

# *Optimization for Data Science*

*Cristian Perez Jensen*

*March 11, 2024*

## Contents

1	<i>Risk minimization</i>	1
1.1	<i>Algorithms in data science</i>	1
1.2	<i>Empirical and expected risk</i>	1
1.3	<i>The map of learning</i>	3
2	<i>Theory of convex functions</i>	4
2.1	<i>Mathematical background</i>	4
2.2	<i>Convex sets</i>	4
2.3	<i>Convex functions</i>	5
2.4	<i>Minimizing convex functions</i>	9
2.5	<i>Convex programming</i>	10
3	<i>Gradient descent</i>	14
3.1	<i>Vanilla analysis</i>	14
3.2	<i>Lipschitz convex functions</i>	15
3.3	<i>Smooth functions</i>	16
3.4	<i>Strongly convex functions</i>	18

*List of symbols*

$\doteq$	Equality by definition
$\mathbb{R}$	Set of real numbers
$f : A \rightarrow B$	Function $f$ that maps elements of set $A$ to elements of set $B$
$\mathbf{v} \in \mathbb{R}^n$	$n$ -dimensional vector
$\mathbf{M} \in \mathbb{R}^{m \times n}$	$m \times n$ matrix
$\mathbf{M}^\top$	Transpose of matrix $\mathbf{M}$
$\mathbf{M}^{-1}$	Inverse of matrix $\mathbf{M}$
$\det(\mathbf{M})$	Determinant of $\mathbf{M}$
$\frac{d}{dx}f(x)$	Ordinary derivative of $f(x)$ w.r.t. $x$ at point $x \in \mathbb{R}$
$\frac{\partial}{\partial x}f(x)$	Partial derivative of $f(x)$ w.r.t. $x$ at point $x \in \mathbb{R}^n$
$\nabla f(x) \in \mathbb{R}^n$	Gradient of $f : \mathbb{R}^n \rightarrow \mathbb{R}$ at point $x \in \mathbb{R}^n$
$Df(x) \in \mathbb{R}^{n \times m}$	Jacobian of $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ at point $x \in \mathbb{R}^n$
$\nabla^2 f(x) \in \mathbb{R}^{n \times n}$	Hessian of $f : \mathbb{R}^n \rightarrow \mathbb{R}$ at point $x \in \mathbb{R}^n$



## 1 Risk minimization

### 1.1 Algorithms in data science

In classical algorithm theory, an optimization problem solves a well-defined problem. For example, Kruskal's algorithm computes the minimum spanning tree of a graph. In data science, it is not as well-defined. The starting point is a learning problem, and the optimization typically happens on training data. However, even a perfect result may fail to solve the learning problem, which is a failure of the model in which the optimization algorithm was applied, rather than the optimization algorithm itself.

### 1.2 Empirical and expected risk

Typically, we have a data source  $\mathcal{X}$ , equipped with an unknown probability distribution. However, we do have access to independent samples  $X_1, \dots, X_n \sim \mathcal{X}$ . The goal is to "explain"  $\mathcal{X}$  through these samples. More specifically, we have a class  $\mathcal{H}$  of hypotheses (possible explanations of  $\mathcal{X}$ ). The goal is then to select the hypothesis  $H \in \mathcal{H}$  that best explains  $\mathcal{X}$ , which we measure by a *risk* (or *loss*) function  $\ell : \mathcal{H} \times \mathcal{X} \rightarrow \mathbb{R}$ .

The expected risk can be computed by

$$\ell(H) \doteq \mathbb{E}_{\mathcal{X}}[\ell(H, X)].$$

The best explanation is the *expected risk minimizer*,

$$H^* = \operatorname{argmin}_{H \in \mathcal{H}} \ell(H).$$

However, since we do not have access to the distribution over  $\mathcal{X}$ , we cannot compute  $\ell(H)$  or  $H^*$ . Thus, we need to compute the *probably approximately correct* (PAC) hypothesis. This means that given  $\delta, \epsilon > 0$ , we want to produce, with probability  $1 - \delta$ , a hypothesis  $\tilde{H} \in \mathcal{H}$  such that

$$\ell(\tilde{H}) \leq \inf_{H \in \mathcal{H}} \ell(H) + \epsilon,$$

meaning that with high probability, we approximately solve the expected risk minimization problem.

However, we can still not compute this, thus we must use our training data to compute the *empirical risk*,

$$\ell_n(H) = \frac{1}{n} \sum_{i=1}^n \ell(H, X_i).$$

This is a random variable, because it depends on the training data, which are all random variables, distributed according to the probability distribution of  $\mathcal{X}$ .

**Lemma 1** (Weak law of large numbers). Let  $H \in \mathcal{H}$  be a hypothesis. For any  $\delta, \epsilon > 0$ , there exists  $n_0 \in \mathbb{N}$  such that for  $n \geq n_0$ ,

$$|\ell_n(H) - \ell(H)| \leq \epsilon,$$

with probability at least  $1 - \delta$ .

Given  $n \in \mathbb{N}$  and training data  $X_1, \dots, X_n$  from  $\mathcal{X}$ , we want to produce a hypothesis  $\tilde{H}_n$  such that

$$\ell_n(\tilde{H}_n) \leq \inf_{H \in \mathcal{H}} \ell_n(H) + \epsilon.$$

In an ideal world,  $\tilde{H}_n$  is also almost the best explanation for the data source  $\mathcal{X}$ ,

$$\ell(\tilde{H}_n) \leq \inf_{H \in \mathcal{H}} \ell(H) + \epsilon.$$

**Remark.** Note that the weak law of large numbers can only be applied to a *fixed* hypothesis, but not to the data-dependent hypothesis  $\tilde{H}_n$ . Thus, we are not always in an ideal world scenario.

A sufficient condition for an ideal world scenario is that the weak law of large numbers uniformly holds for all hypotheses with high probability. This leads us to the following theorem.

**Theorem 2.** Assume that for any  $\delta, \epsilon > 0$ , there exists  $n_0 \in \mathbb{N}$  such that for  $n \geq n_0$ ,

$$\sup_{H \in \mathcal{H}} |\ell_n(H) - \ell(H)| \leq \epsilon,$$

with probability at least  $1 - \delta$ . Then, for  $n \geq n_0$ , an approximate empirical risk minimizer  $\tilde{H}_n$  is PAC for expected risk minimization, meaning that it satisfies

$$\ell(\tilde{H}_n) \leq \inf_{H \in \mathcal{H}} \ell(H) + 3\epsilon,$$

with probability at least  $1 - \delta$ .

*Proof.*

$$\begin{aligned} \ell(\tilde{H}_n) &\leq \ell_n(\tilde{H}_n) + \epsilon \\ &\leq \inf_{H \in \mathcal{H}} \ell_n(H) + 2\epsilon \\ &\leq \inf_{H \in \mathcal{H}} \ell(H) + 3\epsilon, \end{aligned}$$

with probability at least  $1 - \delta$ . ■

It turns out that the assumption made by Theorem 2 holds in many relevant cases, but it is not always true.

Follows from  $\sup_{H \in \mathcal{H}} |\ell_n(H) - \ell(H)| \leq \epsilon$ .

$\tilde{H}_n$  is an almost best explanation of the training data.

Follows from  $\sup_{H \in \mathcal{H}} |\ell_n(H) - \ell(H)| \leq \epsilon$ .

In this course, we will not learn how to pick the theory ( $\mathcal{H}$  and  $\ell$ ), but rather how to solve the optimization problems that arise in empirical risk minimization after the theory has been chosen.

### 1.3 The map of learning

The map of learning can be seen in Figure 1. It plots the empirical risk ( $\ell_n(H_n)$ , training data) against the expected risk ( $\ell(H_n)$ , estimated by a validation set). We can only be in the area denoted by “empirical risk minimization”, because

$$\begin{aligned}\ell_n(\tilde{H}_n) &\leq \inf_{H \in \mathcal{H}} \ell_n + \epsilon \\ &\leq \ell_n(\tilde{H}) + \epsilon \\ &\leq \ell(\tilde{H}) + 2\epsilon \\ &\leq \ell(\tilde{H}_n) + \epsilon\epsilon.\end{aligned}$$

A model is overfit when we have low empirical risk, while having high expected risk. This means that the explanation quality on the data source is much worse than on the training data. The main reason for this is that the theory ( $\mathcal{H}$  and  $\ell$ ) is so complex that it can explain almost anything.

