

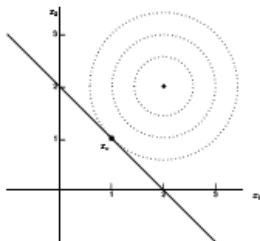
Optimization

Màster de Fonaments de Ciència de Dades

Lecture 1. Optimization. First examples and background

First introductory examples

Problem. Find the point on the line $x_1 + x_2 = 2$ that is closest to the point $(2, 2)^T$.



The problem can be written as

$$\begin{array}{ll} \text{minimize} & f(x, y) = (x - 2)^2 + (y - 2)^2, \\ \text{subject to} & x + y = 2. \end{array}$$

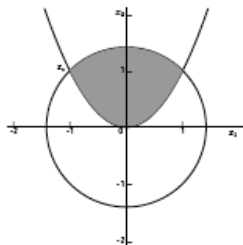
The solution is $(x, y)^T = (1, 1)^T$.

First introductory examples

Problem. Find the point such that:

$$\begin{array}{ll}\text{minimize} & f(x, y) = x, \\ \text{subject to} & x^2 \leq y, \\ & x^2 + y^2 \leq 2.\end{array}$$

In this example, the (feasible) set where we must look for the solution is defined by multiple constraints.



The solution (optimal point) is $(x, y)^T = (-1, 1)^T$.

First introductory examples

Problem. *Find the solution of:*

$$\text{minimize } f(x, y) = (e^x - 1)^2 + (y - 1)^2.$$

This is an example of an **unconstrained** optimization problem.

The **feasible set** here is the entire two-dimensional space.

The solution is $(x, y)^T = (0, 1)^T$, since the function value is zero only at this point and is positive elsewhere.

Introduction

- ▶ **What is Optimization?** Given a system or process, find the **best solution** to this process **within (or not) constraints**.
- ▶ **Objective Function:** Indicator of “goodness” of the solution, e.g., cost, profit, time, etc. In the above examples, the function f .
- ▶ **Decision Variables:** Variables that influence process behavior and can be adjusted for optimization. In the above examples, the variables x and y .
- ▶ We are interested in a **systematic approach** to the optimization process, and to make it as efficient as possible.
- ▶ **Optimization** is also called: *Mathematical Programming*, or *Operations Research*.

Current applications

- ▶ In modern times, (linear and nonlinear) optimization is used in optimal engineering design, finance, statistics and many other fields.
- ▶ Think of:
 - ▶ designing a car with minimal air resistance,
 - ▶ designing a bridge of minimal weight that still meets essential specifications,
 - ▶ defining a stock portfolio where the risk is minimal and the expected return high,...
- ▶ **Rule of thumb:** If you can make a mathematical model of your decision problem, then you can *try to optimize* it!

Optimization viewpoints

- ▶ **Mathematician** - characterization of theoretical properties of optimization, convergence, existence, local convergence rates.
- ▶ **Numerical Analyst** - implementation of optimization method for efficient and "practical" use. Concerned with fast computations, numerical stability, performance.
- ▶ **User** - applies optimization method to real problems. Concerned with reliability, robustness, efficiency, diagnosis, and recovery from failure.
- ▶ Optimization is a **fast moving research field**. Currently, there are over 30 journals devoted to optimization with roughly 200 published papers/month.
- ▶ In **this course**, we will see only the most **basic concepts, results, and procedures**.

Some classical optimization problems - I

1. **Dido's (or isoperimetric) problem.** Among all closed plain curves of a given length, find the one that encloses the largest area.
2. **Heron's problem.** Given two points A and B on the same side of a line L , find a point D on L such that the sum of the distances from A to D and from D to B is a minimum.
3. **Snell's law of refraction.** Given two points A and B on either side of a horizontal line L separating two (homogeneous) different media, find a point D on L such that the time it takes for a light ray to traverse the path ADB is a minimum.
Note: In an inhomogeneous medium, light travels from one point to another along the path requiring the shortest time ($v_i = c/n_i$).
4. **Euclid (Elements, 4th cent. B.C.).** In a given triangle ABC inscribe a parallelogram $ADEF$ ($EF \parallel AB, DE \parallel AC$) of maximal area.
5. **Steiner.** In the plane of a triangle, find a point (Fermat point) such that the sum of its distances to the vertices of the triangle is minimal

Some classical optimization problems - II

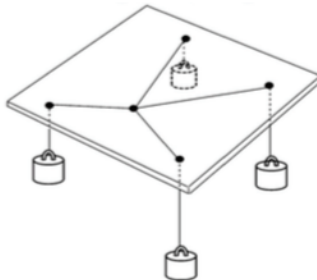
6. Find the maximum of the product of two numbers whose sum is given.
7. Find the maximal area of a right triangle whose small sides have constant sum.
8. In a given circle find a rectangle of maximal area.
9. In a given sphere find a cylinder of maximal volume.
10. Of all rectangular parallelepipeds inscribed in a sphere find the one of maximal volume.
11. Of all rectangular parallelepipeds with square base inscribed in a sphere find the one of maximal volume.
12. **The Brachistochrone.** Let two points A and B be given in a vertical plane. Find the curve that a point M , moving on a path AMB must follow such that, starting from A with zero velocity, it reaches B in the shortest time under its own gravity.

Some classical optimization problems - III

13. **The Fermat point of a set of points.** Given set of points y_1, \dots, y_m in the plane, find a point x^* whose sum of weighted distances to the given set of points is minimized. Mathematically, the problem is

$$\min \sum_{i=1}^m w_i \|x^* - y_i\|, \quad \text{subject to } x^* \in \mathbb{R}^2,$$

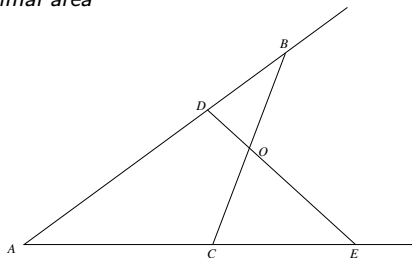
where w_1, \dots, w_m are given positive real numbers.



