# Statistical Learning II

## Lecture 1 - Recap of maths

_____

**Bruno Loureiro**

@ CSD, DI-ENS & CNRS

brloureiro@gmail.com

*DL3 IASO, Université Paris Dauphine-PSL*
*10.09.2025*

# Recap of Linear Algebra

The bread of statistical learning

# The Euclidean space

The Euclidean space $\mathbb{R}^d$ is the vector space of $d$-tuples:

$$\boldsymbol{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix} \in \mathbb{R}^d \ \ (\mathbb{R}^{d \times 1})$$

"column vector"

# The Euclidean space

The Euclidean space $\mathbb{R}^d$ is the vector space of $d$-tuples:

$$\boldsymbol{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix} \in \mathbb{R}^d \ (\mathbb{R}^{d \times 1})$$

"column vector"

Recall, $\mathbb{R}^d$ is a vector space of dimension $d$ with basis:

$$\boldsymbol{e}_i = \begin{bmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix}$$

Position $i$

# The Euclidean space

The Euclidean space is endowed with an inner (or scalar) product

$$\boldsymbol{u}, \boldsymbol{v} \in \mathbb{R}^d \qquad\qquad \langle \boldsymbol{u}, \boldsymbol{v} \rangle = \sum_{i=1}^{d} u_i v_i$$

# The Euclidean space

The Euclidean space is endowed with an inner (or scalar) product

$$\boldsymbol{u}, \boldsymbol{v} \in \mathbb{R}^d \qquad \langle \boldsymbol{u}, \boldsymbol{v} \rangle = \sum_{i=1}^{d} u_i v_i$$

Which induces a natural notion of distance and size:

$$||\boldsymbol{u}||_2^2 = \langle \boldsymbol{u}, \boldsymbol{u} \rangle = \sum_{i=1}^{d} u_i^2 \qquad d(\boldsymbol{u}, \boldsymbol{v}) = ||\boldsymbol{u} - \boldsymbol{v}||_2^2$$

"Euclidean or $\ell_2$ norm"        "Euclidean distance"

We say two vectors $\boldsymbol{u}, \boldsymbol{v} \in \mathbb{R}^d$ are orthogonal if $\langle \boldsymbol{u}, \boldsymbol{v} \rangle = 0$

# Euclidean geometry

$$||u||_2^2 = \langle u, u \rangle = \sum_{i=1}^{d} u_i^2 \qquad d(u, v) = ||u - v||_2^2$$

"Euclidean or $\ell_2$ norm"          "Euclidean distance"

They correspond to our intuitive notion of geometry in the plane



$$d(u, v) = ||u - v||_2^2$$

$$\cos(\theta) = \frac{\langle u, v \rangle}{||u||_2 ||v||_2}$$

# Euclidean geometry

They correspond to our intuitive notion of geometry in the plane

$$u$$

$$d(u, v) = ||u - v||_2^2$$

$$\theta$$

$$\mathbf{0}_d$$

$$v$$

$$\cos(\theta) = \frac{\langle u, v \rangle}{||u||_2 ||v||_2}$$

In particular, we  say two vectors $u, v \in \mathbb{R}^d$ are orthogonal if

$$\langle u, v \rangle = 0$$

$$u$$

$$v$$

# Other norms

One can define other notions of size in $\mathbb{R}^d$

$$||\boldsymbol{u}||_p = \left( \sum_{i=1}^{d} u_i^p \right)^{1/p} \qquad p \geq 1$$

"$\ell_p$ norm"

# Other norms

One can define other notions of size in $\mathbb{R}^d$

$$||\boldsymbol{u}||_p = \left( \sum_{i=1}^{d} u_i^p \right)^{1/p} \qquad p \geq 1$$

"$\ell_p$ norm"

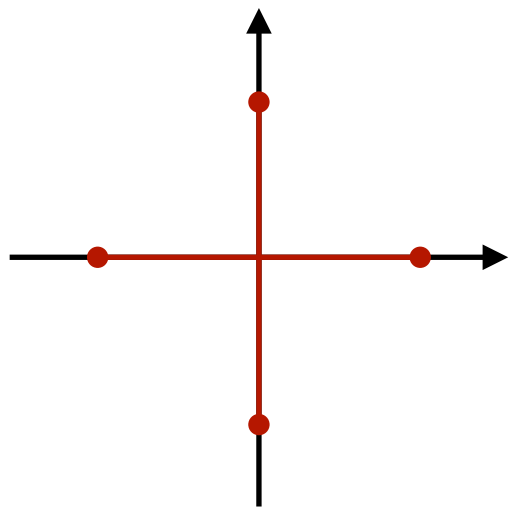⚠️ $||\cdot||_p$ is not associated to an inner product for $p \neq 2$

# Other norms

One can define other notions of size in $\mathbb{R}^d$

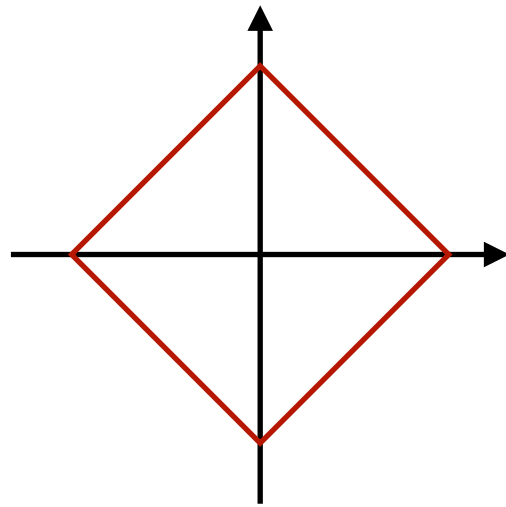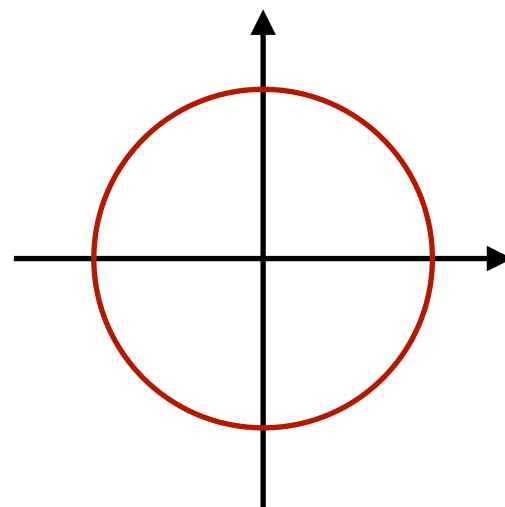$$||\boldsymbol{u}||_p = \left( \sum_{i=1}^{d} u_i^p \right)^{1/p} \qquad p \geq 1$$

