

IFT 6085 - Lecture 2

Basics of convex analysis and gradient descent

Scribes

Winter 2020: Joss Rakotobe

Winter 2019: Andrew Williams, Ankit Vani, Maximilien Le Clei

Winter 2018: Assya Trofimov, Mohammad Pezeshki, Reyhane Askari

Instructor: Ioannis Mitliagkas

1 Introduction

Many machine learning problems involve learning parameters $\theta \in \Theta$ of a function f towards achieving an objective better. Typically, such objectives are characterized by a loss function $L : \Theta \rightarrow \mathbb{R}$, also called the empirical risk, and training the model corresponds to searching the optimal parameters θ^* that minimize this loss.

For example, in supervised learning, θ parameterizes a function $f : X \rightarrow Y$, where any $x \in X$ is an input and any $y \in Y$ is a target label. Then,

$$L(\theta) = \frac{1}{n} \sum_{i=1}^n \ell(f(x_i; \theta), y_i)$$

represents the loss function for a dataset containing n training examples $(x_1, y_1), \dots, (x_n, y_n)$. Here, ℓ is a deterministic function that determines the distance between a target label y_i and the predicted label $y'_i = f(x_i; \theta)$. In this setting, learning is carried out by performing *empirical risk minimization*, which involves optimizing to find parameters $\theta^* \in \arg \min_{\theta \in \Theta} L(\theta)$.

In the first few lectures, we will dive deeper into the basics and theory of optimization that lie at the heart of machine learning. We will step back from the notation we see in machine learning, and start by considering the most general unconstrained optimization problem¹ for a real valued function $f : \mathcal{X} \rightarrow \mathbb{R}$:

$$\min_{x \in \mathcal{X}} f(x)$$

In most of our discussions, we will consider \mathcal{X} or $\text{dom}(f)$ to be the d -dimensional Euclidean space \mathbb{R}^d .

The optimization problem formulated above is NP-hard in general (see [2]. Section 6.6). However, for certain classes of functions f , strong theoretical guarantees and efficient optimization algorithms exist. In this lecture, we consider such a class of functions, called *convex functions* and prove convergence guarantees for an algorithm for convex optimization called *gradient descent*.

2 Convex optimization

This section introduces some concepts in convexity, and then uses them to prove convergence of gradient descent for convex functions.

Although in practice people commonly use the same algorithms for non-convex optimization as they do for convex optimization (e.g. gradient descent), it is important to note that the strong theory for convex optimization algorithms

¹More generally, this would involve an inf instead of min, but in this lecture we keep the notation simple and stick with min.

often breaks down without the convexity assumption. However, ideas from convex analysis and the weakening of certain results can give partial guarantees and offer generalizations for non-convex analysis.

2.1 Background

2.1.1 Lipschitz continuity

Definition 1 (Lipschitz continuity). Let $L \geq 0$. A real-valued function f is L -Lipschitz continuous iff $\forall x, y \in \text{dom } f$,

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

Intuitively, a Lipschitz continuous function is bounded in how fast it can change. Figure 1 illustrates two Lipschitz continuous functions with different Lipschitz constants.

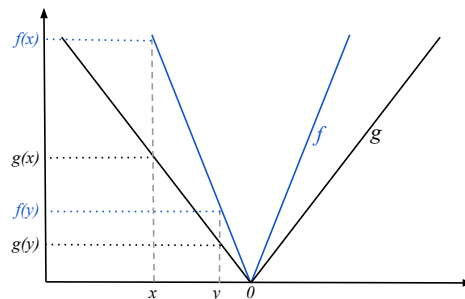


Figure 1: Consider a L_f -Lipschitz continuous function f and a L_g -Lipschitz continuous function g . If f and g are changing as fast as they can, then $L_f > L_g$.

As another example, consider the following function:

$$f(x) = \begin{cases} \exp(-\lambda x), & \text{if } x > 0 \\ 1, & \text{otherwise} \end{cases}$$

$f(x)$ here is L -Lipschitz, and the value of L increases with λ . As the value of λ increases, the function gets closer to discontinuity. In the limit of λ going to ∞ , we recover a step function, which is not Lipschitz continuous for any L . In fact it is not continuous at $x = 0$. This function is illustrated in Figure 2.

