

Lecture 4 — 09/15

Lecturer: Dimitris Papailiopoulos

Scribe: Muxuan Liang

4.1 Properties of Functions

4.1.1 Some useful properties of functions

Suppose that $f : \mathbf{X} \rightarrow \mathbb{R}$, $\nabla f(\mathbf{x})$ exists and is a column-based gradient vector. We will discuss about the following properties of a function

- L -Lipschitz: We say a function is L -Lipschitz if $|f(\mathbf{x}) - f(\mathbf{y})| \leq L\|\mathbf{x} - \mathbf{y}\|_2$.
- β -smooth: We say a function is β -smooth if $\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\| \leq \beta\|\mathbf{x} - \mathbf{y}\|_2$.
- λ -strong convex: We say a function is λ -strong convex if $f(\mathbf{x}) \geq f(\mathbf{y}) + \langle \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle + \frac{\lambda}{2}\|\mathbf{x} - \mathbf{y}\|_2^2$. This implies that strong convex function is always bounded by a quadratic function from below.

4.1.2 Examples

Table 4.1.2 is a summary for some common functions and their properties regarding as convex, L -Lipschitz, and β -smooth.

	Convexity	L -Lipschitz	β -smooth
$ x $	Convex	1	
x^2	Convex	$\max_{x \in D} 2x$	2
$\log(1 + e^x)$	Convex	1	1/4
$\mathbf{w}^\top \mathbf{x} + b$	Convex	$\ \mathbf{w}\ _2$	0
$g(\mathbf{w}^\top \mathbf{x} + b)$	if g convex	$\ \mathbf{w}\ _2 L_g$, if g L_g -Lipschitz	$\ \mathbf{w}\ _2^2 \beta_g$, if g β_g -smooth

4.1.3 Some useful results

Suppose that we have a family of convex function \mathcal{F} , then $\sup_{f \in \mathcal{F}} f(\mathbf{x})$ is also a convex function. If a function f is L -Lipschitz, then $\|\nabla f(\mathbf{x})\|_2 \leq L$. If a function f is β -smooth, then $\|\nabla f(\mathbf{x})\|_2^2 \leq \beta f(\mathbf{x})$. If a function f is λ -strong convex, then $f(\mathbf{x}) - \frac{\lambda}{2}\|\mathbf{x}\|_2^2$ is also convex.

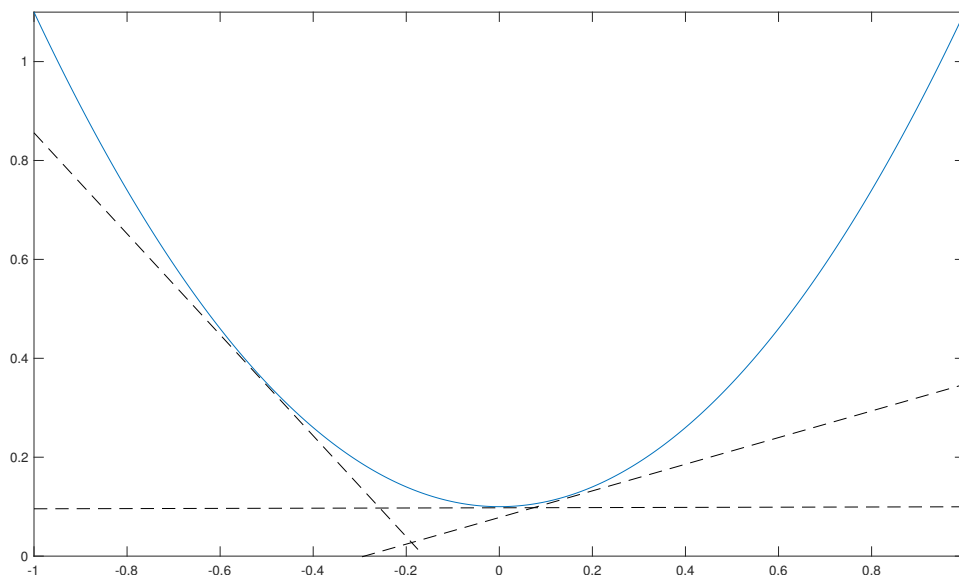


Figure 4.1. The figure show the linearization of a quadratic function. The dashed lines are linearization and solid curve is the quadratic function. The horizontal dashed line implies the gradient of the minimizer is 0.