"$\ell_p$ norm"

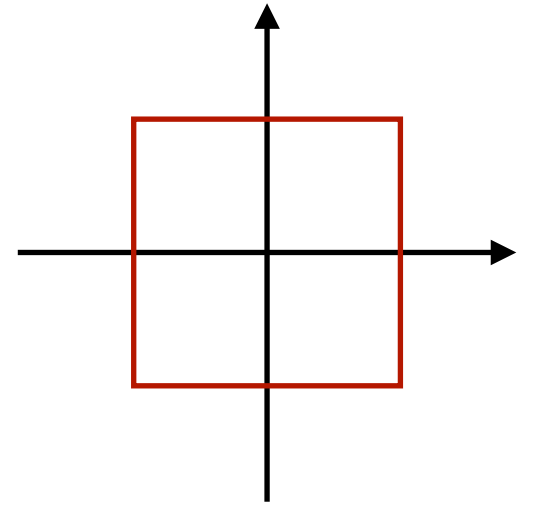⚠ $||\cdot||_p$ is not associated to an inner product for $p \neq 2$



$\ell_0$        $\ell_1$        $\ell_2$        $\ell_\infty$

⚠ Not a norm

# Matrices

A real-valued matrix $A \in \mathbb{R}^{n \times d}$ is a table of real numbers.

$$A = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nd} \end{bmatrix} \in \mathbb{R}^{n \times d}$$

# Matrices

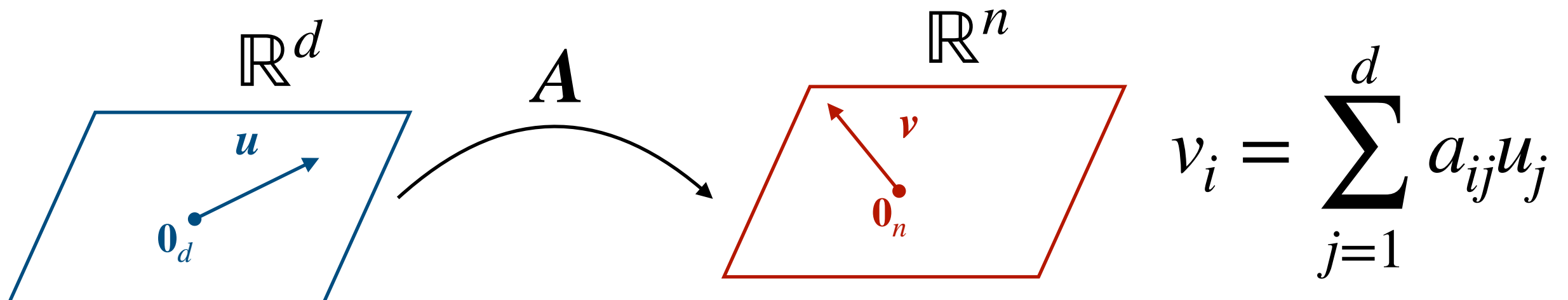A real-valued matrix $A \in \mathbb{R}^{n \times d}$ is a table of real numbers.

$$A = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nd} \end{bmatrix} \in \mathbb{R}^{n \times d}$$

It is most often used to describe the coordinates of linear transformations $A : \mathbb{R}^d \to \mathbb{R}^n$ with respect to a basis.



$$v_i = \sum_{j=1}^{d} a_{ij} u_j$$

# Matrices

A real-valued matrix $A \in \mathbb{R}^{n \times d}$ is a table of real numbers.

$$A = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nd} \end{bmatrix} \in \mathbb{R}^{n \times d}$$

$$A_1 \qquad A_2 \qquad\qquad A_d$$

- The columns of $A \in \mathbb{R}^{n \times d}$ are vectors $A_i \in \mathbb{R}^n$ with $(A_i)_j = a_{ij}$

  "Column space" $\mathrm{col}(A) = \mathrm{span}(A_1, \cdots, A_d) \subset \mathbb{R}^n$

# Matrices

A real-valued matrix $A \in \mathbb{R}^{n \times d}$ is a table of real numbers.

$$A = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nd} \end{bmatrix} \begin{matrix} a_1 \\ a_2 \\ \\ a_n \end{matrix} \in \mathbb{R}^{n \times d}$$

- The columns of $A \in \mathbb{R}^{n \times d}$ are vectors $A_i \in \mathbb{R}^n$ with $(A_i)_j = a_{ij}$

  "Column space" $\mathrm{col}(A) = \mathrm{span}(A_1, \cdots, A_d) \subset \mathbb{R}^n$

- The rows of $A \in \mathbb{R}^{n \times d}$ are vectors $a_j \in \mathbb{R}^d$ with $(a_j)_i = a_{ij}$

  "Row space" of $\mathrm{row}(A) = \mathrm{span}(a_1, \cdots, a_n) \subset \mathbb{R}^d$

# Flattening matrices

The space of matrices $A \in \mathbb{R}^{n \times d}$ is itself a vector space of dimension $nd$. Therefore we can identify:

$$\mathbb{R}^{n \times d} \simeq \mathbb{R}^{nd}$$

By flattening the matrices into vectors.

$$A = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nd} \end{bmatrix} \mapsto \begin{bmatrix} a_{11} \\ a_{12} \\ \vdots \\ a_{1n} \\ a_{21} \\ \vdots \end{bmatrix} \in \mathbb{R}^{nd}$$

# Rank of a matrix

- The rank of a matrix $A \in \mathbb{R}^{n \times d}$ is the dimension of column space

$$\mathrm{rank}(A) = \dim(\mathrm{col}(A))$$

This is equivalent to the number of independent columns.

# Rank of a matrix

- The rank of a matrix $A \in \mathbb{R}^{n \times d}$ is the dimension of column space

$$\text{rank}(A) = \dim(\text{col}(A))$$

This is equivalent to the number of independent columns.

**Proposition**

$$\text{rank}(A) = \dim(\text{col}(A)) = \dim(\text{row}(A))$$

# Rank of a matrix

- The rank of a matrix $A \in \mathbb{R}^{n \times d}$ is the dimension of column space

$$\mathrm{rank}(A) = \dim(\mathrm{col}(A))$$

This is equivalent to the number of independent columns.

