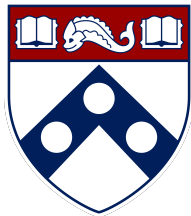


Stochastic gradient methods



Yuxin Chen

Wharton Statistics & Data Science, Fall 2023

Outline

- Stochastic gradient descent (stochastic approximation)
- Convergence analysis
- Reducing variance via iterate averaging

Stochastic programming

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

- $\boldsymbol{\xi}$: randomness in problem
- suppose $f(\cdot, \boldsymbol{\xi})$ is convex for every $\boldsymbol{\xi}$ (and hence $F(\cdot)$ 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” technical conditions

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

issues:

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

Stochastic gradient descent (stochastic approximation)

Stochastic gradient descent (SGD)

— Robbins, Monro '51

stochastic approximation / stochastic gradient descent (SGD)

$$\mathbf{x}^{t+1} = \mathbf{x}^t - \eta_t \mathbf{g}(\mathbf{x}^t; \boldsymbol{\xi}^t) \quad (11.1)$$

where $\mathbf{g}(\mathbf{x}^t; \boldsymbol{\xi}^t)$ is *unbiased* estimate of $\nabla F(\mathbf{x}^t)$, i.e.

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

Stochastic gradient descent (SGD)

— Robbins, Monro '51

stochastic approximation / stochastic gradient descent (SGD)

$$\mathbf{x}^{t+1} = \mathbf{x}^t - \eta_t \mathbf{g}(\mathbf{x}^t; \boldsymbol{\xi}^t) \quad (11.1)$$

- 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 $G(\mathbf{x}) := \mathbb{E}[\mathbf{g}(\mathbf{x}; \boldsymbol{\xi})]$

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

$$\mathbf{x}^{t+1} = \mathbf{x}^t - \eta_t \nabla_{\mathbf{x}} f_{i_t}(\mathbf{x}^t; \{\mathbf{a}_i, y_i\})$$

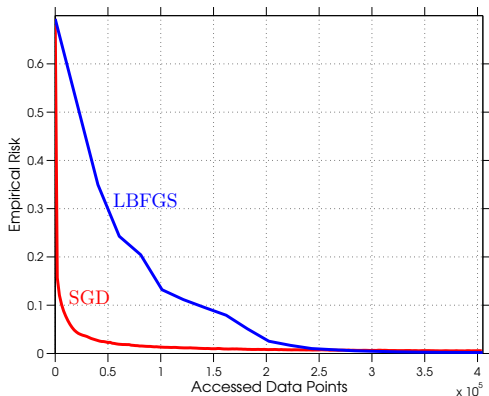
Example: SGD for empirical risk minimization

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 ($\eta_t \equiv 4$)

Example: temporal difference (TD) learning

Reinforcement learning studies a Markov decision process (MDP) with unknown model

core problem: estimate the so-called “value function” under a stationary policy π

$$V^\pi(s) = \mathbb{E}\left[r_0 + \gamma V^\pi(s_1) \mid s_0 = s\right] \quad (11.2)$$

for all $s \in \mathcal{S}$, without knowing the transition probabilities of the MDP

Example: temporal difference (TD) learning

We won't explain what Equation (11.2) means, but remark that ...

- $V^\pi(\cdot)$: value function under policy π
- s_t : state at time t
- \mathcal{S} : state space
- $0 < \gamma < 1$: discount factor
- r_t : reward at time t

Example: temporal difference (TD) learning

The definition of the value function is equivalent to

$$\mathbb{E} \left[\underbrace{V^\pi(s) - r_0 - \gamma V^\pi(s_1)}_{g(V^\pi)} \mid s_0 = s \right] = 0$$

TD(0) algorithm: for $t = 0, 1, \dots$

draw a new state s_{t+1} , collect a reward r_t , then update

$$\begin{aligned} \hat{V}^\pi(s_t) &\leftarrow \hat{V}^\pi(s_t) - \eta_t g(\hat{V}^\pi) && \text{or} \\ \hat{V}^\pi(s_t) &\leftarrow \hat{V}^\pi(s_t) - \eta_t \left\{ \hat{V}^\pi(s_t) - r_t - \hat{V}^\pi(s_{t+1}) \right\} \end{aligned}$$

