

LECTURE #12: GRADIENT DESCENT VARIANTS

Instructor: Aditya Bhaskara

Scribe: Yosuke Mizutani

CS 5966/6966: Theory of Machine Learning

February 17th, 2022

Abstract

This lecture reviewed the basic theorem of the gradient descent for a convex, L -Lipschitz function and introduced its variants, online convex optimization and stochastic gradient descent. Further, we learned that more structure on the function, such as smoothness and strong convexity, provides a faster convergence rate.

1 RECAP: GRADIENT DESCENT ANALYSIS

Gradient descent is an iterative algorithm for finding a (local) minimum of a differentiable function. Particularly, we are interested in finding the minimizer (or the minimum value) for a convex function.

Gradient Descent for Convex Functions

Given a convex function $f : \mathcal{D} \rightarrow \mathbb{R}$ defined over a convex domain $\mathcal{D} \subseteq \mathbb{R}^d$, the algorithm starts with some feasible point $w_0 \in \mathcal{D}$. Then, for $t = 0, 1, \dots, T-1$ iteratively set $w_{t+1} = \Pi(w_t - \eta_t \nabla f(w_t))$, where η_t is the learning rate at step t and $\Pi : \mathbb{R}^d \rightarrow \mathcal{D}$ is a projection to the feasible set.

A function $f : \mathcal{D} \rightarrow \mathbb{R}$ is convex if for all $x, y \in \mathcal{D}$ and $0 \leq t \leq 1$, $f(tx + (1-t)y) \leq tf(x) + (1-t)f(y)$. Any local minimum of a convex function is the global minimum.

The gradient of a scalar-valued multi-variable function $f(x_1, x_2, \dots, x_d)$, denoted by ∇f , is given by $\nabla f = \left(\frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2}, \dots, \frac{\partial f}{\partial x_d} \right)$.

Projection to the feasible set can be “hard” if the domain is not “simple”.

1.1 Basic theorem

For analysis, we assume that the convex function f is L -Lipschitz and its domain is all of \mathbb{R}^d . A function is called L -Lipschitz when:

$$(1) \quad |f(x) - f(y)| \leq L\|x - y\| \quad \text{for every } x, y$$

This implies the following, although we left the proof as an exercise.

$$(2) \quad \|\nabla f(x)\| \leq L \quad \text{for every } x$$

Let $w^* \in \mathbb{R}^d$ be the optimum minimizer, that is, $w^* = \arg \min_{x \in \mathbb{R}^d} f(x)$. Then, we introduce a parameter $B \in \mathbb{R}$, which bounds how bad the starting point is.

$$(3) \quad |w_0 - w^*| \leq B$$

For the given step t , how far is $f(w_t)$ from $f(w^*)$? By the convexity of f , we have:

$$(4) \quad f(w^*) \geq f(w_t) + \langle w^* - w_t, \nabla f(w_t) \rangle$$

One may use the multi-variable Mean Value Theorem, which states that for every x, y , there exists $z \in [x, y]$ such that $f(x) - f(y) = \langle \nabla f(z), x - y \rangle$.

For a differentiable function f , the tangent at x is the linear function $\ell(y) = f(x) + \langle \nabla f(x), y - x \rangle$. If f is convex, it must hold that $f(y) \geq \ell(y)$.

And we get:

$$(5) \quad \langle w_t - w^*, \nabla f(w_t) \rangle \geq f(w_t) - f(w^*)$$

Now, let us define the *potential function* $\Phi_t := \|w_t - w^*\|^2$ for the given t , and consider the *potential drop* $\Phi_t - \Phi_{t+1}$.