Some classical optimization problems - III

Exercise 1. To be delivered before 29-X-2019 as: Ex01-YourSurname.pdf.

13. **Smallest area problem.** *Given an angle with vertex A and a point O in its interior. Pass a line BC through the point O that cuts off from the angle a triangle of minimal area*



Hint: proof that for a triangle of minimal area the segments OB and OC should be equal.

The general optimization problem

Definition:

The **general nonlinear optimization** (NLO) problem can be written as follows:

$$\begin{array}{ll}\min & f(x), \\ \text{subject to} & g_i(x) = 0, \quad i \in I = \{1, \dots, m\}, \\ & h_j(x) \leq 0, \quad j \in J = \{1, \dots, p\}, \\ & x \in \mathcal{C},\end{array}$$

where $x \in \mathbb{R}^n$, $\mathcal{C} \subset \mathbb{R}^n$ is a certain set, and $f, g_1, \dots, g_m, h_1, \dots, h_p$ are real-valued functions defined on \mathcal{C} .

Terminology:

- ▶ The function f is called the **objective function** of the NLO.
- ▶ The set \mathcal{F} defined by:

$$\mathcal{F} = \{x \in \mathcal{C} : g_i(x) = 0, i = 1, \dots, m, h_j(x) \leq 0, j = 1, \dots, p\},$$

is called the **feasible set** (or **feasible region**).

- ▶ If $\mathcal{F} = \emptyset$ then we say that the optimization problem is **infeasible**.
- ▶ If the infimum of f over \mathcal{F} is attained at $x^* \in \mathcal{F}$, then we call x^* an **optimal solution** of the NLO, and $f(x^*)$ the **the optimal (objective) value of the NLO**.

Classification of optimization problems

- **Unconstrained Optimization:** The index sets I and J are empty:

$$g_1 = \dots = g_m = h_1 = \dots = h_p = 0,$$

and $\mathcal{C} = \mathbb{R}^n$.

- **Linear Optimization (LO)** (Linear programming): The functions $f, g_1, \dots, g_m, h_1, \dots, h_p$ are linear (affine: $F(\mathbf{x}) = A\mathbf{x} + \mathbf{b}$) and the set \mathcal{C} either equals to \mathbb{R}^n , the positive (negative) orthant \mathbb{R}_+^n , or is polyhedral.
- **Quadratic Optimization (QO):** The objective function f is quadratic:

$$f(\mathbf{x}) = \mathbf{x}^T Q \mathbf{x} + \mathbf{c}^T \mathbf{x} + \mathbf{d},$$

all the constraint functions $g_1, \dots, g_m, h_1, \dots, h_p$ are linear and the set \mathcal{C} is \mathbb{R}^n or the positive (negative) orthant \mathbb{R}_+^n , and Q is a $n \times n$ real matrix ($Q \in \mathbb{R}^{n \times n}$).

- **Quadratically Constrained Quadratic Optimization:** Same as QO, except that the constraint functions are quadratic.
- **Convex Quadratic Optimization (CQO).**
- **Convex Quadratically Constrained Quadratic Optimization:**
- ...

A well known application of Quadratic Optimization: Regression problems

- If a system

$$Ax = b, \quad A \in \mathbb{R}^{m \times n}, \quad x \in \mathbb{R}^n, \quad b \in \mathbb{R}^m,$$

has **more equations than unknowns** ($m > n$), then, in general, it has no solution, but we can compute the **least squares solution**

$$x^* = \min_{x \in \mathbb{R}^n} \|Ax - b\|,$$

for the **Euclidean norm** $\|x\| = \sqrt{\sum_{i=1}^n x_i^2} = \sqrt{x^T x} \geq 0$.

- Note that

$$\begin{aligned} \|Ax - b\|^2 &= (Ax - b)^T (Ax - b) \\ &= x^T A^T A x - 2b^T A x + \|b\|^2. \end{aligned}$$

- Note also that if $A \in \mathbb{R}^{m \times n}$, then $A^T A \in \mathbb{R}^{n \times n}$, $b^T A \in \mathbb{R}^n$, and introducing $z = Ax$:

$$x^T A^T A x = z^T z = \|z\|^2 \geq 0, \quad \forall x \in \mathbb{R}^n.$$

According to this last inequality, $A^T A$ will be **positive definite** if and only if for all $x \neq 0$ then $Ax \neq 0$, which is equivalent to say that the rank of A is n .

Example of regression problem: Concrete mixing

Mix concrete using n different gravel sizes s_1, s_2, \dots, s_n .

- ▶ The ideal mixture is given by $\mathbf{c} = (c_1, c_2, \dots, c_n)$, where c_i ($0 \leq c_i \leq 1$) is the fraction of size s_i in the mix, and $\sum_{i=1}^n c_i = 1$.
- ▶ Gravel mixtures come from m different mines.
- ▶ The gravel composition at each mine M_j given by $A_j = (a_{1j}, \dots, a_{nj})$ where $0 \leq a_{ij} \leq 1$ for all $i = 1, \dots, n$ and $\sum_{i=1}^n a_{ij} = 1$

	s_1	\dots	s_n	
M_1	a_{11}	\dots	a_{n1}	x_1 = fraction from M_1 in the mix
\vdots	\vdots		\vdots	\vdots
M_m	a_{1m}	\dots	a_{nm}	x_m = fraction from M_m in the mix

In the mix, the amount of gravel with size k should be close to c_k .