**Proposition**
$$\mathrm{rank}(A) = \dim(\mathrm{col}(A)) = \dim(\mathrm{row}(A))$$

- A matrix $A \in \mathbb{R}^{n \times d}$ is said to be full-rank if

$$\mathrm{rank}(A) = \min(n, d)$$

# Another point of view

- Alternatively, we can see the column space $\mathrm{col}(A) \subset \mathbb{R}^n$ as The image of the associated linear map.

$$\mathrm{im}(A) = \mathrm{col}(A) = \{v \in \mathbb{R}^n : Au = v \text{ for some } u \in \mathbb{R}^d\}$$

# Another point of view

- Alternatively, we can see the column space $\mathrm{col}(A) \subset \mathbb{R}^n$ as The image of the associated linear map.

$$\mathrm{im}(A) = \mathrm{col}(A) = \{v \in \mathbb{R}^n : Au = v \text{ for some } u \in \mathbb{R}^d\}$$

- The null-space or kernel of a matrix $A \in \mathbb{R}^{n \times d}$ is defined as:

$$\ker(A) = \{u \in \mathbb{R}^d : Au = 0\}$$

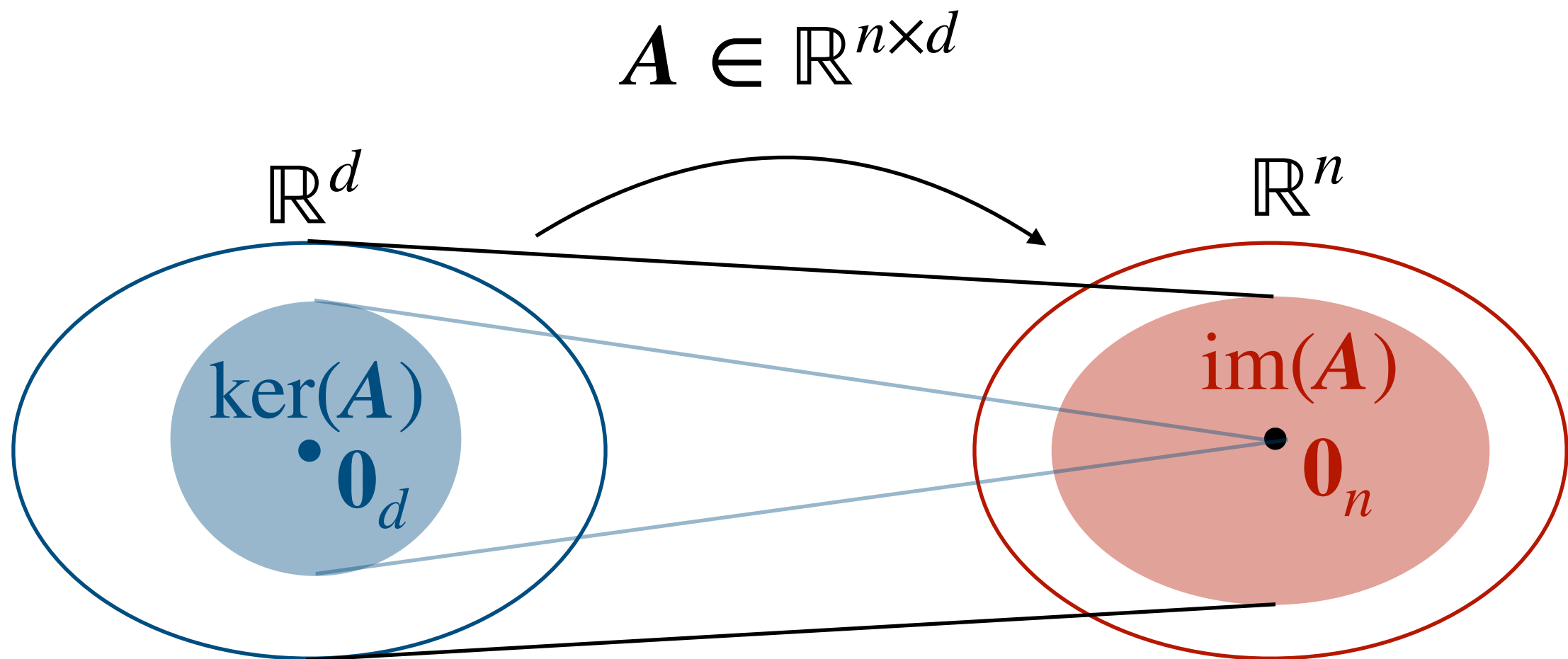Note that $\ker(A) \subset \mathbb{R}^d$
and $\mathbf{0} \in \ker(A)$

# Image and null-space

Proposition

Let $A \in \mathbb{R}^{n \times d}$ denote a linear map. We have:
$$\text{rank}(A) + \dim(\ker(A)) = n$$

$$A \in \mathbb{R}^{n \times d}$$



$\mathbb{R}^d$

$\ker(A)$
$\mathbf{0}_d$

$\mathbb{R}^n$

$\text{im}(A)$
$\mathbf{0}_n$

# Matrix inverse

A square matrix $A \in \mathbb{R}^{d \times d}$ is said to be invertible if there exists $B \in \mathbb{R}^{d \times d}$ such that:

$$AB = I_d$$

In this case, we denote $B = A^{-1}$.

# Matrix inverse

A square matrix $A \in \mathbb{R}^{d \times d}$ is said to be invertible if there exists $B \in \mathbb{R}^{d \times d}$ such that:

$$AB = I_d$$

In this case, we denote $B = A^{-1}$.

⚠️ For any invertible matrix $A \in \mathbb{R}^{d \times d}$
$(A^{-1})^{-1} = A$.

# Matrix inverse

A square matrix $A \in \mathbb{R}^{d \times d}$ is said to be invertible if there exists $B \in \mathbb{R}^{d \times d}$ such that:

$$AB = I_d$$

In this case, we denote $B = A^{-1}$.

⚠️ For any invertible matrix $A \in \mathbb{R}^{d \times d}$ $(A^{-1})^{-1} = A$.

