



Statistical Learning II

Lecture 1 - Introduction & preliminaries

Bruno Loureiro
@ CSD, DI-ENS & CNRS

brloureiro@gmail.com

Course Organisation

- 12 classes of 3h, divided in:
 - 1h30 lectures
 - 1h30 exercises & lab (Python)
(with **Leonardo De Filippis**)



Course Organisation

- 12 classes of 3h, divided in:
 - 1h30 lectures
 - 1h30 exercises & lab (Python)
(with **Leonardo De Filippis**)



- Contact us: via **Microsoft Teams** or **e-mail**:

bruno.loureiro@di.ens.fr and leonardo.defilippis99@gmail.com

Course Organisation

- 12 classes of 3h, divided in:
 - 1h30 lectures
 - 1h30 exercises & lab (Python)
(with **Leonardo De Filippis**)



- Contact us: via **Microsoft Teams** or **e-mail**:
bruno.loureiro@di.ens.fr and leonardo.defilippis99@gmail.com
- Evaluation: 2 x exams (**midterm** and **final**)

Menu for the semester

Goal: Develop a *mathematical* understanding of *classical* and *modern* machine learning models

Menu for the semester

Goal: Develop a *mathematical* understanding of *classical* and *modern* machine learning models

- Classical methods:
 - Ridge regression
 - LASSO
 - Generalised linear models
 - Kernel methods
 - Principal component analysis (PCA)

Menu for the semester

Goal: Develop a *mathematical* understanding of *classical* and *modern* machine learning models

- Classical methods:
 - Ridge regression
 - LASSO
 - Generalised linear models
 - Kernel methods
 - Principal component analysis (PCA)
- Modern methods:
 - Neural networks
 - Diffusion models
 - Your suggestions?

Introduction & Motivation

Or: why should I care about Statistical Learning?

A brief history of ML

W. McCulloch W. Pitts



McCulloch-Pitts
neuron
1943

F. Rosenblat



Rosenblat's
perceptron /
MLP
1957

S. Linnainmaa



Backprop
1970

K. Fukushima



Neocognitron
1980

Y. LeCun



CNNs
1989



2016
AlphaGo



ChatGPT
2022

1940

1st digital
computer

1947

Invention of the
transistor

1969

ARPANET
(Internet)

1982

Hopfield
Model

1997

IBM Deep
Blue

2017

Transformers

DALL-E
2021

2020

AlphaFold2



J. Bardeen W.H. Brattain W. Shockley



J. Hopfield



What about the maths?

Yet, on the mathematical side...

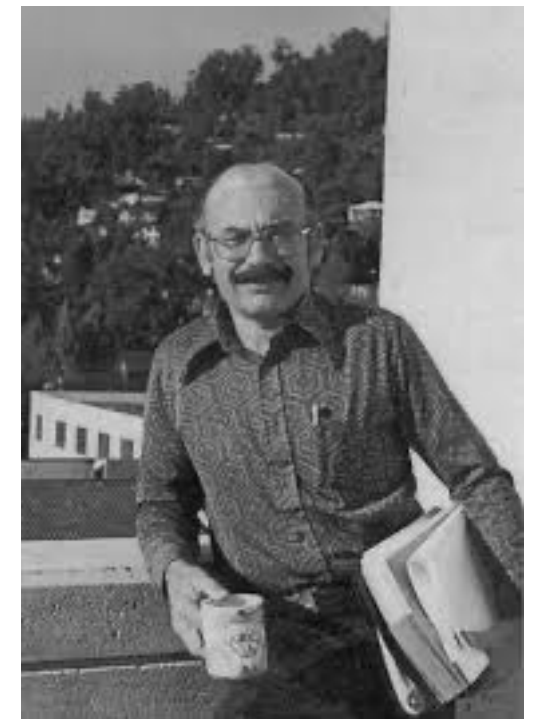
Leo Breiman

Statistics Department, University of California, Berkeley, CA 94305;
e-mail: leo@stat.berkeley.edu

Reflections After Refereeing Papers for NIPS

For instance, there are many important questions regarding neural networks which are largely unanswered. There seem to be conflicting stories regarding the following issues:

- Why don't heavily parameterized neural networks overfit the data?
- What is the effective number of parameters?
- Why doesn't backpropagation head for a poor local minima?
- When should one stop the backpropagation and use the current parameters?



Leo Breiman
1928

What about the maths?

Yet, on the mathematical side...

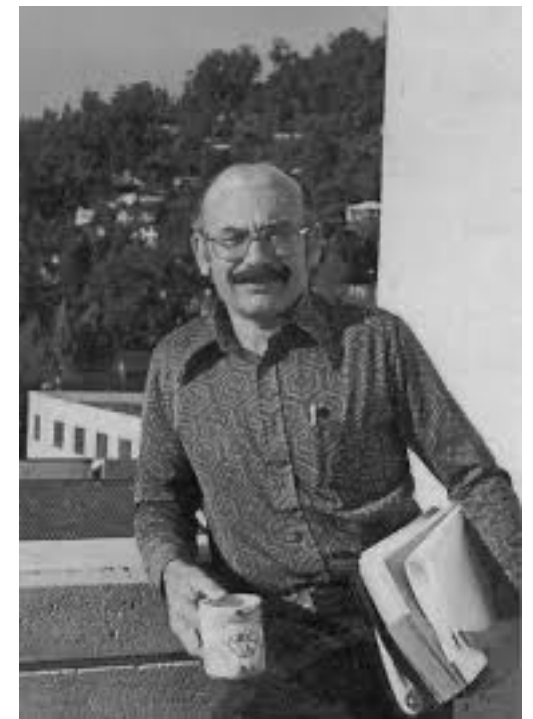
Leo Breiman

Statistics Department, University of California, Berkeley, CA 94305;
e-mail: leo@stat.berkeley.edu

Reflections After Refereeing Papers for NIPS

For instance, there are many important questions regarding neural networks which are largely unanswered. There seem to be conflicting stories regarding the following issues:

- Why don't heavily parameterized neural networks overfit the data?
- What is the effective number of parameters?
- Why doesn't backpropagation head for a poor local minima?
- When should one stop the backpropagation and use the current parameters?



