

# The Five Miracles of Mirror Descent

**Hadar Tal**

hadar.tal@mail.huji.ac.il

This paper is a summary of the educational materials and lectures from Professor Sebastian Bubeck, enhanced by Claire Boyer's comprehensive notes, and structured according to Tomer Koren's course on Optimization for Computer Science.

Winter 2024



# Contents

<b>0</b>	<b>Mathematical Background</b>	<b>1</b>
0.1	Multivariable Calculus . . . . .	1
0.2	Taylor series . . . . .	3
0.3	Important subsets of $\mathbb{R}^n$ . . . . .	4
0.4	Convexity . . . . .	6
0.4.1	Definitions and Fundamental Theorems . . . . .	6
0.4.2	Inequalities and Characterizations . . . . .	6
0.4.3	Optimization and Projection . . . . .	7
0.5	Properties of Convex Functions . . . . .	8
0.6	Important Inequalities . . . . .	10
<b>1</b>	<b>The First Miracle: Robustness</b>	<b>11</b>
1.1	Gradient Descent . . . . .	11
1.1.1	Analysis of the Gradient Descent Algorithm . . . . .	11
<b>2</b>	<b>The Second Miracle: Potential Based</b>	<b>13</b>
2.1	Experts Problem . . . . .	13
	Approach 1: Gradient Descent . . . . .	13
<b>3</b>	<b>The Third Miracle:</b>	<b>15</b>
<b>4</b>	<b>The Fourth Miracle:</b>	<b>17</b>
<b>5</b>	<b>The Fifth Miracle:</b>	<b>19</b>



# Chapter 0

## Mathematical Background

### 0.1 Multivariable Calculus

**Definition 0.1.1.** *Differentiability, single variable*

Let  $f : (a, b) \rightarrow \mathbb{R}$  be a function. We say that  $f$  is differentiable at  $x_0 \in (a, b)$  if

$$\lim_{h \rightarrow 0} \frac{f(x_0 + h) - f(x_0)}{h} \quad (1)$$

exists. If  $f$  is differentiable at  $x_0$ , then  $f'(x_0)$  is the derivative of  $f$  at  $x_0$ .

**Definition 0.1.2.** *Differentiability, single variable (alternative)*

Let  $f : (a, b) \rightarrow \mathbb{R}$  be a function. We say that  $f$  is differentiable at  $x_0 \in (a, b)$  if there exists a number  $m$  such that:

$$f(x_0 + h) = f(x_0) + m \cdot h + E(h) \text{ where } \lim_{h \rightarrow 0} \frac{E(h)}{h} = 0 \quad (2)$$

If  $f$  is differentiable at  $x_0$ , then  $f'(x_0) = m$  is the derivative of  $f$  at  $x_0$ .

**Definition 0.1.3.** *Differentiability, multivariable*

Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  be a function. We say that  $f$  is differentiable at  $x_0$  if there exists a vector  $m \in \mathbb{R}^n$  such that:

$$\lim_{h \rightarrow 0} \frac{f(x_0 + h) - f(x_0) - m \cdot h}{||h||} = 0 \quad (3)$$

If  $f$  is differentiable at  $x_0$ , then  $m$  is the gradient of  $f$  at  $x_0$ , denoted  $\nabla f(x_0)$ .

Suppose the  $S \subseteq \mathbb{R}^n$  and  $f : S \rightarrow \mathbb{R}$  is a function.

**Definition 0.1.4.** *Limit, multivariate function*

We say that the limit of  $f$  at  $x_0$  is  $L$  if for all  $\epsilon > 0$ , there exists  $\delta > 0$  such that for all  $x$  such that  $||x - x_0|| < \delta$ , we have  $|f(x) - L| < \epsilon$ .

**Definition 0.1.5.** *Differentiability, multivariable (alternative)*

We say that  $f$  is differentiable at  $x_0$  if there exists a vector  $m \in \mathbb{R}^n$  such that:

$$f(x_0 + h) = f(x_0) + m^T \cdot h + E(h) \text{ where } \lim_{h \rightarrow 0} \frac{E(h)}{||h||} = 0 \quad (4)$$

If  $f$  is differentiable at  $x_0$ , then  $m$  is the gradient of  $f$  at  $x_0$ , denoted  $\nabla f(x_0)$ .

**Definition 0.1.6.** *Partial Derivative*

The partial derivative of  $f$  with respect to the  $i$ -th variable at  $x$  is:

$$\frac{\partial f}{\partial x_i}(x) = \lim_{h \rightarrow 0} \frac{f(x + h \cdot e_i) - f(x)}{h} \quad (5)$$

where  $e_i$  is the  $i$ -th standard basis vector.

