

# Assignment 2 - Numerical Linear Algebra

Arthur Rabello Oliveira

29/04/2025

## Abstract

We derive linear and polynomial regression in subsets of  $\mathbb{R}$  and discuss the condition number of the associated matrices, numerical algorithms for the SVD and QR factorization are built and used on an efficiency analysis of the 3 methods to do linear or polynomial regression, stability of these algorithms is mentioned and

## Contents

1. Introduction .....	2
2. Linear Regression: The Least Squares Method (a) .....	2
3. Condition of a Problem .....	5
3.1. The Condition number of a problem .....	5
3.2. The relative Conditioning Number .....	6
3.3. Condition Number of Matrices .....	7
3.4. Application (b) .....	8
3.5. More Regression: A Polynomial Perspective (c) .....	9
3.6. Finding A, given (m,n) (d) .....	10
3.7. How Perturbations Affect The Condition Number of A (e) .....	10
3.8. A different dataset .....	11
3.8.1. How Conditioning Changes (f) .....	12
4. Some Algorithms .....	12
4.1. The SVD and QR factorizations .....	12
4.1.1. QR .....	12
4.1.2. SVD .....	12
4.2. A Python-Implementation and a Conditioning Analysis (b) .....	14
4.3. The polynomial approach: An efficiency analysis (c) .....	14

## 1. Introduction

Given  $D \subset \mathbb{R}^2$ , a dataset, approximating this set through a *continuous*  $f : \mathbb{R} \rightarrow \mathbb{R}$  is a very important problem in statistics, we will derive the 2 most important and most used methods to do this: linear and polynomial regression. Both are based on the least squares minimization problem, so an algebraical approach is used and on later sections, we will build some algorithms to better visualize both methods

## 2. Linear Regression: The Least Squares Method (a)

Given a dataset of equally spaced points  $D := \{t_i = \frac{i}{m}\}, i = 0, 1, \dots, m \in \mathbb{R}$ , linear regression consists of finding the best *line*  $f(t) = \alpha + \beta t$  that approximates the points  $(t_i, b_i) \in \mathbb{R}^2$ , where  $b_i$  are arbitrary

Approximating 2 points in  $\mathbb{R}^2$  by a line is trivial, now approximating more points is a task that requires linear algebra. To see this, we will analyze the following example to build intuition for the general case:

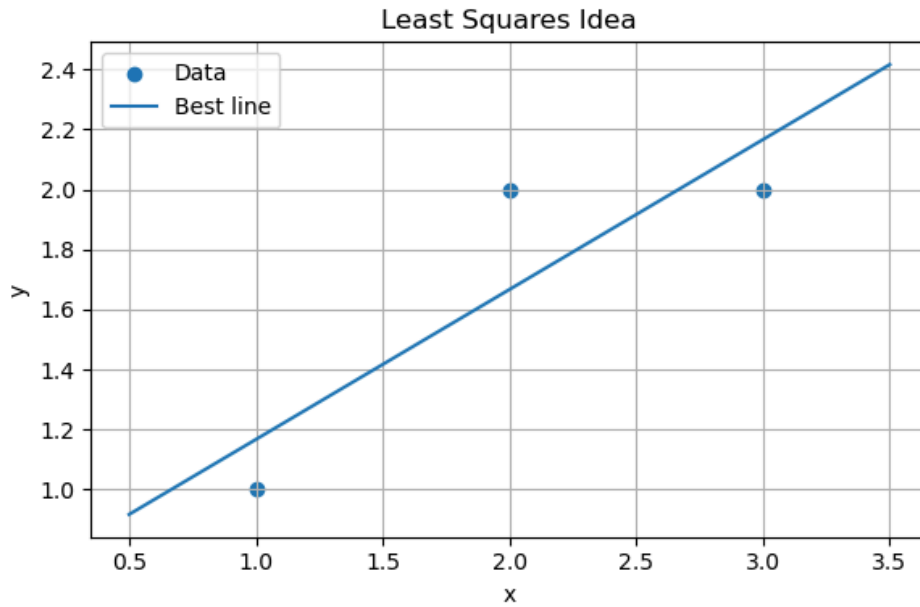


Figure 1: A glimpse into what we want to see

Given the points  $(1, 1), (2, 2), (3, 2) \in \mathbb{R}^2$ , we have  $(t_1, b_1) = (1, 1), (t_2, b_2) = (2, 2), (t_3, b_3) = (3, 2)$  we would like a *line*  $f(t) = y(t) = \alpha + \beta t$  that best approximates  $(t_i, b_i)$ , in other words, since we know that the line does not pass through all 3 points, we would like to find the *closest* line to **each point** of the dataset  $D$ , so the system:

$$\begin{aligned} f(1) &= \alpha + \beta = 1 \\ f(2) &= \alpha + 2\beta = 2 \\ f(3) &= \alpha + 3\beta = 2 \end{aligned} \tag{1}$$

Which is:

$$\underbrace{\begin{pmatrix} 1 & 1 \\ 1 & 2 \\ 1 & 3 \end{pmatrix}}_A \cdot \underbrace{\begin{pmatrix} \alpha \\ \beta \end{pmatrix}}_x = \underbrace{\begin{pmatrix} 1 \\ 2 \\ 2 \end{pmatrix}}_b \tag{2}$$

Clearly has no solution, (the line does not cross the 3 points), but it has a *closest solution*, which we can find through **minimizing** the errors produced by this approximation.

