

ECE 595: Machine Learning I

Lecture 5.2: Stochastic Gradient Descent

Spring 2020

Stanley Chan

School of Electrical and Computer Engineering Purdue University



Outline

Mathematical Background

- Lecture 4: Intro to Optimization
- **Lecture 5: Gradient Descent**

Lecture 5: Gradient Descent

- Gradient Descent
 - Descent Direction
 - Step Size
 - Convergence
- **Stochastic Gradient Descent**
 - **Difference between GD and SGD**
 - **Why does SGD work?**

Stochastic Gradient Descent

Most loss functions in machine learning problems are **separable**:

$$J(\boldsymbol{\theta}) = \frac{1}{N} \sum_{n=1}^N \mathcal{L}(g_{\boldsymbol{\theta}}(\mathbf{x}^n), y^n) = \frac{1}{N} \sum_{n=1}^N J_n(\boldsymbol{\theta}). \quad (1)$$

For example,

- Square-loss:

$$J(\boldsymbol{\theta}) = \sum_{n=1}^N (g_{\boldsymbol{\theta}}(\mathbf{x}^n) - y^n)^2$$

- Cross-entropy loss:

$$J(\boldsymbol{\theta}) = - \sum_{n=1}^N \left\{ y^n \log g_{\boldsymbol{\theta}}(\mathbf{x}^n) + (1 - y^n) \log(1 - g_{\boldsymbol{\theta}}(\mathbf{x}^n)) \right\}$$

- Logistic loss:

$$J(\boldsymbol{\theta}) = \sum_{n=1}^N \log(1 + e^{-y^n \boldsymbol{\theta}^T \mathbf{x}^n})$$

Full Gradient VS Partial Gradient

Vanilla **gradient descent**:

$$\boldsymbol{\theta}^{t+1} = \boldsymbol{\theta}^t - \underbrace{\eta^t \nabla J(\boldsymbol{\theta}^t)}_{\text{main computation}} . \quad (2)$$

The full gradient of the loss is

$$\nabla J(\boldsymbol{\theta}) = \frac{1}{N} \sum_{n=1}^N \nabla J_n(\boldsymbol{\theta}) \quad (3)$$

Stochastic gradient descent:

$$\nabla J(\boldsymbol{\theta}) \approx \frac{1}{|\mathcal{B}|} \sum_{n \in \mathcal{B}} \nabla J_n(\boldsymbol{\theta}) \quad (4)$$

where $\mathcal{B} \subseteq \{1, \dots, N\}$ is a random subset. $|\mathcal{B}| = \text{batch size}$.

SGD Algorithm

Algorithm (Stochastic Gradient Descent)

- 1 Given $\{(\mathbf{x}^n, y^n) \mid n = 1, \dots, N\}$.
- 2 Initialize θ (zero or random)
- 3 For $t = 1, 2, 3, \dots$
 - Draw a random subset $\mathcal{B} \subseteq \{1, \dots, N\}$.
 - Update

$$\theta^{t+1} = \theta^t - \eta^t \frac{1}{|\mathcal{B}|} \sum_{n \in \mathcal{B}} \nabla J_n(\theta) \quad (5)$$

- If $|\mathcal{B}| = 1$, then use only one sample at a time.
- The approximate gradient is **unbiased**: (See Appendix for Proof)

$$\mathbb{E} \left[\frac{1}{|\mathcal{B}|} \sum_{n \in \mathcal{B}} \nabla J_n(\theta) \right] = \nabla J(\theta).$$

Interpreting SGD

- Just showed that the SGD step is unbiased:

$$\mathbb{E} \left[\frac{1}{|\mathcal{B}|} \sum_{n \in \mathcal{B}} \nabla J_n(\boldsymbol{\theta}) \right] = \nabla J(\boldsymbol{\theta}).$$

- Unbiased gradient implies that each update is

gradient + zero-mean noise

- Step size: SGD with constant step size **does not converge**.
- If $\boldsymbol{\theta}^*$ is a minimizer, then $J(\boldsymbol{\theta}^*) = \frac{1}{N} \sum_{n=1}^N J_n(\boldsymbol{\theta}^*) = 0$. But

$$\frac{1}{|\mathcal{B}|} \sum_{n \in \mathcal{B}} J_n(\boldsymbol{\theta}^*) \neq 0, \quad \text{since } \mathcal{B} \text{ is a subset.}$$

- Typical strategy: Start with large step size and gradually decrease: $\eta^t \rightarrow 0$, e.g., $\eta^t = t^{-a}$ for some constant a .

Perspectives of SGD

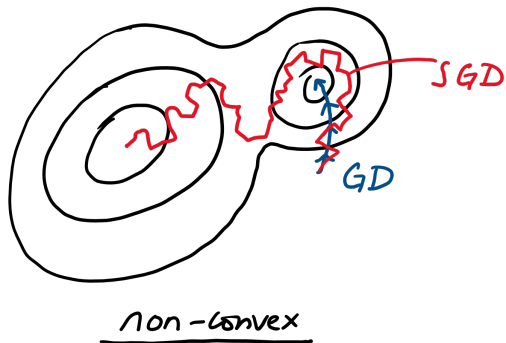
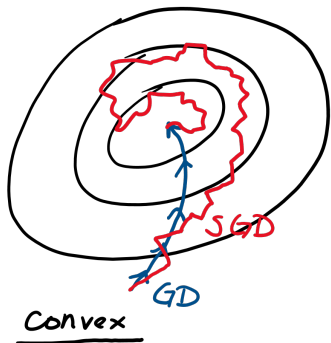
Classical optimization literature have the following observations.

