

# Fundamentals of Mathematics for Machine Learning

Oxford ML Summer School  
8 May 2023

Rasul Tutunov  
Juliusz Ziomek

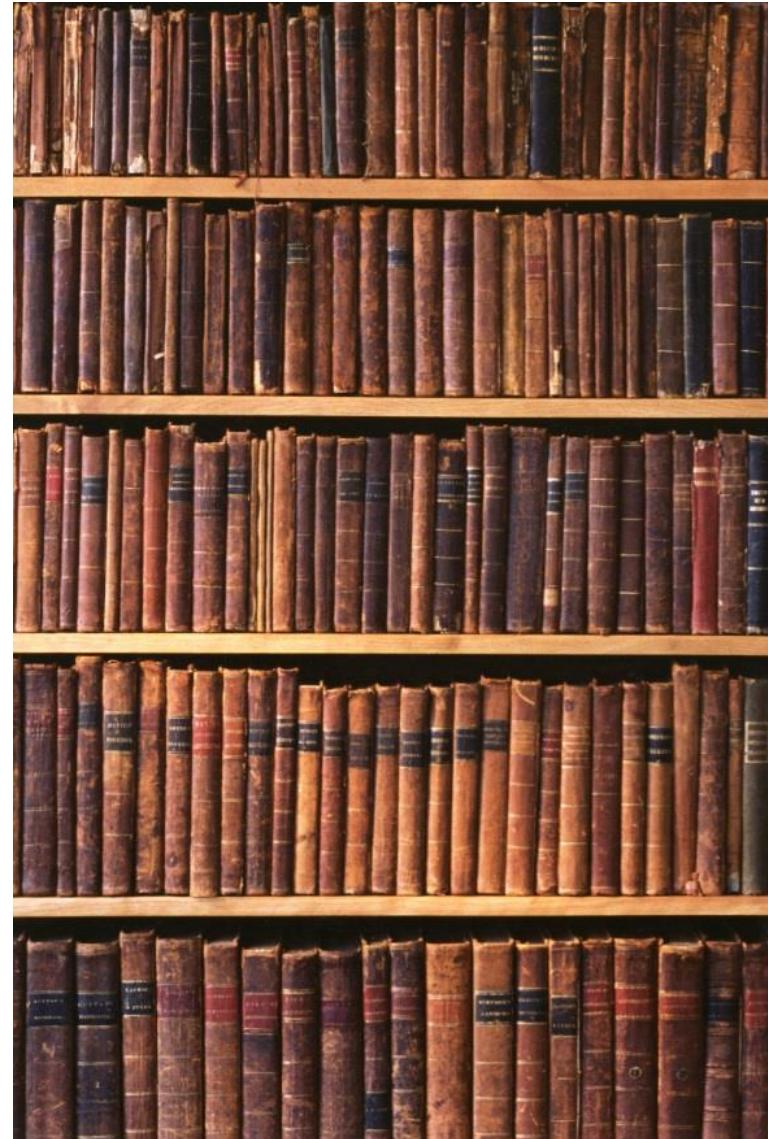
# Plan of the Talk

Linear algebra

Vector Calculus

Probability Theory

Loss functions in ML



Collection of 511 American Mathematical Textbooks, 1760–1850



Courtesy of Robin Halwas

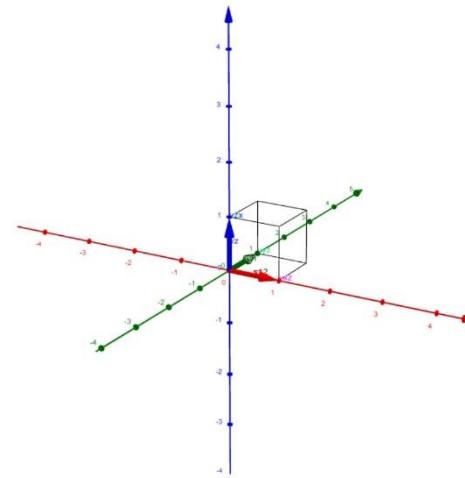
# Linear Algebra

## Vectors in Euclidian Spaces.

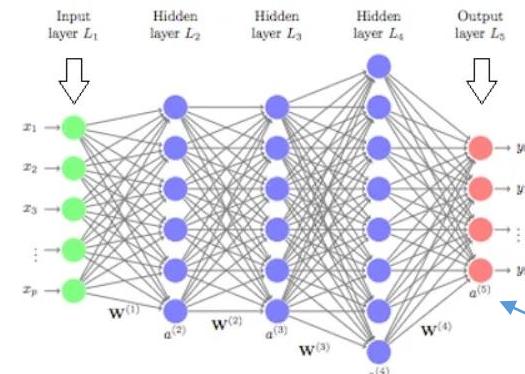
- Definition, properties, and examples.
- Length, inner products, norms, metrics.
- Linear Independence.

## Matrices

- Definition, properties, and examples.
- Trace, inner products, norms, metrics.
- Symmetric matrices, positive semidefinite matrices.
- Eigenvalues, eigenvalue decomposition.
- Cholesky Decomposition.



$$\begin{matrix} & 1 & 2 & \dots & n \\ 1 & a_{11} & a_{12} & \dots & a_{1n} \\ 2 & a_{21} & a_{22} & \dots & a_{2n} \\ 3 & a_{31} & a_{32} & \dots & a_{3n} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ m & a_{m1} & a_{m2} & \dots & a_{mn} \end{matrix}$$



Courtesy of Vishal Yadav

# Euclidian Space

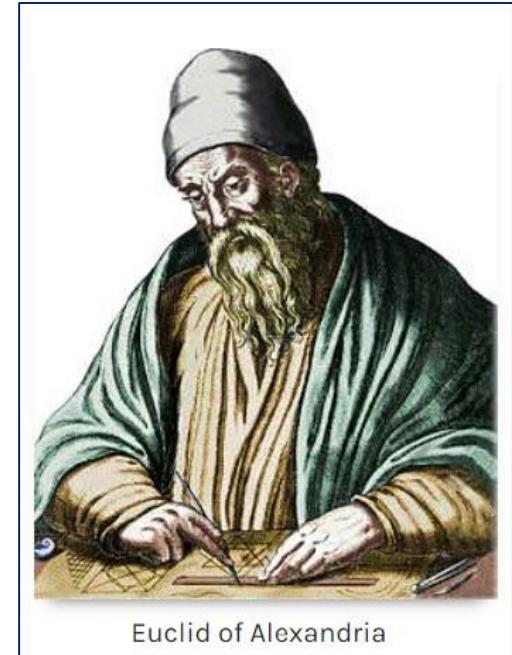
## Euclidian Space:

A space in any finite number of dimensions  $n$ , in which points are designated by coordinates (one for each dimension) and the distance between two points is given by a distance formula.

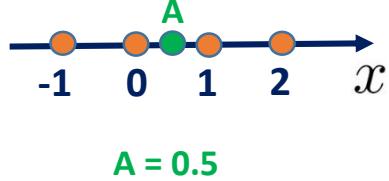
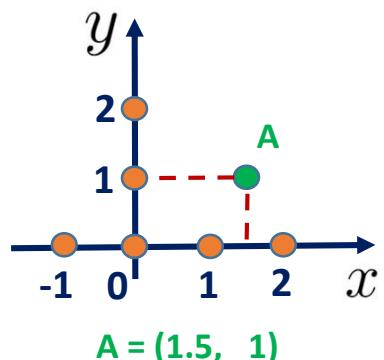
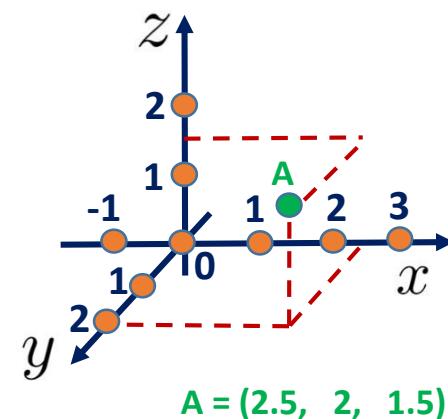
Typical notation:  $\mathbb{R}^n$

## Coordinates of Points

An **ordered collection** of  $n$  numbers, identifying the location of points in space with respect to predefined reference lines (called **coordinate axes**)



## Examples:

 $\mathbb{R}^1$  $\mathbb{R}^2$  $\mathbb{R}^3$

# Euclidian Space

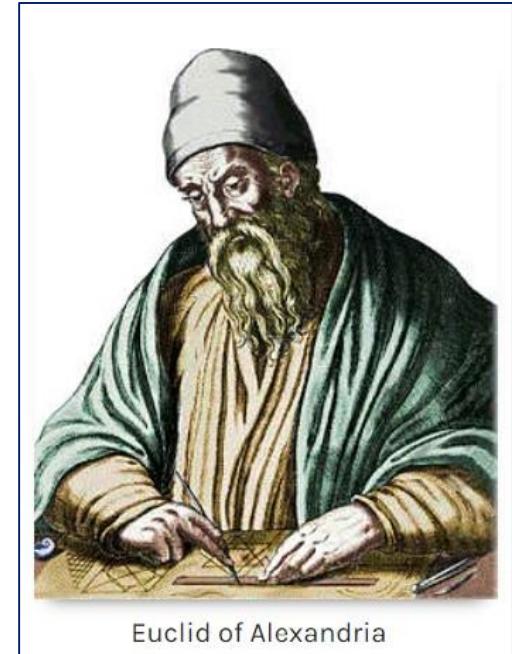
## Euclidian Space:

A space in any finite number of dimensions  $n$ , in which points are designated by coordinates (one for each dimension) and the distance between two points is given by a distance formula.