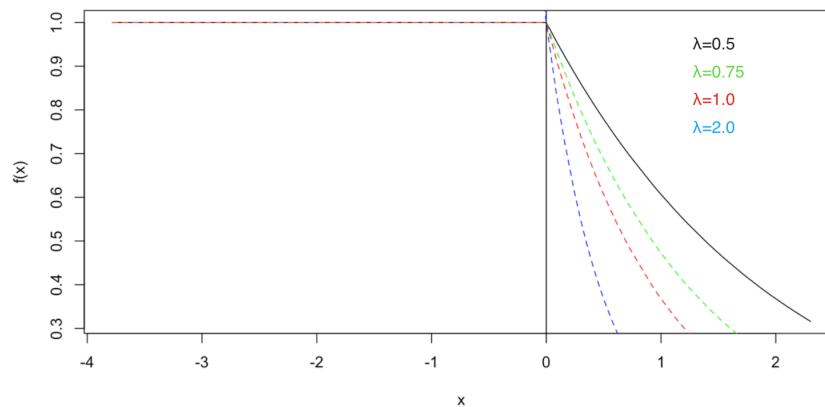


Figure 2: As λ increases in $f(x)$, the Lipschitz constant L increases, and the function gets closer to being discontinuous.

A Lipschitz continuous function does not need to be differentiable. For example, we can have integer-valued and non-smooth Lipschitz continuous functions. However, a corollary of f being L -Lipschitz continuous is that if it is differentiable, the norm of its gradient is bounded by L .

Lemma 1 (First order condition for Lipschitz continuity). *A differentiable function f is L -Lipschitz continuous iff the norm of its gradient is bounded by L :*

$$\forall x \in \text{dom}(f), \|\nabla f(x)\| \leq L$$

Proof. We first give the idea in \mathbb{R} , then generalize in \mathbb{R}^d .

(\implies) For $\text{dom}(f) \subseteq \mathbb{R}$,

$$f'(x) = \lim_{y \rightarrow x} \frac{f(x) - f(y)}{x - y} \quad (\text{definition of derivative})$$

$$\implies |f'(x)| = \lim_{y \rightarrow x} \frac{|f(x) - f(y)|}{|x - y|} \leq L \quad (\text{definition of Lipschitz continuity})$$

In general, if $\text{dom}(f) \subseteq \mathbb{R}^d$ for a differentiable L -Lipschitz function f , then

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

In fact, let $u \in \mathbb{S}^d$ a unit vector, then the directional derivative with respect to u is given by:

$$\nabla_u f(x) = \lim_{h \rightarrow 0} \frac{f(x + hu) - f(x)}{h} \quad (\text{definition of directional derivative})$$

$$\implies |\nabla_u f(x)| = \lim_{h \rightarrow 0} \frac{|f(x + hu) - f(x)|}{\|hu\|} \leq L \quad (\text{definition of Lipschitz continuity})$$

Since $\nabla_u f(x) = \langle \nabla f(x), u \rangle$, taking $u = \frac{\nabla f(x)}{\|\nabla f(x)\|} \implies \nabla_u f(x) = \|\nabla f(x)\|$, thus the result.

(\impliedby) By contraposition, suppose there exist $A, B \in \text{dom}(f)$ such that $|f(A) - f(B)| > L\|A - B\|$.

In the multivariate case, the mean value theorem states that there exists a point C in the line segment \overline{AB} such that $f(A) - f(B) = \langle \nabla f(C), A - B \rangle$.

Hence, by Cauchy-Schwarz inequality, $|f(A) - f(B)| \leq \|\nabla f(C)\| \|A - B\|$, which implies that

$$\|\nabla f(C)\| \geq \frac{|f(A) - f(B)|}{\|A - B\|} > L$$

□

Note that Lipschitz continuity is a special case of continuity: all Lipschitz continuous functions are continuous, but not all continuous functions are Lipschitz continuous (for more information, see [1]).

2.1.2 Convex sets

Before we define convex sets, let us first define a convex combination, which is a constrained version of a linear combination, illustrated in Figure 3 for two points.

