

Optimization for Machine Learning

Lecture 3: Stochastic Gradient Descent

Sebastian Stich

CISPA – <https://cms.cispa.saarland/optml24/>

April 30, 2024

Quiz Week 2 (1)

What does $f(n) \in \mathcal{O}(g(n))$ (for $n \rightarrow \infty$) mean?

Examples:

- ▶ $10n^2 \in \mathcal{O}(n^2)$?
- ▶ $n^2 \in \mathcal{O}(n^3)$?
- ▶ $n^3 + n^2 + n + 1 \in \mathcal{O}(n^3)$?

Formally:

What about $\epsilon \rightarrow 0$?

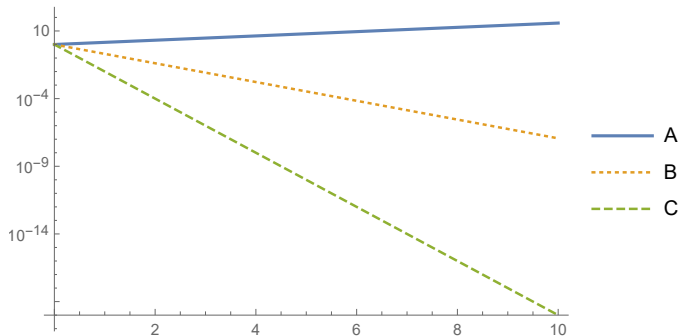
- ▶ Consider $n = \frac{1}{\epsilon}$.
- ▶ $\frac{1}{\epsilon^2} \in \mathcal{O}\left(\frac{1}{\epsilon^3}\right)$?

Quiz Week 2 (2)

Consider gradient descent on a smooth and convex function $f: \mathbb{R}^d \rightarrow \mathbb{R}$,

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \gamma \nabla f(\mathbf{x}_t),$$

for a stepsize $\gamma > 0$.



The figure shows three runs of gradient descent, with the stepsizes $\{\gamma, \gamma/2, \gamma/4\}$, for a (fixed) value of γ . Which curve does correspond to which stepsize?

Chapter 6

Stochastic Gradient Descent

Stochastic gradient descent

Many objective functions are **sum structured**:

$$f(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{x}).$$

Example: f_i is the cost function of the i -th observation, taken from a training set of n observation.

Evaluating $\nabla f(\mathbf{x})$ of a sum-structured function is expensive (sum of n gradients).

Stochastic gradient descent: the algorithm

choose $\mathbf{x}_0 \in \mathbb{R}^d$

sample $i \in [n]$ uniformly at random

$$\mathbf{x}_{t+1} := \mathbf{x}_t - \gamma_t \nabla f_i(\mathbf{x}_t).$$

for **iterations** $t = 0, 1, \dots$, and **stepsizes** $\gamma_t \geq 0$.

Only update with the gradient of f_i instead of the full gradient!

Iteration is n times cheaper than in full gradient descent.

The vector $\mathbf{g}_t := \nabla f_i(\mathbf{x}_t)$ is called a **stochastic gradient**.

\mathbf{g}_t is a vector of d random variables, but we will also simply call this a random variable.

Stochastic Optimization

The finite sum structure is not necessary. All results we discuss in this course do also hold for stochastic optimization problems:

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

- ▶ \mathcal{D} a distribution
- ▶ for every ξ , access to stochastic gradients $\nabla F(\mathbf{x}, \xi)$
- ▶ finite-sum is a special case:

- ▶ algorithm:

sample $\xi_t \sim \mathcal{D}$ uniformly at random
 $\mathbf{x}_{t+1} := \mathbf{x}_t - \gamma_t \nabla F(\mathbf{x}_t, \xi_t).$

Unbiasedness

Consider a stochastic gradient \mathbf{g}_t , for a random index $i_t \in [n]$.

$$\mathbf{g}_t = \nabla f_{i_t}(\mathbf{x}_t),$$

We **cannot** use our previous inequalities as they might not hold, depending on how the stochastic gradient \mathbf{g}_t turns out.