**Theorem 0.1.1.** *(Differentiability vs. Partial Derivatives)*

If  $f$  is differentiable at  $x$ , then all partial derivatives of  $f$  exist at  $x$  and:

$$\nabla f(x) = \left( \frac{\partial f}{\partial x_1}(x), \dots, \frac{\partial f}{\partial x_n}(x) \right) \quad (6)$$

- If any partial derivative of  $f$  does not exist at  $x$ , then  $f$  is not differentiable at  $x$ .
- If all partial derivatives of  $f$  exist at  $x$ , then  $f$  may still not be differentiable at  $x$  and the vector  $m = \nabla f(x)$  is the only possible vector that satisfies the definition of differentiability.

**Definition 0.1.7.** *Continuously Differentiable*

We say that  $f$  is continuously differentiable or of class  $C^1$  if all partial derivatives of  $f$  exist and are continuous at every point in  $S$ .

**Theorem 0.1.2.** If  $f$  is continuously differentiable, then  $f$  is differentiable.

**Definition 0.1.8.** *The directional derivative*

For a given  $x \in S$  and a unit vector  $u \in \mathbb{R}^n$ , the directional derivative of  $f$  at  $x$  in the direction of  $u$  is:

$$\partial_u f(x) = \lim_{h \rightarrow 0} \frac{f(x + h \cdot u) - f(x)}{h} \quad (7)$$

Equivalently,  $\partial_u f(x) = g'(0)$  where  $g(h) = f(x + h \cdot u)$ .

**Theorem 0.1.3.** If  $f$  is differentiable at  $x$ , then for all  $u \in \mathbb{R}^n$ , the directional derivative of  $f$  at  $x$  in the direction of  $u$  exists and is given by:

$$\partial_u f(x) = \nabla f(x) \cdot u \quad (8)$$

**Theorem 0.1.4.** *Fermat's Theorem*

If  $f$  is differentiable at  $x$  and  $x$  is a local minimum of  $f$ , then  $\nabla f(x) = 0$ .

**Theorem 0.1.5.** Suppose that  $f : S \rightarrow \mathbb{R}$  is differentiable at  $x$ . Then  $\nabla f(x)$  is orthogonal to the level set of  $f$  that passes through  $x$ .

**Theorem 0.1.6.** *The mean value theorem*

If  $f : S \rightarrow \mathbb{R}$  is differentiable on the open interval between  $a$  and  $b$ , then there exists  $c \in [a, b]$  such that:

$$f(b) - f(a) = \nabla f(c) \cdot (b - a) \quad (9)$$

where  $[a, b] = a + t(b - a) | t \in [0, 1]$ .

**Definition 0.1.9.** *Second-order partial derivatives*

Suppose that  $f$  is a  $C^1$  function. If the partial derivatives of  $f$  are differentiable, then the second-order partial derivatives of  $f$  are:

$$\frac{\partial^2 f}{\partial x_i \partial x_j} = \frac{\partial}{\partial x_i} \left( \frac{\partial f}{\partial x_j} \right) \quad (10)$$

Equivalently,  $\frac{\partial^2 f}{\partial i \partial j} = \partial_j \partial_i f$ . If  $i = j$  we denote  $\frac{\partial^2 f}{\partial x_i^2}$  or  $(\partial_i^2 f)$

**Definition 0.1.10.** *The  $C^2$  class*

We say that  $f$  is of class  $C^2$  if all second-order partial derivatives of  $f$  exist and are continuous.

**Theorem 0.1.7.** *Clairaut's Theorem*

If  $f$  is of class  $C^2$ , then  $\frac{\partial^2 f}{\partial x_i \partial x_j} = \frac{\partial^2 f}{\partial x_j \partial x_i}$ .

**Definition 0.1.11.** *Hessian Matrix*

The Hessian matrix of  $f$  at  $x$  is the matrix of second-order partial derivatives of  $f$  at  $x$ :

$$\nabla^2 f(x) = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} & \cdots & \frac{\partial^2 f}{\partial x_2 \partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \frac{\partial^2 f}{\partial x_n \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_n^2} \end{bmatrix} \quad (11)$$

**Corollary.** *The interpretation of the Hessian matrix*

Let  $u \in \mathbb{R}^n$  be a unit vector. then

$$\partial_{uu}^2 f(x) = \sum_{i,j=1}^n \partial_{ij}^2 f(x) u_i u_j = u^T \nabla^2 f(x) u \quad (12)$$

## 0.2 Taylor series

**Definition 0.2.1.** *Taylor Series*

Let  $f : \mathbb{R} \rightarrow \mathbb{R}$  be a function that is  $k$  times differentiable at  $x_0$ . Then the Taylor series of  $f$  at  $x_0$  is given by:

$$f(x) = f(x_0) + f'(x_0)(x - x_0) + \frac{f''(x_0)}{2!}(x - x_0)^2 + \cdots + \frac{f^{(k)}(x_0)}{k!}(x - x_0)^k + R_k(x) \quad (13)$$

where  $R_k(x) = \frac{f^{(k+1)}(c)}{(k+1)!}(x - x_0)^{k+1}$  for some  $c$  between  $x$  and  $x_0$ .