Definition 2 (Convex combination). *If $z \in \mathbb{R}^d$ is a linear combination of $x_1, x_2, \dots, x_n \in \mathbb{R}^d$ and the coefficients are non-negative and sum to 1, then z is a convex combination of x_1, x_2, \dots, x_n :*

$$z = \sum_{i=1}^n \theta_i x_i, \quad \text{where } \forall i \in (1, \dots, n), \theta_i \geq 0 \text{ and } \sum_{i=1}^n \theta_i = 1$$

Definition 3 (Convex set). *\mathcal{X} is a convex set if the convex combination any two points in \mathcal{X} is also in \mathcal{X} . That is, for a convex set \mathcal{X} :*

$$\forall x, y \in \mathcal{X}, \forall \theta \in [0, 1], \quad z = \theta x + (1 - \theta)y \in \mathcal{X}$$

Figure 4 gives examples of a convex set and a non-convex set.



Figure 3: All convex combinations $z = \theta x + (1 - \theta)y$ of two points x and y lie on the line segment from x to y . When $\theta = 1$, we get x and when $\theta = 0$, we get y .

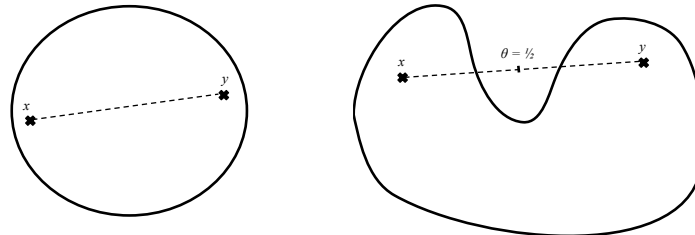


Figure 4: Examples of a convex and a non-convex set. **Left:** Convex set, **Right:** Non-convex set.

2.1.3 Convex functions

Armed with the definition of convex sets, we can finally define convex functions.

Definition 4 (Convex function). A function $f(x)$ is convex iff the domain $\text{dom}(f)$ is a convex set and $\forall x, y \in \text{dom}(f), \forall \theta \in [0, 1]$,

$$f(\theta x + (1 - \theta)y) \leq \theta f(x) + (1 - \theta)f(y)$$

The condition above says that for any two members in the domain of f , the function's value on a convex combination does not exceed the convex combination of those values. It can be viewed as a *zeroth order condition*, as it doesn't involve any condition of differentiability of the function. We will see results with higher order conditions later.

Graphically, when f is convex, for any points x and y in f 's domain, the chord connecting $f(x)$ and $f(y)$ lies above the function between those points. This is illustrated in Figure 5.

For a convex function f , its opposite function $-f$ is defined as a *concave* function. In other words, for a concave function, the inequality condition in Definition 4 is reversed.

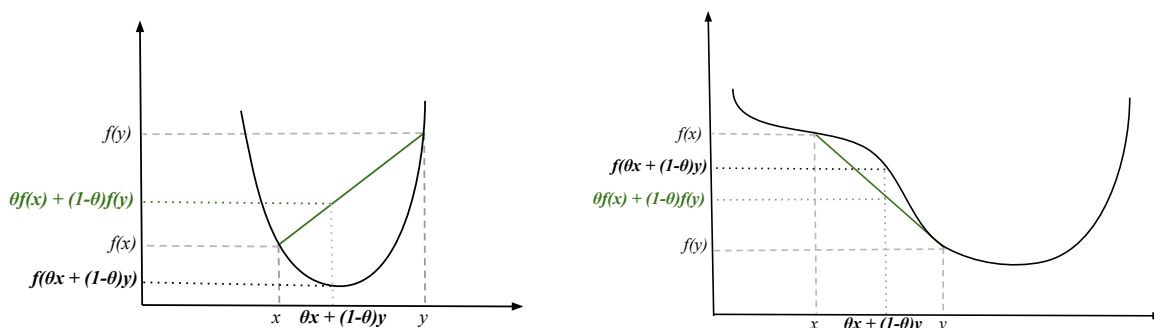


Figure 5: **Left:** Example of a convex function. For any two points $x, y \in \text{dom}(f)$, the chord $\theta f(x) + (1 - \theta)f(y)$, $\theta \in [0, 1]$ lies above the function value $f(\theta x + (1 - \theta)y)$, $\theta \in [0, 1]$. **Right:** Example of a non-convex function². We see here that there exist points x and y for which the chord lies below the curve between $f(x)$ and $f(y)$.

²In fact, this is an example of a quasiconvex function f , meaning that all sublevel sets $S_\alpha(f) = \{x \mid f(x) \leq \alpha\}$ are convex sets.

As we can see in Figure 5, convexity is not a necessary condition for a function to have a unique global minimum. The right plot shows an example of non-convex function having a unique global minimum.