Example: Q-learning

What if we also want to find an optimal policy?

core problem: solve the so-called “Bellman equation”

$$V(s) = \max_{a \in \mathcal{A}} \left\{ R(s, a) + \gamma \mathbb{E}[V(s_1) \mid s_0 = s, a_0 = a] \right\} \quad (11.3)$$

for all $s \in \mathcal{S}$, without knowing the transition probabilities of the MDP

Example: Q-learning

Again we won't explain what the Bellman equation means, but remark that ...

- $V(\cdot)$: value function
- s_t : state at time t
- \mathcal{S} : state space
- a_t : action at time t
- \mathcal{A} : action space
- $0 < \gamma < 1$: discount factor
- $R(\cdot, \cdot)$: reward function

Example: Q-learning

$$V(s) = \max_{a \in \mathcal{A}} \left\{ R(s, a) + \gamma \mathbb{E}[V(s_1) \mid s_0 = s, a_0 = a] \right\}$$

- since the transition probabilities are unknown, it is natural to resort to stochastic approximation methods
- **issue:** the Bellman equation has \mathbb{E} inside the max operator
- **a very cute idea:** introduce the so-called “Q function”

Example: Q-learning

$$V(s) = \max_{a \in \mathcal{A}} \left\{ R(s, a) + \gamma \mathbb{E}[V(s_1) \mid s_0 = s, a_0 = a] \right\}$$

Define the Q function as

$$\begin{aligned} Q(s, a) &:= R(s, a) + \gamma \mathbb{E}[V(s_1) \mid s_0 = s, a_0 = a] \\ &= R(s, a) + \gamma \mathbb{E} \left[\underbrace{\max_{\tilde{a} \in \mathcal{A}} Q(s_1, \tilde{a})}_{=V(s_1)} \mid s_0 = s, a_0 = a \right] \end{aligned} \quad (11.4)$$

- **Q learning:** use stochastic approximation methods to estimate the Q function (rather than the value function $V(\cdot)$)

Example: Q-learning

Definition of the Q-function is equivalent to

$$\mathbb{E} \left[\underbrace{Q(s, a) - R(s, a) - \gamma \max_{\tilde{a} \in \mathcal{A}} Q(s_1, \tilde{a})}_{g(Q)} \mid s_0 = s, a_0 = a \right] = 0$$

Q-learning algorithm: for $t = 0, 1, \dots$

draw a new state s_{t+1} using an action a_t , then update

$$\hat{Q}(s_t, a_t) \leftarrow \hat{Q}(s_t, a_t) - \eta_t g(\hat{Q}) \quad \text{or}$$

$$\hat{Q}(s_t, a_t) \leftarrow \hat{Q}(s_t, a_t) - \eta_t \left\{ \hat{Q}(s_t, a_t) - R(s_t, a_t) - \gamma \max_{\tilde{a} \in \mathcal{A}} \hat{Q}(s_{t+1}, \tilde{a}) \right\}$$

Convergence analysis

Strongly convex and smooth problems

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

- F : μ -strongly convex, L -smooth
- $g(\mathbf{x}^t; \boldsymbol{\xi}^t)$: an **unbiased** estimate of $\nabla F(\mathbf{x}^t)$ given $\{\boldsymbol{\xi}^0, \dots, \boldsymbol{\xi}^{t-1}\}$
- for all \mathbf{x} ,

$$\mathbb{E}[\|g(\mathbf{x}; \boldsymbol{\xi})\|_2^2] \leq \sigma_g^2 + c_g \|\nabla F(\mathbf{x})\|_2^2 \quad (11.5)$$

Convergence: fixed stepsizes

Theorem 11.1 (Convergence of SGD for strongly convex problems; fixed stepsizes)

Under the assumptions in Page 11-20, if $\eta_t \equiv \eta \leq \frac{1}{Lc_g}$, then SGD (11.1) 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

