

3

Orthogonality (正交性)

3.3

PROJECTIONS AND LEAST SQUARES (投影与最小二乘)

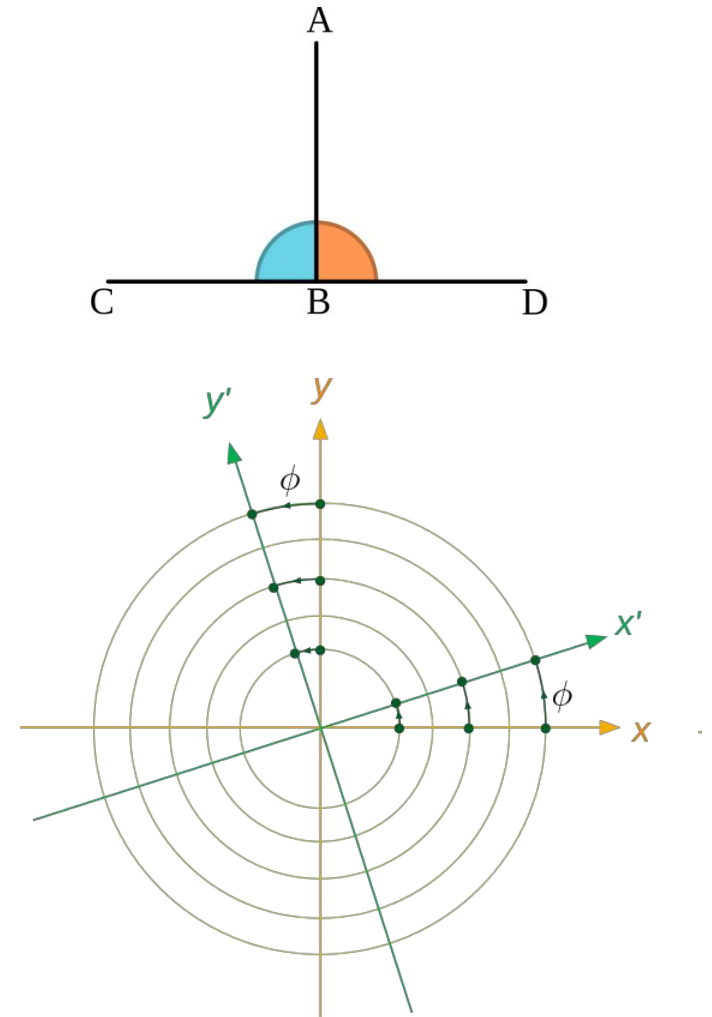
Least squares

Projection Matrix

The matrix $A^T A$

Fitting of Data (数据拟合)

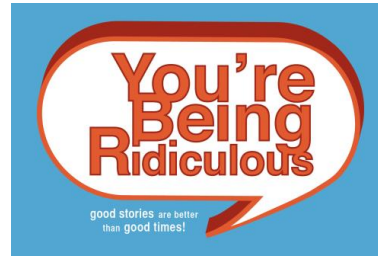
合)



Introduction

Solve $Ax = b$

when there is no solution?



IT HAPPENS ALL THE TIME!

(e.g., it often happens when $m > n$)

In spite of their unsolvability, inconsistent equations arise all the time in practice.

They have to be solved!

Elimination is going to fail.

Throw away some equations? This is hard to justify if all m equations come from the same source.

Better way:

choose the x that minimizes an average error E in the m equations.



We start with an example.

Example 1. The system of linear equations with only one unknown x

$$\begin{cases} 2x = b_1, \\ 3x = b_2, \\ 4x = b_3. \end{cases}$$

has solutions if and only if (b_1, b_2, b_3) is proportional to $(2, 3, 4)$.

If it has no solution, we wish to find a value \hat{x} *which minimizes the difference*

$$E^2 = (2x - b_1)^2 + (3x - b_2)^2 + (4x - b_3)^2.$$

Differentiating it with respect to x , we have

$$\frac{dE^2}{dx} = 4(2x - b_1) + 6(3x - b_2) + 8(4x - b_3)$$

equating 0 leads to

$$\hat{x} = \frac{2b_1 + 3b_2 + 4b_3}{2^2 + 3^2 + 4^2} = \frac{\mathbf{a}^T \mathbf{b}}{\mathbf{a}^T \mathbf{a}},$$

where $\mathbf{a} = (2, 3, 4)^T$ and $\mathbf{b} = (b_1, b_2, b_3)^T$.