Let  $x^* \neq x$  be a solution to the system, let the error produced by approximating the points through a line be  $e = Ax - b$ , we want the smaller error *square* possible (that is why least squares). We square the error to avoid and detect outliers, so:

$$e_1^2 + e_2^2 + e_3^2 \quad 3$$

Is what we want to minimize, where  $e_i$  is the error (distance) from the  $i$ th point to the line:



Figure 2: The errors (distances)

So we will project  $b$  into  $C(A)$ , giving us the closest solution, and the least squares solutions is when  $\hat{x}$  minimizes  $\|Ax - b\|^2$ , this occurs when the residual  $e = Ax - b$  is orthogonal to  $C(A)$ , since  $N(A^T) \perp C(A)$  and the dimensions sum up the left dimension of the matrix, so by the well-known projection formula, we have:

$$\begin{aligned} A^T A \hat{x} &= A^T b \\ &= \begin{pmatrix} 1 & 1 & 1 \\ 1 & 2 & 3 \end{pmatrix} \cdot \begin{pmatrix} 1 & 1 \\ 1 & 2 \\ 1 & 3 \end{pmatrix} \cdot \begin{pmatrix} \alpha \\ \beta \end{pmatrix} = \begin{pmatrix} 3 & 6 \\ 6 & 14 \end{pmatrix} \cdot \begin{pmatrix} \alpha \\ \beta \end{pmatrix} \\ &= \begin{pmatrix} 3 & 6 \\ 6 & 14 \end{pmatrix} \cdot \begin{pmatrix} 1 \\ 2 \\ 3 \end{pmatrix} = \begin{pmatrix} 5 \\ 11 \end{pmatrix} \end{aligned} \quad 4$$

So the system to find  $\hat{x} = \begin{pmatrix} \hat{\alpha} \\ \hat{\beta} \end{pmatrix}$  becomes:

$$\begin{aligned} 3\alpha + 5\beta &= 5 \\ 6\alpha + 14\beta &= 11 \end{aligned} \quad 5$$

Notice that with the errors  $e_i^2$  as:

$$\begin{aligned}
e_1^2 &= (f(t_1) - b_1)^2 = (f(1) - 1)^2 = (\alpha + \beta - 1)^2 \\
e_2^2 &= (f(t_2) - b_2)^2 = (f(2) - 2)^2 = (\alpha + 2\beta - 2)^2 \\
e_3^2 &= (f(t_3) - b_2)^2 = (f(3) - 2)^2 = (\alpha + 3\beta - 2)^2
\end{aligned} \tag{6}$$

The system in eq. (5) is *precisely* what is obtained after using partial derivatives to minimize the errors sum as a function of  $(\alpha, \beta)$ :

$$\begin{aligned}
f(\alpha, \beta) &= (\alpha + \beta - 1)^2 + (\alpha + 2\beta - 2)^2 + (\alpha + 3\beta - 2)^2 \\
&= 3\alpha^2 + 14\beta^2 + 12\alpha\beta - 10\alpha - 22\beta + 9, \\
\frac{\partial f}{\partial \alpha} = \frac{\partial f}{\partial \beta} &= 0 \Leftrightarrow 6\alpha + 12\beta - 10 = 28\beta + 12\alpha - 22 = 0 \Leftrightarrow \begin{cases} 3\alpha + 6\beta = 11 \\ 6\alpha + 14\beta = 11 \end{cases}
\end{aligned} \tag{7}$$

This new system has a solution in  $\hat{\alpha} = \frac{2}{3}, \hat{\beta} = \frac{1}{2}$ , so the equation of the optimal line, obtained through *linear regression* (or least squares) is:

$$y(t) = \frac{2}{3} + \frac{1}{2}t. \tag{8}$$

If we have  $n > 3$  points to approximate through a line, the reasoning is analogous:

Going back to  $D$ , we want to find the extended system as we did in eq. (7), so let the best line be:

$$f(t) = \alpha + \beta t \tag{9}$$

That best approximates the points  $(0, b_0), (\frac{1}{m}, b_1), \dots, (1, b_m)$ . The system is:

$$\begin{aligned}
f(0) &= b_0 = \alpha, \\
f\left(\frac{1}{m}\right) &= b_1 = \alpha + \frac{\beta}{m}, \\
f\left(\frac{2}{m}\right) &= b_2 = \alpha + \frac{2}{m}\beta \\
&\dots \\
f(1) &= b_m = \alpha + \beta
\end{aligned} \tag{10}$$

Or:

$$\underbrace{\begin{pmatrix} 1 & 0 \\ 1 & \frac{1}{m} \\ \vdots & \vdots \\ 1 & 1 \end{pmatrix}}_A \cdot \underbrace{\begin{pmatrix} \alpha \\ \beta \end{pmatrix}}_x = \underbrace{\begin{pmatrix} b_0 \\ \vdots \\ b_m \end{pmatrix}}_b \tag{11}$$

Projecting into  $C(A)$ , we have:

$$\begin{aligned}
A^T A x &= A^T b \\
&= \begin{pmatrix} 1 & 1 & \dots & 1 \\ 0 & \frac{1}{m} & \dots & 1 \end{pmatrix} \cdot \begin{pmatrix} 1 & 0 \\ 1 & \frac{1}{m} \\ \vdots & \vdots \\ 1 & 1 \end{pmatrix} = \begin{pmatrix} m+1 & \frac{m+1}{2} \\ \frac{m+1}{2} & \frac{(m+1)(2m+2)}{6m} \end{pmatrix} \cdot \begin{pmatrix} \hat{\alpha} \\ \hat{\beta} \end{pmatrix} \\
&= \begin{pmatrix} 1 & 1 & \dots & 1 \\ 0 & \frac{1}{m} & \dots & 1 \end{pmatrix} \cdot \begin{pmatrix} b_0 \\ \vdots \\ b_m \end{pmatrix} = \begin{pmatrix} b_0 + b_2 + \dots + b_m \\ \frac{1}{m}[b_1 + 2b_2 + \dots + (m-1)b_{m-1} + b_m] \end{pmatrix}
\end{aligned} \tag{12}$$

So the system to find the optimal vector  $\begin{pmatrix} \hat{\alpha} \\ \hat{\beta} \end{pmatrix}$  is:

$$\begin{pmatrix} m+1 & \frac{m+1}{2} \\ \frac{m+1}{2} & \frac{(m+1)(2m+2)}{6m} \end{pmatrix} \cdot \begin{pmatrix} \hat{\alpha} \\ \hat{\beta} \end{pmatrix} = \begin{pmatrix} b_0 + b_1 + \dots + b_m \\ \frac{1}{m}[b_1 + 2b_2 + \dots + (m-1)b_{m-1} + b_m] \end{pmatrix} \quad 13$$