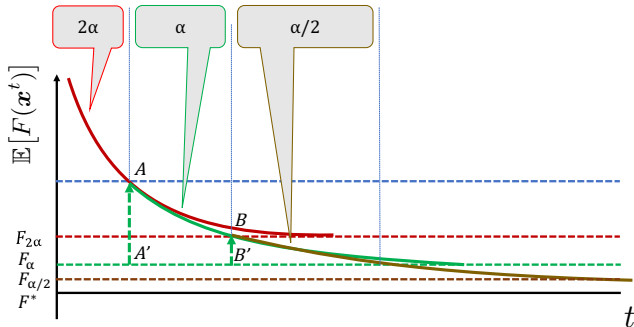
$$\mathbb{E}[F(\mathbf{x}^t) - F(\mathbf{x}^*)] \leq \frac{\eta L \sigma_g^2}{2\mu} + (1 - \eta\mu)^t (F(\mathbf{x}^0) - F(\mathbf{x}^*))$$

- 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 stepsizes η yield better converging points

One practical strategy

Run SGD with fixed stepsizes; whenever progress stalls, reduce stepsizes and continue SGD

— Bottou, Curtis, Nocedal '18



whenever progress stalls, we half the stepsizes and repeat

Convergence with diminishing stepsizes

Theorem 11.2 (Convergence of SGD for strongly convex problems; diminishing stepsizes)

Suppose F is μ -strongly convex, and (11.5) holds with $c_g = 0$. If $\eta_t = \frac{\theta}{t+1}$ for some $\theta > \frac{1}{2\mu}$, then SGD (11.1) 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 $\eta_t \asymp 1/t$

Proof of Theorem 11.2

Using the SGD update rule, we have

$$\begin{aligned}\|\mathbf{x}^{t+1} - \mathbf{x}^*\|_2^2 &= \|\mathbf{x}^t - \eta_t \mathbf{g}(\mathbf{x}^t; \boldsymbol{\xi}^t) - \mathbf{x}^*\|_2^2 \\ &= \|\mathbf{x}^t - \mathbf{x}^*\|_2^2 - 2\eta_t (\mathbf{x}^t - \mathbf{x}^*)^\top \mathbf{g}(\mathbf{x}^t; \boldsymbol{\xi}^t) + \eta_t^2 \|\mathbf{g}(\mathbf{x}^t; \boldsymbol{\xi}^t)\|_2^2\end{aligned}\quad (11.6)$$

Since \mathbf{x}^t is indep. of $\boldsymbol{\xi}^t$, apply the law of total expectation to obtain

$$\begin{aligned}\mathbb{E}[(\mathbf{x}^t - \mathbf{x}^*)^\top \mathbf{g}(\mathbf{x}^t; \boldsymbol{\xi}^t)] &= \mathbb{E}[\mathbb{E}[(\mathbf{x}^t - \mathbf{x}^*)^\top \mathbf{g}(\mathbf{x}^t; \boldsymbol{\xi}^t) \mid \boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_{t-1}]] \\ &= \mathbb{E}[(\mathbf{x}^t - \mathbf{x}^*)^\top \mathbb{E}[\mathbf{g}(\mathbf{x}^t; \boldsymbol{\xi}^t) \mid \boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_{t-1}]] \\ &= \mathbb{E}[(\mathbf{x}^t - \mathbf{x}^*)^\top \nabla F(\mathbf{x}^t)]\end{aligned}\quad (11.7)$$

Proof of Theorem 11.2 (cont.)