The general case is the same. We “solve” $\mathbf{ax} = \mathbf{b}$ by minimizing

$$E^2 = \|\mathbf{ax} - \mathbf{b}\|^2 = (a_1x - b_1)^2 + (a_2x - b_2)^2 + \dots + (a_mx - b_m)^2.$$

The derivative of E^2 is zero at the point \hat{x} , if

$$(a_1x - b_1)a_1 + (a_2x - b_2)a_2 + \dots + (a_mx - b_m)a_m = 0.$$

We are *minimizing the distance* from \mathbf{b} to the line through \mathbf{a} , and calculus gives the same answer,

$$\hat{x} = \frac{a_1b_1 + a_2b_2 + \dots + a_mb_m}{a_1^2 + a_2^2 + \dots + a_m^2} = \frac{\mathbf{a}^T \mathbf{b}}{\mathbf{a}^T \mathbf{a}},$$

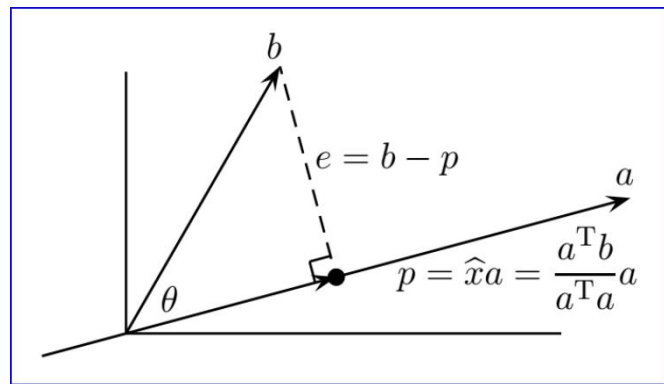
where $\mathbf{a} = (a_1, a_2, \dots, a_m)^T$ and $\mathbf{b} = (b_1, b_2, \dots, b_m)^T$.

Lemma

The least-squares solution to a problem $\mathbf{ax} = \mathbf{b}$ in one unknown is

$$\hat{x} = \frac{\mathbf{a}^T \mathbf{b}}{\mathbf{a}^T \mathbf{a}}.$$

(The error vector \mathbf{e} connecting \mathbf{b} to \mathbf{p} must be perpendicular to \mathbf{a})



I. Best Approximation – Least Squares Solution

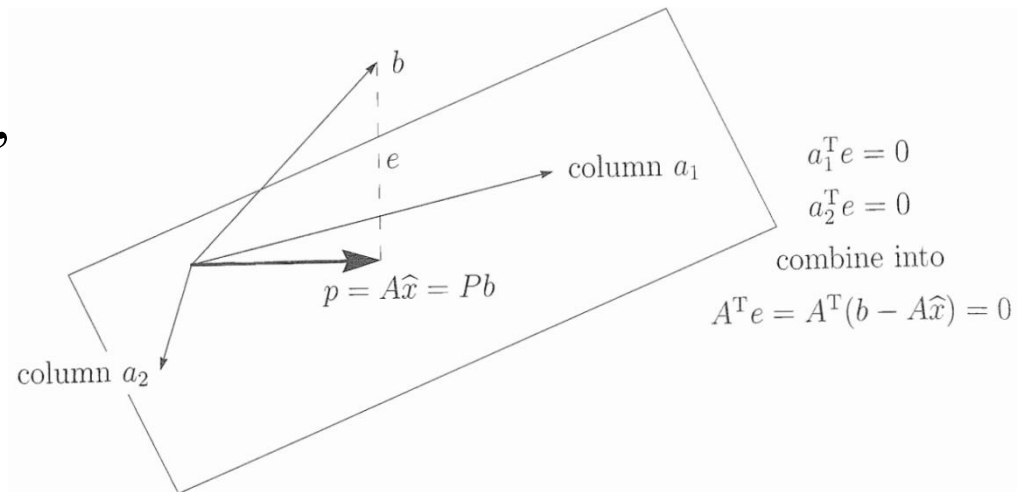
Our next task is to establish a similar formula in the general case (Let A be an $(m \times n)$ -matrix).

Suppose that $Ax = b$ is inconsistent. We wish to find \hat{x} such that $\|A\hat{x} - b\|$ is as small as possible.