A model is underfit when we have high empirical risk, while having high expected risk. This means that we neither explain the training data nor the data source. The main reason for this is that the theory is too simple to capture the nature of the data.

The model is learning when we have low empirical risk and low expected risk. This means that the training was successful. Generalization occurs when the expected risk is close to the empirical risk. Note that this does not mean that the explanation is good, since any “blind explanation” will generalize well due to the weak law of large numbers. Ideally, we want generalization and learning.



Figure 1. The map of learning.

## 2 Theory of convex functions

### 2.1 Mathematical background

**Theorem 3** (Cauchy-Schwarz inequality). Let  $\mathbf{u}, \mathbf{v} \in \mathbb{R}^d$ , then

$$|\mathbf{u}^\top \mathbf{v}| \leq \|\mathbf{u}\| \|\mathbf{v}\|.$$

For non-zero vectors, this is equivalent to

$$-1 \leq \frac{\mathbf{u}^\top \mathbf{v}}{\|\mathbf{u}\| \|\mathbf{v}\|} \leq 1.$$

**Definition 4** (Spectral norm). Let  $A$  be an  $m \times d$  matrix, then

$$\|A\| \doteq \max_{\mathbf{v} \in \mathbb{R}^d, \mathbf{v} \neq \mathbf{0}} \frac{\|A\mathbf{v}\|}{\|\mathbf{v}\|} = \max_{\|\mathbf{v}\|=1} \|A\mathbf{v}\|.$$

Intuitively, this means that  $\|A\|$  is the largest factor by which a vector can be stretched in length under the mapping  $\mathbf{v} \mapsto A\mathbf{v}$ .

**Theorem 5** (Mean value theorem). Let  $a < b$  be real numbers, and let  $h : [a, b] \rightarrow \mathbb{R}$  be a continuous function that is differentiable on  $(a, b)$ . Then, there exists  $c \in (a, b)$  such that

$$h'(c) = \frac{h(b) - h(a)}{b - a}.$$

Geometrically, this means that between  $a$  and  $b$ , there is a tangent to the graph of  $h$  that has the same slope; see Figure 2.

**Theorem 6** (Fundamental theorem of calculus). Let  $a < b$  be real numbers, and let  $h : \text{dom}(h) \rightarrow \mathbb{R}$  be a differentiable function on an open domain  $\text{dom}(h) \supset [a, b]$ , and such that  $h'$  is continuous on  $[a, b]$ . Then,

$$h(b) - h(a) = \int_a^b h'(t) dt.$$

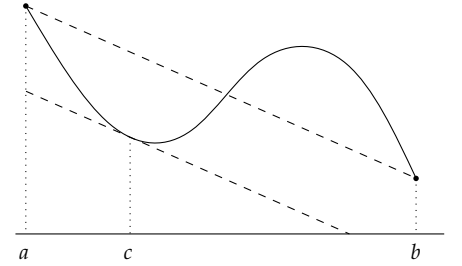


Figure 2. Illustration of the mean value theorem.

### 2.2 Convex sets

**Definition 7** (Convex set). A set  $\mathcal{C} \subseteq \mathbb{R}^d$  is convex if the line segment between any two points of  $\mathcal{C}$  lies in  $\mathcal{C}$ . I.e., if for any  $\mathbf{x}, \mathbf{y} \in \mathcal{C}$  and any  $\lambda$  with  $0 \leq \lambda \leq 1$ , we have

$$\lambda \mathbf{x} + (1 - \lambda) \mathbf{y} \in \mathcal{C}.$$

**Observation.** Let  $\mathcal{C}_i, i \in I$  be convex sets, where  $I$  is a (possibly infinite)

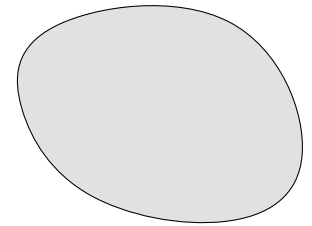
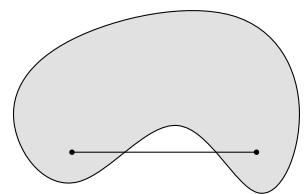


Figure 3. Example of a convex set in  $\mathbb{R}^2$ .





index set. Then,

$$\mathcal{C} = \bigcap_{i \in I} \mathcal{C}_i,$$

is a convex set.

### 2.3 Convex functions

**Definition 8** (Convexity). A function  $f : \text{dom}(f) \rightarrow \mathbb{R}$  is convex if (i)  $\text{dom}(f)$  is convex and (ii) for all  $x, y \in \text{dom}(f)$  and all  $\lambda \in [0, 1]$ , we have

$$f(\lambda x + (1 - \lambda)y) \leq \lambda f(x) + (1 - \lambda)f(y).$$

Geometrically, this condition means that the line segment connecting the two points  $(x, f(x))$  and  $(y, f(y)) \in \mathbb{R}^{d+1}$  lies pointwise above the graph of  $f$ ; see Figure 5.

**Definition 9** (Epigraph). The epigraph of a function  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  is defined as

$$\text{epi}(f) \doteq \left\{ (x, \alpha) \in \mathbb{R}^{d+1} \mid x \in \text{dom}(f), \alpha \geq f(x) \right\}.$$

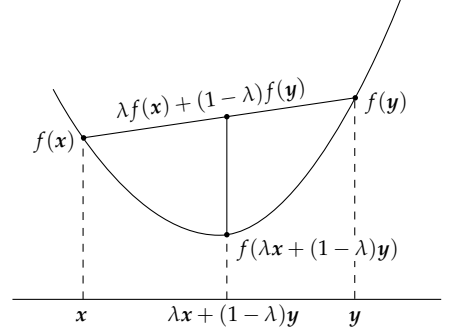
**Observation.**  $f$  is a convex function if and only if  $\text{epi}(f)$  is a convex set.

**Lemma 10** (Jensen's inequality). Let  $f$  be convex,  $x_1, \dots, x_m \in \text{dom}(f)$ ,  $\lambda_1, \dots, \lambda_m \in \mathbb{R}_+$  such that  $\sum_{i=1}^m \lambda_i = 1$ , then

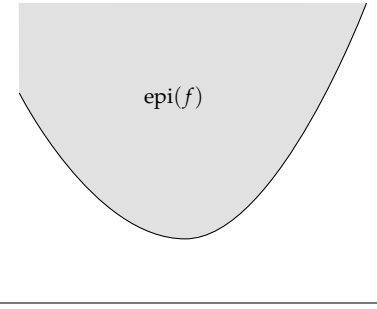
$$f\left(\sum_{i=1}^m \lambda_i x_i\right) \leq \sum_{i=1}^m \lambda_i f(x_i).$$

*Proof.* We will prove Jensen's inequality by induction. The base case ( $m = 2$ ) is true by definition of convexity. Let Jensen's inequality hold for  $m = k$ . Consider  $m = k + 1$ , then

$$\begin{aligned} f\left(\sum_{i=1}^{k+1} \lambda_i x_i\right) &= f\left(\sum_{i=1}^k \lambda_i x_i + \lambda_{k+1} x_{k+1}\right) \\ &= f\left((1 - \lambda_{k+1}) \sum_{i=1}^k \frac{\lambda_i}{1 - \lambda_{k+1}} x_i + \lambda_{k+1} x_{k+1}\right) \\ &\leq (1 - \lambda_{k+1}) f\left(\sum_{i=1}^k \frac{\lambda_i}{1 - \lambda_{k+1}} x_i\right) + \lambda_{k+1} f(x_{k+1}) \\ &\leq (1 - \lambda_{k+1}) \sum_{i=1}^k \frac{\lambda_i}{1 - \lambda_{k+1}} f(x_i) + \lambda_{k+1} f(x_{k+1}) \\ &= \sum_{i=1}^{k+1} \lambda_i f(x_i). \end{aligned}$$



**Figure 5.** Illustration of the classic definition of convexity.



**Figure 6.** Epigraph of a convex function.

By definition of convexity.