Leo Breiman
1928

This was written in 1995!!!

What about the maths?

Yet, on the mathematical side...

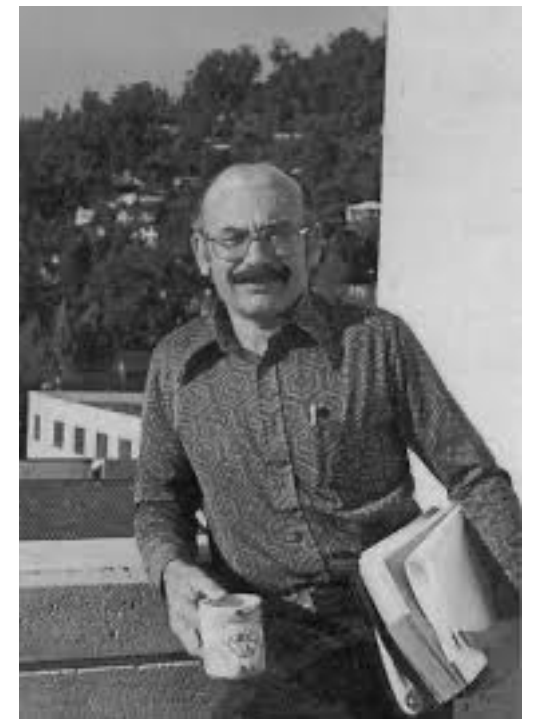
Leo Breiman

Statistics Department, University of California, Berkeley, CA 94305;
e-mail: leo@stat.berkeley.edu

Reflections After Refereeing Papers for NIPS

For instance, there are many important questions regarding neural networks which are largely unanswered. There seem to be conflicting stories regarding the following issues:

- Why don't heavily parameterized neural networks overfit the data?
- What is the effective number of parameters?
- Why doesn't backpropagation head for a poor local minima?
- When should one stop the backpropagation and use the current parameters?



Leo Breiman
1928

This was written in 1995!!!

But why should I care?

Some good reasons to care

- Reliability and Liability

If a model does something unexpected, **who is responsible?**

Crucial in **sensitive applications**, e.g. medicine, law, self-driving cars/planes...

Some good reasons to care

- Reliability and Liability

If a model does something unexpected, **who is responsible?**

Crucial in **sensitive applications**, e.g. medicine, law, self-driving cars/planes...

- Efficient design

Data centres are responsible for **4% of the energy consumption in the US.**

Can we **design models** and **algorithms** that learn more **efficiently** from data?

Some good reasons to care

- Reliability and Liability

If a model does something unexpected, **who is responsible?**

Crucial in **sensitive applications**, e.g. medicine, law, self-driving cars/planes...

- Efficient design

Data centres are responsible for **4% of the energy consumption in the US.**

Can we **design models** and **algorithms** that learn more **efficiently** from data?

- Scientific curiosity

Some good reasons to care

- Reliability and Liability

If a model does something unexpected, **who is responsible?**

Crucial in **sensitive applications**, e.g. medicine, law, self-driving cars/planes...

- Efficient design

Data centres are responsible for **4% of the energy consumption in the US.**

Can we **design models** and **algorithms** that learn more **efficiently** from data?

- Scientific curiosity

- And in the worst case, understanding the maths will make you a **better engineer / data scientist.**

Our expectations

My expectations: By the end of the term, I expect you to:

- Get acquainted with the most popular machine learning algorithms

Our expectations

My expectations: By the end of the term, I expect you to:

- Get acquainted with the most popular machine learning algorithms
- Appreciate (some) of the mathematics behind the methods.

Our expectations

My expectations: By the end of the term, I expect you to:

- Get acquainted with the most popular machine learning algorithms
- Appreciate (some) of the mathematics behind the methods.
- Be able to implement these methods from scratch.

99% of books in statistical learning:

*Let $(\mathbf{x}_i, y_i) \in \mathbb{R}^d \times \mathbb{R}, i = 1, \dots, n$,
denote independently drawn samples from a
probability distribution....*

99% of books in statistical learning:

Let $(\mathbf{x}_i, y_i) \in \mathbb{R}^d \times \mathbb{R}, i = 1, \dots, n$,
denote *independently drawn samples from a*
probability distribution....

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:

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

“column vector”

The Euclidean space

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


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

“column vector”

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

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

Position i



The Euclidean space

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

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

The Euclidean space

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

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

Which induces a natural notion of **distance** and **size**:

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

“Euclidean or ℓ_2 norm” “Euclidean distance”

