

# 18

---

## Stochastic Gradient Descent

---

In this lecture we introduce **Stochastic Gradient Descent** (SGD), a probabilistic version of gradient descent that has been around since the 1950s, and that has become popular in the context of data science and machine learning. To motivate the algorithm, consider a set of functions  $\mathcal{H} = \{h_{\mathbf{w}} : \mathbf{w} \in \mathbb{R}^d\}$ , where each such function depends on  $d$  parameters. Also consider a smooth loss functions  $L$ , or a smooth approximation of a loss function. Given samples  $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^n \subset \mathcal{X} \times \mathcal{Y}$ , define the functions

$$f_i(\mathbf{w}) = L(h_{\mathbf{w}}(\mathbf{x}_i), \mathbf{y}_i), \quad i \in \{1, \dots, n\}.$$

The problem of finding functions that minimize the empirical risk is

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimize}} \quad \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{w}).$$

The  $f_i$  are often assumed to be convex and smooth. In addition one often considers a regularization term  $R(\mathbf{w})$ . In what follows, we abstract from the machine learning context and consider purely the associated optimization problem.

### Stochastic Gradient Descent

We consider an objective function of the form

$$f(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{w}). \tag{18.1}$$

In what follows we assume that the functions  $f_i$  are convex and differentiable. If  $n$  is large, then computing the gradient can be very expensive. However, and considering the machine learning context, where  $f(\mathbf{w})$  is an estimator of the generalization risk  $\mathbb{E}_{\xi}[f_{\xi}(\mathbf{w})]$  of a family of functions  $f_{\xi}$  parametrized by a random vector  $\xi$ , we can shift the focus to finding an **unbiased estimator** of the gradient. Quite trivially, choosing

an index  $j$  uniformly at random and computing the gradient of  $f_j(\mathbf{w})$  gives such an unbiased estimator by *definition*:

$$\mathbb{E}_U[f_U(\mathbf{w})] = \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{w}),$$

where  $\mathbb{P}\{U = j\} = 1/n$  for  $j \in [n] = \{1, \dots, n\}$ . The Stochastic Gradient Descent (SGD) algorithm proceeds as follows. Begin with  $\mathbf{w}^0 \in \mathbb{R}^d$ . At each step  $k$ , compute an unbiased estimator  $\mathbf{g}^k$  of the gradient at  $\mathbf{w}^k$ :

$$\mathbb{E}[\mathbf{g}^k | \mathbf{w}^k] = \nabla f(\mathbf{w}^k).$$

Next, take a step in direction  $\mathbf{g}^k$ :

$$\mathbf{w}^{k+1} = \mathbf{w}^k - \eta_k \mathbf{g}^k,$$

where  $\eta_k$  is a step length (or **learning rate** in machine learning jargon).

While there are many variants of stochastic gradient descent, we consider the simplest version in which  $\mathbf{g}^k$  is chosen by picking one of the gradients  $\nabla f_i(\mathbf{w})$  uniformly at random, and we refer to this as SGD with *uniform sampling*. A commonly used generalization is **mini-batch** sampling, where one chooses a small set of indices  $I \subset \{1, \dots, n\}$  at random, instead of only one. We also restrict to the smooth setting without a regularization term; in the non-smooth setting one would apply a proximal operator. Since SGD involves random choices, convergence results are stated in terms of the expected value. Let  $U$  be a random variable with distribution  $\mathbb{P}\{U = i\} = 1/n$  for  $i \in [n]$ . Then

$$\mathbb{E}_U[\nabla f_U(\mathbf{w})] = \frac{1}{n} \sum_{i=1}^n \nabla f_i(\mathbf{w}) = \nabla f(\mathbf{w}),$$

so that  $\nabla f_U$  is an unbiased estimator of  $\nabla f$ . Assuming that  $f$  has a unique minimizer  $\mathbf{w}^*$ , we define the empirical **variance** at the optimal point  $\mathbf{w}^*$  as

$$\sigma^2 = \mathbb{E}_U[\|\nabla f_U(\mathbf{w}^*)\|^2] = \frac{1}{n} \sum_{i=1}^n \|\nabla f_i(\mathbf{w}^*)\|^2. \quad (18.2)$$

We can now state the convergence result for stochastic gradient descent.