$$\begin{aligned}
(6) \quad \Phi_t - \Phi_{t+1} &= \|w_t - w^*\|^2 - \|w_{t+1} - w^*\|^2 \\
&= \|w_t - w^*\|^2 - \|w_t - \eta \nabla f(w_t) - w^*\|^2 && w_{t+1} = w_t - \eta \nabla f(w_t) \\
&= \|w_t - w^*\|^2 - \|(w_t - w^*) - \eta \nabla f(w_t)\|^2 \\
&= \|w_t - w^*\|^2 - && \text{(norm of a sum of vectors)} \\
&\quad \left[\|w_t - w^*\|^2 - 2 \langle w_t - w^*, \eta \nabla f(w_t) \rangle + \|\eta \nabla f(w_t)\|^2 \right] \\
&= 2\eta \langle w_t - w^*, \nabla f(w_t) \rangle - \eta^2 \|\nabla f(w_t)\|^2 \\
&\geq 2\eta [f(w_t) - f(w^*)] - \eta^2 \|\nabla f(w_t)\|^2 && \text{from (5)} \\
&\geq 2\eta [f(w_t) - f(w^*)] - \eta^2 L^2 && \text{from (2)}
\end{aligned}$$

Therefore, we get:

$$(7) \quad f(w_t) - f(w^*) \leq \frac{\Phi_t - \Phi_{t+1}}{2\eta} + \frac{L^2\eta}{2}$$

Finally, consider the sum over these values after T steps.

$$\begin{aligned}
(8) \quad \sum_{t=0}^{T-1} f(w_t) - f(w^*) &\leq \sum_{t=0}^{T-1} \frac{\Phi_t - \Phi_{t+1}}{2\eta} + \frac{L^2\eta}{2} \\
&= \frac{1}{2\eta} \left(\sum_{t=0}^{T-1} \Phi_t - \Phi_{t+1} \right) + T \cdot \frac{L^2\eta}{2} \\
&= \frac{\Phi_0 - \Phi_T}{2\eta} + T \cdot \frac{L^2\eta}{2} && \text{(telescoping sum)} \\
&= \frac{\|w_0 - w^*\|^2 - \|w_T - w^*\|^2}{2\eta} + T \cdot \frac{L^2\eta}{2} \\
&\leq \frac{\|w_0 - w^*\|^2}{2\eta} + T \cdot \frac{L^2\eta}{2} \\
&\leq \frac{B^2}{2\eta} + T \cdot \frac{L^2\eta}{2} && \text{from (3)}
\end{aligned}$$

This gives the following theorem.

1.1 THEOREM. Consider running T steps of gradient descent with a fixed learning rate η . Then we have:

$$\frac{1}{T} \sum_{t=1}^T f(w_t) - f(w^*) \leq \frac{B^2}{2\eta T} + \frac{L^2\eta}{2}$$

Here we use indices $t \in [1, T]$, and (8) still holds.

When we set $\eta = \frac{B}{L\sqrt{T}}$, we get a nicer form:

$$(9) \quad \frac{1}{T} \sum_{t=1}^T f(w_t) - f(w^*) \leq \frac{BL}{\sqrt{T}}$$

We say this is $\frac{1}{\sqrt{T}}$ -convergence.

1.2 Constrained domain

The same proof works if we had a constrained domain.

Let $w_{t+\frac{1}{2}} = w_t - \eta \nabla f(w_t)$, which may or may not be feasible. Then, we have the following due to the convex domain.

$$(10) \quad \|w^* - \Pi(w_{t+\frac{1}{2}})\|^2 \leq \|w^* - w_{t+\frac{1}{2}}\|^2$$

It turns out that the projection to the feasible set does not violate any inequalities in the main proof.

1.3 Different functions at time steps

The proof works even if functions at different time steps were different. Let f_1, f_2, \dots, f_T be all convex, L -Lipschitz functions.

We update our gradient descent to $w_{t+1} = w_t - \nabla f_t(w_t)$. Now the potential drop becomes:

$$(11) \quad \begin{aligned} \Phi_t - \Phi_{t+1} &= 2\eta \langle \nabla f_t(w_t), w_t - w^* \rangle - \eta^2 \|\nabla f_t(w_t)\|^2 \\ &\geq 2\eta [f_t(w_t) - f_t(w^*)] - \eta^2 L^2 \end{aligned}$$

This implies the following result.

$$(12) \quad \sum_{t=1}^T f_t(w_t) - f_t(w^*) \leq \frac{B^2}{2\eta} + \frac{L^2 \eta T}{2}$$

When $\eta = \frac{B}{L\sqrt{T}}$, we have:

$$(13) \quad \sum_{t=1}^T f_t(w_t) - f_t(w^*) \leq \sqrt{T}$$

2 ONLINE CONVEX OPTIMIZATION

Now, let us see some application of gradient descent. The previous result is directly applicable to the online convex optimization, where a learner makes a series of decisions to minimize the total loss, and loss functions (convex, L -Lipschitz) f_1, \dots, f_T over the same domain \mathcal{D} are given sequentially.