2.2 First-order conditions for convexity

We will see that for differentiable and twice differentiable functions, it is possible to define convexity in terms of first- and second-order conditions for convexity. Note that all the definitions of convexity are equivalent when the appropriate level of differentiability holds (for more information, see [2]).

Lemma 2 (First order condition for convexity). *A differentiable function f is convex iff $\text{dom}(f)$ is convex and $\forall x, y \in \text{dom}(f)$,*

$$f(y) \geq f(x) + \nabla f(x)^\top (y - x)$$

Proof. (\implies) Suppose f is convex. By definition of a convex function, $\text{dom}(f)$ is convex. Let $x, y \in \text{dom}(f)$ and $\theta \in [0, 1]$,

$$\begin{aligned} & f((1 - \theta)x + \theta y) \leq (1 - \theta)f(x) + \theta f(y) && \text{(definition of convexity)} \\ \implies & f(x + \theta(y - x)) \leq f(x) + \theta(f(y) - f(x)) \\ \implies & \frac{f(x + \theta(y - x)) - f(x)}{\theta} \leq f(y) - f(x) \\ \implies & \frac{f(x + \theta u) - f(x)}{\theta} \leq f(y) - f(x) && (u = y - x) \\ \implies & \nabla_u f(x) \leq f(y) - f(x) && \text{(taking } \theta \rightarrow 0) \end{aligned}$$

Since $\nabla_u f(x) = \langle \nabla f(x), u \rangle = \nabla f(x)^\top u = \nabla f(x)^\top (y - x)$, the result follows.

(\impliedby) Suppose $\text{dom}(f)$ is convex and

$$\forall x, y \in \text{dom}(f), f(y) \geq f(x) + \nabla f(x)^\top (y - x)$$

Let $x, y \in \text{dom}(f), \theta \in [0, 1]$.

By the convexity of $\text{dom}(f)$, $z = \theta x + (1 - \theta)y \in \text{dom}(f)$

Applying the above inequation using z :

$$f(x) \geq f(z) + \nabla f(z)^\top (x - z) \tag{1}$$

$$f(y) \geq f(z) + \nabla f(z)^\top (y - z) \tag{2}$$

Multiplying inequation (1) by θ , inequation (2) by $1 - \theta$ and adding, we get:

$$\begin{aligned} \theta f(x) + (1 - \theta)f(y) & \geq \theta(f(z) + \nabla f(z)^\top (x - z)) + (1 - \theta)(f(z) + \nabla f(z)^\top (y - z)) \\ & = \theta f(z) + \nabla f(z)^\top (\theta x - \theta z) + (1 - \theta)f(z) + \nabla f(z)^\top ((1 - \theta)y - (1 - \theta)z) \\ & = f(z) + \nabla f(z)^\top (\theta x + (1 - \theta)y - z) \\ & = f(z) \\ & = f(\theta x + (1 - \theta)y) \end{aligned}$$

□

Intuitively, this says that for a convex function, a tangent of its graph at any point must lie below the graph. This is illustrated in Figure 6.

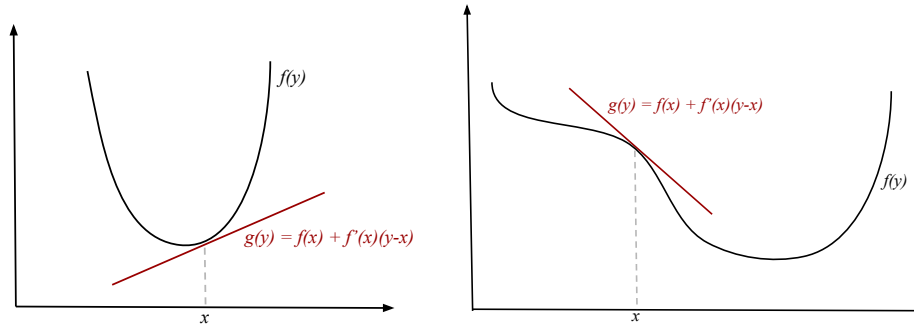


Figure 6: Example of a convex and a non-convex function illustrating the first-order condition for convexity. For the convex function on the left, all possible tangents will lie below the graph. However, for the non-convex function on the right, there exists a tangent such that it lies above the graph for some points in the function's domain.