Or:

$$\begin{pmatrix} m+1 & \sum_{i=1}^m t_i \\ \sum_{i=1}^m t_i & \sum_{i=1}^m t_i^2 \end{pmatrix} \cdot \begin{pmatrix} \hat{\alpha} \\ \hat{\beta} \end{pmatrix} = \begin{pmatrix} \sum_{i=1}^m b_i \\ \sum_{i=1}^m \frac{i}{m} \cdot b_i \end{pmatrix} \quad 14$$

And the least squares optimal solution is:

$$\begin{pmatrix} \hat{\alpha} \\ \hat{\beta} \end{pmatrix} = \begin{pmatrix} m+1 & \sum_{i=1}^m t_i \\ \sum_{i=1}^m t_i & \sum_{i=1}^m t_i^2 \end{pmatrix}^{-1} \cdot \begin{pmatrix} \sum_{i=1}^m b_i \\ \sum_{i=1}^m \frac{i}{m} \cdot b_i \end{pmatrix} \quad 15$$

Now we will discuss if this method is efficient, using **Conditioning Numbers** as base.

### 3. Condition of a Problem

A *problem* is usually described as a function  $f : X \rightarrow Y$  from a **normed** vector space  $X$  of data (it has to be normed so we can *quantify* data) and a *normed* vector space  $Y$  of solutions,  $f$  is not always a well-behaved continuous function, which is why we are interested in **well-conditioned** problems and not in **ill-conditioned** problems, which we define:

**Definition 3.1:** (Well-Conditioned Problem) A problem  $f : X \rightarrow Y$  is *well-conditioned* at  $x_0 \in X \Leftrightarrow \forall \varepsilon > 0, \exists \delta > 0 \mid \|x - x_0\| < \delta \Rightarrow \|f(x) - f(x_0)\| < \varepsilon$ .

This means that small perturbations in  $x$  lead to small changes in  $f(x)$ , a problem is **ill-conditioned** if  $f(x)$  can suffer huge changes with small changes in  $x$ .

We usually say  $f$  is well-conditioned if it is well-conditioned  $\forall x \in X$ , if there is at least one  $x_i$  in which the problem is ill-conditioned, then we can use that whole problem is ill-conditioned.

#### 3.1. The Condition number of a problem

Conditioning numbers are a tool to quantify how well/ill conditioned a problem is:

**Definition 3.1.1:** (Absolute Conditioning Number) Let  $\delta x$  be a small perturbation of  $x$ , so  $\delta f = f(x + \delta x) - f(x)$ . The **absolute** conditioning number of  $f$  is:

$$\hat{\kappa} = \lim_{\delta \rightarrow 0} \sup_{\|\delta x\| \leq \delta} \frac{\|\delta f\|}{\|\delta x\|} \quad 16$$

The limit of the supremum can be seen as the supremum of all *infinitesimal* perturbations, so this can be rewritten as:

$$\hat{\kappa} = \sup_{\delta x} \frac{\|\delta f\|}{\|\delta x\|} \quad 17$$

If  $f$  is differentiable, we can evaluate the abs.conditioning number using its derivative, if  $J$  is the matrix whose  $i \times j$  entry is the derivative  $\frac{\partial f_i}{\partial x_j}$  (jacobian of  $f$ ), then we know that  $\delta f \approx J(x)\delta x$ , with equality in the limit  $\|\delta x\| \rightarrow 0$ . So the absolute conditioning number of  $f$  becomes:

$$\hat{\kappa} = \|J(x)\|, \quad 18$$

### 3.2. The relative Conditioning Number

When, instead of analyzing the whole set  $X$  of data, we are interested in *relative* changes, we use the **relative condition number**:

**Definition 3.2.1:** (Relative Condition Number) Given  $f : X \rightarrow Y$  a problem, the *relative condition number*  $\kappa(x)$  at  $x \in X$  is:

$$\kappa(x) = \lim_{\delta \rightarrow 0} \sup_{\|\delta x\| \leq \delta} \left( \frac{\|\delta f\|}{\|f(x)\|} \right) \cdot \left( \frac{\|\delta x\|}{\|x\|} \right)^{-1} \quad 19$$

Or, as we did in Definition 3.1.1, assuming that  $\delta f$  and  $\delta x$  are infinitesimal:

$$\kappa(x) = \sup_{\delta x} \left( \frac{\|\delta f\|}{\|f(x)\|} \right) \cdot \left( \frac{\|\delta x\|}{\|x\|} \right)^{-1} \quad 20$$

If  $f$  is differentiable:

$$\kappa(x) = (\|J(x)\|) \cdot \left( \frac{\|f(x)\|}{\|x\|} \right)^{-1} \quad 21$$

Relative condition numbers are more useful than absolute conditioning numbers because the **floating point arithmetic** used in many computers produces *relative* errors, the latter is not a highlight of this discussion.

Here are some examples of conditioning:

*Example 3.2.1.:* Consider the problem of obtaining the scalar  $\frac{x}{2}$  from  $x \in \mathbb{R}$ . The function  $f(x) = \frac{x}{2}$  is differentiable, so by eq. (21):

$$\kappa(x) = (\|J\|) \cdot \left( \frac{\|f(x)\|}{\|x\|} \right)^{-1} = \left( \frac{1}{2} \right) \cdot \left( \frac{\frac{x}{2}}{x} \right)^{-1} = 1. \quad 22$$

This problem is well-conditioned ( $\kappa$  is small).