Concrete mixing: mathematical formulation

Exercise 2. To be delivered before 29-X-2019 as: Ex02-YourSurname.pdf.

Find the best possible approximation x_1, \dots, x_m of the ideal mixture, c_1, \dots, c_n , by using the material from the m mines.

Show that the optimal mixture will be the point x such that:

$$\begin{aligned} \min \quad & (Ax - c)^T (Ax - c), \\ \text{s.t.} \quad & \sum_{j=1}^m x_j = 1, \quad \text{and} \quad x_j \geq 0, \end{aligned}$$

where the matrix $A = (A_1, \dots, A_m)$ has A_j as columns.

Some mathematical notation and background

Notation and background

- ▶ Scalar and cross product
- ▶ Lines and planes
- ▶ Continuity
- ▶ Derivatives
- ▶ Gradients
- ▶ Approximation of functions

Scalar and cross product

Let $\mathbf{x} = (x_1, \dots, x_n)^T, \mathbf{y} = (y_1, \dots, y_n)^T \in \mathbb{R}^n$, we define:

- ▶ **Scalar (dot) product:** $\mathbf{x} \cdot \mathbf{y} = \mathbf{x}^T \mathbf{y} = x_1 y_1 + \dots + x_n y_n \in \mathbb{R}$.
- ▶ **Euclidean norm:** $\|\mathbf{x}\| = \sqrt{\mathbf{x} \cdot \mathbf{x}} = \sqrt{\mathbf{x}^T \mathbf{x}} = \sqrt{x_1^2 + \dots + x_n^2}$.
- ▶ **Euclidean distance:** $d(\mathbf{x}, \mathbf{y}) = \|\mathbf{y} - \mathbf{x}\| = \sqrt{(y_1 - x_1)^2 + \dots + (y_n - x_n)^2}$.
- ▶ **Cosinus of the angle:** $\cos(\widehat{\mathbf{x}, \mathbf{y}}) = \frac{\mathbf{x} \cdot \mathbf{y}}{\|\mathbf{x}\| \|\mathbf{y}\|}$.
- ▶ **Perpendicularity (orthogonality):** $\mathbf{x} \perp \mathbf{y} \Leftrightarrow \mathbf{x} \cdot \mathbf{y} = 0$.

Let $\mathbf{x} = (x_1, x_2, x_3), \mathbf{y} = (y_1, y_2, y_3) \in \mathbb{R}^3$, we define:

- ▶ **Cross product:**

$$\mathbf{x} \times \mathbf{y} = \begin{vmatrix} \mathbf{i} & \mathbf{j} & \mathbf{k} \\ x_1 & x_2 & x_3 \\ y_1 & y_2 & y_3 \end{vmatrix} = \begin{pmatrix} x_2 y_3 - x_3 y_2 \\ x_3 y_1 - x_1 y_3 \\ x_1 y_2 - x_2 y_1 \end{pmatrix}.$$

Note that

$$\mathbf{x} \times \mathbf{y} \perp \mathbf{x} \quad \text{and} \quad \mathbf{x} \times \mathbf{y} \perp \mathbf{y}.$$

Lines and planes

- In \mathbb{R}^2 : The **line** determined by the **point** $\mathbf{a} = (a_1, a_2)^T$ and the **vector** $\mathbf{v} = (v_1, v_2)^T$ is

$$\mathbf{x} = \mathbf{a} + t\mathbf{v}, \quad t \in \mathbb{R} \quad \Leftrightarrow \quad \begin{pmatrix} x \\ y \end{pmatrix} = \begin{pmatrix} a_1 \\ a_2 \end{pmatrix} + t \begin{pmatrix} v_1 \\ v_2 \end{pmatrix}, \quad t \in \mathbb{R},$$

that can also be written as

$$\frac{x - a_1}{v_1} = \frac{y - a_2}{v_2} \quad \Leftrightarrow \quad Ax + By + C = 0,$$

with $A = v_2$, $B = -v_1$, $C = -a_1 v_2 + a_2 v_1$.

- In \mathbb{R}^3 : The **line** determined by the **point** $\mathbf{a} = (a_1, a_2, a_3)^T$ and the **vector** $\mathbf{v} = (v_1, v_2, v_3)^T$ is

$$\mathbf{x} = \mathbf{a} + t\mathbf{v}, \quad t \in \mathbb{R} \quad \Leftrightarrow \quad \begin{pmatrix} x \\ y \\ z \end{pmatrix} = \begin{pmatrix} a_1 \\ a_2 \\ a_3 \end{pmatrix} + t \begin{pmatrix} v_1 \\ v_2 \\ v_3 \end{pmatrix}, \quad t \in \mathbb{R},$$

that can also be written as

$$\frac{x - a_1}{v_1} = \frac{y - a_2}{v_2} = \frac{z - a_3}{v_3}.$$

Lines and planes

- In \mathbb{R}^3 : The plane determined by the point $\mathbf{a} = (a_1, a_2, a_3)^T$ and the vectors $\mathbf{u} = (u_1, u_2, u_3)^T$ and $\mathbf{v} = (v_1, v_2, v_3)^T$ is

$$\mathbf{x} = \mathbf{a} + t\mathbf{u} + s\mathbf{v} \quad \Leftrightarrow \quad \begin{pmatrix} x \\ y \\ z \end{pmatrix} = \begin{pmatrix} a_1 \\ a_2 \\ a_3 \end{pmatrix} + t \begin{pmatrix} u_1 \\ u_2 \\ u_3 \end{pmatrix} + s \begin{pmatrix} v_1 \\ v_2 \\ v_3 \end{pmatrix},$$

with $t, s \in \mathbb{R}$.