2.3 First- and second-order conditions for convexity

Before discussing the second-order condition for convexity, let us review the multivariate generalization of a second derivative, namely a *Hessian*:

$$\nabla^2 f(x) = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_d} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_d \partial x_1} & \frac{\partial^2 f}{\partial x_d \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_d^2} \end{bmatrix}$$

For a function $f : \mathbb{R} \rightarrow \mathbb{R}$, $f(x) = \frac{h}{2}x^2$, the second derivative $f''(x) = h$ corresponds to a measure on how quickly the slope of the function can change. Similarly, the Hessian represents the curvature of a function f with $\text{dom}(f) \subseteq \mathbb{R}^d$. A multivariate quadratic function f can be written as $f(x) = \frac{1}{2}x^\top Hx$, where H is the Hessian. The eigenvalues of the Hessian determine the curvature of the function along its eigenvectors. Consider the eigendecomposition³

$$H = Q\Lambda Q^\top \quad (3)$$

where

$$\Lambda = \begin{bmatrix} \lambda_1 & & & \\ & \lambda_2 & & \\ & & \ddots & \\ & & & \lambda_d \end{bmatrix}$$

Changing the basis to Q , we can focus on the directions described by $Q = [q_1, q_2, \dots, q_d]$. Then, along the direction q_i , we get the curvature λ_i . Figure 7 illustrates the curvature of a convex quadratic function. Note that not every quadratic function is convex (see Lemma 4).

Definition 5 (Positive semi-definite). A symmetric matrix $M \in \mathbb{R}^{d \times d}$ is positive semi-definite if

$$\forall x \in \mathbb{R}^d, x^\top Mx \geq 0$$

In that case, we will use the notation $M \succeq 0$.

A positive semi-definite matrix does not necessarily have all positive elements⁵. However, as we will see in the next theorem, its eigenvalues λ_i are non-negative.

³The Schwarz Theorem implies that the Hessian is always a symmetric matrix. Applying the spectral theorem, it can always be decomposed in the mentioned form, with Q an orthogonal matrix, i.e., $Q^{-1} = Q^\top$.

⁴Most eigendecomposition algorithms return eigenvalues in non-decreasing order.

⁵There exist matrices with all positive entries, that are not positive semi-definite, and there exist positive semi-definite matrices that have negative elements

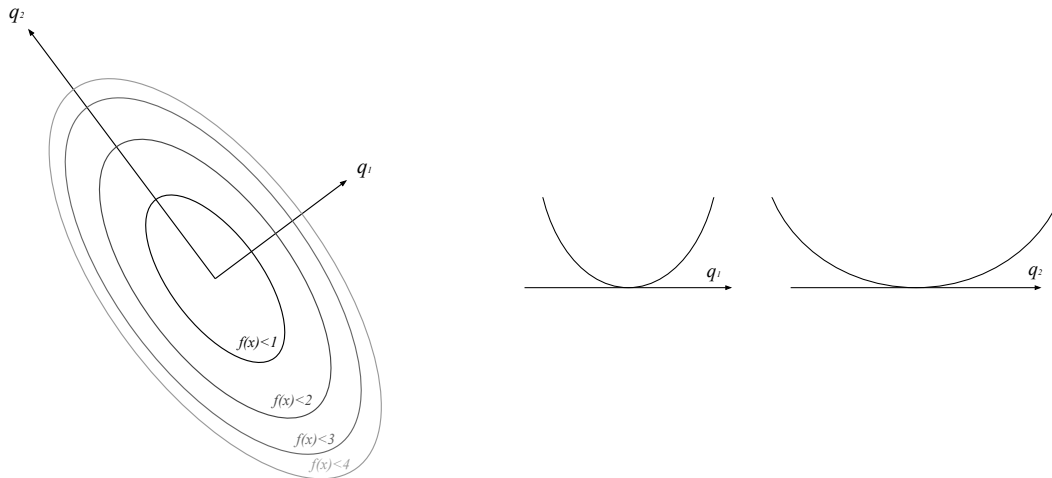


Figure 7: **Left:** Looking along the principle directions of the quadratic function $f(x) = \frac{1}{2}x^\top Hx$, we see that the curve changes faster along q_1 than q_2 . Here⁴, $\lambda_1 > \lambda_2$. **Right:** Cross-sections along q_1 and q_2 , showing that f has higher curvature along q_1 than q_2 .