*Example 3.2.2.:* Consider the problem of computing the scalar  $x_1 - x_2$  from  $(x_1, x_2) \in \mathbb{R}^2$  (Use the  $\infty$ -norm in  $\mathbb{R}^2$  for simplicity). The function associated is differentiable and the jacobian is:

$$J = \left( \frac{\partial f}{\partial x_1} \quad \frac{\partial f}{\partial x_2} \right) = (1 \quad -1) \quad 23$$

With  $\|J\|_\infty = 2$ , so the condition number is:

$$\kappa = (\|J\|_\infty) \cdot \left( \frac{\|f(x)\|}{\|x\|} \right)^{-1} = \frac{2}{|x_1 - x_2| \cdot \max\{|x_1|, |x_2|\}} \quad 24$$

This problem can be ill-conditioned if  $|x_1 - x_2| \approx 0$  ( $\kappa$  gets huge), and well-conditioned otherwise

### 3.3. Condition Number of Matrices

We will deduce the conditioning number of a matrix from the conditioning number of *matrix-vector* multiplication:

Consider the problem of obtaining  $Ax$  given  $A \in \mathbb{C}^{m \times n}$ , we will calculate the relative condition number with respect to perturbations in  $x$ , directly from Definition 3.2.1, we have:

$$\kappa = \sup_{\delta x} \frac{\|A(x + \delta x) - Ax\|}{\|Ax\|} \cdot \left( \frac{\|\delta x\|}{\|x\|} \right)^{-1} = \sup_{\delta x} \frac{\|A\delta x\|}{\|\delta x\|} \cdot \left( \frac{\|Ax\|}{\|x\|} \right)^{-1} \quad 25$$

Since  $\sup_{\delta x} \frac{\|A\delta x\|}{\|\delta x\|} = \|A\|$ , we have:

$$\kappa = \|A\| \cdot \frac{\|x\|}{\|Ax\|} \quad 26$$

This is a precise formula as a function of  $(A, x)$ .

Suppose for a moment that  $A$  is square and non-singular:

**Theorem 3.3.1:**  $\forall x \in \mathbb{C}^n, A \in \mathbb{C}^{n \times n}, \det(A) \neq 0$ , the following holds:

$$\frac{\|x\|}{\|Ax\|} \leq \|A^{-1}\| \quad 27$$

*Proof:* Since  $\forall A, B \in \mathbb{C}^{n \times n}, \|AB\| \leq \|A\| \|B\|$ , we have:

$$\|AA^{-1}x\| \leq \|Ax\| \|A^{-1}\| \Leftrightarrow \frac{\|x\|}{\|Ax\|} \leq \|A^{-1}\| \quad 28$$

□

So using this in eq. (26), we can write:

$$\kappa \leq \|A\| \cdot \|A^{-1}\| \quad 29$$

Or:

$$\kappa = \alpha \|A\| \cdot \|A^{-1}\| \quad 30$$

With

$$\alpha = \frac{\|x\|}{\|Ax\|} \cdot (\|A^{-1}\|)^{-1} \quad 31$$

From Theorem 3.3.1, we can choose  $x$  to make  $\alpha = 1$ , and therefore  $\kappa = \|A\| \cdot \|A^{-1}\|$ .

Consider now the problem of calculating  $A^{-1}b$  given  $A \in \mathbb{C}^{n \times n}$ . This is mathematically identical to the problem we just analyzed, so the following theorem has already been proven:

**Theorem 3.3.2:** Let  $A \in \mathbb{C}^{n \times n}$ ,  $\det(A) \neq 0$ , and consider the problem of computing  $b$ , from  $Ax = b$ , by perturbing  $x$ . Then the following holds:

$$\kappa = \|A\| \frac{\|x\|}{\|b\|} \leq \|A\| \cdot \|A^{-1}\| \quad 32$$

Where  $\kappa$  is the condition number of the problem.

*Proof:* Read from eq. (25) to eq. (31). □

Finally,  $\|A\| \cdot \|A^{-1}\|$  is so useful it has a name: **the condition number of A** (relative to the norm  $\|\cdot\|$ )

If  $A$  is singular, usually we write  $\kappa(A) = \infty$ . Notice that if  $\|\cdot\| = \|\cdot\|_2$ , then  $\|A\| = \sigma_1$  and  $\|A^{-1}\| = \frac{1}{\sigma_m}$ , so:

$$\kappa(A) = \frac{\sigma_1}{\sigma_m} \quad 33$$

This is very useful, as we show an interesting application:

### 3.4. Application (b)

Still on least squares, using  $\|\cdot\|_2$  the conditioning number of eq. (14) is:

$$\kappa(A) = \|A\|_2 \cdot \|A^{-1}\|_2 = \frac{\sigma_1}{\sigma_m} \quad 34$$

And the singular values  $\sigma_i$  are the square roots of the eigenvalues  $\lambda_i$  of  $B = A^T A$  (non-crescent order), so from eq. (13):

$$\begin{aligned} \det(B - \lambda I) = 0 &\Leftrightarrow \det \left( \begin{pmatrix} m+1-\lambda & \frac{m+1}{2} \\ \frac{m+1}{2} & \frac{(m+1)(2m+1)}{6} - \lambda \end{pmatrix} \right) = 0 \\ &\Leftrightarrow (m+1-\lambda) \left[ \frac{(m+1)(2m+1)}{6} - \lambda \right] - \left( \frac{m+1}{2} \right)^2 = 0 \\ &\Leftrightarrow \lambda^2 - \frac{(m+1)(8m+1)}{6m} \lambda + \frac{(m+1)^2(m+2)}{12m} = 0 \\ &\Leftrightarrow \lambda = \frac{m+1}{12m} \left[ (8m+1) \pm \sqrt{52m^2 - 8m + 1} \right] \end{aligned} \quad 35$$

So the singular values are:

$$\begin{aligned} \sigma_1 = \sqrt{\lambda_1} &= \sqrt{\frac{m+1}{12m} \left[ (8m+1) + \sqrt{52m^2 - 8m + 1} \right]}, \\ \sigma_2 = \sqrt{\lambda_2} &= \sqrt{\frac{m+1}{12m} \left[ (8m+1) - \sqrt{52m^2 - 8m + 1} \right]} \end{aligned} \quad 36$$

And the condition number of the matrix  $A$  is:



$$\begin{aligned}\kappa(A) &= \frac{\sigma_1}{\sigma_m} = \frac{\sqrt{\frac{m+1}{12m} [(8m+1) + \sqrt{52m^2 - 8m + 1}]}}{\sqrt{\frac{m+1}{12m} [(8m+1) - \sqrt{52m^2 - 8m + 1}]}} \\ &= \sqrt{\frac{(8m+1) + \sqrt{52m^2 - 8m + 1}}{(8m+1) - \sqrt{52m^2 - 8m + 1}}}\end{aligned}\tag{37}$$

