

Course Notes for EE227C (Spring 2018): Convex Optimization and Approximation

Instructor: Moritz Hardt

Email: `hardt+ee227c@berkeley.edu`

Graduate Instructor: Max Simchowitz

Email: `msimchow+ee227c@berkeley.edu`

January 23, 2018

2 Lecture 2: Gradient Descent

2.1 Gradient Descent

The procedure of gradient descent is defined by the recursion:

$$x_{t+1} = x_t - \eta \nabla f(x_t)$$

where η is the step size. This works to solve the problem

$$\min_{x \in \Omega} f(x)$$

for f convex, differentiable, and L -Lipschitz

Definition 2.1 (L -Lipschitz). A function is said to be L -Lipschitz if its gradient is bounded,

$$\|\nabla f(x)\| \leq L$$

Fact 2.2. $f(x)$ is L -Lipschitz implies that the difference between two points in the range is bounded,

$$|f(x) - f(y)| \leq L\|x - y\|$$

Question 2.1. How do we ensure that $x_{t+1} \in \Omega$?

Solution: Project onto Ω

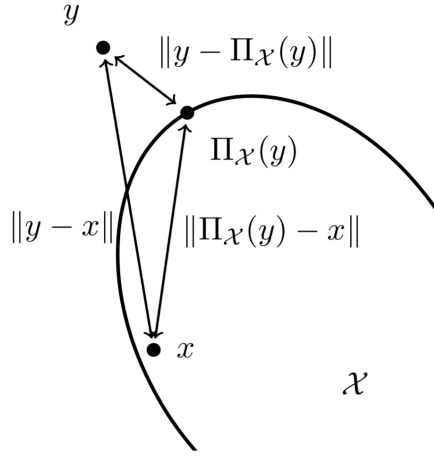


Figure 1: Projection of y onto set \mathcal{X} .

Definition 2.3 (Projection). The *projection* of a point y onto a set Ω is defined as

$$\Pi_{\Omega}(x) = \underset{y \in \Omega}{\operatorname{argmin}} \|x - y\|$$

Example 2.4. A projection onto the Euclidean ball B_2 is just normalization:

$$\Pi_{B_2}(x) = \frac{x}{\|x\|}$$

The crucial property of projections is that they satisfy the following condition:

$$\|\Pi_{\Omega}(y) - x\|^2 \leq \|y - x\|^2$$

i.e. the projection of y onto a convex set containing x is closer to x . See [Figure ??](#) for a geometric picture.

Lemma 2.5.

$$\|\Pi_{\Omega}(y) - x\|^2 \leq \|y - x\|^2 - \|y - \Pi_{\Omega}(y)\|^2$$

Which follows from the Pythagorean theorem. Note that this lemma implies the above property.

2.1.1 Modifying Gradient Descent with Projections

So now we can modify our original procedure to use two steps.

$$y_{t+1} = x_t - \eta \nabla f(x_t)$$

$$x_{t+1} = \Pi_{\Omega}(y_{t+1})$$

And we are guaranteed that $x_{t+1} \in \Omega$. Note that computing the projection may be the hardest part of your problem, as you are computing an argmin. However, there are convex sets for which we know explicitly how to compute the projection (see [Example 2.4](#)).

Theorem 2.6 (Projected Gradient Descent for Lipchitz Functions). *Assume that function f is convex, differentiable, and closed with bounded gradients. Let L be the Lipchitz constant of f over the convex domain Ω . Let R be the upper bound on the distance from the initial point x_1 to the optimal point $x^* = \arg \min_{x \in \Omega} f(x)$ (i.e. $\|x_1 - x^*\|_2$). Let t be the number of iterations of project gradient descent. If the learning rate η is set to $\eta = \frac{R}{L\sqrt{t}}$, then*

$$f\left(\frac{1}{t} \sum_{s=1}^t x_s\right) - f(x^*) \leq \frac{RL}{\sqrt{t}}.$$

This means that the difference between the functional value of the average point during the optimization process from the optimal value is bounded above by a constant proportional to $\frac{1}{\sqrt{t}}$.