- Compared to GD in **convex** problems:
- SGD offers a **trade-off** between accuracy and efficiency
- More iterations
- Less gradient evaluation per iteration
- Noise is a by-product

Recent studies of SGD for **non-convex** problems found that

- SGD for training deep neural networks works
- SGD finds solution faster
- SGD find a better local minima
- Noise matters

GD compared to SGD



Smoothing the Landscape

Analyzing SGD is an active research topic. Here is one by Kleinberg et al. (<https://arxiv.org/pdf/1802.06175.pdf> ICML 2018)

- The SGD step can be written as GD + noise:

$$\begin{aligned}\mathbf{x}^{t+1} &= \mathbf{x}^t - \eta(\nabla f(\mathbf{x}^t) + \mathbf{w}^t) \\ &= \underbrace{\mathbf{x}^t - \eta\nabla f(\mathbf{x}^t)}_{\stackrel{\text{def}}{=} \mathbf{y}^t} - \eta\mathbf{w}^t.\end{aligned}$$

- \mathbf{y}^t is the “ideal” location returned by GD.
- Let us analyze \mathbf{y}^{t+1} :

$$\begin{aligned}\mathbf{y}^{t+1} &\stackrel{\text{def}}{=} \mathbf{x}^{t+1} - \eta\nabla f(\mathbf{x}^{t+1}) \\ &= (\mathbf{y}^t - \eta\mathbf{w}^t) - \eta\nabla f(\mathbf{y}^t - \eta\mathbf{w}^t)\end{aligned}$$

- Assume $\mathbb{E}[\mathbf{w}] = 0$, then

$$\mathbb{E}[\mathbf{y}^{t+1}] = \mathbf{y}^t - \eta\nabla\mathbb{E}[f(\mathbf{y}^t - \eta\mathbf{w}^t)]$$

Smoothing the Landscape

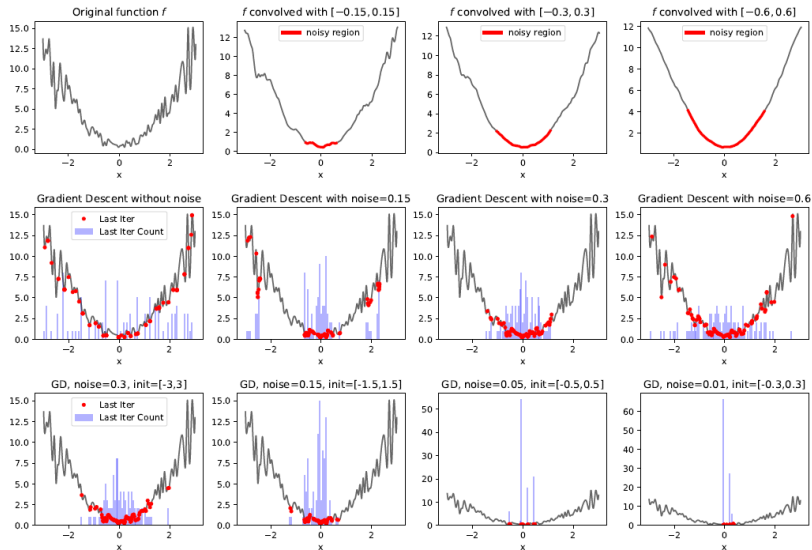
- Let us look at $\mathbb{E}[f(\mathbf{y}^t - \eta \mathbf{w}^t)]$:

$$\mathbb{E}[f(\mathbf{y} - \eta \mathbf{w})] = \int f(\mathbf{y} - \eta \mathbf{w}) p(\mathbf{w}) d\mathbf{w},$$

where $p(\mathbf{w})$ is the distribution of \mathbf{w} .

- $\int f(\mathbf{y} - \eta \mathbf{w}) p(\mathbf{w}) d\mathbf{w}$ is the **convolution** between f and p .
- $p(\mathbf{w}) \geq 0$ for all \mathbf{w} , so the convolution always **smooths** the function.
- Learning rate controls the smoothness
- Too small: Under-smooth. You have not yet escaped from bad local minimum.
- Too large: Over-smooth. You may miss a local minimum.

Smoothing the Landscape



Reading List

Gradient Descent

- S. Boyd and L. Vandenberghe, “Convex Optimization”, Chapter 9.2-9.4.
- J. Nocedal and S. Wright, “Numerical Optimization”, Chapter 3.1-3.3.
- Y. Nesterov, “Introductory lectures on convex optimization”, Chapter 2.
- CMU 10.725 Lecture <https://www.stat.cmu.edu/~ryantibs/convexopt/lectures/grad-descent.pdf>

Stochastic Gradient Descent

- CMU 10.725 Lecture <https://www.stat.cmu.edu/~ryantibs/convexopt/lectures/stochastic-gd.pdf>
- Kleinberg et al. (2018) “When Does SGD Escape Local Minima”, <https://arxiv.org/pdf/1802.06175.pdf>