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

Euclidean geometry

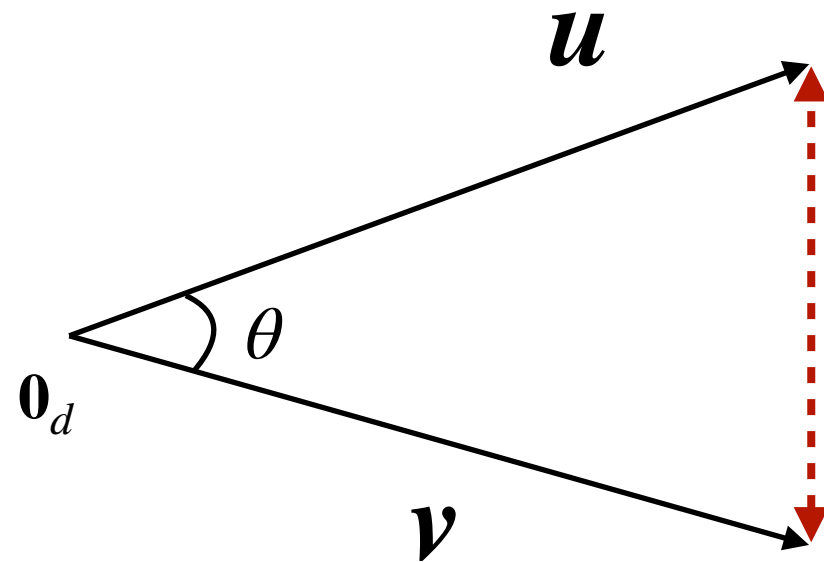
$$||\mathbf{u}||_2^2 = \langle \mathbf{u}, \mathbf{u} \rangle = \sum_{i=1}^d u_i^2$$

“Euclidean or ℓ_2 norm”

$$d(\mathbf{u}, \mathbf{v}) = ||\mathbf{u} - \mathbf{v}||_2^2$$

“Euclidean distance”

They correspond to our intuitive notion of geometry in the plane

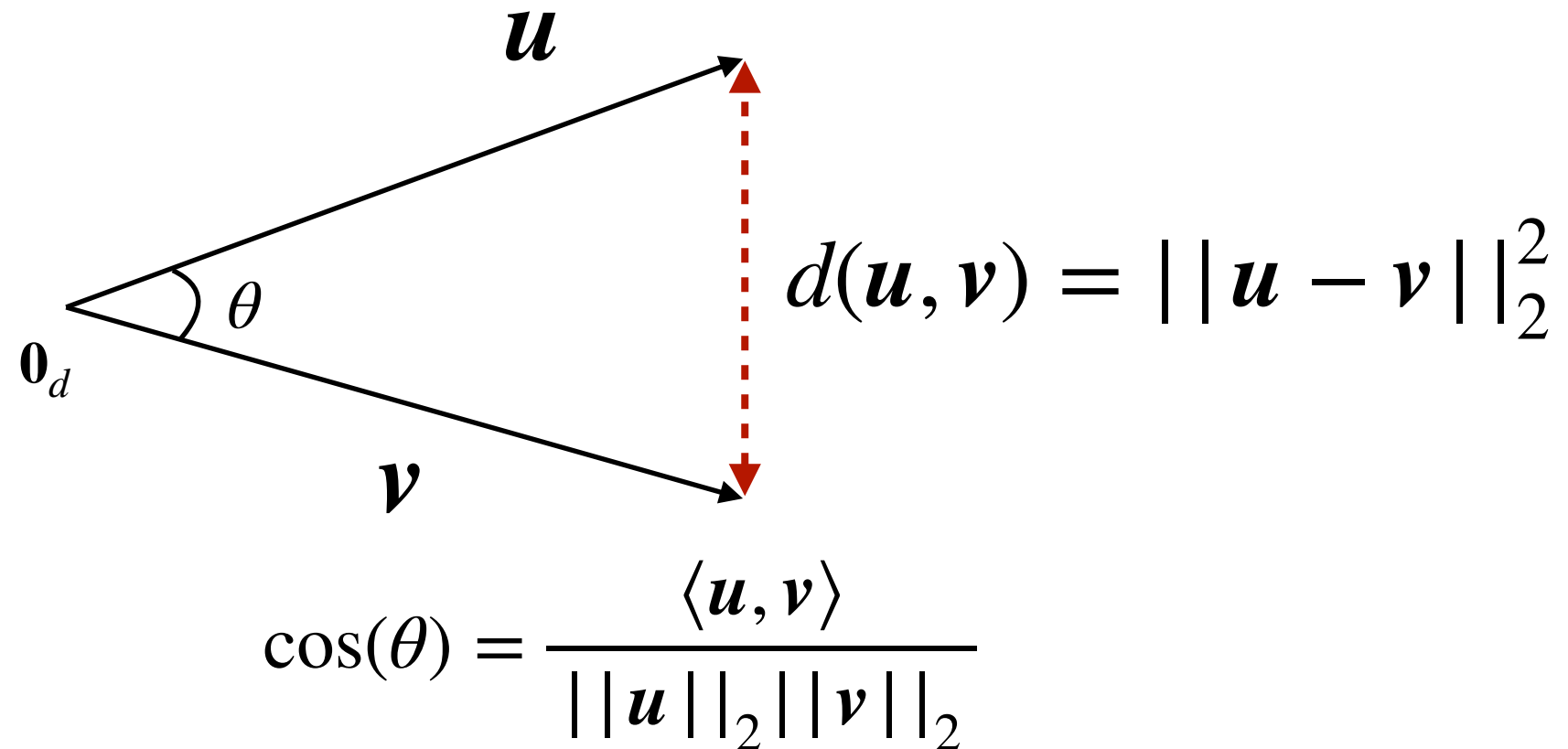


$$d(\mathbf{u}, \mathbf{v}) = ||\mathbf{u} - \mathbf{v}||_2^2$$

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

Euclidean geometry

They correspond to our intuitive notion of geometry in the plane



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

$$\langle \mathbf{u}, \mathbf{v} \rangle = 0$$



Other norms

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

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

" ℓ_p norm"

Other norms

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

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

" ℓ_p norm"