We will show (and exploit): many inequalities holds **in expectation**.

For this, we use that by definition, \mathbf{g}_t is an **unbiased estimate** of $\nabla f(\mathbf{x}_t)$:

$$\mathbb{E}[\mathbf{g}_t] = \frac{1}{n} \sum_{i=1}^n \nabla f_i(\mathbf{x}_t) = \nabla f(\mathbf{x}_t).$$

Convexity in expectation

Note, for any fixed vector $\mathbf{y} \in \mathbb{R}^d$:

$$\mathbb{E}[\mathbf{g}_t^\top \mathbf{y}] = \mathbb{E}[\mathbf{g}_t]^\top \mathbf{y} = \nabla f(\mathbf{x}_t)^\top \mathbf{y}.$$

Hence, for a convex function $f: \mathbb{R}^d \rightarrow \mathbb{R}$:

$$\mathbb{E}[\mathbf{g}_t^\top (\mathbf{x}_t - \mathbf{x}^*)] = \nabla f(\mathbf{x}_t)^\top (\mathbf{x}_t - \mathbf{x}^*) \geq f(\mathbf{x}_t) - f(\mathbf{x}^*).$$

Quadratic upper with stochastic updates?

Can we also use expectation with the quadratic upper bound?

Recall, a step of SGD: $\mathbf{x}_{t+1} = \mathbf{x}_t - \gamma \mathbf{g}_t$.

$$\begin{aligned} & \mathbb{E} \left[f(\mathbf{x}_t) + \nabla f(\mathbf{x}_t)^\top (\mathbf{x}_{t+1} - \mathbf{x}_t) + \frac{L}{2} \|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 \right] \\ &= \mathbb{E} \left[f(\mathbf{x}_t) + \nabla f(\mathbf{x}_t)^\top (-\gamma \mathbf{g}_t) + \frac{L}{2} \|\gamma \mathbf{g}_t\|^2 \right] \\ &= f(\mathbf{x}_t) - \gamma \nabla f(\mathbf{x}_t)^\top \nabla f(\mathbf{x}_t) + \frac{\gamma^2 L}{2} \mathbb{E} [\|\mathbf{g}_t\|^2] \end{aligned}$$

What is $\mathbb{E} [\|\mathbf{g}_t\|^2]$? **We need one more assumption!**

Case 1: **Bounded Gradients**

Bounded Gradient Assumption

Assume that there exists a constant $B \geq 0$, such that:

$$\mathbb{E} \left[\|\mathbf{g}_t\|^2 \right] \leq B^2$$

for all t .

- + This simplifies the proofs to a certain degree, while still comprehensively addressing most of the additional complexity presented by stochastic gradients..
- Might not hold. (Example: quadratic functions)

Bounded stochastic gradients: $\mathcal{O}(1/\varepsilon^2)$ steps

Theorem (Lecture-3).1

Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be convex and differentiable, \mathbf{x}^* a global minimum; furthermore, suppose that $\|\mathbf{x}_0 - \mathbf{x}^*\| \leq R$, and that $\mathbb{E}[\|\mathbf{g}_t\|^2] \leq B^2$ for all t . Choosing the constant stepsize

$$\gamma := \frac{R}{B\sqrt{T}}$$

stochastic gradient descent yields

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}[f(\mathbf{x}_t)] - f(\mathbf{x}^*) \leq \frac{RB}{\sqrt{T}}.$$

- ▶ we assume bounded stochastic gradients **in expectation**;
- ▶ error bound holds **in expectation**.