Typical notation:  $\mathbb{R}^n$

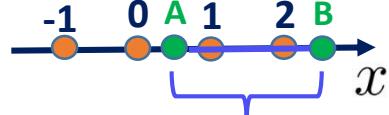
## Distance Formula

Algebraic expression that gives the distances between pairs of points in terms of their coordinates



## Examples:

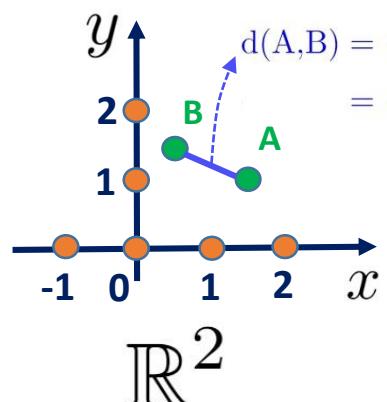
$$A = 0.5 \quad B = 2.5$$



$$d(A,B) = |A - B| = 2$$

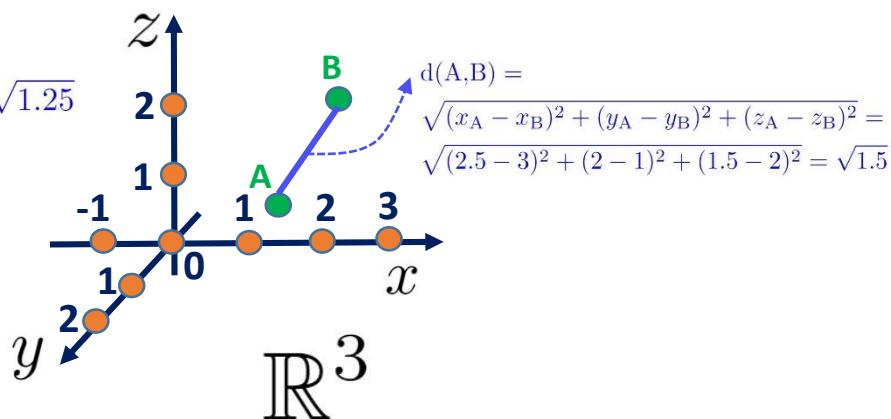
$$\mathbb{R}^1$$

$$A = (1.5, 1) \quad B = (0.5, 1.5)$$



$$\mathbb{R}^2$$

$$A = (2.5, 2, 1.5) \quad B = (3, 1, 2)$$



$$\mathbb{R}^3$$

# Euclidian Spaces

In  $n$ -dimensional Euclidian space  $\mathbb{R}^n$ , each point is designated by  $n$  coordinates:  $A = (x_1^{(A)}, x_2^{(A)}, \dots, x_n^{(A)})$

Distance (Euclidian) between two points  $A = (x_1^{(A)}, x_2^{(A)}, \dots, x_n^{(A)})$  and  $B = (x_1^{(B)}, x_2^{(B)}, \dots, x_n^{(B)})$  is defined as:

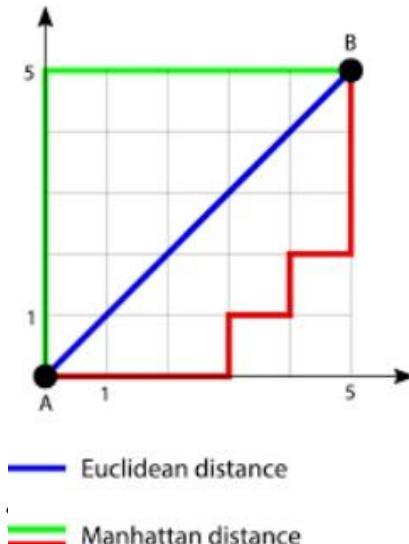
$$d(A, B) = \sqrt{\sum_{i=1}^n (x_i^{(A)} - x_i^{(B)})^2}$$

## Remark:

Distance function  $d(A, B)$  in Euclidian space, can be defined in many ways. Actually any function  $\rho(A, B)$  satisfying the following metric axioms is a valid distance:

## Example:

- For any distinct points A and B it is positive:  $\rho(A, B) > 0$
- For any point A:  $\rho(A, A) = 0$
- For any point A and B:  $\rho(A, B) = \rho(B, A)$
- For any points A and B it is true that  $\rho(A, B) \leq \rho(A, C) + \rho(C, B)$  for any point C.



# Vectors in Euclidian Spaces

Vector in Euclidian space  $\mathbb{R}^n$  is a geometric object that has magnitude and direction. It characterized by two points:

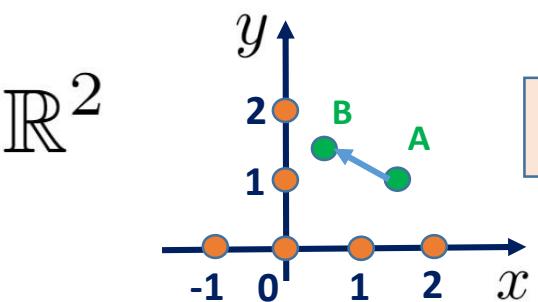
the beginning point  $A = (a_1, \dots, a_n)$  and the end point  $B = (b_1, \dots, b_n)$

The diagram shows a 2D coordinate system with x and y axes. A vector  $\vec{AB}$  originates from point  $A$  and ends at point  $B$ . A right-angle bracket on the left indicates the vector's direction. To the right, the vector is represented as a column matrix  $\vec{AB} = \begin{bmatrix} b_1 - a_1 \\ b_2 - a_2 \\ \vdots \\ b_n - a_n \end{bmatrix}$ . Blue arrows point from the labels  $b_1 - a_1$ ,  $b_2 - a_2$ , and  $b_n - a_n$  to the corresponding entries in the matrix.

**Magnitude** of vector  $\vec{AB}$  is its length. Typically, is defined via distance measure between points  $A$  and  $B$

## Examples:

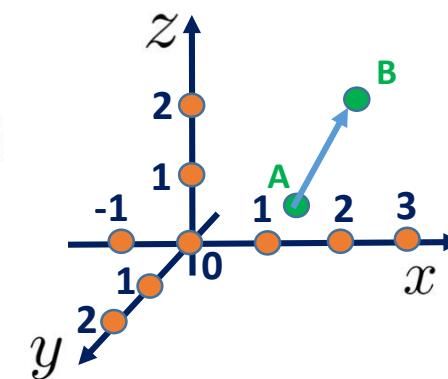
$$A = (1.5, 1) \quad B = (0.5, 1.5)$$



$$\vec{AB} = \begin{bmatrix} 0.5 - 1.5 \\ 1.5 - 1 \end{bmatrix} = \begin{bmatrix} -1 \\ 0.5 \end{bmatrix}$$

$$A = (2.5, 2, 1.5) \quad B = (3, 1, 2)$$

$$\mathbb{R}^3$$



$$\vec{AB} = \begin{bmatrix} 3 - 2.5 \\ 1 - 2 \\ 2 - 1.5 \end{bmatrix} = \begin{bmatrix} 0.5 \\ -1 \\ 0.5 \end{bmatrix}$$

# Vectors in Euclidian Spaces

Typically, the beginning point  $A$  of vector  $\overrightarrow{AB}$  is chosen to be  $A = (0, \dots, 0)$ . Hence:

$$\overrightarrow{AB} = \overrightarrow{B} = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{bmatrix} = \underbrace{[b_1, b_2, \dots, b_n]}_{\text{row vector representation}}^T$$

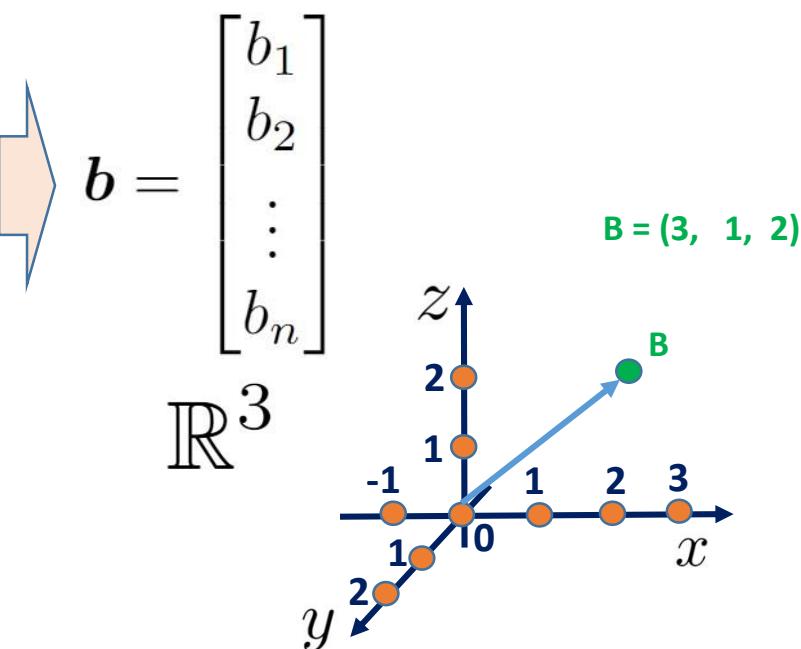
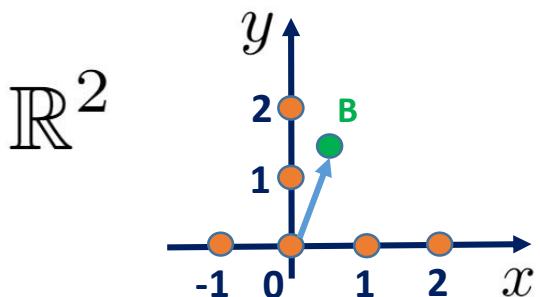
Transpose operation, allows to write row vector as column vector

column vector representation

Typical notation for vectors uses bold lowercase letters:

## Examples:

$$B = (0.5, 1.5)$$



# Important Operations on Vectors

## Linear Operations:

$$\mathbf{v} = \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_n \end{bmatrix}, \quad \mathbf{u} = \begin{bmatrix} u_1 \\ u_2 \\ \vdots \\ u_n \end{bmatrix}, \quad \alpha\mathbf{v} + \beta\mathbf{u} = \begin{bmatrix} \alpha v_1 + \beta u_1 \\ \alpha v_2 + \beta u_2 \\ \vdots \\ \alpha v_n + \beta u_n \end{bmatrix}$$

scalars  
(any real numbers)

both  $\mathbf{v}$  and  $\mathbf{u}$  must be the same dimensionality

Geometric interpretation:

## Examples:

$$\mathbf{v} = \begin{bmatrix} 1 \\ 2 \\ 5 \end{bmatrix}, \quad \mathbf{u} = \begin{bmatrix} 3 \\ 9 \\ 15 \end{bmatrix}, \quad 2\mathbf{v} + 0.5\mathbf{u} = \begin{bmatrix} 2 \cdot 1 + 0.5 \cdot 3 \\ 2 \cdot 2 + 0.5 \cdot 9 \\ 2 \cdot 5 + 0.5 \cdot 15 \end{bmatrix} = \begin{bmatrix} 3.5 \\ 8.5 \\ 17.5 \end{bmatrix}$$

$\alpha\mathbf{v}$   
 $\beta\mathbf{u}$   
 $0 < \beta < 1$

# Important Operations on Vectors

## Inner product:

$$\mathbf{v} = \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_n \end{bmatrix}, \quad \mathbf{u} = \begin{bmatrix} u_1 \\ u_2 \\ \vdots \\ u_n \end{bmatrix}, \quad \langle \mathbf{u}, \mathbf{v} \rangle = \mathbf{v}^T \mathbf{u} = \underbrace{v_1 u_1 + v_2 u_2 + \dots + v_n u_n}_{\text{Result of this product}}$$

is NOT a vector, but a scalar.

both  $\mathbf{v}$  and  $\mathbf{u}$  must be the same dimensionality

## Application in ML: Basic properties:

Maximal margin hyperplane for support vector machines has form:  $\mathbf{w}^T \mathbf{x} - b = 0$