**Definition 0.2.2.** *Taylor Series for Multivariable Functions ( $k=2$ )*

Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  be a function that is  $C^2$  at  $x_0$ . Then for any  $h$  such that  $x_0 + h \in S$ , there exists  $\theta \in [0, 1]$  such that:

$$f(x_0 + h) = f(x_0) + \nabla f(x_0) \cdot h + \frac{1}{2} h^T \nabla^2 f(x_0 + \theta h) h \quad (14)$$

### 0.3 Important subsets of $\mathbb{R}^n$

**Definition 0.3.1.** *Open set*

A set  $S \subseteq \mathbb{R}^n$  is open if for all  $x \in S$ , there exists  $\epsilon > 0$  such that  $B(x, \epsilon) \subseteq S$ .

**Definition 0.3.2.** *Closed set*

A set  $S \subseteq \mathbb{R}^n$  is closed if its complement is open.

**Definition 0.3.3.** *Interior point*

A point  $x \in S$  is an interior point of  $S$  if there exists  $\epsilon > 0$  such that  $B(x, \epsilon) \subseteq S$ .

**Corollary 0.3.1.** *Open set characterization*

A set  $S \subseteq \mathbb{R}^n$  is open if and only if every point in  $S$  is an interior point of  $S$ .

**Definition 0.3.4.** *Boundary point*

A point  $x \in S$  is a boundary point of  $S$  if for all  $\epsilon > 0$ ,  $B(x, \epsilon) \cap S \neq \emptyset$  and  $B(x, \epsilon) \cap S^c \neq \emptyset$ .

**Definition 0.3.5.** *Half-space*

A half-space in  $\mathbb{R}^n$  is a set of the form  $\{x \in \mathbb{R}^n : a^T x \leq b\}$  for some  $a \in \mathbb{R}^n$  and  $b \in \mathbb{R}$ .

**Definition 0.3.6.** *Hyperplane*

A hyperplane in  $\mathbb{R}^n$  is a set of the form  $\{x \in \mathbb{R}^n : a^T x = b\}$  for some  $a \in \mathbb{R}^n$  and  $b \in \mathbb{R}$ .

**Definition 0.3.7.** *Polyhedron (Polyhedra)*

A polyhedron in  $\mathbb{R}^n$  is a set of the form  $\{x \in \mathbb{R}^n : Ax \leq b\}$  for some  $A \in \mathbb{R}^{m \times n}$  and  $b \in \mathbb{R}^m$ . Equivalently, a polyhedron is the intersection of finitely many half-spaces.

**Definition 0.3.8.** *Polytope*

A polytope in  $\mathbb{R}^n$  is a bounded polyhedron - i.e., there exists  $r > 0$  such that  $\forall x \in \{x \in \mathbb{R}^n : Ax \leq b\} \implies \|x\| \leq r$ . Equivalently, a polytope is the convex hull of finitely many points.

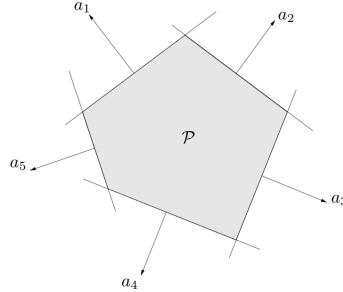


Figure 1: Polytope

**Definition 0.3.9.** *Convex set*

A set  $S \subseteq \mathbb{R}^n$  is convex if for all  $x, y \in S$  and  $\lambda \in [0, 1]$ , we have  $\lambda x + (1 - \lambda)y \in S$ .

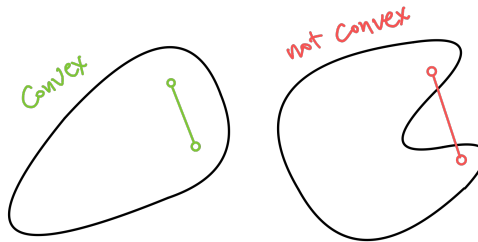


Figure 2: Convex set



**Definition 0.3.10.** *Convex hull*

The convex hull of a set  $S \subseteq \mathbb{R}^n$  is the smallest convex set that contains  $S$ .

**Definition 0.3.11.** *Conic combination*

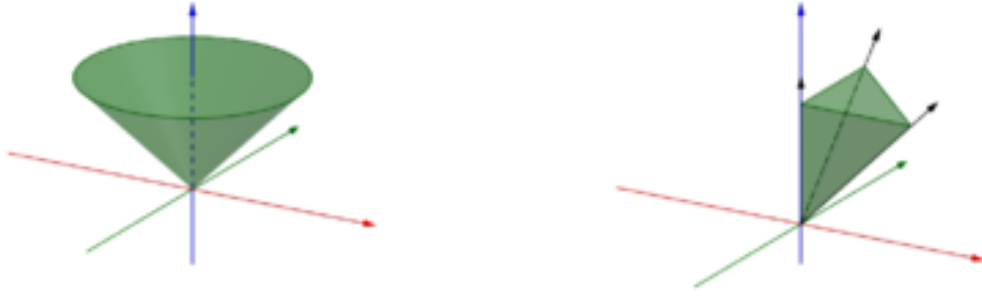
A point  $x \in \mathbb{R}^n$  is a conic combination of  $y_1, \dots, y_k \in \mathbb{R}^n$  if there exist  $\lambda_1, \dots, \lambda_k \geq 0$  such that  $x = \sum_{i=1}^k \lambda_i y_i$ .