The column space $C(A)$ consists of all vectors of the form Ax with $x \in \mathbb{R}^n$, i.e., $C(A) = \{Ax \mid x \in \mathbb{R}^n\}$.

Since $Ax = b$ has no solution, we have $b \neq Ax$ for any $x \in \mathbb{R}^n$, that is, the vector b is not in the column space $C(A)$.

Among the vectors of $C(A)$, the vector b is nearest to the projection of b in $C(A)$.



Projection onto the column space of a 3 by 2 matrix

Namely, $A\hat{x}$ should be the projection of b in $C(A)$.

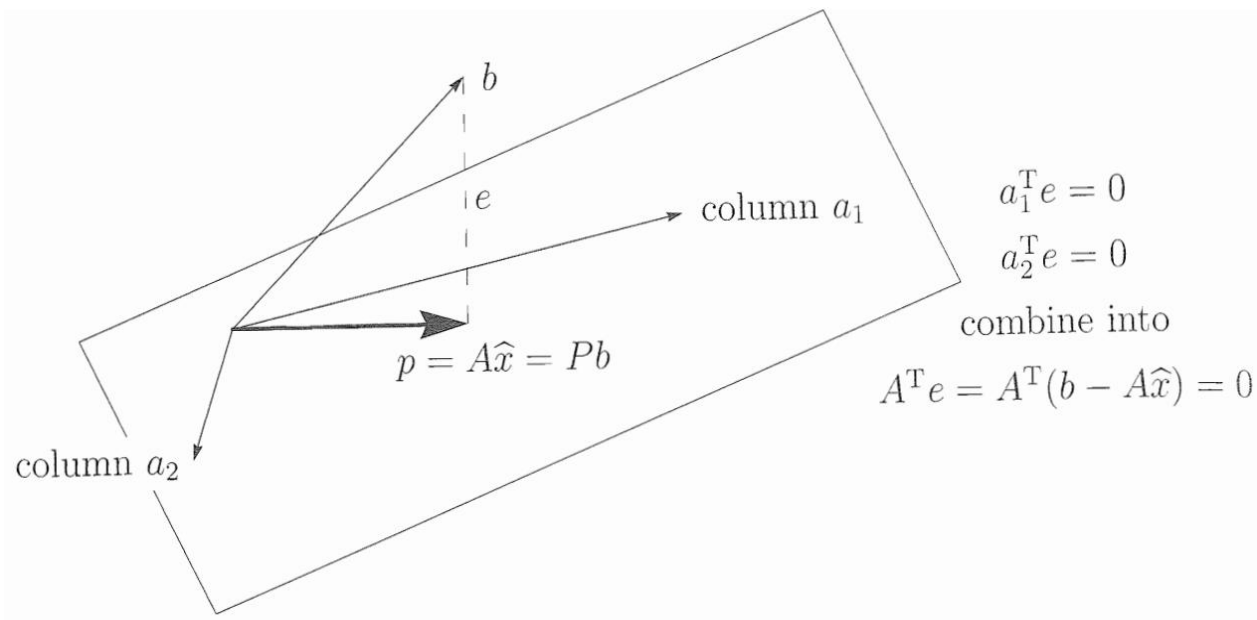
In other words, the difference $b - A\hat{x}$ is perpendicular to $C(A)$, i.e.,

$$b - A\hat{x} \perp C(A).$$

Equivalently, $b - A\hat{x}$ is perpendicular to all columns of A .

Since all vectors perpendicular to the column space lie in the *left nullspace*, thus $b - A\hat{x}$ lies in the left nullspace of A , so

$$A^T(b - A\hat{x}) = 0.$$



Projection onto the column space of a 3 by 2 matrix

Theorem 1. *If a system $A\mathbf{x} = \mathbf{b}$ is inconsistent (has no solution), its least-squares solution minimizes $\|A\mathbf{x} - \mathbf{b}\|^2$:*