We want our total loss to be comparable with the best *fixed* minimizer x in

hindsight. That is, $x^* = \arg \min_{x \in \mathcal{D}} \sum_{t=1}^T f_t(x)$.

Think about learning game where adversary chooses functions.

From (13), we obtain:

$$(14) \quad \left(\sum_{t=1}^T f_t(x_t) \right) - \min_{x \in \mathcal{D}} \sum_{t=1}^T f_t(x) \leq \sqrt{T}$$

When there are k switches, the right-hand side of (14) becomes \sqrt{kT} .

The left-hand side of (14) is called *regret*. If we want to compare our decisions to a dynamic minimizer, it is called *switching regret* or *dynamic regret*.

3 STOCHASTIC GRADIENT DESCENT

In stochastic gradient descent (SGD), we do not require to use the same function f for every iteration, but we randomly choose a function g such that the expected value of g over the possible functions is equal to f (i.e. $f = \mathbb{E}[g]$).

For example, ERM (empirical risk minimization) for training data $(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)$ can be viewed as loss:

$$(15) \quad f : \frac{1}{m} \sum_{i=1}^m \ell(h_w(x_i), y_i)$$

$[m]$ denotes $[1, m] = \{1, \dots, m\}$.

Imagine we sample an index $i \sim [m]$ at random (with replacement). If we define $g_i(w) = \ell(h_w(x_i), y_i)$, then we have:

$$(16) \quad f(w) = \frac{1}{m} \sum_{i=1}^m g_i(w)$$

From the linearity of the gradient, $\nabla f = \frac{1}{m} \sum_{i=1}^m \nabla g_i$.

$f = \mathbb{E}[h_t]$ holds even if we had larger batch sizes.

Now, consider T iterations where at each step t , we pick $i_t \sim [m]$. We write $h_t(w) = g_{i_t}(w)$ for the function chosen at step t as a random variable. Observe that at every step t , $\mathbb{E}[h_t] = f$.

The stochastic gradient descent works as follows.

Stochastic Gradient Descent

Given a convex function $f : \mathcal{D} \rightarrow \mathbb{R}$ defined over a convex domain $\mathcal{D} \subseteq \mathbb{R}^d$ and a distribution \mathcal{X} of functions $(\mathcal{D} \rightarrow \mathbb{R})$ such that $f = \mathbb{E}_{h \sim \mathcal{X}}[h]$, the algorithm starts with some feasible point $w_0 \in \mathcal{D}$.

Then, for $t = 0, 1, \dots, T-1$ iteratively pick a random function $h_t \sim \mathcal{X}$ and set $w_{t+1} = \Pi(w_t - \eta \nabla h_t(w_t))$, where η is the learning rate and $\Pi : \mathbb{R}^d \rightarrow \mathcal{D}$ is a projection to the feasible set.

Intuitively, we can decompose a function f into a number of simple functions g_1, \dots, g_m so that (16) holds. And for each step we randomly pick one of the g_i 's and perform gradient descent as usual.

3.1 Analysis of Stochastic Gradient Descent

Here we want to prove the following theorem.

3.1 THEOREM. Let \bar{w} be the average point over w_1, \dots, w_T chosen by stochastic gradient descent, that is, $\bar{w} = \frac{1}{T} \sum_{i=1}^T w_i$, and w^* be the optimum minimizer $w^* =$

$\arg \min_{x \in \mathcal{D}} f(x)$. Then,

$$\mathbb{E}[f(\bar{w})] - f(w^*) \leq \frac{B^2}{2\eta T} + \frac{L^2\eta}{2},$$

where $|w_0 - w^*| \leq B$ and L is a bound on the Lipschitz constant of all the g_i 's.