Proof I

$$\begin{aligned}\mathbb{E} \left[\|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2 \right] &= \mathbb{E} \left[\|\mathbf{x}_t - \gamma \mathbf{g}_t - \mathbf{x}^*\|^2 \right] \\ &= \mathbb{E} \left[\|\mathbf{x}_t - \mathbf{x}^*\|^2 - 2\gamma \mathbf{g}_t^\top (\mathbf{x}_t - \mathbf{x}^*) + \gamma^2 \|\mathbf{g}_t\|^2 \right] \\ &\leq \|\mathbf{x}_t - \mathbf{x}^*\|^2 - 2\gamma \nabla f(\mathbf{x}_t)^\top (\mathbf{x}_t - \mathbf{x}^*) + \gamma^2 B^2 \\ &\leq \|\mathbf{x}_t - \mathbf{x}^*\|^2 - 2\gamma (f(\mathbf{x}_t) - f(\mathbf{x}^*)) + \gamma^2 B^2\end{aligned}\tag{1}$$

Proof II

We re-arrange and prepare to apply the telescoping sum trick:

$$2(f(\mathbf{x}_t) - f(\mathbf{x}^*)) \leq \frac{\|\mathbf{x}_t - \mathbf{x}^*\|^2 - \mathbb{E} \left[\|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2 \right]}{\gamma} + \gamma B^2$$

This does not seem to work! However, we can also take expectation over \mathbf{x}_t :

$$2\mathbb{E} [f(\mathbf{x}_t) - f(\mathbf{x}^*)] \leq \frac{\mathbb{E} \|\mathbf{x}_t - \mathbf{x}^*\|^2 - \mathbb{E} \|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2}{\gamma} + \gamma B^2$$

Note: this argument can be made more rigorous. See lecture notes or other sources for details.

Proof III

By telescoping (and dividing by T):

$$\frac{2}{T} \sum_{i=0}^{T-1} \mathbb{E} f(\mathbf{x}_t) - f(\mathbf{x}^*) \leq \frac{\|\mathbf{x}_0 - \mathbf{x}^*\|^2}{\gamma T} + \gamma B^2 \leq \frac{R^2}{\gamma T} + \gamma B^2$$

We now observe that the choice $\gamma = \frac{R}{B\sqrt{T}}$ indeed implies the theorem.

Tame strong convexity: $\mathcal{O}(1/\varepsilon)$ steps

Theorem (Lecture-3).2

Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be differentiable and strongly convex with parameter $\mu > 0$; let \mathbf{x}^* be the unique global minimum of f and assume that $\mathbb{E}[\|\mathbf{g}_t\|^2] \leq B^2$ for all t . With decreasing step size

$$\gamma_t := \frac{2}{\mu(t+1)}$$

stochastic gradient descent yields

$$\mathbb{E} \left[f \left(\frac{2}{T(T+1)} \sum_{t=1}^T t \cdot \mathbf{x}_t \right) - f(\mathbf{x}^*) \right] \leq \frac{2B^2}{\mu(T+1)}.$$

- weighted averaging puts more importance on recent iterates!

Proof I

The proof is starting in the same way. Except that we can use strong convexity:

$$\nabla f(\mathbf{x}_t)^\top (\mathbf{x}_t - \mathbf{x}^*) \geq f(\mathbf{x}_t) - f(\mathbf{x}^*) + \frac{\mu}{2} \|\mathbf{x}_t - \mathbf{x}^*\|^2$$

Equation (1) will change into:

$$\mathbb{E} \|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2 \leq (1 - \mu\gamma_t/2) \|\mathbf{x}_t - \mathbf{x}^*\|^2 - 2\gamma_t (f(\mathbf{x}_t) - f(\mathbf{x}^*)) + \gamma_t^2 B^2$$

And therefore

$$\mathbb{E} f(\mathbf{x}_t) - f(\mathbf{x}^*) \leq \frac{\gamma_t B^2}{2} + \frac{1 - \mu\gamma_t/2}{2\gamma_t} \mathbb{E} \|\mathbf{x}_t - \mathbf{x}^*\|^2 - \frac{1}{2\gamma_t} \mathbb{E} \|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2$$

Proof II

Plug in $\gamma_t^{-1} = \mu(1+t)/2$ and multiply with t on both sides:

$$\begin{aligned} t \cdot \mathbb{E}(f(\mathbf{x}_t) - f(\mathbf{x}^*)) &\leq \frac{B^2 t}{\mu(t+1)} + \frac{\mu}{4} \left(t(t-1) \mathbb{E} \|\mathbf{x}_t - \mathbf{x}^*\|^2 - (t+1)t \mathbb{E} \|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2 \right) \\ &\leq \frac{B^2}{\mu} + \frac{\mu}{4} \left(t(t-1) \mathbb{E} \|\mathbf{x}_t - \mathbf{x}^*\|^2 - (t+1)t \mathbb{E} \|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2 \right). \end{aligned}$$