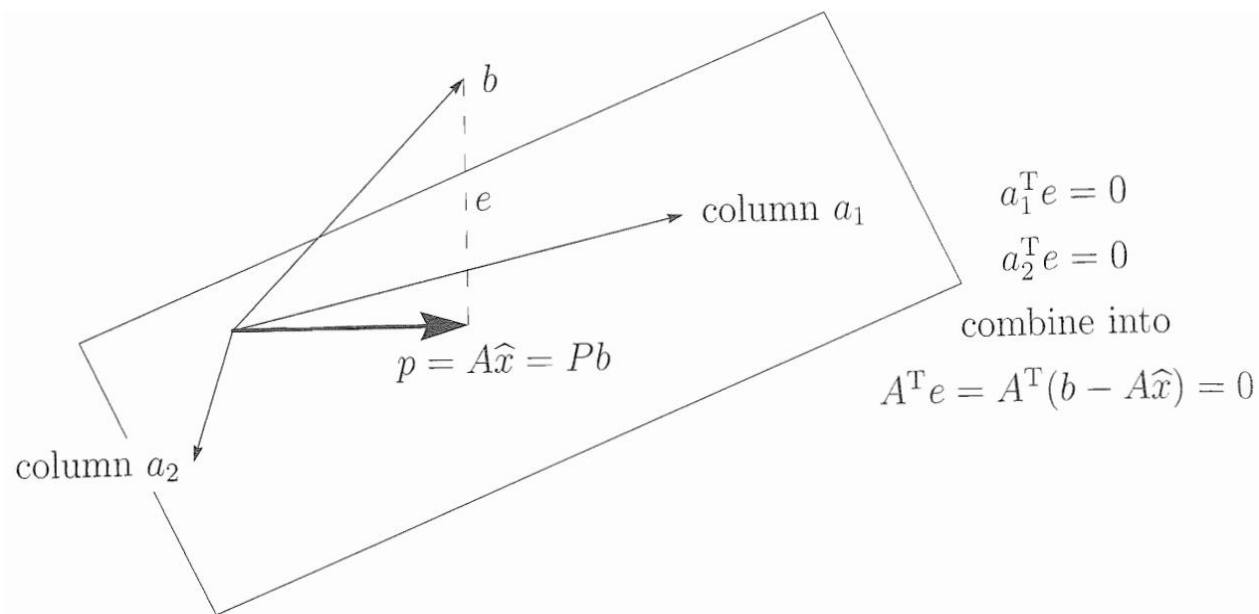
$$A^T A \hat{\mathbf{x}} = A^T \mathbf{b}. \quad \text{(Normal equations)}$$

Moreover, if $A^T A$ is invertible, then

$$\hat{\mathbf{x}} = (A^T A)^{-1} A^T \mathbf{b}. \quad \text{(Best estimate)}$$

The projection of \mathbf{b} onto the column space is the nearest point $A\hat{\mathbf{x}}$:

$$\mathbf{p} = A\hat{\mathbf{x}} = A(A^T A)^{-1} A^T \mathbf{b}. \quad \text{(Projection)}$$



Projection onto the column space of a 3 by 2 matrix

Example 2. The following system of linear equations $A\mathbf{x} = \mathbf{b}$ has no solution

$$\begin{bmatrix} 1 & 1 & 0 \\ 1 & 2 & 3 \\ 0 & 1 & 3 \end{bmatrix} \begin{bmatrix} x \\ y \\ z \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \\ 5 \end{bmatrix}.$$

We wish to find an approximation of \mathbf{x} , denoted by $\hat{\mathbf{x}}$, such that $A\hat{\mathbf{x}}$ is very close to \mathbf{b} , that is, such that the ‘error’ $\|\mathbf{b} - A\hat{\mathbf{x}}\|$ is as small as possible.

There is a formula for finding such an approximation $\hat{\mathbf{x}}$:

$$A^T A \hat{\mathbf{x}} = A^T \mathbf{b}.$$

Let $\hat{\mathbf{x}} = (x_0, y_0, z_0)^T$. Then

$$\begin{bmatrix} 1 & 1 & 0 \\ 1 & 2 & 1 \\ 0 & 3 & 3 \end{bmatrix} \begin{bmatrix} 1 & 1 & 0 \\ 1 & 2 & 3 \\ 0 & 1 & 3 \end{bmatrix} \begin{bmatrix} x_0 \\ y_0 \\ z_0 \end{bmatrix} = \begin{bmatrix} 2 & 3 & 3 \\ 3 & 6 & 9 \\ 3 & 9 & 18 \end{bmatrix} \begin{bmatrix} x_0 \\ y_0 \\ z_0 \end{bmatrix} = \begin{bmatrix} 1 & 1 & 0 \\ 1 & 2 & 1 \\ 0 & 3 & 3 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \\ 5 \end{bmatrix} = \begin{bmatrix} 1 \\ 6 \\ 15 \end{bmatrix},$$

Using the typical method, we can find solutions for (x_0, y_0, z_0) , *one of* which is $(-4, 3, 0)$.