When  $m \rightarrow \infty$ : COLOCAR O LIMITE

Here are some python examples:

(COLOCA AS PORRA DOS EXEMPLO)

### 3.5. More Regression: A Polynomial Perspective (c)

In this section we will discuss what changes when we decide to use **polynomials** instead of **lines** to approximate our dataset:

$$f(t) = \alpha + \beta t \rightarrow p(t) = \varphi_0 + \varphi_1 t + \dots + \varphi_n t^n \tag{38}$$

From a first perspective, it seems way more efficient to describe a dataset with many variables then to do so with a simple line  $\alpha + \beta t$ , so let's use the same dataset  $S := \{(t_i, b_i), t_i = \frac{i}{m}\}, i = 0, 1, \dots, m$ . Where  $b_i$  is arbitrary, as we did in Section 2. Finding the new system to be solved gives us:

$$\begin{aligned}p(t_0 = 0) &= b_0 = \varphi_0, \\ p\left(t_1 = \frac{1}{m}\right) &= b_1 = \varphi_0 + \varphi_1 \frac{1}{m} + \dots + \varphi_n \left(\frac{1}{m}\right)^n \\ p\left(t_2 = \frac{2}{m}\right) &= b_2 = \varphi_0 + \varphi_1 \frac{2}{m} + \varphi_2 \left(\frac{2}{m}\right)^2 + \dots + \varphi_n \left(\frac{2}{m}\right)^n \\ &\vdots \\ p(t_m = 1) &= b_m = \varphi_0 + \dots + \varphi_n\end{aligned}\tag{39}$$

And:

$$\underbrace{\begin{pmatrix} 1 & 0 & 0 & \dots & 0 \\ 1 & \frac{1}{m} & \left(\frac{1}{m}\right)^2 & \dots & \left(\frac{1}{m}\right)^n \\ 1 & \frac{2}{m} & \left(\frac{2}{m}\right)^2 & \dots & \left(\frac{2}{m}\right)^n \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & 1 & 1 & \dots & 1 \end{pmatrix}}_{A_{n+1 \times n+1}} \cdot \underbrace{\begin{pmatrix} \varphi_0 \\ \varphi_1 \\ \varphi_2 \\ \vdots \\ \varphi_n \end{pmatrix}}_{\Phi_{n+1 \times 1}} = \underbrace{\begin{pmatrix} b_0 \\ b_1 \\ b_2 \\ \vdots \\ b_m \end{pmatrix}}_{b_{m+1 \times 1}} \tag{40}$$

Projecting into  $C(A)$ :

$$\begin{aligned}
A^T A \hat{\Phi} &= \begin{pmatrix} 1 & 1 & 1 & \dots & 1 \\ 0 & \frac{1}{m} & \frac{2}{m} & \dots & 1 \\ 0 & (\frac{1}{m})^2 & (\frac{2}{m})^2 & \dots & 1 \\ \vdots & \vdots & \vdots & \dots & \vdots \\ 0 & (\frac{1}{m})^n & (\frac{2}{m})^n & \dots & 1 \end{pmatrix} \cdot \begin{pmatrix} 1 & 0 & 0 & \dots & 0 \\ 1 & \frac{1}{m} & (\frac{1}{m})^2 & \dots & (\frac{1}{m})^n \\ 1 & \frac{2}{m} & (\frac{2}{m})^2 & \dots & (\frac{2}{m})^n \\ \vdots & \vdots & \vdots & \dots & \vdots \\ 1 & 1 & 1 & \dots & 1 \end{pmatrix} \cdot \begin{pmatrix} \hat{\varphi}_0 \\ \hat{\varphi}_1 \\ \hat{\varphi}_2 \\ \vdots \\ \hat{\varphi}_n \end{pmatrix} \\
&= \begin{pmatrix} n+1 & \sum_{i=1}^m \frac{i}{m} & \sum_{i=1}^m (\frac{i}{m})^2 & \dots & \sum_{i=1}^m (\frac{i}{m})^n \\ \sum_{i=1}^m \frac{i}{m} & \sum_{i=1}^m (\frac{i}{m})^2 & \sum_{i=1}^m (\frac{i}{m})^3 & \dots & \sum_{i=1}^m (\frac{i}{m})^{n+1} \\ \sum_{i=1}^m (\frac{i}{m})^2 & \sum_{i=1}^m (\frac{i}{m})^3 & \sum_{i=1}^m (\frac{i}{m})^4 & \dots & \sum_{i=1}^m (\frac{i}{m})^n + 2 \\ \vdots & \vdots & \vdots & \dots & \vdots \\ \sum_{i=1}^m (\frac{i}{m})^n & \sum_{i=1}^m (\frac{i}{m})^{n+1} & \sum_{i=1}^m (\frac{i}{m})^{n+2} & \dots & \sum_{i=1}^m (\frac{i}{m})^{2n} \end{pmatrix} \cdot \begin{pmatrix} \hat{\varphi}_0 \\ \hat{\varphi}_1 \\ \hat{\varphi}_2 \\ \vdots \\ \hat{\varphi}_n \end{pmatrix} \\
&= \begin{pmatrix} 1 & 1 & 1 & \dots & 1 \\ 0 & \frac{1}{m} & \frac{2}{m} & \dots & 1 \\ 0 & (\frac{1}{m})^2 & (\frac{2}{m})^2 & \dots & 1 \\ \vdots & \vdots & \vdots & \dots & \vdots \\ 0 & (\frac{1}{m})^n & (\frac{2}{m})^n & \dots & 1 \end{pmatrix} \cdot \begin{pmatrix} b_0 \\ b_1 \\ b_2 \\ \vdots \\ b_m \end{pmatrix} = \begin{pmatrix} \sum_{i=0}^m b_i \\ \sum_{i=0}^m \frac{ib_i}{m} \\ \sum_{i=0}^m (\frac{i}{m})^2 m \\ \vdots \\ \sum_{i=0}^m (\frac{i}{m})^n b_i \end{pmatrix}
\end{aligned} \tag{41}$$