**Definition 0.3.12.** *Conic hull*

The conic hull of a finite set  $S \subseteq \mathbb{R}^n$  is the set of all conic combinations of points in  $S$ .

**Definition 0.3.13.** *Convex cone*

A set  $S \subseteq \mathbb{R}^n$  is a convex cone if for all  $x \in S$  and  $\lambda \geq 0$ , we have  $\lambda x \in S$ .



(a) Convex cone that is not a conic hull of finitely many generators. (b) Convex cone generated by the conic combination of three black vectors (conic hull).

**Definition 0.3.14.** *Normal cone*

The normal cone to a set  $S$  at a point  $x$  is defined as

$$N_S(x) = \{v \in \mathbb{R}^n : \langle v, y - x \rangle \leq 0 \text{ for all } y \in S\} \quad (15)$$

**Definition 0.3.15.** *Tangent cone*

The tangent cone to a set  $S$  at a point  $x$  is defined as

$$T_S(x) = \{v \in \mathbb{R}^n : \lim_{t \rightarrow 0^+} \frac{x + tv - x}{t} \in S\} \quad (16)$$

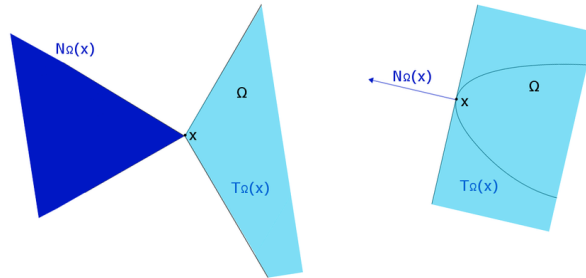


Figure 4: Normal and tangent cones

**Theorem 0.3.1.** *Normal cone of polyhedron*

The normal cone to a polyhedron  $S = \{x \in \mathbb{R}^n : \forall j \in [m] \quad a_j \cdot x \leq b_j\}$  at a point  $x$  is given by

$$N_S(x) = \left\{ \sum_j \lambda_j a_j : \lambda_j \geq 0 \text{ and } a_j \cdot x = b_j \right\} \quad (17)$$

## 0.4 Convexity

### 0.4.1 Definitions and Fundamental Theorems

**Definition 0.4.1.** (*Convex function*): A function  $f : S \rightarrow \mathbb{R}$  defined on a convex set  $S$  is convex if, for all  $x, y \in S$  and  $\lambda \in [0, 1]$ ,

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

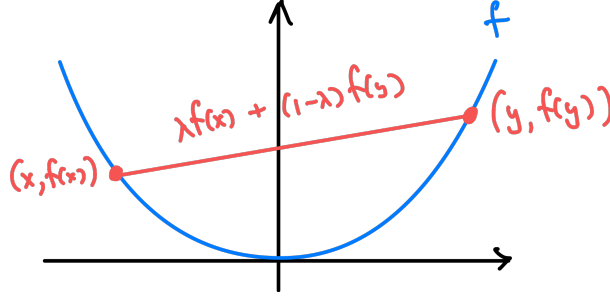


Figure 5: Convex function

**Theorem 0.4.1.** (*Characterization via epigraph*): A function  $f : S \rightarrow \mathbb{R}$  is convex if and only if its epigraph  $\{(x, t) \in S \times \mathbb{R} : f(x) \leq t\}$  is a convex set.

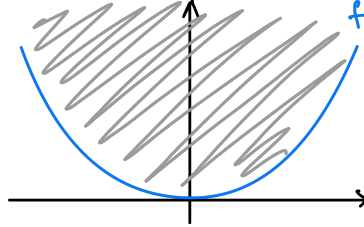


Figure 6: Epigraph of a convex function

**claim 0.4.1.** (*Convexity of sublevel sets*): If  $f : S \rightarrow \mathbb{R}$  is convex, then the sublevel set  $S_t = \{x \in S : f(x) \leq t\}$  is convex for any  $t \in \mathbb{R}$ .