Remark 1. Suppose \mathbf{b} is actually in the column space of \mathbf{A} —it is a combination $\mathbf{b} = \mathbf{A}\mathbf{x}$ of the columns. Then the projection of \mathbf{b} is still \mathbf{b} :

$$\mathbf{b} \text{ in column space} \quad \mathbf{p} = \mathbf{A}(\mathbf{A}^T\mathbf{A})^{-1}\mathbf{A}^T\mathbf{A}\mathbf{x} = \mathbf{A}\mathbf{x} = \mathbf{b}.$$

The closest point \mathbf{p} is just \mathbf{b} itself—which is obvious.

Remark 2. At the other extreme, suppose \mathbf{b} is *perpendicular* to every column, so $\mathbf{A}^T\mathbf{b} = \mathbf{0}$. In this case \mathbf{b} projects to the zero vector:

$$\mathbf{b} \text{ in left nullspace} \quad \mathbf{p} = \mathbf{A}(\mathbf{A}^T\mathbf{A})^{-1}\mathbf{A}^T\mathbf{b} = \mathbf{A}(\mathbf{A}^T\mathbf{A})^{-1}\mathbf{0} = \mathbf{0}.$$

Remark 3. When \mathbf{A} is square and invertible, the column space is the whole space. Every vector projects to itself, \mathbf{p} equals \mathbf{b} , and $\hat{\mathbf{x}} = \mathbf{x}$:

$$\text{If } \mathbf{A} \text{ is invertible} \quad \mathbf{p} = \mathbf{A}(\mathbf{A}^T\mathbf{A})^{-1}\mathbf{A}^T\mathbf{b} = \mathbf{A}\mathbf{A}^{-1}(\mathbf{A}^T)^{-1}\mathbf{A}^T\mathbf{b} = \mathbf{b}.$$

When \mathbf{A} is rectangular that is not possible.

Remark 4. Suppose \mathbf{A} has only one column, containing \mathbf{a} . Then the matrix $\mathbf{A}^T\mathbf{A}$ is the number $\mathbf{a}^T\mathbf{a}$ and $\hat{\mathbf{x}}$ is $\frac{\mathbf{a}^T\mathbf{b}}{\mathbf{a}^T\mathbf{a}}$. We return to the earlier formula.

II. Projection Matrices and the Matrix $A^T A$

For a matrix A of size $m \times n$, the column space $C(A)$ is a subspace of \mathbf{R}^m . Is there a projection from \mathbf{R}^m to $C(A)$?

As mentioned above, $A\hat{\mathbf{x}}$ is the vector in $C(A)$ that is the projection of \mathbf{b} in $C(A)$, and is nearest to \mathbf{b} , which is

$$A\hat{\mathbf{x}} = A(A^T A)^{-1} A^T \mathbf{b}$$

So the matrix $A(A^T A)^{-1} A^T$ is a transformation: $\mathbf{b} \mapsto A\hat{\mathbf{x}}$.

Corollary. *Let A be a matrix of size $m \times n$ such that $A^T A$ is invertible. Then the projection of \mathbf{R}^m to $C(A)$ has matrix*

$$A(A^T A)^{-1} A^T.$$

It is easy to see that the matrix $P = A(A^T A)^{-1} A^T$ is symmetric, and

$$P^2 = A(A^T A)^{-1} A^T A(A^T A)^{-1} A^T = A(A^T A)^{-1} A^T = P.$$

Theorem 2. Any symmetric matrix \mathbf{P} with $\mathbf{P}^2 = \mathbf{P}$ represents a projection. (任何对称的幂等矩阵都对应一个投影变换)

Proof. Let \mathbf{P} have size $m \times m$.

We claim that \mathbf{P} projects \mathbf{R}^m to the column space $C(\mathbf{P})$.

Let $\mathbf{v} \in \mathbf{R}^m$. Consider the difference $\mathbf{v} - \mathbf{P}\mathbf{v}$. A vector of $C(\mathbf{P})$ has the form $\mathbf{P}\mathbf{x}$, and

$$(\mathbf{P}\mathbf{x})^T(\mathbf{v} - \mathbf{P}\mathbf{v}) = \mathbf{x}^T\mathbf{P}^T\mathbf{v} - \mathbf{x}^T\mathbf{P}^T\mathbf{P}\mathbf{v} = \mathbf{x}^T(\mathbf{P} - \mathbf{P}^2)\mathbf{v} = 0.$$

Thus $\mathbf{v} - \mathbf{P}\mathbf{v}$ is perpendicular to $\mathbf{P}\mathbf{x}$, and so $\mathbf{P}\mathbf{v}$ is the projection of \mathbf{v} in $C(\mathbf{P})$.