Induction assumption,  $\sum_{i=1}^k \frac{\lambda_i}{1 - \lambda_{k+1}} = 1$ .

Since it holds for the base case and the induction step, Jensen's inequality must hold for all  $m$ . ■

**Lemma 11.** Let  $f$  be convex and suppose that  $\text{dom}(f) \subseteq \mathbb{R}^d$  is open, then  $f$  is continuous.

**Definition 12** (Differentiable functions). Let  $f : \text{dom}(f) \rightarrow \mathbb{R}^m$  where  $\text{dom}(f) \subseteq \mathbb{R}^d$  is open.  $f$  is called differentiable at  $x \in \text{dom}(f)$  if there exists an  $m \times d$  matrix  $A$  and an error function  $r : \mathbb{R}^d \rightarrow \mathbb{R}^m$  defined around  $0 \in \mathbb{R}^d$  such that for all  $y$  in some neighborhood of  $x$ ,

$$f(y) = f(x) + A(y - x) + r(y - x),$$

where  $\lim_{v \rightarrow 0} \|r(v)\|/\|v\| = 0$  (error function  $r$  is sublinear around  $0$ ).  $A$  is unique and called the Jacobian matrix of  $f$  at  $x$ .

**Lemma 13** (First-order convexity). Suppose that  $\text{dom}(f)$  is open and that  $f$  is differentiable. In particular, the gradient

$$\nabla f(x) \doteq \left[ \frac{\partial f(x)}{\partial x_1}, \dots, \frac{\partial f(x)}{\partial x_d} \right]$$

exists at every point  $x \in \text{dom}(f)$ .

Then,  $f$  is convex if and only if (i)  $\text{dom}(f)$  is convex and (ii)

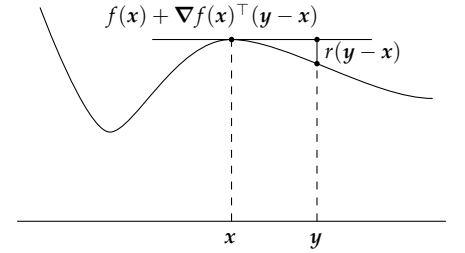
$$\forall x, y \in \text{dom}(f) : f(y) \geq f(x) + \nabla f(x)^\top (y - x).$$

Geometrically, this means that the graph is above all tangent hyperplanes; see Figure 8.

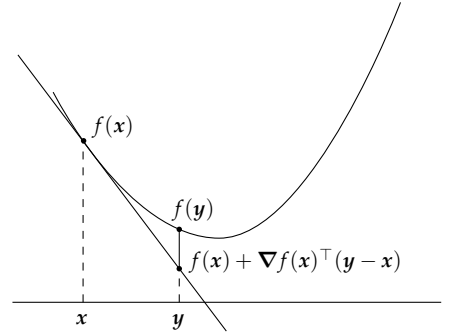
*Proof.* We will first prove that convexity implies that the first-order definition holds. Then, we will prove that the first-order definition implies convexity, making the definitions equivalent for differentiable functions  $f$ .

$\Rightarrow$ : Suppose  $f$  is convex. Then, for all  $t \in (0, 1)$ ,

$$\begin{aligned} f((1-t)x + ty) &\leq (1-t)f(x) + tf(y) \\ \Leftrightarrow f(x + t(y-x)) &\leq f(x) + t(f(y) - f(x)) \\ \Leftrightarrow f(y) &\geq f(x) + \frac{f(x + t(y-x)) - f(x)}{t} \\ &= f(x) + \frac{\nabla f(x)^\top t(y-x) + r(t(y-x))}{t} \\ &= f(x) + \nabla f(x)^\top (y-x) + \underbrace{\frac{r(t(y-x))}{t}}_{\rightarrow 0 \text{ for } t \rightarrow 0}. \end{aligned}$$



**Figure 7.** Graph of the affine function  $f(x) + \nabla f(x)^\top (y - x)$  is a tangent hyperplane to the graph of  $f$  at  $(x, f(x))$ .



**Figure 8.** Illustration of the first-order characterization of convexity (Lemma 13).

Definition of convexity

$\Leftarrow$ : Define  $\mathbf{z} \doteq \lambda \mathbf{x} + (1 - \lambda) \mathbf{y} \in \text{dom}(f)$  for  $\lambda \in [0, 1]$  by convexity of  $\text{dom}(f)$ . Then, we have the following inequalities,

$$\begin{aligned} f(\mathbf{x}) &\geq f(\mathbf{z}) + \nabla f(\mathbf{z})^\top (\mathbf{x} - \mathbf{z}) \\ \lambda f(\mathbf{x}) &\geq \lambda \left( f(\mathbf{z}) + \nabla f(\mathbf{z})^\top (\mathbf{x} - \mathbf{z}) \right) \\ f(\mathbf{y}) &\geq f(\mathbf{z}) + \nabla f(\mathbf{z})^\top (\mathbf{y} - \mathbf{z}) \\ (1 - \lambda) f(\mathbf{y}) &\geq (1 - \lambda) \left( f(\mathbf{z}) + \nabla f(\mathbf{z})^\top (\mathbf{y} - \mathbf{z}) \right). \end{aligned}$$

This implies the following,

$$\begin{aligned} \lambda f(\mathbf{x}) + (1 - \lambda) f(\mathbf{y}) &\geq \lambda \left( f(\mathbf{z}) + \nabla f(\mathbf{z})^\top (\mathbf{x} - \mathbf{z}) \right) + (1 - \lambda) \left( f(\mathbf{z}) + \nabla f(\mathbf{z})^\top (\mathbf{y} - \mathbf{z}) \right) \\ &= f(\mathbf{z}) + \lambda \nabla f(\mathbf{z})^\top (\mathbf{x} - \mathbf{z}) + (1 - \lambda) \nabla f(\mathbf{z})^\top (\mathbf{y} - \mathbf{z}) \\ &= f(\mathbf{z}) + \nabla f(\mathbf{z})^\top (\lambda (\mathbf{x} - \mathbf{z}) + (1 - \lambda) (\mathbf{y} - \mathbf{z})) \\ &= f(\mathbf{z}) + \nabla f(\mathbf{z})^\top (\lambda \mathbf{x} + (1 - \lambda) \mathbf{y} - (\lambda \mathbf{z} + (1 - \lambda) \mathbf{z})) \\ &= f(\mathbf{z}) + \nabla f(\mathbf{z})^\top (\mathbf{z} - \mathbf{z}) \\ &= f(\lambda \mathbf{x} + (1 - \lambda) \mathbf{y}). \end{aligned}$$

■

**Lemma 14** (First-order convexity alternative). Suppose that  $\text{dom}(f)$  is open and that  $f$  is differentiable. Then,  $f$  is convex if and only if (i)  $\text{dom}(f)$  is convex and (ii)

$$\forall \mathbf{x}, \mathbf{y} \in \text{dom}(f) : (\nabla f(\mathbf{y}) - \nabla f(\mathbf{x}))^\top (\mathbf{y} - \mathbf{x}) \geq 0.$$

Geometrically, this means that the gradient is monotonic.

*Proof.* We will first prove that it holds from left to right, and then from right to left.

