

# Stochastic Gradient Descent



Cong Ma

University of Chicago, Winter 2026

# Outline

---

- Stochastic gradient descent (stochastic approximation)
- Convergence analysis

# Stochastic Programming

---

$$\text{minimize}_{\mathbf{x}} \quad \underbrace{F(\mathbf{x}) = \mathbb{E}[f(\mathbf{x}; \Xi)]}_{\text{expected risk, population risk, ...}}$$

- $\Xi$  denotes the random source in the objective.
- Assume  $f(\cdot, \Xi)$  is convex for every realization of  $\Xi$ ; then  $F$  is convex.

## Example: Empirical Risk Minimization

---

Let  $\{\mathbf{a}_i, y_i\}_{i=1}^n$  be  $n$  random samples, and consider

$$\text{minimize}_{\mathbf{x}} \quad \underbrace{F(\mathbf{x}) := \frac{1}{n} \sum_{i=1}^n f(\mathbf{x}; \{\mathbf{a}_i, y_i\})}_{\text{empirical risk}}$$

e.g. quadratic loss  $f(\mathbf{x}; \{\mathbf{a}_i, y_i\}) = (\mathbf{a}_i^\top \mathbf{x} - y_i)^2$

If one draws index  $j \sim \text{Unif}(1, \dots, n)$  uniformly at random, then

$$F(\mathbf{x}) = \mathbb{E}_j[f(\mathbf{x}; \{\mathbf{a}_j, y_j\})]$$

## A Natural Solution

---

Under mild regularity conditions, we can interchange gradient and expectation:

$$\begin{aligned}\mathbf{x}^{t+1} &= \mathbf{x}^t - \eta_t \nabla F(\mathbf{x}^t) \\ &= \mathbf{x}^t - \eta_t \nabla \mathbb{E}[f(\mathbf{x}^t; \Xi)] \\ &= \mathbf{x}^t - \eta_t \mathbb{E}[\nabla_{\mathbf{x}} f(\mathbf{x}^t; \Xi)]\end{aligned}$$

### issues:

- distribution of  $\Xi$  may be unknown
- even if it is known, evaluating high-dimensional expectation is often expensive

# **Stochastic gradient descent (stochastic approximation)**

# Stochastic Approximation / Stochastic Gradient Descent

---

— Robbins, Monro '51

## stochastic approximation / stochastic gradient descent (SGD)

$$\mathbf{x}^{t+1} = \mathbf{x}^t - \eta_t \mathbf{g}(\mathbf{x}^t; \Xi^t) \quad (1)$$

where  $\mathbf{g}(\mathbf{x}^t; \Xi^t)$  is an *unbiased* estimate of  $\nabla F(\mathbf{x}^t)$ :

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

- a stochastic algorithm for finding a critical point  $\mathbf{x}$  obeying  $\nabla F(\mathbf{x}) = \mathbf{0}$
- more generally, a stochastic algorithm for finding the roots of  $G(\mathbf{x}) := \mathbb{E}[\mathbf{g}(\mathbf{x}; \Xi)]$

## Example: SGD for Empirical Risk Minimization

---

$$\text{minimize}_{\mathbf{x}} \quad \underbrace{F(\mathbf{x}) := \frac{1}{n} \sum_{i=1}^n f(\mathbf{x}; \{\mathbf{a}_i, y_i\})}_{\text{empirical risk}}$$

**for**  $t = 0, 1, \dots$

choose  $i_t$  uniformly at random, and run

$$\mathbf{x}^{t+1} = \mathbf{x}^t - \eta_t \nabla_{\mathbf{x}} f(\mathbf{x}^t; \{\mathbf{a}_{i_t}, y_{i_t}\}).$$

## Example: SGD for Empirical Risk Minimization

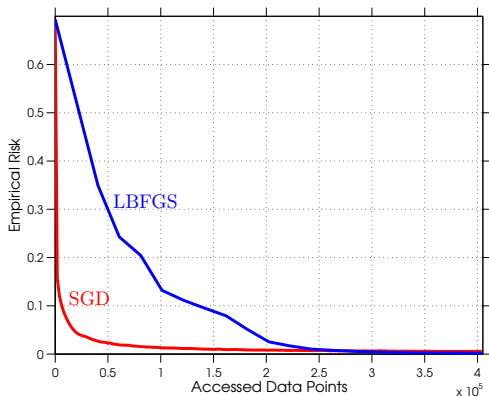
---

**benefits:** SGD exploits information more efficiently than batch methods

- practical data usually involve lots of redundancy; using all data simultaneously in each iteration might be inefficient
- SGD is particularly efficient at the very beginning, as it achieves fast initial improvement with very low per-iteration cost

