

# MATH 532: Linear Algebra

## Chapter 4: Vector Spaces

Greg Fasshauer

Department of Applied Mathematics  
Illinois Institute of Technology

Spring 2015



# Outline

- 1 Spaces and Subspaces
- 2 Four Fundamental Subspaces
- 3 Linear Independence
- 4 Bases and Dimension
- 5 More About Rank
- 6 Classical Least Squares
- 7 Kriging as best linear unbiased predictor



# Spaces and Subspaces

While the discussion of vector spaces can be rather dry and abstract, they are an essential tool for describing the world we work in, and to understand many practically relevant consequences.

After all, linear algebra is pretty much the workhorse of modern applied mathematics.

Moreover, many concepts we discuss now for traditional “vectors” apply also to vector spaces of functions, which form the foundation of functional analysis.



# Vector Space

## Definition

A set  $\mathcal{V}$  of elements (**vectors**) is called a **vector space** (or linear space) over the scalar field  $\mathcal{F}$  if

- (A1)  $\mathbf{x} + \mathbf{y} \in \mathcal{V}$  for any  $\mathbf{x}, \mathbf{y} \in \mathcal{V}$   
(closed under addition),
- (A2)  $(\mathbf{x} + \mathbf{y}) + \mathbf{z} = \mathbf{x} + (\mathbf{y} + \mathbf{z})$  for all  $\mathbf{x}, \mathbf{y}, \mathbf{z} \in \mathcal{V}$ ,
- (A3)  $\mathbf{x} + \mathbf{y} = \mathbf{y} + \mathbf{x}$  for all  $\mathbf{x}, \mathbf{y} \in \mathcal{V}$ ,
- (A4) There exists a **zero vector**  $\mathbf{0} \in \mathcal{V}$  such that  $\mathbf{x} + \mathbf{0} = \mathbf{x}$  for every  $\mathbf{x} \in \mathcal{V}$ ,
- (A5) For every  $\mathbf{x} \in \mathcal{V}$  there is a **negative**  $(-\mathbf{x}) \in \mathcal{V}$  such that  $\mathbf{x} + (-\mathbf{x}) = \mathbf{0}$ ,
- (M1)  $\alpha\mathbf{x} \in \mathcal{V}$  for every  $\alpha \in \mathcal{F}$  and  $\mathbf{x} \in \mathcal{V}$  (closed under scalar multiplication),
- (M2)  $(\alpha\beta)\mathbf{x} = \alpha(\beta\mathbf{x})$  for all  $\alpha\beta \in \mathcal{F}$ ,  $\mathbf{x} \in \mathcal{V}$ ,
- (M3)  $\alpha(\mathbf{x} + \mathbf{y}) = \alpha\mathbf{x} + \alpha\mathbf{y}$  for all  $\alpha \in \mathcal{F}$ ,  $\mathbf{x}, \mathbf{y} \in \mathcal{V}$ ,
- (M4)  $(\alpha + \beta)\mathbf{x} = \alpha\mathbf{x} + \beta\mathbf{x}$  for all  $\alpha, \beta \in \mathcal{F}$ ,  $\mathbf{x} \in \mathcal{V}$ ,
- (M5)  $1\mathbf{x} = \mathbf{x}$  for all  $\mathbf{x} \in \mathcal{V}$ .

# Examples of vector spaces

- $\mathcal{V} = \mathbb{R}^m$  and  $\mathcal{F} = \mathbb{R}$  (traditional **real vectors**)
- $\mathcal{V} = \mathbb{C}^m$  and  $\mathcal{F} = \mathbb{C}$  (traditional **complex vectors**)
- $\mathcal{V} = \mathbb{R}^{m \times n}$  and  $\mathcal{F} = \mathbb{R}$  (**real matrices**)
- $\mathcal{V} = \mathbb{C}^{m \times n}$  and  $\mathcal{F} = \mathbb{C}$  (**complex matrices**)

But also

- $\mathcal{V}$  is **polynomials** of a certain degree with real coefficients,  $\mathcal{F} = \mathbb{R}$
- $\mathcal{V}$  is **continuous functions** on an interval  $[a, b]$ ,  $\mathcal{F} = \mathbb{R}$



# Subspaces

## Definition

Let  $\mathcal{S}$  be a nonempty **subset** of  $\mathcal{V}$ . If  $\mathcal{S}$  is a vector space, then  $\mathcal{S}$  is called a **subspace** of  $\mathcal{V}$ .

**Q:** What is the difference between a **subset** and a **subspace**?

**A:** The **structure** provided by the axioms (A1)–(A5), (M1)–(M5)

## Theorem

The **subset**  $\mathcal{S} \subseteq \mathcal{V}$  is a **subspace** of  $\mathcal{V}$  if and only if

$$\alpha \mathbf{x} + \beta \mathbf{y} \in \mathcal{S} \quad \text{for all } \mathbf{x}, \mathbf{y} \in \mathcal{S}, \alpha, \beta \in \mathcal{F}. \quad (1)$$

## Remark

$\mathcal{Z} = \{\mathbf{0}\}$  is called the **trivial subspace**.

## Proof.

“ $\implies$ ”: Clear, since we actually have

$$(1) \iff (A1) \text{ and } (M1)$$

“ $\impliedby$ ”: Only (A1), (A4), (A5) and (M1) need to be checked (why?).

In fact, we see that (A1) and (M1) imply (A4) and (A5):

If  $\mathbf{x} \in \mathcal{S}$ , then — using (M1) —  $-1\mathbf{x} = -\mathbf{x} \in \mathcal{S}$ , i.e., (A5) holds.

Using (A1),  $\mathbf{x} + (-\mathbf{x}) = \mathbf{0} \in \mathcal{S}$ , so that (A4) holds. □



## Definition

Let  $\mathcal{S} = \{\mathbf{v}_1, \dots, \mathbf{v}_r\} \subseteq \mathcal{V}$ . The **span** of  $\mathcal{S}$  is

$$\text{span}(\mathcal{S}) = \left\{ \sum_{i=1}^r \alpha_i \mathbf{v}_i : \alpha_i \in \mathcal{F} \right\}.$$

## Remark

- $\text{span}(\mathcal{S})$  contains *all possible linear combinations of vectors in  $\mathcal{S}$* .
- One can easily show that  $\text{span}(\mathcal{S})$  is a *subspace* of  $\mathcal{V}$ .

## Example (Geometric interpretation)

- 1 If  $\mathcal{S} = \{\mathbf{v}_1\} \subseteq \mathbb{R}^3$ , then  $\text{span}(\mathcal{S})$  is the line through the origin with direction  $\mathbf{v}_1$ .
- 2 If  $\mathcal{S} = \{\mathbf{v}_1, \mathbf{v}_2 : \mathbf{v}_1 \neq \alpha \mathbf{v}_2, \alpha \neq 0\} \subseteq \mathbb{R}^3$ , then  $\text{span}(\mathcal{S})$  is the plane through the origin “spanned by”  $\mathbf{v}_1$  and  $\mathbf{v}_2$ .



## Definition

Let  $\mathcal{S} = \{\mathbf{v}_1, \dots, \mathbf{v}_r\} \subseteq \mathcal{V}$ . If  $\text{span } \mathcal{S} = \mathcal{V}$  then  $\mathcal{S}$  is called a **spanning set** for  $\mathcal{V}$ .

## Remark

- A spanning set is sometimes referred to as a (finite) **frame**.
- A spanning set is **not the same as a basis** since the spanning set **may include redundancies**.

## Example

- $\left\{ \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix} \right\}$  is a spanning set for  $\mathbb{R}^3$ .
- $\left\{ \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix}, \begin{pmatrix} 2 \\ 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 2 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 0 \\ 2 \end{pmatrix} \right\}$  is also a spanning set for  $\mathbb{R}^3$ .

## Connection to linear systems

### Theorem

Let  $\mathcal{S} = \{\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n\}$  be the set of columns of an  $m \times n$  matrix  $A$ .  $\text{span}(\mathcal{S}) = \mathbb{R}^m$  if and only if for every  $\mathbf{b} \in \mathbb{R}^m$  there exists an  $\mathbf{x} \in \mathbb{R}^n$  such that  $A\mathbf{x} = \mathbf{b}$  (i.e., if and only if  $A\mathbf{x} = \mathbf{b}$  is consistent for every  $\mathbf{b} \in \mathbb{R}^m$ ).

### Proof.

By definition,  $\mathcal{S}$  is a spanning set for  $\mathbb{R}^m$  if and only if for every  $\mathbf{b} \in \mathbb{R}^m$  there exist  $\alpha_1, \dots, \alpha_n \in \mathbb{R}$  such that

$$\mathbf{b} = \alpha_1 \mathbf{a}_1 + \dots + \alpha_n \mathbf{a}_n = A\mathbf{x},$$

where  $A = \begin{pmatrix} \mathbf{a}_1 & \mathbf{a}_2 & \cdots & \mathbf{a}_n \end{pmatrix}_{m \times n}$  and  $\mathbf{x} = \begin{pmatrix} \alpha_1 \\ \vdots \\ \alpha_n \end{pmatrix}$ . □

**Remark**

The *sum*

$$\mathcal{X} + \mathcal{Y} = \{\mathbf{x} + \mathbf{y} : \mathbf{x} \in \mathcal{X}, \mathbf{y} \in \mathcal{Y}\}$$

is a subspace of  $\mathcal{V}$  provided  $\mathcal{X}$  and  $\mathcal{Y}$  are subspaces.

If  $S_{\mathcal{X}}$  and  $S_{\mathcal{Y}}$  are spanning sets for  $\mathcal{X}$  and  $\mathcal{Y}$ , respectively, then  $S_{\mathcal{X}} \cup S_{\mathcal{Y}}$  is a spanning set for  $\mathcal{X} + \mathcal{Y}$ .



# Four Fundamental Subspaces

Recall that a **linear function**  $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$  satisfies

$$f(\alpha \mathbf{x} + \beta \mathbf{y}) = \alpha f(\mathbf{x}) + \beta f(\mathbf{y}) \quad \forall \alpha, \beta \in \mathbb{R}, \mathbf{x}, \mathbf{y} \in \mathbb{R}^n.$$

## Example

Let  $A$  be a real  $m \times n$  matrix and

$$f(\mathbf{x}) = A\mathbf{x}, \quad \mathbf{x} \in \mathbb{R}^n.$$

The function  $f$  is linear since  $A(\alpha \mathbf{x} + \beta \mathbf{y}) = \alpha A\mathbf{x} + \beta A\mathbf{y}$ .

Moreover, the **range** of  $f$ ,

$$\mathcal{R}(f) = \{A\mathbf{x} : \mathbf{x} \in \mathbb{R}^n\} \subseteq \mathbb{R}^m,$$

is a **subspace** of  $\mathbb{R}^m$  since for all  $\alpha, \beta \in \mathbb{R}$  and  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$

$$\alpha \underbrace{(A\mathbf{x})}_{\in \mathcal{R}(f)} + \beta \underbrace{(A\mathbf{y})}_{\in \mathcal{R}(f)} = A(\alpha \mathbf{x} + \beta \mathbf{y}) \in \mathcal{R}(f).$$

## Remark

For the situation in this example we can also use the terminology *range of A* (or *image of A*), i.e.,

$$R(\mathbf{A}) = \{\mathbf{Ax} : \mathbf{x} \in \mathbb{R}^n\} \subseteq \mathbb{R}^m$$

Similarly,

$$R(\mathbf{A}^T) = \{\mathbf{A}^T \mathbf{y} : \mathbf{y} \in \mathbb{R}^m\} \subseteq \mathbb{R}^n$$

is called the *range of  $\mathbf{A}^T$* .