$\Rightarrow$ : If  $f$  is convex, the first-order characterization of convexity (Lemma 13) yields the following,

$$\begin{aligned} f(\mathbf{x}) &\geq f(\mathbf{y}) + \nabla f(\mathbf{y})^\top (\mathbf{x} - \mathbf{y}) \\ f(\mathbf{y}) &\geq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}), \end{aligned}$$

for all  $\mathbf{x}, \mathbf{y} \in \text{dom}(f)$ . Adding up these two inequalities yields

$$\begin{aligned} f(\mathbf{x}) + f(\mathbf{y}) &\geq f(\mathbf{y}) + \nabla f(\mathbf{y})^\top (\mathbf{x} - \mathbf{y}) + f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) \\ \Leftrightarrow 0 &\geq \nabla f(\mathbf{y})^\top (\mathbf{x} - \mathbf{y}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) \\ &= (\nabla f(\mathbf{y}) - \nabla f(\mathbf{x}))^\top (\mathbf{x} - \mathbf{y}) \\ \Leftrightarrow 0 &\leq (\nabla f(\mathbf{y}) - \nabla f(\mathbf{x}))^\top (\mathbf{y} - \mathbf{x}). \end{aligned}$$

$\Leftarrow$ : Define  $\mathbf{z} \doteq (1 - t) \mathbf{x} + t \mathbf{y} \in \text{dom}(f)$  for  $\mathbf{x}, \mathbf{y} \in \text{dom}(f), t \in (0, 1)$  by convexity of  $\text{dom}(f)$ . Observe that  $\mathbf{z} = \mathbf{x} + t(\mathbf{y} - \mathbf{x})$ . Then, we have the

following inequality, according to the monotonicity of the gradient,

$$\begin{aligned} (\nabla f(z) - \nabla f(x))^\top (z - x) &\geq 0 \\ (\nabla f(x + t(y - x)) - \nabla f(x))^\top (x + t(y - x) - x) &\geq 0 \\ (\nabla f(x + t(y - x)) - \nabla f(x))^\top (y - x) &\geq 0. \end{aligned}$$

Divide by  $t$

Let  $h(t) = f(x + t(y - x))$ , then  $h'(t) = \nabla f(x + t(y - x))^\top (y - x)$ . Hence, we can rewrite the inequality as the following,

$$h'(t) \geq \nabla f(x)^\top (y - x).$$

By the mean value theorem, there exists  $c \in (0, 1)$  such that  $h'(c) = h(1) - h(0)$ . I.e.,

$$h'(c) = f(y) - f(x).$$

Thus, we can rewrite the inequality to the following,

$$\begin{aligned} f(y) &= f(x) + h'(c) \\ &\geq f(x) + \nabla f(x)^\top (y - x), \end{aligned}$$

which is the first-order characterization of convexity (Lemma 13). ■

**Lemma 15** (Second-order convexity). Suppose that  $\text{dom}(f)$  is open and that  $f$  is twice differentiable. In particular, the Hessian

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

exists at every point  $x \in \text{dom}(f)$  and is symmetric.

Then,  $f$  is convex if and only if (i)  $\text{dom}(f)$  is convex and (ii) for all  $x \in \text{dom}(f)$ , we have

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

I.e.,  $\nabla^2 f(x)$  is positive semidefinite ( $M$  is positive semidefinite if  $x^\top M x \geq 0$  for all  $x$  and  $x^\top M x > 0$  for all  $x \neq 0$ ).

Geometrically, this means that  $f$  has non-negative curvature everywhere. I.e., the growth rate should be growing.

**Observation** (Operations that preserve convexity). Let  $f_1, \dots, f_m$  be convex functions,  $\lambda_1, \dots, \lambda_m \in \mathbb{R}_+$ , then

$$f \doteq \max_{i=1}^m f_i,$$

and

$$f \doteq \sum_{i=1}^m \lambda_i f_i.$$

are convex on  $\text{dom}(f) \doteq \bigcap_{i=1}^m \text{dom}(f_i)$ .

Let  $f$  be a convex function with  $\text{dom}(f) \subseteq \mathbb{R}^d$ ,  $g: \mathbb{R}^m \rightarrow \mathbb{R}^d$  an affine function, i.e.  $g(x) = Ax + b$  for some  $A \in \mathbb{R}^{d \times m}$ ,  $b \in \mathbb{R}^d$ . Then, the function  $f \circ g$  is convex on  $\text{dom}(f \circ g) \doteq \{x \in \mathbb{R}^m \mid g(x) \in \text{dom}(f)\}$ .

## 2.4 Minimizing convex functions

**Definition 16** (Local minimum). A local minimum of  $f: \text{dom}(f) \rightarrow \mathbb{R}$  is a point  $x$  such that there exists  $\epsilon > 0$  with

$$f(x) \leq f(y) \quad \forall y \in \text{dom}(f) \text{ s.t. } \|y - x\| < \epsilon.$$

**Remark.** This definition of local minima means that in some small neighborhood around  $x$ ,  $x$  is the smallest point.

**Lemma 17.** Let  $x^*$  be a local minimum of a convex function  $f: \text{dom}(f) \rightarrow \mathbb{R}$ . Then,  $x^*$  is a global minimum, meaning that  $f(x^*) \leq f(y) \quad \forall y \in \text{dom}(f)$ .

*Proof.* Proof by contradiction. Suppose there exists  $y \in \text{dom}(f)$  such that  $f(y) < f(x^*)$ . Define  $y' \doteq \lambda x^* + (1 - \lambda)y$  for  $\lambda \in (0, 1)$ . From convexity, we get that  $f(y') < f(x^*)$ , because  $f(y) < f(x^*)$ . If we choose  $\lambda$  so close to 1 such that  $\|y' - x^*\| < \epsilon$  gives the inequality  $f(x^*) \leq f(y')$ . This yields a contradiction, thus there cannot exist a  $y \in \text{dom}(f)$  such that  $f(y) < f(x^*)$ , meaning that  $f(x^*) \leq f(y)$  for all  $y \in \text{dom}(f)$ . ■

**Lemma 18.** Suppose that  $f$  is convex and differentiable over an open domain  $\text{dom}(f)$ . Let  $x \in \text{dom}(f)$ . If  $\nabla f(x) = 0$  (critical point), then  $x$  is a global minimum.

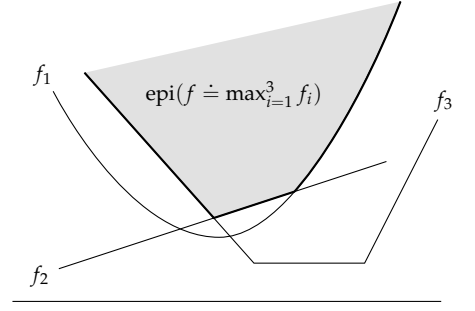
*Proof.* Suppose that  $\nabla f(x) = 0$ . According to the first-order characterization of convexity, we have

$$f(y) \geq f(x) + \underbrace{\nabla f(x)^\top (y - x)}_{=0} = f(x)$$

for all  $y \in \text{dom}(f)$ , so  $x$  is a global minimum. ■

**Definition 19** (Strict convexity). A function  $f: \text{dom}(f) \rightarrow \mathbb{R}$  is strictly convex if (i)  $\text{dom}(f)$  is convex and (ii) for all  $x \neq y \in \text{dom}(f)$  and all  $\lambda \in (0, 1)$ , we have

$$f(\lambda x + (1 - \lambda)y) < \lambda f(x) + (1 - \lambda)f(y).$$



**Figure 9.** The maximum operator over  $m$  convex functions is a convex function. As can be seen, the epigraph of  $f$  is convex.

An example of a strictly convex function can be found in Figure 5, while an example of a function that is not strictly convex can be found in Figure 10.

**Lemma 20.** Let  $f : \text{dom}(f) \rightarrow \mathbb{R}$  be strictly convex. Then,  $f$  has at most one global minimum.

But, how do we know a minimizer exists? To be able to find this out, we must use the Weierstrass theorem.

**Definition 21** (Sublevel set).  $f : \mathbb{R}^d \rightarrow \mathbb{R}, \alpha \in \mathbb{R}$ . The set  $f^{\leq \alpha} \doteq \{x \in \mathbb{R}^d \mid f(x) \leq \alpha\}$  is the  $\alpha$ -sublevel set of  $f$ .

**Theorem 22** (Weierstrass). Let  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  be a continuous function, and suppose there is a non-empty and bounded sublevel set  $f^{\leq \alpha}$ . Then,  $f$  has a global minimum.

*Proof.* We know that  $f$  attains a minimum over the closed and bounded set  $f^{\leq \alpha}$  at some  $x^*$ , because  $f$  is continuous. This  $x^*$  is also a global minimum as it has value  $f(x^*) \leq \alpha$ , while any  $x \notin f^{\leq \alpha}$  has value  $f(x) > \alpha \geq f(x^*)$ . ■

Since convex functions are continuous by Lemma 11, the Weierstrass theorem also applies to convex functions. Thus, to figure out whether a convex function has a minimizer, we need to find any non-empty and bounded sublevel set of this function.