Furthermore, strong convexity gives

$$\langle \nabla F(\mathbf{x}^t), \mathbf{x}^t - \mathbf{x}^* \rangle = \langle \nabla F(\mathbf{x}^t) - \underbrace{\nabla F(\mathbf{x}^*)}_{=0}, \mathbf{x}^t - \mathbf{x}^* \rangle \geq \mu \|\mathbf{x}^t - \mathbf{x}^*\|_2^2$$

$$\implies \mathbb{E}[\langle \nabla F(\mathbf{x}^t), \mathbf{x}^t - \mathbf{x}^* \rangle] \geq \mu \mathbb{E}[\|\mathbf{x}^t - \mathbf{x}^*\|_2^2] \quad (11.8)$$

Combine (11.6), (11.7), (11.8) and (11.5) (with $c_g = 0$) to obtain

$$\mathbb{E}[\|\mathbf{x}^{t+1} - \mathbf{x}^*\|_2^2] \leq (1 - 2\mu\eta_t)\mathbb{E}[\|\mathbf{x}^t - \mathbf{x}^*\|_2^2] + \underbrace{\eta_t^2 \sigma_g^2}_{\text{does not vanish unless } \eta_t \rightarrow 0} \quad (11.9)$$

Take $\eta_t = \frac{\theta}{t+1}$ and use induction to conclude the proof (exercise!)

Optimality

— Nemirovski, Yudin '83, Agarwal et al. '11, Raginsky, Rakhlin '11

Informally, when minimizing strongly convex functions, no algorithm performing t queries to noisy first-order oracles can achieve an accuracy better than the order of $1/t$

\implies SGD with stepsizes $\eta_t \asymp 1/t$ is optimal

Optimality

— Nemirovski, Yudin '83

More precisely, consider a class of problems in which f is μ -strongly convex and L -smooth, and $\text{Var}(\|\mathbf{g}(\mathbf{x}^t; \boldsymbol{\xi}^t)\|_2) \leq \sigma^2$. Then the worst-case iteration complexity for (stochastic) first-order methods:

$$\sqrt{\frac{L}{\mu}} \log \left(\frac{L \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2}{\varepsilon} \right) + \frac{\sigma^2}{\mu \varepsilon}$$

- for deterministic cases: $\sigma = 0$, and hence the lower bound is

$$\sqrt{\frac{L}{\mu}} \log \left(\frac{L \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2}{\varepsilon} \right) \quad (\text{achievable by Nesterov's method})$$

Optimality

— Nemirovski, Yudin '83

More precisely, consider a class of problems in which f is μ -strongly convex and L -smooth, and $\text{Var}(\|\mathbf{g}(\mathbf{x}^t; \boldsymbol{\xi}^t)\|_2) \leq \sigma^2$. Then the worst-case iteration complexity for (stochastic) first-order methods:

$$\sqrt{\frac{L}{\mu}} \log \left(\frac{L \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2}{\varepsilon} \right) + \frac{\sigma^2}{\mu \varepsilon}$$

- for noisy cases with **large** σ , the lower bound is dominated by

$$\frac{\sigma^2}{\mu} \cdot \frac{1}{\varepsilon}$$

Comparisons with batch GD

Empirical risk minimization with n samples:

	iteration complexity	per-iteration cost	total comput. cost
batch GD	$\log \frac{1}{\varepsilon}$	n	$n \log \frac{1}{\varepsilon}$
SGD	$\frac{1}{\varepsilon}$	1	$\frac{1}{\varepsilon}$

SGD is more appealing for large n and moderate accuracy ε (in which case $\frac{1}{\varepsilon} < n \log \frac{1}{\varepsilon}$)

— which often arises in the *big data* regime!

Convex problems

What if we lose strong convexity?

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

- F : convex
- $\mathbb{E}[\|g(\mathbf{x}; \boldsymbol{\xi})\|_2^2] \leq \sigma_g^2$ for all \mathbf{x}
- $g(\mathbf{x}^t; \boldsymbol{\xi}^t)$ is an unbiased estimate of $\nabla F(\mathbf{x}^t)$ given $\{\boldsymbol{\xi}^0, \dots, \boldsymbol{\xi}^{t-1}\}$

Convex problems

Suppose we return a **weighted average**

$$\tilde{\mathbf{x}}^t := \sum_{k=0}^t \frac{\eta_k}{\sum_{j=0}^t \eta_j} \mathbf{x}^k$$