- The above equation of the plane can also be written as

$$\begin{vmatrix} x - a_1 & y - a_2 & z - a_3 \\ u_1 & u_2 & u_3 \\ v_1 & v_2 & v_3 \end{vmatrix} = 0$$

or as

$$Ax + By + Cz + D = 0,$$

with $(A, B, C)^T = \mathbf{u} \times \mathbf{v}$.

Continuity

Consider the function

$$f : \mathcal{C} \subset \mathbb{R}^n \longrightarrow \mathbb{R},$$

we define:

- ▶ The **domain \mathcal{C} of f** as the set of points $\mathbf{x} \in \mathbb{R}^n$ where f is defined.
- ▶ The **graph of f** , as the subset of \mathbb{R}^{n+1} defined by:

$$\{(\mathbf{x}, z) \in \mathbb{R}^{n+1} : \mathbf{x} = (x_1, \dots, x_n)^T \in \mathcal{C} \subset \mathbb{R}^n, z = f(\mathbf{x}) \in \mathbb{R}\} \subset \mathbb{R}^{n+1}.$$

- ▶ For each $c \in \mathbb{R}$, the **level set c of f** as:

$$f^{-1}(c) = \{\mathbf{x} \in \mathcal{C} : f(\mathbf{x}) = c\} \subset \mathbb{R}^n.$$

- ▶ We say that f is **continuous at a point $\mathbf{a} \in \mathcal{C}$** if and only if

$$\lim_{\mathbf{x} \rightarrow \mathbf{a}} f(\mathbf{x}) = f(\mathbf{a}).$$

Continuity

Some **fundamental properties** of continuous functions are:

- ▶ The elementary functions of one variable $e^x, \log x, \sin x, \cos x, \dots$ and the coordinate functions

$$\begin{array}{lll} x_i : & \mathbb{R}^n & \longrightarrow \mathbb{R} \\ & \mathbf{x} = (x_1, \dots, x_n)^T & \longrightarrow x_i \end{array}$$

are continuous in their domain.

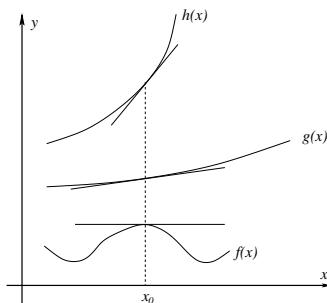
- ▶ Addition, subtraction, product, division (except at the points where the denominator vanishes) and composition of continuous functions are also continuous functions.
- ▶ Given a **continuous** function

$$f : \mathcal{C} \subset \mathbb{R}^n \longrightarrow \mathbb{R},$$

such that \mathcal{C} is **compact** (closed and bounded), then f is **bounded** and f **attains its maximum and minimum values on \mathcal{C} .**

Derivatives

- ▶ The **derivative of a function** $y = f(x)$ of a variable x is a measure of the rate at which the value y of the function changes with respect to the change of the variable x .
- ▶ If x and y are **real numbers**, and if the **graph** of f is plotted against x , the **derivative** is the **slope** of this graph at each point.



Derivatives

Let f be a real valued function defined in an open neighborhood of a real number a , then:

- ▶ The **derivative** of $y = f(x)$ with respect to x at a is, geometrically, the **slope of the tangent line** to the graph of f at $(a, f(a))^T$.
- ▶ The slope of the tangent line is very close to the slope of the line through $(a, f(a))$ and a nearby point on the graph, for example $(a + h, f(a + h))^T$.
- ▶ The slope m of the secant line is

$$m = \frac{\Delta f(a)}{\Delta a} = \frac{f(a + h) - f(a)}{(a + h) - (a)} = \frac{f(a + h) - f(a)}{h}.$$

- ▶ A value of h close to zero gives, in general, a good approximation to the slope of the tangent line

Derivatives. Rigorous definition

- ▶ Geometrically, the **limit of the secant lines is the tangent line**. Therefore, the limit of the difference quotient as h approaches zero, if it exists, should represent the slope of the tangent line to $(a, f(a))$.
- ▶ This limit is defined to be the **derivative of the function f at a** :

$$f'(a) = \lim_{h \rightarrow 0} \frac{f(a+h) - f(a)}{h}.$$

- ▶ When the limit exists, f is said to be **differentiable at a** .
- ▶ Equivalently, the derivative satisfies the property that

$$\lim_{h \rightarrow 0} \frac{f(a+h) - f(a) - f'(a) \cdot h}{h} = 0,$$

which has the intuitive interpretation that the tangent line to f at a gives the **best linear approximation**

$$f(a+h) \approx f(a) + f'(a)h,$$

to f near a .

Derivatives in higher dimensions

- ▶ A vector-valued function $\mathbf{y}(t)$ of a real variable sends real numbers to vectors in some vector space (\mathbb{R}^n).

$$\begin{array}{ccc} \mathbf{y} : & \mathbb{R} & \longrightarrow \mathbb{R}^n \\ & t & \longrightarrow \mathbf{y}(t). \end{array}$$

- ▶ A vector-valued function can be split up into its coordinate functions

$$\mathbf{y}(t) = (y_1(t), \dots, y_n(t))^T.$$

- ▶ The derivative of the **curve** $\mathbf{y}(t)$ is defined to be the vector, called the **tangent vector**, whose coordinates are the derivatives of the coordinate functions

$$\mathbf{y}'(t) = (y_1'(t), \dots, y_n'(t))^T, \quad \text{or equivalently} \quad \mathbf{y}'(t) = \lim_{h \rightarrow 0} \frac{\mathbf{y}(t+h) - \mathbf{y}(t)}{h},$$

if the limit exists.