**Theorem 18.1.** *Assume the function  $f$  is  $\alpha$ -strongly convex and that the  $f_i$  are convex and  $\beta$ -smooth for  $i \in [n]$  and  $4\beta > \alpha$ . Assume  $f$  has a unique minimizer  $\mathbf{w}^*$  and define the variance as in (18.2). Then for any starting point  $\mathbf{w}^0$ , the sequence of iterates  $\{\mathbf{w}^k\}$  generated by SGD with uniform sampling and step length  $\eta = 1/(2\beta)$  satisfies*

$$\mathbb{E}[\|\mathbf{w}^k - \mathbf{w}^*\|^2] \leq \left(1 - \frac{\alpha}{4\beta}\right)^k \|\mathbf{w}^0 - \mathbf{w}^*\|^2 + \frac{2\sigma^2}{\alpha\beta}.$$

*Proof.* As in the analysis of gradient descent, we get

$$\|\mathbf{w}^{k+1} - \mathbf{w}^*\|^2 = \|\mathbf{w}^k - \mathbf{w}^*\|^2 - 2\eta \nabla f_U(\mathbf{w}^k)^\top (\mathbf{w}^k - \mathbf{w}^*) + \eta^2 \|\nabla f_U(\mathbf{w}^k)\|^2.$$

Taking the expectation conditional on  $\mathbf{w}^k$ , we get

$$\begin{aligned} \mathbb{E}[\|\mathbf{w}^{k+1} - \mathbf{w}^*\|^2 | \mathbf{w}^k] &= \|\mathbf{w}^k - \mathbf{w}^*\|^2 - 2\eta \nabla f(\mathbf{w}^k)^\top (\mathbf{w}^k - \mathbf{w}^*) \\ &\quad + \eta^2 \mathbb{E}[\|\nabla f_U(\mathbf{w}^k)\|^2 | \mathbf{w}^k], \end{aligned} \quad (18.3)$$

where we used the fact that the expectation satisfies  $\mathbb{E}[\nabla f_U(\mathbf{w})] = \nabla f(\mathbf{w})$ . For the last term we use the bound

$$\begin{aligned} \mathbb{E}[\|\nabla f_U(\mathbf{w}^k)\|^2 | \mathbf{w}^k] &= \mathbb{E}[\|\nabla f_U(\mathbf{w}^k) - \nabla f_U(\mathbf{w}^*) + \nabla f_U(\mathbf{w}^*)\|^2] \\ &\leq 2\mathbb{E}[\|\nabla f_U(\mathbf{w}^k) - \nabla f_U(\mathbf{w}^*)\|^2] + 2\mathbb{E}[\|\nabla f_U(\mathbf{w}^*)\|^2] \\ &= 2\mathbb{E}[\|\nabla f_U(\mathbf{w}^k) - \nabla f_U(\mathbf{w}^*)\|^2] + 2\sigma^2. \end{aligned}$$

Using the characterization of  $\beta$  smoothness from the previous chapter, namely

$$\frac{1}{\beta} \|\nabla f_i(\mathbf{w}) - \nabla f_i(\mathbf{v})\|^2 \leq (\nabla f_i(\mathbf{w}) - \nabla f_i(\mathbf{v}))^\top (\mathbf{w} - \mathbf{v}), \quad (18.4)$$

we get that

$$\begin{aligned} \mathbb{E}[\|\nabla f_U(\mathbf{w}^k) - \nabla f_U(\mathbf{w}^*)\|^2 | \mathbf{w}^k] &= \frac{1}{n} \sum_{i=1}^n \|\nabla f_i(\mathbf{w}^k) - \nabla f_i(\mathbf{w}^*)\|^2 \\ &\stackrel{(18.4)}{\leq} \frac{\beta}{n} \sum_{i=1}^n (\nabla f_i(\mathbf{w}^k) - \nabla f_i(\mathbf{w}^*))^\top (\mathbf{w}^k - \mathbf{w}^*) \\ &= \frac{\beta}{n} \sum_{i=1}^n \nabla f_i(\mathbf{w}^k)^\top (\mathbf{w}^k - \mathbf{w}^*) \\ &\quad - \frac{\beta}{n} \sum_{i=1}^n \nabla f_i(\mathbf{w}^*)^\top (\mathbf{w}^k - \mathbf{w}^*) \\ &= \beta \nabla f(\mathbf{w})^\top (\mathbf{w}^k - \mathbf{w}^*), \end{aligned}$$