Before proving the theorem, recall that

- First order characterization of convexity: $f(y) = f(x) + \nabla f(x)^\top (y - x)$
- "Fundamental Theorem of Optimization": An inner product can be written as a sum of norms: $u^\top v = \frac{1}{2}(\|u\|^2 + \|v\|^2 - \|u - v\|^2)$. This property can be seen by writing $\|u - v\|^2$ as $\|u - v\|^2 = \|u\|^2 + \|v\|^2 - 2u^\top v$.
- L -Lipchitz: For all x , $\|\nabla f(x)\| \leq L$.
- Pythagorean Theorem: $\|\pi_{\Omega}(y) - x\|^2 \leq \|y - x\|^2 - \|y - \pi_{\Omega}(x)\|^2$

Proof of Theorem 2.6 for compact sets. The proof begins by first bounding the difference in function values $f(x_s) - f(x^*)$.

$$f(x_s) - f(x^*) \leq \nabla f(x_s)^\top (x_s - x^*) \quad (1)$$

$$= \frac{1}{\eta} (x_s - y_{s+1})^\top (x_s - x^*) \quad (2)$$

$$= \frac{1}{2\eta} \left(\|x_s - x^*\|^2 + \|x_s - y_{s+1}\|^2 - \|y_{s+1} - x^*\|^2 \right) \quad (3)$$

$$= \frac{1}{2\eta} \left(\|x_s - x^*\|^2 - \|y_{s+1} - x^*\|^2 \right) + \frac{\eta}{2} \|\nabla f(x_s)\|^2 \quad (4)$$

$$\leq \frac{1}{2\eta} \left(\|x_s - x^*\|^2 - \|y_{s+1} - x^*\|^2 \right) + \frac{\eta L^2}{2} \quad (5)$$

$$\leq \frac{1}{2\eta} \left(\|x_s - x^*\|^2 - \|x_{s+1} - x^*\|^2 \right) + \frac{\eta L^2}{2} \quad (6)$$

$$(7)$$

Equation 1 comes from the definition of convexity. Equation 2 comes from the update rule for projected gradient descent. Equation 3 comes from the “Fundamental Theorem of Optimization.” Equation 4 comes from the update rule for projected gradient descent. Equation 5 is because f is L -Lipchitz. Equation 6 comes from the Pythagorean Theorem.

Now, sum these differences from $s = 1$ to $s = T$:

$$\sum_{s=1}^t f(x_s) - f(x^*) \leq \frac{1}{2\eta} \sum_{s=1}^t \left(\|x_s - x^*\|^2 - \|x_{s+1} - x^*\|^2 \right) + \frac{\eta L^2 t}{2} \quad (8)$$

$$= \frac{1}{2\eta} \left(\|x_1 - x^*\|^2 - \|x_t - x^*\|^2 \right) + \frac{\eta L^2 t}{2} \quad (9)$$

$$\leq \frac{1}{2\eta} \left(\|x_1 - x^*\|^2 + \frac{\eta L^2 t}{2} \right) \quad (10)$$

$$\leq \frac{R^2}{2\eta} + \frac{\eta L^2 t}{2} \quad (11)$$

$$(12)$$

Equation 9 is because Equation 8 is a telescoping sum. Equation 10 is because $\|x_t - x^*\|^2 \geq 0$. Equation 11 is by the assumption that $\|x_1 - x^*\|^2 \leq R^2$.

Then bound $f\left(\frac{1}{t} \sum_{s=1}^t x_s\right) - f(x^*)$ by the above sum:

$$f\left(\frac{1}{t}\sum_{s=1}^t x_s\right) \leq \frac{1}{t}\sum_{s=1}^t f(x_s) \quad (13)$$

$$\left(\frac{1}{t}\sum_{s=1}^t x_s\right) - f(x^*) \leq \frac{1}{t}\sum_{s=1}^t f(x_s) - f(x^*) \quad (14)$$

$$\leq \frac{R^2}{2\eta} + \frac{\eta L^2 t}{2} \quad (15)$$

Equation 13 is by convexity. $\frac{R^2}{2\eta} + \frac{\eta L^2 t}{2}$, the upper bound of the difference between $f\left(\frac{1}{t}\sum_{s=1}^t x_s\right)$ and $f(x^*)$ is minimized when η is set to be $\frac{RL}{\sqrt{t}}$. ■