Now we get telescoping...

$$\sum_{t=0}^{T-1} t \cdot \mathbb{E}(f(\mathbf{x}_t) - f(\mathbf{x}^*)) \leq \frac{TB^2}{\mu} + \frac{\mu}{4} \left(0 - T(T+1) \mathbb{E} \|\mathbf{x}_{T+1} - \mathbf{x}^*\|^2 \right) \leq \frac{TB^2}{\mu}.$$

Finally, use $\frac{2}{T(T+1)} \sum_{t=1}^T t = 1$, and Jensen's inequality.

Discussion

- ▶ strong convexity helps: $\mathcal{O}(\frac{1}{\epsilon})$ convergence, vs. $\mathcal{O}(\frac{1}{\epsilon^2})$
- ▶ stochastic gradients make the convergence more difficult: $\mathcal{O}(\frac{1}{\epsilon})$ convergence vs. $\mathcal{O}(\log(\frac{1}{\epsilon}))$ in the deterministic setting for gradient descent!
(recall Exercise Sheet 2)
- ▶ Note: The $\mathcal{O}(\frac{1}{\epsilon})$ convergence is optimal!
- ▶ Weighted averaging is a common & useful trick to adapt telescoping sum proofs to the strongly-convex case!

Case 2: **Bounded Variance**

Bounded Variance Assumption

Assume that there exists a constant $\sigma \geq 0$, such that:

$$\mathbb{E} \left[\|\mathbf{g}_t - \nabla f(\mathbf{x}_t)\|^2 \right] \leq \sigma^2$$

for all t .

- + Standard and widely-accepted model in complexity theory.
- Might not hold on all (but much fewer) problems of interest.
- ▶ (Convergence proof: we will cover some examples next week—you could try yourself as an exercise!)

Mini-batch SGD

Mini-batch SGD

Instead of using a single element f_i , use an average of several of them:

$$\tilde{\mathbf{g}}_t := \frac{1}{m} \sum_{j=1}^m \mathbf{g}_t^j.$$

where \mathbf{g}_t^j denotes a stochastic gradient drawn uniformly and independently at random.
 m denotes the **batch size**.

Extreme cases:

$m = 1 \Leftrightarrow$ SGD as originally defined

$m = n \Leftrightarrow$ full gradient descent

Benefit: Gradient computation can be naively parallelized

Mini-batch SGD

Variance Intuition: Taking an average of many independent random variables reduces the variance. So for larger size of the mini-batch m , $\tilde{\mathbf{g}}_t$ will be closer to the true gradient, in expectation:

$$\begin{aligned}\mathbb{E}\left[\left\|\tilde{\mathbf{g}}_t - \nabla f(\mathbf{x}_t)\right\|^2\right] &= \mathbb{E}\left[\left\|\frac{1}{m} \sum_{j=1}^m \mathbf{g}_t^j - \nabla f(\mathbf{x}_t)\right\|^2\right] \\ &= \frac{1}{m} \mathbb{E}\left[\left\|\mathbf{g}_t^1 - \nabla f(\mathbf{x}_t)\right\|^2\right] \leq \frac{\sigma^2}{m}.\end{aligned}$$

- variance reduction by a factor of at least m

Lecture 3 Recap

- ▶ SGD: the most important building block in ML/DL optimization!
 - ▶ low per-iteration cost
 - ▶ ideal if **low-accuracy** approximations suffice (say, $\epsilon \geq 0.01$)
- ▶ SGD convergence proof under the bounded gradient assumption
 - ▶ we will discuss next week a proof with the bounded variance assumption
- ▶ variance-reduction effect of mini-batches
- ▶ weighted averaging to make telescoping work

Discussion

Discussion

Discussion