where we used that  $\nabla f(\mathbf{w}^*) = \mathbf{0}$  for the last equality. Hence, we get the bound

$$\mathbb{E}[\|\nabla f_U(\mathbf{w}^k)\|^2 | \mathbf{w}^k] \leq 2\beta \nabla f(\mathbf{w})^\top (\mathbf{w}^k - \mathbf{w}^*) + 2\sigma^2.$$

Plugging this into (18.3), we get

$$\mathbb{E}[\|\mathbf{w}^{k+1} - \mathbf{w}^*\|^2 | \mathbf{w}^k] \leq \|\mathbf{w}^k - \mathbf{w}^*\|^2 - (2\eta - 2\eta^2\beta) \nabla f(\mathbf{w}^k)^\top (\mathbf{w}^k - \mathbf{w}^*) + 2\eta^2\sigma^2.$$

Using the  $\alpha$ -strong convexity, we get the bound

$$\nabla f(\mathbf{w}^k)^\top \mathbf{w}^k - \mathbf{w}^* \geq f(\mathbf{w}^k) - f(\mathbf{w}^*) + \frac{\alpha}{2} \|\mathbf{w}^k - \mathbf{w}^*\|^2 \geq \frac{\alpha}{2} \|\mathbf{w}^k - \mathbf{w}^*\|^2,$$

since  $f(\mathbf{w}^k) \geq f(\mathbf{w}^*)$ , so that we get

$$\mathbb{E}[\|\mathbf{w}^{k+1} - \mathbf{w}^*\|^2 | \mathbf{w}^k] \leq (1 - \eta\alpha(1 - \eta\beta))\|\mathbf{w}^k - \mathbf{w}^*\|^2 + 2\eta^2\sigma^2.$$

With step length  $\eta = 1/(2\beta)$ , and taking the expected value over all previous iterates, we get

$$\mathbb{E}[\|\mathbf{w}^{k+1} - \mathbf{w}^*\|^2] \leq \left(1 - \frac{\alpha}{4\beta}\right) \mathbb{E}[\|\mathbf{w}^k - \mathbf{w}^*\|^2] + \frac{\sigma^2}{2\beta^2}.$$

Applying this bound recursively (and moving the index down), we get

$$\begin{aligned} \mathbb{E}[\|\mathbf{w}^k - \mathbf{w}^*\|^2] &\leq \left(1 - \frac{\alpha}{4\beta}\right)^k \|\mathbf{w}^0 - \mathbf{w}^*\|^2 + \frac{\sigma^2}{2\beta^2} \sum_{j=0}^{k-1} \left(1 - \frac{\alpha}{4\beta}\right)^j \\ &\leq \left(1 - \frac{\alpha}{4\beta}\right)^k \|\mathbf{w}^0 - \mathbf{w}^*\|^2 + \frac{2\sigma^2}{\alpha\beta}, \end{aligned}$$

where we used that  $4\beta > \alpha$ .  $\square$

**Example 18.2.** Consider the problem of **logistic regression**, where the aim is to minimize the objective function

$$f(\mathbf{w}) = \sum_{i=1}^n (\log(1 + \exp(\mathbf{x}_i^T \mathbf{w})) - y_i \mathbf{x}_i^T \mathbf{w})$$

over a vector of weights  $\mathbf{w}$ . This problem arises in the context of a binary classification problem with data pairs  $(\mathbf{x}_i, y_i)$  and  $y_i \in \{0, 1\}$ . Setting

$$p := \frac{e^{\mathbf{x}^T \mathbf{w}}}{1 + e^{\mathbf{x}^T \mathbf{w}}},$$

the resulting classifier is the function

$$h(\mathbf{w}) = \begin{cases} 1 & p > 1/2 \\ 0 & p \leq 1/2 \end{cases}$$

The function  $f$  is convex, and the gradient is

$$\nabla f(\mathbf{w}) = -\mathbf{X}^T(\mathbf{y} - \mathbf{p}(\mathbf{w})),$$

where  $\mathbf{X} \in \mathbb{R}^{n \times d}$  is the matrix with the  $\mathbf{x}_i^T$  as rows,  $\mathbf{y} = (y_1, \dots, y_n)^T$ , and  $\mathbf{p}(\mathbf{w}) \in \mathbb{R}^n$  has coordinates

$$p_i(\mathbf{w}) = \frac{\exp(\mathbf{x}_i^T \mathbf{w})}{1 + \exp(\mathbf{x}_i^T \mathbf{w})}, \quad 1 \leq i \leq n.$$

