

## Chapter 8: Least squared error

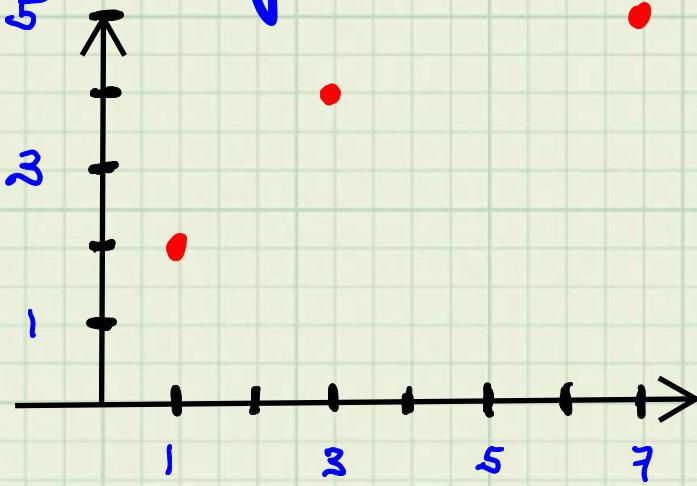
solutions to  $Ax = b$ ,  $b$  not a linear combination of the columns of  $A$  Big

**Application: Linear Regression =**

Fitting Functions to Data 

Machine Learning (AI) does a nonlinear version of "regression" = fitting functions to data!

### Summary



$x_i$	$y_i$
1	2
3	4
7	5

Linear model  $\rightarrow$  two unknowns  $\left[ \begin{matrix} m \\ b \end{matrix} \right]$

$$y_1 = mx_1 + b$$

$$y_2 = mx_2 + b$$

$$y_3 = mx_3 + b$$

$$\left[ \begin{matrix} y_1 \\ y_2 \\ y_3 \end{matrix} \right] = \left[ \begin{matrix} x_1 & 1 \\ x_2 & 1 \\ x_3 & 1 \end{matrix} \right] \left[ \begin{matrix} m \\ b \end{matrix} \right]$$

$$\underbrace{\begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix}}_Y = \underbrace{\begin{bmatrix} x_1 & 1 \\ x_2 & 1 \\ x_3 & 1 \end{bmatrix}}_{\Phi} \underbrace{\begin{bmatrix} m \\ b \end{bmatrix}}_x$$

$$\left\{ \begin{matrix} Y = \Phi \cdot \alpha \\ \end{matrix} \right.$$

$$\begin{bmatrix} 2 \\ 4 \\ 5 \end{bmatrix} = \begin{bmatrix} 1 & 1 \\ 3 & 1 \\ 7 & 1 \end{bmatrix} \begin{bmatrix} m \\ b \end{bmatrix}$$

$Y$        $\Phi$        $m$

When you are first learning something unfamiliar notation throws you off your game!

Re-write as

$$Y = \Phi \cdot \alpha$$

$$A \cdot x = b$$

$A = \Phi$  (regressor matrix)

$x = \alpha = \begin{bmatrix} m \\ b \end{bmatrix}$  unknowns

$b = Y$  measured data

- There is no line that goes through all data points
- There is no solution to our system of linear equations
- We want the line that "best fits"

the data !!!

Write  $e(x) := Ax - b$

= error in the solution,  
and it depends on the  
 $x$  we select

• For some  $x$ ,  $e(x) := Ax - b$  might  
be HUGE

For other  $x$ ,  $e(x) := Ax - b$  might be  
smaller.

**Question:** Can we find an  
 $x$  that makes the norm of the  
error as small as possible ???

If so, is it unique ?

**Answer:** Suppose the columns of

$A$  are linearly independent. Then  
there exists  $x^* \in \mathbb{R}^m$  such that

$\|e(x^*)\| = \|Ax^*-b\| \leq \|Ax-b\|$  for all  $x \in \mathbb{R}^m$ .

Moreover,  $x^*$  is unique.

Meaning:

$$A = n \times m, \quad b = n \times 1$$

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

$x^*$  minimizes the norm of the error.

One writes  $x^* = \arg \min_{x \in \mathbb{R}^m} \|Ax-b\|$

$x^*$  is that vector in  $\mathbb{R}^m$  achieving the minimum error.

How to Compute it?

Fact : Suppose the columns of  $A$  are linearly independent.

Then TFAE

$$\textcircled{1} \quad x^* = \arg \min_{x \in \mathbb{R}^m} \|Ax - b\|$$

$$\textcircled{2} \quad A^T A x^* = A^T b$$

$$\textcircled{3} \quad x^* = (A^T A)^{-1} A^T b$$

Recall :  $\det(A^T A) \neq 0 \iff$  columns of  $A$  are lin. indep.