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

Other norms

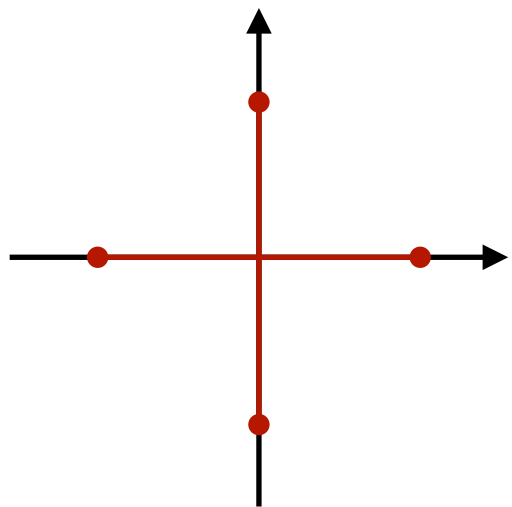
One can define other notions of size in \mathbb{R}^d

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

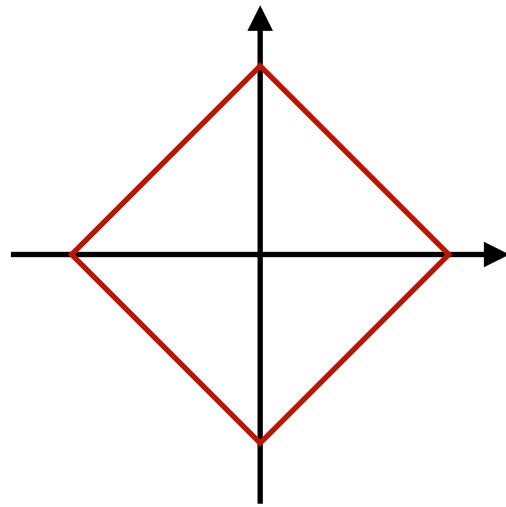
" ℓ_p norm"



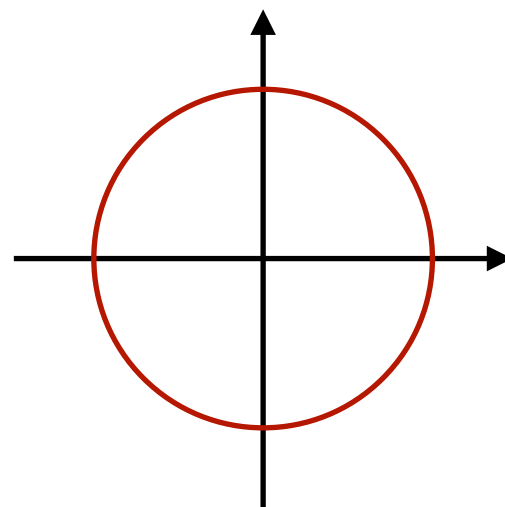
$||\cdot||_p$ is not associated to an inner product for $p \neq 2$



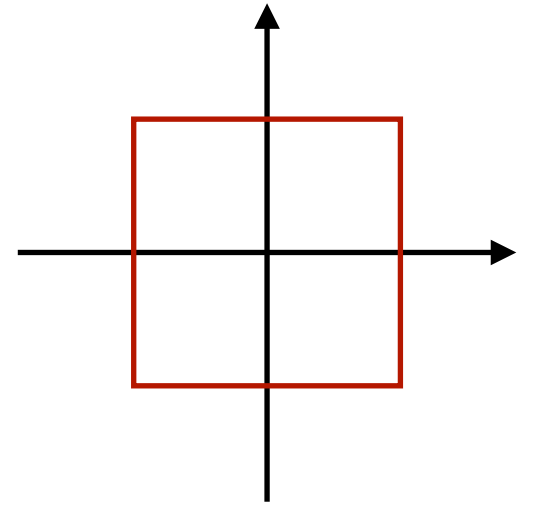
ℓ_0



ℓ_1



ℓ_2



ℓ_∞



Not a norm

Matrices

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

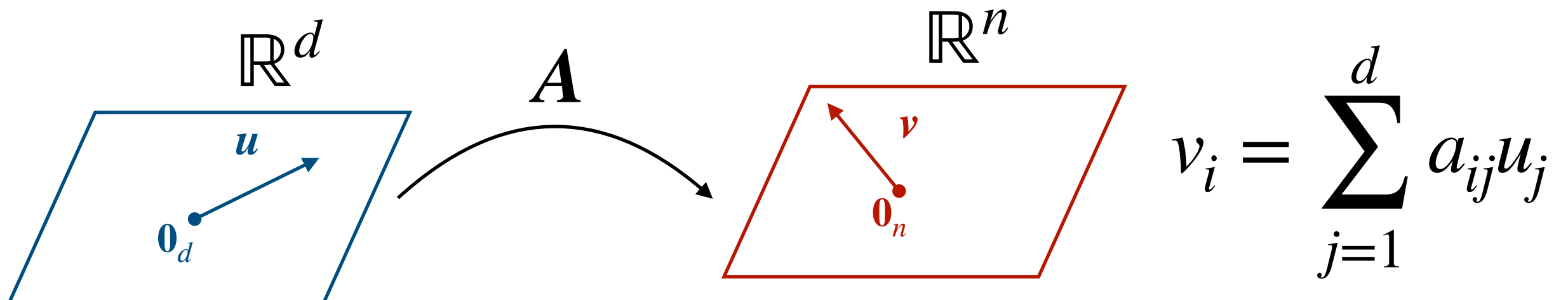
$$\mathbf{A} = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nd} \end{bmatrix} \in \mathbb{R}^{n \times d}$$

Matrices

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

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

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



Matrices

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

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

$\mathbf{A}_1 \quad \mathbf{A}_2 \quad \cdots \quad \mathbf{A}_d$

- The columns of $\mathbf{A} \in \mathbb{R}^{n \times d}$ are vectors $\mathbf{A}_i \in \mathbb{R}^n$ with $(\mathbf{A}_i)_j = a_{ij}$
- “Column space” $\text{col}(\mathbf{A}) = \text{span}(\mathbf{A}_1, \cdots, \mathbf{A}_d) \subset \mathbb{R}^n$