### 0.4.2 Inequalities and Characterizations

**Theorem 0.4.2.** (*Jensen's inequality*): If  $f$  is a convex function, then for any  $x_1, x_2, \dots, x_n \in S$  and any non-negative weights  $\alpha_i$  such that  $\sum_{i=1}^n \alpha_i = 1$ ,

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

**Theorem 0.4.3.** (*First-order characterization, aka "the gradient inequality"*): If  $f$  is a differentiable convex function on an open set  $S$ , then for all  $x, y \in S$ ,

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

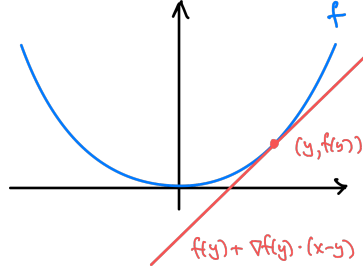


Figure 7: First-order characterization of convexity

**Definition 0.4.2.** *Bergman divergence (distance)*

The Bergman divergence between two points  $x, y \in \mathbb{R}^n$  is defined as

$$D_f(x, y) = f(x) - f(y) - \nabla f(y)^\top (x - y) \quad (18)$$

**Theorem 0.4.4.** *(Jensen's inequality, generalized for expectation): If  $f$  is a convex function and  $X$  is a random variable over  $S$ , then*

$$f(\mathbb{E}[X]) \leq \mathbb{E}[f(X)].$$

**Theorem 0.4.5.** *(Second-order characterization of convexity): A twice differentiable function  $f$  is convex on an open set  $S$  if and only if the Hessian matrix of  $f$  is positive semidefinite at every point in  $S$ .*

### 0.4.3 Optimization and Projection

**Definition 0.4.3.** *(Convex optimization): The problem of minimizing a convex function over a convex set.*

**Theorem 0.4.6.** *(Optimality conditions, unconstrained): If  $f$  is convex and differentiable,  $x^*$  is a local minimum of  $f \Leftrightarrow x^*$  is a global minimum of  $f \Leftrightarrow \nabla f(x^*) = 0$ .*

**Theorem 0.4.7.** *(Optimality conditions, constrained): If  $f$  is differentiable and  $C$  is a convex set,  $x^*$  is a local minimum of  $f$  on  $C$  if and only if  $\langle \nabla f(x^*), x - x^* \rangle \geq 0$  for all  $x \in C$ .*

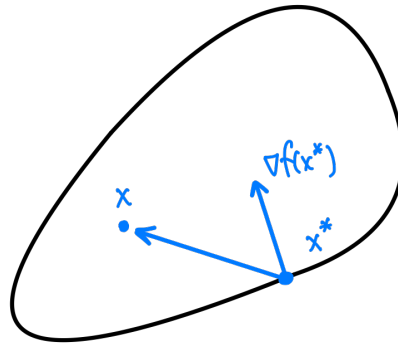


Figure 8: Optimality conditions, constrained

**Corollary 0.4.1.** *Optimality conditions, constrained (alternative)*

If  $f$  is differentiable and  $C$  is a convex set, then  $x^*$  is a local minimum of  $f$  on  $C$  if and only if  $-\nabla f(x^*) \in N_C(x^*)$ .

**Definition 0.4.4.** (Projection): The projection of a point  $x$  onto a convex set  $S$  is defined as  $\Pi_S(x) = \arg \min_{y \in S} \|y - x\|$ .

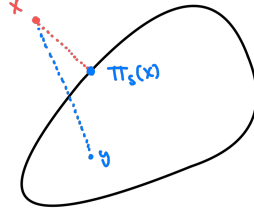


Figure 9: Projection

**Theorem 0.4.8.** *Generalized cosine theorem*

Let  $S \subseteq \mathbb{R}^d$  be convex and  $x \in \mathbb{R}^d$ . Then the projection  $\Pi_S[x]$  is unique and satisfies:

$$\|x - \Pi_S[x]\|^2 + \|\Pi_S[x] - y\|^2 \leq \|x - y\|^2, \quad \forall y \in S. \quad (19)$$

In particular:

$$\|\Pi_S[x] - y\| \leq \|x - y\|, \quad \forall y \in S. \quad (20)$$

## 0.5 Properties of Convex Functions

**Definition 0.5.1.** *L - Lipschitz continuous*

A function  $f : S \rightarrow \mathbb{R}$  is L-Lipschitz continuous if for all  $x, y \in S$ ,

$$|f(x) - f(y)| \leq L\|x - y\| \quad (21)$$

**Theorem 0.5.1.** *Convexity and Lipschitz continuity*

If  $f$  is convex, differentiable and L-Lipschitz continuous, then  $\|\nabla f(x)\| \leq L$  for all  $x \in S$ .