Theorem 1 (Characterization of a positive semi-definite matrix). A symmetric matrix $M \in \mathbb{R}^{d \times d}$ is positive semi-definite iff all its eigenvalues λ_i are non-negative: $\forall i \in \{1, \dots, d\}$,

$$\lambda_i \geq 0$$

Proof. Suppose $M \succeq 0$ and denote v_i an eigenvector associated to λ_i . Then, using the definition of an eigenvalue and eigenvector, we have $\forall i \in \{1, \dots, d\}$,

$$0 \leq v_i^\top M v_i = v_i^\top \lambda_i v_i = \lambda_i \|v_i\|^2$$

Hence, $\lambda_i \geq 0$.

Suppose $\lambda_i \geq 0$ for $i = 1, \dots, d$. Then, using the eigendecomposition of M , $M = Q\Lambda Q^\top$ as mentioned previously in equation (3), $\forall x \in \mathbb{R}^d$,

$$x^\top M x = x^\top Q \underbrace{\Lambda Q^\top}_{y} x = y^\top \Lambda y = \sum_{i=1}^d \lambda_i y_i^2 \geq 0$$

□

Before stating the second order condition for convexity, we will need a small but useful other lemma.

Lemma 3. A function f is convex iff $\text{dom}(f)$ is a convex set and $\forall x \in \text{dom}(f), v \in \mathbb{R}$, the real function $g_{x,v}$ is convex, where

$$g_{x,v}(t) := f(x + tv)$$

for all $t \in \mathbb{R}$ such that $x + tv \in \text{dom}(f)$.

Proof. (\implies) Suppose f is convex, then by definition of a convex function, $\text{dom}(f)$ is a convex set.

First, we prove that $\text{dom}(g_{x,v})$ is a convex set.

$\forall s, t \in \text{dom}(g_{x,v}), \theta \in [0, 1]$, we have $x + sv \in \text{dom}(f)$ and $x + tv \in \text{dom}(f)$.

Hence by convexity of $\text{dom}(f)$,

$$\theta(x + sv) + (1 - \theta)(x + tv) = x + (\theta s + (1 - \theta)t)v \in \text{dom}(f)$$

Therefore, $\theta s + (1 - \theta)t \in \text{dom}(g_{x,v})$, which means that $\text{dom}(g_{x,v})$ is a convex set.

Let $s, t \in \text{dom}(g_{x,v})$, $\theta \in [0, 1]$, we have :

$$\begin{aligned} g_{x,v}(\theta s + (1 - \theta)t) &= f(x + (\theta s + (1 - \theta)t)v) \\ &= f(\theta(x + sv) + (1 - \theta)(x + tv)) \\ &\leq \theta f(x + sv) + (1 - \theta)f(x + tv) \\ &= \theta g_{x,v}(s) + (1 - \theta)g_{x,v}(t) \end{aligned}$$

Hence, $g_{x,v}$ is convex.

(\Leftarrow) By contraposition, suppose f is not convex and $\text{dom}(f)$ is a convex set. Then, $\exists a, b \in \text{dom}(f)$, $\theta \in [0, 1]$ such that

$$f(\theta a + (1 - \theta)b) > \theta f(a) + (1 - \theta)f(b)$$

Therefore,

$$g_{a,b-a}(\theta) = f(a + \theta(b - a)) = f(\theta b + (1 - \theta)a) > \theta f(a) + (1 - \theta)f(b) = \theta g_{a,b-a}(0) + (1 - \theta)g_{a,b-a}(1)$$

Hence, $g_{a,b-a}$ is not convex. □

The Lemma 3 gives us an equivalence between the convexity of a function in \mathbb{R}^d and the convexity of this function along any line. It will be very helpful in the next proof, as it allows us to derive results in \mathbb{R}^d from results in \mathbb{R} .

Lemma 4 (Second order condition for convexity). *A twice differentiable function f is convex iff $\text{dom}(f)$ is a convex set and $\forall x \in \text{dom}(f)$,*

$$\nabla^2 f(x) \succeq 0$$

Proof. [3] We begin the proof with the case where $\text{dom}(f) \subseteq \mathbb{R}$ and then generalize for $\text{dom}(f) \subseteq \mathbb{R}^d$

Let's suppose $\text{dom}(f) \subseteq \mathbb{R}$. First, notice that the positive semi-definite condition in \mathbb{R} is simply a non-negative condition.