Theorem 11.3

Under the assumptions in Page 11-30, one has

$$\mathbb{E}[F(\tilde{\mathbf{x}}^t) - F(\mathbf{x}^*)] \leq \frac{\frac{1}{2}\mathbb{E}[\|\mathbf{x}^0 - \mathbf{x}^*\|_2^2] + \frac{1}{2}\sigma_g^2 \sum_{k=0}^t \eta_k^2}{\sum_{k=0}^t \eta_k}$$

- if $\eta_t \asymp 1/\sqrt{t}$, then

$$\mathbb{E}[F(\tilde{\mathbf{x}}^t) - F(\mathbf{x}^*)] \lesssim \frac{\log t}{\sqrt{t}}$$

Proof of Theorem 11.3

Remark: very similar to the convergence analysis for subgradient methods

By convexity of F , we have $F(\mathbf{x}) \geq F(\mathbf{x}^t) + (\mathbf{x} - \mathbf{x}^t)^\top \nabla F(\mathbf{x}^t)$

$$\implies \mathbb{E}[(\mathbf{x}^t - \mathbf{x}^*)^\top \nabla F(\mathbf{x}^t)] \geq \mathbb{E}[F(\mathbf{x}^t) - F(\mathbf{x}^*)]$$

This together with (11.6) and (11.7) implies

$$2\eta_k \mathbb{E}[F(\mathbf{x}^k) - F(\mathbf{x}^*)] \leq \mathbb{E}[\|\mathbf{x}^k - \mathbf{x}^*\|_2^2] - \mathbb{E}[\|\mathbf{x}^{k+1} - \mathbf{x}^*\|_2^2] + \eta_k^2 \sigma_g^2$$

Sum over $k = 0, \dots, t$ to obtain

$$\begin{aligned} \sum_{k=0}^t 2\eta_k \mathbb{E}[F(\mathbf{x}^k) - F(\mathbf{x}^*)] &\leq \mathbb{E}[\|\mathbf{x}^0 - \mathbf{x}^*\|_2^2] - \mathbb{E}[\|\mathbf{x}^{t+1} - \mathbf{x}^*\|_2^2] + \sigma_g^2 \sum_{k=0}^t \eta_k^2 \\ &\leq \mathbb{E}[\|\mathbf{x}^0 - \mathbf{x}^*\|_2^2] + \sigma_g^2 \sum_{k=0}^t \eta_k^2 \end{aligned}$$

Proof of Theorem 11.3 (cont.)

Setting $v_t = \frac{\eta_t}{\sum_{k=0}^t \eta_k}$ yields

$$\sum_{k=0}^t v_k \mathbb{E}[F(\mathbf{x}^k) - F(\mathbf{x}^*)] \leq \frac{\frac{1}{2} \mathbb{E}[\|\mathbf{x}^0 - \mathbf{x}^*\|_2^2] + \frac{1}{2} \sigma_g^2 \sum_{k=0}^t \eta_k^2}{\sum_{k=0}^t \eta_k}$$

By convexity of F , we arrive at

$$\mathbb{E}[F(\tilde{\mathbf{x}}^t) - F(\mathbf{x}^*)] \leq \frac{\frac{1}{2} \mathbb{E}[\|\mathbf{x}^0 - \mathbf{x}^*\|_2^2] + \frac{1}{2} \sigma_g^2 \sum_{k=0}^t \eta_k^2}{\sum_{k=0}^t \eta_k}$$

Reducing variance via iterate averaging

Stepsize choice $O(1/t)$?

Two conflicting regimes

- the noiseless case (i.e. $\mathbf{g}(\mathbf{x}; \boldsymbol{\xi}) = \nabla F(\mathbf{x})$): stepsizes $\eta_t \asymp 1/t$ are way too conservative
- the general noisy case: longer stepsizes ($\eta_t \gg 1/t$) might fail to suppress noise (and hence slow down convergence)

Can we modify SGD so as to allow for larger stepsizes without compromising convergence rates?