4.2 Importance of convexity

Suppose $f(\mathbf{x})$ is convex and differentiable, for any point \mathbf{x}_0 , the linearization of $f(\mathbf{x})$ is a minorization at point \mathbf{x}_0 . To be specific,

$$f(\mathbf{x}) \geq f(\mathbf{x}_0) + \langle \nabla f(\mathbf{x}_0), \mathbf{x} - \mathbf{x}_0 \rangle, \quad (4.1)$$

which implies that $\nabla f(\mathbf{x}^*) = 0$, where $f(\mathbf{x}^*) = \min f(\mathbf{x})$. Figure (??) show the relationship of the linearization and the minimizer.

4.3 Gradient Descent

Our goal is to get $\min_{\mathbf{x}} f(\mathbf{x})$ and minimizer \mathbf{x}^* . Suppose that we have a guess or last iteration of the minimizer \mathbf{x}_k , we can write our next guess or iteration $\mathbf{x}_{k+1} = \mathbf{x}_k + \mathbf{u}_k$. We hope that when $k \rightarrow +\infty$, $\|\nabla f(\mathbf{x}_k)\| \rightarrow 0$. Our idea is to get \mathbf{x}_{k+1} through the following optimization

$$\mathbf{x}_{k+1} = \arg \min_{\mathbf{x}} \{f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{x} - \mathbf{x}_k \rangle + \frac{1}{2\gamma} \|\mathbf{x} - \mathbf{x}_k\|_2^2\}. \quad (4.2)$$

We denote the right-hand side of equation (4.2) as $g(\mathbf{x})$. The first two term is the linearization of $f(\mathbf{x})$ at \mathbf{x}_k , which gives the descent direction. The last term is the regularization to keep x close to \mathbf{x}_k . Note that equation (4.2) is a quadratic programming without

constraints. We can take the derivative and set it to be 0. Then we can get

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \gamma \nabla f(\mathbf{x}_k). \quad (4.3)$$

4.3.1 Convergence of Gradient Descent

We will investigate the convergence result of the Gradient Descent method (4.3).

Theorem 4.1. *If we set $\gamma = \frac{\Delta_1}{L\sqrt{T}}$, where L is the Lipschitz constant of $f(\mathbf{x})$ and $\Delta_k = \|\mathbf{x}_k - \mathbf{x}^*\|_2$, then*

$$f\left(\frac{1}{T} \sum_{k=1}^T \mathbf{x}_k\right) - f(\mathbf{x}^*) \leq \frac{\Delta_1 L}{\sqrt{T}}. \quad (4.4)$$

Proof: We start from the convexity of f . We know that

$$\begin{aligned} f(\mathbf{x}_k) - f(\mathbf{x}^*) &\leq \langle \nabla f(\mathbf{x}_k), \mathbf{x}_k - \mathbf{x}^* \rangle \\ &= \langle \frac{1}{\gamma}(\mathbf{x}_{k+1} - \mathbf{x}_k), \mathbf{x}_k - \mathbf{x}^* \rangle \\ &= \frac{1}{2\gamma} (\Delta_k^2 - \Delta_{k+1}^2 + \gamma^2 \|\nabla f(\mathbf{x}_k)\|_2^2). \end{aligned}$$

For f is L -Lipschitz, $\|\nabla f(\mathbf{x}_k)\|_2^2 \leq L^2$, then we have the following inequality

$$f(\mathbf{x}_k) - f(\mathbf{x}^*) \leq \frac{1}{2\gamma} (\Delta_k^2 - \Delta_{k+1}^2) + \frac{\gamma}{2} L^2.$$

We take $k = 1, \dots, T$ and add those inequality up. Then we can get

$$\sum_{k=1}^T f(\mathbf{x}_k) - f(\mathbf{x}^*) \leq \frac{1}{2\gamma} \Delta_1^2 + \frac{\gamma}{2} T L^2.$$

By convexity and optimize over γ , when we set $\gamma = \frac{\Delta_1}{L\sqrt{T}}$, we have

$$f\left(\frac{1}{T} \sum_{k=1}^T \mathbf{x}_k\right) - f(\mathbf{x}^*) \leq \frac{\Delta_1 L}{\sqrt{T}}.$$

□