(\Rightarrow) If f is convex, then by Lemma 2, $\forall x \in \text{dom}(f)$, $h > 0$:

$$f(x + h) - f(x) \geq f'(x)h \quad \text{and} \quad f(x) - f(x + h) \geq f'(x + h)(-h)$$

Hence, $f'(x + h) \geq \frac{f(x+h) - f(x)}{h} \geq f'(x)$, which means that f' is non-decreasing, hence f'' is non-negative.

(\Leftarrow) Suppose f'' is non-negative. By Taylor's theorem of order 2, we have $\forall x, y \in \text{dom}(f)$, $\exists z \in [x, y]$ such that:

$$f(y) = f(x) + f'(x)(y - x) + \frac{1}{2}f''(z)(y - x)^2$$

Therefore, $f(y) \geq f(x) + f'(x)(y - x)$, which is the first order condition of convexity.

Now let's move on to the general case, where $\text{dom}(f) \subseteq \mathbb{R}^d$.

Applying Lemma 3, f is convex iff $\forall x \in \text{dom}(f)$, $v \in \mathbb{R}^d$, the real function $g_{x,v}$, defined by $g_{x,v}(t) = f(x + tv)$ for t such that $x + tv \in \text{dom}(f)$, is convex. By the above proof in \mathbb{R} , this is true iff

$$g''_{x,v}(t) = v^\top \nabla^2 f(x + tv)v \geq 0$$

Since this is for all $v \in \mathbb{R}^d$, $x + tv \in \text{dom}(f)$, by taking $t = 0$, we can see that it is equivalent to $\nabla^2 f(x) \succeq 0, \forall x \in \text{dom}(f)$. □

The Lemma states that for a twice differentiable function to be convex, its Hessian must be positive semi-definite.

In general, for any non-negative eigenvalue of the Hessian, the curvature of the function is non-negative along the corresponding eigenvector, and thus the function is convex in that direction. On the other hand, a non-positive eigenvalue represents non-positive curvature along the eigenvector, and the function is concave in that direction. Then, we can see that for a function to be convex, it has to have non-negative curvature, and thus non-negative eigenvalues, in all directions.

2.4 Gradient descent

We will now study our first optimization method. Gradient descent is an iterative optimization algorithm that starts from an initial point, and moves in the direction of the steepest descent. This makes the algorithm especially useful for convex optimization, since for a convex function, any local minimum is also a global minimum. It can be proved that this direction we're looking for is given by the opposite of the gradient, hence the name of the algorithm.

Specifically, starting from an initial guess x_1 , the algorithm generates the sequence $x_1, x_2, \dots, x_T \in \mathbb{R}^d$ to approach the minimum of a differentiable function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ using the following update rule:

$$x_{k+1} = x_k - \gamma \nabla f(x_k),$$

Here, γ is called the *step size*, also known as the *learning rate* in the machine learning literature. If f is convex and γ decays at an appropriate rate, then it is guaranteed that as $T \rightarrow \infty$, $x_T \rightarrow x^*$, where $x^* \in \arg \min_{x \in \text{dom}(f)} f(x)$ is an optimal value.

Lemma 5. *If f is a L -Lipschitz continuous differentiable function, then*

$$\|\nabla f(x_k)\|_2^2 \leq L^2.$$

Proof. See Lemma 1. □

Theorem 2. *Let f be convex and L -Lipschitz continuous⁶. If we take T steps of gradient descent with the step size*

$$\gamma = \frac{\|x_1 - x^*\|_2}{L\sqrt{T}}$$

Then the following holds:

$$f\left(\frac{1}{T} \sum_{k=1}^T x_k\right) - f(x^*) \leq \frac{\|x_1 - x^*\|L}{\sqrt{T}}$$

(see [2]. Theorem 3.2)

⁶Although gradient descent and the theorem rely on the notion of using gradients, the function does not need to be differentiable. The algorithm, theorem and proof still hold when any of the subgradients is used in place of a gradient at points where the function is not differentiable.