# Column space and row space

Since

$$A\mathbf{x} = \alpha_1\mathbf{a}_1 + \dots + \alpha_n\mathbf{a}_n,$$

we have  $R(A) = \text{span}\{\mathbf{a}_1, \dots, \mathbf{a}_n\}$ , i.e.,

$R(A)$  is the **column space** of  $A$ .

Similarly,

$R(A^T)$  is the **row space** of  $A$ .



## Example

Consider

$$A = \begin{pmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{pmatrix}$$

By definition

- the columns of  $A$  span  $R(A)$ , i.e., they form a spanning set of  $R(A)$ ,
- the rows of  $A$  span  $R(A^T)$ , i.e., they form a spanning set of  $R(A^T)$ ,

However, since

$$(A)_{*3} = 2(A)_{*2} - (A)_{*1} \quad \text{and} \quad (A)_{3*} = 2(A)_{2*} - (A)_{1*}$$

we also have

- $R(A) = \text{span}\{(A)_{*1}, (A)_{*2}\}$
- $R(A^T) = \text{span}\{(A)_{1*}, (A)_{2*}\}$

In general, **how do we find such minimal spanning sets** as in the previous example?

An important tool is

### Lemma

Let  $A, B$  be  $m \times n$  matrices. Then

$$(1) \quad R(A^T) = R(B^T) \iff A \overset{\text{row}}{\sim} B \quad (\iff E_A = E_B).$$

$$(2) \quad R(A) = R(B) \iff A \overset{\text{col}}{\sim} B \quad (\iff E_{A^T} = E_{B^T}).$$





## Proof

- ① “ $\Leftarrow$ ”: Assume  $A \stackrel{\text{row}}{\sim} B$ , i.e., there exists a nonsingular matrix  $P$  such that

$$PA = B \iff A^T P^T = B^T.$$

Now  $\mathbf{a} \in R(A^T) \iff \mathbf{a} = A^T \mathbf{y}$  for some  $\mathbf{y}$ .

We rewrite this as

$$\begin{aligned} \mathbf{a} &= \underbrace{A^T P^T}_{=B^T} P^{-T} \mathbf{y} \\ \iff \mathbf{a} &= B^T \mathbf{x} \quad \text{for } \mathbf{x} = P^{-T} \mathbf{y} \\ \iff \mathbf{a} &\in R(B^T). \end{aligned}$$



(cont.)

“ $\implies$ ”: Assume  $R(A^T) = R(B^T)$ , i.e.,

$$\text{span}\{(A)_{1*}, \dots, (A)_{m*}\} = \text{span}\{(B)_{1*}, \dots, (B)_{m*}\},$$

i.e., the rows of A are linear combinations of rows of B and vice versa.

Now **apply row operations to A** (all collected in P) to obtain

$$PA = B, \quad \text{i.e., } A \overset{\text{row}}{\sim} B.$$

2 Let  $A = A^T$  and  $B = B^T$  in (1). □



## Theorem

Let  $A$  be an  $m \times n$  matrix and  $U$  any row echelon form obtained from  $A$ . Then

- 1  $R(A^T) = \text{span of nonzero rows of } U.$
- 2  $R(A) = \text{span of basic columns of } A.$

## Remark

Later we will see that any *minimal* span of the columns of  $A$  forms a *basis* for  $R(A)$ .



## Proof

- 1 This follows from (1) in the Lemma since  $A \stackrel{\text{row}}{\sim} U$ .
- 2 Assume the columns of  $A$  are permuted (with a matrix  $Q_1$ ) such that

$$AQ_1 = (B \ N),$$

where  $B$  contains the **basic columns**, and  $N$  the **nonbasic columns**.

By definition, the **nonbasic columns are linear combinations of the basic columns**, i.e., **there exists a nonsingular  $Q_2$  such that**

$$(B \ N) Q_2 = (B \ O),$$

where  $O$  is a zero matrix.



(cont.)

Putting this together, we have

$$A \underbrace{Q_1 Q_2}_{=Q} = (B \ O),$$

so that  $A \overset{\text{col}}{\sim} (B \ O)$ .

(2) in the Lemma says that

$$R(A) = \text{span}\{B_{*1}, \dots, B_{*r}\},$$

where  $r = \text{rank}(A)$ . □



So far, we have **two of the four fundamental subspaces**:

$$R(A) \quad \text{and} \quad R(A^T).$$

**Third fundamental subspace:**  $N(A) = \{\mathbf{x} : A\mathbf{x} = \mathbf{0}\} \subseteq \mathbb{R}^n$ ,

$N(A)$  is the **nullspace of  $A$**

(also called the **kernel of  $A$** )

**Fourth fundamental subspace:**  $N(A^T) = \{\mathbf{y} : A^T\mathbf{y} = \mathbf{0}\} \subseteq \mathbb{R}^m$ ,

$N(A^T)$  is the **left nullspace of  $A$**

### Remark

$N(A)$  is a **linear space**, i.e., a **subspace** of  $\mathbb{R}^n$ .

To see this, assume  $\mathbf{x}, \mathbf{y} \in N(A)$ , i.e.,  $A\mathbf{x} = A\mathbf{y} = \mathbf{0}$ .

Then

$$A(\alpha\mathbf{x} + \beta\mathbf{y}) = \alpha A\mathbf{x} + \beta A\mathbf{y} = \mathbf{0},$$

so that  $\alpha\mathbf{x} + \beta\mathbf{y} \in N(A)$ .

## How to find a (minimal) spanning set for $N(A)$

Find a row echelon form  $U$  of  $A$  and solve  $U\mathbf{x} = \mathbf{0}$ .

### Example

We can compute  $A = \begin{pmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{pmatrix} \rightarrow U = \begin{pmatrix} 1 & 2 & 3 \\ 0 & -3 & -6 \\ 0 & 0 & 0 \end{pmatrix}$ .

So that  $U\mathbf{x} = \mathbf{0} \implies \begin{cases} x_2 = -2x_3 \\ x_1 = -2x_2 - 3x_3 = x_3 \end{cases}$ , or

$$\begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} x_3 \\ -2x_3 \\ x_3 \end{pmatrix} = x_3 \begin{pmatrix} 1 \\ -2 \\ 1 \end{pmatrix}.$$

Therefore

$$N(A) = \text{span} \left\{ \begin{pmatrix} 1 \\ -2 \\ 1 \end{pmatrix} \right\}.$$

## Remark

We will see later that — as in the example — *if  $\text{rank}(A) = r$ , then  $N(A)$  is spanned by  $n - r$  vectors.*

## Theorem

Let  $A$  be an  $m \times n$  matrix. Then

- 1  $N(A) = \{\mathbf{0}\} \iff \text{rank}(A) = n.$
- 2  $N(A^T) = \{\mathbf{0}\} \iff \text{rank}(A) = m.$

## Proof.

- 1 We know  $\text{rank}(A) = n \iff A\mathbf{x} = \mathbf{0}$ , but that implies  $\mathbf{x} = \mathbf{0}$ .
- 2 Repeat (1) with  $A = A^T$  and use  $\text{rank}(A^T) = \text{rank}(A)$ .





## How to find a spanning set of $N(A^T)$

### Theorem

Let  $A$  be an  $m \times n$  matrix with  $\text{rank}(A) = r$ , and let  $P$  be a nonsingular matrix so that  $PA = U$  (row echelon form). Then *the last  $m - r$  rows of  $P$  span  $N(A^T)$* .

### Remark

*We will later see that this spanning set is also a basis for  $N(A^T)$ .*



## Proof

Partition  $P$  as  $P = \begin{pmatrix} P_1 \\ P_2 \end{pmatrix}$ , where  $P_1$  is  $r \times m$  and  $P_2$  is  $m - r \times m$ .

The claim of the theorem implies that we should **show that**  
 $R(P_2^T) = N(A^T)$ .

We do this in two parts:

- 1 Show that  $R(P_2^T) \subseteq N(A^T)$ .
- 2 Show that  $N(A^T) \subseteq R(P_2^T)$ .



(cont.)

- 1 Partition  $U_{m \times n} = \begin{pmatrix} C \\ O \end{pmatrix}$  with  $C \in \mathbb{R}^{r \times n}$  and  $O \in \mathbb{R}^{m-r \times n}$  (a zero matrix).  
Then

$$PA = U \iff \begin{pmatrix} P_1 \\ P_2 \end{pmatrix} A = \begin{pmatrix} C \\ O \end{pmatrix} \implies P_2 A = O.$$

This also means that

$$A^T P_2^T = O^T,$$

i.e., every column of  $P_2^T$  is in  $N(A^T)$  so that  $R(P_2^T) \subseteq N(A^T)$ .



(cont.)

2 Now, show  $N(A^T) \subseteq R(P_2^T)$ .

We **assume**  $\mathbf{y} \in N(A^T)$  and **show that then**  $\mathbf{y} \in R(P_2^T)$ .

By definition,

$$\mathbf{y} \in N(A^T) \implies A^T \mathbf{y} = \mathbf{0} \iff \mathbf{y}^T A = \mathbf{0}^T.$$

Since  $PA = U \implies A = P^{-1}U$ , and so

$$\mathbf{0}^T = \mathbf{y}^T P^{-1} U = \mathbf{y}^T P^{-1} \begin{pmatrix} C \\ \mathbf{0} \end{pmatrix}$$

or

$$\mathbf{0}^T = \mathbf{y}^T Q_1 C, \quad \text{where } P^{-1} = \begin{pmatrix} \underbrace{Q_1}_{m \times r} & \underbrace{Q_2}_{m \times m-r} \end{pmatrix}.$$



(cont.)

However, since  $\text{rank}(C) = r$  and  $C$  is  $m \times n$  we get (using  $m = r$  in our earlier theorem)

$$N(C^T) = \{\mathbf{0}\}$$

and therefore  $\mathbf{y}^T Q_1 = \mathbf{0}^T$ .

Obviously, this implies that we also have

$$\mathbf{y}^T Q_1 P_1 = \mathbf{0}^T \tag{2}$$



(cont.)

Now  $P = \begin{pmatrix} P_1 \\ P_2 \end{pmatrix}$  and  $P^{-1} = (Q_1 \quad Q_2)$  so that

$$I = P^{-1}P = Q_1P_1 + Q_2P_2$$

or

$$Q_1P_1 = I - Q_2P_2. \quad (3)$$

Now we **insert (3) into (2)** and get

Therefore  $\mathbf{y} \in R(P_2^T)$ .



Finally,

### Theorem

Let  $A, B$  be  $m \times n$  matrices.

$$\textcircled{1} \quad N(A) = N(B) \iff A \overset{\text{row}}{\sim} B.$$

$$\textcircled{2} \quad N(A^T) = N(B^T) \iff A \overset{\text{col}}{\sim} B.$$

### Proof.

See [Mey00, Section 4.2].



# Linear Independence

## Definition

A set of vectors  $S = \{\mathbf{v}_1, \dots, \mathbf{v}_n\}$  is called **linearly independent** if

$$\alpha_1 \mathbf{v}_1 + \alpha_2 \mathbf{v}_2 + \dots + \alpha_n \mathbf{v}_n = \mathbf{0} \implies \alpha_1 = \alpha_2 = \dots = \alpha_n = 0.$$

Otherwise  $S$  is **linearly dependent**.

## Remark

*Linear independence is a property of a **set**, not of vectors.*





## Example

Is  $S = \left\{ \begin{pmatrix} 1 \\ 4 \\ 7 \end{pmatrix}, \begin{pmatrix} 2 \\ 5 \\ 8 \end{pmatrix}, \begin{pmatrix} 3 \\ 6 \\ 9 \end{pmatrix} \right\}$  linearly independent?