Motivation for iterate averaging

SGD with long stepsizes poorly suppresses noise, which tends to oscillate around the global minimizers due to the noisy nature of gradient computation

One may, however, average iterates to mitigate oscillation and reduce variance

Acceleration by averaging

— Ruppert '88, Polyak '90, Polyak, Juditsky '92

$$\text{return } \bar{\mathbf{x}}^t := \frac{1}{t} \sum_{i=0}^{t-1} \mathbf{x}^i \quad (11.10)$$

with larger stepsizes $\eta_t \asymp t^{-\alpha}$, $\alpha < 1$

Key idea: average the iterates to reduce variance and improve convergence

Example: a toy quadratic problem

$$\text{minimize}_{\mathbf{x} \in \mathbb{R}^d} \quad \frac{1}{2} \|\mathbf{x}\|_2^2$$

- constant stepsizes: $\eta_t \equiv \eta < 1$
- $\mathbf{g}(\mathbf{x}^t; \boldsymbol{\xi}^t) = \mathbf{x}^t + \boldsymbol{\xi}^t$ with
 - $\mathbb{E}[\boldsymbol{\xi}^t \mid \boldsymbol{\xi}^0, \dots, \boldsymbol{\xi}^{t-1}] = \mathbf{0}$
 - $\mathbb{E}[\boldsymbol{\xi}^t \boldsymbol{\xi}^{t\top} \mid \boldsymbol{\xi}^0, \dots, \boldsymbol{\xi}^{t-1}] = \mathbf{I}$

Example: a toy quadratic problem

$$\text{minimize}_{\mathbf{x} \in \mathbb{R}^d} \quad \frac{1}{2} \|\mathbf{x}\|_2^2$$

SGD iterates:

$$\mathbf{x}^1 = \mathbf{x}^0 - \eta(\mathbf{x}^0 + \boldsymbol{\xi}^0) = (1 - \eta)\mathbf{x}^0 - \eta\boldsymbol{\xi}^0$$

$$\mathbf{x}^2 = \mathbf{x}^1 - \eta(\mathbf{x}^1 + \boldsymbol{\xi}^1) = (1 - \eta)^2\mathbf{x}^0 - \eta(1 - \eta)\boldsymbol{\xi}^0 - \eta\boldsymbol{\xi}^1$$

\vdots

$$\mathbf{x}^t = (1 - \eta)^t\mathbf{x}^0 - \eta(1 - \eta)^{t-1}\boldsymbol{\xi}^0 - \eta(1 - \eta)^{t-2}\boldsymbol{\xi}^1 - \dots$$

Example: a toy quadratic problem

$$\text{minimize}_{\mathbf{x} \in \mathbb{R}^d} \quad \frac{1}{2} \|\mathbf{x}\|_2^2$$

$$\begin{aligned} \bar{\mathbf{x}}^t &\approx \underbrace{\frac{1}{t} \sum_{k=0}^{t-1} (1-\eta)^k \mathbf{x}^0}_{= \frac{1}{t} \frac{1-(1-\eta)^t}{\eta} \mathbf{x}^0 \xrightarrow{t \rightarrow \infty} \mathbf{0}} - \underbrace{\eta \{1 + (1-\eta) + \dots\} \frac{1}{t} \sum_{k=0}^{t-1} \boldsymbol{\xi}^k}_{\text{imprecise; but close enough for large } t} \\ &\approx -\frac{1}{t} \sum_{k=0}^{t-1} \boldsymbol{\xi}^k \quad (\text{since } 1 + (1-\eta) + \dots = \eta^{-1}) \\ &\xrightarrow{t \rightarrow \infty} \frac{1}{\sqrt{t}} \mathcal{N}(\mathbf{0}, \mathbf{I}) \quad (\text{the central limit theorem for martingale}) \end{aligned}$$

Example: more general quadratic problems

$$\text{minimize}_{\mathbf{x} \in \mathbb{R}^d} \quad \frac{1}{2} \mathbf{x}^\top \mathbf{A} \mathbf{x} - \mathbf{b}^\top \mathbf{x}$$