## Example: SGD for Empirical Risk Minimization

— Bottou, Curtis, Nocedal '18



binary classification with logistic loss and RCV1 dataset ( $\eta_t \equiv 4$ )

# Convergence analysis

## Convergence Roadmap

---

- We first study strongly convex and smooth objectives under a second-moment bound on stochastic gradients.
- With fixed stepsizes, SGD converges linearly to a noise-dominated neighborhood.
- With diminishing stepsizes, SGD converges to  $x^*$  at rate  $O(1/t)$ .

## Strongly Convex and Smooth Problems

---

$$\text{minimize}_{\mathbf{x}} \quad F(\mathbf{x}) := \mathbb{E}[f(\mathbf{x}; \Xi)]$$

- $F$ :  $\mu$ -strongly convex,  $L$ -smooth
- $\mathbf{g}(\mathbf{x}^t; \Xi^t)$  is conditionally unbiased:

$$\mathbb{E}[\mathbf{g}(\mathbf{x}^t; \Xi^t) \mid \Xi^0, \dots, \Xi^{t-1}] = \nabla F(\mathbf{x}^t)$$

- for all  $\mathbf{x}$ ,

$$\mathbb{E}[\|\mathbf{g}(\mathbf{x}; \Xi)\|_2^2] \leq \sigma_g^2 + c_g \|\nabla F(\mathbf{x})\|_2^2 \quad (2)$$

## Convergence: Fixed Stepsizes

---

### Theorem 1 (Convergence of SGD for strongly convex problems; fixed stepsizes)

Under the assumptions in Frame 13, if  $\eta_t \equiv \eta \leq \frac{1}{Lc_g}$ , then SGD (1) achieves

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

- check Bottou, Curtis, Nocedal '18 (Theorem 4.6) for the proof

## Implications: SGD with Fixed Stepsizes

---

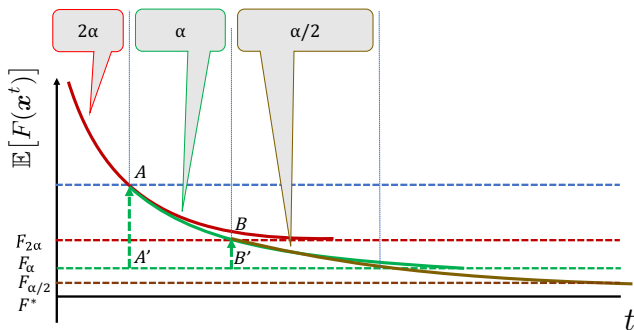
$$\mathbb{E}[F(\mathbf{x}^t) - F(\mathbf{x}^*)] \leq \frac{\eta L \sigma_g^2}{2\mu} + (1 - \eta\mu)^t (F(\mathbf{x}^0) - F(\mathbf{x}^*))$$

- fast (linear) convergence at the very beginning
- converges to some neighborhood of  $\mathbf{x}^*$  — variation in gradient computation prevents further progress
- when gradient computation is noiseless (i.e.  $\sigma_g = 0$ ), it converges linearly to optimal points
- smaller fixed stepsizes  $\eta$  yield better terminal accuracy

## One Practical Strategy

Run SGD with a fixed stepsize; whenever progress stalls, reduce the stepsize and continue.

— Bottou, Curtis, Nocedal '18



Whenever progress stalls, halve the stepsize and repeat.

## Convergence with Diminishing Stepsizes

---

### Theorem 2 (Convergence of SGD for strongly convex problems; diminishing stepsizes)

Suppose  $F$  is  $\mu$ -strongly convex, and (2) holds with  $c_g = 0$ . If  $\eta_t = \frac{\theta}{t+1}$  for some  $\theta > \frac{1}{2\mu}$ , then SGD (1) achieves

$$\mathbb{E}[\|\mathbf{x}^t - \mathbf{x}^*\|_2^2] \leq \frac{c_\theta}{t+1}$$

where  $c_\theta = \max \left\{ \frac{2\theta^2 \sigma_g^2}{2\mu\theta - 1}, \|\mathbf{x}^0 - \mathbf{x}^*\|_2^2 \right\}$

- convergence rate  $O(1/t)$  with diminishing stepsize  $\eta_t \asymp 1/t$