Consider

$$\alpha_1 \begin{pmatrix} 1 \\ 4 \\ 7 \end{pmatrix} + \alpha_2 \begin{pmatrix} 2 \\ 5 \\ 8 \end{pmatrix} + \alpha_3 \begin{pmatrix} 3 \\ 6 \\ 9 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix}$$

$$\iff \mathbf{Ax} = \mathbf{0}, \quad \text{where } \mathbf{A} = \begin{pmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{pmatrix}, \quad \mathbf{x} = \begin{pmatrix} \alpha_1 \\ \alpha_2 \\ \alpha_3 \end{pmatrix}$$



## Example ((cont.))

Since

$$A \stackrel{\text{row}}{\sim} E_A = \begin{pmatrix} 1 & 2 & 3 \\ 0 & 1 & 2 \\ 0 & 0 & 0 \end{pmatrix}$$

we know that  $N(A)$  is nontrivial, i.e., the system  $A\mathbf{x} = \mathbf{0}$  has a nonzero solution, and therefore  $\mathcal{S}$  is linearly dependent.



More generally,

### Theorem

Let  $A$  be an  $m \times n$  matrix.

- 1 The columns of  $A$  are linearly independent if and only if  $N(A) = \{\mathbf{0}\} \iff \text{rank}(A) = n$ .
- 2 The rows of  $A$  are linearly independent if and only if  $N(A^T) = \{\mathbf{0}\} \iff \text{rank}(A) = m$ .

### Proof.

See [Mey00, Section 4.3]. □



## Definition

A square matrix  $A$  is called **diagonally dominant** if

$$|a_{ii}| > \sum_{\substack{j=1 \\ j \neq i}}^n |a_{ij}|, \quad i = 1, \dots, n.$$

## Remark

- *Aside from being nonsingular (see next slide), diagonally dominant matrices are important since they ensure that **Gaussian elimination will succeed without pivoting**.*
- *Also, diagonal dominance ensures convergence of certain iterative solvers (more later).*



## Theorem

Let  $A$  be an  $n \times n$  matrix. If  $A$  is diagonally dominant then  $A$  is nonsingular.

## Proof

We will show that  $N(A) = \{\mathbf{0}\}$  since then we know that  $\text{rank}(A) = n$  and  $A$  is nonsingular.

We will do this with a **proof by contradiction**.

We **assume** that there exists an  $\mathbf{x} (\neq \mathbf{0}) \in N(A)$  and we will **conclude** that  $A$  cannot be diagonally dominant.



(cont.)

If  $\mathbf{x} \in N(\mathbf{A})$  then  $\mathbf{A}\mathbf{x} = \mathbf{0}$ .

Now we take  $k$  so that  $x_k$  is the maximum (in absolute value) component of  $\mathbf{x}$  and consider

$$A_{k*}\mathbf{x} = 0.$$

We can rewrite this as

$$\sum_{j=1}^n a_{kj}x_j = 0 \iff a_{kk}x_k = -\sum_{\substack{j=1 \\ j \neq k}}^n a_{kj}x_j.$$



(cont.)

Now we take absolute values:

$$\begin{aligned}
 |a_{kk}x_k| &= \left| \sum_{\substack{j=1 \\ j \neq k}}^n a_{kj}x_j \right| \leq \sum_{\substack{j=1 \\ j \neq k}}^n |a_{kj}| |x_j| \\
 &\leq \underbrace{|x_k|}_{\text{max. component}} \sum_{\substack{j=1 \\ j \neq k}}^n |a_{kj}|
 \end{aligned}$$

Finally, dividing both sides by  $|x_k|$  yields

$$|a_{kk}| \leq \sum_{\substack{j=1 \\ j \neq k}}^n |a_{kj}|,$$

which shows that **A cannot be diagonally dominant** (which is a contradiction since **A was assumed to be diagonally dominant**).  $\square$

## Example

Consider  $m$  real numbers  $x_1, \dots, x_m$  such that  $x_i \neq x_j$ ,  $i \neq j$ .  
 Show that the **columns** of the **Vandermonde matrix**

$$V = \begin{pmatrix} 1 & x_1 & x_1^2 & \cdots & x_1^{n-1} \\ 1 & x_2 & x_2^2 & \cdots & x_2^{n-1} \\ & & & \vdots & \\ 1 & x_m & x_m^2 & \cdots & x_m^{n-1} \end{pmatrix}$$

**form a linearly independent set provided  $n \leq m$ .**

From above, the columns of  $V$  are linearly independent if and only if  $N(V) = \{\mathbf{0}\}$

$$\iff V\mathbf{z} = \mathbf{0} \implies \mathbf{z} = \mathbf{0}, \quad \mathbf{z} = \begin{pmatrix} \alpha_0 \\ \vdots \\ \alpha_{n-1} \end{pmatrix}.$$



## Example

(cont.)

Now  $V\mathbf{z} = \mathbf{0}$  if and only if

$$\alpha_0 + \alpha_1 x_i + \alpha_2 x_i^2 + \dots + \alpha_{n-1} x_i^{n-1} = 0, \quad i = 1, \dots, m.$$

In other words,  $x_1, x_2, \dots, x_m$  are all (distinct) roots of

$$p(x) = \alpha_0 + \alpha_1 x + \alpha_2 x^2 + \dots + \alpha_{n-1} x^{n-1}.$$

This is a polynomial of degree at most  $n - 1$ .It can have  $m$  distinct roots only if  $m \leq n - 1$ .

Otherwise,  $p$  is the zero polynomial, i.e.,  $\alpha_0 = \alpha_1 = \dots = \alpha_{n-1} = 0$ , so that the columns of  $V$  are linearly dependent.



The example implies that in the special case  $m = n$  there is a unique polynomial of degree (at most)  $m - 1$  that interpolates the data  $\{(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)\} \subset \mathbb{R}^2$ .

We see this by writing the polynomial in the form

$$\ell(t) = \alpha_0 + \alpha_1 t + \alpha_2 t^2 + \dots + \alpha_{m-1} t^{m-1}.$$

Then, interpolation of the data implies

$$\ell(x_i) = y_i, \quad i = 1, \dots, m$$

or

$$\begin{pmatrix} 1 & x_1 & x_1^2 & \cdots & x_1^{m-1} \\ 1 & x_2 & x_2^2 & \cdots & x_2^{m-1} \\ & & & \vdots & \\ 1 & x_m & x_m^2 & \cdots & x_m^{m-1} \end{pmatrix} \begin{pmatrix} \alpha_0 \\ \alpha_1 \\ \vdots \\ \alpha_{m-1} \end{pmatrix} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{pmatrix}.$$

Since the columns of  $V$  are linearly independent it is nonsingular, and the coefficients  $\alpha_0, \dots, \alpha_{m-1}$  are uniquely determined.



In fact,

$$\ell(t) = \sum_{i=1}^m y_i L_i(t) \quad (\text{Lagrange interpolation polynomial})$$

$$\text{with } L_i(t) = \prod_{\substack{k=1 \\ k \neq i}}^m (t - x_k) / \prod_{\substack{k=1 \\ k \neq i}}^m (x_i - x_k) \quad (\text{Lagrange functions}).$$

To verify (4) we note that the **degree of  $\ell$  is  $m - 1$**  (since each  $L_i$  is of degree  $m - 1$ ) and

$$L_i(x_j) = \delta_{ij}, \quad i, j = 1, \dots, m,$$

so that

$$\ell(x_j) = \sum_{i=1}^m y_i \underbrace{L_i(x_j)}_{=\delta_{ij}} = y_j, \quad j = 1, \dots, m.$$



## Theorem

Let  $S = \{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_n\} \subseteq \mathcal{V}$  be nonempty. Then

- 1 If  $S$  contains a linearly dependent subset, then  $S$  is linearly dependent.
- 2 If  $S$  is linearly independent, then every subset of  $S$  is also linearly independent.
- 3 If  $S$  is linearly independent and if  $\mathbf{v} \in \mathcal{V}$ , then  $S_{\text{ext}} = S \cup \{\mathbf{v}\}$  is linearly independent if and only if  $\mathbf{v} \notin \text{span}(S)$ .
- 4 If  $S \subseteq \mathbb{R}^m$  and  $n > m$ , then  $S$  must be linearly dependent.



## Proof

- 1 If  $S$  contains a linearly dependent subset,  $\{\mathbf{u}_1, \dots, \mathbf{u}_k\}$  say, then there exist nontrivial coefficients  $\alpha_1, \dots, \alpha_k$  such that

$$\alpha_1 \mathbf{u}_1 + \dots + \alpha_k \mathbf{u}_k = \mathbf{0}.$$

Clearly, then

$$\alpha_1 \mathbf{u}_1 + \dots + \alpha_k \mathbf{u}_k + \mathbf{0} \mathbf{u}_{k+1} + \dots + \mathbf{0} \mathbf{u}_n = \mathbf{0}$$

and  $S$  is also linearly dependent.

- 2 Follows from (1) by contraposition.



(cont.)

③ “ $\implies$ ”: Assume  $\mathcal{S}_{\text{ext}}$  is linearly independent. Then  $\mathbf{v}$  can't be a linear combination of  $\mathbf{u}_1, \dots, \mathbf{u}_n$ .

“ $\impliedby$ ”: Assume  $\mathbf{v} \notin \text{span}(\mathcal{S})$  and consider

$$\alpha_1 \mathbf{u}_1 + \alpha_2 \mathbf{u}_2 + \dots + \alpha_n \mathbf{u}_n + \alpha_{n+1} \mathbf{v} = \mathbf{0}.$$

First,  $\alpha_{n+1} = 0$  since otherwise  $\mathbf{v} \in \text{span}(\mathcal{S})$ .

That leaves

$$\alpha_1 \mathbf{u}_1 + \alpha_2 \mathbf{u}_2 + \dots + \alpha_n \mathbf{u}_n = \mathbf{0}.$$

However, the linear independence of  $\mathcal{S}$  implies  $\alpha_j = 0$ ,  $i = 1, \dots, n$ , and therefore  $\mathcal{S}_{\text{ext}}$  is linearly independent.



(cont.)

- 4 We know that the **columns** of an  $m \times n$  matrix  $A$  are **linearly independent** if and only if  $\text{rank}(A) = n$ .

Here  $A = (\mathbf{u}_1 \quad \mathbf{u}_2 \quad \cdots \quad \mathbf{u}_n)$  with  $\mathbf{u}_i \in \mathbb{R}^m$ .

If  $n > m$ , then  $\text{rank}(A) \leq m$  and  $S$  must be linearly dependent.  $\square$



# Bases and Dimension

Earlier we introduced the concept of a **spanning set** of a vector space  $\mathcal{V}$ , i.e.,

$$\mathcal{V} = \text{span}\{\mathbf{v}_1, \dots, \mathbf{v}_n\}$$

Now

## Definition

Consider a vector space  $\mathcal{V}$  with **spanning set**  $\mathcal{S}$ . If  $\mathcal{S}$  is also **linearly independent** then we call it a **basis of  $\mathcal{V}$** .

## Example

- 1  $\{\mathbf{e}_1, \dots, \mathbf{e}_n\}$  is the **standard basis** for  $\mathbb{R}^n$ .
- 2 The columns/rows of an  $n \times n$  matrix  $A$  with  $\text{rank}(A) = n$  form a basis for  $\mathbb{R}^n$ .