So the system to be solved is:

$$\begin{pmatrix} n+1 & \sum_{i=1}^m \frac{i}{m} & \dots & \sum_{i=1}^m (\frac{i}{m})^n \\ \sum_{i=1}^m \frac{i}{m} & \sum_{i=1}^m (\frac{i}{m})^2 & \dots & \sum_{i=1}^m (\frac{i}{m})^{n+1} \\ \sum_{i=1}^m (\frac{i}{m})^2 & \sum_{i=1}^m (\frac{i}{m})^3 & \dots & \sum_{i=1}^m (\frac{i}{m})^n + 2 \\ \vdots & \vdots & \vdots & \vdots \\ \sum_{i=1}^m (\frac{i}{m})^n & \sum_{i=1}^m (\frac{i}{m})^{n+1} & \dots & \sum_{i=1}^m (\frac{i}{m})^{2n} \end{pmatrix} \cdot \begin{pmatrix} \hat{\varphi}_0 \\ \hat{\varphi}_1 \\ \hat{\varphi}_2 \\ \vdots \\ \hat{\varphi}_n \end{pmatrix} = \begin{pmatrix} \sum_{i=0}^m b_i \\ \sum_{i=0}^m \frac{ib_i}{m} \\ \sum_{i=0}^m (\frac{i}{m})^2 m \\ \vdots \\ \sum_{i=0}^m (\frac{i}{m})^n b_i \end{pmatrix} \tag{42}$$

Therefore the least squares *polynomial regression* solution is:

$$\begin{pmatrix} \hat{\varphi}_0 \\ \hat{\varphi}_1 \\ \hat{\varphi}_2 \\ \vdots \\ \hat{\varphi}_n \end{pmatrix} = \begin{pmatrix} n+1 & \sum_{i=1}^m \frac{i}{m} & \dots & \sum_{i=1}^m (\frac{i}{m})^n \\ \sum_{i=1}^m \frac{i}{m} & \sum_{i=1}^m (\frac{i}{m})^2 & \dots & \sum_{i=1}^m (\frac{i}{m})^{n+1} \\ \sum_{i=1}^m (\frac{i}{m})^2 & \sum_{i=1}^m (\frac{i}{m})^3 & \dots & \sum_{i=1}^m (\frac{i}{m})^n + 2 \\ \vdots & \vdots & \vdots & \vdots \\ \sum_{i=1}^m (\frac{i}{m})^n & \sum_{i=1}^m (\frac{i}{m})^{n+1} & \dots & \sum_{i=1}^m (\frac{i}{m})^{2n} \end{pmatrix}^{-1} \cdot \begin{pmatrix} \sum_{i=0}^m b_i \\ \sum_{i=0}^m \frac{ib_i}{m} \\ \sum_{i=0}^m (\frac{i}{m})^2 m \\ \vdots \\ \sum_{i=0}^m (\frac{i}{m})^n b_i \end{pmatrix} \tag{43}$$

### 3.6. Finding A, given (m,n) (d)

Here we show a python-implemented function that calculates the matrix  $A$  as a function of the dimensions  $(m, n)$ :

BOTAR A FUNÇÃO PYTHON

### 3.7. How Perturbations Affect The Condition Number of A (e)

In this section we analyze what happens to  $\kappa(A)$ , when  $A$  is perturbed with  $m = 100$  and  $n = 1, \dots, 20$ , the following graphs have been produced by the algorithm shown in BOTAR O ALGORITMO:

### 3.8. A different dataset

If we change  $S := \{(t_i, b_i) \mid t_i = \frac{i}{m}, i = 0, 1, \dots, m\}$  to  $\hat{S} = \{(t_i, b_i) \mid t_i = \frac{i}{m} - \frac{1}{2}\}$ , the polynomial regression becomes:

$$\begin{aligned} p\left(t_0 = -\frac{1}{2}\right) &= \varphi_0 + \varphi_1\left(-\frac{1}{2}\right) + \dots + \varphi_n\left(-\frac{1}{2}\right)^n = b_0 \\ p\left(t_1 = \frac{1}{m} - \frac{1}{2}\right) &= \varphi_0 + \varphi_1\left(\frac{1}{m} - \frac{1}{2}\right) + \varphi_2\left(\frac{1}{m} - \frac{1}{2}\right)^2 + \dots + \varphi_n\left(\frac{1}{m} - \frac{1}{2}\right)^n \\ &\vdots \\ p\left(t_m = 1 - \frac{1}{2}\right) &= \varphi_0 + \varphi_1\left(1 - \frac{1}{2}\right) + \dots + \varphi_n\left(1 - \frac{1}{2}\right)^n \end{aligned} \quad 44$$

So:

$$\underbrace{\begin{pmatrix} 1 & -\frac{1}{2} & \left(-\frac{1}{2}\right)^2 & \dots & \left(-\frac{1}{2}\right)^n \\ 1 & \left(\frac{1}{m} - \frac{1}{2}\right) & \left(\frac{1}{m} - \frac{1}{2}\right)^2 & \dots & \left(\frac{1}{m} - \frac{1}{2}\right)^n \\ 1 & \left(\frac{2}{m} - \frac{1}{2}\right) & \left(\frac{2}{m} - \frac{1}{2}\right)^2 & \dots & \left(\frac{2}{m} - \frac{1}{2}\right)^n \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & \left(-\frac{1}{2}\right) & \left(-\frac{1}{2}\right)^2 & \dots & \left(-\frac{1}{2}\right)^n \end{pmatrix}}_A \cdot \underbrace{\begin{pmatrix} \varphi_0 \\ \varphi_1 \\ \vdots \\ \varphi_n \end{pmatrix}}_{\Phi} = \underbrace{\begin{pmatrix} b_0 \\ b_1 \\ \vdots \\ b_m \end{pmatrix}}_b \quad 45$$