## 2.5 Convex programming

In standard form, optimization problems look like the following,

$$\begin{aligned} & \text{minimize} && f_0(x) \\ & \text{subject to} && f_i(x) \leq 0, \quad i = 1, \dots, m \\ & && h_j(x) = 0, \quad j = 1, \dots, p. \end{aligned}$$

The problem domain is then  $\mathcal{D} = \bigcap_{i=0}^m \text{dom}(f_i) \cap \bigcap_{j=1}^p \text{dom}(h_j)$ .

In a convex program, all  $f_i$  are convex functions, and all  $h_j$  are affine functions with domain  $\mathbb{R}^d$ .

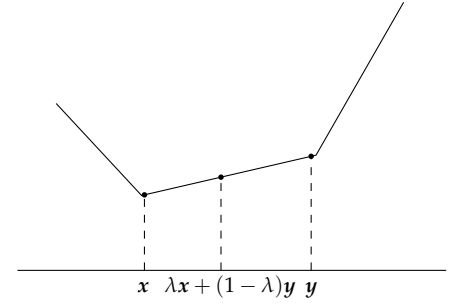


Figure 10. A non-strictly convex function.

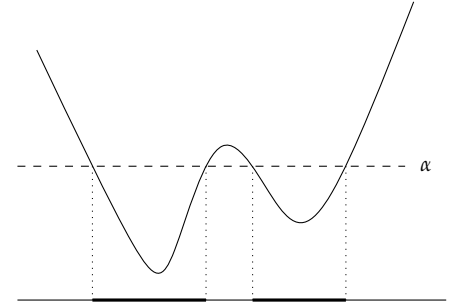


Figure 11. Illustration of a sublevel set of a non-convex function.

**Definition 23** (Lagrange dual function). Given an optimization problem in standard form, its Lagrangian is the function  $L : \mathcal{D} \times \mathbb{R}^m \times \mathbb{R}^p \rightarrow \mathbb{R}$  given by

$$L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\nu}) = f_0(\mathbf{x}) + \sum_{i=1}^m \lambda_i f_i(\mathbf{x}) + \sum_{j=1}^p \nu_j h_j(\mathbf{x}).$$

The  $\lambda_i, \nu_j$  are called Lagrange multipliers. The Lagrange dual function is the function  $g : \mathbb{R}^m \times \mathbb{R}^p \rightarrow \mathbb{R} \cup \{-\infty\}$  defined by

$$g(\boldsymbol{\lambda}, \boldsymbol{\nu}) = \inf_{\mathbf{x} \in \mathcal{D}} L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\nu}).$$

Only the  $(\boldsymbol{\lambda}, \boldsymbol{\nu})$  pairs with  $g(\boldsymbol{\lambda}, \boldsymbol{\nu}) > -\infty$  are interesting.

**Lemma 24** (Weak Lagrange duality). Let  $\mathbf{x}$  be a feasible solution ( $f_i(\mathbf{x}) \leq 0$  for  $i = 1, \dots, m$  and  $h_j(\mathbf{x}) = 0$  for  $j = 1, \dots, p$ ). Let  $g$  be the Lagrange dual function and  $\boldsymbol{\lambda} \in \mathbb{R}^m, \boldsymbol{\nu} \in \mathbb{R}^p$  such that  $\boldsymbol{\lambda} \geq \mathbf{0}$ . Then

$$g(\boldsymbol{\lambda}, \boldsymbol{\nu}) \leq f_0(\mathbf{x}).$$

*Proof.*

$$g(\boldsymbol{\lambda}, \boldsymbol{\nu}) \leq L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\nu}) = f_0(\mathbf{x}) + \underbrace{\sum_{i=1}^m \lambda_i f_i(\mathbf{x})}_{\leq 0} + \underbrace{\sum_{j=1}^p \nu_j h_j(\mathbf{x})}_{=0} \leq f_0(\mathbf{x}).$$

■

However, we want to know what the best lower bound is that we can get in this way. For this, we must choose  $\boldsymbol{\lambda} \geq \mathbf{0}, \boldsymbol{\nu}$  such that  $g(\boldsymbol{\lambda}, \boldsymbol{\nu})$  is maximized. This can be phrased as another optimization problem, called the Lagrange dual,

$$\begin{aligned} & \text{maximize} && g(\boldsymbol{\lambda}, \boldsymbol{\nu}) \\ & \text{subject to} && \boldsymbol{\lambda} \geq \mathbf{0}. \end{aligned}$$

**Observation.** The Lagrange dual is a convex program, even if the primal is not.

**Theorem 25.** Suppose that a convex program has a feasible solution  $\tilde{\mathbf{x}}$  that in addition satisfies  $f_i(\tilde{\mathbf{x}}) < 0, i = 1, \dots, m$  (Slater point). Then, the infimum value of the primal equals the supremum value of its Lagrange dual. Moreover, if the value is finite, it is attained by a feasible solution of the dual,

$$\inf f_0(\mathbf{x}) = \max g(\boldsymbol{\lambda}, \boldsymbol{\nu}).$$

A case of particular interest is that strong duality holds and the joint value is attained in both the primal and dual problem.<sup>1</sup> If all  $f_i$  and  $h_j$  are differentiable, then the Karush-Kuhn-Tucker (KKT) conditions provide necessary and, under convexity, also sufficient conditions for this case to occur.

<sup>1</sup> This would mean that  $\min f_0(\mathbf{x}) = \max g(\boldsymbol{\lambda}, \boldsymbol{\nu})$ .

**Definition 26** (Zero duality gap). Let  $\tilde{\mathbf{x}}$  be feasible for the primal and  $(\tilde{\boldsymbol{\lambda}}, \tilde{\boldsymbol{\nu}})$  feasible for the Lagrange dual. The primal and dual solutions  $\tilde{\mathbf{x}}$  and  $(\tilde{\boldsymbol{\lambda}}, \tilde{\boldsymbol{\nu}})$  are said to have zero duality gap if  $f_0(\tilde{\mathbf{x}}) = g(\tilde{\boldsymbol{\lambda}}, \tilde{\boldsymbol{\nu}})$ .

If the primal and dual have zero duality gap, then  $\min f_0(\mathbf{x}) = \max g(\boldsymbol{\lambda}, \boldsymbol{\nu})$ .

The consequence of zero duality gap is the *master equation*,

$$\begin{aligned} f_0(\tilde{\mathbf{x}}) &= g(\tilde{\boldsymbol{\lambda}}, \tilde{\boldsymbol{\nu}}) \\ &= \inf_{\mathbf{x} \in \mathcal{D}} \left( f_0(\mathbf{x}) + \sum_{i=1}^m \tilde{\lambda}_i f_i(\mathbf{x}) + \sum_{j=1}^p \tilde{\nu}_j h_j(\mathbf{x}) \right) \\ &\leq f_0(\tilde{\mathbf{x}}) + \underbrace{\sum_{i=1}^m \tilde{\lambda}_i f_i(\tilde{\mathbf{x}})}_{\leq 0} + \underbrace{\sum_{j=1}^p \tilde{\nu}_j h_j(\tilde{\mathbf{x}})}_{=0} \\ &\leq f_0(\tilde{\mathbf{x}}), \end{aligned}$$

which means that all inequalities turn into equalities. Thus,

$$\tilde{\lambda}_i f_i(\tilde{\mathbf{x}}) = 0, \quad i = 1, \dots, m,$$

which is called complementary slackness, because if  $\tilde{\lambda}_i \neq 0$ , then  $f_i(\tilde{\mathbf{x}}) = 0$ , and vice versa. Furthermore, if all  $f_i$  and  $h_j$  are differentiable, then

$$\nabla f_0(\tilde{\mathbf{x}}) + \sum_{i=1}^m \tilde{\lambda}_i \nabla f_i(\tilde{\mathbf{x}}) + \sum_{j=1}^p \tilde{\nu}_j \nabla h_j(\tilde{\mathbf{x}}) = \mathbf{0},$$

which is called the vanishing Lagrangian gradient condition.

In summary, we get the following result.