## Our Line Fitting Problem

$$A = \begin{bmatrix} 1 & 1 \\ 3 & 1 \\ 7 & 1 \end{bmatrix}$$

$$b = \begin{bmatrix} 2 \\ 4 \\ 5 \end{bmatrix}$$

$$x = \begin{bmatrix} m \\ b \\ \alpha \end{bmatrix}$$

$$Ax = b \iff Y = \Phi \cdot \alpha$$

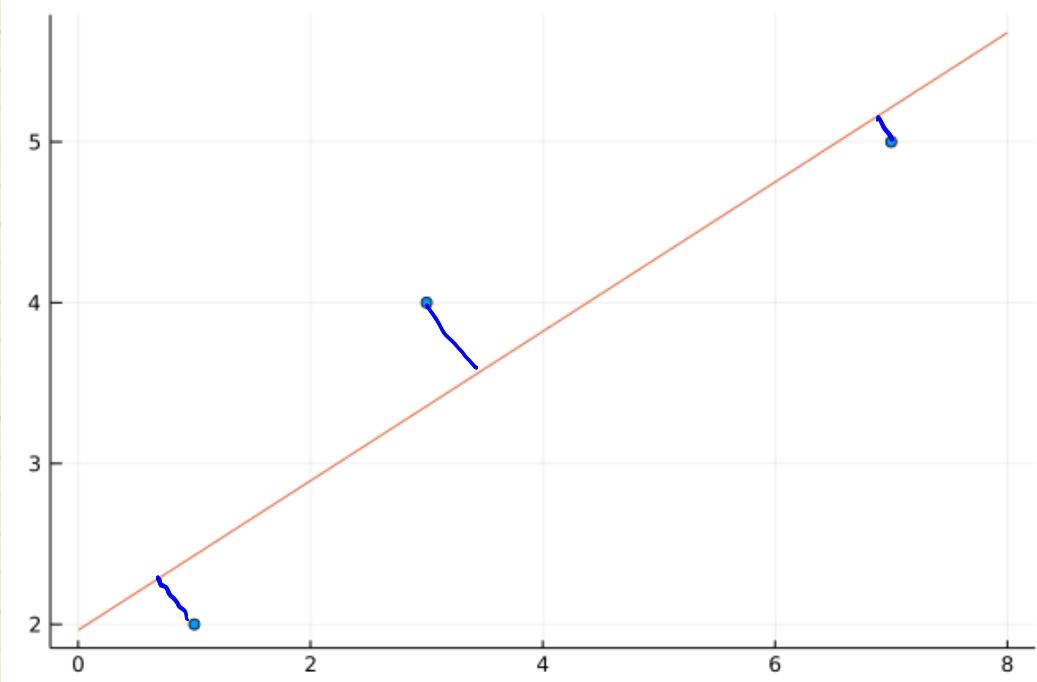
$$(A^T A)x^* = A^T b$$

$$A^T A = \begin{bmatrix} 1 & 3 & 7 \\ 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} 1 & 1 \\ 3 & 1 \\ 7 & 1 \end{bmatrix} = \begin{bmatrix} 59 & 11 \\ 11 & 3 \end{bmatrix}$$

$$A^T b = \begin{bmatrix} 1 & 3 & 7 \\ 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} 2 \\ 4 \\ 5 \end{bmatrix} = \begin{bmatrix} 49 \\ 11 \end{bmatrix}$$

$$x^* = \begin{bmatrix} 0.4b \\ 1.9b \end{bmatrix}$$

$$\begin{aligned} m &= 0.4b \\ b &= 1.9b \end{aligned}$$



# Why True?

$$\textcircled{1} \quad x^* = \arg \min_{x \in \mathbb{R}^m} \|Ax - b\|$$

$$\Updownarrow \quad x^* = \arg \min_{x \in \mathbb{R}^m} \|Ax - b\|^2$$

$$\|Ax - b\|^2 = (Ax - b)^T (Ax - b)$$

Some painful algebra to "complete the square"

$$(Ax - b)^T (Ax - b) = (A^T A x - A^T b)^T (A^T A)^{-1} (A^T A x - A^T b) + \\ + \underbrace{b^T b - b^T A (A^T A)^{-1} A^T b}_{\text{does not depend on } x}$$

minimize with

$$A^T A x - A^T b = 0$$

Someone asked: Why not minimize with a large negative value instead of zero?

Answer:  $(A^T A)^{-1}$  is what we call a positive definite matrix (see Appendix of Book)

Such matrices have the property that  $y^T (A^T A) y \geq 0$  for all  $y \in \mathbb{R}^m$  and  $y^T (A^T A) y = 0 \iff y = 0$

This a bit advanced for ROB 101.