Remark. In this case, $\mathbf{I} - \mathbf{P}$ also represents a projection.

There is a simple way to decide *whether $A^T A$ is invertible*.

Theorem 3. *The matrices $A^T A$ and A have the same nullspace. In particular, if A has full column rank, then $A^T A$ is invertible.*

Proof. If $A\mathbf{x} = \mathbf{0}$, then $A^T A\mathbf{x} = \mathbf{0}$, and $N(A) \subseteq N(A^T A)$.

Conversely, suppose that $A^T A\mathbf{x} = \mathbf{0}$. Then

$$\|A\mathbf{x}\|^2 = (A\mathbf{x})^T (A\mathbf{x}) = \mathbf{x}^T A^T A\mathbf{x} = \mathbf{x}^T \mathbf{0} = 0,$$

which implies that $A\mathbf{x} = \mathbf{0}$.

Thus $N(A^T A) \subseteq N(A)$, and so $N(A^T A) = N(A)$.

Let A be of size $m \times n$. Then $A^T A$ is of size $n \times n$.

If A has column rank n , then $N(A) = \{\mathbf{0}\}$, and so $N(A^T A) = \{\mathbf{0}\}$.

Therefore, $A^T A$ is invertible.

Note: $A^T A$ is square and symmetric.

Moreover, if A has independent columns, then $A^T A$ is *square, symmetric, and invertible*.

(如果 A 的列线性无关, 则 $A^T A$ 是方阵、对称而且可逆)

For Example,

$$\mathbf{A} = \begin{bmatrix} 1 & 1 \\ 1 & 2 \\ 1 & 5 \end{bmatrix}, \text{ then}$$

$$\mathbf{A}^T \mathbf{A} = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 2 & 5 \end{bmatrix} \begin{bmatrix} 1 & 1 \\ 1 & 2 \\ 1 & 5 \end{bmatrix} = \begin{bmatrix} 3 & 8 \\ 8 & 30 \end{bmatrix} \text{ is invertible.}$$

$$\mathbf{B} = \begin{bmatrix} 1 & 3 \\ 1 & 3 \\ 1 & 3 \end{bmatrix}, \text{ then}$$

$$\mathbf{B}^T \mathbf{B} = \begin{bmatrix} 1 & 1 & 1 \\ 3 & 3 & 3 \end{bmatrix} \begin{bmatrix} 1 & 3 \\ 1 & 3 \\ 1 & 3 \end{bmatrix} = \begin{bmatrix} 3 & 9 \\ 9 & 27 \end{bmatrix} \text{ is *not* invertible.}$$

III. Least Squares Fitting of Data (数据的最小二乘拟合)

Look at a line

$$C + Dt = b.$$

If there is no experimental error, then two measurements of b will determine the line.

But if there is error, by a series of experiments, we obtain a system with **two unknowns** C , D :

$$\begin{cases} C + Dt_1 = b_1, \\ C + Dt_2 = b_2, \\ \dots \\ C + Dt_m = b_m. \end{cases}$$

In matrix form $A\mathbf{x} = \mathbf{b}$, where

$$A = \begin{bmatrix} 1 & t_1 \\ 1 & t_2 \\ \vdots & \vdots \\ 1 & t_m \end{bmatrix}, \quad \mathbf{x} = \begin{bmatrix} C \\ D \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_m \end{bmatrix}.$$

If the system $\mathbf{Ax} = \mathbf{b}$ is inconsistent, then we wish to find $\hat{\mathbf{x}} = (\hat{C}, \hat{D})^T$ to minimize the squared error E^2 :

$$E^2 = \|\mathbf{b} - \mathbf{Ax}\|^2 = (b_1 - C - Dt_1)^2 + \dots + (b_m - C - Dt_m)^2.$$