## Remark

*Linear algebra deals with **finite-dimensional** linear spaces.*

*Functional analysis can be considered as **infinite-dimensional linear algebra**, where the linear spaces are usually **function spaces** such as*

- *infinitely differentiable functions with **Taylor (polynomial) basis***

$$\{1, x, x^2, x^3, \dots\}$$

- *square integrable functions with **Fourier basis***

$$\{1, \sin(x), \cos(x), \sin(2x), \cos(2x), \dots\}$$



Earlier we mentioned the idea of **minimal spanning sets**.

### Theorem

Let  $\mathcal{V}$  be a subspace of  $\mathbb{R}^m$  and let

$$\mathcal{B} = \{\mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_n\} \subseteq \mathcal{V}.$$

The following are equivalent:

- 1  $\mathcal{B}$  is a basis for  $\mathcal{V}$ .
- 2  $\mathcal{B}$  is a minimal spanning set for  $\mathcal{V}$ .
- 3  $\mathcal{B}$  is a maximal linearly independent subset of  $\mathcal{V}$ .

### Remark

We say “a basis” here since  $\mathcal{V}$  can have many different bases.



## Proof

Since it is **difficult to directly relate (2) and (3)**, our strategy will be to prove

- Show  $(1) \implies (2)$  and  $(2) \implies (1)$ , so that  $(1) \iff (2)$ .
- Show  $(1) \implies (3)$  and  $(3) \implies (1)$ , so that  $(1) \iff (3)$ .

Then — by transitivity — we will also have  $(2) \iff (3)$ .



Proof (cont.)

(1)  $\implies$  (2): Assume  $\mathcal{B}$  is a basis (i.e., a linearly independent spanning set) of  $\mathcal{V}$  and show that it is minimal.

Assume  $\mathcal{B}$  is not minimal, i.e., we can find a smaller spanning set  $\{\mathbf{x}_1, \dots, \mathbf{x}_k\}$  for  $\mathcal{V}$  with  $k \leq n$  elements.

But then

$$\mathbf{b}_j = \sum_{i=1}^k \alpha_{ij} \mathbf{x}_i, \quad j = 1, \dots, n,$$

or

$$\mathbf{B} = \mathbf{X}\mathbf{A},$$

where

$$\mathbf{B} = (\mathbf{b}_1 \quad \mathbf{b}_2 \quad \cdots \quad \mathbf{b}_n) \in \mathbb{R}^{m \times n},$$

$$\mathbf{X} = (\mathbf{x}_1 \quad \mathbf{x}_2 \quad \cdots \quad \mathbf{x}_k) \in \mathbb{R}^{m \times k},$$

$$[\mathbf{A}]_{ij} = \alpha_{ij}, \quad \mathbf{A} \in \mathbb{R}^{k \times n}.$$

**Proof** (cont.)

Now,  $\text{rank}(A) \leq k < n$ , which implies  $N(A)$  is nontrivial, i.e., there exists a  $\mathbf{z} \neq \mathbf{0}$  such that

$$A\mathbf{z} = \mathbf{0}.$$

But then

$$B\mathbf{z} = XA\mathbf{z} = \mathbf{0},$$

and therefore  $N(B)$  is nontrivial.

However, since  $\mathcal{B}$  is a basis, the columns of  $B$  are linearly independent (i.e.,  $N(B) = \{\mathbf{0}\}$ ) — and that is a **contradiction**.

Therefore,  $\mathcal{B}$  has to be minimal.



Proof (cont.)

(2)  $\implies$  (1): Assume  $\mathcal{B}$  is a minimal spanning set and show that it must also be linearly independent.

This is clear since

- if  $\mathcal{B}$  were linearly dependent,
- then we would be able to remove at least one vector from  $\mathcal{B}$  and still have a spanning set
- but then it would not have been minimal.



Proof (cont.)

(3)  $\implies$  (1): Assume  $\mathcal{B}$  is a maximal linearly independent subset of  $\mathcal{V}$  and show that  $\mathcal{B}$  is a basis of  $\mathcal{V}$ .

Assume that  $\mathcal{B}$  is not a basis, i.e., there exists a  $\mathbf{v} \in \mathcal{V}$  such that  $\mathbf{v} \notin \text{span}\{\mathbf{b}_1, \dots, \mathbf{b}_n\}$ .

Then — by an earlier theorem — the extension set  $\mathcal{B} \cup \{\mathbf{v}\}$  is linearly independent.

But this contradicts the maximality of  $\mathcal{B}$ , so that  $\mathcal{B}$  has to be a basis.



Proof (cont.)

(1)  $\implies$  (3): Assume  $\mathcal{B}$  is a basis, but not a maximal linearly independent subset of  $\mathcal{V}$ , and show that this leads to a contradiction.

Let

$$\mathcal{Y} = \{\mathbf{y}_1, \dots, \mathbf{y}_k\} \subseteq \mathcal{V}, \quad \text{with } k > n$$

be a maximal linearly independent subset of  $\mathcal{V}$  (note that such a set always exists).

But then  $\mathcal{Y}$  must be a basis for  $\mathcal{V}$  by our “(1)  $\implies$  (3)” argument.

On the other hand,  $\mathcal{Y}$  has more vectors than  $\mathcal{B}$  and a basis has to be a minimal spanning set.

Therefore  $\mathcal{B}$  has to already be a maximal linearly independent subset of  $\mathcal{V}$ .  $\square$





## Remark

Above we remarked that  $\mathcal{B}$  is not unique, i.e., a vector space  $\mathcal{V}$  can have many different bases.

However, the number of elements in all of these bases is unique.

## Definition

The dimension of the vector space  $\mathcal{V}$  is given by

$$\dim \mathcal{V} = \text{the number of elements in any basis of } \mathcal{V}.$$

Special case: by convention

$$\dim\{\mathbf{0}\} = 0.$$



### Example

Consider

$$\mathcal{P} = \text{span} \left\{ \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} \right\} \subset \mathbb{R}^3.$$

Geometrically,  $\mathcal{P}$  corresponds to the plane  $z = 0$ , i.e., the  $xy$ -plane.

Note that  $\dim \mathcal{P} = 2$ .

Moreover, any subspace of  $\mathbb{R}^3$  has dimension at most 3.



In general,

### Theorem

Let  $\mathcal{M}$  and  $\mathcal{N}$  be vector spaces such that  $\mathcal{M} \subseteq \mathcal{N}$ . Then

- 1  $\dim \mathcal{M} \leq \dim \mathcal{N}$ ,
- 2  $\dim \mathcal{M} = \dim \mathcal{N} \implies \mathcal{M} = \mathcal{N}$ .

### Proof.

See [Mey00]. □



## Back to the 4 fundamental subspaces

Consider an  $m \times n$  matrix  $A$  with  $\text{rank}(A) = r$ .

$R(A)$  We know that

$$R(A) = \text{span}\{\text{columns of } A\}.$$

If  $\text{rank}(A) = r$ , then only  $r$  columns of  $A$  are linearly independent, i.e.,

$$\dim R(A) = r.$$

A basis of  $R(A)$  is given by the basic columns of  $A$  (determined via a row echelon form  $U$ ).



$R(A^T)$  We know that

$$R(A^T) = \text{span}\{\text{rows of } A\}.$$

Again,  $\text{rank}(A) = r$  implies that only  $r$  rows of  $A$  are linearly independent, i.e.,

$$\dim R(A^T) = r.$$

A basis of  $R(A^T)$  is given by the nonzero rows of  $U$  (from the LU factorization of  $A$ ).



$N(A^T)$  One of our earlier theorems states that the **last  $m - r$  rows of  $P$  span  $N(A^T)$**  (where  $P$  is nonsingular such that  $PA = U$  is in row echelon form).

Since  $P$  is nonsingular **these rows are linearly independent** and so

$$\dim N(A^T) = m - r.$$

A **basis of  $N(A^T)$**  is given by the **last  $m - r$  rows of  $P$** .



$N(A)$  Replace  $A$  by  $A^T$  above so that

$$\dim N\left((A^T)^T\right) = n - \text{rank}(A^T) = n - r$$

so that

$$\dim N(A) = n - r.$$

A basis of  $N(A)$  is given by the  $n - r$  linearly independent solutions of  $A\mathbf{x} = \mathbf{0}$ .



**Theorem**

For any  $m \times n$  matrix  $A$  we have

$$\dim R(A) + \dim N(A) = n.$$

This follows directly from the above discussion of  $R(A)$  and  $N(A)$ .

The theorem shows that there is always a **balance between the rank of  $A$  and the dimension of its nullspace.**





### Example

Find the **dimension** and a **basis** for

$$S = \text{span} \left\{ \begin{pmatrix} 1 \\ 2 \\ 3 \\ 1 \end{pmatrix}, \begin{pmatrix} 2 \\ 4 \\ 6 \\ 2 \end{pmatrix}, \begin{pmatrix} 2 \\ 4 \\ 6 \\ 4 \end{pmatrix}, \begin{pmatrix} 3 \\ 6 \\ 9 \\ 5 \end{pmatrix}, \begin{pmatrix} 1 \\ 2 \\ 6 \\ 3 \end{pmatrix} \right\}.$$

Before we even do any calculations we know that

$$S \subseteq \mathbb{R}^4, \quad \text{so that } \dim S \leq 4.$$

We will now answer this question in **two different ways** using

$$A = \begin{pmatrix} 1 & 2 & 2 & 3 & 1 \\ 2 & 4 & 4 & 6 & 2 \\ 3 & 6 & 6 & 9 & 6 \\ 1 & 2 & 4 & 5 & 3 \end{pmatrix}.$$

## Example (cont.)

Via  $R(A)$ , i.e., by finding the basic columns of  $A$ :

$$A = \begin{pmatrix} 1 & 2 & 2 & 3 & 1 \\ 2 & 4 & 4 & 6 & 2 \\ 3 & 6 & 6 & 9 & 6 \\ 1 & 2 & 4 & 5 & 3 \end{pmatrix} \xrightarrow{\text{G.-J.}} E_A = \begin{pmatrix} 1 & 2 & 0 & 1 & 0 \\ 0 & 0 & 1 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 \end{pmatrix}$$

Therefore,  $\dim \mathcal{S} = 3$  and

$$\mathcal{S} = \text{span} \left\{ \begin{pmatrix} 1 \\ 2 \\ 3 \\ 1 \end{pmatrix}, \begin{pmatrix} 2 \\ 4 \\ 6 \\ 4 \end{pmatrix}, \begin{pmatrix} 1 \\ 2 \\ 6 \\ 3 \end{pmatrix} \right\}$$

since the **basic columns of  $E_A$**  are the first, third and fifth columns.



## Example (cont.)

Via  $R(A^T)$ , i.e.,  $R(A) = \text{span}\{\text{rows of } A^T\}$ , i.e., we need the nonzero rows of  $U$  (from the LU factorization of  $A^T$ ):

$$A^T = \begin{pmatrix} 1 & 2 & 3 & 1 \\ 2 & 4 & 6 & 2 \\ 2 & 4 & 6 & 4 \\ 3 & 6 & 9 & 4 \\ 1 & 2 & 6 & 3 \end{pmatrix} \xrightarrow{\text{zero out } [A^T]_{*,1}} \begin{pmatrix} 1 & 2 & 3 & 1 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 2 \\ 0 & 0 & 0 & 2 \\ 0 & 0 & 3 & 2 \end{pmatrix} \xrightarrow{\text{permute}} \underbrace{\begin{pmatrix} 1 & 2 & 3 & 1 \\ 0 & 0 & 3 & 2 \\ 0 & 0 & 0 & 2 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix}}_{=U}$$

Therefore,  $\dim \mathcal{S} = 3$  and

$$\mathcal{S} = \text{span} \left\{ \begin{pmatrix} 1 \\ 2 \\ 3 \\ 1 \end{pmatrix}, \begin{pmatrix} 0 \\ 0 \\ 3 \\ 2 \end{pmatrix}, \begin{pmatrix} 0 \\ 0 \\ 0 \\ 2 \end{pmatrix} \right\}$$

since the nonzero rows of  $U$  are the first, second and third rows.