Notice that in SGD w_i 's and \bar{w} are random variables. From Jensen's inequality, we have:

See Section 14.1.1 (page 186) of Shalev-Shwartz and Ben-David's book.

$$\begin{aligned} (17) \quad f(\bar{w}) - f(w^*) &= f\left(\frac{1}{T} \sum_{t=1}^T w_t\right) - f(w^*) \\ &\leq \frac{1}{T} \left(\sum_{t=1}^T f(w_t)\right) - f(w^*) \\ &= \frac{1}{T} \sum_{t=1}^T f(w_t) - f(w^*) \end{aligned}$$

A key idea to proceed the proof is to keep track of the *expected* potential drop.

$$\begin{aligned} (18) \quad \mathbb{E}_{h_t}[\Phi_t - \Phi_{t+1} \mid w_t] &= \mathbb{E}_{h_t}[2\eta \langle \nabla h_t(w_t), w_t - w^* \rangle] - \eta^2 \|\nabla h_t(w_t)\|^2 \\ &\geq 2\eta \cdot \mathbb{E}_{h_t}[\langle \nabla h_t(w_t), w_t - w^* \rangle] - \eta^2 L^2 && \text{from (2)} \\ &= 2\eta \cdot \langle \mathbb{E}_{h_t}[\nabla h_t(w_t)], w_t - w^* \rangle - \eta^2 L^2 \\ &= 2\eta \cdot \langle \nabla f(w_t), w_t - w^* \rangle - \eta^2 L^2 \\ &\geq 2\eta(f(w_t) - f(w^*)) - \eta^2 L^2 && \text{from (5)} \end{aligned}$$

The expectations are over the choice of h_t .

Now, the remaining task is to accumulate the left-hand side of (18). Notice that it still telescopes; what we condition on doesn't matter!

Convince yourself with $T = 2$.

$$(19) \quad \sum_{t=1}^T \mathbb{E}_{h_t}[\Phi_t - \Phi_{t+1} \mid w_t] = \mathbb{E}_{h_1, h_2, \dots, h_T}[\Phi_1 - \Phi_{T+1}] \leq \Phi_0 \leq B^2$$

The accumulated right-hand side would be:

$$\begin{aligned} \sum_{t=1}^T 2\eta(f(w_t) - f(w^*)) - \eta^2 L^2 &= 2\eta \left[\sum_{t=1}^T f(w_t) - f(w^*) \right] - T\eta^2 L^2 \\ &\geq 2\eta T (f(\bar{w}) - f(w^*)) - T\eta^2 L^2 && \text{from (17)} \\ (20) \quad &= \eta T \left[2(\mathbb{E}[f(\bar{w})] - f(w^*)) - \eta L^2 \right] \end{aligned}$$

By combining (19) and (20), we obtain the following, which completes the proof.

$$\begin{aligned} (21) \quad B^2 &\geq \eta T \left[2(\mathbb{E}[f(\bar{w})] - f(w^*)) - \eta L^2 \right] \\ \mathbb{E}[f(\bar{w})] - f(w^*) &\leq \frac{B^2}{2\eta T} + \frac{L^2\eta}{2} \end{aligned}$$

4 MORE STRUCTURE ON FUNCTION

4.1 Smooth functions

A function is *smooth* if its gradient is also Lipschitz (i.e. the gradient does not change rapidly). The gradient descent for smooth functions achieves $\frac{1}{T}$ -convergence. This topic was covered in Lecture #13.

4.2 Strongly convex functions

A function f is *strongly convex* if $f(y) \geq f(x) + \langle \nabla f(x), y - x \rangle + \mu \|y - x\|^2$ for some $\mu \geq 0$ and for all x, y . In this case, an error is bounded by e^{-T} after T steps.

A function is convex when $\mu = 0$.

Refer to Section 14.4.4 (page 195) of Shalev-Shwartz and Ben-David's book.