**Proposition**

A square matrix $A \in \mathbb{R}^{d \times d}$ is invertible if and only if it is full-rank
$$\mathrm{rank}(A) = d$$

# Matrix inverse

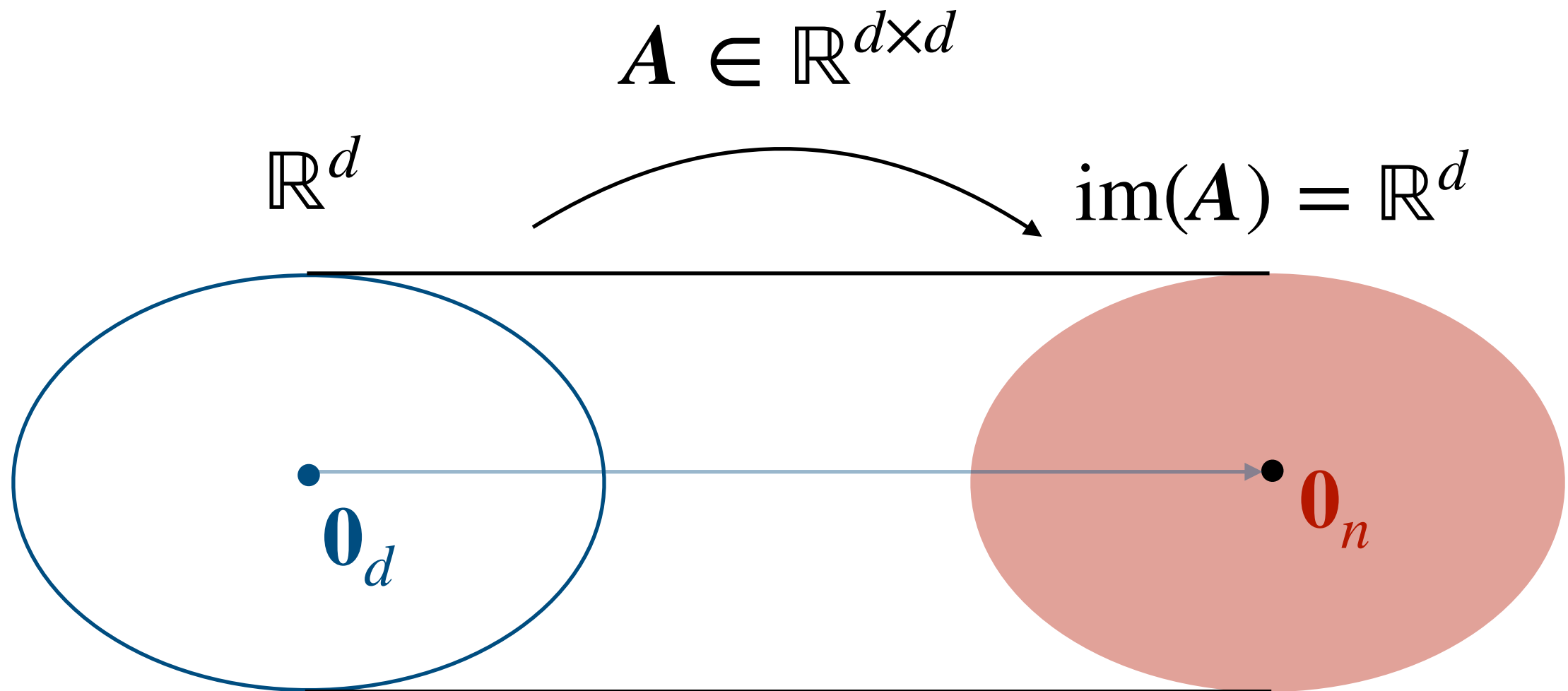**Proposition**

A square matrix $A \in \mathbb{R}^{d \times d}$ is invertible if and only if it is full-rank

$$\text{rank}(A) = d$$

$$A \in \mathbb{R}^{d \times d}$$

$\mathbb{R}^d$

$\text{im}(A) = \mathbb{R}^d$

$\mathbf{0}_d$

$\mathbf{0}_n$

# Matrix transpose

- The transpose of a matrix $A \in \mathbb{R}^{n \times d}$ with elements $a_{ij}$ the matrix with $A^{\top} \in \mathbb{R}^{d \times n}$ with elements $a_{ji}$

$$A = \quad\quad\quad\quad A^{\top} = $$

# Matrix transpose

- The transpose of a matrix $A \in \mathbb{R}^{n \times d}$ with elements $a_{ij}$ the matrix with $A^\top \in \mathbb{R}^{d \times n}$ with elements $a_{ji}$

$$A = \qquad A^\top =$$

- We have:

$$(A^\top)^\top = A$$

$$(aA + bB)^\top = aA^\top + bB^\top$$

$$(A^{-1})^\top = (A^\top)^{-1}$$

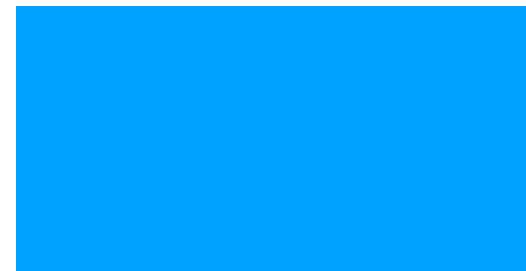$$(AB)^\top = B^\top A^\top \qquad 🧐 \ \underline{\text{Exercise}}\text{: check this.}$$

# Matrix transpose

- The transpose of a matrix $A \in \mathbb{R}^{n \times d}$ with elements $a_{ij}$ the matrix with $A^\top \in \mathbb{R}^{d \times n}$ with elements $a_{ji}$

$$A = \qquad \qquad A^\top = $$

- Note that by seeing $u, v \in \mathbb{R}^{d \times 1}$ as column vectors, we can also write the Euclidean inner product as:

$$\langle u, v \rangle = u^\top v \qquad \text{🤨 Exercise: check this.}$$

# Matrix trace

- The trace of a square matrix $A \in \mathbb{R}^{d \times d}$ is the sum of its diagonal:

$$\text{Tr } A = \sum_{i=1}^{d} a_{ii}$$

# Matrix trace

- The trace of a square matrix $A \in \mathbb{R}^{d \times d}$ is the sum of its diagonal:

$$\text{Tr}\, A = \sum_{i=1}^{d} a_{ii}$$

- It satisfies:

$$\text{Tr}\, AB = \text{Tr}\, BA$$

$$\text{Tr}\, (aA + bB) = a\text{Tr}\, A + b\text{Tr}\, B$$

$$\text{Tr}\, A^{\top} = \text{Tr}\, A$$

🧐 Exercise: check this.