- ▶ If $\mathbf{e}_1, \dots, \mathbf{e}_n$ is the standard basis for \mathbb{R}^n , then

$$\mathbf{y}(t) = y_1(t)\mathbf{e}_1 + \dots + y_n(t)\mathbf{e}_n,$$

and since each of the basis vectors is a constant

$$\mathbf{y}'(t) = y_1'(t)\mathbf{e}_1 + \dots + y_n'(t)\mathbf{e}_n.$$

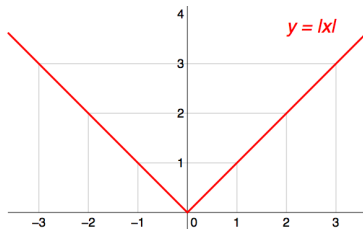
Continuity and differentiability

- **Property:** If

$$\begin{array}{rcl} f : \mathbb{R}^n & \longrightarrow & \mathbb{R} \\ x & \longrightarrow & f(x) \end{array}$$

is **differentiable** at **a** , then **f** **must** also be **continuous** at **a** .

- **Property:** If a function is **continuous** at a point it **may not** be **differentiable** there.
- **Example:** The absolute value function $f(x) = |x|$ is continuous at $x = 0$, but it is not differentiable there, since the tangent slopes do not approach the same value from the left as they do from the right.



Partial derivatives

- ▶ If f is a real value function that depends on n variables

$$\begin{aligned} f : \mathbb{R}^n &\longrightarrow \mathbb{R} \\ \mathbf{x} &\longrightarrow f(\mathbf{x}) = f(x_1, \dots, x_n), \end{aligned}$$

the **partial derivative** of $f(\mathbf{x})$ in the direction x_i at the point $\mathbf{a} = (a_1, \dots, a_n)^T$ is defined to be:

$$\frac{\partial f}{\partial x_i}(\mathbf{a}) = \lim_{h \rightarrow 0} \frac{f(a_1, \dots, a_i + h, \dots, a_n) - f(a_1, \dots, a_i, \dots, a_n)}{h}.$$

- ▶ In the above difference quotient, all the variables except x_i are held fixed. That choice of fixed values determines a function of one variable

$$f_{a_1, \dots, a_{i-1}, a_{i+1}, \dots, a_n}(x_i) = f(a_1, \dots, a_{i-1}, x_i, a_{i+1}, \dots, a_n),$$

and, by definition:

$$\frac{df_{a_1, \dots, a_{i-1}, a_{i+1}, \dots, a_n}}{dx_i}(a_i) = \frac{\partial f}{\partial x_i}(\mathbf{a}).$$

First and second partial derivatives

Let $\mathbf{a} \in \mathcal{C} \subset \mathbb{R}^n$ be a point where the real function

$$f : \mathcal{C} \longrightarrow \mathbb{R},$$

is differentiable.

- ▶ **Property:** If a real-valued function f is differentiable at an interior point $\mathbf{a} \in \mathcal{C}$, then its first partial derivatives exist at \mathbf{a} .
- ▶ **Definition:** If the partial derivatives are continuous at \mathbf{a} , then f is said to be **continuously differentiable** at \mathbf{a} .
- ▶ **Property:** If f is **twice differentiable** at $\mathbf{a} \in \mathcal{C}$, then the second partial derivatives exist there.
- ▶ **Definition:** If the second partial derivatives are continuous at \mathbf{a} , then f is said to be **twice continuously differentiable** at \mathbf{a} .
- ▶ **Definition:** If f is twice continuously differentiable at \mathbf{a} we define the **Hessian** matrix of f at \mathbf{a} as the $n \times n$ symmetric matrix $\nabla^2 f(\mathbf{a})$ given by:

$$\nabla^2 f(\mathbf{a}) = \left(\frac{\partial^2 f(\mathbf{a})}{\partial x_i \partial x_j} \right), \quad i, j = 1, \dots, n.$$

Directional derivatives

- ▶ If f is a real-valued function on \mathbb{R}^n , then the partial derivatives of f measure its variation in the direction of the coordinate axes.
- ▶ If f is a function of x and y ($x, y \in \mathbb{R}$), then its partial derivatives measure the variation in f in the x direction and the y direction. They do not, however, directly measure the variation of f in any other direction, such as along the diagonal line $y = x$.
- ▶ These are measured using directional derivatives. Choose a vector

$$\mathbf{v} = (v_1, \dots, v_n)^T.$$

The **directional derivative** of f in the direction of \mathbf{v} at the point \mathbf{x} is defined by

$$D_{\mathbf{v}}f(\mathbf{x}) = \lim_{h \rightarrow 0} \frac{f(\mathbf{x} + h\mathbf{v}) - f(\mathbf{x})}{h} = \left. \frac{d}{dt} \right|_{t=0} f(\mathbf{x} + t\mathbf{v}) = \sum_{j=1}^n \frac{\partial f(\mathbf{x})}{\partial x_j} v_j,$$

where we have used the chain rule to get the last equality.

The chain rule

► Let

$$\begin{array}{ccc} \alpha : & I \subset \mathbb{R} & \longrightarrow & C \\ & t & \longrightarrow & \alpha(t) = (x_1(t), \dots, x_n(t))^T, \end{array}$$

be a **differentiable curve** in $C \subset D \subset \mathbb{R}^n$ and

$$\begin{array}{ccc} f : & D \subset \mathbb{R}^n & \longrightarrow & \mathbb{R}^m \\ & \mathbf{x} & \longrightarrow & f(\mathbf{x}) \end{array}$$

be a differentiable function. Then

$$f(\alpha(t)) = f(x_1(t), \dots, x_n(t)),$$

and

$$\frac{d}{dt} f(\alpha(t)) = \frac{\partial f}{\partial x_1}(\alpha(t)) x_1'(t) + \dots + \frac{\partial f}{\partial x_n}(\alpha(t)) x_n'(t).$$

Directional derivatives

- ▶ We want to compute the directional derivative after **changing the length of the vector \mathbf{v}** .
- ▶ Suppose that $\mathbf{v} = \lambda \mathbf{u}$. If in

$$\frac{f(\mathbf{x} + h\mathbf{v}) - f(\mathbf{x})}{h},$$

we substitute $h = k/\lambda$ and $\mathbf{v} = \lambda \mathbf{u}$, we get

$$\frac{f(\mathbf{x} + (k/\lambda)(\lambda \mathbf{u})) - f(\mathbf{x})}{k/\lambda} = \lambda \cdot \frac{f(\mathbf{x} + k\mathbf{u}) - f(\mathbf{x})}{k}.$$

This is λ times the difference quotient that we had for the directional derivative of f with respect to \mathbf{u} .

- ▶ Taking the limit as h tends to zero is the same as taking the limit as k tends to zero, because h and k are multiples of each other.
- ▶ Therefore, $D_{\mathbf{v}}(f) = \lambda D_{\mathbf{u}}(f)$. Because of this rescaling property, **directional derivatives are considered only for unit vectors**: $\|\mathbf{v}\| = 1$.

The gradient

Consider the function

$$f : \mathcal{C} \subset \mathbb{R}^n \longrightarrow \mathbb{R}.$$

- ▶ If f has a partial derivatives $\partial f / \partial x_j$ with respect to each variable x_j , then at any point $\mathbf{a} \in \mathcal{C}$, these partial derivatives define the vector

$$\nabla f(\mathbf{a}) = \left(\frac{\partial f}{\partial x_1}(\mathbf{a}), \dots, \frac{\partial f}{\partial x_n}(\mathbf{a}) \right)^T.$$

This vector is called the **gradient of f at \mathbf{a}** .

- ▶ **Theorem:** If all the partial derivatives of f exist and are **continuous** at \mathbf{a} , then the function f is **differentiable** at \mathbf{a} and the gradient of f at \mathbf{a} exists
- ▶ From

$$D_{\mathbf{v}}f(\mathbf{a}) = \sum_{j=1}^n \frac{\partial f(\mathbf{a})}{\partial x_j} v_j,$$

we get

$$D_{\mathbf{v}}f(\mathbf{a}) = (\nabla f(\mathbf{a})) \cdot \mathbf{v}.$$

Properties of the gradient

- **Property:** If $f : D \subset \mathbb{R}^n \rightarrow \mathbb{R}$ is differentiable, $\mathbf{a} \in D$, and $\mathbf{u} \in \mathbb{R}^n$ is a unitary vector ($\|\mathbf{u}\| = 1$), then

$$D_{\mathbf{u}}f(\mathbf{a}) = (\nabla f(\mathbf{a})) \cdot \mathbf{u} = \|\nabla f(\mathbf{a})\| \cos \theta,$$

where θ is the angle between \mathbf{u} and $\nabla f(\mathbf{a})$.

- **Property:** The gradient vector $\nabla f(\mathbf{a})$ gives the maximum direction variation of f at the point \mathbf{a} (since $\cos \theta$ is maximum $\Leftrightarrow \theta = 0$).
- **Property:** Gradients are orthogonal to the level curves and the level surfaces of a function f .

Proof. Let $\mathbf{r}(t) = (x_1(t), x_2(t), \dots, x_n(t))^T$ be a level curve (or a curve on a level surface) this means that $f(\mathbf{r}(t))$ is constant for any value of t . Then

$$\frac{d}{dt}f(\mathbf{r}(t)) = 0.$$

Using the chain rule for the computation of the derivative, we get

$$\begin{aligned} \frac{d}{dt}f(\mathbf{r}(t)) &= \frac{d}{dt}f(x_1(t), x_2(t), \dots, x_n(t)) \\ &= \frac{\partial f}{\partial x_1}(\mathbf{r}(t))x_1'(t) + \dots + \frac{\partial f}{\partial x_n}(\mathbf{r}(t))x_n'(t) = \nabla f(\mathbf{r}(t))^T \mathbf{r}'(t), \end{aligned}$$

and since $\mathbf{r}'(t)$ is the tangent vector to the curve, the property follows.

Properties of the gradient. Examples

- ▶ **Property:** The equations of the tangent plane and the normal line of the level set of f at \mathbf{a} are:

- ▶ Tangent plane

$$(\nabla f(\mathbf{a})) \cdot (\mathbf{x} - \mathbf{a}) = 0 \quad \Leftrightarrow \quad \frac{\partial f}{\partial x_1}(x_1 - a_1) + \cdots + \frac{\partial f}{\partial x_n}(x_n - a_n) = 0.$$

- ▶ Normal line

$$\mathbf{x} = \mathbf{a} + \lambda \nabla f(\mathbf{a}), \quad \lambda \in \mathbb{R}.$$

- ▶ **Example:** Compute the tangent plane to the surface $3x^2y + z^2 - 4 = 0$ at the point $(1, 1, 1)^T$.

Let $f(\mathbf{x}) = 3x^2y + z^2 - 4$, since

$$\begin{aligned}\nabla f(\mathbf{x})^T &= (6xy, 3x^2, 2z)^T, \\ \nabla f(1, 1, 1)^T &= (6, 3, 2)^T,\end{aligned}$$

the plane is

$$6(x - 1) + 3(y - 1) + 2(z - 1) = 0 \quad \Leftrightarrow \quad 6x + 3y + 2z = 11.$$

Linear approximation of functions

- ▶ We have already seen that if f is a real function in one variable, the linear approximation of the function $f(x)$ at a point x_0 is defined by the linear function

$$L(x) = f(x_0) + f'(x_0)(x - x_0).$$

- ▶ In two dimensions, the linear approximation of the function $f(x, y)$ at the point $(x_0, y_0)^T$ is defined as the linear function

$$\begin{aligned} L(x, y) &= f(x_0, y_0) + \frac{\partial f}{\partial x}(x_0, y_0)(x - x_0) + \frac{\partial f}{\partial y}(x_0, y_0)(y - y_0) \\ &= f(x_0, y_0) + \left(\frac{\partial f}{\partial x}(x_0, y_0), \frac{\partial f}{\partial y}(x_0, y_0) \right) \begin{pmatrix} x - x_0 \\ y - y_0 \end{pmatrix} \\ &= f(x_0, y_0) + (\nabla f(x_0, y_0))^T \begin{pmatrix} x - x_0 \\ y - y_0 \end{pmatrix}. \end{aligned}$$

- ▶ In dimension n

$$L(\mathbf{x}) = f(\mathbf{x}_0) + (\nabla f(\mathbf{x}_0))^T (\mathbf{x} - \mathbf{x}_0).$$

Linear approximation of functions

- **Example:** Estimate the value of $f(0.01, 24.8, 1.02)$ for $f(x, y, z) = e^x \sqrt{y} z$.

We take $\mathbf{x}_0 = (0, 5, 1)^T$ and we use the linear approximation of f to compute an estimation of $f(0.01, 24.8, 1.02)$.

Clearly

$$\begin{aligned} f(\mathbf{x}_0) &= 5, \\ \nabla f(\mathbf{x})^T &= \left(e^x \sqrt{y} z, \frac{e^x z}{2\sqrt{y}}, e^x \sqrt{y} \right)^T, \\ \nabla f(\mathbf{x}_0)^T &= (5, 1/10, 5)^T, \\ L(\mathbf{x}) &= f(\mathbf{x}_0) + (\nabla f(\mathbf{x}_0))^T (\mathbf{x} - \mathbf{x}_0) \\ &= 5 + (5, 1/10, 5) \begin{pmatrix} x - 0 \\ y - 5 \\ z - 1 \end{pmatrix} = 5 + 5x + \frac{y - 5}{10} + 5(z - 1) \end{aligned}$$

We approximate $f(0.01, 24.8, 1.02) = 5.1306$ by $L(0.01, 24.8, 1.02) = 5.13$

The differential matrix

Let

$$\begin{aligned} f : \mathcal{C} \subset \mathbb{R}^n &\longrightarrow \mathbb{R}^m \\ \mathbf{x} &\longrightarrow f(\mathbf{x}) = (f_1(\mathbf{x}), \dots, f_m(\mathbf{x})). \end{aligned}$$

- ▶ We say that f is differentiable if f_1, \dots, f_m are differentiable.
- ▶ The differential of f at an interior point $\mathbf{a} \in \mathcal{C}$ is

$$Df(\mathbf{a}) = \begin{pmatrix} \nabla f_1(\mathbf{a}) \\ \vdots \\ \nabla f_m(\mathbf{a}) \end{pmatrix} = \begin{pmatrix} \frac{\partial f_1(\mathbf{a})}{\partial x_1} & \cdots & \frac{\partial f_1(\mathbf{a})}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m(\mathbf{a})}{\partial x_1} & \cdots & \frac{\partial f_m(\mathbf{a})}{\partial x_n} \end{pmatrix}.$$

- ▶ If $g : \mathcal{D} \subset \mathbb{R}^p \longrightarrow \mathcal{C} \subset \mathbb{R}^n$ and $f : \mathcal{C} \subset \mathbb{R}^n \longrightarrow \mathbb{R}^m$ are both differentiable, then the composition $h = f \circ g$

$$\begin{aligned} h : \mathcal{D} &\longrightarrow \mathcal{C} && \longrightarrow \mathbb{R}^m \\ \mathbf{x} &\longrightarrow g(\mathbf{x}) = (g_1(\mathbf{x}), \dots, g_n(\mathbf{x}))^T && \longrightarrow h(\mathbf{x}) = f(g_1(\mathbf{x}), \dots, g_n(\mathbf{x})) \end{aligned}$$

is also differentiable.

The differential matrix

If $g : \mathcal{D} \subset \mathbb{R}^p \longrightarrow \mathcal{C} \subset \mathbb{R}^n$ and $f : \mathcal{C} \subset \mathbb{R}^n \longrightarrow \mathbb{R}^m$ are both differentiable, then the differential of the composition $h = f \circ g$ at an interior point $\mathbf{a} \in \mathcal{D}$ is the product of the differentials

$$Dh(\mathbf{a}) = Df(g(\mathbf{a}))Dg(\mathbf{a})$$
$$Dh(\mathbf{a}) = \begin{pmatrix} \frac{\partial f_1(g(\mathbf{a}))}{\partial x_1} & \cdots & \frac{\partial f_1(g(\mathbf{a}))}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m(g(\mathbf{a}))}{\partial x_1} & \cdots & \frac{\partial f_m(g(\mathbf{a}))}{\partial x_n} \end{pmatrix} \begin{pmatrix} \frac{\partial g_1(\mathbf{a})}{\partial x_1} & \cdots & \frac{\partial g_1(\mathbf{a})}{\partial x_p} \\ \vdots & \ddots & \vdots \\ \frac{\partial g_n(\mathbf{a})}{\partial x_1} & \cdots & \frac{\partial g_n(\mathbf{a})}{\partial x_p} \end{pmatrix}.$$

The differential matrix. Linear approximations

- ▶ If $f : I \subset \mathbb{R} \longrightarrow \mathbb{R}$ is differentiable, then for $dx \approx 0$

$$f(x + dx) \approx f(x) + f'(x)dx$$

- ▶ If $f : \mathcal{C} \subset \mathbb{R}^n \longrightarrow \mathbb{R}$ is differentiable, $\mathbf{x} = (x_1, \dots, x_n)^T$, $d\mathbf{x} = (dx_1, \dots, dx_n)^T \approx \mathbf{0}$, then

$$f(\mathbf{x} + d\mathbf{x}) \approx f(\mathbf{x}) + (\nabla f(\mathbf{x})) \cdot d\mathbf{x}$$

- ▶ If $f : \mathcal{C} \subset \mathbb{R}^n \longrightarrow \mathbb{R}^m$ is differentiable, $\mathbf{x} = (x_1, \dots, x_n)^T$, $d\mathbf{x} = (dx_1, \dots, dx_n)^T \approx \mathbf{0}$, then

$$f(\mathbf{x} + d\mathbf{x}) \approx f(\mathbf{x}) + DF(\mathbf{x}) d\mathbf{x}$$

Critical points

- **Definition.** Given a differentiable function $f : \mathcal{C} \subset \mathbb{R}^n \rightarrow \mathbb{R}$, \mathbf{a} is a **critical point of f** is

$$\nabla f(\mathbf{a}) = \mathbf{0} \quad \Leftrightarrow \quad \left\{ \begin{array}{l} \frac{\partial f(\mathbf{a})}{\partial x_1} = 0, \\ \vdots \\ \frac{\partial f(\mathbf{a})}{\partial x_n} = 0. \end{array} \right.$$

- If \mathbf{a} is not a critical point of f , then $\nabla f(\mathbf{a})$ gives the direction along which f increases or decreases faster. In particular, if \mathbf{a} is not a critical point of f then it can be not a maximum or minimum of f .
- The critical points of f are the candidates to be the **local extrema** (relative extrema) of f .

Quadratic approximation of functions

- ▶ We have already seen that, in dimension n , the linear approximation of the function $f(\mathbf{x})$ at a point \mathbf{a} is defined by the function

$$L(\mathbf{x}) = f(\mathbf{a}) + \nabla f(\mathbf{a})(\mathbf{x} - \mathbf{a}).$$

- ▶ Is \mathbf{a} is a critical point of f , then $\nabla f(\mathbf{a}) = 0$, and the linear approximation of f at \mathbf{a} is constant.
- ▶ The second order approximation is obtained using Taylor's formula

$$f(\mathbf{x}) = f(\mathbf{a}) + \nabla f(\mathbf{a})(\mathbf{x} - \mathbf{a}) + \frac{1}{2}\nabla^2 f(\mathbf{a})(\mathbf{x} - \mathbf{a})^2 + \dots$$

where the value of $\nabla^2 f(\mathbf{a})(\mathbf{x} - \mathbf{a})^2 \in \mathbb{R}$ is given by

$$(x_1 - a_1, \dots, x_n - a_n) \begin{pmatrix} \frac{\partial^2 f(\mathbf{a})}{\partial x_1^2} & \dots & \frac{\partial^2 f(\mathbf{a})}{\partial x_n \partial x_1} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 f(\mathbf{a})}{\partial x_1 \partial x_n} & \dots & \frac{\partial^2 f(\mathbf{a})}{\partial x_n^2} \end{pmatrix} \begin{pmatrix} x_1 - a_1 \\ \vdots \\ x_n - a_n \end{pmatrix}.$$

- ▶ Denoting the Hessian $\nabla^2 f(\mathbf{a})$ by $H(\mathbf{a})$, the quadratic approximation of f at the point \mathbf{a} is written as

$$Q(\mathbf{x}) = f(\mathbf{a}) + \nabla f(\mathbf{a})(\mathbf{x} - \mathbf{a}) + \frac{1}{2}(\mathbf{x} - \mathbf{a})^T H(\mathbf{a})(\mathbf{x} - \mathbf{a}).$$

Quadratic functions

- ▶ For any $n \times n$ matrix Q ($Q \in \mathbb{R}^{n \times n}$) we have

$$Q \text{ is symmetric} \Leftrightarrow Q^T = Q$$

$$Q \text{ is skew-symmetric} \Leftrightarrow Q^T = -Q$$

$$Q \text{ is positive semidefinite (PSD)} \Leftrightarrow x^T Q x \geq 0 \text{ for all } x \in \mathbb{R}^n$$

$$Q \text{ is positive definite (PD)} \Leftrightarrow \begin{aligned} &x^T Q x \geq 0 \text{ for all } x \in \mathbb{R}^n \\ &\text{and } x^T Q x = 0 \text{ if and only if } x = 0 \end{aligned}$$

- ▶ Let f be the quadratic function given by

$$f(x) = x^T Q x + c^T x + d$$

where $Q \in \mathbb{R}^{n \times n}$, $c \in \mathbb{R}^n$ and $d \in \mathbb{R}$. Then f is:

$$\text{▶ linear} \quad \Leftrightarrow \quad Q = 0 \text{ and } d = 0 \quad \Rightarrow \quad f(x) = c^T x$$

$$\text{▶ affine} \quad \Leftrightarrow \quad Q = 0 \quad \Rightarrow \quad f(x) = c^T x + d$$

$$\text{▶ convex} \quad \Leftrightarrow \quad Q \text{ is PSD} \quad \Rightarrow \quad f(x) = x^T Q x + c^T x + d$$