- $\mathbf{A} \succeq \mu \mathbf{I} \succ \mathbf{0}$ (strongly convex)
- constant stepsizes: $\eta_t \equiv \eta < 1/\mu$
- $g(\mathbf{x}^t; \boldsymbol{\xi}^t) = \mathbf{A} \mathbf{x}^t - \mathbf{b} + \boldsymbol{\xi}^t$ with
 - $\mathbb{E}[\boldsymbol{\xi}^t \mid \boldsymbol{\xi}^0, \dots, \boldsymbol{\xi}^{t-1}] = \mathbf{0}$
 - $\mathbf{S} := \lim_{t \rightarrow \infty} \mathbb{E}[\boldsymbol{\xi}^t \boldsymbol{\xi}^{t\top} \mid \boldsymbol{\xi}^0, \dots, \boldsymbol{\xi}^{t-1}]$ is finite

Example: more general quadratic problems

$$\text{minimize}_{\mathbf{x} \in \mathbb{R}^d} \quad \frac{1}{2} \mathbf{x}^\top \mathbf{A} \mathbf{x} - \mathbf{b}^\top \mathbf{x}$$

Theorem 11.4

Fix d . Then as $t \rightarrow \infty$, the iterate average $\bar{\mathbf{x}}^t$ obeys

$$\sqrt{t}(\bar{\mathbf{x}}^t - \mathbf{x}^*) \underbrace{\xrightarrow{\mathcal{D}}}_{\text{convergence in distribution}} \mathcal{N}(\mathbf{0}, \mathbf{A}^{-1} \mathbf{S} \mathbf{A}^{-1})$$

Example: quadratic problems

$$\sqrt{t}(\bar{\mathbf{x}}^t - \mathbf{x}^*) \xrightarrow{\mathcal{D}} \mathcal{N}(\mathbf{0}, \mathbf{A}^{-1} \mathbf{S} \mathbf{A}^{-1}), \quad t \rightarrow \infty$$

- asymptotically, $\|\bar{\mathbf{x}}^t - \mathbf{x}^*\|_2^2 \asymp 1/t$, which matches the convergence rate in Theorem 11.2
- much longer stepsizes ($\eta_t \asymp 1$)
 \implies faster convergence for less noisy cases (e.g. $\xi^t = \mathbf{0}$)

Proof sketch of Theorem 11.5

(1) Let $\Delta^t = \mathbf{x}^t - \mathbf{x}^*$ and $\bar{\Delta}^t = \bar{\mathbf{x}}^t - \mathbf{x}^*$. SGD update rule gives

$$\begin{aligned}\Delta^{t+1} &= \Delta^t - \eta(\mathbf{A}\Delta^t + \boldsymbol{\xi}^t) = (\mathbf{I} - \eta\mathbf{A})\Delta^t - \eta\boldsymbol{\xi}^t \\ \implies \Delta^{t+1} &= (\mathbf{I} - \eta\mathbf{A})^{t+1}\Delta^0 - \eta \sum_{k=0}^t (\mathbf{I} - \eta\mathbf{A})^{t-k}\boldsymbol{\xi}^k\end{aligned}$$

(2) Simple calculation gives (check Polyak, Juditsky '92)

$$\bar{\Delta}^t = \frac{1}{t\eta}\mathbf{G}_0^t\Delta^0 + \frac{1}{t}\sum_{j=0}^{t-2}\mathbf{A}^{-1}\boldsymbol{\xi}^j + \frac{1}{t}\sum_{j=0}^{t-2}(\mathbf{G}_j^t - \mathbf{A}^{-1})\boldsymbol{\xi}^j$$

$$\text{where } \mathbf{G}_j^t := \eta \sum_{i=0}^{t-1-j} (\mathbf{I} - \eta\mathbf{A})^i$$

Proof sketch of Theorem 11.5 (cont.)