Projecting onto  $C(A)$ :

$$\begin{aligned} &\underbrace{\begin{pmatrix} 1 & 1 & 1 & \dots & 1 \\ -\frac{1}{2} & \left(\frac{1}{m} - \frac{1}{2}\right) & \left(\frac{2}{m} - \frac{1}{2}\right) & \dots & -\frac{1}{2} \\ \left(-\frac{1}{2}\right)^2 & \left(\frac{1}{m} - \frac{1}{2}\right)^2 & \left(\frac{2}{m} - \frac{1}{2}\right)^2 & \dots & \left(-\frac{1}{2}\right)^2 \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \left(-\frac{1}{2}\right)^n & \left(\frac{1}{m} - \frac{1}{2}\right)^n & \left(\frac{2}{m} - \frac{1}{2}\right)^n & \dots & \left(-\frac{1}{2}\right)^n \end{pmatrix}}_{A^T} \cdot \underbrace{\begin{pmatrix} 1 & -\frac{1}{2} & \left(-\frac{1}{2}\right)^2 & \dots & \left(-\frac{1}{2}\right)^n \\ 1 & \left(\frac{1}{m} - \frac{1}{2}\right) & \left(\frac{1}{m} - \frac{1}{2}\right)^2 & \dots & \left(\frac{1}{m} - \frac{1}{2}\right)^n \\ 1 & \left(\frac{2}{m} - \frac{1}{2}\right) & \left(\frac{2}{m} - \frac{1}{2}\right)^2 & \dots & \left(\frac{2}{m} - \frac{1}{2}\right)^n \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & -\frac{1}{2} & \left(-\frac{1}{2}\right)^2 & \dots & \left(-\frac{1}{2}\right)^n \end{pmatrix}}_A \cdot \underbrace{\begin{pmatrix} \widehat{\varphi}_0 \\ \widehat{\varphi}_1 \\ \vdots \\ \widehat{\varphi}_n \end{pmatrix}}_{\hat{\Phi}} \\ &= \begin{pmatrix} n+1 & \sum_{i=0}^n \left(\frac{i}{m} - \frac{1}{2}\right) & \sum_{i=0}^n \left(\frac{i}{m} - \frac{1}{2}\right)^2 & \dots & \sum_{i=0}^n \left(\frac{i}{m} - \frac{1}{2}\right)^n \\ \sum_{i=0}^n \frac{i}{m} - \frac{1}{2} & \sum_{i=0}^n \left(\frac{i}{m} - \frac{1}{2}\right)^2 & \sum_{i=0}^n \left(\frac{i}{m} - \frac{1}{2}\right)^3 & \dots & \sum_{i=0}^n \left(\frac{i}{m} - \frac{1}{2}\right)^{n+1} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \sum_{i=0}^n \left(\frac{i}{m} - \frac{1}{2}\right)^n & \sum_{i=0}^n \left(\frac{i}{m} - \frac{1}{2}\right)^{n+1} & \sum_{i=0}^n \left(\frac{i}{m} - \frac{1}{2}\right)^{n+2} & \dots & \sum_{i=0}^n \left(\frac{i}{m} - \frac{1}{2}\right)^{2n} \end{pmatrix} \cdot \begin{pmatrix} \widehat{\varphi}_0 \\ \widehat{\varphi}_1 \\ \vdots \\ \widehat{\varphi}_n \end{pmatrix} \quad 46 \\ &= \begin{pmatrix} 1 & 1 & 1 & \dots & 1 \\ -\frac{1}{2} & \left(\frac{1}{m} - \frac{1}{2}\right) & \left(\frac{2}{m} - \frac{1}{2}\right) & \dots & -\frac{1}{2} \\ \left(-\frac{1}{2}\right)^2 & \left(\frac{1}{m} - \frac{1}{2}\right)^2 & \left(\frac{2}{m} - \frac{1}{2}\right)^2 & \dots & \left(-\frac{1}{2}\right)^2 \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \left(-\frac{1}{2}\right)^n & \left(\frac{1}{m} - \frac{1}{2}\right)^n & \left(\frac{2}{m} - \frac{1}{2}\right)^n & \dots & \left(-\frac{1}{2}\right)^n \end{pmatrix} \cdot \begin{pmatrix} b_0 \\ b_1 \\ b_2 \\ \vdots \\ b_m \end{pmatrix} = \begin{pmatrix} \sum_{i=0}^n b_i \\ \sum_{i=0}^n \left(\frac{i}{m} - \frac{1}{2}\right) b_i \\ \sum_{i=0}^n \left(\frac{i}{m} - \frac{1}{2}\right)^2 b_i \\ \vdots \\ \sum_{i=0}^n \left(\frac{i}{m} - \frac{1}{2}\right)^n b_i \end{pmatrix} \end{aligned}$$

$\therefore$  The least squares optimal solution  $(\widehat{\varphi}_0, \widehat{\varphi}_1, \dots, \widehat{\varphi}_n)$  is:

$$\begin{pmatrix} n+1 & \sum_{i=0}^n \left(\frac{i}{m} - \frac{1}{2}\right) & \dots & \sum_{i=0}^n \left(\frac{i}{m} - \frac{1}{2}\right)^n \\ \sum_{i=0}^n \frac{i}{m} - \frac{1}{2} & \sum_{i=0}^n \left(\frac{i}{m} - \frac{1}{2}\right)^2 & \dots & \sum_{i=0}^n \left(\frac{i}{m} - \frac{1}{2}\right)^{n+1} \\ \vdots & \vdots & \ddots & \vdots \\ \sum_{i=0}^n \left(\frac{i}{m} - \frac{1}{2}\right)^n & \sum_{i=0}^n \left(\frac{i}{m} - \frac{1}{2}\right)^{n+1} & \dots & \sum_{i=0}^n \left(\frac{i}{m} - \frac{1}{2}\right)^{2n} \end{pmatrix}^{-1} \cdot \begin{pmatrix} \sum_{i=0}^n b_i \\ \sum_{i=0}^n \left(\frac{i}{m} - \frac{1}{2}\right) b_i \\ \sum_{i=0}^n \left(\frac{i}{m} - \frac{1}{2}\right)^2 b_i \\ \vdots \\ \sum_{i=0}^n \left(\frac{i}{m} - \frac{1}{2}\right)^n b_i \end{pmatrix} \quad 47$$

We provide numerical examples in the next section for a better visualization of eq. (47).