Proof. Using the first order condition of convexity and rearranging terms, we can write:

$$\begin{aligned}
f(x_k) - f(x^*) &\leq \langle \nabla f(x_k), x_k - x^* \rangle \\
&= \left\langle \frac{1}{\gamma} (x_k - x_{k+1}), x_k - x^* \right\rangle && \text{(using gradient descent update rule)} \\
&= \frac{1}{2\gamma} (-\|x_k - x_{k+1} - (x_k - x^*)\|_2^2 + \|x_k - x_{k+1}\|_2^2 + \|x_k - x^*\|_2^2) \\
&\hspace{10em} \text{(using } \|a - b\|^2 = \|a\|^2 + \|b\|^2 - 2\langle a, b \rangle \text{ and rearranging terms)} \\
&= \frac{1}{2\gamma} (-\|x_{k+1} - x^*\|_2^2 + \|\gamma \nabla f(x_k)\|_2^2 + \|x_k - x^*\|_2^2) && \text{(using gradient descent update rule)} \\
&= \frac{1}{2\gamma} (\|x_k - x^*\|_2^2 - \|x_{k+1} - x^*\|_2^2) + \frac{\gamma}{2} \|\nabla f(x_k)\|_2^2
\end{aligned}$$

Using Lemma 5, we can thus write:

$$f(x_k) - f(x^*) \leq \frac{1}{2\gamma} (\|x_k - x^*\|_2^2 - \|x_{k+1} - x^*\|_2^2) + \frac{\gamma}{2} L^2 \quad (4)$$

Let us now perform the change of variables by defining $D_k = \|x_k - x^*\|$. Then, from Equation 4, we have:

$$\begin{aligned}
f(x_1) - f(x^*) &\leq \frac{1}{2\gamma} (D_1^2 - D_2^2) + \frac{\gamma}{2} L^2 \\
f(x_2) - f(x^*) &\leq \frac{1}{2\gamma} (D_2^2 - D_3^2) + \frac{\gamma}{2} L^2 \\
&\vdots \\
f(x_{T-1}) - f(x^*) &\leq \frac{1}{2\gamma} (D_{T-1}^2 - D_T^2) + \frac{\gamma}{2} L^2 \\
f(x_T) - f(x^*) &\leq \frac{1}{2\gamma} (D_T^2 - D_{T+1}^2) + \frac{\gamma}{2} L^2 \leq \frac{1}{2\gamma} D_T^2 + \frac{\gamma}{2} L^2
\end{aligned}$$

Adding all the terms above, we get a telescopic sum where most of the D_k terms cancel. We get:

$$\begin{aligned}
\sum_{k=1}^T (f(x_k) - f(x^*)) &\leq \frac{1}{2\gamma} D_1^2 + \frac{T\gamma L^2}{2} \\
\implies \left(\frac{1}{T} \sum_{k=1}^T f(x_k) \right) - f(x^*) &\leq \frac{1}{2\gamma T} D_1^2 + \frac{\gamma L^2}{2} \quad (5)
\end{aligned}$$

Since f is convex, Jensen's inequality tells us that $f\left(\frac{1}{T} \sum_{k=1}^T x_k\right) \leq \frac{1}{T} \sum_{k=1}^T f(x_k)$. Thus, we can rewrite Equation 5 as:

$$f\left(\frac{1}{T} \sum_{k=1}^T x_k\right) - f(x^*) \leq \frac{1}{2\gamma T} D_1^2 + \frac{\gamma L^2}{2} \quad (6)$$

Plugging in $\gamma = \frac{\|x_1 - x^*\|_2}{L\sqrt{T}}$ gives us the result

$$f\left(\frac{1}{T} \sum_{k=1}^T x_k\right) - f(x^*) \leq \frac{\|x_1 - x^*\|_2}{\sqrt{T}}$$

□

To understand how we derived the optimal γ in the above theorem, notice that the RHS of Equation 6 is a convex function of γ in the domain $\mathbb{R}_{\geq 0}$. The minimizing γ would give us the tightest bound, which can be found analytically

using convexity by setting the gradient to zero and solving for γ .

The convergence rate we derived here for gradient descent is $O(1/\sqrt{T})$, which is quite slow. We will see in the following lectures how stronger assumptions on the function f can guarantee significantly faster convergence rates for gradient descent. However, we obtain a convergence in average, which is very useful for some type of functions, such as the ones that are not differentiable at their minimum.

References

- [1] S. Boyd and L. Vandenberghe. *Convex optimization*. Cambridge university press, 2004.
- [2] S. Bubeck et al. Convex optimization: Algorithms and complexity. *Foundations and Trends® in Machine Learning*, 8(3-4):231–357, 2015.
- [3] G. Hall. Lecture on convex functions, 2016. URL http://www.princeton.edu/~aaa/Public/Teaching/ORF523/S16/ORF523_S16_Lec7_gh.pdf.