**Definition 0.5.2.** *Smooth function*

A differentiable function  $f$  is  $\beta$ -smooth over  $S \subseteq \text{dom} f$  if for all  $x, y \in S$ :

$$-\frac{\beta}{2}\|y - x\|^2 \leq f(y) - f(x) - \nabla f(x) \cdot (y - x) \leq \frac{\beta}{2}\|y - x\|^2.$$

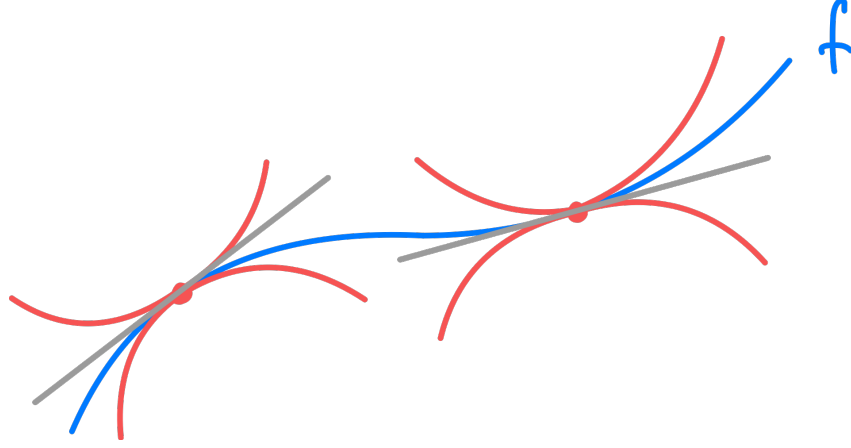


Figure 10: Smooth function

**Theorem 0.5.2.** *Lipschitz gradient interpretation*

Let  $f$  be differentiable and let  $S \subseteq \text{dom}f$  be convex and closed. Suppose that

$$\|\nabla f(x) - \nabla f(y)\| \leq \beta \|x - y\|, \quad \forall x, y \in S.$$

Then  $f$  is  $\beta$ -smooth over  $S$ .

**Theorem 0.5.3.** *Second-order characterization of smoothness*

Let  $f$  be  $C^2$  and let  $S \subseteq \text{dom}f$  be convex and closed. Then  $f$  is  $\beta$ -smooth over  $S$  if and only if

$$-\beta I \preceq \nabla^2 f(x) \preceq \beta I, \quad \forall x \in S.$$

**Lemma 0.5.1.** *The Descent Lemma*

Let  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  be  $\beta$ -smooth, and let  $x \in \mathbb{R}^d$ .

- For  $\eta \leq \frac{1}{\beta}$ ,  $x^+ = x - \eta \nabla f(x)$ , we have

$$f(x^+) - f(x) \leq -\frac{\eta}{2} \|\nabla f(x)\|^2.$$

- For  $x^* \in \arg \min_x f(x)$ , we have

$$\frac{1}{2\beta} \|\nabla f(x)\|^2 \leq f(x) - f(x^*).$$

**Basic Facts:**

- An affine function  $f : \mathbb{R}^d \rightarrow \mathbb{R}$ ,  $f(x) = a^\top x + b$ , is 0-smooth.
- A quadratic function  $f : \mathbb{R}^d \rightarrow \mathbb{R}$ ,  $f(x) = \frac{1}{2}x^\top A x + b^\top x + c$ , is  $\lambda_{\max}(A)$ -smooth.
- A linear combination of smooth functions is smooth with an appropriate parameter.
- A convex combination of  $\beta$ -smooth functions is  $\beta$ -smooth.

**Definition 0.5.3.** *Strong convexity*

A function  $f$  is  $\alpha$ -strongly convex (for  $\alpha \geq 0$ ) over a convex and closed set  $S \subseteq \text{dom}f$  if for any  $x \in S$ , there exists  $g_x \in \partial f(x)$  such that:

$$\forall y \in S, \quad f(y) \geq f(x) + g_x \cdot (y - x) + \frac{\alpha}{2} \|y - x\|^2.$$

In particular, a differentiable  $f$  is  $\alpha$ -strongly convex over  $S$  if for any  $x \in S$ ,

$$\forall y \in S, \quad f(y) \geq f(x) + \nabla f(x) \cdot (y - x) + \frac{\alpha}{2} \|y - x\|^2.$$

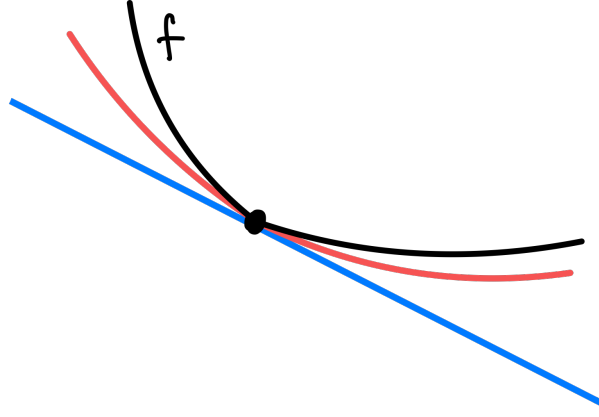


Figure 11: Strongly convex function

**Theorem 0.5.4.** *Strong convexity, second-order characterization*

Let  $f$  be  $C^2$  and let  $S \subseteq \text{dom} f$  be convex and closed. Then  $f$  is  $\alpha$ -strongly convex over  $S$  if and only if

$$\forall x \in S, \quad \nabla^2 f(x) \succeq \alpha I.$$

**Theorem 0.5.5.** *Usage of strong convexity*

If a differentiable  $f$  is  $\alpha$ -strongly convex over a convex and closed  $S \subseteq \text{dom} f$  with a minimum at  $x^* \in S$ , then

$$\forall x \in S, \quad \frac{\alpha}{2} \|x - x^*\|^2 \leq f(x) - f(x^*) \leq \frac{1}{2\alpha} \|\nabla f(x)\|^2.$$