## Example

## Extend

$$\mathcal{S} = \text{span} \left\{ \begin{pmatrix} 1 \\ 2 \\ 3 \\ 1 \end{pmatrix}, \begin{pmatrix} 1 \\ 2 \\ 6 \\ 3 \end{pmatrix} \right\}$$

to a basis for  $\mathbb{R}^4$ .

The procedure will be to **augment the columns of  $\mathcal{S}$  by an identity matrix**, i.e., to form

$$A = \begin{pmatrix} 1 & 1 & 1 & 0 & 0 & 0 \\ 2 & 2 & 0 & 1 & 0 & 0 \\ 3 & 6 & 0 & 0 & 1 & 0 \\ 1 & 3 & 0 & 0 & 0 & 1 \end{pmatrix}$$

and then to **get a basis via the basic columns of  $U$** .

## Example (cont.)

$$\begin{aligned}
 A = \begin{pmatrix} 1 & 1 & 1 & 0 & 0 & 0 \\ 2 & 2 & 0 & 1 & 0 & 0 \\ 3 & 6 & 0 & 0 & 1 & 0 \\ 1 & 3 & 0 & 0 & 0 & 1 \end{pmatrix} &\longrightarrow \begin{pmatrix} 1 & 1 & 1 & 0 & 0 & 0 \\ 0 & 0 & -2 & 1 & 0 & 0 \\ 0 & 3 & -3 & 0 & 1 & 0 \\ 0 & 2 & -1 & 0 & 0 & 1 \end{pmatrix} \\
 \longrightarrow \begin{pmatrix} 1 & 1 & 1 & 0 & 0 & 0 \\ 0 & 2 & -1 & 0 & 0 & 1 \\ 0 & 0 & -\frac{3}{2} & 0 & 1 & -\frac{3}{2} \\ 0 & 0 & -2 & 1 & 0 & 0 \end{pmatrix} &\longrightarrow \begin{pmatrix} 1 & 1 & 1 & 0 & 0 & 0 \\ 0 & 2 & -1 & 0 & 0 & 1 \\ 0 & 0 & -\frac{3}{2} & 0 & 1 & -\frac{3}{2} \\ 0 & 0 & 0 & 1 & -\frac{4}{3} & 2 \end{pmatrix}
 \end{aligned}$$

so that the basic columns are  $[A]_{*1}$ ,  $[A]_{*2}$ ,  $[A]_{*3}$ ,  $[A]_{*4}$  and

$$\mathbb{R}^4 = \text{span} \left\{ \begin{pmatrix} 1 \\ 2 \\ 3 \\ 1 \end{pmatrix}, \begin{pmatrix} 1 \\ 2 \\ 6 \\ 3 \end{pmatrix}, \begin{pmatrix} 1 \\ 0 \\ 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 1 \\ 0 \\ 0 \end{pmatrix} \right\}.$$

Earlier we defined the **sum of subspaces**  $\mathcal{X}$  and  $\mathcal{Y}$  as

$$\mathcal{X} + \mathcal{Y} = \{\mathbf{x} + \mathbf{y} : \mathbf{x} \in \mathcal{X}, \mathbf{y} \in \mathcal{Y}\}$$

### Theorem

*If  $\mathcal{X}, \mathcal{Y}$  are subspaces of  $\mathcal{V}$ , then*

$$\dim(\mathcal{X} + \mathcal{Y}) = \dim \mathcal{X} + \dim \mathcal{Y} - \dim(\mathcal{X} \cap \mathcal{Y}).$$

### Proof.

See [Mey00], but the basic idea is pretty clear.  
We want to avoid double counting. □



## Corollary

Let  $A$  and  $B$  be  $m \times n$  matrices. Then

$$\text{rank}(A + B) \leq \text{rank}(A) + \text{rank}(B).$$

## Proof

First we note that

$$R(A + B) \subseteq R(A) + R(B) \quad (4)$$

since for any  $\mathbf{b} \in R(A + B)$  we have

$$\mathbf{b} = (A + B)\mathbf{x} = A\mathbf{x} + B\mathbf{x} \in R(A) + R(B).$$



(cont.)

Now,

$$\text{rank}(A + B) = \dim R(A + B)$$

$$\stackrel{(4)}{\leq} \dim(R(A) + R(B))$$

$$\stackrel{\text{Thm}}{=} \dim R(A) + \dim R(B) - \dim (R(A) \cap R(B))$$

$$\leq \dim R(A) + \dim R(B)$$

$$= \text{rank}(A) + \text{rank}(B)$$





## More About Rank

We know that  $A \sim B$  if and only if  $\text{rank}(A) = \text{rank}(B)$ .

Thus (for invertible  $P, Q$ ),  $PAQ = B$  implies  $\text{rank}(A) = \text{rank}(PAQ)$ .

As we now show, it is a general fact that **multiplication by a nonsingular matrix does not change the rank of a given matrix.**

Moreover, **multiplication by an arbitrary matrix can only lower the rank.**

### Theorem

*Let  $A$  be an  $m \times n$  matrix, and let  $B$  be  $n \times p$ . Then*

$$\text{rank}(AB) = \text{rank}(B) - \dim(N(A) \cap R(B)).$$

### Remark

*Note that if  $A$  is nonsingular, then  $N(A) = \{\mathbf{0}\}$  so that  $\dim(N(A) \cap R(B)) = 0$  and  $\text{rank}(AB) = \text{rank}(B)$ .*

## Proof

Let  $\mathcal{S} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_s\}$  be a basis for  $N(A) \cap R(B)$ .

Since  $N(A) \cap R(B) \subseteq R(B)$  we know that

$$\dim(R(B)) = s + t, \quad \text{for some } t \geq 0.$$

We can construct an extension set such that

$$\mathcal{B} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_s, \mathbf{z}_1, \dots, \mathbf{z}_2, \dots, \mathbf{z}_t\}$$

is a basis for  $R(B)$ .



(cont.)

If we can show that  $\dim(R(AB)) = t$  then

$$\text{rank}(B) = \dim(R(B)) = s + t = \dim(N(A) \cap R(B)) + \dim(R(AB)),$$

and we are done.

Therefore, we now show that  $\dim(R(AB)) = t$ .

In particular, we show that

$$\mathcal{T} = \{A\mathbf{z}_1, A\mathbf{z}_2, \dots, A\mathbf{z}_t\}$$

is a basis for  $R(AB)$ .

We do this by showing that

- 1  $\mathcal{T}$  is a spanning set for  $R(AB)$ ,
- 2  $\mathcal{T}$  is linearly independent.

(cont.)

**Spanning set:** Consider an arbitrary  $\mathbf{b} \in R(AB)$ . It can be written as

$$\mathbf{b} = A\mathbf{B}\mathbf{y} \quad \text{for some } \mathbf{y}.$$

But then  $\mathbf{B}\mathbf{y} \in R(B)$ , so that

$$\mathbf{B}\mathbf{y} = \sum_{i=1}^s \xi_i \mathbf{x}_i + \sum_{j=1}^t \eta_j \mathbf{z}_j$$

and

$$\mathbf{b} = A\mathbf{B}\mathbf{y} = \sum_{i=1}^s \xi_i A\mathbf{x}_i + \sum_{j=1}^t \eta_j A\mathbf{z}_j = \sum_{j=1}^t \eta_j A\mathbf{z}_j$$

since  $\mathbf{x}_i \in N(A)$ .

(cont.)

**Linear independence:** Let's use the definition of linear independence and look at

$$\sum_{i=1}^t \alpha_i \mathbf{A} \mathbf{z}_i = \mathbf{0} \iff \mathbf{A} \sum_{i=1}^t \alpha_i \mathbf{z}_i = \mathbf{0}.$$

The identity on the right implies that  $\sum_{i=1}^t \alpha_i \mathbf{z}_i \in N(\mathbf{A})$ .

But we also have  $\mathbf{z}_i \in \mathcal{B}$ , i.e.,  $\sum_{i=1}^t \alpha_i \mathbf{z}_i \in R(\mathcal{B})$ .

And so together

$$\sum_{i=1}^t \alpha_i \mathbf{z}_i \in N(\mathbf{A}) \cap R(\mathcal{B}).$$

(cont.)

Now, since  $\mathcal{S} = \{\mathbf{x}_1, \dots, \mathbf{x}_s\}$  is a basis for  $N(\mathbf{A}) \cap R(\mathbf{B})$  we have

$$\sum_{i=1}^t \alpha_i \mathbf{z}_i = \sum_{j=1}^s \beta_j \mathbf{x}_j \iff \sum_{i=1}^t \alpha_i \mathbf{z}_i - \sum_{j=1}^s \beta_j \mathbf{x}_j = \mathbf{0}.$$

But  $\mathcal{B} = \{\mathbf{x}_1, \dots, \mathbf{x}_s, \mathbf{z}_1, \dots, \mathbf{z}_t\}$  is linearly independent, so that  $\alpha_1 = \dots = \alpha_t = \beta_1 = \dots = \beta_s = 0$  and therefore  $\mathcal{T}$  is also linearly independent.  $\square$



It turns out that  $\dim(N(A) \cap R(B))$  is relatively difficult to determine.

Therefore, the following upper and lower bounds for  $\text{rank}(AB)$  are useful.

### Theorem

Let  $A$  be an  $m \times n$  matrix, and let  $B$  be  $n \times p$ . Then

- 1  $\text{rank}(AB) \leq \min\{\text{rank}(A), \text{rank}(B)\},$
- 2  $\text{rank}(AB) \geq \text{rank}(A) + \text{rank}(B) - n.$



## Proof of (1)

We show that  $\text{rank}(AB) \leq \text{rank}(A)$  and  $\text{rank}(AB) \leq \text{rank}(B)$ .

The previous theorem states

$$\text{rank}(AB) = \text{rank}(B) - \underbrace{\dim(N(A) \cap R(B))}_{\geq 0} \leq \text{rank}(B).$$

Similarly,

$$\text{rank}(AB) = \text{rank}(AB)^T = \text{rank}(B^T A^T) \stackrel{\text{as above}}{\leq} \text{rank}(A^T) = \text{rank}(A).$$

To make things as tight as possible we take the smaller of the two upper bounds.





## Proof of (2)

We begin by noting that  $N(A) \cap R(B) \subseteq N(A)$ .

Therefore,

$$\dim(N(A) \cap R(B)) \leq \dim(N(A)) = n - \text{rank}(A).$$

But then (using the previous theorem)

$$\begin{aligned} \text{rank}(AB) &= \text{rank}(B) - \dim(N(A) \cap R(B)) \\ &\geq \text{rank}(B) - n + \text{rank}(A). \end{aligned}$$



To prepare for our study of **least squares solutions**, where the matrices  $A^T A$  and  $AA^T$  are important, we prove

### Lemma

Let  $A$  be a real  $m \times n$  matrix. Then

- 1  $\text{rank}(A^T A) = \text{rank}(AA^T) = \text{rank}(A)$ .
- 2  $R(A^T A) = R(A^T)$ ,  $R(AA^T) = R(A)$ .
- 3  $N(A^T A) = N(A)$ ,  $N(AA^T) = N(A^T)$ .