(3) From the central limit theorem for martingales,

$$\frac{1}{\sqrt{t}} \sum_{j=0}^{t-2} \mathbf{A}^{-1} \boldsymbol{\xi}^j \xrightarrow{\mathcal{D}} \mathcal{N}(\mathbf{0}, \mathbf{A}^{-1} \mathbf{S} \mathbf{A}^{-1})$$

(4) With proper stepsizes, one has (check Polyak, Juditsky '92)

$$\|\mathbf{G}_0^t\| < \infty, \|\mathbf{G}_j^t - \mathbf{A}^{-1}\| < \infty \text{ and } \lim_{t \rightarrow \infty} \frac{1}{t} \sum_{j=0}^{t-1} \|\mathbf{G}_j^t - \mathbf{A}^{-1}\|^2 = 0$$

(5) Combining these bounds establishes Theorem 11.5

More general strongly convex problems

$$\text{minimize}_{\mathbf{x} \in \mathbb{R}^d} F(\mathbf{x})$$

- F : strongly convex
- stepsizes: $\eta_t \asymp t^{-\alpha}$ with $\alpha \in (0.5, 1)$
- $\mathbf{g}(\mathbf{x}^t; \boldsymbol{\xi}^t) = \nabla F(\mathbf{x}^t) + \boldsymbol{\xi}^t$
 - $\mathbb{E}[\boldsymbol{\xi}^t \mid \boldsymbol{\xi}^0, \dots, \boldsymbol{\xi}^{t-1}] = \mathbf{0}$
 - $\mathbf{S} := \lim_{t \rightarrow \infty} \mathbb{E}[\boldsymbol{\xi}^t \boldsymbol{\xi}^{t\top} \mid \boldsymbol{\xi}^0, \dots, \boldsymbol{\xi}^{t-1}]$ is finite

More general strongly convex problems

$$\text{minimize}_{\mathbf{x} \in \mathbb{R}^d} F(\mathbf{x})$$

Theorem 11.5 (informal, Polyak, Juditsky '92)

Fix d and let $t \rightarrow \infty$. For a large class of strongly convex problems, if $\eta_t \asymp t^{-\alpha}$ for some $1/2 < \alpha < 1$, then

$$\sqrt{t}(\bar{\mathbf{x}}^t - \mathbf{x}^*) \xrightarrow{\mathcal{D}} \mathcal{N}\left(\mathbf{0}, (\nabla^2 F(\mathbf{x}^*))^{-1} \mathbf{S}(\nabla^2 F(\mathbf{x}^*))^{-1}\right)$$

- depend on the local curvature at / around minimizers
- allow the stepsize η_t to be longer than $1/t$

Reference

- "*A stochastic approximation method*, H. Robbins, S. Monro, *the annals of mathematical statistics*, 1951.
- "*Robust stochastic approximation approach to stochastic programming*," A. Nemirovski et al., *SIAM Journal on optimization*, 2009.
- "*Optimization methods for large-scale machine learning*," L. Bottou et al., arXiv, 2016.
- "*New stochastic approximation type procedures*," B. Polyak, *Automat. Remote Control*, 1990.
- "*Acceleration of stochastic approximation by averaging*," B. Polyak, A. Juditsky, *SIAM Journal on Control and Optimization*, 1992.
- "*Efficient estimations from a slowly convergent Robbins-Monro process*," D. Ruppert, 1988.

Reference

- "*Problem complexity and method efficiency in optimization*," A. Nemirovski, D. Yudin, Wiley, 1983.
- "*Information-theoretic lower bounds on the oracle complexity of stochastic convex optimization*," A. Agarwal, P. Bartlett, P. Ravikumar, M. Wainwright, *IEEE Transactions on Information Theory*, 2011.
- "*Information-based complexity, feedback and dynamics in convex programming*," M. Raginsky, A. Rakhlin, *IEEE Transactions on Information Theory*, 2011.
- "*First-order methods in optimization*," A. Beck, Vol. 25, SIAM, 2017.
- "*Acceleration of stochastic approximation by averaging*," B. Polyak, A. Juditsky, *SIAM Journal on Control and Optimization*, 1992.
- "*A convergence theorem for nonnegative almost supermartingales and some applications*," H. Robbins, D. Siegmund, *Optimizing methods in statistics*, 1971.