s_k}{B_k} \sum_{i=1}^{B_k} \nabla f(\mathbf{x}, \xi_i)$$

- For a fixed B_k : sublinear convergence rate
 - ▶ Fixed step-size: sublinear convergence to an error ball around a stationary point
 - ▶ Diminishing step-size: sublinear convergence to a stationary point
- Can grow B_k to achieve faster rate:
 - ▶ Early iterations: cheap SG iterations
 - ▶ Later iterations: Use larger batch-sizes (no need to play with step-sizes)

Next Class



Variance-Reduced First-Order Methods