## Proof

From our earlier theorem we know

$$\text{rank}(A^T A) = \text{rank}(A) - \dim(N(A^T) \cap R(A)).$$

For (1) to be true we need to show  $\dim(N(A^T) \cap R(A)) = 0$ , i.e.,  $N(A^T) \cap R(A) = \{\mathbf{0}\}$ .

This is true since

$$\mathbf{x} \in N(A^T) \cap R(A) \implies A^T \mathbf{x} = \mathbf{0} \text{ and } \mathbf{x} = A\mathbf{y} \text{ for some } \mathbf{y}.$$

Therefore (using  $\mathbf{x}^T = \mathbf{y}^T A^T$ )

$$\mathbf{x}^T \mathbf{x} = \mathbf{y}^T A^T \mathbf{x} = 0.$$

But

$$\mathbf{x}^T \mathbf{x} = 0 \iff \sum_{i=1}^m x_i^2 = 0 \implies \mathbf{x} = \mathbf{0}.$$

(cont.)

$\text{rank}(AA^T) = \text{rank}(A^T)$  obtained by **switching  $A$  and  $A^T$** , and then use  $\text{rank}(A^T) = \text{rank}(A)$ .

The first part of (2) follows from  $R(A^T A) \subseteq R(A^T)$  (see HW) and

$$\dim(R(A^T A)) = \text{rank}(A^T A) \stackrel{(1)}{=} \text{rank}(A^T) = \dim(R(A^T))$$

since for  $\mathcal{M} \subseteq \mathcal{N}$  with  $\dim \mathcal{M} = \dim \mathcal{N}$  one has  $\mathcal{M} = \mathcal{N}$  (from an earlier theorem).

The other part of (2) follows by switching  $A$  and  $A^T$ .



(cont.)

The first part of (3) follows from  $N(A) \subseteq N(A^T A)$  (see HW) and

$$\dim(N(A)) = n - \text{rank}(A) = n - \text{rank}(A^T A) = \dim(N(A^T A))$$

using the same reasoning as above.

The other part of (3) follows by switching  $A$  and  $A^T$ .  $\square$



# Connection to least squares and normal equations

Consider a — possibly inconsistent — linear system

$$\mathbf{Ax} = \mathbf{b}$$

with  $m \times n$  matrix  $A$  (and  $\mathbf{b} \notin R(A)$  if inconsistent).

To find a “solution” we multiply both sides by  $A^T$  to get the **normal equations**:

$$A^T A \mathbf{x} = A^T \mathbf{b},$$

where  $A^T A$  is an  $n \times n$  matrix.



## Theorem

Let  $A$  be an  $m \times n$  matrix,  $\mathbf{b}$  an  $m$ -vector, and consider the normal equations

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

associated with  $A \mathbf{x} = \mathbf{b}$ .

- 1 The *normal equations are always consistent*, i.e., for every  $A$  and  $\mathbf{b}$  there exists at least one  $\mathbf{x}$  such that  $A^T A \mathbf{x} = A^T \mathbf{b}$ .
- 2 If  $A \mathbf{x} = \mathbf{b}$  is consistent, then  $A^T A \mathbf{x} = A^T \mathbf{b}$  has the same solution set (the *least squares solution* of  $A \mathbf{x} = \mathbf{b}$ ).
- 3  $A^T A \mathbf{x} = A^T \mathbf{b}$  has a *unique solution if and only if*  $\text{rank}(A) = n$ .  
Then

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

regardless of whether  $A \mathbf{x} = \mathbf{b}$  is consistent or not.

- 4 If  $A \mathbf{x} = \mathbf{b}$  is consistent and has a unique solution, then the same holds for  $A^T A \mathbf{x} = A^T \mathbf{b}$  and  $\mathbf{x} = (A^T A)^{-1} A^T \mathbf{b}$ .

## Proof

(1) follows from our previous lemma, i.e.,

$$A^T \mathbf{b} \in R(A^T) = R(A^T A).$$

To show (2) we assume the  $\mathbf{p}$  is some particular solution of  $A\mathbf{x} = \mathbf{b}$ , i.e.,  $A\mathbf{p} = \mathbf{b}$ .

If we multiply by  $A^T$ , then

$$A^T A\mathbf{p} = A^T \mathbf{p},$$

so that  $\mathbf{p}$  is also a solution of the normal equations.





(cont.)

Now, the **general solution of  $A\mathbf{x} = \mathbf{b}$**  is from the set (see Problem 2 on HW#4)

$$\mathcal{S} = \mathbf{p} + N(A).$$

Moreover, the general solution of  $A^T A\mathbf{x} = A^T \mathbf{b}$  is of the form

$$\mathbf{p} + N(A^T A) \stackrel{\text{lemma}}{=} \mathbf{p} + N(A) = \mathcal{S}.$$



(cont.)

For (3) we want to show that  $A^T A \mathbf{x} = A^T \mathbf{b}$  has a unique solution if and only if  $\text{rank}(A) = n$ .

What we know immediately is that  $A^T A \mathbf{x} = A^T \mathbf{b}$  has a unique solution if and only if  $\text{rank}(A^T A) = n$ .

Since we showed earlier that  $\text{rank}(A^T A) = \text{rank}(A)$  this part is done.

Now, if  $\text{rank}(A^T A) = n$  we know that  $A^T A$  is invertible (even though  $A^T$  and  $A$  may not be) and therefore

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

To show (4) we note that  $A \mathbf{x} = \mathbf{b}$  has a unique solution if and only if  $\text{rank}(A) = n$ . But  $\text{rank}(A^T A) = \text{rank}(A)$  and the rest follows from (3).  $\square$



## Remark

The normal equations are *not recommended for serious computations* since they are often rather *ill-conditioned* since one can show that

$$\text{cond}(A^T A) = \text{cond}(A)^2.$$

*There's an example in [Mey00] that illustrates this fact.*



## Historical definition of rank

Let  $A$  be an  $m \times n$  matrix. Then  $A$  has *rank*  $r$  if there exists at least one nonsingular  $r \times r$  submatrix of  $A$  (and none larger).

### Example

The matrix

$$A = \begin{pmatrix} 1 & 2 & 2 & 3 & 1 \\ 2 & 4 & 4 & 6 & 2 \\ 3 & 6 & 6 & 9 & 6 \\ 1 & 2 & 4 & 5 & 3 \end{pmatrix}$$

cannot have rank 4 since rows one and two are linearly dependent.

But  $\text{rank}(A) \geq 2$  since  $\begin{pmatrix} 9 & 6 \\ 5 & 3 \end{pmatrix}$  is nonsingular.



### Example (cont.)

In fact,  $\text{rank}(A) = 3$  since

$$\begin{pmatrix} 4 & 6 & 2 \\ 6 & 9 & 6 \\ 4 & 5 & 3 \end{pmatrix}$$

is nonsingular.

Note that other singular  $3 \times 3$  submatrices are allowed, such as

$$\begin{pmatrix} 1 & 2 & 2 \\ 2 & 4 & 4 \\ 3 & 6 & 6 \end{pmatrix}.$$



Earlier we showed that

$$\text{rank}(AB) \leq \text{rank}(A),$$

i.e., multiplication by another matrix does not increase the rank of a given matrix, i.e., we can't "fix" a singular system by multiplication.

Now

### Theorem

Let  $A$  and  $E$  be  $m \times n$  matrices. Then

$$\text{rank}(A + E) \geq \text{rank}(A),$$

*provided the entries of  $E$  are "sufficiently small".*



This theorem has at least two **fundamental consequences of practical importance**:

- **Beware!!** A theoretically singular system may become nonsingular, i.e., have a “solution” — just due to round-off error.
- We may want to intentionally “fix” a singular system, so that it has a “solution”. One such strategy is known as **Tikhonov regularization**, i.e.,

$$A\mathbf{x} = \mathbf{b} \longrightarrow (A + \mu I)\mathbf{x} = \mathbf{b},$$

where  $\mu$  is a (small) regularization parameter.



## Proof

We assume that  $\text{rank}(A) = r$  and that we have nonsingular  $P$  and  $Q$  such that we can convert  $A$  to **rank normal form**, i.e.,

$$PAQ = \begin{pmatrix} I_r & O \\ O & O \end{pmatrix}.$$

Then — formally —  $PEQ = \begin{pmatrix} E_{11} & E_{12} \\ E_{21} & E_{22} \end{pmatrix}$  with appropriate blocks  $E_{ij}$ .

This allows us to write

$$P(A + E)Q = \begin{pmatrix} I_r + E_{11} & E_{12} \\ E_{21} & E_{22} \end{pmatrix}.$$





(cont.)

Now, we note that

$$(I - B)(I + B + B^2 + \dots + B^{k-1}) = I - B^k \\ \rightarrow I,$$

provided the entries of  $B$  are “sufficiently small” (i.e., so that  $B^k \rightarrow O$  for  $k \rightarrow \infty$ ).

Therefore  $(I - B)^{-1}$  exists.

This technique is known as the **Neumann series** expansion of the inverse of  $I - B$ .



(cont.)

Now, letting  $B = -E_{11}$ , we know that  $(I_r + E_{11})^{-1}$  exists and we can write

$$\begin{pmatrix} I_r & O \\ -E_{21}(I_r + E_{11})^{-1} & I_{m-r} \end{pmatrix} \begin{pmatrix} I_r + E_{11} & E_{12} \\ E_{21} & E_{22} \end{pmatrix} \begin{pmatrix} I_r & -(I_r + E_{11})^{-1}E_{12} \\ O & I_{n-r} \end{pmatrix} \\ = \begin{pmatrix} I_r + E_{11} & O \\ O & S \end{pmatrix},$$

where  $S = E_{22} - E_{21}(I_r + E_{11})^{-1}E_{12}$  is the **Schur complement** of  $I + E_{11}$  in PAQ.



(cont.)

The Schur complement calculation shows that

$$A + E \sim \begin{pmatrix} I_r + E_{11} & O \\ O & S \end{pmatrix}.$$

But then this rank normal form with invertible diagonal blocks tells us

$$\begin{aligned} \text{rank}(A + E) &= \text{rank}(I_r + E_{11}) + \text{rank}(S) \\ &= \text{rank}(A) + \text{rank}(S) \\ &\geq \text{rank}(A). \end{aligned}$$

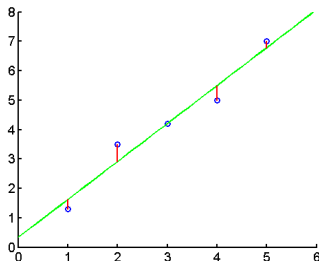


## Linear least squares (linear regression)

**Given:** data  $\{(t_1, b_1), (t_2, b_2), \dots, (t_m, b_m)\}$

**Find:** “best fit” by a line

t	1	2	3	4	5
b	1.3	3.5	4.2	5.0	7.0



### Idea for best fit

Minimize the sum of the squares of the vertical distances of line from the data points.

