

# Lecture 3.1: The rank

Optimization and Computational Linear Algebra for Data Science

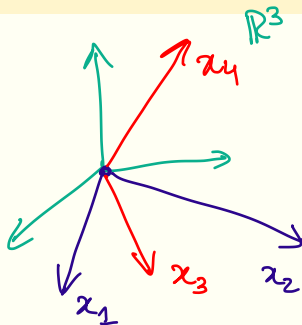
# Rank of a family of vectors

## Definition

We define the rank of a family  $x_1, \dots, x_k$  of vectors of  $\mathbb{R}^n$  as the dimension of its span:

$$\text{rank}(x_1, \dots, x_k) \stackrel{\text{def}}{=} \dim(\text{Span}(x_1, \dots, x_k)) = 0, 1, 2, 3, \dots$$

Example



$$\begin{aligned}\text{rank}(x_1, x_2) &= 2 \\ \text{rank}(x_1, x_2, x_3) &= 2 \\ \text{rank}(x_1, \dots, x_4) &= 4\end{aligned}$$

# Rank of a matrix

## Definition

Let  $M \in \mathbb{R}^{n \times m}$ . Let  $c_1, \dots, c_m \in \mathbb{R}^n$  be its columns. We define

$$\text{rank}(M) \stackrel{\text{def}}{=} \text{rank}(c_1, \dots, c_m) = \dim(\text{Im}(M)).$$

$$M = \begin{pmatrix} | & & | \\ c_1 & \dots & c_m \\ | & & | \end{pmatrix}$$

$\text{Span}(c_1, \dots, c_m)$

# Example

$$M = \begin{pmatrix} 1 & 2 & 3 \\ 0 & 4 & 5 \end{pmatrix}$$

$c_1 \quad c_2 \quad c_3$

$$\underline{\underline{\text{Span}(c_1, c_2, c_3)}} = ?$$

- $c_1, c_2$  are two lin. indep vectors of  $\mathbb{R}^2$   
 $\rightarrow (c_1, c_2)$  is a basis of  $\mathbb{R}^2$   
 $\rightarrow \text{Span}(c_1, c_2) = \mathbb{R}^2$
- $\text{Span}(c_1, c_2, c_3) = \underline{\mathbb{R}^2}$

$$\text{rank}(M) = 2.$$

# « Rank of columns = rank of rows »

## Proposition

Let  $M \in \mathbb{R}^{n \times m}$ . Let  $r_1, \dots, r_n \in \mathbb{R}^m$  be the rows of  $M$  and  $c_1, \dots, c_m \in \mathbb{R}^n$  be its columns. Then we have

$$\underline{\text{rank}(r_1, \dots, r_n)} = \underline{\text{rank}(c_1, \dots, c_m)} = \underline{\text{rank}(M)}.$$

$$M = \begin{pmatrix} \text{---} r_1 \text{---} \\ \vdots \\ \text{---} r_n \text{---} \end{pmatrix} = \begin{pmatrix} | & & | \\ c_1 & \dots & c_m \\ | & & | \end{pmatrix}$$

# The rank in Data Science

Consider a matrix  $M$  of size  $1000 \times 500$ :

$$M = \begin{pmatrix} - & r_1 & - \\ & \vdots & \\ - & r_{1000} & - \end{pmatrix} \quad \updownarrow \quad \underline{1000 \text{ rows}}$$

What does it mean to say that «  $\text{rank}(M) = 5$  »?

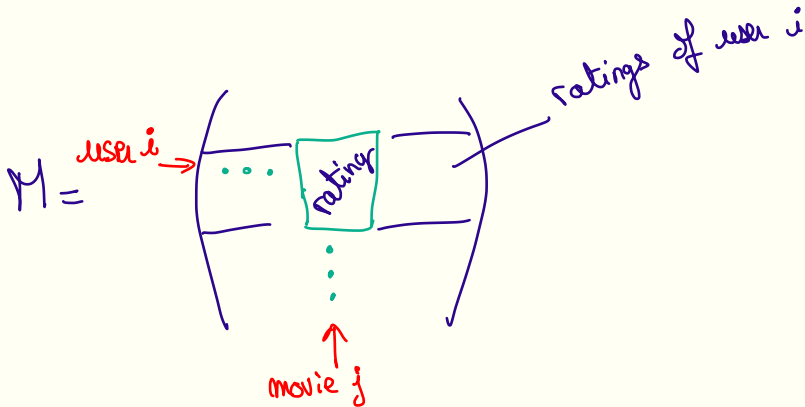
- $\dim \text{Span}(\underline{r_1, \dots, r_{1000}}) = 5$
- Consequently one can generate any row of  $M$  using only 5 vectors.

using  
lin. comb.

# The rank in Data Science

Imagine now that

- ❑ The rows of  $M$  corresponds to Netflix's users.
- ❑ The columns of  $M$  corresponds to Netflix's movies.
- ❑ The entry  $M_{i,j}$  is rating of the movie  $j$  by the user  $i$ , assuming that all the users have rated all the movies.



# The rank in Data Science

Imagine now that

- ❖ The rows of  $M$  corresponds to Netflix's users.
- ❖ The columns of  $M$  corresponds to Netflix's movies.
- ❖ The entry  $M_{i,j}$  is rating of the movie  $j$  by the user  $i$ , assuming that all the users have rated all the movies.

**Claim:** the rank of  $M$  is "small".

- ❖ The ratings of a user can be obtained as a linear combination of a small number of « profiles ».
- ❖ In practice, we do not have access to the full matrix, so we can use this assumption to predict the missing entries.