# Symmetric matrices

- A square matrix $\boldsymbol{A} \in \mathbb{R}^{d \times d}$ is <span style="color:red">symmetric</span> if $\boldsymbol{A}^{\top} = \boldsymbol{A}$

# Symmetric matrices

- A square matrix $\boldsymbol{A} \in \mathbb{R}^{d \times d}$ is <span style="color:red">symmetric</span> if $\boldsymbol{A}^\top = \boldsymbol{A}$

⚠️ For any $\boldsymbol{A} \in \mathbb{R}^{n \times d}$, $\boldsymbol{A}^\top \boldsymbol{A} \in \mathbb{R}^{d \times d}$ and $\boldsymbol{A}\boldsymbol{A}^\top \in \mathbb{R}^{n \times n}$ are symmetric matrices.

# Symmetric matrices

- A square matrix $A \in \mathbb{R}^{d \times d}$ is symmetric if $A^\top = A$

⚠️ For any $A \in \mathbb{R}^{n \times d}$, $A^\top A \in \mathbb{R}^{d \times d}$ and $AA^\top \in \mathbb{R}^{n \times n}$ are symmetric matrices.

Letting $a_i \in \mathbb{R}^d$ denote the rows of $A \in \mathbb{R}^{n \times d}$, we have:

$$(AA^\top)_{ij} = \langle a_i, a_j \rangle$$

🧐 Exercise: check this.

Note: a similar representation holds for columns of $A$

# Orthogonal matrices

- A square matrix $A \in \mathbb{R}^{d \times d}$ is orthogonal if $A^\top = A^{-1}$

# Orthogonal matrices

- A square matrix $A \in \mathbb{R}^{d \times d}$ is orthogonal if $A^\top = A^{-1}$

Orthogonal matrices preserve the norm and distance between vectors (they are isometries):

$$||Au||_2 = ||u||_2$$

🧐 Exercise: check this.

# Orthogonal matrices

- A square matrix $A \in \mathbb{R}^{d \times d}$ is orthogonal if $A^{\top} = A^{-1}$

Orthogonal matrices preserve the norm and distance between vectors (they are isometries):

$$||Au||_2 = ||u||_2$$

🧐 <u>Exercise</u>: check this.

Geometrically, they define rotations

$$e_2 = \begin{bmatrix} 0 \\ 1 \end{bmatrix} \qquad A = \begin{bmatrix} 0 & -1 \\ 1 & 0 \end{bmatrix}$$

$$\mathbf{0} \qquad e_1 = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

# Projection matrix

- A square matrix $A \in \mathbb{R}^{d \times d}$ is a projection if $A^2 = A$

Moreover, if $A$ is also orthogonal, we call it a orthogonal projection.



$$P = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

# Projection matrix

- A square matrix $A \in \mathbb{R}^{d \times d}$ is a projection if $A^2 = A$

Moreover, if $A$ is also orthogonal, we call it a orthogonal projection.

<div style="background-color:#b6f0b0;">

**Proposition**

Any $v \in \mathbb{R}^d$ can be uniquely written as:

$$v = u + Av \qquad u \in \ker(A)$$

</div>

# Projection matrix

- A square matrix $A \in \mathbb{R}^{d \times d}$ is a projection if $A^2 = A$

Moreover, if $A$ is also orthogonal, we call it a orthogonal projection.

> **Proposition**
>
> Any $v \in \mathbb{R}^d$ can be uniquely written as:
>
> $$v = u + Av \qquad\qquad u \in \ker(A)$$

⚠️ The only projection matrix which is invertible is the identity.

# Eigen-(values, vectors)

Let $A \in \mathbb{R}^{d \times d}$ denote a square matrix. An eigenvector is a vector that is only re-scaled under the action of $A$:

$$Av = \lambda v$$

Where $\lambda \in \mathbb{R}$ is known as an eigenvalue.

# Eigen-(values, vectors)

Let $A \in \mathbb{R}^{d \times d}$ denote a square matrix. An eigenvector is a vector that is only re-scaled under the action of $A$:

$$Av = \lambda v$$

Where $\lambda \in \mathbb{R}$ is known as an eigenvalue.

We call the set of eigenvalues the spectrum of $A$:

$$\text{spec}(A) = \{\lambda \in \mathbb{R} : Av = \lambda v\}$$

# Eigen-(values, vectors)

Let $A \in \mathbb{R}^{d \times d}$ denote a square matrix. An <span style="color:#b33">eigenvector</span> is a vector that is only re-scaled under the action of $A$:

$$Av = \lambda v$$

Where $\lambda \in \mathbb{R}$ is known as an <span style="color:#b33">eigenvalue</span>.

We call the set of eigenvalues the spectrum of $A$:

$$\operatorname{spec}(A) = \{\lambda \in \mathbb{R} : Av = \lambda v\}$$

- A square matrix $A \in \mathbb{R}^{d \times d}$ can have at most $d$ independent eigenvectors.

- An eigenvalue $\lambda$ can be associated to more than one independent eigenvector.

# Positive matrices

- A square matrix $A \in \mathbb{R}^{d \times d}$ is called <span style="color:#a00">positive definite</span> if all eigenvalues are positive:

$$\lambda \in \operatorname{spec}(A) \Rightarrow \lambda > 0$$

# Positive matrices

- A square matrix $A \in \mathbb{R}^{d \times d}$ is called <span style="color:darkred">positive definite</span> if all eigenvalues are positive:

$$\lambda \in \operatorname{spec}(A) \Rightarrow \lambda > 0$$

- A square matrix $A \in \mathbb{R}^{d \times d}$ is called <span style="color:darkred">positive semi-definite</span> if all eigenvalues are non-negative:

$$\lambda \in \operatorname{spec}(A) \Rightarrow \lambda \geq 0$$

# Positive matrices

- A square matrix $A \in \mathbb{R}^{d \times d}$ is called positive definite if all eigenvalues are positive:

$$\lambda \in \mathrm{spec}(A) \Rightarrow \lambda > 0$$

- A square matrix $A \in \mathbb{R}^{d \times d}$ is called positive semi-definite if all eigenvalues are non-negative:

$$\lambda \in \mathrm{spec}(A) \Rightarrow \lambda \geq 0$$

**Proposition**

Symmetric matrices $A \in \mathbb{R}^{d \times d}$ are positive semi-definite

🧐 Exercise: prove this.