### 3.8.1. How Conditioning Changes (f)

The following implementation uses code built in Section 3.6 to showcase eq. (47):

BOTA A PORRA DO CODIGO

## 4. Some Algorithms

We have shown the solutions to the least squares problem  $Ax = b$ , but this problem could be solved with factorizations of  $A$ , such as the QR and SVD, in the following sections we will define these factorizations and use them to solve the least squares problem.

### 4.1. The SVD and QR factorizations

#### 4.1.1. QR

The QR factorization of  $A \in \mathbb{C}^{m \times n}$  consists of decomposing the co

#### 4.1.2. SVD

The *singular value decomposition* of a matrix is based on the fact that the image of the unit sphere under a  $m \times n$  matrix is a **hyperellipse**:

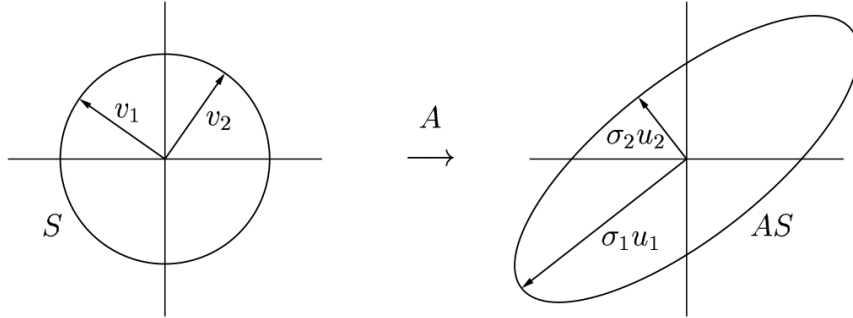


Figure 3: SVD of a  $2 \times 2$  matrix

So the independent directions  $v_1, v_2$  have been mapped to another set of orthogonal directions  $\sigma_1 v_1, \sigma_2 v_2$ , so with  $S := \{v \in \mathbb{C}^n \mid \|v\| = 1\}$  as the unit ball, let's define:

**Definition 4.1.2.1:** (Singular Values) The  $n$  *singular values*  $\sigma_i$  of  $A \in \mathbb{C}(m \times n)$  are the lengths of the  $n$  new axes of  $AS$ , written in non-crescent order  $\sigma_1 \geq \dots \geq \sigma_n$ .

**Definition 4.1.2.2:** (Left Singular Vectors) The  $n$  **left** singular vectors of  $A$  are the unit vectors  $u_i$  laying in  $AS$ , oriented to correspond and number the singular values  $\sigma_i$ , respectively

**Definition 4.1.2.3:** (Right Singular Vectors) The **right** singular vectors of  $A$  are the  $v_i$  in  $S$  that are the preimages of  $\sigma_i u_i \in AS$ , such that  $Av_i = \sigma_i u_i$

The equation  $Av_i = \sigma_i u_i$  is equivalent to:

$$A \cdot \begin{pmatrix} v_1 & v_2 & \dots & v_n \end{pmatrix} = \begin{pmatrix} \sigma_1 u_1 & \sigma_2 u_2 & \dots & \sigma_n u_n \end{pmatrix} \quad 48$$

Better:

$$A \cdot \begin{pmatrix} v_1 & v_2 & \dots & v_n \end{pmatrix} = \begin{pmatrix} u_1 & u_2 & \dots & u_n \end{pmatrix} \cdot \begin{pmatrix} \sigma_1 & 0 & \dots & 0 \\ 0 & \sigma_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma_n \end{pmatrix} \quad 49$$

Or simple  $AV = U\Sigma$ , but since  $V$  has orthonormal columns:

$$A = U\Sigma V^* \quad 50$$

The SVD is a very particular factorization for matrices, as the following theorem states:

**Theorem 4.1.2.1:** (Existence of SVD) Every matrix  $A \in \mathbb{C}^{m \times n}$  has a singular value decomposition

*Proof:* We proceed by fixing the largest image of  $A$  and using induction on the dimension of  $A$ :

Let  $\sigma_1 = \|A\|_2$ . There must exist unitary vectors  $u_1, v_1 \in \mathbb{C}^n$  such that  $Av_1 = \sigma_1 u_1$ . PROVAR ESSA PORRA DIREITO  $\square$

Is the SVD factorization of  $A$ . There are more about the SVD on computing  $U, \Sigma, V^*$ , as we will show below:

**Theorem 4.1.2.2:**  $\forall A \in \mathbb{C}^{m \times n}$ , the following holds:

- The eigenvalues of  $A^*A$  are the singular values *squared* of  $A$ , and the column-eigenvectors of  $A^*A$  form the matrix  $V$ .
- The eigenvalues of  $AA^*$  are the singular values *squared* of  $A$ , and the column-eigenvectors of  $AA^*$  form the matrix  $U$ .

*Proof:* PROVA ESSA PORRA DIREITO  $\square$

By Theorem 4.1.2.2, calculating the SVD of  $A$  has been reduced to calculating the eigenvalues and eigenvectors of  $A^*A$  and  $AA^*$ , here are some examples of singular value decompositions:

BOTA AS PORRA DOS EXEMPLO

*Example 4.1.2.1.:*

*Example 4.1.2.2.:*

*Example 4.1.2.3.:*

## 4.2. A Python-Implementation and a Conditioning Analysis (b)

Here we will solve the least squares problem using the 2 factorizations shown in Section 4.1.1 and Section 4.1.2, as well as the ordinary approach to least squares shown in Section 2.

We will also use these algorithms to do linear regression on the simple functions  $f, g, h : \mathbb{R} \rightarrow \mathbb{R}$  defined as:

$$f(t) = \sin(t)$$

$$g(t) = e^t$$

$$h(t) = \cos(3t)$$

51

BOTA OS CODIGO

## 4.3. The polynomial approach: An efficiency analysis (c)

Here we will analyse what happens when we do *polynomial* regression with the tools shown in Section 4.2. The same functions of eq. (51) will be used here, with polynomials of degree up to  $n = 15$ :

BOTA AS PORRA DOS CODIGO