Matrices

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

$$\mathbf{A} = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nd} \end{bmatrix} \begin{matrix} \mathbf{a}_1 \\ \mathbf{a}_2 \\ \\ \mathbf{a}_n \end{matrix} \in \mathbb{R}^{n \times d}$$

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

“Column space” $\text{col}(\mathbf{A}) = \text{span}(\mathbf{A}_1, \dots, \mathbf{A}_d) \subset \mathbb{R}^n$

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

“Row space” of $\text{row}(\mathbf{A}) = \text{span}(\mathbf{a}_1, \dots, \mathbf{a}_n) \subset \mathbb{R}^d$

Flattening matrices

The space of matrices $\mathbf{A} \in \mathbb{R}^{n \times d}$ 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.

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

Rank of a matrix

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

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

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

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

This is equivalent to the **number of independent columns**.

Proposition

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

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

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

Another point of view

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

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

Another point of view

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

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

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

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

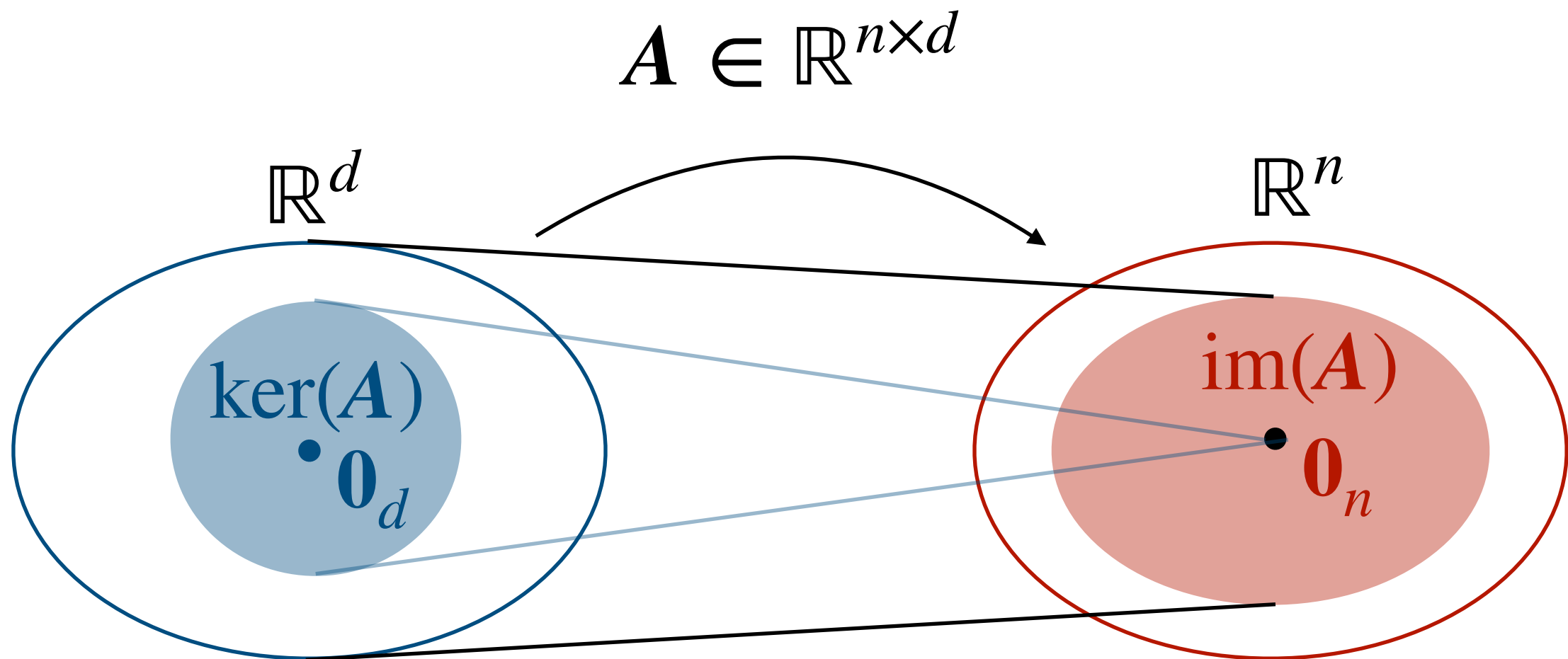


Note that $\ker(\mathbf{A}) \subset \mathbb{R}^d$
and $\mathbf{0} \in \ker(\mathbf{A})$