# Positive matrices

- A square matrix $A \in \mathbb{R}^{d \times d}$ is called <span style="color:red">positive definite</span> if all eigenvalues are positive:

$$\lambda \in \mathrm{spec}(A) \Rightarrow \lambda > 0$$

- A square matrix $A \in \mathbb{R}^{d \times d}$ is called <span style="color:red">positive semi-definite</span> if all eigenvalues are non-negative:

$$\lambda \in \mathrm{spec}(A) \Rightarrow \lambda \geq 0$$

**Proposition**

Symmetric matrices $A \in \mathbb{R}^{d \times d}$ are positive semi-definite

⚠️ not necessarily positive definite. 🧐 <u>Exercise</u>: prove this.

# Spectral theorem

Any symmetric matrix $A \in \mathbb{R}^{d \times d}$ can be decomposed as

$$A = UDU^\top$$

$U \in \mathbb{R}^{d \times d}$ are orthogonal matrices and $D$ is a diagonal matrix with elements given by the eigenvalues.

# Spectral theorem

**Theorem**

Any symmetric matrix $A \in \mathbb{R}^{d \times d}$ can be decomposed as

$$A = UDU^{\top}$$

$U \in \mathbb{R}^{d \times d}$ are orthogonal matrices and $D$ is a diagonal matrix with elements given by the eigenvalues.

We can equivalently write the spectral decomposition as:

$$A = \sum_{i=1}^{\text{rank}(A)} \lambda_i v_i v_i^{\top}$$

Where $v_i \in \mathbb{R}^d$ are orthonormal eigenvectors.

# Important facts

- The trace of a symmetric matrix $A \in \mathbb{R}^{d \times d}$ is the sum of its eigenvalues

$$\mathrm{Tr}\, A = \sum_{i=1}^{d} \lambda_i$$

# Important facts

- The trace of a symmetric matrix $A \in \mathbb{R}^{d \times d}$ is the sum of its eigenvalues

$$\mathrm{Tr}\, A = \sum_{i=1}^{d} \lambda_i$$

- A square matrix $A \in \mathbb{R}^{d \times d}$ is invertible i.f.f. $0 \notin \mathrm{spec}(A)$

# Important facts

- The trace of a symmetric matrix $A \in \mathbb{R}^{d \times d}$ is the sum of its eigenvalues

$$\mathrm{Tr}\, A = \sum_{i=1}^{d} \lambda_i$$

- A square matrix $A \in \mathbb{R}^{d \times d}$ is invertible i.f.f. $0 \notin \mathrm{spec}(A)$

- The eigenvalues of a projection matrix $P \in \mathbb{R}^{d \times d}$ are $0$ or $1$

$$P = \sum_{i=1}^{\mathrm{rank}(P)} v_i v_i^{\top}$$

🧐 <u>Exercise</u>: show this.

Moreover, $P \in \mathbb{R}^{d \times d}$ is orthogonal if $v_i$ are orthogonal vectors.

# Singular value decomposition

Note that for any real matrix $\boldsymbol{A} \in \mathbb{R}^{n \times d}$, $\boldsymbol{A}^\top \boldsymbol{A} \in \mathbb{R}^{d \times d}$ and $\boldsymbol{A} \boldsymbol{A}^\top \in \mathbb{R}^{n \times n}$ are a symmetric matrices.

# Singular value decomposition

Note that for any real matrix $\boldsymbol{A} \in \mathbb{R}^{n \times d}$, $\boldsymbol{A}^\top \boldsymbol{A} \in \mathbb{R}^{d \times d}$ and $\boldsymbol{A}\boldsymbol{A}^\top \in \mathbb{R}^{n \times n}$ are a symmetric matrices.

Therefore, $\boldsymbol{A}^\top \boldsymbol{A}$ and $\boldsymbol{A}\boldsymbol{A}^\top$ can be diagonalised:

$$A^\top A = \sum_{i=1}^{r} \lambda_i \boldsymbol{v}_i \boldsymbol{v}_i^\top \qquad AA^\top = \sum_{i=1}^{r} \lambda_i \boldsymbol{u}_i \boldsymbol{u}_i^\top$$

Where:
$$r = \mathrm{rank}(\boldsymbol{A}^\top A) = \mathrm{rank}(\boldsymbol{A}\boldsymbol{A}^\top)$$

$\boldsymbol{u}_i \in \mathbb{R}^n, \boldsymbol{v}_i \in \mathbb{R}^d$ are orthonormal vectors.

$\lambda_i \geq 0$

# Singular value decomposition

Note that for any real matrix $A \in \mathbb{R}^{n \times d}$, $A^\top A \in \mathbb{R}^{d \times d}$ and $AA^\top \in \mathbb{R}^{n \times n}$ are a symmetric matrices.

Therefore, $A^\top A$ and $AA^\top$ can be diagonalised:

$$A^\top A = \sum_{i=1}^{r} \lambda_i v_i v_i^\top \qquad AA^\top = \sum_{i=1}^{r} \lambda_i u_i u_i^\top$$

Where:  $r = \mathrm{rank}(A^\top A) = \mathrm{rank}(AA^\top)$

$u_i \in \mathbb{R}^n, v_i \in \mathbb{R}^d$ are orthonormal vectors.

$\lambda_i \geq 0$

Therefore, defining the singular values $\sigma_i = \sqrt{\lambda_i}$

# Singular value decomposition

**Theorem**

Any real matrix $A \in \mathbb{R}^{n \times d}$ can be decomposed as

$$A = \sum_{i=1}^{\text{rank}(A)} \sigma_i u_i v_i^{\top}$$

# Singular value decomposition

Any real matrix $A \in \mathbb{R}^{n \times d}$ can be decomposed as

$$A = \sum_{i=1}^{\text{rank}(A)} \sigma_i u_i v_i^\top$$

This can be equivalently written as:

$$A = UDV^\top$$

With:  $U \in \mathbb{R}^{n \times n}$ and $V \in \mathbb{R}^{d \times d}$ orthogonal matrices

$D \in \mathbb{R}^{n \times d}$ a rectangular matrix with the singular values $\sigma_i$

# Singular value decomposition

**Theorem**

Any real matrix $A \in \mathbb{R}^{n \times d}$ can be decomposed as

$$A = \sum_{i=1}^{\text{rank}(A)} \sigma_i u_i v_i^\top$$

This can be equivalently written as:

$$A = UDV^\top$$

With: $U \in \mathbb{R}^{n \times n}$ and $V \in \mathbb{R}^{d \times d}$ orthogonal matrices