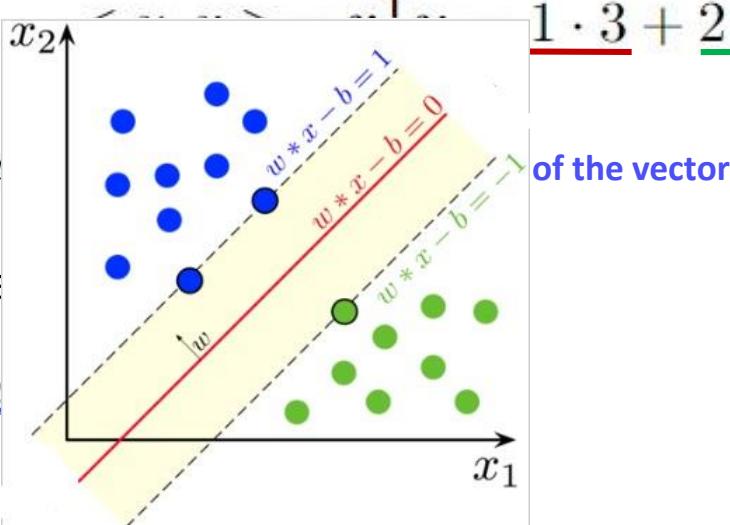
$$\mathbf{v} = \begin{bmatrix} 1 \\ 2 \\ 5 \end{bmatrix}, \quad \mathbf{u} = \begin{bmatrix} 3 \\ -1 \\ 0.5 \end{bmatrix}$$

for support vector machines has form:  $\mathbf{w}^T \mathbf{x} - b = 0$

$\mathbf{w}^T \mathbf{x} - b = 1$   $\mathbf{w}^T \mathbf{x} - b = 0$   $\mathbf{w}^T \mathbf{x} - b = -1$

parameters of the vector

- Multiplication on itself:  $\mathbf{v}^T \mathbf{v}$
- Orthogonality of vectors  $\mathbf{u}$  and  $\mathbf{v}$ :  $\mathbf{u}^T \mathbf{v} = 0$
- Linearity:  $(\alpha \mathbf{v} + \beta \mathbf{u})^T \mathbf{w} = \alpha \mathbf{v}^T \mathbf{w} + \beta \mathbf{u}^T \mathbf{w}$



# Important Operations on Vectors

## Component wise-product:

$$\mathbf{v} = \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_n \end{bmatrix}, \quad \mathbf{u} = \begin{bmatrix} u_1 \\ u_2 \\ \vdots \\ u_n \end{bmatrix}, \quad \mathbf{v} \odot \mathbf{u} = \begin{bmatrix} v_1 \cdot u_1 \\ v_2 \cdot u_2 \\ \vdots \\ v_n \cdot u_n \end{bmatrix}$$

both  $\mathbf{v}$  and  $\mathbf{u}$  must be the same dimensionality

Result of this product  
is a vector.

## Application in ML:

### Examples:

ADAM update equations:

$$\mathbf{v} = \begin{bmatrix} 1 \\ 2 \\ 5 \end{bmatrix}, \quad \mathbf{u} = \begin{bmatrix} 3 \\ -1 \\ 0.2 \end{bmatrix}$$
$$\mathbf{m}_t = \beta_1 \mathbf{m}_{t-1} + (1 - \beta_1) \mathbf{g}_t = \begin{bmatrix} 1 \cdot 3 \\ \underline{2} \cdot (-1) \\ 5 \cdot 0.2 \end{bmatrix} = \begin{bmatrix} 3 \\ -2 \\ 1 \end{bmatrix}$$

$$\mathbf{v}_t = \beta_2 \mathbf{v}_{t-1} + (1 - \beta_2) \mathbf{g}_t \odot \mathbf{g}_t$$

$$\theta_t = \theta_{t-1} - \alpha \frac{\mathbf{m}_t}{\sqrt{\mathbf{v}_t} + \epsilon \mathbf{1}}$$

Component-wise product and division.

# Important Operations on Vectors

## Kronecker-product:

$$\mathbf{v} = \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_n \end{bmatrix}, \quad \mathbf{u} = \begin{bmatrix} u_1 \\ u_2 \\ \vdots \\ u_m \end{bmatrix}, \quad \mathbf{v} \otimes \mathbf{u} =$$

$$\begin{bmatrix} v_1\mathbf{u} \\ v_2\mathbf{u} \\ \vdots \\ v_n\mathbf{u} \end{bmatrix} = \begin{bmatrix} v_1 \cdot u_1 \\ v_1 \cdot u_2 \\ \vdots \\ v_1 \cdot u_m \\ v_2 \cdot u_1 \\ v_2 \cdot u_2 \\ \vdots \\ v_2 \cdot u_m \\ \vdots \\ v_n \cdot u_1 \\ v_n \cdot u_2 \\ \vdots \\ v_n \cdot u_m \end{bmatrix}$$

Result of this product  
is a vector in  $\mathbb{R}^{nm}$

Vectors  $\mathbf{v}$  and  $\mathbf{u}$  can have different dimensionality

Not symmetric:  $\mathbf{v} \otimes \mathbf{u} \neq \mathbf{u} \otimes \mathbf{v}$

## Example:

$$\mathbf{v} = \begin{bmatrix} 1 \\ 2 \end{bmatrix}, \quad \mathbf{u} = \begin{bmatrix} 3 \\ -5 \\ 20 \end{bmatrix},$$

$$\mathbf{v} \otimes \mathbf{u} = \begin{bmatrix} 1 \cdot 3 \\ 1 \cdot (-5) \\ 1 \cdot 20 \\ 2 \cdot 3 \\ 2 \cdot (-5) \\ 2 \cdot 20 \end{bmatrix} = \begin{bmatrix} 3 \\ -5 \\ 20 \\ 6 \\ -10 \\ 40 \end{bmatrix}$$

# Vector Norms

Norm of vector  $v$  is any univariate function  $\|\cdot\|$ , satisfying norm axioms:

- definiteness:  $\|v\| = 0 \rightarrow v = \mathbf{0}$
- absolute homogeneity:  $\|\alpha v\| = |\alpha| \|v\|$  for any vectors  $v$  and scalar  $\alpha$
- subadditivity or triangular inequality:  $\|u + v\| \leq \|u\| + \|v\|$  for any vectors  $v$  and  $u$
- Non-negativity:  $\|v\| \geq 0$  for any vector  $v$

## Application in ML:

Train ML model:

$$\min_w \mathcal{L}_{\text{loss}}(w, \mathcal{D}) + \|w\|$$


Regularisation term

# Some Vector Norms and some Not Vector Norms

- Euclidian (or  $L_2$ -norm):

$$\|v\|_2 = \sqrt{v_1^2 + v_2^2 + \dots + v_n^2} \quad \Rightarrow \quad v = [1, 3, 5]^\top, \|v\|_2 = \sqrt{1^2 + 3^2 + 5^2} = \sqrt{1 + 9 + 25} = \sqrt{35}.$$

- $L_1$ -norm:

$$\|v\|_1 = |v_1| + \dots + |v_n| \quad \Rightarrow \quad u = [1, -3, 4]^\top, \|u\|_1 = |1| + |-3| + |4| = 1 + 3 + 4 = 8$$

- Infinity norm (or  $L_\infty$ -norm):

$$\|v\|_\infty = \max\{|v_1|, |v_2|, \dots, |v_n|\} \quad \Rightarrow \quad w = [-20, -3, 8]^\top, \|w\|_\infty = \max\{|-20|, |-3|, |8|\} = \max\{20, 3, 8\} = 20$$

- Sparsity norm (or  $L_0$ -norm):

$$\|v\|_0 = |\{i \in [1, \dots, n] \text{ such that } v_i \neq 0\}| \quad \Rightarrow \quad h = [0, -3, 0, 4]^\top, \|h\|_0 = 2$$

One of them, actually is not a norm... Which one and why?

# Vector Norms (Equivalence)

For the same vector  $\mathbf{v}$  different norms gives different values....

However, for any two norm  $\|\cdot\|_p$  and  $\|\cdot\|_q$  there are constants  $0 < a \leq b$  such that:

$$a\|\mathbf{v}\|_p \leq \|\mathbf{v}\|_q \leq b\|\mathbf{v}\|_p \quad \text{for any vector } \mathbf{v}$$

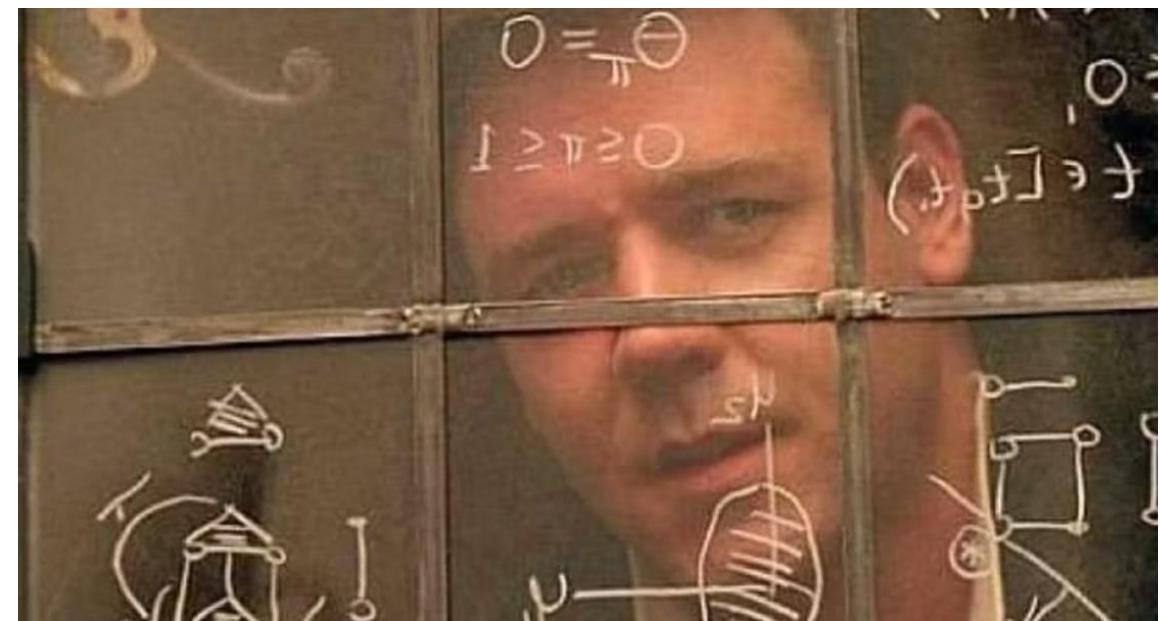
**Equivalence of norms**

## Exercise:

Prove the following statements:

- If  $\|\cdot\|$  is a valid norm, then  $\|\mathbf{u} - \mathbf{v}\|$  is a valid distance function between points  $\mathbf{u}$  and  $\mathbf{v}$
- For any vectors  $\mathbf{u}$  and  $\mathbf{v}$  any scalar product function:

$$|\langle \mathbf{u}, \mathbf{v} \rangle| \leq \|\mathbf{u}\|_2 \|\mathbf{v}\|_2$$



# Linear Independence

Collection of vectors  $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_K\}$  is called linear independent if:

$$\alpha_1 \mathbf{v}_1 + \alpha_2 \mathbf{v}_2 + \dots + \alpha_K \mathbf{v}_K = \mathbf{0} \text{ implies } \alpha_1 = 0, \alpha_2 = 0, \dots, \alpha_K = 0$$

## Example:

$$e_1 = \begin{bmatrix} 1 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}, \quad e_2 = \begin{bmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}, \quad e_3 = \begin{bmatrix} 0 \\ 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix}, \quad \dots, \quad e_n = \begin{bmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 1 \end{bmatrix}$$

Standard Basis in  $\mathbb{R}^n$

are linearly independent

any vector in  $\mathbb{R}^n$  can be written as their linear combination:

$$\begin{bmatrix} 2.7 \\ -13 \\ 5 \\ \frac{2}{7} \end{bmatrix} = 2.7 \underbrace{\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}}_{e_1} - 13 \underbrace{\begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}}_{e_2} + 5 \underbrace{\begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}}_{e_3} + \frac{2}{7} \underbrace{\begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}}_{e_4}$$

$$\mathbf{v}_1 = \begin{bmatrix} 1 \\ 1 \end{bmatrix}, \quad \mathbf{v}_2 = \begin{bmatrix} 0 \\ -1 \end{bmatrix}, \quad \mathbf{v}_3 = \begin{bmatrix} -1 \\ 1 \end{bmatrix}$$

$$\mathbf{v}_1 + \mathbf{v}_3 - 2\mathbf{v}_2 = \begin{bmatrix} 1 \\ 1 \end{bmatrix} + \begin{bmatrix} -1 \\ 1 \end{bmatrix} + 2 \begin{bmatrix} 0 \\ -1 \end{bmatrix} = \begin{bmatrix} 1 - 1 + 0 \\ 1 + 1 - 2 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

are NOT linearly independent

Now, we are moving to the next topic....

matrices...



Reality

Given

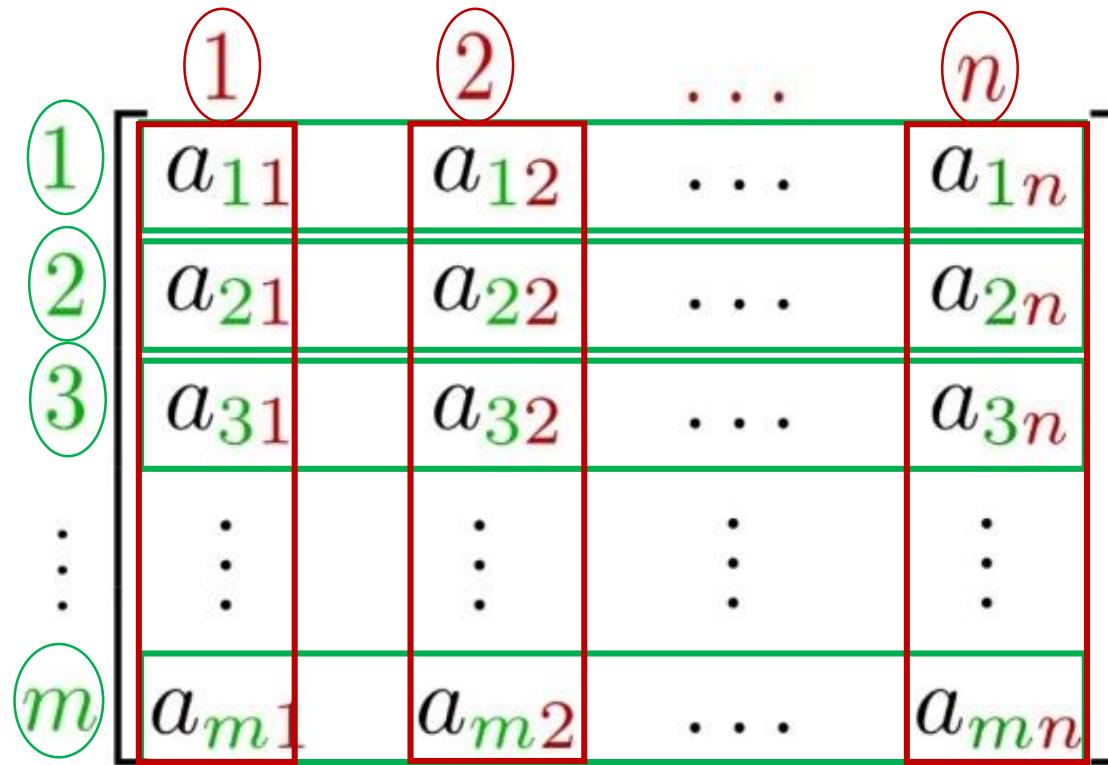
$$A = \begin{bmatrix} 2 & 3 & 0 \\ 4 & 3 & 7 \end{bmatrix} \quad B = \begin{bmatrix} 5 & 7 \\ 6 & 4 \end{bmatrix}$$

Find AB and BA.



# Matrices

Matrix of size  $m$  by  $n$  of real numbers has the following form:



## Examples:

$$A = \begin{bmatrix} 0 & 2 \\ 3 & 4 \\ 8 & -0.65 \\ \frac{1}{3} & 5 \end{bmatrix} \in \mathbb{R}^{4 \times 2}, \quad B = \begin{bmatrix} 0 & 2 & 3 & 7 \\ 8 & 4 & 0.4 & -5 \\ 2 & -0.65 & \frac{3}{7} & -\log 3 \end{bmatrix} \in \mathbb{R}^{3 \times 4} \quad C = \begin{bmatrix} 0 & 2 & 6 & \frac{2}{3} \\ 3 & 4 & -0.75 & 43 \\ 8 & -0.65 & \log 7 & 1 \\ \frac{1}{3} & 5 & 0 & -54 \end{bmatrix} \in \mathbb{R}^{4 \times 4}$$

If  $m = n$  matrix is called square matrix

# Matrix Operations (Transpose)

Transpose operation on matrix swaps rows with columns:

$$A = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \dots & a_{mn} \end{bmatrix} \in \mathbb{R}^{m \times n} \quad A^T = \begin{bmatrix} a_{11} & a_{21} & \dots & a_{m1} \\ a_{12} & a_{22} & \dots & a_{m2} \\ \vdots & \vdots & \ddots & \vdots \\ a_{1n} & a_{2n} & \dots & a_{mn} \end{bmatrix} \in \mathbb{R}^{n \times m}$$

Exam  $A$  is called **symmetric** if:  $A^T = A$

$$A = \begin{bmatrix} 1 & 5 & 0 \\ 5 & -2 & 9 \\ 0 & 9 & 91 \end{bmatrix} \rightarrow A^T = \begin{bmatrix} 1 & 5 & 0 \\ 5 & -2 & 9 \\ 0 & 9 & 91 \end{bmatrix} \Rightarrow A = A^T \in \mathbb{R}^{3 \times 3} \text{ symmetric:}$$

$$B = \begin{bmatrix} 1 & 5 & 0 \\ -7 & 2.3 & 9 \\ -12 & 9 & 91 \end{bmatrix} \rightarrow B^T = \begin{bmatrix} 1 & -7 & -12 \\ 5 & 2.3 & 9 \\ 0 & 9 & 91 \end{bmatrix} \Rightarrow B \neq B^T \in \mathbb{R}^{3 \times 3} \text{ NOT symmetric:}$$

Square matrix remains square after transposition

# Matrix Operations (Linear Operations)

Linear operation for matrices are element-wise:

$$A = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \cdots & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mn} \end{bmatrix} \quad B = \begin{bmatrix} b_{11} & b_{12} & \cdots & b_{1n} \\ b_{21} & b_{22} & \cdots & b_{2n} \\ \vdots & \vdots & \cdots & \vdots \\ b_{m1} & b_{m2} & \cdots & b_{mn} \end{bmatrix}$$

for any  $\alpha, \beta \in \mathbb{R}$

$$\alpha A + \beta B = \begin{bmatrix} \alpha a_{11} + \beta b_{11} & \alpha a_{12} + \beta b_{12} & \cdots & \alpha a_{1n} + \beta b_{1n} \\ \alpha a_{21} + \beta b_{21} & \alpha a_{22} + \beta b_{22} & \cdots & \alpha a_{2n} + \beta b_{2n} \\ \vdots & \vdots & \cdots & \vdots \\ \alpha a_{m1} + \beta b_{m1} & \alpha a_{m2} + \beta b_{m2} & \cdots & \alpha a_{mn} + \beta b_{mn} \end{bmatrix}$$

both  $A, B \in \mathbb{R}^{m \times n}$  must be the same dimensionality

## Examples:

$$A = \begin{bmatrix} 2 & 3 & 5 \\ -1 & 2 & 3 \\ -4 & 10 & 7 \end{bmatrix} \quad B = \begin{bmatrix} 10 & 100 & 1000 \\ 1 & -3 & 5 \\ 5 & -5 & 3 \end{bmatrix}$$



$$2A - 0.5B = \begin{bmatrix} 2 \cdot 2 - 0.5 \cdot 10 & 2 \cdot 3 - 0.5 \cdot 100 & 2 \cdot 5 - 0.5 \cdot 1000 \\ 2 \cdot (-1) - 0.5 \cdot 1 & 2 \cdot 2 - 0.5 \cdot (-3) & 2 \cdot 3 - 0.5 \cdot 5 \\ 2 \cdot (-4) - 0.5 \cdot 5 & 2 \cdot 10 - 0.5 \cdot (-5) & 2 \cdot 7 - 0.5 \cdot 3 \end{bmatrix} = \begin{bmatrix} -1 & -44 & -490 \\ -2.5 & 5.5 & 3.5 \\ -10.5 & 22.5 & 12.5 \end{bmatrix}$$

# Matrix Operations (Multiplication)

For two matrices  $A \in \mathbb{R}^{m \times n}$  and  $B \in \mathbb{R}^{n \times p}$ :

$$A = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ \boxed{a_{i1} & a_{i2} & \dots & a_{in}} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \dots & a_{mn} \end{bmatrix}$$

and  $B = \begin{bmatrix} b_{11} & b_{12} & \dots & \boxed{b_{1j}} & \dots & b_{1p} \\ b_{21} & b_{22} & \dots & \boxed{b_{2j}} & \dots & b_{2p} \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ b_{n1} & b_{n2} & \dots & \boxed{b_{nj}} & \dots & b_{np} \end{bmatrix}$