**Theorem 27** (Karush-Kuhn-Tucker necessary conditions). Let  $\tilde{\mathbf{x}}$  and  $(\tilde{\boldsymbol{\lambda}}, \tilde{\boldsymbol{\nu}})$  be feasible solutions of the primal optimization problem and its Lagrange dual, respectively, with zero duality gap. If all  $f_i$  and  $h_j$  are differentiable, then

$$\begin{aligned} \tilde{\lambda}_i f_i(\tilde{\mathbf{x}}) &= 0, \quad i = 1, \dots, m \\ \nabla f_0(\tilde{\mathbf{x}}) + \sum_{i=1}^m \tilde{\lambda}_i \nabla f_i(\tilde{\mathbf{x}}) + \sum_{j=1}^p \tilde{\nu}_j \nabla h_j(\tilde{\mathbf{x}}) &= \mathbf{0}. \end{aligned}$$

Complementary slackness

Vanishing Lagrangian gradient

If we have a convex program, then the KKT conditions are sufficient for zero duality gap. The motivation behind this is that they may be easier to solve for than solving the primal optimization problem. However, we cannot always count on the KKT conditions to be solvable, because they are only guaranteed if there are primal and dual solutions of zero duality



gap. But, if the primal has a Slater point<sup>2</sup>, then the KKT conditions are equivalent to the existence of a primal minimizer.

<sup>2</sup> A point  $\tilde{x}$  such that  $f_i(\tilde{x}) < 0$  for all  $i = 1, \dots, m$ .

### 3 Gradient descent

Gradient descent is an optimization algorithm that aims to find the global minimizer. We assume that  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  is convex, differentiable, and has a global minimizer  $\mathbf{x}^*$ . Then, the goal of gradient descent is, for  $\epsilon > 0$ , to find a  $\mathbf{x} \in \mathbb{R}^d$  such that

$$f(\mathbf{x}) - f(\mathbf{x}^*) \leq \epsilon.$$

It works by first choosing a  $\mathbf{x}_0 \in \mathbb{R}^d$  and then iteratively updating this value by

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \gamma \nabla f(\mathbf{x}_t),$$

for timesteps  $t = 0, 1, \dots$  and stepsize  $\gamma > 0$ .

#### 3.1 Vanilla analysis

We want to be able to bound  $f(\mathbf{x}_t) - f(\mathbf{x}^*)$ , which tells us how far off we are from the minimum at timestep  $t$ . For this, we can use the first-order characterization of convexity,

$$f(\mathbf{x}^*) \geq f(\mathbf{x}_t) + \nabla f(\mathbf{x}_t)^\top (\mathbf{x}^* - \mathbf{x}_t) \iff f(\mathbf{x}_t) - f(\mathbf{x}^*) \leq \nabla f(\mathbf{x}_t)^\top (\mathbf{x}_t - \mathbf{x}^*).$$

Using the gradient descent update rule, we get

$$\nabla f(\mathbf{x}_t) = \frac{\mathbf{x}_t - \mathbf{x}_{t+1}}{\gamma},$$

which gives the following bound

$$\begin{aligned} f(\mathbf{x}_t) - f(\mathbf{x}^*) &\leq \nabla f(\mathbf{x}_t)^\top (\mathbf{x}_t - \mathbf{x}^*) \\ &= \frac{1}{\gamma} (\mathbf{x}_t - \mathbf{x}_{t+1})^\top (\mathbf{x}_t - \mathbf{x}^*) \\ &= \frac{1}{2\gamma} \left( \|\mathbf{x}_t - \mathbf{x}_{t+1}\|^2 + \|\mathbf{x}_t - \mathbf{x}^*\|^2 - \|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2 \right) \\ &= \frac{\gamma}{2} \|\nabla f(\mathbf{x}_t)\|^2 + \frac{1}{2\gamma} \left( \|\mathbf{x}_t - \mathbf{x}^*\|^2 - \|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2 \right). \end{aligned}$$

$2\mathbf{v}^\top \mathbf{w} = \|\mathbf{v}\|^2 + \|\mathbf{w}\|^2 - \|\mathbf{v} - \mathbf{w}\|^2$  (cosine theorem).

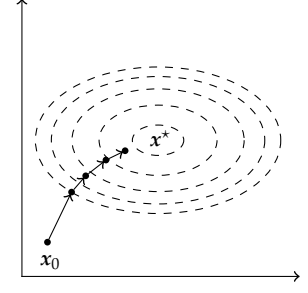
$$\mathbf{x}_t - \mathbf{x}_{t+1} = \gamma \nabla f(\mathbf{x}_t).$$

Summing this up over the first  $T$  iterations, we get an upper bound on the summed error,

$$\begin{aligned} \sum_{t=0}^{T-1} (f(\mathbf{x}_t) - f(\mathbf{x}^*)) &\leq \sum_{t=0}^{T-1} \nabla f(\mathbf{x}_t)^\top (\mathbf{x}_t - \mathbf{x}^*) \\ &= \frac{\gamma}{2} \sum_{t=0}^{T-1} \|\nabla f(\mathbf{x}_t)\|^2 + \frac{1}{2\gamma} \left( \|\mathbf{x}_0 - \mathbf{x}^*\|^2 - \|\mathbf{x}_T - \mathbf{x}^*\|^2 \right) && \text{Telescoping sum.} \\ &\leq \frac{\gamma}{2} \sum_{t=0}^{T-1} \|\nabla f(\mathbf{x}_t)\|^2 + \frac{1}{2\gamma} \|\mathbf{x}_0 - \mathbf{x}^*\|^2. && \|\mathbf{x}_T - \mathbf{x}^*\|^2 \geq 0. \end{aligned}$$

We can also use this to get a bound on the average error,

$$\frac{1}{T} \sum_{t=0}^{T-1} (f(\mathbf{x}_t) - f(\mathbf{x}^*)) \leq \frac{1}{T} \left( \frac{\gamma}{2} \sum_{t=0}^{T-1} \|\nabla f(\mathbf{x}_t)\|^2 + \frac{1}{2\gamma} \|\mathbf{x}_0 - \mathbf{x}^*\|^2 \right).$$



**Figure 12.** Gradient descent updates.

Take a small step into the direction of the negative gradient to move toward the minimum.

### 3.2 Lipschitz convex functions

**Theorem 28.** Let  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  be convex and differentiable with a global minimum  $\mathbf{x}^*$ . Furthermore, suppose that  $\|\mathbf{x}_0 - \mathbf{x}^*\| \leq R$  and  $\|\nabla f(\mathbf{x})\| \leq B$  for all  $\mathbf{x}$ . Choosing the stepsize,

$$\gamma \doteq \frac{R}{B\sqrt{T}},$$

gradient descent yields

$$\frac{1}{T} \sum_{t=0}^{T-1} (f(\mathbf{x}_t) - f(\mathbf{x}^*)) \leq \frac{RB}{\sqrt{T}}.$$

A function  $f : \mathbb{R} \rightarrow \mathbb{R}$  is Lipschitz continuous if there exists an  $M$ , such that  $\forall x, y \in \mathbb{R}$  :

$$\frac{|f(x) - f(y)|}{|x - y|} \leq M.$$

This holds if the gradient is bounded.

*Proof.* We can plug our bounds into the vanilla analysis,

$$\begin{aligned} \sum_{t=0}^{T-1} (f(\mathbf{x}_t) - f(\mathbf{x}^*)) &\leq \frac{\gamma}{2} \sum_{t=0}^{T-1} \|\nabla f(\mathbf{x}_t)\|^2 + \frac{1}{2\gamma} \|\mathbf{x}_0 - \mathbf{x}^*\|^2 \\ &\leq \frac{\gamma}{2} B^2 T + \frac{1}{2\gamma} R^2 \\ &\doteq q(\gamma). \end{aligned}$$

Now, we want to choose  $\gamma$  such that the bound  $q(\gamma)$  is minimized. We can compute the minimum by computing the derivative,

$$q'(\gamma) = \frac{1}{2} B^2 T - \frac{R^2}{2\gamma^2},$$

and solving for  $q'(\gamma) = 0$ , which yields

$$\gamma = \frac{R}{B\sqrt{T}}.$$

Then, we can compute the bound by

$$q\left(\frac{R}{B\sqrt{T}}\right) = RB\sqrt{T}.$$

Dividing by  $T$  yields the result. ■

We want to find out how many iterations we would need to ensure that the average error is bounded by  $\epsilon$ . We can use Theorem 28 to compute a lower bound on the number of iterations,

$$\frac{RB}{\sqrt{T}} \leq \epsilon \implies T \geq \frac{R^2 B^2}{\epsilon^2}.$$

So, the amount of iterations until convergence is of the order  $\mathcal{O}(1/\epsilon^2)$ . This means that we need at most  $10000 \cdot R^2 B^2$  iterations to achieve an error of 0.01.