More precisely, let

$$f(t) = \alpha + \beta t$$

with  $\alpha, \beta$  such that

$$\begin{aligned} \sum_{i=1}^m \varepsilon_i^2 &= \sum_{i=1}^m (f(t_i) - b_i)^2 \\ &= \sum_{i=1}^m (\alpha + \beta t_i - b_i)^2 = G(\alpha, \beta) \quad \rightarrow \quad \text{min} \end{aligned}$$

From calculus, **necessary (and sufficient) condition for minimum**

$$\frac{\partial G(\alpha, \beta)}{\partial \alpha} = 0, \quad \frac{\partial G(\alpha, \beta)}{\partial \beta} = 0.$$

where

$$\frac{\partial G(\alpha, \beta)}{\partial \alpha} = 2 \sum_{i=1}^m (\alpha + \beta t_i - b_i), \quad \frac{\partial G(\alpha, \beta)}{\partial \beta} = 2 \sum_{i=1}^m (\alpha + \beta t_i - b_i) t_i$$



Equivalently,

$$\begin{aligned} \left( \sum_{i=1}^m 1 \right) \alpha + \left( \sum_{i=1}^m t_i \right) \beta &= \sum_{i=1}^m b_i \\ \left( \sum_{i=1}^m t_i \right) \alpha + \left( \sum_{i=1}^m t_i^2 \right) \beta &= \sum_{i=1}^m b_i t_i \end{aligned}$$

which can be written as

$$\mathbf{Q}\mathbf{x} = \mathbf{y}$$

with

$$\mathbf{Q} = \begin{pmatrix} \sum_{i=1}^m 1 & \sum_{i=1}^m t_i \\ m & m \\ \sum_{i=1}^m t_i & \sum_{i=1}^m t_i^2 \end{pmatrix}, \quad \mathbf{x} = \begin{pmatrix} \alpha \\ \beta \end{pmatrix}, \quad \mathbf{y} = \begin{pmatrix} \sum_{i=1}^m b_i \\ m \\ \sum_{i=1}^m b_i t_i \end{pmatrix}$$



We can write each of these sums as inner products:

$$\sum_{i=1}^m 1 = \mathbf{1}^T \mathbf{1}, \quad \sum_{i=1}^m t_i = \mathbf{1}^T \mathbf{t} = \mathbf{t}^T \mathbf{1}, \quad \sum_{i=1}^m t_i^2 = \mathbf{t}^T \mathbf{t}$$

$$\sum_{i=1}^m b_i = \mathbf{1}^T \mathbf{b} = \mathbf{b}^T \mathbf{1}, \quad \sum_{i=1}^m b_i t_i = \mathbf{b}^T \mathbf{t} = \mathbf{t}^T \mathbf{b},$$

where

$$\mathbf{1}^T = (1 \quad \cdots \quad 1), \quad \mathbf{t}^T = (t_1 \quad \cdots \quad t_m), \quad \mathbf{b}^T = (b_1 \quad \cdots \quad b_m)$$

With this notation we have

$$\mathbf{Qx} = \mathbf{y} \iff \begin{pmatrix} \mathbf{1}^T \mathbf{1} & \mathbf{1}^T \mathbf{t} \\ \mathbf{t}^T \mathbf{1} & \mathbf{t}^T \mathbf{t} \end{pmatrix} \mathbf{x} = \begin{pmatrix} \mathbf{1}^T \mathbf{b} \\ \mathbf{t}^T \mathbf{b} \end{pmatrix}$$

$$\iff \mathbf{A}^T \mathbf{A} \mathbf{x} = \mathbf{A}^T \mathbf{b}, \quad \mathbf{A}^T = \begin{pmatrix} \mathbf{1}^T \\ \mathbf{t}^T \end{pmatrix}, \quad \mathbf{A} = (\mathbf{1} \quad \mathbf{t})$$



Therefore we can find the parameters of the line,  $\mathbf{x} = \begin{pmatrix} \alpha \\ \beta \end{pmatrix}$ , by solving the square linear system

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

Also note that since  $\varepsilon_i = \alpha + \beta t_i - b_i$  we have

$$\begin{aligned} \boldsymbol{\varepsilon} = \begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_m \end{pmatrix} &= \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} \alpha + \begin{pmatrix} t_1 \\ \vdots \\ t_m \end{pmatrix} \beta - \begin{pmatrix} b_1 \\ \vdots \\ b_m \end{pmatrix} \\ &= \mathbf{1}\alpha + \mathbf{t}\beta - \mathbf{b} = \mathbf{A}\mathbf{x} - \mathbf{b}. \end{aligned}$$

This implies that

$$G(\alpha, \beta) = \sum_{i=1}^m \varepsilon_i^2 = \boldsymbol{\varepsilon}^T \boldsymbol{\varepsilon} = (\mathbf{A}\mathbf{x} - \mathbf{b})^T (\mathbf{A}\mathbf{x} - \mathbf{b}).$$





## Example

Data:

t	-1	0	1	2	3	4	5	6
b	10	9	7	5	4	3	0	-1

$$\begin{aligned}
 A^T A \mathbf{x} &= A^T \mathbf{b} & \iff & \begin{pmatrix} \sum_{i=1}^8 1 & \sum_{i=1}^8 t_i \\ \sum_{i=1}^8 t_i & \sum_{i=1}^8 t_i^2 \end{pmatrix} \begin{pmatrix} \alpha \\ \beta \end{pmatrix} = \begin{pmatrix} \sum_{i=1}^8 b_i \\ \sum_{i=1}^8 b_i t_i \end{pmatrix} \\
 & & \iff & \begin{pmatrix} 8 & 20 \\ 20 & 92 \end{pmatrix} \begin{pmatrix} \alpha \\ \beta \end{pmatrix} = \begin{pmatrix} 37 \\ 25 \end{pmatrix} \\
 & & \implies & \alpha \approx 8.643, \beta \approx -1.607
 \end{aligned}$$

So that the best fit line to the given data is

$$f(t) \approx 8.643 - 1.607t.$$

## General Least Squares

The general least squares problem behaves analogously to the linear example.

### Theorem

Let  $A$  be a real  $m \times n$  matrix and  $\mathbf{b}$  an  $m$ -vector. *Any vector  $\mathbf{x}$  that minimizes the square of the residual  $A\mathbf{x} - \mathbf{b}$ , i.e.,*

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

*is called a least squares solution of  $A\mathbf{x} = \mathbf{b}$ .*

*The set of all least squares solutions is obtained by solving the normal equations*

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

*Moreover, a unique solution exists if and only if  $\text{rank}(A) = n$  so that*

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

## Proof

The statement about **uniqueness** follows directly from our earlier **theorem** on p. 92 on the normal equations.

To characterize the least squares solutions we **first show that if  $\mathbf{x}$  minimizes  $G(\mathbf{x})$  then  $\mathbf{x}$  satisfies  $A^T A \mathbf{x} = A^T \mathbf{b}$ .**

As in our earlier example, a necessary condition for the minimum is:  
 $\frac{\partial G(\mathbf{x})}{\partial x_i} = 0, i = 1, \dots, n.$

Let's first work out what  $G(\mathbf{x})$  looks like:

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

since  $\mathbf{b}^T A \mathbf{x} = (\mathbf{b}^T A \mathbf{x})^T = \mathbf{x}^T A^T \mathbf{b}$  is a scalar.

(cont.)

Therefore

$$\begin{aligned}
 \frac{\partial G(\mathbf{x})}{\partial x_i} &= \frac{\partial \mathbf{x}^T}{\partial x_i} \mathbf{A}^T \mathbf{A} \mathbf{x} + \mathbf{x}^T \mathbf{A}^T \mathbf{A} \frac{\partial \mathbf{x}}{\partial x_i} - 2 \frac{\partial \mathbf{x}^T}{\partial x_i} \mathbf{A}^T \mathbf{b} \\
 &= \mathbf{e}_i^T \mathbf{A}^T \mathbf{A} \mathbf{x} + \mathbf{x}^T \mathbf{A}^T \mathbf{A} \mathbf{e}_i - 2 \mathbf{e}_i^T \mathbf{A}^T \mathbf{b} \\
 &= 2 \mathbf{e}_i^T \mathbf{A}^T \mathbf{A} \mathbf{x} - 2 \mathbf{e}_i^T \mathbf{A}^T \mathbf{b}
 \end{aligned}$$

since  $\mathbf{x}^T \mathbf{A}^T \mathbf{A} \mathbf{e}_i = (\mathbf{x}^T \mathbf{A}^T \mathbf{A} \mathbf{e}_i)^T = \mathbf{e}_i^T \mathbf{A}^T \mathbf{A} \mathbf{x}$  is a scalar.

This means that

$$\frac{\partial G(\mathbf{x})}{\partial x_i} = 0 \iff (\mathbf{A}^T)_{i*} \mathbf{A} \mathbf{x} = (\mathbf{A}^T)_{i*} \mathbf{b}.$$

If we collect all such conditions (for  $i = 1, \dots, n$ ) in one linear system we get

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

(cont.)

To verify that we indeed have a minimum we show that if  $\mathbf{z}$  is a solution of the normal equations then  $G(\mathbf{z})$  is minimal.

$$\begin{aligned} G(\mathbf{z}) &= (\mathbf{Az} - \mathbf{b})^T (\mathbf{Az} - \mathbf{b}) \\ &= \mathbf{z}^T \mathbf{A}^T \mathbf{A} \mathbf{z} - 2\mathbf{z}^T \mathbf{A}^T \mathbf{b} + \mathbf{b}^T \mathbf{b} \\ &= \mathbf{z}^T (\underbrace{\mathbf{A}^T \mathbf{A} \mathbf{z} - \mathbf{A}^T \mathbf{b}}_{=0}) - \mathbf{z}^T \mathbf{A}^T \mathbf{b} + \mathbf{b}^T \mathbf{b} = -\mathbf{z}^T \mathbf{A}^T \mathbf{b} + \mathbf{b}^T \mathbf{b}. \end{aligned}$$

Now, for any other  $\mathbf{y} = \mathbf{z} + \mathbf{u}$  we have

$$\begin{aligned} G(\mathbf{y}) &= (\mathbf{z} + \mathbf{u})^T \mathbf{A}^T \mathbf{A} (\mathbf{z} + \mathbf{u}) - 2(\mathbf{z} + \mathbf{u})^T \mathbf{A}^T \mathbf{b} + \mathbf{b}^T \mathbf{b} \\ &= G(\mathbf{z}) + \mathbf{u}^T \mathbf{A}^T \mathbf{A} \mathbf{u} + \underbrace{\mathbf{z}^T \mathbf{A}^T \mathbf{A} \mathbf{u}}_{=\mathbf{u}^T \mathbf{A}^T \mathbf{A} \mathbf{z}} + \mathbf{u}^T \mathbf{A}^T \mathbf{A} \mathbf{z} - 2\mathbf{u}^T \underbrace{\mathbf{A}^T \mathbf{b}}_{\mathbf{A}^T \mathbf{A} \mathbf{z}} \\ &= G(\mathbf{z}) + \mathbf{u}^T \mathbf{A}^T \mathbf{A} \mathbf{u} \geq G(\mathbf{z}) \end{aligned}$$

since  $\mathbf{u}^T \mathbf{A}^T \mathbf{A} \mathbf{u} = \sum_{i=1}^m (\mathbf{A}\mathbf{u})_i^2 \geq 0$ .  $\square$

## Remark

Using this framework *we can compute least squares fits from any linear function space.*

## Example

- 1 Let  $f(t) = \alpha_0 + \alpha_1 t + \alpha_2 t^2$ , i.e., we can use quadratic polynomials (or any other degree).
- 2 Let  $f(t) = \alpha_0 + \alpha_1 \sin t + \alpha_2 \cos t$ , i.e., we can use trigonometric polynomials.
- 3 Let  $f(t) = \alpha e^t + \beta \sqrt{t}$ , i.e., we can use just about anything we want.



## Regression in Statistics (BLUE)

One assumes that there is a **random process** that generates **data as a random variable  $Y$**  of the form

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n,$$

where  $X_1, \dots, X_n$  are (input) **random variables** and  $\beta_1, \dots, \beta_n$  are **unknown parameters**.

Now the actually **observed data may be affected by noise**, i.e.,

$$y = Y + \varepsilon = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n + \varepsilon,$$

where  $\varepsilon \sim \mathcal{N}(0, \sigma^2)$  (normally distributed with mean zero and variance  $\sigma^2$ ) is **another random variable denoting the noise**.