We can apply different versions of gradient descent to this problem. Figure 1 shows the typical paths of gradient descent and of stochastic gradient descent for a problem with 100 data points. Note that using a naive approach to computing the gradient, one would need to compute 100 gradients at each step. Stochastic gradient descent, on the other hand, fails to converge due to the variance of the gradient estimator (see Figure 2).

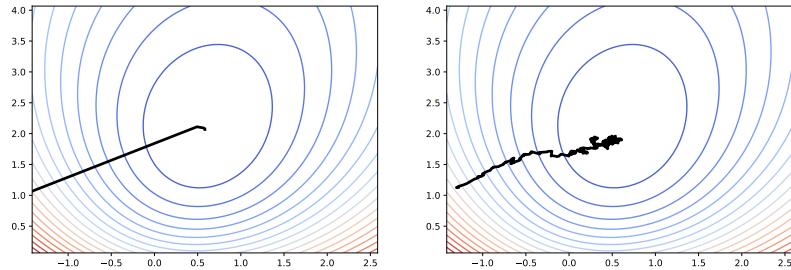


Figure 18.1: The path of gradient descent and of stochastic gradient descent with constant step length.

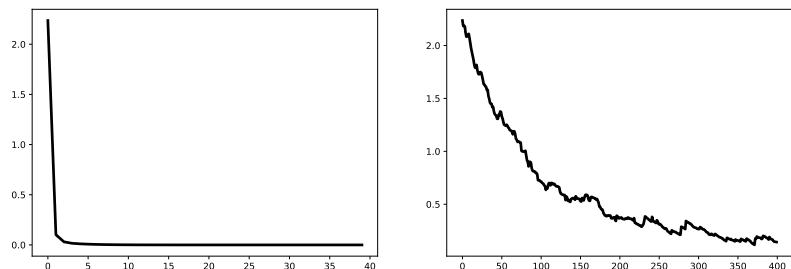


Figure 18.2: Convergence of gradient descent and SGD.

## Extensions

The version of SGD described here is the most basic one. There are many possible extensions to the methods. They include considering different sampling schemes, including **mini-batching** and **importance sampling**. These sampling strategies have the effect of reducing the variance  $\sigma^2$ . In addition, improvements can be made in the step length selection and when dealing with non-smooth functions, where the proximal operator comes into play.

## Notes

The origins of stochastic gradient descent go back to the work of Robbins and Monro in 1951 [4]. The algorithm has been rediscovered many times, and gained popularity due to its effectiveness in training deep neural networks on large data sets, where gradient computations are very expensive. Despite its simplicity, a systematic and rigorous analysis has not been available until recently. The presentation in this chapter is based loosely on the papers [3] and [1]. A more general and systematic analysis of SGD that includes non-smooth objectives is given in [2]. These works also discuss general sampling techniques, not just uniform sampling.

- [1] Léon Bottou, Frank E Curtis, and Jorge Nocedal. Optimization methods for large-scale machine learning. *Siam Review*, 60(2):223–311, 2018.
- [2] Eduard Gorbunov, Filip Hanzely, and Peter Richtárik. A unified theory of SGD: Variance reduction, sampling, quantization and coordinate descent. *arXiv preprint arXiv:1905.11261*, 2019.
- [3] Robert Mansel Gower, Nicolas Loizou, Xun Qian, Alibek Sailanbayev, Egor Shulgin, and Peter Richtárik. Sgd: General analysis and improved rates. *arXiv preprint arXiv:1901.09401*, 2019.
- [4] Herbert Robbins and Sutton Monro. A stochastic approximation method. *The Annals of Mathematical Statistics*, pages 400–407, 1951.