Image and null-space

Proposition

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



Matrix inverse

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

$$\mathbf{AB} = \mathbf{I}_d$$

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

Matrix inverse

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

$$\mathbf{AB} = \mathbf{I}_d$$

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



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

Matrix inverse

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

$$\mathbf{AB} = \mathbf{I}_d$$

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



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

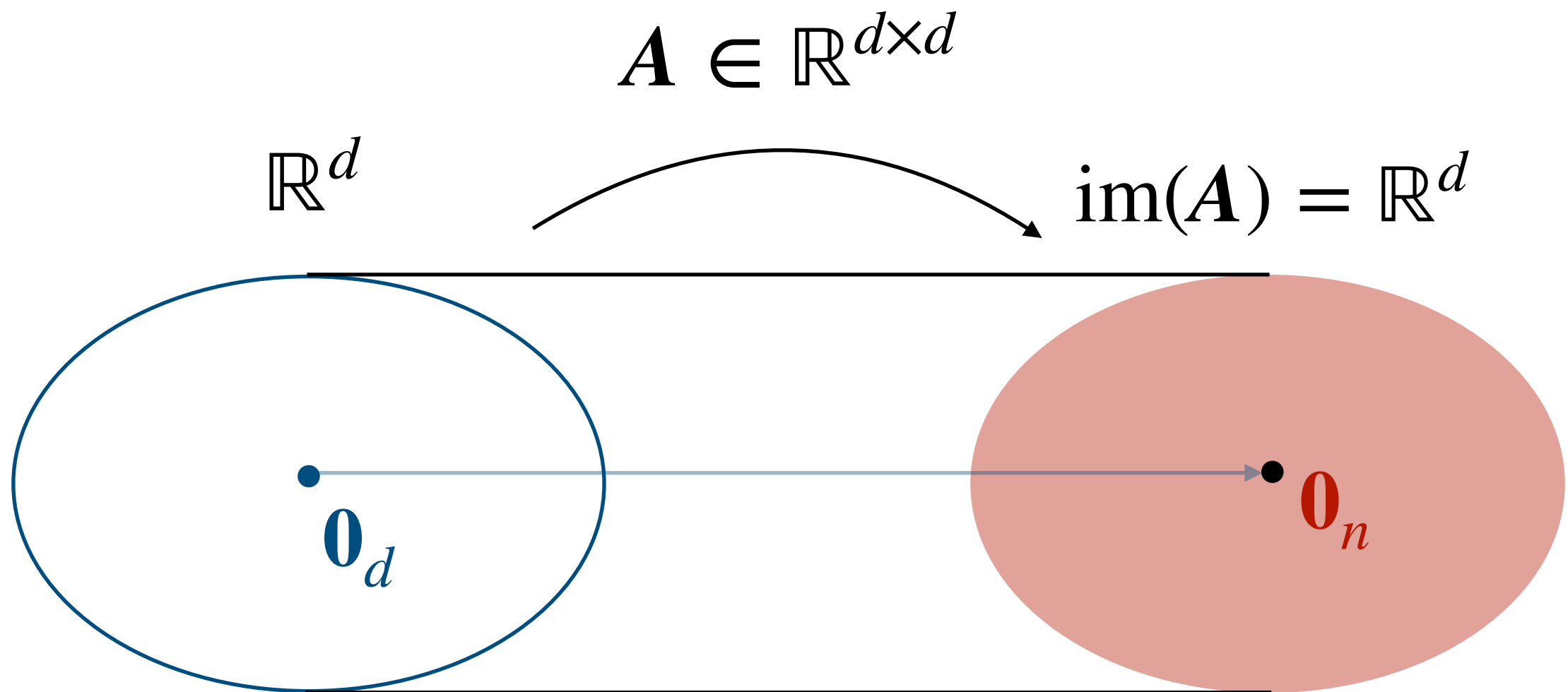
Proposition

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

Matrix inverse

Proposition

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



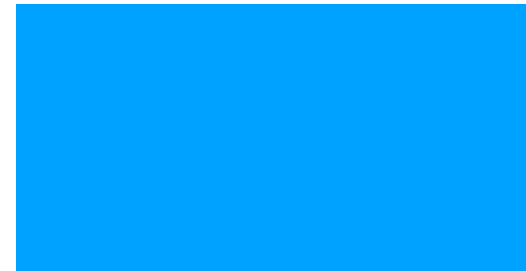
Matrix transpose

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

$$\mathbf{A} =$$



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



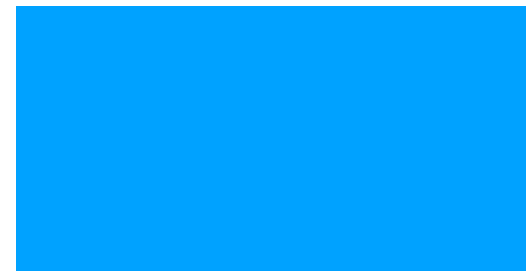
Matrix transpose

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

$$\mathbf{A} =$$



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



- We have:

$$(\mathbf{A}^\top)^\top = \mathbf{A}$$

$$(a\mathbf{A} + b\mathbf{B})^\top = a\mathbf{A}^\top + b\mathbf{B}^\top$$

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

$$(\mathbf{AB})^\top = \mathbf{B}^\top \mathbf{A}^\top$$



Exercise: check this.

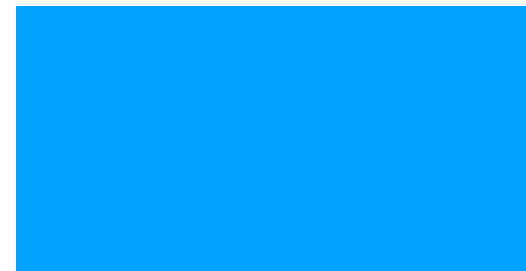
Matrix transpose

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

$$\mathbf{A} =$$



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



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

$$\langle \mathbf{u}, \mathbf{v} \rangle = \mathbf{u}^\top \mathbf{v}$$



Exercise: check this.

Matrix trace

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

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

Matrix trace

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

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

- It satisfies: $\text{Tr } \mathbf{AB} = \text{Tr } \mathbf{BA}$

$$\text{Tr } (a\mathbf{A} + b\mathbf{B}) = a\text{Tr } \mathbf{A} + b\text{Tr } \mathbf{B}$$

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



Exercise: check this.

Symmetric matrices

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

Symmetric matrices

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



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

Symmetric matrices

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



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

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

$$(\mathbf{A} \mathbf{A}^\top)_{ij} = \langle \mathbf{a}_i, \mathbf{a}_j \rangle$$



Exercise: check this.

Note: a similar representation holds for columns of \mathbf{A}

Orthogonal matrices

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

Orthogonal matrices

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

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

$$||\mathbf{A}\mathbf{u}||_2 = ||\mathbf{u}||_2$$



Exercise: check this.