Number of columns in  $A$  must be equal to number of rows in  $B$



$$AB = C = \begin{bmatrix} c_{11} & c_{12} & \dots & c_{1j} & \dots & c_{1p} \\ c_{21} & c_{22} & \dots & c_{2j} & \dots & c_{2p} \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ c_{i1} & c_{i2} & \dots & \boxed{c_{ij}} & \dots & \vdots \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ c_{m1} & c_{m2} & \dots & c_{mj} & \dots & c_{mp} \end{bmatrix} \in \mathbb{R}^{m \times p}$$

where

$$\boxed{c_{ij}} =$$

$$\begin{bmatrix} a_{i1} & a_{i2} & \dots & a_{in} \end{bmatrix}$$

$$\begin{bmatrix} b_{1j} \\ b_{2j} \\ \vdots \\ b_{nj} \end{bmatrix}$$

$$= \sum_{k=1}^n a_{ik} b_{nk}$$

inner product

# Matrix Operations (Multiplication)

## Examples:

$$A = \begin{bmatrix} 1 & 2 \\ -3 & 4 \\ 5 & 0.5 \end{bmatrix} \quad B = \begin{bmatrix} -1 & 0 & 0 & -1 \\ 1 & 1 & -2 & 0 \end{bmatrix}$$

Notice, number of columns in  $A$  is 2  
and number of rows in  $B$  is 2

$$C = \begin{bmatrix} 1 \cdot (-1) + 2 \cdot 1 & 1 \cdot 0 + 2 \cdot 1 & 1 \cdot 0 + 2 \cdot (-2) & 1 \cdot (-1) + 2 \cdot 0 \\ (-3) \cdot (-1) + 4 \cdot 1 & (-3) \cdot 0 + 4 \cdot 1 & (-3) \cdot 0 + 4 \cdot (-2) & (-3) \cdot (-1) + 4 \cdot 0 \\ 5 \cdot (-1) + 0.5 \cdot 1 & 5 \cdot 0 + 0.5 \cdot 1 & 5 \cdot 0 + 0.5 \cdot (-2) & 5 \cdot (-1) + 0.5 \cdot 0 \end{bmatrix} = \begin{bmatrix} 1 & 2 & -4 & -1 \\ 7 & 4 & -8 & 4 \\ -4.5 & 0.5 & -1 & -5 \end{bmatrix}$$

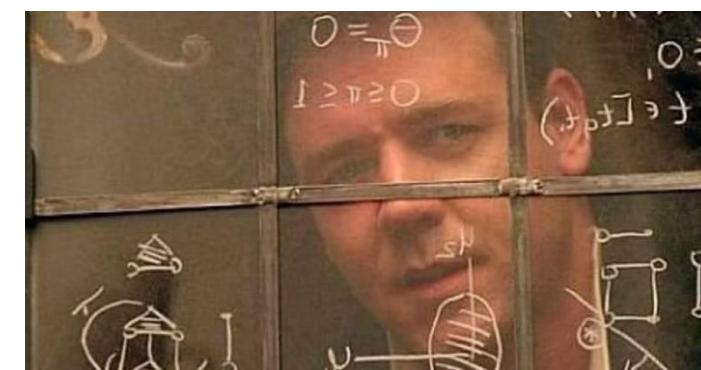
$$A = \begin{bmatrix} 1 & 2 & 0 \\ -3 & 4 & 8 \\ 5 & 0.5 & -3 \\ 0 & -1 & 1 \end{bmatrix}, \quad b = \begin{bmatrix} 10 \\ 100 \\ 100 \end{bmatrix} \quad Ab = \begin{bmatrix} 1 \cdot 10 + 2 \cdot 100 + 0 \cdot 1000 \\ (-3) \cdot 10 + 4 \cdot 100 + 8 \cdot 1000 \\ 5 \cdot 10 + 0.5 \cdot 100 + (-3) \cdot 1000 \\ 0 \cdot 10 + (-1) \cdot 100 + 1 \cdot 1000 \end{bmatrix} = \begin{bmatrix} 210 \\ 8370 \\ -2800 \\ 900 \end{bmatrix}$$

Matrix vector multiplication

## Exercise:

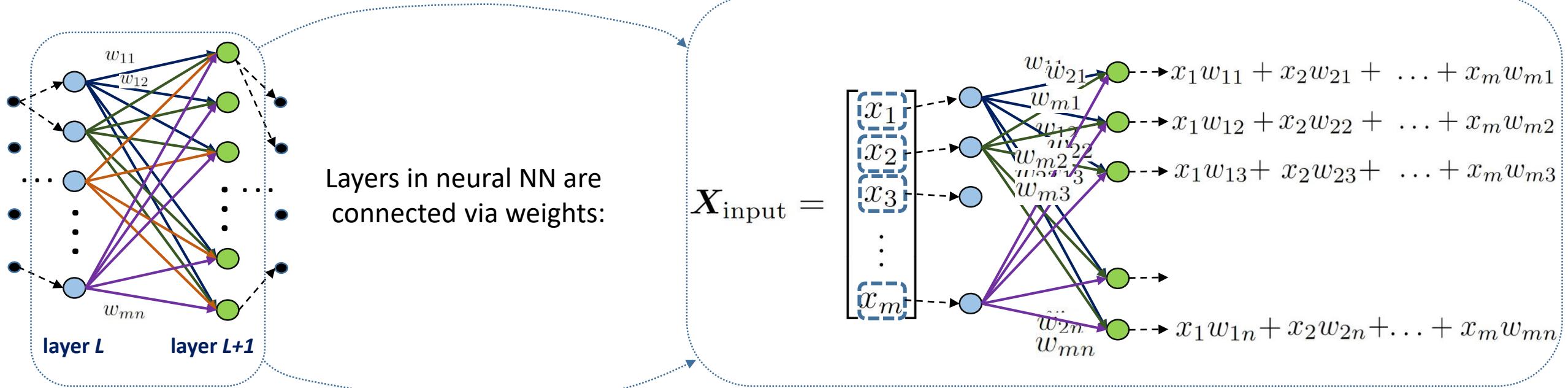
Give example of two square matrices  $A, B$  such that  $AB \neq BA$

Show that for any square matrix:  $AI = IA$  where  $I = \begin{bmatrix} 1 & 0 & 0 & \dots & 0 \\ 0 & 1 & 0 & \dots & 0 \\ 0 & 0 & 1 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & 1 \end{bmatrix} \in \mathbb{R}^{n \times n}$



# Matrix Operations (Multiplication)

## ML application:



Results mathematical expressions:

$$x_1w_{11} + x_2w_{21} + x_3w_{31} + \dots + x_mw_{m1}$$

$$x_1w_{12} + x_2w_{22} + x_3w_{32} + \dots + x_mw_{m2}$$

$$x_1w_{13} + x_2w_{23} + x_3w_{33} + \dots + x_mw_{m3}$$

$\vdots$

$$x_1w_{1n} + x_2w_{2n} + x_3w_{3n} + \dots + x_mw_{mn}$$

## Exercise:

A portrait of a man with glasses and a mustache, possibly a historical figure related to mathematics or science, is positioned above the exercise area.

$$\equiv \begin{bmatrix} w_{11} & w_{12} & w_{13} & \dots & w_{1n} \\ w_{21} & w_{22} & w_{23} & \dots & w_{2n} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ w_{m1} & w_{m2} & w_{m3} & \dots & w_{mn} \end{bmatrix}^T \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_m \end{bmatrix} = W^T X_{\text{input}}$$

# Matrix Multiplication (Complexity)

How computational expensive is matrix multiplication?

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

$$B = \begin{bmatrix} b_{11} & b_{12} & \dots & b_{1n} \\ b_{21} & b_{22} & \dots & b_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ b_{n1} & b_{n2} & \dots & b_{nn} \end{bmatrix} \in \mathbb{R}^{n \times n}$$

$$C = AB = \begin{bmatrix} c_{11} & c_{12} & \dots & c_{1n} \\ c_{21} & c_{22} & \dots & c_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ c_{n1} & c_{n2} & \dots & c_{nn} \end{bmatrix}$$

$c_{ij} = \sum_{k=1}^n a_{ik} b_{kj}$

Each entry requires  
n multiplications and  
n additions

There are total  $n^2$  entries in matrix  $C$

Total number of arithmetic operations:

$$2n^3 = \mathcal{O}(n^3)$$

## Best Current Theoretical Algorithm:

A Refined Laser Method and Faster Matrix Multiplication

Josh Alman\*

Virginia Vassilevska Williams†

October 13, 2020

$\mathcal{O}(n^w)$ ,  $w \leq 2.37286$   
previous best:  $w \leq 2.37287$   
But... it is galactic algorithm

## Widely Used Practical Algorithm

Published: August 1969

Gaussian elimination is not optimal

Volker Strassen

Numerische Mathematik 13, 354–356 (1969) | Cite this article

$$\mathcal{O}(n^{\log_2 7})$$

## Cool Recent Advances

Discovering faster matrix multiplication algorithms with reinforcement learning

Alhussein Fawzi, Matej Balog, Aja Huang, Thomas Hubert, Bernardino Romera-Paredes, Mohammadamin Barekatian, Alexander Novikov, Francisco J. R. Ruiz, Julian Schrittweis, Grzegorz Swirszcz, David Silver, Demis Hassabis & Pushmeet Kohli

Nature 610, 47–53 (2022) | Cite this article

$$\mathcal{O}(n^{\log_4 47})$$

# Inner Products of Matrices

Given two matrices  $A \in \mathbb{R}^{m \times n}$  and  $B \in \mathbb{R}^{n \times m}$  can we compute inner product  $\langle A, B \rangle$  between them?