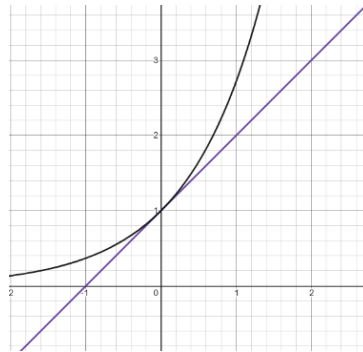
In particular, the minimum of a strongly convex function is unique.

## 0.6 Important Inequalities

**Theorem 0.6.1.**  $1 + x \leq e^x$ 

For all  $x \in \mathbb{R}$ , we have  $1 + x \leq e^x$ .

*Proof.* Let  $f(x) = e^x - 1 - x$ . Then  $f'(x) = e^x - 1$  and  $f''(x) = e^x > 0$ . Thus,  $f$  is convex and  $f(0) = 0$ . Therefore,  $f(x) \geq 0$  for all  $x \in \mathbb{R}$ .  $\square$

Figure 12:  $1 + x \leq e^x$

# Chapter 1

## The First Miracle: Robustness

Let  $f$  be a convex function, and let  $x^*$  be a minimizer of  $f$ .

### 1.1 Gradient Descent

**Definition 1.1.1.** *Gradient Descent*

$$x_{t+1} = x_t - \eta \nabla f(x_t) \quad (1.1)$$

It holds that:

$$f(x^*) \geq f(x_t) + \nabla f(x_t) \cdot (x^* - x_t) \quad (1.2)$$

$$0 \leq f(x_t) - f(x^*) \leq \nabla f(x_t) \cdot (x_t - x^*) \quad (1.3)$$

#### 1.1.1 Analysis of the Gradient Descent Algorithm

$$\begin{aligned} \|a\|^2 &= \|b\|^2 + \|a - b\|^2 \\ \|b\|^2 &= \|a\|^2 - \|a - b\|^2 = \|a\|^2 - (\|a\|^2 - 2a \cdot b + \|b\|^2) = 2a \cdot b - \|b\|^2 \end{aligned}$$

Then we have:

$$\begin{aligned} \|x^* - x_t\|^2 - \|x^* - x_{t+1}\|^2 &= -2\eta(x^* - x_t) \cdot \nabla f(x_t) - \eta^2 \|\nabla f(x_t)\|^2 \\ &= 2\eta(x_t - x^*) \cdot \nabla f(x_t) - \eta^2 \|\nabla f(x_t)\|^2 \\ &\geq 2(f(x_t) - f(x^*)) - \eta^2 L^2 \end{aligned}$$

Where the last inequality follows from the convexity and the Lipschitz continuity of  $f$ .

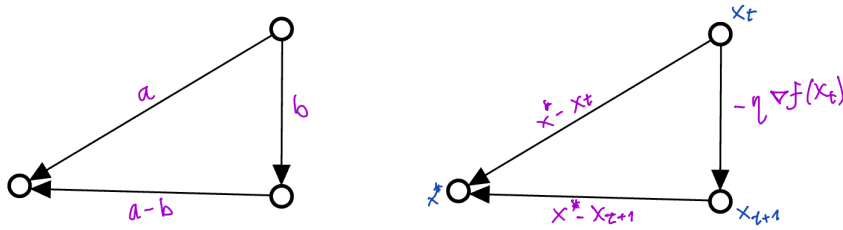


Figure 1.1: Gradient Descent

Then if we sum the above inequality from  $t = 1$  to  $T$ , we get:

$$\sum_{t=1}^T (f(x_t) - f(x^*)) \leq \frac{\|x_1 - x^*\|^2}{2\eta} + \frac{\eta L^2}{2} T$$

In fact, this is a specific case of the Fundamental Inequality of Optimization.