Orthogonal matrices

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

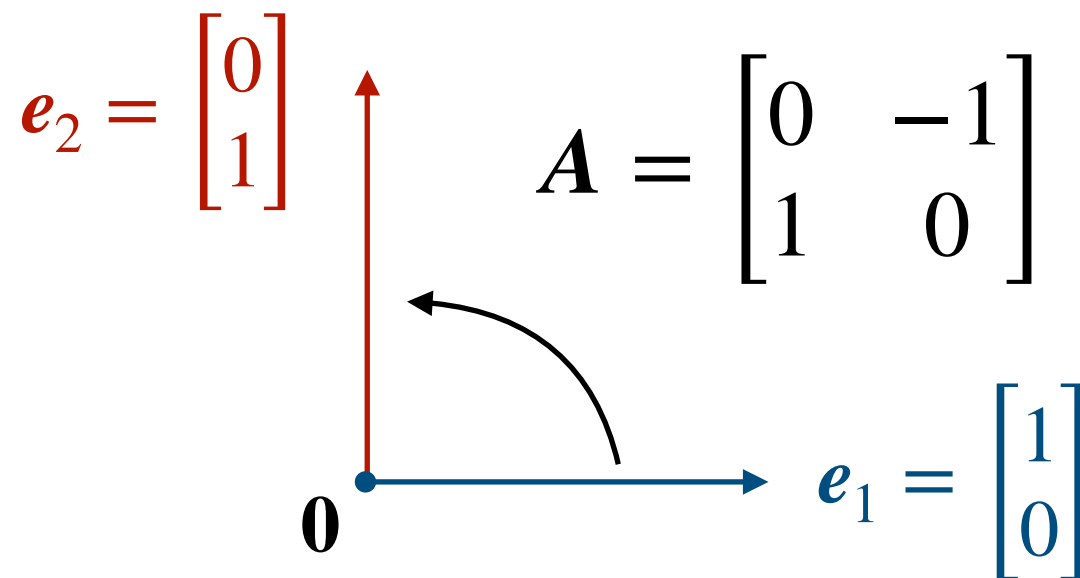
Orthogonal matrices preserve the norm and distance between vectors (they are **isometries**):

$$||\mathbf{A}\mathbf{u}||_2 = ||\mathbf{u}||_2$$



Exercise: check this.

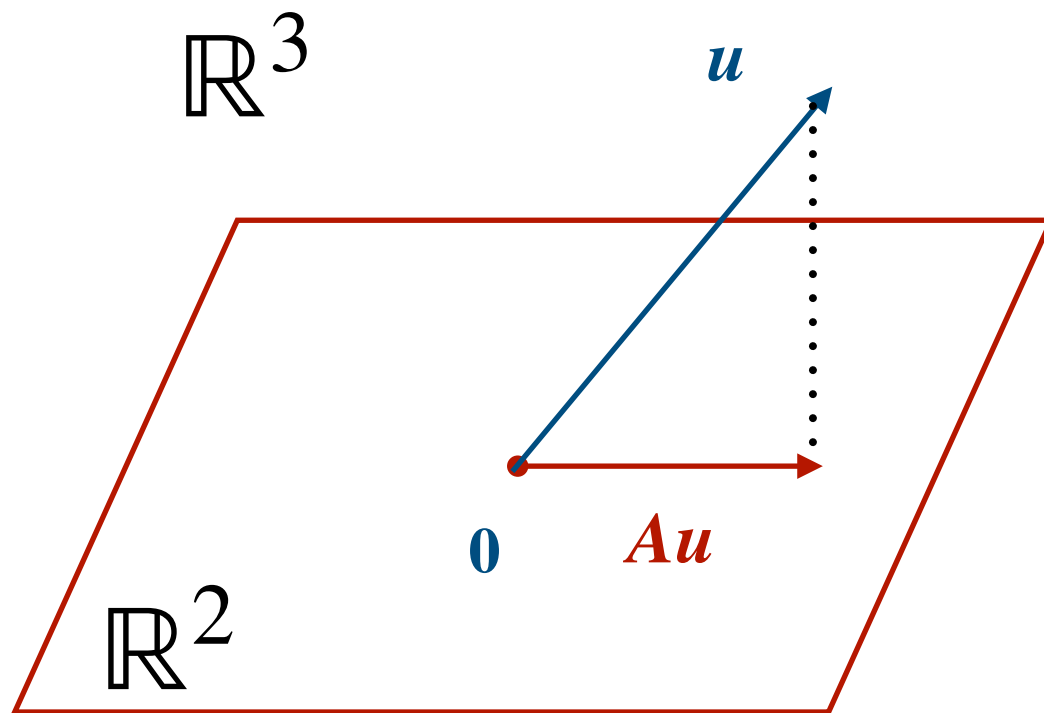
Geometrically, they define **rotations**



Projection matrix

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

Moreover, if \mathbf{A} is also orthogonal, we call it a **orthogonal projection**.



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

Projection matrix

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

Moreover, if \mathbf{A} is also orthogonal, we call it a **orthogonal projection**.

Proposition

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

$$\mathbf{v} = \mathbf{u} + \mathbf{A}\mathbf{v} \qquad \mathbf{u} \in \ker(\mathbf{A})$$

Projection matrix

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

Moreover, if \mathbf{A} is also orthogonal, we call it a **orthogonal projection**.

Proposition

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

$$\mathbf{v} = \mathbf{u} + \mathbf{A}\mathbf{v} \qquad \mathbf{u} \in \ker(\mathbf{A})$$



The only projection matrix which is invertible is the identity.