**To determine the model parameters  $\beta_1, \dots, \beta_n$**  we now look at measurements, i.e.,

$$y_i = \beta_0 + \beta_1 x_{i,1} + \beta_2 x_{i,2} + \dots + \beta_n x_{i,n} + \varepsilon, \quad i = 1, \dots, m.$$



In matrix-vector form this gives us

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$$

Now, the **least squares solution** of  $\mathbf{X}\boldsymbol{\beta} = \mathbf{y}$ , i.e.,  $\hat{\boldsymbol{\beta}} = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y}$  is in fact the **best linear unbiased estimator (BLUE)** for  $\boldsymbol{\beta}$ .

To show this one needs an **assumption that the error is unbiased**, i.e.,  $\mathbb{E}[\boldsymbol{\varepsilon}] = \mathbf{0}$ .

Then

$$\mathbb{E}[\mathbf{y}] = \mathbb{E}[\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}] = \mathbb{E}[\mathbf{X}\boldsymbol{\beta}] + \mathbb{E}[\boldsymbol{\varepsilon}] = \mathbf{X}\boldsymbol{\beta}$$

and therefore

$$\begin{aligned}\mathbb{E}[\hat{\boldsymbol{\beta}}] &= \mathbb{E}[(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y}] = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbb{E}[\mathbf{y}] \\ &= (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{X}\boldsymbol{\beta} = \boldsymbol{\beta},\end{aligned}$$

so that the **estimator is indeed unbiased**.





## Remark

*One can also show (maybe later) that  $\hat{\beta}$  has minimal variance among all unbiased linear estimators, so it is the best linear unbiased estimator of the model parameters.*

*In fact, the theorem ensuring this is the so-called **Gauss-Markov theorem**.*



## Kriging: a regression approach

**Assume:** the approximate value of a realization of a **zero-mean** (Gaussian) random field is given by a **linear predictor** of the form

$$\hat{Y}_{\mathbf{x}} = \sum_{j=1}^N Y_{\mathbf{x}_j} w_j(\mathbf{x}) = \mathbf{w}(\mathbf{x})^T \mathbf{Y},$$

where  $\hat{Y}_{\mathbf{x}}$  and  $Y_{\mathbf{x}_j}$  are **random variables**,  $\mathbf{Y} = (Y_{\mathbf{x}_1} \ \cdots \ Y_{\mathbf{x}_N})^T$ , and  $\mathbf{w}(\mathbf{x}) = (w_1(\mathbf{x}) \ \cdots \ w_N(\mathbf{x}))^T$  is a vector of **weight functions** at  $\mathbf{x}$ . Since all of the  $Y_{\mathbf{x}_j}$  have zero mean the predictor  $\hat{Y}_{\mathbf{x}}$  is automatically **unbiased**.

**Goal:** to compute “optimal” weights  $w_j^*(\cdot)$ ,  $j = 1, \dots, N$ . To this end, consider the **mean-squared error (MSE)** of the predictor, i.e.,

$$\text{MSE}(\hat{Y}_{\mathbf{x}}) = \mathbb{E} \left[ \left( Y_{\mathbf{x}} - \mathbf{w}(\mathbf{x})^T \mathbf{Y} \right)^2 \right].$$

We now present some details (see [FM15]).



## Covariance Kernel

We need the **covariance kernel**  $K$  of a random field  $Y$  with mean  $\mu(\mathbf{x})$ . It is defined via

$$\begin{aligned}
 \sigma^2 K(\mathbf{x}, \mathbf{z}) &= \text{Cov}(Y_{\mathbf{x}}, Y_{\mathbf{z}}) = \mathbb{E}[(Y_{\mathbf{x}} - \mu(\mathbf{x}))(Y_{\mathbf{z}} - \mu(\mathbf{z}))] \\
 &= \mathbb{E}[(Y_{\mathbf{x}} - \mathbb{E}[Y_{\mathbf{x}}])(Y_{\mathbf{z}} - \mathbb{E}[Y_{\mathbf{z}}])] \\
 &= \mathbb{E}[Y_{\mathbf{x}}Y_{\mathbf{z}} - Y_{\mathbf{x}}\mathbb{E}[Y_{\mathbf{z}}] - \mathbb{E}[Y_{\mathbf{x}}]Y_{\mathbf{z}} + \mathbb{E}[Y_{\mathbf{x}}]\mathbb{E}[Y_{\mathbf{z}}]] \\
 &= \mathbb{E}[Y_{\mathbf{x}}Y_{\mathbf{z}}] - \mathbb{E}[Y_{\mathbf{x}}]\mathbb{E}[Y_{\mathbf{z}}] - \mathbb{E}[Y_{\mathbf{x}}]\mathbb{E}[Y_{\mathbf{z}}] + \mathbb{E}[Y_{\mathbf{x}}]\mathbb{E}[Y_{\mathbf{z}}] \\
 &= \mathbb{E}[Y_{\mathbf{x}}Y_{\mathbf{z}}] - \mathbb{E}[Y_{\mathbf{x}}]\mathbb{E}[Y_{\mathbf{z}}] = \mathbb{E}[Y_{\mathbf{x}}Y_{\mathbf{z}}] - \mu(\mathbf{x})\mu(\mathbf{z}).
 \end{aligned}$$

Therefore, the variance of the random field,

$$\text{Var}(Y_{\mathbf{x}}) = \mathbb{E}[Y_{\mathbf{x}}^2] - \mathbb{E}[Y_{\mathbf{x}}]^2 = \mathbb{E}[Y_{\mathbf{x}}^2] - \mu^2(\mathbf{x}),$$

corresponds to the “diagonal” of the covariance, i.e.,

$$\text{Var}(Y_{\mathbf{x}}) = \sigma^2 K(\mathbf{x}, \mathbf{x}).$$



Let's now work out the MSE:

$$\begin{aligned} \text{MSE}(\hat{Y}_{\mathbf{x}}) &= \mathbb{E} \left[ \left( Y_{\mathbf{x}} - \mathbf{w}(\mathbf{x})^T \mathbf{Y} \right)^2 \right] \\ &= \mathbb{E}[Y_{\mathbf{x}} Y_{\mathbf{x}}] - 2\mathbb{E}[Y_{\mathbf{x}} \mathbf{w}(\mathbf{x})^T \mathbf{Y}] + \mathbb{E}[\mathbf{w}(\mathbf{x})^T \mathbf{Y} \mathbf{Y}^T \mathbf{w}(\mathbf{x})] \end{aligned}$$

Now use  $\mathbb{E}[Y_{\mathbf{x}} Y_{\mathbf{z}}] = K(\mathbf{x}, \mathbf{z})$  (the **covariance**, since  $\mathbf{Y}$  is centered):

$$\text{MSE}(\hat{Y}_{\mathbf{x}}) = \sigma^2 K(\mathbf{x}, \mathbf{x}) - 2\mathbf{w}(\mathbf{x})^T (\sigma^2 \mathbf{k}(\mathbf{x})) + \mathbf{w}(\mathbf{x})^T (\sigma^2 \mathbf{K}) \mathbf{w}(\mathbf{x}),$$

where

$$\sigma^2 \mathbf{k}(\mathbf{x}) = \sigma^2 (k_1(\mathbf{x}) \ \cdots \ k_N(\mathbf{x}))^T: \text{ with}$$

$$\sigma^2 k_j(\mathbf{x}) = \sigma^2 K(\mathbf{x}, \mathbf{x}_j) = \mathbb{E}[Y_{\mathbf{x}} Y_{\mathbf{x}_j}]$$

**K**: the covariance matrix has entries  $\sigma^2 K(\mathbf{x}_i, \mathbf{x}_j) = \mathbb{E}[Y_{\mathbf{x}_i} Y_{\mathbf{x}_j}]$

Finding the minimum MSE is straightforward. Differentiation and equating to zero yields

$$-2\mathbf{k}(\mathbf{x}) + 2\mathbf{K}\mathbf{w}(\mathbf{x}) = 0,$$

and so the optimum weight vector is

$$\mathbf{w}^*(\mathbf{x}) = \mathbf{K}^{-1} \mathbf{k}(\mathbf{x}).$$



We have shown that the (simple) kriging predictor

$$\hat{Y}_{\mathbf{x}} = \mathbf{k}(\mathbf{x})^T \mathbf{K}^{-1} \mathbf{Y}$$

is the **best** (in the MSE sense) **linear unbiased predictor** (BLUP).

Since we are given the observations  $\mathbf{y}$  as realizations of  $\mathbf{Y}$  we can compute the **prediction**

$$\hat{y}_{\mathbf{x}} = \mathbf{k}(\mathbf{x})^T \mathbf{K}^{-1} \mathbf{y}.$$



The MSE of the kriging predictor with optimal weights  $\hat{\mathbf{w}}^*(\cdot)$ ,

$$\mathbb{E} \left[ \left( Y_{\mathbf{x}} - \hat{Y}_{\mathbf{x}} \right)^2 \right] = \sigma^2 \left( K(\mathbf{x}, \mathbf{x}) - \mathbf{k}(\mathbf{x})^T \mathbf{K}^{-1} \mathbf{k}(\mathbf{x}) \right),$$

is known as the **kriging variance**.

It allows us to give **confidence intervals** for our prediction. It also gives rise to a criterion for choosing an **optimal parametrization** of the family of covariance kernels used for prediction.

### Remark

*For Gaussian random fields the BLUP is also the best **nonlinear** unbiased predictor (see, e.g., [BTA04, Chapter 2]).*



## Remark

- 1 The *simple kriging approach just described is precisely how Krige [Kri51] introduced the method:*
  - The unknown value to be predicted is given by a *weighted average of the observed values*, where the *weights depend on the prediction location*.
  - Usually *one assigns a smaller weight to observations further away from  $\mathbf{x}$* .

*The latter statement implies that one should be using kernels whose associated weights decay away from  $\mathbf{x}$ . Positive definite translation invariant kernels have this property.*

- 2 *More advanced kriging variants are discussed in papers such as [SWMW89, SSS13], or books such as [Cre93, Ste99, BTA04].*



# References I

- [BTA04] A. Berlinet and C. Thomas-Agnan, *Reproducing Kernel Hilbert Spaces in Probability and Statistics*, Kluwer, Dordrecht, 2004.
- [Cre93] N. Cressie, *Statistics for Spatial Data*, revised ed., Wiley–Interscience, New York, 1993.
- [FM15] G. E. Fasshauer and M. J. McCourt, *Kernel-based Approximation Methods using MATLAB*, Interdisciplinary Mathematical Sciences, vol. 19, World Scientific Publishing, Singapore, 2015.
- [Kri51] D. G. Krige, *A statistical approach to some basic mine valuation problems on the Witwatersrand*, J. Chem. Met. & Mining Soc., S. Africa **52** (1951), no. 6, 119–139.
- [Mey00] Carl D. Meyer, *Matrix Analysis and Applied Linear Algebra*, SIAM, Philadelphia, PA, 2000.
- [SSS13] M. Scheuerer, R. Schaback, and M. Schlather, *Interpolation of spatial data — a stochastic or a deterministic problem?*, Eur. J. Appl. Math. **24** (2013), no. 4, 601–629.





# References II

- [Ste99] M. L. Stein, *Interpolation of Spatial Data: Some Theory for Kriging*, Springer, Berlin; New York, 1999.
- [SWMW89] Jerome Sacks, W. J. Welch, T. J. Mitchell, and H. P. Wynn, *Design and analysis of computer experiments*, Stat. Sci. **4** (1989), no. 4, 409–423.