## Method 1: Via the Trace

Trace of the square matrix  $G \in \mathbb{R}^{k \times k}$ :

$$G = \begin{bmatrix} g_{11} & g_{12} & \dots & g_{1k} \\ g_{21} & g_{22} & \dots & g_{2k} \\ \vdots & \vdots & \ddots & \vdots \\ g_{k1} & g_{k2} & \dots & g_{kk} \end{bmatrix} \quad \text{Tr}(G) = \sum_{i=1}^k g_{ii}$$

**HELL YEAH**

## Example:

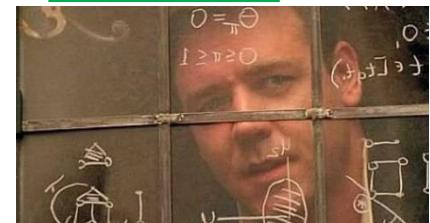
$$G = \begin{bmatrix} 1 & 3 & 4 \\ 41 & -3 & 80 \\ -0.9 & 65 & 21 \end{bmatrix} \quad \text{Tr}(G) = 1 + (-3) + 21 = 19$$

The inner product now can be defined as:  $\langle A, B \rangle = \text{Tr}(AB) = \sum_{i=1}^m c_{ii} = \sum_{i=1}^m \sum_{j=1}^n a_{ik} b_{kj}$

## Example:

$$A = \begin{bmatrix} 1 & 3 \\ 0 & -2 \\ 9 & -5 \\ 5 & -0.5 \end{bmatrix}, \quad B = \begin{bmatrix} 0 & 1 & 6 & 2 \\ -1 & 5 & -2.5 & 1 \end{bmatrix} \quad AB = \begin{bmatrix} 3 & 16 & -1.5 & 5 \\ 2 & -10 & 5 & -2 \\ 95 & -16 & 66.5 & 13 \\ 0.5 & 2.5 & 31.25 & 9.5 \end{bmatrix} \quad \Rightarrow \text{Tr}(AB) = 3 + (-10) + 66.5 + 9.5 = 69$$

## Exercise:



Show symmetric property:  $\text{Tr}(AB) = \text{Tr}(BA)$

## Hint:

Use definition of trace and try reordering of terms

# Inner Products of Matrices

Given two matrices  $A \in \mathbb{R}^{m \times n}$  and  $B \in \mathbb{R}^{n \times m}$  can we compute inner product  $\langle A, B \rangle$  between them?

## Method 2: Via Matrix Vectorization:

Any matrix  $G \in \mathbb{R}^{n \times m}$  can be represented using vector notation:

$$G = \begin{bmatrix} g_{11} & g_{12} & \dots & g_{1m} \\ g_{21} & g_{22} & \dots & g_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ g_{n1} & g_{n2} & \dots & g_{nm} \end{bmatrix}$$

$$\text{vec}(G) = \begin{bmatrix} g_{11} \\ g_{12} \\ \vdots \\ g_{1m} \\ g_{21} \\ g_{22} \\ \vdots \\ g_{2m} \\ \vdots \\ g_{n1} \\ g_{n2} \\ \vdots \\ g_{nm} \end{bmatrix}$$

**concatenating rows :**

$$\text{vec}(G) = \begin{bmatrix} g_{11} \\ g_{21} \\ \vdots \\ g_{n1} \\ g_{12} \\ g_{22} \\ \vdots \\ g_{n2} \\ \vdots \\ g_{1m} \\ g_{2m} \\ \vdots \\ g_{nm} \end{bmatrix}$$

**concatenating columns :**

## Example:

$$G = \begin{bmatrix} 1 & 2 \\ 3 & 4 \\ 5 & 6 \end{bmatrix}$$

$$\text{vec}(G) = \begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \\ 5 \\ 6 \end{bmatrix} \quad \text{vec}(G) = \begin{bmatrix} 1 \\ 3 \\ 5 \\ 2 \\ 4 \\ 6 \end{bmatrix}$$

Now, we can use any inner product between vectors:

$$\langle A, B \rangle = \langle \text{vec}(A), \text{vec}(B) \rangle$$

# Matrix Inverse

Square matrix  $A \in \mathbb{R}^{n \times n}$  is called inverse to another square matrix  $B \in \mathbb{R}^{n \times n}$  if:

$$AB = BA = \underline{\underline{I}}$$

$$I = \begin{bmatrix} 1 & 0 & 0 & \dots & 0 \\ 0 & 1 & 0 & \dots & 0 \\ 0 & 0 & 1 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & 1 \end{bmatrix} \in \mathbb{R}^{n \times n}$$

Identity matrix:

- Only square matrices might have inverse.
- Not all square matrices have inverse:  
 $A = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}$  - is not invertible:

## Example:

$$A = \begin{bmatrix} 1 & 2 & 1 \\ 4 & 4 & 5 \\ 6 & 7 & 7 \end{bmatrix}, \quad B = \begin{bmatrix} -7 & -7 & 6 \\ 2 & 1 & -1 \\ 4 & 5 & -4 \end{bmatrix} \quad AB = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

## Matrix Inverse Computation:

Gaussian Elimination Algorithm

LU Decomposition Algorithms

$\mathcal{O}(n^3)$

In general, big matrices  
are expensive to invert...

But, if matrix has special structure:

- Symmetric diagonal dominant
  - Sparse lower/upper triangular
  - ⋮
- can be inverted in a much faster way

# Solving Linear Equations

System of  $m$  linear equations with  $n$  variables have general form:

$$\begin{cases} a_{11}x_1 + a_{12}x_2 + \dots + a_{1n}x_n = b_1 \\ a_{21}x_1 + a_{22}x_2 + \dots + a_{2n}x_n = b_2 \\ a_{31}x_1 + a_{32}x_2 + \dots + a_{3n}x_n = b_3 \\ \dots \\ a_{m1}x_1 + a_{m2}x_2 + \dots + a_{mn}x_n = b_m \end{cases}$$

Unknown variables :  $x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^n$

Known variables :

$$A = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \dots & a_{mn} \end{bmatrix} \in \mathbb{R}^{m \times n}$$
$$b = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_m \end{bmatrix} \in \mathbb{R}^m$$

In a matrix vector form can be written compactly:

$$Ax = b$$

Example:

$$\begin{aligned} 2x_1 - 3x_2 + 5x_3 &= 2 \\ 6x_1 + 0x_2 - 7x_3 &= -3 \\ 5x_1 - 8x_2 + 0.5x_3 &= 0 \end{aligned}$$

A =  $\begin{bmatrix} 2 & -3 & 5 \\ 6 & 1 & -7 \\ 5 & -8 & 0.5 \end{bmatrix}$

$b = \begin{bmatrix} 2 \\ -3 \\ 0 \end{bmatrix}$

# Solving Linear Equations

Given a system of linear equations:

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

Can be reduced to matrix inversion problem

( by multiplying both sides on  $\mathbf{A}^T$  ):  $\underbrace{\mathbf{A}\mathbf{x}}_{\in \mathbb{R}^{m \times n}} = \mathbf{b} \iff \underbrace{\mathbf{A}^T \mathbf{A}\mathbf{x}}_{\text{square matrix: } \in \mathbb{R}^{n \times n}} = \mathbf{A}^T \mathbf{b}$

$$\xrightarrow{\hspace{1cm}} \mathbf{x}^* = \underbrace{(\mathbf{A}^T \mathbf{A})^{-1}}_{\text{Inverse of } \mathbf{A}^T \mathbf{A}} \mathbf{A}^T \mathbf{b}$$

Sometimes, instead of multiplying on  $\mathbf{A}^T$ , we also multiply on matrix  $\mathbf{G} \in \mathbb{R}^{n \times m}$  such that  $\mathbf{G}\mathbf{A}$  has desirable (for inversion) structure:

$$\mathbf{A}\mathbf{x} = \mathbf{b} \iff \mathbf{G}\mathbf{A}\mathbf{x} = \mathbf{G}\mathbf{b}$$
$$\mathbf{x}^* = (\mathbf{G}\mathbf{A})^{-1} \mathbf{G}\mathbf{b}$$

$\mathbf{G} \in \mathbb{R}^{n \times m}$   
is called preconditioner

Since this reduction, linear equation solver has complexity bounded by inversion complexity:  $\mathcal{O}(n^3)$

In practice, features like sparsity of matrix of the system  $\mathbf{A}$  along with proper preprocessing can have huge effect on performance.

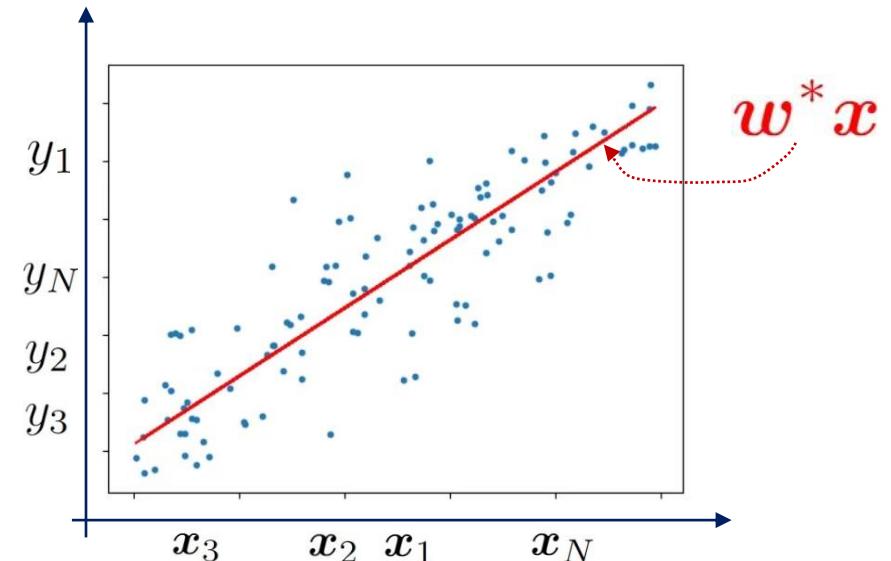
# Solving Linear Equations (ML Application)

In linear regression task, we are given a data-set  $\mathcal{D} = \{\mathbf{x}_i, y_i\}_{i=1}^N$

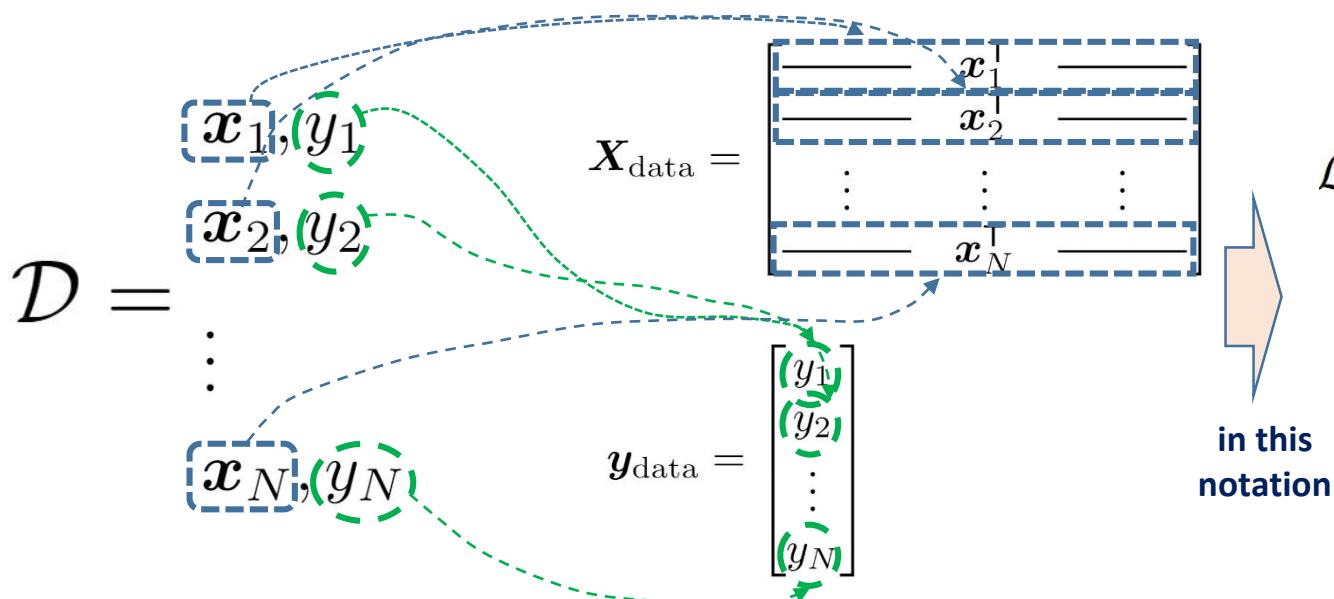
and our goal is to find ML model  $\mathbf{w}^\top \mathbf{x}$  that fits this data.

To find the optimal model parameter we aim optimization problem:

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \mathcal{L}_{\text{loss}}(\mathbf{w}) = \sum_{i=1}^N (\mathbf{w}^\top \mathbf{x}_i - y_i)^2$$



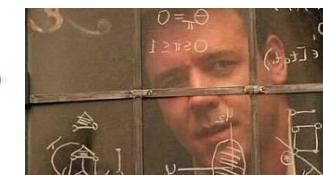
How can we write it in matrix vector form?



in this  
notation

$$\mathcal{L}_{\text{loss}}(\mathbf{w}) = \sum_{i=1}^N (\mathbf{w}^\top \mathbf{x}_i - y_i)^2$$

Why?



Optimal parameter  
is the solution of  
the following system

$$X_{\text{data}}^\top X_{\text{data}} \mathbf{w}^* = X_{\text{data}}^\top \mathbf{y}_{\text{data}}$$

system of linear equations

Lets have 10 min break



# Eigenvalues/Eigenvectors

For square matrix  $A \in \mathbb{R}^{n \times n}$  a non-zero vector  $x \in \mathbb{R}^n$  such that:

$$Ax = \lambda x$$

for some scalar  $\lambda \in \mathbb{R}$

is called eigenvector of matrix  $A$  corresponding to eigenvalue  $\lambda$ . For same value  $\lambda$  you might have several eigenvectors.

## Example:

$$A = \begin{bmatrix} 2 & 1 \\ 0 & 2 \end{bmatrix} \quad x = \begin{bmatrix} 1 \\ 0 \end{bmatrix} \quad Ax = \begin{bmatrix} 2 & 1 \\ 0 & 2 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 2 \cdot 1 + 1 \cdot 0 \\ 0 \cdot 1 + 2 \cdot 0 \end{bmatrix} = 2 \begin{bmatrix} 1 \\ 0 \end{bmatrix} = 2x$$

$x = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$  - is eigenvector  
with eigenvalue  
 $\lambda = 2$

## Geometric meaning:

Linear mapping  $\Phi : \mathbb{R}^n \rightarrow \mathbb{R}^m$  transforms any vector  $x \in \mathbb{R}^n$  into a vector  $\Phi(x) \in \mathbb{R}^m$  such that

$$\Phi(\alpha x + \beta z) = \alpha \Phi(x) + \beta \Phi(z) \quad \alpha, \beta \in \mathbb{R}$$

For any linear mapping  $\Phi : \mathbb{R}^n \rightarrow \mathbb{R}^m$  there exist matrix  $A \in \mathbb{R}^{m \times n}$  such that  $\Phi(x) = Ax$

$Ax = \lambda x$  implies that vector  $x \in \mathbb{R}^n$  does not change after transformation defined by  $A$

# Eigendecomposition and Diagonalisation

Square matrix  $D \in \mathbb{R}^{n \times n}$  is called diagonal if it has zero entries on all off-diagonal positions:

$$D = \begin{bmatrix} d_{11} & 0 & \cdots & 0 \\ 0 & d_{22} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & d_{nn} \end{bmatrix}$$

These positions are the diagonals and called off-diagonal.

$$D = \begin{bmatrix} 4 & 0 & 0 & 0 \\ 0 & -3 & 0 & 0 \\ 0 & 0 & 0.5 & 0 \\ 0 & 0 & 0 & \frac{2}{5} \end{bmatrix} \quad \text{OR} \quad D = \begin{bmatrix} -2 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 3 \end{bmatrix} \quad \text{OR} \quad I = \begin{bmatrix} 1 & 0 & 0 & \cdots & 0 \\ 0 & 1 & 0 & \cdots & 0 \\ 0 & 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \cdots & 1 \end{bmatrix}$$

A matrix  $A \in \mathbb{R}^{n \times n}$  is called diagonalizable if there exist an invertible matrix  $P \in \mathbb{R}^{n \times n}$  such that:

$$A = P^{-1}DP$$

for some diagonal matrix  $D \in \mathbb{R}^{n \times n}$

## Theorem (Eigendecomposition of symmetric matrices):

Let  $\{u_1, u_2, \dots, u_n\}$  be a collection of eigenvectors of symmetric matrix  $A \in \mathbb{R}^{n \times n}$  corresponding to eigenvalues  $\{\lambda_1, \lambda_2, \dots, \lambda_n\}$ ,

i.e.  $Au_i = \lambda_i u_i$  for  $i \in \{1, \dots, n\}$ . Then, matrix  $A$  can be factorized as:

$$A = \begin{bmatrix} | & | & & | \\ u_1 & u_2 & \cdots & u_n \\ | & | & \cdots & | \end{bmatrix} \begin{bmatrix} \lambda_1 & 0 & \cdots & 0 \\ 0 & \lambda_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \lambda_n \end{bmatrix} \begin{bmatrix} | & & & | \\ u_1^\top & & & u_n^\top \\ | & & \cdots & | \\ u_1^\top & & \cdots & u_n^\top \\ | & & \vdots & | \\ \vdots & & \vdots & \vdots \\ | & & & | \end{bmatrix}$$

# Spectral Decomposition of Symmetric Matrix

Let us prove that any symmetric matrix  $A \in \mathbb{R}^{n \times n}$  can be decomposed as:

$$A = \sum_{i=1}^n \lambda_i u_i u_i^\top$$

where  $\{u_1, u_2, \dots, u_n\}$  are eigenvectors of  $A$  corresponding to eigenvalues  $\{\lambda_1, \lambda_2, \dots, \lambda_n\}$

How to show that two equal size matrices  $A = B$ ? We can show that their entries are equal:  $a_{ij} = b_{ij}$  for  $i, j \in \{1, \dots, n\}$

Using **Eigendecomposition Theorem**:

$$A = \begin{bmatrix} | & | & \cdots & | \\ u_1 & u_2 & \cdots & u_n \\ | & | & \cdots & | \end{bmatrix} \begin{bmatrix} \lambda_1 & 0 & \cdots & 0 \\ 0 & \lambda_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \lambda_n \end{bmatrix} \begin{bmatrix} | & | & \cdots & | \\ u_1^\top & u_2^\top & \cdots & u_n^\top \\ | & | & \cdots & | \end{bmatrix}$$

We'll show:  $[U \Lambda U^\top]_{ij} = \left[ \sum_{m=1}^n \lambda_m u_m u_m^\top \right]_{ij} = \sum_{m=1}^n \lambda_m [u_m u_m^\top]_{ij}$    
 for arbitrary  $i, j \in \{1, \dots, n\}$

**Step 1:** Computing  $\left[ \sum_{m=1}^n \lambda_m [u_m u_m^\top]_{ij} \right]$

Notice:  $\begin{bmatrix} | & | & \cdots & | & m=1 & 0 & \cdots & 0 & | & | \\ & 0 & \lambda_2 & \cdots & 0 & | & | & u_1^\top & \cdots & | \\ & & \lambda_2 u_2^\top & \cdots & & | & | & u_2^\top & \cdots & | \\ & & & \ddots & & | & | & u_n^\top & \cdots & | \end{bmatrix} = \lambda \begin{bmatrix} | & | & \cdots & | & [u_m]_1 & [u_m]_2 & \cdots & [u_m]_n \\ & [u_m]_1 & [u_m]_1^2 & [u_m]_1 [u_m]_2 & \cdots & [u_m]_1 [u_m]_n \\ & [u_m]_2 & [u_m]_2^2 & [u_m]_2 [u_m]_1 & \cdots & [u_m]_2 [u_m]_n \\ & \vdots & \vdots & \vdots & \ddots & \vdots \\ & [u_m]_n & [u_m]_n^2 & [u_m]_n [u_m]_1 & \cdots & [u_m]_n [u_m]_n \end{bmatrix}$

$\lambda_m u_m u_m^\top = \lambda \begin{bmatrix} | & | & \cdots & | & [u_m]_1 & [u_m]_2 & \cdots & [u_m]_n \\ & [u_m]_1 & [u_m]_1^2 & [u_m]_1 [u_m]_2 & \cdots & [u_m]_1 [u_m]_n \\ & [u_m]_2 & [u_m]_2^2 & [u_m]_2 [u_m]_1 & \cdots & [u_m]_2 [u_m]_n \\ & \vdots & \vdots & \vdots & \ddots & \vdots \\ & [u_m]_n & [u_m]_n^2 & [u_m]_n [u_m]_1 & \cdots & [u_m]_n [u_m]_n \end{bmatrix}$

Compare:  $[U \Lambda U^\top]_{ij} = \lambda_m [u_m]_i [u_m]_j$

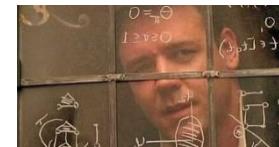
Hence:  $\sum_{m=1}^n \lambda_m [u_m u_m^\top]_{ij} = \sum_{m=1}^n \lambda_m [u_m]_i [u_m]_j$

# Positive Semidefinite (PSD) Matrices

A symmetric matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$  is called positive semidefinite (positive definite) if for any non-zero vector  $\mathbf{x} \in \mathbb{R}^n$ :

$$\mathbf{x}^\top \mathbf{A} \mathbf{x} \geq 0 \quad (\mathbf{x}^\top \mathbf{A} \mathbf{x} > 0) \quad \text{Notation: } \mathbf{A} \succeq \mathbf{0} \quad (\mathbf{A} \succ \mathbf{0})$$

**Exercise:** Matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$  is positive semidefinite (positive definite) if all its eigenvalues are non-negative (positive).