By the theorem obtained above, we have

$$\mathbf{A}^T \mathbf{A} \hat{\mathbf{x}} = \mathbf{A}^T \mathbf{b}.$$

which is

$$\begin{bmatrix} m & t_1 + \dots + t_m \\ t_1 + \dots + t_m & t_1^2 + \dots + t_m^2 \end{bmatrix} \begin{bmatrix} \hat{C} \\ \hat{D} \end{bmatrix} = \begin{bmatrix} b_1 + \dots + b_m \\ t_1 b_1 + \dots + t_m b_m \end{bmatrix}.$$

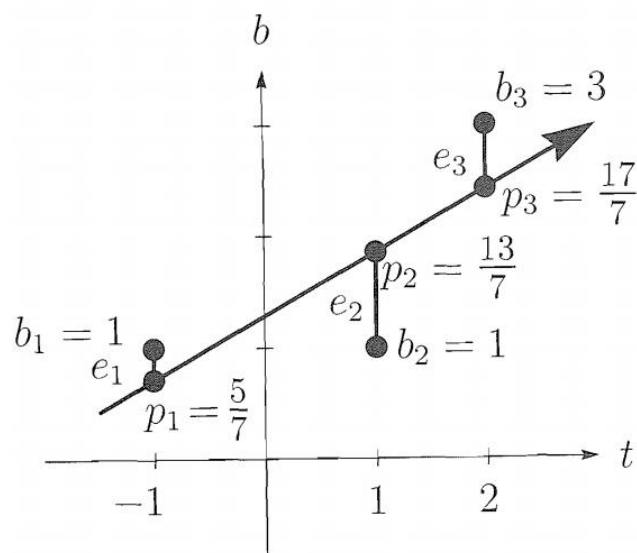
Theorem 4. *Let the measurements b_1, b_2, \dots, b_m be given at distinct points t_1, t_2, \dots, t_m . Then the line*

$$\hat{C} + \hat{D}t = b$$

with $\hat{\mathbf{x}} = (\hat{C}, \hat{D})^T$ which minimizes E^2 comes from least squares:

$$\mathbf{A}^T \mathbf{A} \hat{\mathbf{x}} = \mathbf{A}^T \mathbf{b}.$$

Example 3 Three measurements b_1, b_2, b_3 are marked on the figure:



$b = 1$ at $t = -1$; $b = 1$ at $t = 1$; $b = 3$ at $t = 2$.

$$\begin{bmatrix} 1 & -1 \\ 1 & 1 \\ 1 & 2 \end{bmatrix} \begin{bmatrix} C \\ D \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \\ 3 \end{bmatrix}.$$

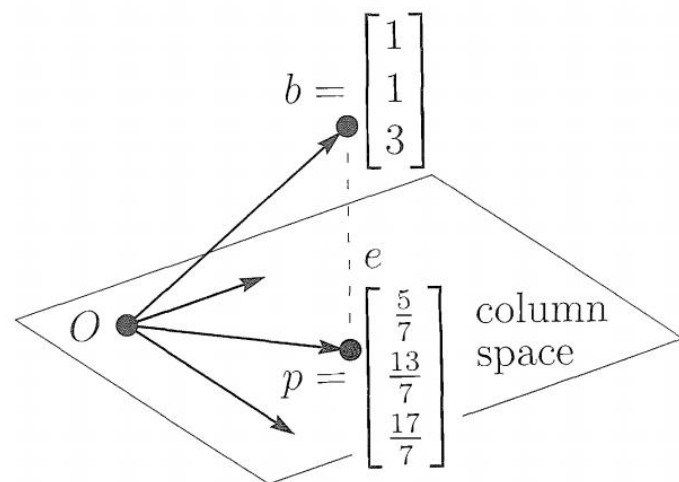
$A\mathbf{x} = \mathbf{b}$ can't be solved because the points are not on a line.

Therefore they are solved by least squares:

$$A^T A \hat{\mathbf{x}} = A^T \mathbf{b}.$$

is $\begin{bmatrix} 3 & 2 \\ 2 & 6 \end{bmatrix} \begin{bmatrix} C \\ D \end{bmatrix} = \begin{bmatrix} 5 \\ 6 \end{bmatrix} \Rightarrow \hat{C} = \frac{9}{7}, \hat{D} = \frac{4}{7}.$

The best line is $\frac{9}{7} + \frac{4}{7}t$.



Remark. The mathematics of least squares is not limited to fitting the data by straight lines—*nonlinear least squares*.

Key words:

Least Squares

The matrix $A^T A$; Projection matrix

Fitting of Data

Homework

See Blackboard