**Theorem 1.1.1.** *Fundamental Inequality of Optimization (unconstrained version)*

Suppose  $x_{t+1} = x_t - \eta g_t$  for all  $t$ , where  $g_1, \dots, g_T \in \mathbb{R}^d$  are arbitrary vectors. Then for all  $x^* \in \mathbb{R}^d$  it holds that

$$\sum_{t=1}^T g_t \cdot (x_t - x^*) \leq \frac{\|x_1 - x^*\|^2}{2\eta} + \frac{\eta}{2} \sum_{t=1}^T \|g_t\|^2.$$

*Proof.* Fundamental Inequality of Optimization

The proof tracks  $\|x_t - x^*\|^2$  as a “potential”. First write

$$\|x_{t+1} - x^*\|^2 = \|(x_t - x^*) - \eta g_t\|^2 = \|x_t - x^*\|^2 - 2\eta g_t \cdot (x_t - x^*) + \eta^2 \|g_t\|^2,$$

that is,

$$\|x_t - x^*\|^2 - \|x_{t+1} - x^*\|^2 = 2\eta g_t \cdot (x_t - x^*) - \eta^2 \|g_t\|^2.$$

Summing over  $t = 1, \dots, T$  and telescoping terms, we obtain

$$\|x_1 - x^*\|^2 - \|x_{T+1} - x^*\|^2 = 2\eta \sum_{t=1}^T g_t \cdot (x_t - x^*) - \eta^2 \sum_{t=1}^T \|g_t\|^2.$$

Organizing terms, we conclude:

$$\sum_{t=1}^T g_t \cdot (x_t - x^*) \leq \frac{\|x_1 - x^*\|^2 - \|x_{T+1} - x^*\|^2}{2\eta} + \frac{\eta}{2} \sum_{t=1}^T \|g_t\|^2.$$

□



## Chapter 2

# The Second Miracle: Potential Based

### 2.1 Experts Problem

At each time step, the player picks an action  $I_t \in [n]$  (we have  $n$  experts) and the adversary picks a loss vector  $l_t \in [0, 1]^n$ . The player incurs loss  $l_t(I_t)$  and the goal is to minimize the regret:

$$\text{Regret}_T(i) = \sum_{t=1}^T (l_t(I_t) - l_t(i)) \quad (2.1)$$

We consider the case where in each time step the player chooses an action from a distribution  $\vec{p}$  over the  $n$  experts (a vector from the simplex):

$$\vec{p} \in \Delta_n = \{\vec{p} \in \mathbb{R}_+^n : p_i \geq 0, \sum_{i=1}^n p_i = 1\}$$

#### Approach 1: Gradient Descent

We can use gradient descent on  $f_t(\vec{p}_t) = \vec{l}_t \cdot \vec{p}$ , where  $\vec{l}_t$  is the loss vector at time  $t$ . It holds that  $\nabla f_t(\vec{p}_t) = \vec{l}_t$ . We can use the analysis of the gradient descent algorithm for gradient descent of convex functions varying in time.

Let  $q \in \Delta_n$  be any distribution. Then we have:

$$\begin{aligned} f_t(q) &\geq f_t(\vec{p}_t) + \nabla f_t(q) \cdot (q - \vec{p}_t) \implies \\ f_t(\vec{p}_t) - f_t(q) &\leq \nabla f_t(q) \cdot (\vec{p}_t - q) \end{aligned}$$

Then:

$$\begin{aligned} \|q - p_t\|^2 - \|q - p_{t+1}\|^2 &= -2\eta(q - p_t) \cdot \nabla f_t(p_t) - \eta^2 \|\nabla f_t(p_t)\|^2 \implies \\ f_t(\vec{p}_t) - f_t(q) &\leq \nabla f_t(q) \cdot (\vec{p}_t - q) = \frac{1}{2\eta} (\|q - \vec{p}_t\|^2 - \|q - \vec{p}_{t+1}\|^2) + \frac{\eta}{2} \|\nabla f_t(\vec{p}_t)\|^2 \implies \\ \sum_{t=1}^T (f_t(\vec{p}_t) - f_t(q)) &\leq \frac{1}{2\eta} (\|q - \vec{p}_1\|^2 - \|q - \vec{p}_{T+1}\|^2) + \frac{\eta}{2} \sum_{t=1}^T \|\nabla f_t(\vec{p}_t)\|^2 \\ &\leq \frac{1}{2\eta} \|q - \vec{p}_1\|^2 + \frac{\eta}{2} \sum_{t=1}^T \|\nabla f_t(\vec{p}_t)\|^2 \\ &\leq \frac{1}{\eta} + \frac{\eta}{2} Tn = \mathcal{O}(\sqrt{Tn}) \end{aligned}$$

We have used the facts that:

- Both  $q$  and  $\vec{p}_1$  are distributions, so  $\|q - \vec{p}_1\|^2 \leq 2$ .
- $\|\nabla f_t(\vec{p}_t)\|^2 \leq n$  (as the loss vector is in  $0, 1^n$ ).

we can see that in this case, the rate of convergence DO depend on the dimension of the problem, in contrast to the non-varying case. The fact that the rate of convergence DO NOT depend on the dimension of the problem in GD is one of the reasons why GD is so useful in practice.

## Chapter 3

### The Third Miracle:



## Chapter 4

### The Fourth Miracle:



## Chapter 5

### The Fifth Miracle: