

# Mathematical modelling

Lecture notes

Faculty of Computer and Information Science  
University of Ljubljana

2022/23

## Chapter 0:

# What is Mathematical Modelling?

- ▶ Types of models
- ▶ Modelling cycle
- ▶ Numerical errors

# Introduction

The task of mathematical modelling is to find and evaluate solutions to real world problems with the use of mathematical concepts and tools.

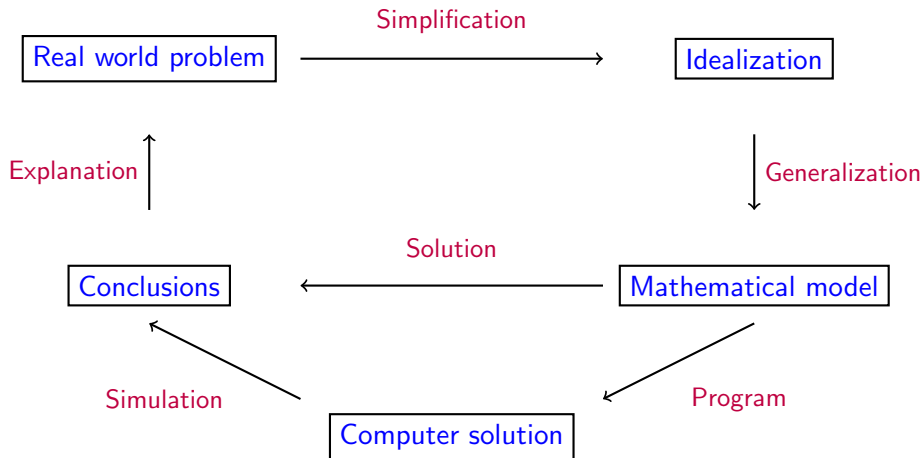
In this course we will introduce some (by far not all) mathematical tools that are used in setting up and solving mathematical models.

We will (together) also solve specific problems, study examples and work on projects.

# Contents

- ▶ Introduction
- ▶ Linear models: systems of linear equations, matrix inverses, SVD decomposition, PCA
- ▶ Nonlinear models: vector functions, linear approximation, solving systems of nonlinear equations
- ▶ Geometric models: curves and surfaces
- ▶ Dynamical models: differential equations, dynamical systems

# Modelling cycle



## What should we pay attention to?

- ▶ Simplification: relevant assumptions of the model (distinguish important features from irrelevant)
- ▶ Generalization: choice of mathematical representations and tools (for example: how to represent an object - as a point, a geometric shape, ...)
- ▶ Solution: as simple as possible and well documented
- ▶ Conclusions: are the results within the expected range, do they correspond to “facts” and experimental results?

A mathematical model is not universal, it is an approximation of the real world that works only within a certain scale where the assumptions are at least approximately realistic.

## Example

An object (ball) with mass  $m$  is thrown vertically into the air. What should we pay attention to when modelling its motion?

- ▶ The assumptions of the model: relevant forces and parameters (gravitation, friction, wind, ...), how to model the object (a point, a homogeneous or nonhomogeneous geometric object, angle and rotation in the initial thrust, ...)
- ▶ Choice of the mathematical model: differential equation, discrete model, ...
- ▶ Computation: analytic or numeric, choice of method, ...
- ▶ Do the results make sense?

## Errors

An important part of modelling is estimating the errors!

Errors are an integral part of every model.

Errors come from: assumptions of the model, imprecise data, mistakes in the model, computational precision, errors in numerical and computational methods, mistakes in the computations, mistakes in the programs, ...

Absolute error = Approximate value - Correct value

$$\Delta x = \bar{x} - x$$

Relative error =  $\frac{\text{Absolute error}}{\text{Correct value}}$

$$\delta_x = \frac{\Delta x}{x}$$



## Example: quadratic equation

$$x^2 + 2a^2x - q = 0$$

Analytic solutions are

$$x_1 = -a^2 - \sqrt{a^4 + q} \quad \text{and} \quad x_2 = -a^2 + \sqrt{a^4 + q}.$$

What happens if  $a^2 = 10000$ ,  $q = 1$ ? **Problem with stability in calculating  $x_2$ .**

More stable way for computing  $x_2$  (so that we do not subtract numbers which are nearly the same) is

$$\begin{aligned} x_2 &= -a^2 + \sqrt{a^4 + q} = \frac{(-a^2 + \sqrt{a^4 + q})(a^2 + \sqrt{a^4 + q})}{a^2 + \sqrt{a^4 + q}} \\ &= \frac{q}{a^2 + \sqrt{a^4 + q}}. \end{aligned}$$

## Example of real life disasters

- ▶ Disasters caused because of numerical errors:  
(<http://www-users.math.umn.edu/~arnold//disasters/>)
  - ▶ **The Patriot Missile failure, Dharan, Saudi Arabia, February 25 1991**, 28 deaths: **bad analysis of rounding errors.**
  - ▶ **The exploding of the Ariane 5 rocket, French Guiana, June 4, 1996**: **the consequence of overflow in the horizontal velocity.**  
[https://www.youtube.com/watch?v=PK\\_yguLapgA](https://www.youtube.com/watch?v=PK_yguLapgA)  
<https://www.youtube.com/watch?v=W3YJeoYgozw>  
<https://www.arianespace.com/vehicle/ariane-5/>
  - ▶ **The sinking of the Sleipner offshore platform, Stavanger, Norway, August 12, 1991**, billions of dollars of the loss: **inaccurate finite element analysis, i.e., the method for solving partial differential equations.**  
<https://www.youtube.com/watch?v=eGdiPs4THW8>

## Chapter 1:

# Linear model

- ▶ Definition
- ▶ Systems of linear equations
- ▶ Generalized inverses
- ▶ The Moore-Penrose (MP) inverse
- ▶ Singular value decomposition
- ▶ Principal component analysis
- ▶ MP inverse and solving linear systems

# 1. Linear mathematical models

Given points

$$\{(x_1, y_1), \dots, (x_m, y_m)\}, \quad x_i \in \mathbb{R}^n, \quad y_i \in \mathbb{R},$$

the task is to find a function  $F(x, a_1, \dots, a_p)$  that is a good fit for the data.

The values of the parameters  $a_1, \dots, a_p$  should be chosen so that the equations

$$y_i = F(x, a_1, \dots, a_p), \quad i = 1, \dots, m,$$

are satisfied or, if this is not possible, that the error is as small as possible.

Least squares method: the parameters are determined so that the sum of squared errors

$$\sum_{i=1}^m (F(x_i, a_1, \dots, a_p) - y_i)^2$$

is as small as possible.

The mathematical model is linear, when the function  $F$  is a linear function of the parameters:

$$F(x, a_1, \dots, a_p) = a_1\varphi_1(x) + \varphi_2(x) + \dots + a_p\varphi_p(x),$$

where  $\varphi_1, \varphi_2, \dots, \varphi_p$  are functions of a specific type.

Examples of linear models:

1. linear regression:  $x, y \in \mathbb{R}$ ,  $\varphi_1(x) = 1, \varphi_2(x) = x$ ,
2. polynomial regression:  $x, y \in \mathbb{R}$ ,  $\varphi_1(x) = 1, \dots, \varphi_p(x) = x^{p-1}$ ,
3. multivariate linear regression:  $x = (x_1, \dots, x_n) \in \mathbb{R}^n, y \in \mathbb{R}$ ,

$$\varphi_1(x) = 1, \varphi_2(x) = x_1, \dots, \varphi_n(x) = x_n,$$

4. frequency or spectral analysis:

$$\varphi_1(x) = 1, \varphi_2(x) = \cos \omega x, \varphi_3(x) = \sin \omega x, \varphi_4(x) = \cos 2\omega x, \dots$$

(there can be infinitely many functions  $\varphi_i(x)$  in this case)

Examples of nonlinear models:  $F(x, a, b) = ae^{bx}$  and  $F(x, a, b, c) = \frac{a + bx}{c + x}$ .

Given the data points  $\{(x_1, y_1), \dots, (x_m, y_m)\}$ ,  $x_i \in \mathbb{R}^n$ ,  $y_i \in \mathbb{R}$ , the parameters of a linear model

$$y = a_1\varphi_1(x) + a_2\varphi_2(x) + \dots + a_p\varphi_p(x)$$

should satisfy the system of linear equations

$$y_i = a_1\varphi_1(x_i) + a_2\varphi_2(x_i) + \dots + a_p\varphi_p(x_i), \quad i = 1, \dots, m,$$

or, in a matrix form,

$$\begin{bmatrix} \varphi_1(x_1) & \varphi_2(x_1) & \dots & \varphi_p(x_1) \\ \varphi_1(x_2) & \varphi_2(x_2) & \dots & \varphi_p(x_2) \\ \dots & \dots & \dots & \dots \\ \varphi_1(x_m) & \varphi_2(x_m) & \dots & \varphi_p(x_m) \end{bmatrix} \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_p \end{bmatrix} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_p \end{bmatrix}.$$

## 1.1 Systems of linear equations and generalized inverses

A system of linear equations in the matrix form is given by

$$Ax = b,$$

where

- ▶  $A$  is the matrix of coefficients of order  $m \times n$  where  $m$  is the number of equations and  $n$  is the number of unknowns,
- ▶  $x$  is the vector of unknowns and
- ▶  $b$  is the right side vector.

## Existence of solutions:

Let  $A = [a_1, \dots, a_n]$ , where  $a_i$  are vectors representing the columns of  $A$ .

For any vector  $x = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}$  the product  $Ax$  is a linear combination

$$Ax = \sum_i x_i a_i.$$

The system is **solvable** if and only if the vector  $b$  can be expressed as a linear combination of the columns of  $A$ , that is, it is in the column space  $\mathcal{C}(A)$  of  $A$ , i.e.,  $b \in \mathcal{C}(A)$ .



By adding  $b$  to the columns of  $A$  we obtain the extended matrix of the system

$$[A \mid b] = [a_1, \dots, a_n \mid b],$$

### Theorem

*The system  $Ax = b$  is solvable if and only if the rank of  $A$  equals the rank of the extended matrix  $[A \mid b]$ , i.e.,*

$$\text{rank } A = \text{rank } [A \mid b] =: r.$$

*The solution is unique if the rank of the two matrices equals the number of unknowns, i.e.,  $r = n$ .*

A generic case is the following:

If  $A$  is a square matrix ( $n = m$ ) that has an inverse matrix  $A^{-1}$ , the system has a unique solution

$$x = A^{-1}b.$$

Let  $A \in \mathbb{R}^{n \times n}$  be a square matrix. The following conditions are equivalent and characterize when a matrix  $A$  is invertible or nonsingular:

- ▶ The matrix  $A$  has an inverse.
- ▶ The rank of  $A$  equals  $n$ .
- ▶  $\det(A) \neq 0$ .
- ▶ The null space  $N(A) = \{x : Ax = 0\}$  is trivial.
- ▶ All eigenvalues of  $A$  are nonzero.
- ▶ For each  $b$  the system of equations  $Ax = b$  has precisely one solution.

A square matrix that does not satisfy the above conditions does not have an inverse.

### Example

$$A = \begin{bmatrix} 1 & 0 & 1 \\ 0 & 1 & -1 \\ 1 & 1 & 1 \end{bmatrix}, \quad B = \begin{bmatrix} 1 & 0 & 1 \\ 0 & 1 & -1 \\ 1 & 1 & 0 \end{bmatrix}$$

$A$  is invertible and is of rank 3,  $B$  is not invertible and is of rank 2.

For a rectangular matrix  $A$  of dimension  $m \times n$ ,  $m \neq n$ , its inverse is not defined (at least in the above sense...).

## Definition

A generalized inverse of a matrix  $A \in \mathbb{R}^{n \times m}$  is a matrix  $G \in \mathbb{R}^{m \times n}$  such that

$$AGA = A. \quad (1)$$

## Remark

*Note that the dimension of  $A$  and its generalized inverse are transposed to each other. This is the only way which enables the multiplication  $A \cdot * \cdot A$ .*

## Proposition

*If  $A$  is invertible, it has a unique generalized inverse, which is equal to  $A^{-1}$ .*

## Proof.

Let  $G$  be a generalized inverse of  $A$ , i.e., (1) holds. Multiplying (1) with  $A^{-1}$  from the left and the right side we obtain:

$$\text{Left hand side (LHS): } A^{-1}AGAA^{-1} = IGI = G,$$

$$\text{Right hand side (RHS): } A^{-1}AA^{-1} = IA^{-1} = A^{-1},$$

where  $I$  is the identity matrix. The equality LHS=RHS implies that  $G = A^{-1}$ .

## Theorem

Every matrix  $A \in \mathbb{R}^{n \times m}$  has a generalized inverse.

## Proof.

Let  $r$  be the rank of  $A$ .

**Case 1.**  $\text{rank } A = \text{rank } A_{11}$ , where

$$A = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix}$$

and  $A_{11} \in \mathbb{R}^{r \times r}$ ,  $A_{12} \in \mathbb{R}^{r \times (m-r)}$ ,  $A_{21} \in \mathbb{R}^{(n-r) \times r}$ ,  $A_{22} \in \mathbb{R}^{(n-r) \times (m-r)}$ .

We claim that

$$G = \begin{bmatrix} A_{11}^{-1} & 0 \\ 0 & 0 \end{bmatrix},$$

where 0s denote zero matrices of appropriate sizes, is the generalized inverse of  $A$ . To prove this claim we need to check that

$$AGA = A.$$

$$\begin{aligned}
 AGA &= \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} \begin{bmatrix} A_{11}^{-1} & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} = \begin{bmatrix} I & 0 \\ A_{21}A_{11}^{-1} & 0 \end{bmatrix} \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} \\
 &= \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{21}A_{11}^{-1}A_{12} \end{bmatrix}.
 \end{aligned}$$

For  $AGA$  to be equal to  $A$  we must have

$$A_{21}A_{11}^{-1}A_{12} = A_{22}. \quad (2)$$

It remains to prove (2). Since we are in Case 1, it follows that every column of  $\begin{bmatrix} A_{12} \\ A_{22} \end{bmatrix}$  is in the column space of  $\begin{bmatrix} A_{11} \\ A_{21} \end{bmatrix}$ . Hence, there is a coefficient matrix  $W \in \mathbb{R}^{r \times (m-r)}$  such that

$$\begin{bmatrix} A_{12} \\ A_{22} \end{bmatrix} = \begin{bmatrix} A_{11} \\ A_{21} \end{bmatrix} W = \begin{bmatrix} A_{11}W \\ A_{21}W \end{bmatrix}.$$

We obtain the equations  $A_{11}W = A_{12}$  and  $A_{21}W = A_{22}$ . Since  $A_{11}$  is invertible, we get  $W = A_{11}^{-1}A_{12}$  and hence  $A_{21}A_{11}^{-1}A_{12} = A_{22}$ , which is (2).

**Case 2.** *The upper left  $r \times r$  submatrix of  $A$  is not invertible.*

One way to handle this case is to use permutation matrices  $P$  and  $Q$ , such

that  $PAQ = \begin{bmatrix} \tilde{A}_{11} & \tilde{A}_{12} \\ \tilde{A}_{21} & \tilde{A}_{22} \end{bmatrix}$ ,  $\tilde{A}_{11} \in \mathbb{R}^{r \times r}$  and  $\text{rank } \tilde{A}_{11} = r$ . By Case 1 we

have that the generalized inverse  $(PAQ)^g$  of  $PAQ$  equals to  $\begin{bmatrix} \tilde{A}_{11}^{-1} & 0 \\ 0 & 0 \end{bmatrix}$ .

Thus,

$$(PAQ) \begin{bmatrix} \tilde{A}_{11}^{-1} & 0 \\ 0 & 0 \end{bmatrix} (PAQ) = PAQ. \quad (3)$$

Multiplying (3) from the left by  $P^{-1}$  and from the right by  $Q^{-1}$  we get

$$A \left( Q \begin{bmatrix} \tilde{A}_{11}^{-1} & 0 \\ 0 & 0 \end{bmatrix} P \right) A = A.$$

So,  $Q \begin{bmatrix} \tilde{A}_{11}^{-1} & 0 \\ 0 & 0 \end{bmatrix} P = \left( P^T \begin{bmatrix} (\tilde{A}_{11}^{-1})^T & 0 \\ 0 & 0 \end{bmatrix} Q^T \right)^T$  is a generalized inverse of  $A$ . □

## Algorithm for computing a generalized inverse of $A$

Let  $r$  be the rank of  $A$ .

1. Find any nonsingular submatrix  $B$  in  $A$  of order  $r \times r$ ,
2. in  $A$  substitute
  - ▶ elements of the submatrix  $B$  for corresponding elements of  $(B^{-1})^T$ ,
  - ▶ all other elements with 0,
3. the transpose of the obtained matrix is a generalized inverse  $G$ .

### Example

Compute at least one generalized inverse of

$$A = \begin{bmatrix} 0 & 0 & 2 & 0 \\ 0 & 0 & 1 & 0 \\ 2 & 0 & 1 & 4 \end{bmatrix}.$$



- Note that  $\text{rank } A = 2$ . For  $B$  from the algorithm one of the possibilities is

$$B = \begin{bmatrix} 1 & 0 \\ 1 & 4 \end{bmatrix},$$

i.e., the submatrix in the right lower corner.

- Computing  $B^{-1}$  we get  $B^{-1} = \begin{bmatrix} 1 & 0 \\ -\frac{1}{4} & \frac{1}{4} \end{bmatrix}$  and hence

$$(B^{-1})^T = \begin{bmatrix} 1 & -\frac{1}{4} \\ 0 & \frac{1}{4} \end{bmatrix}.$$

- A generalized inverse of  $A$  is then

$$G = \begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & -\frac{1}{4} \\ 0 & 0 & 0 & \frac{1}{4} \end{bmatrix}^T = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & -\frac{1}{4} & \frac{1}{4} \end{bmatrix}.$$

Generalized inverses of a matrix  $A$  play a similar role as the usual inverse (when it exists) in solving a linear system  $Ax = b$ .

### Theorem

Let  $A \in \mathbb{R}^{n \times m}$  and  $b \in \mathbb{R}^m$ . If the system

$$Ax = b \tag{4}$$

is solvable (that is,  $b \in \mathcal{C}(A)$ ) and  $G$  is a generalized inverse of  $A$ , then

$$x = Gb \tag{5}$$

is a solution of the system (4).

Moreover, all solutions of the system (4) are exactly vectors of the form

$$x_z = Gb + (GA - I)z, \tag{6}$$

where  $z$  varies over all vectors from  $\mathbb{R}^m$ .

## Proof.

We write  $A$  in the column form

$$A = [a_1 \quad a_2 \quad \dots \quad a_m],$$

where  $a_i$  are column vectors of  $A$ . Since the system (4) is solvable, there exist real numbers  $\alpha_1, \dots, \alpha_m \in \mathbb{R}$  such that

$$\sum_{i=1}^m \alpha_i a_i = b. \quad (7)$$

First we will prove that  $Gb$  also solves (4). Multiplying (7) with  $G$  we get

$$Gb = \sum_{i=1}^m \alpha_i Ga_i. \quad (8)$$

Multiplying (9) with  $A$  the left side becomes  $A(Gb)$ , so we have to check that

$$\sum_{i=1}^m \alpha_i AGa_i = b. \quad (9)$$

Since  $G$  is a generalized inverse of  $A$ , we have that  $AGA = A$  or restricting to columns of the left hand side we get

$$AGa_i = a_i \quad \text{for every } i = 1, \dots, m.$$

Plugging this into the left side of (9) we get exactly (7), which holds and proves (9).

For the moreover part we have to prove two facts:

- (i) Any  $x_z$  of the form (6) solves (4).
- (ii) If  $A\tilde{x} = b$ , then  $\tilde{x}$  is of the form  $x_z$  for some  $z \in \mathbb{R}^m$ .

(i) is easy to check:

$$\begin{aligned} Ax_z &= A(Gb + (GA - I)z) = AGb + A(GA - I)z \\ &= b + (AGA - A)z = b. \end{aligned}$$

To prove (ii) note that

$$A(\tilde{x} - Gb) = 0,$$

which implies that

$$\tilde{x} - Gb \in \ker A.$$

It remains to check that

$$\ker A = \{(GA - I)z : z \in \mathbb{R}^m\}. \quad (10)$$

The inclusion ( $\supseteq$ ) of (10) is straightforward:

$$A((GA - I)z) = (AGA - A)z = 0.$$

For the inclusion ( $\subseteq$ ) of (10) we have to notice that any  $v \in \ker A$  is equal to  $(GA - I)z$  for  $z = -v$ :

$$(GA - I)(-v) = -GA v + v = 0 + v = v. \quad \square$$

## Example

Find all solutions of the system

$$Ax = b,$$

where  $A = \begin{bmatrix} 0 & 0 & 2 & 0 \\ 0 & 0 & 1 & 0 \\ 2 & 0 & 1 & 4 \end{bmatrix}$  and  $b = \begin{bmatrix} 2 \\ 1 \\ 4 \end{bmatrix}$ .

- ▶ Recall from the example a few slides above that  $G = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & -\frac{1}{4} & \frac{1}{4} \end{bmatrix}$ .
- ▶ Calculating  $Gb$  and  $GA - I$  we get

$$Gb = \begin{bmatrix} 0 \\ 0 \\ 1 \\ \frac{3}{4} \end{bmatrix} \quad \text{and} \quad A = \begin{bmatrix} -1 & 0 & 0 & 0 \\ 0 & -1 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ \frac{1}{2} & 0 & 0 & 0 \end{bmatrix}.$$

- ▶ Hence,

$$x_z = \begin{bmatrix} -z_1 & -z_2 & 1 & \frac{3}{4} + \frac{1}{2}z_1 \end{bmatrix}^T$$

where  $z_1, z_2$  vary over  $\mathbb{R}$ .

## 1.2 The Moore-Penrose generalized inverse

Among all generalized inverses of a matrix  $A$ , one has especially nice properties.

### Definition

The Moore-Penrose generalized inverse, or shortly the MP inverse of  $A \in \mathbb{R}^{n \times m}$  is any matrix  $A^+ \in \mathbb{R}^{m \times n}$  satisfying the following four conditions:

1.  $A^+$  is a generalized inverse of  $A$ :  $AA^+A = A$ .
2.  $A$  is a generalized inverse of  $A^+$ :  $A^+AA^+ = A^+$ .
3. The square matrix  $AA^+ \in \mathbb{R}^{n \times n}$  is symmetric:  $(AA^+)^T = AA^+$ .
4. The square matrix  $A^+A \in \mathbb{R}^{m \times m}$  is symmetric:  $(A^+A)^T = A^+A$ .

### Remark

*There are two natural questions arising after defining the MP inverse:*

- ▶ *Does every matrix admit a MP inverse? **Yes.***
- ▶ *Is the MP inverse unique? **Yes.***

## Theorem

The MP inverse  $A^+$  of a matrix  $A$  is unique.

## Proof.

Assume that there are two matrices  $M_1$  and  $M_2$  that satisfy the four conditions in the definition of MP inverse of  $A$ . Then,

$$\begin{aligned} AM_1 &= (AM_2A)M_1 && \text{by property (1)} \\ &= (AM_2)(AM_1) = (AM_2)^T(AM_1)^T && \text{by property (3)} \\ &= M_2^T(AM_1A)^T = M_2^T A^T && \text{by property (1)} \\ &= (AM_2)^T = AM_2 && \text{by property (3)} \end{aligned}$$

A similar argument involving properties (2) and (4) shows that

$$M_1A = M_2A,$$

and so

$$M_1 = M_1AM_1 = M_1AM_2 = M_2AM_2 = M_2.$$





## Remark

*Let us assume that  $A^+$  exists (we will shortly prove this fact). Then the following properties are true:*

- ▶ *If  $A$  is a square invertible matrix, then  $A^+ = A^{-1}$ .*
- ▶  $(A^+)^+ = A$ .
- ▶  $(A^T)^+ = (A^+)^T$ .

In the rest of this chapter we will be interested in two obvious questions:

- ▶ How do we compute  $A^+$ ?
- ▶ Why would we want to compute  $A^+$ ?

To answer the first question, we will begin by three special cases.

## Construction of the MP inverse of $A \in \mathbb{R}^{n \times m}$ :

**Case 1:**  $A^T A \in \mathbb{R}^{m \times m}$  is an invertible matrix. (In particular,  $m \leq n$ .)

In this case  $A^+ = (A^T A)^{-1} A^T$ .

To see this, we have to show that the matrix  $(A^T A)^{-1} A^T$  satisfies properties (1) to (4):

1.  $AMA = A(A^T A)^{-1} A^T A = A(A^T A)^{-1} (A^T A) = A$ .
2.  $MAM = (A^T A)^{-1} A^T A (A^T A)^{-1} A^T = (A^T A)^{-1} A^T = M$ .
- 3.

$$\begin{aligned} (AM)^T &= \left( A(A^T A)^{-1} A^T \right)^T = A \left( (A^T A)^{-1} \right)^T A^T = \\ &= A \left( (A^T A)^T \right)^{-1} A^T = A(A^T A)^{-1} A^T = AM. \end{aligned}$$

4. Analogous to the previous fact.

**Case 2:**  $AA^T$  is an invertible matrix. (In particular,  $n \leq m$ .)

In this case  $A^T$  satisfies the condition for Case 1, so  $(A^T)^+ = (AA^T)^{-1}A$ .

Since  $(A^T)^+ = (A^+)^T$  it follows that

$$\begin{aligned}A^+ &= \left( (A^T)^+ \right)^T = \left( (AA^T)^{-1}A \right)^T = A^T \left( (AA^T)^{-1} \right)^T \\ &= A^T \left( (AA^T)^{-T} \right)^{-1} = A^T (AA^T)^{-1}.\end{aligned}$$

Hence,  $A^+ = A^T (AA^T)^{-1}$ .

**Case 3:**  $\Sigma \in \mathbb{R}^{n \times m}$  is a diagonal matrix of the form

$$\Sigma = \begin{bmatrix} \sigma_1 & & & \\ & \sigma_2 & & \\ & & \ddots & \\ & & & \sigma_n \end{bmatrix} \quad \text{or} \quad \tilde{\Sigma} = \begin{bmatrix} \sigma_1 & & & \\ & \sigma_2 & & \\ & & \ddots & \\ & & & \sigma_m \end{bmatrix}.$$

The MP inverse is

$$\Sigma^+ = \begin{bmatrix} \sigma_1^+ & & & \\ & \sigma_2^+ & & \\ & & \ddots & \\ & & & \sigma_n^+ \end{bmatrix} \quad \text{or} \quad \tilde{\Sigma}^+ = \begin{bmatrix} \sigma_1^+ & & & \\ & \sigma_2^+ & & \\ & & \ddots & \\ & & & \sigma_m^+ \end{bmatrix},$$

$$\text{where } \sigma_i^+ = \begin{cases} \frac{1}{\sigma_i}, & \sigma_i \neq 0, \\ 0, & \sigma_i = 0. \end{cases}$$

**Case 4:** A general matrix  $A$ . (using SVD)

Theorem (Singular value decomposition - SVD)

Let  $A \in \mathbb{R}^{n \times m}$  be a matrix. Then it can be expressed as a product

$$A = U\Sigma V^T,$$

where

- ▶  $U \in \mathbb{R}^{n \times n}$  is an orthogonal matrix with left singular vectors  $u_i$  as its columns,
- ▶  $V \in \mathbb{R}^{m \times m}$  is an orthogonal matrix with right singular vectors  $v_i$  as its columns,

▶  $\Sigma = \left[ \begin{array}{ccc|c} \sigma_1 & & & 0 \\ & \ddots & & \vdots \\ & & \sigma_r & 0 \\ \hline & & & 0 \\ & 0 & & 0 \end{array} \right] = \begin{bmatrix} S & 0 \\ 0 & 0 \end{bmatrix} \in \mathbb{R}^{n \times m}$  is a diagonal matrix

with singular values

$$\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > 0$$

on the diagonal.

## Derivations for computing SVD

If  $A = U\Sigma V^T$ , then

$$A^T A = (V\Sigma^T U^T)(U\Sigma V^T) = V\Sigma^T \Sigma V^T = V \begin{bmatrix} S^2 & 0 \\ 0 & 0 \end{bmatrix} V^T \in \mathbb{R}^{m \times m},$$

$$AA^T = (U\Sigma V^T)(U\Sigma V^T)^T = U\Sigma \Sigma^T U^T = U \begin{bmatrix} S^2 & 0 \\ 0 & 0 \end{bmatrix} U^T \in \mathbb{R}^{n \times n}.$$

Let

$$V = [v_1 \quad v_2 \quad \cdots \quad v_m] \quad \text{and} \quad U = [u_1 \quad u_2 \quad \cdots \quad u_n]$$

be the column decompositions of  $V$  and  $U$ .

Let  $e_1, \dots, e_m \in \mathbb{R}^m$  and  $f_1, \dots, f_n \in \mathbb{R}^n$  be the standard coordinate vectors of  $\mathbb{R}^m$  and  $\mathbb{R}^n$ , i.e., the only nonzero component of  $e_i$  (resp.  $f_j$ ) is the  $i$ -th one (resp.  $j$ -th one), which is 1. Then

$$A^T A v_i = V\Sigma^T \Sigma V^T v_i = V\Sigma^T \Sigma e_i = \begin{cases} \sigma_i^2 v_i, & \text{if } i \leq r, \\ 0, & \text{if } i > r, \end{cases}$$

$$AA^T u_j = U\Sigma \Sigma^T U^T u_j = U\Sigma \Sigma^T f_j = \begin{cases} \sigma_j^2 u_j, & \text{if } j \leq r, \\ 0, & \text{if } j > r. \end{cases}$$

Further on,

$$(AA^T)(Av_i) = A(A^T A)v_i = \begin{cases} \sigma_i^2 Av_i, & \text{if } i \leq r, \\ 0, & \text{if } i > r, \end{cases}$$

$$(A^T A)(A^T u_j) = A^T(AA^T)u_j = \begin{cases} \sigma_j^2 A^T u_j, & \text{if } j \leq r, \\ 0, & \text{if } j > r. \end{cases}$$

It follows that:

- ▶  $\Sigma^T \Sigma = \begin{bmatrix} S^2 & 0 \\ 0 & 0 \end{bmatrix} \in \mathbb{R}^{m \times m}$  (resp.  $\Sigma \Sigma^T = \begin{bmatrix} S^2 & 0 \\ 0 & 0 \end{bmatrix} \in \mathbb{R}^{n \times n}$ ) is the diagonal matrix with eigenvalues  $\sigma_i^2$  of  $A^T A$  (resp.  $AA^T$ ) on its diagonal, so the singular values  $\sigma_i$  are their square roots.
- ▶  $V$  has the corresponding eigenvectors (normalized and pairwise orthogonal) of  $A^T A$  as its columns, so the right singular vectors are eigenvectors of  $A^T A$ .
- ▶  $U$  has the corresponding eigenvectors (normalized and pairwise orthogonal) of  $AA^T$  as its columns, so the left singular vectors are eigenvectors of  $AA^T$ .

- $Av_i$  is an eigenvector of  $AA^T$  corresponding to  $\sigma_i^2$  and so

$$u_i = \frac{Av_i}{\|Av_i\|} = \frac{Av_i}{\sigma_i}$$

is a left singular vector corresponding to  $\sigma_i$ , where in the second equality we used that

$$\|Av_i\| = \sqrt{(Av_i)^T(Av_i)} = \sqrt{v_i^T A^T Av_i} = \sqrt{\sigma_i^2 v_i^T v_i} = \sigma_i \|v_i\| = \sigma_i.$$

- $A^T u_j$  is an eigenvector of  $A^T A$  corresponding to  $\sigma_j^2$  and so

$$v_j = \frac{A^T u_j}{\|A^T u_j\|} = \frac{A^T u_j}{\sigma_j}$$

is a right singular vector corresponding to  $\sigma_j$ , where in the second equality we used that

$$\|A^T u_j\| = \sqrt{(A^T u_j)^T(A^T u_j)} = \sqrt{u_j^T A A^T u_j} = \sqrt{\sigma_j^2 u_j^T u_j} = \sigma_j \|u_j\| = \sigma_j.$$



## Algorithm for SVD computation

- ▶ Compute the eigenvalues and an orthonormal basis consisting of eigenvectors of the symmetric matrix  $A^T A$  or  $AA^T$  (depending on which is of them is of smaller size).
- ▶ The singular values of the matrix  $A \in \mathbb{R}^{n \times m}$  are equal to  $\sigma_i = \sqrt{\lambda_i}$ , where  $\lambda_i$  are the nonzero eigenvalues of  $A^T A$  (resp.  $AA^T$ ).
- ▶ The left singular vectors are the corresponding orthonormal eigenvectors of  $AA^T$ .
- ▶ The right singular vector are the corresponding orthonormal eigenvectors of  $A^T A$ .
- ▶ If  $u$  (resp.  $v$ ) is a left (resp. right) singular vector corresponding to the singular value  $\sigma_i$ , then  $v = A^T u$  (resp.  $u = Av$ ) is a right (resp. left) singular vector corresponding to the same singular value.
- ▶ The remaining columns of  $U$  (resp.  $V$ ) consist of an orthonormal basis of the kernel (i.e., the eigenspace of  $\lambda = 0$ ) of  $AA^T$  (resp.  $A^T A$ ).

## General algorithm for computation of $A^+$ (long version)

1. For  $A^T A$  compute its eigenvalues

$$\lambda_1 \geq \lambda_2 \geq \dots, \geq \lambda_r > \lambda_{r+1} = \dots = \lambda_m = 0$$

and the corresponding orthonormal eigenvectors

$$v_1, \dots, v_r, v_{r+1}, \dots, v_m,$$

and form the matrices

$$\Sigma = \text{diag}(\sqrt{\lambda_1}, \dots, \sqrt{\lambda_m}) \in \mathbb{R}^{n \times m},$$

$$V_1 = [v_1 \ \dots \ v_r], \quad V_2 = [v_{r+1} \ \dots \ v_m] \quad \text{and} \quad V = [V_1 \ V_2].$$

2. Let

$$u_1 = \frac{Av_1}{\sigma_1}, \quad u_2 = \frac{Av_2}{\sigma_2}, \quad \dots, \quad u_r = \frac{Av_r}{\sigma_r},$$

and  $u_{r+1}, \dots, u_n$  vectors, such that  $\{u_1, \dots, u_n\}$  is an orthonormal basis for  $\mathbb{R}^n$ . Form the matrices

$$U_1 = [u_1 \ \dots \ u_r], \quad U_2 = [u_{r+1} \ \dots \ u_n] \quad \text{and} \quad U = [U_1 \ U_2].$$

3. Then

$$A^+ = V\Sigma^+U^T.$$

## General algorithm for computation of $A^+$ (short version)

1. For  $A^T A$  compute its **nonzero** eigenvalues

$$\lambda_1 \geq \lambda_2 \geq \dots, \geq \lambda_r > 0$$

and the corresponding orthonormal eigenvectors

$$v_1, \dots, v_r,$$

and form the matrices

$$S = \text{diag}(\sqrt{\lambda_1}, \dots, \sqrt{\lambda_r}) \in \mathbb{R}^{r \times r},$$

$$V_1 = [v_1 \ \dots \ v_r] \in \mathbb{R}^{m \times r}.$$

2. Put the vectors

$$u_1 = \frac{Av_1}{\sigma_1}, \quad u_2 = \frac{Av_2}{\sigma_2}, \quad \dots, \quad u_r = \frac{Av_r}{\sigma_r}$$

in the matrix

$$U_1 = [u_1 \ \dots \ u_r].$$

3. Then

$$A^+ = V_1 \Sigma^+ U_1^T.$$

## Correctness of the computation of $A^+$

**Step 1.**  $V\Sigma^+U^T$  is equal to  $A^+$ .

(i)  $AA^+A = A$ :

$$\begin{aligned}AA^+A &= (U\Sigma V^T)(V\Sigma^+U^T)(U\Sigma V^T) = U\Sigma(V^TV)\Sigma^+(U^TU)\Sigma V^T \\ &= U\Sigma\Sigma^+\Sigma V^T = U\Sigma V^T = A.\end{aligned}$$

(ii)  $A^+AA^+ = A^+$ : Analogous to (i).

(iii)  $(AA^+)^T = AA^+$ :

$$\begin{aligned}(AA^+)^T &= \left( (U\Sigma V^T)(V\Sigma^+U^T) \right)^T = \left( U\Sigma\Sigma^+U^T \right)^T \\ &= \left( U \begin{bmatrix} I_r & 0 \\ 0 & 0 \end{bmatrix} U^T \right)^T = U \begin{bmatrix} I_r & 0 \\ 0 & 0 \end{bmatrix} U^T \\ &= (U\Sigma V^T)(V\Sigma^+U^T) = A^+.\end{aligned}$$

(iv)  $(A^+A)^T = A^+A$ : Analogous to (iii).

**Step 2.**  $V\Sigma^+U^T$  is equal to  $V_1\Sigma^+U_1^T$ .

$$V\Sigma U^T = [V_1 \quad V_2] \begin{bmatrix} S & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} U_1^T \\ U_2^T \end{bmatrix} = [V_1 S \quad 0] \begin{bmatrix} U_1^T \\ U_2^T \end{bmatrix} = V_1 S U_1^T.$$

## Example

Compute the SVD and  $A^+$  of the matrix  $A = \begin{bmatrix} 3 & 2 & 2 \\ 2 & 3 & -2 \end{bmatrix}$ .

- ▶  $AA^T = \begin{bmatrix} 17 & 8 \\ 8 & 17 \end{bmatrix}$  has eigenvalues 25 and 9.
- ▶ The eigenvectors of  $AA^T$  corresponding to the eigenvalues 25, 9 are

$$u_1 = \left[ \frac{1}{\sqrt{2}} \quad \frac{1}{\sqrt{2}} \right]^T, \quad u_2 = \left[ \frac{1}{\sqrt{2}} \quad -\frac{1}{\sqrt{2}} \right]^T.$$

- ▶ The left singular vectors of  $A$  are

$$v_1 = \frac{A^T u_1}{\sigma_1} = \left[ \frac{1}{\sqrt{2}} \quad \frac{1}{\sqrt{2}} \quad 0 \right]^T, \quad v_2 = \frac{A^T u_2}{\sigma_2} = \left[ \frac{1}{3\sqrt{2}} \quad -\frac{1}{3\sqrt{2}} \quad \frac{4}{3\sqrt{2}} \right]^T.$$

$$v_3 = v_1 \times v_2 = \left[ \frac{2}{\sqrt{3}} \quad -\frac{2}{3} \quad -\frac{1}{3} \right]^T.$$

$$A = U\Sigma V^T = \begin{bmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \end{bmatrix} \begin{bmatrix} 5 & 0 & 0 \\ 0 & 3 & 0 \end{bmatrix} \begin{bmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} & 0 \\ \frac{1}{3\sqrt{2}} & -\frac{1}{3\sqrt{2}} & \frac{4}{3\sqrt{2}} \\ \frac{2}{\sqrt{3}} & -\frac{2}{3} & -\frac{1}{3} \end{bmatrix}.$$

$$\begin{aligned} A^+ &= V\Sigma^+U^T = \begin{bmatrix} \frac{1}{\sqrt{2}} & \frac{1}{3\sqrt{2}} & \frac{2}{\sqrt{3}} \\ \frac{1}{\sqrt{2}} & -\frac{1}{3\sqrt{2}} & -\frac{2}{3} \\ 0 & \frac{4}{3\sqrt{2}} & -\frac{1}{3} \end{bmatrix} \begin{bmatrix} \frac{1}{5} & 0 \\ 0 & \frac{1}{3} \\ 0 & 0 \end{bmatrix} \begin{bmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \end{bmatrix} \\ &= \begin{bmatrix} \frac{7}{45} & \frac{2}{45} \\ \frac{2}{45} & \frac{7}{45} \\ \frac{2}{9} & -\frac{2}{9} \end{bmatrix}. \end{aligned}$$

## 1.3 The MP inverse and systems of linear equations

Let  $A \in \mathbb{R}^{n \times m}$ , where  $m > n$ . A system of equations  $Ax = b$  that has more variables than constraints. Typically such system has infinitely many solutions, but it may happen that it has no solutions. We call such system an underdetermined system.

### Theorem

1. *An underdetermined system of linear equations*

$$Ax = b \tag{11}$$

*is solvable if and only if  $AA^+b = b$ .*

2. *If there are infinitely many solutions, the solution  $A^+b$  is the one with the smallest norm, i.e.,*

$$\|A^+b\| = \min \{\|x\| : Ax = b\}.$$

*Moreover, it is the unique solution of smallest norm.*

## Proof of Theorem.

We already know that  $Ax = b$  is solvable iff  $Gb$  is a solution, where  $G$  is any generalized inverse of  $A$ . Since  $A^+$  is one of the generalized inverses, this proves the first part of the theorem.

To prove the second part of the theorem, first recall that all the solutions of the system are precisely the set

$$\{A^+b + (A^+A - I)z : z \in \mathbb{R}^m\}.$$

So we have to prove that for every  $z \in \mathbb{R}^m$ ,

$$\|A^+b\| \leq \|A^+b + (A^+A - I)z\|.$$

We have that:

$$\begin{aligned} \|A^+b + (A^+A - I)z\|^2 &= \\ &= (A^+b + (A^+A - I)z)^T (A^+b + (A^+A - I)z) \\ &= (A^+b)^T (A^+b) + 2(A^+b)^T (A^+A - I)z + ((A^+A - I)z)^T ((A^+A - I)z) \\ &= \|A^+b\|^2 + 2(A^+b)^T (A^+A - I)z + \|(A^+A - I)z\|^2 \end{aligned}$$



Now,

$$\begin{aligned}(A^+b)^T (A^+A - I)z &= b^T (A^+)^T (A^+A - I)z \\ &= b^T (A^+)^T (A^+A)^T z - b^T (A^+)^T z \\ &= b^T ((A^+A)A^+)^T z - b^T (A^+)^T z \\ &= b^T (A^+AA^+)^T z - b^T (A^+)^T z \\ &= b^T (A^+)^T z - b^T (A^+)^T z = 0,\end{aligned}$$

where we used the fact  $(A^+A)^T = A^+A$  in the second equality.

Thus,

$$\|A^+b + (A^+A - I)z\|^2 = \|A^+b\|^2 + \|(A^+A - I)z\|^2 \geq \|A^+b\|^2,$$

with the equality iff  $(A^+A - I)z = 0$ . This proves the second part of the theorem. □

## Example

- ▶ The solutions of the underdetermined system  $x + y = 1$  geometrically represent an affine line. Matricially,  $A = \begin{bmatrix} 1 & 1 \end{bmatrix}$ ,  $b = 1$ . Hence,  $A^+b = A^+1$  is the point on the line, which is the nearest to the origin. Thus, the vector of this point is perpendicular to the line.
- ▶ The solutions of the underdetermined system  $x + 2y + 3z = 5$  geometrically represent an affine hyperplane. Matricially,  $A = \begin{bmatrix} 1 & 2 & 3 \end{bmatrix}$ ,  $b = 5$ . Hence,  $A^+b = A^+5$  is the point on the hyperplane, which is the nearest to the origin. Thus, the vector of this point is normal to the hyperplane.
- ▶ The solutions of the underdetermined system  $x + y + z = 1$  and  $x + 2y + 3z = 5$  geometrically represent an affine line in  $\mathbb{R}^3$ . Matricially,  $A = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 2 & 3 \end{bmatrix}$ ,  $b = \begin{bmatrix} 1 \\ 5 \end{bmatrix}$ . Hence,  $A^+b$  is the point on the line, which is the nearest to the origin. Thus, the vector of this point is perpendicular to the line.

## Example

Find the point on the plane  $3x + y + z = 2$  closest to the origin.

- ▶ In this case,

$$A = \begin{bmatrix} 3 & 1 & 1 \end{bmatrix} \quad \text{and} \quad b = [2].$$

- ▶ We have that  $AA^T = [11]$  and hence its only eigenvalue is  $\lambda = 11$  with eigenvector  $u = [1]$ , implying that

$$U = [1] \quad \text{and} \quad \Sigma = \begin{bmatrix} \sqrt{11} & 0 & 0 \end{bmatrix}.$$

- ▶ Hence,

$$v_1 = \frac{A^T u}{\|A^T u\|} = \frac{A^T u}{\sigma_1} = \frac{1}{\sqrt{11}} \begin{bmatrix} 3 & 1 & 1 \end{bmatrix}^T.$$



$$A^+ = V\Sigma^+U^T = \frac{1}{\sqrt{11}} \begin{bmatrix} 3 \\ 1 \\ 1 \end{bmatrix} \frac{1}{\sqrt{11}} [1] = \begin{bmatrix} \frac{3}{11} \\ \frac{1}{11} \\ \frac{1}{11} \end{bmatrix}.$$



$$x^+ = A^+ b = \begin{bmatrix} \frac{6}{11} & \frac{2}{11} & \frac{2}{11} \end{bmatrix}^T.$$

## Overdetermined systems

Let  $A \in \mathbb{R}^{n \times m}$ , where  $n > m$ . This system is called overdetermined, since here are more constraints than variables. Such a system typically has no solutions, but it might have one or even infinitely many solutions.

Least squares approximation problem: if the system  $Ax = b$  has no solutions, then a best fit for the solution is a vector  $x$  such that the error  $\|Ax - b\|$  or, equivalently in the row decomposition

$$A = \begin{bmatrix} \alpha_1 \\ \vdots \\ \alpha_n \end{bmatrix},$$

its square

$$\|Ax - b\|^2 = \sum_{i=1}^n (\alpha_i x - b_i)^2,$$

is the smallest possible.

## Theorem

If the system  $Ax = b$  has no solutions, then

$$x^+ = A^+ b$$

is the solution to the least squares approximation problem:

$$\min\{\|Ax - b\| : x \in \mathbb{R}^n\}. \quad (12)$$

Moreover, if  $\text{rank } A = m$ , then (12) has a unique solution. If  $\text{rank } A < m$ , then  $x^+$  has the smallest second norm  $\|x^+\|_2$  among all solutions to (12).

## Proof.

Let  $A = U\Sigma V^T$  be the SVD decomposition of  $A$ . We have that

$$\|Ax - b\| = \|U\Sigma V^T x - b\| = \|\Sigma V^T x - U^T b\|,$$

where we used that

$$\|U^T v\| = \|v\|$$

in the second equality (which holds since  $U^T$  is an orthogonal matrix).

Let

$$\Sigma = \begin{bmatrix} S & 0 \\ 0 & 0 \end{bmatrix}, \quad U = [U_1 \quad U_2], \quad V = [V_1 \quad V_2], \quad \text{where}$$

$S \in \mathbb{R}^{r \times r}$ ,  $U_1 \in \mathbb{R}^{n \times r}$ ,  $U_2 \in \mathbb{R}^{n \times (n-r)}$ ,  $V_1 \in \mathbb{R}^{m \times r}$ ,  $V_2 \in \mathbb{R}^{m \times (m-r)}$ . Thus,

$$\begin{aligned} \|\Sigma V^T x - U^T b\| &= \left\| \begin{bmatrix} S & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} V_1^T \\ V_2^T \end{bmatrix} x - \begin{bmatrix} U_1^T \\ U_2^T \end{bmatrix} b \right\| \\ &= \left\| \begin{bmatrix} S V_1^T x - U_1^T b \\ U_2^T b \end{bmatrix} \right\|. \end{aligned}$$

But this norm is minimal iff

$$S V_1^T x - U_1^T b = 0$$

or equivalently

$$V_1^T x = S^{-1} U_1^T b. \tag{13}$$

Further on,

$$V^T V = \begin{bmatrix} V_1^T V_1 & V_1^T V_2 \\ V_2^T V_1 & V_2^T V_2 \end{bmatrix} = I_n,$$

implies that  $V_1^T V_1 = I_r$  and  $V_2^T V_1 = 0$ , where  $I_k$  stands for the  $k \times k$  identity matrix.

If  $\text{rank } A = m$ , then  $V_1 \in \mathbb{R}^{m \times m}$  is invertible with the inverse  $V_1^T$  and hence,

$$V_1 S^{-1} U_1^T b = A^+ b$$

is the unique solution to (12).

If  $r = \text{rank } A < m$ , then all  $x$  which solve (13) are of the form  $A_1^+ b + z$ , for  $z \in \ker V_1^T$ . Since  $\ker V_1^T = \text{im } V_2$  and  $V_2^T V_1 = 0$ , it follows that the norm of  $A_1^+ b + z$  is minimal for  $z = 0$ . □

## Remark

*The closest vector to  $b$  in the column space  $C(A) = \{Ax : x \in \mathbb{R}^m\}$  of  $A$  is the orthogonal projection of  $b$  onto  $C(A)$ . It follows that  $A^+ b$  is this projection. Equivalently,  $b - (A^+ b)$  is orthogonal to any vector  $Ax$ ,  $x \in \mathbb{R}^m$ , which can be proved also directly.*

## Example

Given points  $\{(x_1, y_1), \dots, (x_n, y_n)\}$  in the plane, we are looking for the line  $ax + b = y$  which is the least squares best fit.

If  $n > 2$ , we obtain an overdetermined system

$$\begin{bmatrix} x_1 & 1 \\ \vdots & \\ x_n & 1 \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}.$$

The solution of the least squares approximation problem is given by

$$\begin{bmatrix} a \\ b \end{bmatrix} = A^+ \begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix}.$$

The line  $y = ax + b$  in the [regression line](#).



## 1.4 Principal component analysis (PCA)

SVD is an essential tool for the [PCA](#), which is a very well-known and efficient method for **data compression**, **dimension reduction**, ...

Due to its importance in different fields, it has many other names: discrete Karhunen-Loève transform (KLT), Hotelling transform, empirical orthogonal functions (EOF), ...

Let  $\{X_1, \dots, X_m\}$  be a sample of vectors from  $\mathbb{R}^n$ .

In applications, often  $m \ll n$ , where  $n$  is very large, for example,  $X_1, \dots, X_m$  can be

- ▶ vectors of gene expressions in  $m$  tissue samples or
- ▶ vectors of grayscale in images
- ▶ bag of words vectors, with components corresponding to the numbers of certain words from some dictionary in specific texts, ... ,

or  $n \ll m$  for example if the data represents a point cloud in a low dimensional space  $\mathbb{R}^n$  (for example in the plane).

We will assume that  $m \ll n$ . Also assume that the data is [centralized](#), i.e., the centroid is in the origin

$$\mu = \frac{1}{m} \sum_{i=1}^m X_i = 0 \in \mathbb{R}^n.$$

If not, we subtract  $\mu$  from all vectors in the data set.

A [matrix norm](#)  $\|\cdot\| : \mathbb{R}^{n \times m} \rightarrow \mathbb{R}$  is a function, which generalizes the notion of the absolute value for numbers to matrices. It is used to measure a distance between matrices. In contrast with the absolute value, which is unique up to multiplication with a positive constant, there are many different matrix norms.

Two important matrix norms are the following:

1. [Spectral norm](#)  $\|\cdot\|_2$ :

$$\|A\|_2 := \max_{\|x\|_2=1} \|Ax\|_2 = \max_{j=1, \dots, \min(n,m)} \sigma_j(A).$$

2. [Frobenius norm](#)  $\|\cdot\|_F$ :

$$\|A\|_F := \sqrt{\sum_{i,j} a_{i,j}^2} = \sqrt{\sum_{j=1, \dots, \min(n,m)} \sigma_j(A)^2}.$$

Let

$$X = [X_1 \quad X_2 \quad \cdots \quad X_m]^T$$

be the matrix of dimension  $m \times n$  with data in the rows.

Let  $X^T X \in \mathbb{R}^{m \times m}$  and  $XX^T \in \mathbb{R}^{n \times n}$  be the covariance matrices of the data.

- ▶ The principal values of the data set  $\{X_1, \dots, X_r\}$  are the nonzero eigenvalues  $\lambda_i = \sigma_i^2$  of the covariance matrices (where  $\sigma_i$  are the singular values of  $X$ ).
- ▶ The principal directions in  $\mathbb{R}^n$  are corresponding eigenvectors  $v_1, \dots, v_r$ , i.e. the columns of the matrix  $V$  from the SVD of  $X$ . The remaining columns of  $V$  (i.e. the eigenvectors corresponding to 0) form a basis of the null space of  $X$ .
- ▶ The first column  $v_1$ , the first principal direction, corresponds to the direction in  $\mathbb{R}^n$  with the largest variance in the data  $X_i$ , that is, the most informative direction for the data set, the second the second most important, ...
- ▶ The principal directions in  $\mathbb{R}^m$  are the columns  $u_1, \dots, u_r$  of the matrix  $U$  and represent the coefficients in the linear decomposition of the vectors  $X_1, \dots, X_m$  along the orthonormal basis  $v_1, \dots, v_n$  of  $\mathbb{R}^n$ .

PCA provides a linear dimension reduction method based on a projection of the data from the space  $\mathbb{R}^n$  into a lower dimensional subspace spanned by the first few principal vectors  $v_1, \dots, v_k$  in  $\mathbb{R}^n$ .

The idea is to approximate

$$X_i = \sigma_1 u_{1,i} v_1 + \dots + \sigma_m u_{m,i} v_m \cong \sigma_1 u_{1,i} v_1 + \dots + \sigma_k u_{k,i} v_k$$

with the first  $k$  most informative directions in  $\mathbb{R}^n$  and suppress the last  $m - k$ .

PCA has the following amazing property:

### Theorem

*Among all possible projections of  $p: \mathbb{R}^n \rightarrow \mathbb{R}^k$  onto a  $k$ -dimensional subspace, PCA provides the best in the sense that the errors*

$$\|X - p(X)\|_F^2 \quad \text{and} \quad \|X - p(X)\|_2^2,$$

*where  $p(X) = [p(X_1) \ \dots \ p(X_m)]^T$ , are the smallest possible.*

## Chapter 3:

# Nonlinear models

- ▶ Definition and examples
- ▶ Systems of nonlinear equations
- ▶ Vector functions of vector variables
  - ▶ Derivative and Jacobian matrix
  - ▶ Linear approximation
- ▶ Newton's method for square systems
  - ▶ Univariate case: Tangent method
  - ▶ Use in optimization
- ▶ Gauss-Newton's method for rectangular systems

### 3. Nonlinear models

#### General formulation

Given is a sample of points  $\{(x_1, y_1), \dots, (x_m, y_m)\}$ ,  $x_i \in \mathbb{R}^n$ ,  $y_i \in \mathbb{R}$ .

The mathematical model is nonlinear if the function

$$y = F(x, a_1, \dots, a_p) \quad (14)$$

is a nonlinear function of the parameters  $a_j$ . This means it cannot be written in the form

$$y = a_1 f_1(x) + a_2 f_2(x) + \dots + a_p f_p(x),$$

where each  $f_j : \mathbb{R}^n \rightarrow \mathbb{R}$  is some function.

Plugging each data points into (14) we obtain a **system of nonlinear equations**

$$\begin{aligned} y_1 &= F(x_1, a_1, \dots, a_p), \\ &\vdots \\ y_m &= F(x_m, a_1, \dots, a_p), \end{aligned} \quad (15)$$

in the parameters  $a_1, \dots, a_p \in \mathbb{R}$ .

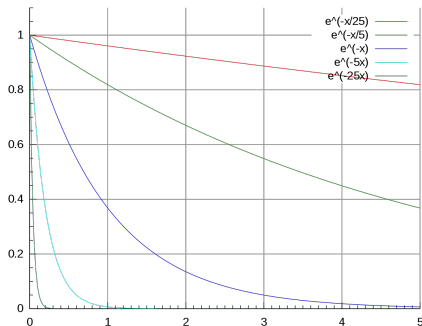
## Examples

1. Exponential decay or growth:  $F(x, a, k) = ae^{kx}$ ,  $a$  and  $k$  are parameters.

A quantity  $y$  changes at a rate proportional to its current value, which can be described by the differential equation

$$\frac{dy}{dx} = ky.$$

The solution to this equation (obtained by the use of separation of variables) is  $y = F(x, a, k)$ .



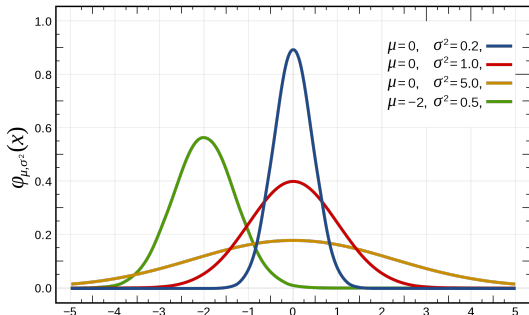
## Examples

2. Gaussian model:  $F(x, a, b, c) = ae^{-\left(\frac{x-b}{c}\right)^2}$ ,  $a, b, c \in \mathbb{R}$  parameters.

$a$  is the value of the maximum obtained at  $x = b$  and  $c$  determines the width of the curve.

It is used in statistics to describe the normal distribution, but also in signal and image processing.

In statistics  $a = \frac{1}{\sigma\sqrt{2\pi}}$ ,  $b = \mu$ ,  $c = \sqrt{2}\sigma$ , where  $\mu$ ,  $\sigma$  are the expected value and the standard deviation of a normally distributed random variable.





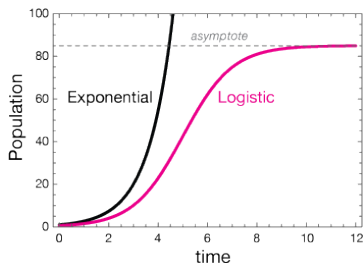
## Examples

### 3. Logistic model: $F(x, a, b, k) = \frac{a}{(1+be^{-kx})}$ , $k > 0$

The logistic function was devised as a model of population size by adjusting the exponential model which also considers the saturation of the environment, hence the growth first changes to linear and then stops.

The logistic function  $F(x, a, b, k)$  is a solution of the first order non-linear differential equation

$$\frac{dy(x)}{dx} = ky(x) \left( 1 - \frac{y(x)}{a} \right).$$



### Examples

4. In the area around a radiotelescope the use of microwave ovens is forbidden, since the radiation interferes with the telescope. We are looking for the location  $(a, b)$  of a microwave oven that is causing problems.

The radiation intensity decreases with the distance  $r$  from the source according to  $u(r) = \frac{\alpha}{1+r}$ . In cartesian coordinates:

$$u(x, y) = \frac{\alpha}{1 + \sqrt{(x - a)^2 + (y - b)^2}},$$

where  $(a, b)$  is a position of the microwave.

**Task:** Find the position of the microwave, if the measured values of the signal at three locations are  $u(0, 0) = 0.27$ ,  $u(1, 1) = 0.36$  in  $u(0, 2) = 0.3$ .

This gives the following system of equations for the parameters  $\alpha, a, b$ :

$$\begin{aligned}\frac{\alpha}{1 + \sqrt{a^2 + b^2}} &= 0.27 \\ \frac{\alpha}{1 + \sqrt{(1 - a)^2 + (1 - b)^2}} &= 0.36 \\ \frac{\alpha}{1 + \sqrt{a^2 + (2 - b)^2}} &= 0.3\end{aligned}$$

## An equivalent, more convenient formulation of the nonlinear system

- ▶ Our goal is to fit the data points

$$\{(x_1, y_1), \dots, (x_m, y_m)\}, \quad x_i \in \mathbb{R}^n, \quad y_i \in \mathbb{R}.$$

- ▶ We choose a fitting function

$$F(x, a_1, \dots, a_p)$$

which depends on the unknown parameters  $a_1, \dots, a_p$ .

- ▶ Equivalent formulation of the system (15) (which will be more suitable for solving with numerical algorithms) is:

1. For  $i = 1, \dots, m$  define the functions

$$g_i : \mathbb{R}^p \rightarrow \mathbb{R} \quad \text{by the rule} \quad g_i(a_1, \dots, a_p) = y_i - F(x_i, a_1, \dots, a_p).$$

2. Solve or approximate the following system by the least squares method

$$\begin{aligned} g_1(a_1, \dots, a_p) &= 0, \\ &\vdots \\ g_m(a_1, \dots, a_p) &= 0. \end{aligned} \tag{16}$$

In a compact way (16) can be expressed by introducing a vector function

$$G: \mathbb{R}^p \rightarrow \mathbb{R}^m, \quad G(a_1, \dots, a_p) = (g_1(a_1, \dots, a_p), \dots, g_m(a_1, \dots, a_p)), \quad (17)$$

and search for the tuples  $(a_1, \dots, a_p)$  that solve the system (or minimize the norm of the left-hand side)

$$G(a_1, \dots, a_p) = (0, \dots, 0). \quad (18)$$

## Remark

*Solving (18) is a difficult problem. Even if the exact solution exists, it is not easy (or even impossible) to compute. For example, there does not even exist an analytic formula to determine roots of a general polynomial of degree 5 or more.*

But we will learn some numerical algorithms to *approximate* the solutions of (18).

## 3.1 Vector functions of a vector variable

Necessary terminology to achieve our plan

$G$  from (17) is an example of

- ▶ a vector function: since it maps into  $\mathbb{R}^m$ , where  $m$  might be bigger than 1.
- ▶ a vector variable: since it maps from  $\mathbb{R}^p$ , where  $p$  might be bigger than 1.

Remark

- ▶ *If  $m = 1$  and  $p > 1$ , then  $G$  is a usual multivariate function.*
- ▶ *If  $m = 1$  and  $p = 1$ , then  $G$  is a usual (univariate) function.*

For easier reference in the continuation we call  $g_1, \dots, g_m$  from (17) the component (or coordinate) functions of  $G$ .

### Examples

1. A linear vector function  $G : \mathbb{R}^n \rightarrow \mathbb{R}^m$  is such that all the component functions  $g_i$  are linear:

$$g_i(x_1, \dots, x_n) = a_{i1} \cdot x_1 + a_{i2} \cdot x_2 + \dots + a_{in} \cdot x_n, \quad \text{where } a_{ij} \in \mathbb{R}. \quad (19)$$

In this case

$$G(x) = Ax,$$

where

$$A = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \dots & a_{mn} \end{bmatrix}.$$

2. Adding constants  $b_i \in \mathbb{R}$  to the left side of (19) we get the definition of an affine linear vector function,

$$g_i(x_1, \dots, x_n) = a_{i1}x_1 + a_{i2}x_2 + \dots + a_{in}x_n + b_i,$$

and then

$$G(x) = Ax + b, \quad \text{where } b = [ b_1 \quad b_2 \quad \dots \quad b_n ]^T.$$

3. Most of the (vector) functions are nonlinear, e.g.,

$$f: \mathbb{R}^3 \rightarrow \mathbb{R}^2, \quad f(x, y, z) = (x^2 + y^2 + z^2 - 1, x + y + z),$$

$$g: \mathbb{R}^2 \rightarrow \mathbb{R}^3, \quad g(z, w) = (zw, \cos z + w^2 - 2, e^{2z}),$$

$$h: \mathbb{R} \rightarrow \mathbb{R}^2, \quad h(t) = (t + 3, e^{-3t}).$$

Derivative of a vector function - is needed in the algorithms we will use

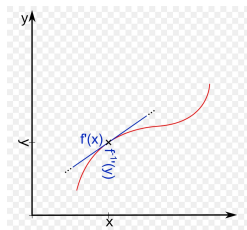
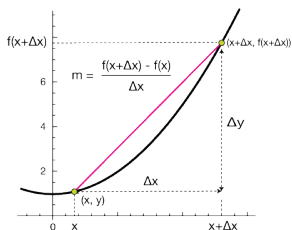
The derivative of a vector function  $F : \mathbb{R}^n \rightarrow \mathbb{R}^m$  in the point

$$a := (a_1, \dots, a_n) \in \mathbb{R}^n$$

is called the Jacobian matrix of  $F$  in  $a$ :

$$J_F(a) = DF(a) = \begin{bmatrix} \frac{\partial f_1}{\partial x_1}(a) & \cdots & \frac{\partial f_1}{\partial x_n}(a) \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1}(a) & \cdots & \frac{\partial f_m}{\partial x_n}(a) \end{bmatrix}.$$

► If  $n = m = 1$ , the  $Df(x) = f'(x)$  is the usual derivative.



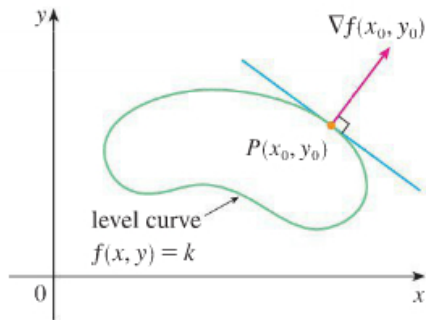


## Derivative - continued

- ▶ For general  $n$  and  $m = 1$ ,  $f$  is a function of  $n$  variables and

$$Df(x) = \text{grad } f(x)$$

is its gradient.



- ▶ For general  $m$  and  $n$ ,  $Df(x) = \begin{bmatrix} \text{grad } f_1 \\ \vdots \\ \text{grad } f_m \end{bmatrix}$  is a vector of gradients of component functions.

## Examples

1. For an affine linear function  $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ , given by  $f(x) = Ax + b$ , it is easy to check that

$$Df(x) = A.$$

2. For a vector function  $f : \mathbb{R}^3 \rightarrow \mathbb{R}^2$ , given by

$$f(x, y, z) = (x^2 + y^2 + z^2 - 1, x + y + z),$$

then

$$Df(x) = \begin{bmatrix} 2x & 2y & 2z \\ 1 & 1 & 1 \end{bmatrix}.$$

### Application of the derivative - linear approximation

A linear approximation of the vector function  $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$  at the point  $a \in \mathbb{R}^n$  is the affine linear function

$$L_a : \mathbb{R}^n \rightarrow \mathbb{R}^m, \quad L_a(x) = Ax + b$$

that satisfies the following conditions:

1. It has the **same value** as  $f$  in  $a$ :  $L_a(a) = f(a)$ .
2. It has the **same derivative** as  $f$  at  $a$ :  $DL_a(a) = Df(a)$ .

It is easy to check that

$$L_a(x) = f(a) + Df(a)(x - a).$$

►  $n = m = 1$ :

$$L_a(x) = f(a) + f'(a)(x - a)$$

The graph  $y = L_a(x)$  is the tangent to the graph  $y = f(x)$  at the point  $a$ .

## Application of the derivative - linear approximation continued

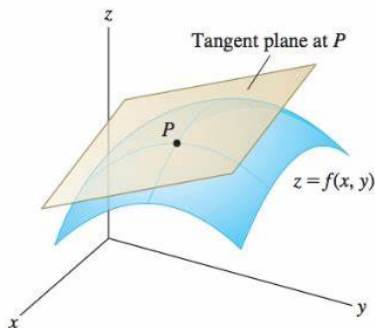
- ▶ If  $n = 2$  and  $m = 1$ , then

$$L_{(a,b)}(x, y) = f(a, b) + \text{grad}f(a, b) \begin{bmatrix} x - a \\ y - b \end{bmatrix}.$$

The graph

$$z = L_{(a,b)}(x, y)$$

is the tangent plane to the surface  $z = f(x, y)$  at the point  $(a, b)$ .



**Example**

The linear approximation of the function

$$f : \mathbb{R}^3 \rightarrow \mathbb{R}^2, \quad f(x, y, z) = (x^2 + y^2 + z^2 - 1, x + y + z)$$

at  $a = (1, -1, 1)$  is the affine linear function

$$\begin{aligned} L_a(x, y, z) &= f(1, -1, 1) + Df(1, -1, 1) \begin{bmatrix} x - 1 \\ y + 1 \\ z - 1 \end{bmatrix} \\ &= \begin{bmatrix} 2 \\ 1 \end{bmatrix} + \begin{bmatrix} 2 & -2 & 2 \\ 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} x - 1 \\ y + 1 \\ z - 1 \end{bmatrix} \\ &= \begin{bmatrix} 2 + 2(x - 1) - 2(y + 1) + 2(z - 1) \\ 1 + (x - 1) + (y + 1) + (z - 1) \end{bmatrix} \\ &= \begin{bmatrix} 2 & -2 & 2 \\ 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} x \\ y \\ z \end{bmatrix} + \begin{bmatrix} -4 \\ 0 \end{bmatrix}. \end{aligned}$$

## 3.2 Solving systems of nonlinear equations

Let  $f : D \rightarrow \mathbb{R}^m$  be a vector function, defined on some set  $D \subset \mathbb{R}^n$ .

We will study the [Gauss-Newton method](#) to solve the system  $f(x) = 0$  in terms of least squares. This is one of the numerical methods for searching approximate solution of this system. It is based on linear approximations of  $f$ .

### Newton's method for $n = m = 1$

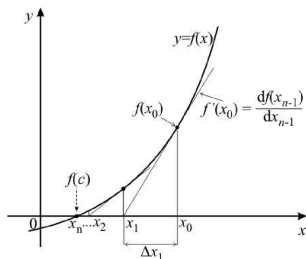
We are searching zeroes of the function  $f : D \rightarrow \mathbb{R}$ ,  $D \subseteq \mathbb{R}$ , i.e., we are solving  $f(x) = 0$ .

### Newton's or tangent method:

We construct a recursive sequence with:

- ▶  $x_0$  is an initial term,
- ▶  $x_{k+1}$  is a solution of

$$L_{x_k}(x) = f(x_k) + f'(x_k)(x - x_k) = 0, \text{ so } x_{k+1} = x_k - \frac{f(x_k)}{f'(x_k)}.$$



## Theorem

The sequence  $x_j$  converges to a solution  $\alpha$ ,  $f(\alpha) = 0$ , if:

- (1)  $0 \neq |f'(x)|$  for all  $x \in I$ , where  $I$  is some interval containing  $\alpha$ ,
- (2)  $x_0$  is sufficiently close to  $\alpha$ .

Under these assumptions the convergence is quadratic, meaning that:

$$\text{If we denote by } \varepsilon_j = |x_j - \alpha|, \text{ then } \varepsilon_{j+1} \leq M\varepsilon_j^2,$$

where  $M$  is some constant. If  $f$  is twice differentiable, then

$$M \leq \max_{x \in I} |f''(x)| / \min_{x \in I} |f'(x)|.$$



Proof.

Condition (1) implies in particular that  $\alpha$  is a simple zero of  $f$ . Plugging  $\alpha$  in the Taylor expansion of  $f$  around  $x_i$  we get

$$\begin{aligned} 0 = f(\alpha) &= f(x_i) + f'(x_i)(\alpha - x_i) + \frac{f''(\eta)}{2}(\alpha - x_i)^2 \\ &= f(x_i) + f'(x_i)(\alpha - x_i) + \frac{f''(\eta)}{2}(\alpha - x_i)^2 \end{aligned} \tag{20}$$

where  $\eta$  is between  $\alpha$  and  $x_i$ . Dividing (20) with  $f'(x_i)$  we get

$$0 = \frac{f(x_i)}{f'(x_i)} - (\alpha - x_i) + \frac{f''(\eta)}{2f'(x_i)}e_i^2$$

and hence

$$\left(x_i - \frac{f(x_i)}{f'(x_i)}\right) - \alpha = x_{i+1} - \alpha = \frac{f''(\eta)}{2f'(x_i)}e_i^2.$$

Thus,

$$e_{i+1} = \left| \frac{f''(\eta)}{2f'(x_i)} \right| e_i^2$$

Now

$$\left| \frac{f''(\eta)}{2f'(x_i)} \right| \leq \frac{\max_{x \in I} |f''(x)|}{\min_{x \in I} |f'(x)|}.$$

To prove that the sequence converges note that there exists  $\delta_0 > 0$  such that

$$M\delta_0 < \frac{1}{2}.$$

Hence, if  $e_i \leq \delta_0$ , then

$$e_{i+1} = \left| \frac{f''(\eta)}{2f'(x_i)} \right| e_i^2 = \frac{1}{2} e_i.$$

Therefore

$$\lim_{n \rightarrow \infty} e_n = \lim_{n \rightarrow \infty} \frac{1}{2^n} \cdot e_0 = 0.$$



### Newton's method for $n = m > 1$

Newton's method generalizes to systems of  $n$  nonlinear equations in  $n$  unknowns:

- ▶  $x_0$  – initial approximation,
- ▶  $x_{k+1}$  – solution of

$$L_{x_k}(x) = f(x_k) + Df(x_k)(x - x_k) = 0,$$

so

$$x_{k+1} = x_k - Df(x_k)^{-1}f(x_k).$$

In practice inverses are difficult to calculate (require too many operations) and the linear system for  $\Delta x_k = x_{k+1} - x_k$

$$Df(x_k)\Delta x_k = -f(x_k)$$

is solved at each step (using  $LU$  decomposition of  $Df(x_k)$ ) and hence

$$x_{k+1} = x_k + \Delta x_k.$$

## Example

Derive Newton's method for solving the system of quadratic equations:

$$\begin{aligned}x^2 + y^2 - 10x + y &= 1, \\x^2 - y^2 - x + 10y &= 25.\end{aligned}$$

We are searching for the zero of the vector function

$$F : \mathbb{R}^2 \rightarrow \mathbb{R}^2, \quad F(x, y) = (x^2 + y^2 - 10x + y - 1, x^2 - y^2 - x + 10y - 25).$$

The Jacobian of  $F$  in  $(x, y)$  is

$$DF(x, y) = \begin{bmatrix} 2x - 10 & 2x - 1 \\ 2y + 1 & -2y + 10 \end{bmatrix}.$$

Using Newton's method we:

- ▶ Choose an initial term  $(x_0, y_0)$ .
- ▶ Calculate  $x_{r+1} = x_r + \Delta x_r$ , where  $DF(x_r, y_r)\Delta x_r = -F(x_r, y_r)^T$ .

## Newton optimization method:

We would like to find the extrema of the function  $F : \mathbb{R}^n \rightarrow \mathbb{R}$ .

Since the extrema are *critical (or stationary) points*, the candidates are zeroes of the gradient, i.e.,

$$G(x) := \text{grad } F(x) = [ F_{x_1}(x) \quad \cdots \quad F_{x_n}(x) ] = 0. \quad (21)$$

(21) is a system of  $n$  equations for  $n$  variables, the Jacobian of the vector function  $G$  is the so called Hessian of  $F$ :

$$DG(x) = H(x) = \begin{bmatrix} F_{x_1x_1} & \cdots & F_{x_1x_n} \\ \vdots & \ddots & \vdots \\ F_{x_nx_1} & \cdots & F_{x_nx_n} \end{bmatrix}.$$

If the sequence of iterates

$$x_0, \quad x_{k+1} = x_k - H^{-1}(x_k)G(x_k)$$

converges, the limit is a critical point of  $F$ , i.e., a candidate for the minimum (or maximum).

# Gradient descent

Optimization methods can also be used to ensure a **sufficiently accurate starting approximation** for the Newton-based techniques. (Like bisection does for a single one-variable equation.)

Finding solutions of the system  $F(x) = 0$ , where

$$F = [F_1, \dots, F_n]^T : \mathbb{R}^n \rightarrow \mathbb{R}^n$$

is equivalent to finding **global minima** of

$$g(x) := \|F\|^2 = F_1(x)^2 + \dots + F_n(x)^2 : \mathbb{R}^n \rightarrow \mathbb{R}.$$

We search for the local minima (**which are not necessarily global minima!**) of  $g$  as follows:

1. Choose  $x_0$ .
2. Determine the constant  $\alpha$  in  $x_r - \alpha \cdot \text{grad}(g(x_r))$  which minimizes

$$h(\alpha) = g(x_r - \alpha \cdot \text{grad}(g(x_r))).$$

(Or is significantly smaller than  $h(0) = g(x_r)$ .)

3.  $x_{r+1} = x_r - \alpha \cdot \text{grad}(g(x_r))$ .

## Quasi-Newtonov methods: Broyden's method

- ▶ For large  $n$ , the Newton's method is very expensive, since we need to evaluate  $n^2$  partial derivatives at each step and use  $\mathcal{O}(n^3)$  flops (+, -, ·, :) to solve the linear system.
- ▶ Broyden's method avoids computing derivatives. For  $n = m = 1$  it replaces the tangent by a secant through the last two iterates. It mimicks this idea also for larger  $n = m$ .

Let  $B_r$  be an approximate for  $J_f(x_r)$ . **Broyden's method** works as follows:

1. Solve  $B_r \Delta x_r = -f(x_r)$ ,
2.  $x_{r+1} = x_r + \Delta x_r$ ,
3. Determine  $B_{r+1}$ .

The last step searches for a matrix  $B_{r+1}$ , which fulfils the **secant condition**:

$$B_{r+1}(x_{r+1} - x_r) = f(x_{r+1}) - f(x_r)$$

and is the closest to  $B_r$  in the spectral norm  $\|\cdot\|_2$ .

It turns out that

$$B_{r+1} = B_r + \frac{f(x_{r+1})(\Delta x_r)^T}{\|\Delta x_r\|_2^2}.$$

### Application on the microwave oven example

Recall from above the microwave oven example. The system of equations for the parameters  $\alpha$ ,  $a$ ,  $b$  is:

$$\begin{aligned}\frac{\alpha}{1 + \sqrt{a^2 + b^2}} - 0.27 &= 0 \\ \frac{\alpha}{1 + \sqrt{(1-a)^2 + (1-b)^2}} - 0.36 &= 0 \\ \frac{\alpha}{1 + \sqrt{a^2 + (2-b)^2}} - 0.3 &= 0.\end{aligned}$$

Newton: [Click](#)

Broyden: [Click](#)

Gradient descent: [Click](#)

Tests: [Click](#)



### Newton's method for $m > n > 0$

We have an overdetermined system

$$f : \mathbb{R}^n \rightarrow \mathbb{R}^m, \quad f(x) = (0, \dots, 0) \quad (22)$$

of  $m$  nonlinear equations for  $n$  unknowns, where  $m > n$ .

The system (22) generally does not have a solution, so we are looking for a solution of (22) by the least squares method, i.e.,  $\alpha \in \mathbb{R}^n$  such that the distance of  $f(\alpha)$  from the origin is the smallest possible:

$$\|f(\alpha)\|^2 = \min\{\|f(x)\|^2\}.$$

The [Gauss-Newton method](#) is a generalization of the Newton's method, where instead of the inverse of the Jacobian its MP inverse is used at each step:

$$x_0 \dots \text{initial approximation}, \quad x_{k+1} = x_k - Df(x_k)^+ f(x_k),$$

where  $Df(x_k)^+$  is the MP inverse of  $Df(x_k)$ . If the matrix

$(Df(x_k)^T Df(x_k))$  is nonsingular at each step  $k$ , then

$$x_{k+1} = x_k - (Df(x_k)^T Df(x_k))^{-1} Df(x_k)^T f(x_k).$$

At each step  $x_{k+1}$  is the least squares approximation to the solution of the overdetermined linear system  $L_{x_k}(x) = 0$ , that is,

$$\|L_{x_k}(x_{k+1})\|^2 = \min\{\|L_{x_k}(x)\|^2, x \in \mathbb{R}^n\}.$$

Convergence is not guaranteed, but:

- ▶ if the sequence  $x_k$  converges, the limit  $x = \lim_k x_k$  is a local (but not necessarily global) minimum of  $\|f(x)\|^2$ .

It follows that the Gauss-Newton method is an algorithm for the local minimum of  $\|f(x)\|^2$ .

## Example

We are given point  $(x_i, y_i) \in \mathbb{R}^2$  for  $i = 1, \dots, m$  and are searching for the function

$$f(x, a, b) = ae^{bx}$$

which fits this data best by the method of least squares.

So we have the overdetermined system  $F(a, b) = 0$ , where

$$F : \mathbb{R}^2 \rightarrow \mathbb{R}^m, \quad F(a, b) = (y_1 - ae^{bx_1}, \dots, y_m - ae^{bx_m}).$$

The Jacobian of  $F$  is

$$DF(a, b) = \begin{bmatrix} -e^{bx_1} & ax_1 e^{bx_1} \\ \vdots & \vdots \\ -e^{bx_m} & ax_m e^{bx_m} \end{bmatrix}.$$

Using the Gauss-Newton method:

- ▶ We choose initial approximation  $(a_0, b_0)$ .
- ▶ Calculate iterates

$$\begin{bmatrix} a_{r+1} \\ b_{r+1} \end{bmatrix} = \begin{bmatrix} a_r \\ b_r \end{bmatrix} - DF(a_r, b_r)^+ F(a_r, b_r)^T.$$