### 3.3 Smooth functions

**Definition 29** (Smooth functions). Let  $f : \text{dom}(f) \rightarrow \mathbb{R}$  be differentiable,  $\mathcal{X} \subseteq \text{dom}(f)$  convex,  $L \in \mathbb{R}_+$ .  $f$  is smooth (with parameter  $L$ ) over  $\mathcal{X}$  if  $\forall \mathbf{x}, \mathbf{y} \in \mathcal{X}$ :

$$\underbrace{f(\mathbf{y}) \leq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x})}_{\text{first-order convexity}} + \frac{L}{2} \|\mathbf{x} - \mathbf{y}\|^2.$$

**Remark.** This definition does not require convexity.

Geometrically, this definition of smoothness means that the graph  $f$  is below a not too steep tangent paraboloid at  $(\mathbf{x}, f(\mathbf{x}))$ ; see Figure 13. Furthermore, we can easily check if  $f$  is smooth with parameter  $L$  if the following function is convex,

$$g(\mathbf{x}) \doteq \frac{L}{2} \mathbf{x}^\top \mathbf{x} - f(\mathbf{x}),$$

over  $\text{dom}(g) \doteq \text{dom}(f)$ .

**Observation** (Operations that preserve smoothness). Let  $f_1, \dots, f_m$  be smooth with parameters  $L_1, \dots, L_m$  and let  $\lambda_1, \dots, \lambda_m \in \mathbb{R}_+$ . Then, the function  $f \doteq \sum_{i=1}^m \lambda_i f_i$  is smooth with parameter  $\sum_{i=1}^m \lambda_i L_i$ .

Let  $f$  be smooth with parameter  $L$  and let  $g(\mathbf{x}) = \mathbf{A}\mathbf{x} + \mathbf{b}$ . Then, the function  $f \circ g$  is smooth with parameter  $L\|\mathbf{A}\|^2$ , where  $\|\mathbf{A}\|$  is the spectral norm of  $\mathbf{A}$ .

**Lemma 30.** Let  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  be convex and differentiable, then  $f$  is smooth with parameter  $L$  if

$$\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\| \leq L\|\mathbf{x} - \mathbf{y}\|.$$

I.e.,  $\nabla f$  is Lipschitz continuous.

**Lemma 31** (Sufficient decrease). Let  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  be differentiable and smooth with parameter  $L$ . With stepsize

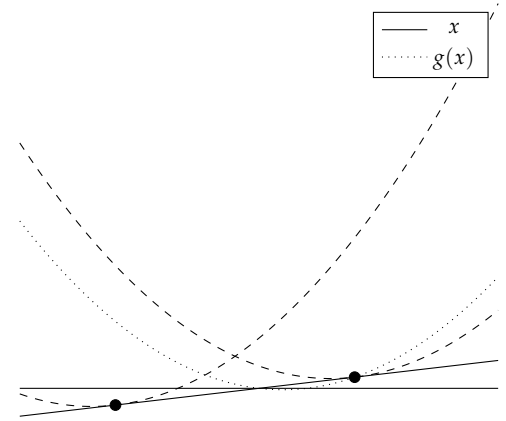
$$\gamma \doteq \frac{1}{L},$$

gradient descent satisfies

$$f(\mathbf{x}_{t+1}) \leq f(\mathbf{x}_t) - \frac{1}{2L} \|\nabla f(\mathbf{x}_t)\|^2.$$

Sufficient decrease implies that, every step, we get closer to the minimum.

“Not too curved.”



**Figure 13.** Plot of  $f(x) = x$  with the tangent paraboloids at  $x = -8, 5$ , showing smoothness with parameter  $L = 1$ . Furthermore, it shows that  $g(x) = \frac{1}{2}x^2 - f(x)$  is convex.

*Proof.*

$$\begin{aligned}
 f(\mathbf{x}_{t+1}) &\leq f(\mathbf{x}_t) + \nabla f(\mathbf{x}_t)^\top (\mathbf{x}_{t+1} - \mathbf{x}_t) + \frac{L}{2} \|\mathbf{x}_t - \mathbf{x}_{t+1}\|^2 \\
 &= f(\mathbf{x}_t) + \nabla f(\mathbf{x}_t)^\top \left( -\frac{\nabla f(\mathbf{x}_t)}{L} \right) + \frac{L}{2} \left\| \frac{\nabla f(\mathbf{x}_t)}{L} \right\|^2 \\
 &= f(\mathbf{x}_t) - \frac{1}{L} \|\nabla f(\mathbf{x}_t)\|^2 + \frac{1}{2L} \|\nabla f(\mathbf{x}_t)\|^2 \\
 &= f(\mathbf{x}_t) - \frac{1}{2L} \|\nabla f(\mathbf{x}_t)\|^2.
 \end{aligned}$$

Smoothness.

$$\text{GD: } \mathbf{x}_{t+1} = \mathbf{x}_t - \frac{1}{L} \nabla f(\mathbf{x}_t).$$

**Theorem 32.** Let  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  be convex and differentiable with a global minimum  $\mathbf{x}^*$ . Furthermore, suppose that  $f$  is smooth with parameter  $L$ . Choosing stepsize

$$\gamma \doteq \frac{1}{L},$$

gradient descent yields

$$f(\mathbf{x}_T) - f(\mathbf{x}^*) \leq \frac{L}{2T} \|\mathbf{x}_0 - \mathbf{x}^*\|^2, \quad T > 0.$$

*Proof.* Due to sufficient decrease, we can bound the sum of squared gradients (which will be useful when bounding the vanilla analysis),

$$\begin{aligned}
 \frac{1}{2L} \sum_{t=0}^{T-1} \|\nabla f(\mathbf{x}_t)\|^2 &\leq \sum_{t=0}^{T-1} f(\mathbf{x}_t) - f(\mathbf{x}_{t+1}) \\
 &= f(\mathbf{x}_0) - f(\mathbf{x}_T).
 \end{aligned}$$

$$\text{SD: } f(\mathbf{x}_{t+1}) \leq f(\mathbf{x}_t) - \frac{1}{2L} \|\nabla f(\mathbf{x}_t)\|^2.$$

Telescoping sum.

Using the vanilla analysis, we can derive a bound on the average error,

$$\begin{aligned}
 \sum_{t=0}^{T-1} (f(\mathbf{x}_t) - f(\mathbf{x}^*)) &\leq \frac{1}{2L} \sum_{t=0}^{T-1} \|\nabla f(\mathbf{x}_t)\|^2 + \frac{L}{2} \|\mathbf{x}_0 - \mathbf{x}^*\|^2 \\
 &\leq f(\mathbf{x}_0) - f(\mathbf{x}_T) + \frac{L}{2} \|\mathbf{x}_0 - \mathbf{x}^*\|^2 \\
 \sum_{t=1}^T (f(\mathbf{x}_t) - f(\mathbf{x}^*)) &\leq \frac{L}{2} \|\mathbf{x}_0 - \mathbf{x}^*\|^2 \\
 \frac{1}{T} \sum_{t=1}^T (f(\mathbf{x}_t) - f(\mathbf{x}^*)) &\leq \frac{L}{2T} \|\mathbf{x}_0 - \mathbf{x}^*\|^2.
 \end{aligned}$$

$$\gamma \doteq \frac{1}{L}.$$

From sufficient decrease, we know that the last iterate must be the best, thus

$$f(\mathbf{x}_T) - f(\mathbf{x}^*) \leq \frac{1}{T} \sum_{t=1}^T (f(\mathbf{x}_t) - f(\mathbf{x}^*)) \leq \frac{L}{2T} \|\mathbf{x}_0 - \mathbf{x}^*\|^2.$$

Using this result, we can compute a bound on  $T$  to achieve an error smaller than  $\epsilon$ ,

$$f(x_T) - f(x^*) \leq \frac{L}{2T} \|x_0 - x^*\|^2 \leq \frac{R^2 L}{2T} \leq \epsilon.$$

The error then becomes,

$$T \geq \frac{R^2 L}{2\epsilon},$$

which is on the order of  $\mathcal{O}(1/\epsilon)$ . This means that we need at most  $50 \cdot R^2 L$  iterations for an error of 0.01 (as opposed to  $10000 \cdot R^2 B^2$  in the Lipschitz case).