Eigen-(values, vectors)

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

$$\mathbf{A}\mathbf{v} = \lambda\mathbf{v}$$

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

Eigen-(values, vectors)

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

$$\mathbf{A}\mathbf{v} = \lambda\mathbf{v}$$

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

We call the set of eigenvalues the spectrum of \mathbf{A} :

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

Eigen-(values, vectors)

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

$$\mathbf{A}\mathbf{v} = \lambda\mathbf{v}$$

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

We call the set of eigenvalues the spectrum of \mathbf{A} :

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

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



- An eigenvalue λ can be associated to more than one independent eigenvector.

Positive matrices

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

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

Positive matrices

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

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

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

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

Positive matrices

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

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

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

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

Proposition

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



Exercise: prove this.

Positive matrices

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

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

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

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

Proposition

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



not necessarily positive definite.



Exercise: prove this.

Spectral theorem

Theorem

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

$$\mathbf{A} = \mathbf{U} \mathbf{D} \mathbf{U}^\top$$

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

Spectral theorem

Theorem

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

$$\mathbf{A} = \mathbf{U} \mathbf{D} \mathbf{U}^\top$$

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

We can equivalently write the spectral decomposition as:

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

Where $\mathbf{v}_i \in \mathbb{R}^d$ are orthonormal eigenvectors.

Important facts

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

$$\text{Tr } \mathbf{A} = \sum_{i=1}^d \lambda_i$$

Important facts

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

$$\text{Tr } \mathbf{A} = \sum_{i=1}^d \lambda_i$$

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

Important facts

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

$$\text{Tr } \mathbf{A} = \sum_{i=1}^d \lambda_i$$

- A square matrix $\mathbf{A} \in \mathbb{R}^{d \times d}$ is invertible i.f.f. $0 \notin \text{spec}(\mathbf{A})$
- The eigenvalues of a projection matrix $\mathbf{P} \in \mathbb{R}^{d \times d}$ are 0 or 1

$$\mathbf{P} = \sum_{i=1}^{\text{rank}(\mathbf{P})} \mathbf{v}_i \mathbf{v}_i^\top$$



Exercise: show this.

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

Singular value decomposition

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

Singular value decomposition

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

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

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

Where: $r = \text{rank}(\mathbf{A}^\top \mathbf{A}) = \text{rank}(\mathbf{A}\mathbf{A}^\top)$

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

$$\lambda_i \geq 0$$

Singular value decomposition

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

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

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

Where: $r = \text{rank}(\mathbf{A}^\top \mathbf{A}) = \text{rank}(\mathbf{A}\mathbf{A}^\top)$

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

$$\lambda_i \geq 0$$

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

Singular value decomposition

Theorem

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

$$\mathbf{A} = \sum_{i=1}^{\text{rank}(\mathbf{A})} \sigma_i \mathbf{u}_i \mathbf{v}_i^{\top}$$

Singular value decomposition

Theorem

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

$$\mathbf{A} = \sum_{i=1}^{\text{rank}(\mathbf{A})} \sigma_i \mathbf{u}_i \mathbf{v}_i^{\top}$$

This can be equivalently written as:

$$\mathbf{A} = \mathbf{U} \mathbf{D} \mathbf{V}^{\top}$$

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

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

Singular value decomposition

Theorem

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

$$\mathbf{A} = \sum_{i=1}^{\text{rank}(\mathbf{A})} \sigma_i \mathbf{u}_i \mathbf{v}_i^\top$$

This can be equivalently written as:

$$\mathbf{A} = \mathbf{U} \mathbf{D} \mathbf{V}^\top$$

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

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



Computationally, it is more efficient to define

$\mathbf{U} \in \mathbb{R}^{n \times r}$, $\mathbf{V} \in \mathbb{R}^{d \times r}$ and $\mathbf{D} \in \mathbb{R}^{r \times r}$

Pseudo-inverse

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

$$\mathbf{A} = \sum_{i=1}^{\text{rank}(\mathbf{A})} \sigma_i \mathbf{u}_i \mathbf{v}_i^\top$$

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

$$\mathbf{A}^+ = \sum_{i=1}^{\text{rank}(\mathbf{A})} \frac{1}{\sigma_i} \mathbf{v}_i \mathbf{u}_i^\top$$

Pseudo-inverse

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

$$\mathbf{A}^+ = \sum_{i=1}^{\text{rank}(\mathbf{A})} \frac{1}{\sigma_i} \mathbf{v}_i \mathbf{u}_i^\top$$

It satisfies: $\mathbf{A}\mathbf{A}^+\mathbf{A} = \mathbf{A}$ $\mathbf{A}^+\mathbf{A}\mathbf{A}^+ = \mathbf{A}^+$

$$(\mathbf{A}^+)^+ = \mathbf{A}$$

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

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



Exercise: show this.

Pseudo-inverse

The pseudo-inverse is useful to define orthogonal projectors

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

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

$$\mathbf{A} \mathbf{A}^+ \in \mathbb{R}^{n \times n}$$



Exercise:
show this.

Define orthogonal projection operators in the column and row space of \mathbf{A} , respectively.

Pseudo-inverse

The pseudo-inverse is useful to define orthogonal projectors

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

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

$$\mathbf{A} \mathbf{A}^+ \in \mathbb{R}^{n \times n}$$



Exercise:
show this.

Define orthogonal projection operators in the column and row space of \mathbf{A} , respectively.

Similarly,

$$\mathbf{I}_d - \mathbf{A}^+ \mathbf{A} \in \mathbb{R}^{d \times d}$$

$$\mathbf{I}_n - \mathbf{A} \mathbf{A}^+ \in \mathbb{R}^{n \times n}$$

Define orthogonal projection operators in the kernel of \mathbf{A} and \mathbf{A}^T , respectively.