$D \in \mathbb{R}^{n \times d}$ a rectangular matrix with the singular values $\sigma_i$

⚠️ Computationally, it is more efficient to define
$U \in \mathbb{R}^{n \times r}$, $V \in \mathbb{R}^{d \times r}$ and $D \in \mathbb{R}^{r \times r}$

# Pseudo-inverse

The SVD allow us to define a generalised notion of matrix inverse. Let $A \in \mathbb{R}^{n \times d}$ with SVD:

$$A = \sum_{i=1}^{\text{rank}(A)} \sigma_i u_i v_i^\top$$

The pseudo-inverse $A^+ \in \mathbb{R}^{d \times n}$ is defined via its SVD:

$$A^+ = \sum_{i=1}^{\text{rank}(A)} \frac{1}{\sigma_i} v_i u_i^\top$$

# Pseudo-inverse

The pseudo-inverse $A^+ \in \mathbb{R}^{d \times n}$ is defined via its SVD:

$$A^+ = \sum_{i=1}^{\text{rank}(A)} \frac{1}{\sigma_i} v_i u_i^\top$$

It satisfies: $\quad AA^+A = A \qquad A^+AA^+ = A^+$

$$(A^+)^+ = A$$

If $A$ is invertible, $A^+ = A^{-1}$

If $A$ is full-rank, $\quad A^+ = \begin{cases} (A^\top A)^{-1}A^\top & \text{if } n \geq d \\ A^\top(AA^\top)^{-1}A^\top & \text{if } n < d \end{cases}$

🧐 <u>Exercise</u>: show this.

# Pseudo-inverse

The pseudo-inverse is useful to define orthogonal projectors

For any real matrix $A \in \mathbb{R}^{n \times d}$:

$$A^{+}A \in \mathbb{R}^{d \times d} \qquad\qquad AA^{+} \in \mathbb{R}^{n \times n}$$

🧐 Exercise: show this.

Define orthogonal projection operators in the column and row space of $A$, respectively.

# Pseudo-inverse

The pseudo-inverse is useful to define orthogonal projectors

For any real matrix $A \in \mathbb{R}^{n \times d}$:

$$A^+ A \in \mathbb{R}^{d \times d} \qquad AA^+ \in \mathbb{R}^{n \times n}$$

🧐 Exercise: show this.

Define orthogonal projection operators in the column and row space of $A$, respectively.

Similarly,
$$I_d - A^+ A \in \mathbb{R}^{d \times d} \qquad I_n - AA^+ \in \mathbb{R}^{n \times n}$$

Define orthogonal projection operators in the kernel of $A$ and $A^\top$, respectively.

# Recap of Probability

The butter of statistical learning

# Random variable

A random variable $X$ mathematically formalises the notion of a "measurement" or "random event".

# Random variable

A random variable $X$ mathematically formalises the notion of a "measurement" or "random event". It can be:

- Discrete: when the possible outcomes are countable.

Examples:
- the outcome of tossing a coin $X \in \{\text{head}, \text{tail}\}$
- rolling a dice $X \in \{1,\ldots,6\}$
- The number of people in France $X \in \mathbb{N}$

# Random variable

A random variable $X$ mathematically formalises the notion of a "measurement" or "random event". It can be:

- Discrete: when the possible outcomes are countable.

Examples:
- the outcome of tossing a coin $X \in \{\mathrm{head}, \mathrm{tail}\}$
- rolling a dice $X \in \{1, \ldots, 6\}$
- The number of people in France $X \in \mathbb{N}$

Discrete r.v.s are described by their probability distribution

$$\mathbb{P}(X = k)$$

A positive "function" that sums to one. $\sum_{k \in \mathrm{supp}(X)} \mathbb{P}(X = k) = 1$

# Random variable

A random variable $X$ mathematically formalises the notion of a "measurement" or "random event". It can be:

- Continuous: when the possible outcomes are uncountable.

# Random variable

A random variable $X$ mathematically formalises the notion of a "measurement" or "random event". It can be:

- Continuous: when the possible outcomes are uncountable.

Examples:
- The temperature in the room $X \in \mathbb{R}$
- The GDP of France next year $X \in \mathbb{R}$

# Random variable

A random variable $X$ mathematically formalises the notion of a "measurement" or "random event". It can be:

- Continuous: when the possible outcomes are uncountable.

Examples:
- The temperature in the room $X \in \mathbb{R}$
- The GDP of France next year $X \in \mathbb{R}$

Continuous r.v.s are described by their probability density function (p.d.f.), which integrates to probabilities:

$$\mathbb{P}(X \in [a, b]) = \int_a^b \mathrm{d}x \; p_X(x)$$

A "function" that integrates to one: $\int_{\mathrm{supp}(X)} \mathrm{d}x \; p_X(x) = 1$

# Random variable

A random variable $X$ mathematically formalises the notion of a "measurement" or "random event". It can be:

- Continuous: when the possible outcomes are uncountable.

Examples:
- The temperature in the room $X \in \mathbb{R}$
- The GDP of France next year $X \in \mathbb{R}$

Continuous r.v.s are described by their probability density function (p.d.f.), which integrates to probabilities:

$$\mathbb{P}(X \in [a, b]) = \int_a^b \mathrm{d}x \; p_X(x)$$

A "function" that integrates to one: $\int_{\mathrm{supp}(X)} \mathrm{d}x \; p_X(x) = 1$

⚠️ The p.d.f. is NOT a probability. It can be negative.

# Normal distribution

A Gaussian r.v. $X \sim \mathcal{N}(\mu, \sigma^2)$ has the following p.d.f.:

$$p_X(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

# Normal distribution

A Gaussian r.v. $X \sim \mathcal{N}(\mu, \sigma^2)$ has the following p.d.f.:

$$p_X(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

High-probability

Low-probability



One-Dimensional Gaussian PDF (Mean=0, Variance=1)

# Expectation and variance

Let $X \sim p_X$ denote a continuous r.v.

- The expectation (or mean) of $X$ is defined as

$$\mathbb{E}[X] = \int \mathrm{d}x \; p_X(x)x$$

For example, for $X \sim \mathcal{N}(\mu, \sigma^2)$, we have $\mathbb{E}[X] = \mu$

# Expectation and variance

Let $X \sim p_X$ denote a continuous r.v.