### 3.4 Strongly convex functions

**Definition 33** (Strong convexity). Let  $f : \text{dom}(f) \rightarrow \mathbb{R}$  be a convex and differentiable function,  $\mathcal{X} \in \text{dom}(f)$  convex and  $\mu \in \mathbb{R}_+, \mu > 0$ .  $f$  is called *strongly convex* with parameter  $\mu$  over  $\mathcal{X}$  if  $\forall x, y \in \mathcal{X}$ :

$$f(y) \geq f(x) + \nabla f(x)^\top (y - x) + \frac{\mu}{2} \|x - y\|^2.$$

Geometrically, this means that, for any  $x$ , the graph of  $f$  is above a not too flat tangent paraboloid at  $(x, f(x))$ ; see Figure 14. Furthermore, we can easily check if  $f$  is strongly convex if and only if the following function is convex,

$$g(x) \doteq f(x) - \frac{\mu}{2} x^\top x.$$

By assuming that a function  $f$  is smooth and strongly convex, we can use a stronger lower bound to derive a bound on the error from the vanilla analysis, TODO

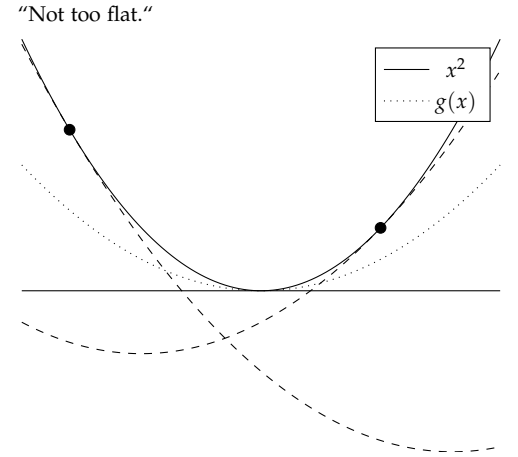
**Theorem 34.** Let  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  be convex and differentiable with a global minimum  $x^*$ . Furthermore, suppose that  $f$  is smooth with parameter  $L$  and strongly convex with parameter  $\mu > 0$ . Choosing  $\gamma \doteq \frac{1}{L}$ , gradient descent with arbitrary  $x_0$  satisfies the following two properties,

1. Squared distances to  $x^*$  are geometrically decreasing,

$$\|x_{t+1} - x^*\|^2 \leq \left(1 - \frac{\mu}{L}\right) \|x_t - x^*\|^2;$$

2. The absolute error after  $T$  iterations is exponentially small in  $T$ ,

$$f(x_T) - f(x^*) \leq \frac{L}{2} \left(1 - \frac{\mu}{L}\right)^T \|x_0 - x^*\|^2.$$



**Figure 14.** Plot of  $f(x) = x^2$  with the tangent paraboloids at  $x = -6, 4$ , showing strong convexity with parameter  $\mu = 1$ . Furthermore, it shows  $g(x) = f(x) - \frac{1}{2}x^2$ , which is convex.

*Proof of 1.* Using the vanilla analysis, we know the following,

$$\begin{aligned} f(\mathbf{x}_t) - f(\mathbf{x}^*) &\leq \nabla f(\mathbf{x}_t)^\top (\mathbf{x}_t - \mathbf{x}^*) \\ &= \frac{\gamma}{2} \|\nabla f(\mathbf{x}_t)\|^2 + \frac{1}{2\gamma} (\|\mathbf{x}_t - \mathbf{x}^*\|^2 - \|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2). \end{aligned}$$

But, using strong convexity, we have a stronger lower bound on the left hand side,

$$\begin{aligned} f(\mathbf{x}_t) - f(\mathbf{x}^*) &\leq \nabla f(\mathbf{x}_t)^\top (\mathbf{x}_t - \mathbf{x}^*) - \frac{\mu}{2} \|\mathbf{x}_t - \mathbf{x}^*\|^2 && \text{Strong convexity.} \\ &= \frac{\gamma}{2} \|\nabla f(\mathbf{x}_t)\|^2 + \frac{1}{2\gamma} (\|\mathbf{x}_t - \mathbf{x}^*\|^2 - \|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2) - \frac{\mu}{2} \|\mathbf{x}_t - \mathbf{x}^*\|^2. \end{aligned}$$

This can be rewritten to the following bound,

$$\|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2 \leq \underbrace{2\gamma(f(\mathbf{x}^*) - f(\mathbf{x}_t)) + \gamma^2 \|\nabla f(\mathbf{x}_t)\|^2}_{\text{"noise"}} + (1 - \mu\gamma) \|\mathbf{x}_t - \mathbf{x}^*\|^2.$$

Now we need to show that the noise is non-positive,

$$\begin{aligned} 2\gamma(f(\mathbf{x}^*) - f(\mathbf{x}_t)) + \gamma^2 \|\nabla f(\mathbf{x}_t)\|^2 &= \frac{2}{L} (f(\mathbf{x}^*) - f(\mathbf{x}_t)) + \frac{1}{L^2} \|\nabla f(\mathbf{x}_t)\|^2 && \gamma \doteq \frac{1}{L}. \\ &\leq \frac{2}{L} (f(\mathbf{x}_{t+1}) - f(\mathbf{x}_t)) + \frac{1}{L^2} \|\nabla f(\mathbf{x}_t)\|^2 && \text{Sufficient decrease.} \\ &\leq -\frac{1}{L^2} \|\nabla f(\mathbf{x}_t)\|^2 + \frac{1}{L^2} \|\nabla f(\mathbf{x}_t)\|^2 && \text{SD: } f(\mathbf{x}_{t+1}) \leq f(\mathbf{x}_t) - \frac{1}{2L} \|\nabla f(\mathbf{x}_t)\|^2. \\ &= 0. \end{aligned}$$

Hence, the noise is non-positive, and we get the following,

$$\|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2 \leq \left(1 - \frac{\mu}{L}\right) \|\mathbf{x}_t - \mathbf{x}^*\|^2.$$

■

*Proof of 2.* From (1), we know the following,

$$\|\mathbf{x}_T - \mathbf{x}^*\|^2 \leq \left(1 - \frac{\mu}{L}\right)^T \|\mathbf{x}_0 - \mathbf{x}^*\|^2.$$

Using smoothness, we can derive a bound on the final iterate error,

$$\begin{aligned} f(\mathbf{x}_T) - f(\mathbf{x}^*) &\leq \nabla f(\mathbf{x}^*)^\top (\mathbf{x}_T - \mathbf{x}^*) + \frac{L}{2} \|\mathbf{x}_T - \mathbf{x}^*\|^2 \\ &= \frac{L}{2} \|\mathbf{x}_T - \mathbf{x}^*\|^2 && \nabla f(\mathbf{x}^*) = 0, \text{ because it is a stationary point.} \\ &\leq \frac{L}{2} \left(1 - \frac{\mu}{L}\right)^T \|\mathbf{x}_0 - \mathbf{x}^*\|^2. && \text{Using (1).} \end{aligned}$$

■

From there, we can derive a lower bound on the number of iterations  $T$  to get an error of at most  $\epsilon$ ,

$$T \geq \frac{L}{\mu} \log\left(\frac{R^2 L}{2\epsilon}\right),$$

This means that we need  $\frac{L}{\mu} \log(50 \cdot R^2 L)$  iterations for an error of at most 0.01, as opposed to  $50 \cdot R^2 L$  in the smooth case. This bound only depends linearly on  $\frac{L}{\mu}$ , which might be very high.

## *References*