- The expectation (or mean) of $X$ is defined as

$$\mathbb{E}[X] = \int \mathrm{d}x \; p_X(x)x$$

For example, for $X \sim \mathcal{N}(\mu, \sigma^2)$, we have $\mathbb{E}[X] = \mu$

- The variance of $X$ is defined as:

$$\mathrm{Var}[X] = \mathbb{E}[(X - \mathbb{E}[X])^2] = \mathbb{E}[X^2] - \mathbb{E}[X]^2$$

For example, for $X \sim \mathcal{N}(\mu, \sigma^2)$, we have $\mathrm{Var}[X] = \sigma^2$

# Change of variables

Let $X \sim p_X$ denote a continuous r.v. and $f : \mathbb{R} \to \mathbb{R}$

# Change of variables

Let $X \sim p_X$ denote a continuous r.v. and $f : \mathbb{R} \to \mathbb{R}$

Then, $Y = f(X)$ is also a random variable, with p.d.f. given by

$$p_Y(y) = \int \mathrm{d}x \; p_X(x)\delta(y - f(x))$$

Where $\delta(x)$ is the "Dirac delta function":

$$\int_{\mathbb{R}} \mathrm{d}x \; \delta(x - y)f(x) = f(y)$$

# Joint distribution

Two random variables $X, Y$ that concern the same random experiment are characterised by their joint p.d.f.

$$p_{X,Y}(x, y)$$

# Joint distribution

Two random variables $X, Y$ that concern the same random experiment are characterised by their joint p.d.f.

$$p_{X,Y}(x, y)$$

The correlation between $X, Y$ is defined by

$$\mathbb{E}[XY] = \int \mathrm{d}x \int \mathrm{d}y \; p_{X,Y}(x, y) xy$$

# Joint distribution

Two random variables $X, Y$ that concern the same random experiment are characterised by their joint p.d.f.

$$p_{X,Y}(x, y)$$

The correlation between $X, Y$ is defined by

$$\mathbb{E}[XY] = \int \mathrm{d}x \int \mathrm{d}y \; p_{X,Y}(x, y)xy$$

We say $X, Y$ are uncorrelated if $\mathbb{E}[XY] = \mathbb{E}[X]\mathbb{E}[Y]$

# Independence

- Given two r.v.s $X, Y \sim p_{X,Y}$, we define the marginal distributions

$$p_X(x) = \int dy \; p_{X,Y}(x, y) \qquad p_Y(y) = \int dx \; p_{X,Y}(x, y)$$

# Independence

- Given two r.v.s $X, Y \sim p_{X,Y}$, we define the marginal distributions

$$p_X(x) = \int dy \; p_{X,Y}(x, y) \qquad p_Y(y) = \int dx \; p_{X,Y}(x, y)$$

- We say the r.v.s. $X, Y$ are independent if

$$p_{X,Y}(x, y) = p_X(x) p_Y(x)$$

# Independence

- Given two r.v.s $X, Y \sim p_{X,Y}$, we define the marginal distributions

$$p_X(x) = \int \mathrm{d}y \; p_{X,Y}(x, y) \qquad p_Y(y) = \int \mathrm{d}x \; p_{X,Y}(x, y)$$

- We say the r.v.s. $X, Y$ are independent if

$$p_{X,Y}(x, y) = p_X(x) p_Y(x)$$

⚠️ Note that independence implies uncorrelated, but not the converse!

🧐 Exercise: Construct a counter-example

# Conditional distribution

- Given two r.v.s $X, Y \sim p_{X,Y}$, we define the conditional p.d.f.

$$p_{X|Y}(x \mid y) = \frac{p_{X,Y}(x, y)}{p_Y(y)}$$

# Conditional distribution

- Given two r.v.s $X, Y \sim p_{X,Y}$, we define the conditional p.d.f.

$$p_{X|Y}(x \mid y) = \frac{p_{X,Y}(x, y)}{p_Y(y)}$$

Note that $X, Y \sim p_{X,Y}$ are independent if and only if:

$$p_{X|Y}(x \mid y) = p_X(x)$$

# Conditional distribution

- Given two r.v.s $X, Y \sim p_{X,Y}$, we define the conditional p.d.f.

$$p_{X|Y}(x|y) = \frac{p_{X,Y}(x, y)}{p_Y(y)}$$

Note that $X, Y \sim p_{X,Y}$ are independent if and only if:

$$p_{X|Y}(x|y) = p_X(x)$$

Theorem (Bayes theorem)

$$p_{X|Y}(x|y) = \frac{p_{Y|X}(y|x)p_X(x)}{p_Y(y)}$$

# Law of large numbers

Let $X_1, \ldots, X_n \sim p_X$ denote i.i.d. r.v.s. with mean $\mathbb{E}[X_i] = \mu$

Define the sample mean (note this is itself a r.v.)

$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^{n} X_i$$

# Law of large numbers

Let $X_1, \ldots, X_n \sim p_X$ denote i.i.d. r.v.s. with mean $\mathbb{E}[X_i] = \mu$

Define the sample mean (note this is itself a r.v.)

$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^{n} X_i$$

Theorem (Weak LLN)

$$\bar{X}_n \xrightarrow{P} \mu \qquad \text{as} \qquad n \to \infty$$

$$\lim_{n\to\infty} \mathbb{P}(|\bar{X}_n - \mu| < \epsilon) = 1$$

⚠️ Be aware there are many variations of the LLN.

# Central limit theorem

Let $X_1, \ldots, X_n \sim p_X$ denote i.i.d. r.v.s. with mean $\mathbb{E}[X_i] = \mu$ and variance $\mathrm{Var}(X_i) = \sigma^2 < \infty$

Again, consider the sample mean

$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^{n} X_i$$

# Central limit theorem

Let $X_1, \dots, X_n \sim p_X$ denote i.i.d. r.v.s. with mean $\mathbb{E}[X_i] = \mu$ and variance $\mathrm{Var}(X_i) = \sigma^2 < \infty$

Again, consider the sample mean

$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^{n} X_i$$

**Theorem (Lindeberg CLT)**

$$\sqrt{n}(\bar{X}_n - \mu) \xrightarrow{d} \mathcal{N}(\mu, \sigma^2)$$

$$\lim_{n \to \infty} \mathbb{P}(\sqrt{n}(\bar{X}_n - \mu) \leq z) = \mathbb{P}(Z \leq z/\sigma) \qquad Z \sim \mathcal{N}(0,1)$$

⚠ Be aware there are many variations of the CLT.