## ML application:

In Gaussian Processes, the symmetric covariance function  $k(\mathbf{x}_1, \mathbf{x}_2)$  defines similarity between function values  $f(\mathbf{x}_1)$  and  $f(\mathbf{x}_2)$

To be proper covariance function, for any collection of points  $\{\mathbf{x}_i\}_{i=1}^N$  their covariance matrix must be PSD:

$$\mathbf{K} = \begin{bmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & k(\mathbf{x}_1, \mathbf{x}_2) & \dots & k(\mathbf{x}_1, \mathbf{x}_n) \\ k(\mathbf{x}_2, \mathbf{x}_1) & k(\mathbf{x}_2, \mathbf{x}_2) & \dots & k(\mathbf{x}_2, \mathbf{x}_n) \\ \vdots & \vdots & \ddots & \vdots \\ k(\mathbf{x}_n, \mathbf{x}_1) & k(\mathbf{x}_n, \mathbf{x}_2) & \dots & k(\mathbf{x}_n, \mathbf{x}_n) \end{bmatrix} \succeq \mathbf{0}$$

Examples of such functions include:

- Square Exponential:  $k_{\text{SE}}(\mathbf{x}_1, \mathbf{x}_2) = \sigma^2 e^{-\frac{\|\mathbf{x}_1 - \mathbf{x}_2\|_2^2}{\ell}}$   $\sigma, \ell \leftarrow$  parameters
- Rational Quadratic:  $k_{\text{RQ}}(\mathbf{x}_1, \mathbf{x}_2) = \left(1 + \frac{\|\mathbf{x}_1 - \mathbf{x}_2\|_2^2}{2\alpha k^2}\right)^{-\alpha}$   $\alpha, k \leftarrow$  parameters

# Cholesky Factorisation

**Theorem (Cholesky Factorisation):**

Any symmetric, positive definite matrix  $A \in \mathbb{R}^{n \times n}$  can be factorized into a product:

$$A = LL^T$$

where  $L$  is lower-triangular matrix (called Cholesky factor):

$$L = \begin{bmatrix} \ell_{11} & 0 & \dots & 0 \\ \ell_{21} & \ell_{22} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ \ell_{n1} & \ell_{n2} & \dots & \ell_{nn} \end{bmatrix}$$



André-Louis Cholesky (1875 - 1918)

**Example:**

$$\begin{bmatrix} 4 & 12 & -16 \\ 12 & 37 & -43 \\ -16 & -43 & 98 \end{bmatrix} = \underbrace{\begin{bmatrix} 2 & 0 & 0 \\ 6 & 1 & 0 \\ -8 & 5 & 3 \end{bmatrix}}_L \underbrace{\begin{bmatrix} 2 & 6 & -8 \\ 0 & 1 & 5 \\ 0 & 0 & 3 \end{bmatrix}}_{L^T}$$

Computational complexity:  $\mathcal{O}(n^3)$

# Matrix Norms

Similar to vectors, we can introduce the norm function for matrices:

Norm of a matrix  $\mathbf{A} \in \mathbb{R}^{m \times n}$  is any univariate function  $\|\cdot\|$ , satisfying norm axioms:

- definiteness:  $\|\mathbf{A}\| = 0 \rightarrow \mathbf{A} = \mathbf{0}$
- absolute homogeneity :  $\|\alpha \mathbf{A}\| = |\alpha| \|\mathbf{A}\|$  for any matrix  $\mathbf{A}$  and scalar  $\alpha$
- subadditivity or triangular inequality:  $\|\mathbf{A} + \mathbf{B}\| \leq \|\mathbf{A}\| + \|\mathbf{B}\|$  for any matrices  $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{m \times n}$
- Non-negativity:  $\|\mathbf{A}\| > 0$  for any matrix  $\mathbf{A}$
- Sub-multiplicativity:  $\|\mathbf{AB}\| \leq \|\mathbf{A}\| \|\mathbf{B}\|$  for any matrices  $\mathbf{A} \in \mathbb{R}^{m \times n}, \mathbf{B} \in \mathbb{R}^{n \times p}$

## ML application:

Non-linear support vector machines is trained using Log-Euclidian kernel:  $k_{\text{LE}}(\mathbf{A}, \mathbf{B}) = \exp\left(-\frac{1}{2\sigma^2} \|\log \mathbf{A} - \log \mathbf{B}\|_F^2\right)$

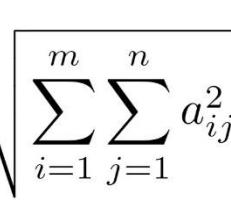
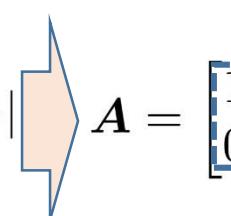
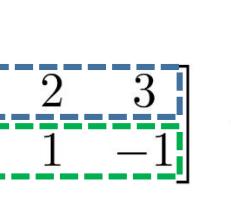
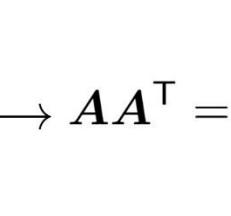
to compare covariance operators  $\mathbf{A}, \mathbf{B}$

Frobenius norm

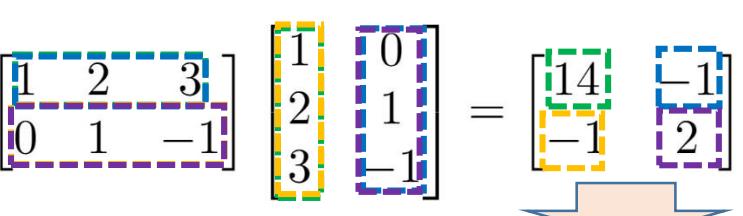
# Matrix Norms (Examples)

Given matrix  $A = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \dots & a_{mn} \end{bmatrix} \in \mathbb{R}^{m \times n}$

## Examples:

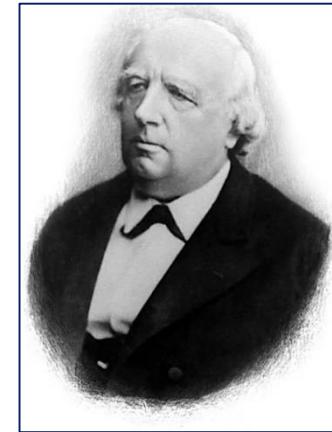
- Frobenius norm:  $\|A\|_F = \sqrt{\sum_{i=1}^m \sum_{j=1}^n a_{ij}^2}$    $A = \begin{bmatrix} 1 & 2 & 3 \\ 0 & 1 & -1 \end{bmatrix} \rightarrow \|A\|_F = \sqrt{1^2 + 2^2 + 3^2 + 0^2 + 1^2 + (-1)^2} = \sqrt{16} = 4$
- Maximum norm:  $\|A\|_1 = \max_{1 \leq j \leq n} \sum_{i=1}^m |a_{ij}|$    $A = \begin{bmatrix} 1 & 2 & 3 \\ 0 & 1 & -1 \end{bmatrix} \rightarrow \|A\|_1 = \max\{1+0, 2+1, 3+1\} = \max\{1, 3, 4\} = 4$
- Infinity norm:  $\|A\|_\infty = \max_{1 \leq i \leq m} \sum_{j=1}^n |a_{ij}|$    $A = \begin{bmatrix} 1 & 2 & 3 \\ 0 & 1 & -1 \end{bmatrix} \rightarrow \|A\|_\infty = \max\{1+2+3, 0+1+1\} = \max\{6, 2\} = 6$
- Spectral norm:  $\|A\|_2 = \sqrt{\lambda_{\max}(AA^\top)}$  

Maximum eigenvalue of  $AA^\top$

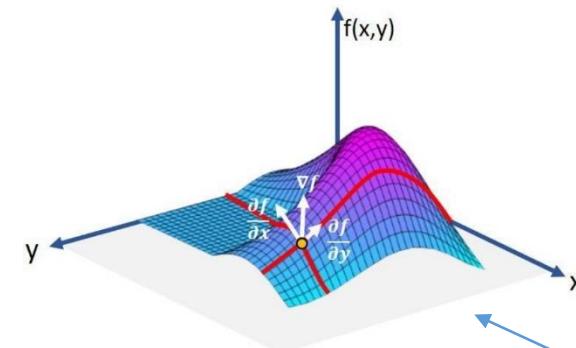
 $A = \begin{bmatrix} 1 & 2 & 3 \\ 0 & 1 & -1 \end{bmatrix} \rightarrow AA^\top = \begin{bmatrix} 1 & 2 & 3 \\ 0 & 1 & -1 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 2 & 1 \\ 3 & -1 \end{bmatrix} = \begin{bmatrix} 14 & -1 \\ -1 & 2 \end{bmatrix}$ 

 $\|A\|_2 = \sqrt{8 + \sqrt{37}}$

# Vector Calculus

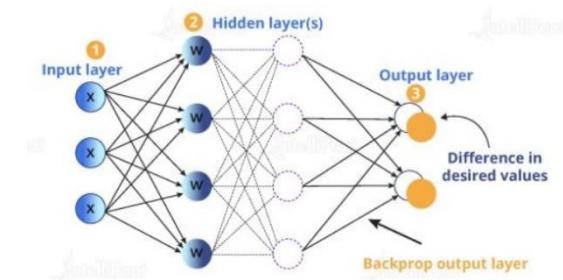
- Derivatives, Differentiation Rules.
- Partial Derivatives, Gradient, Hessian, Taylor Series.
- Jacobian of a smooth mapping.
- Gradients with respect to matrices.
- Backpropagation.



Karl Theodor Wilhelm Weierstrass  
(1815 - 1897)



Courtesy of Brad Moffat



Courtesy of IntelliPaat

# Derivatives, Differentiation Rules.

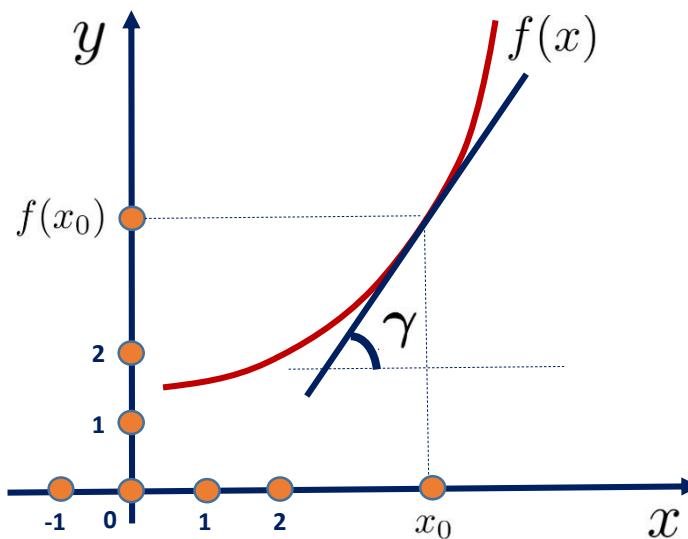
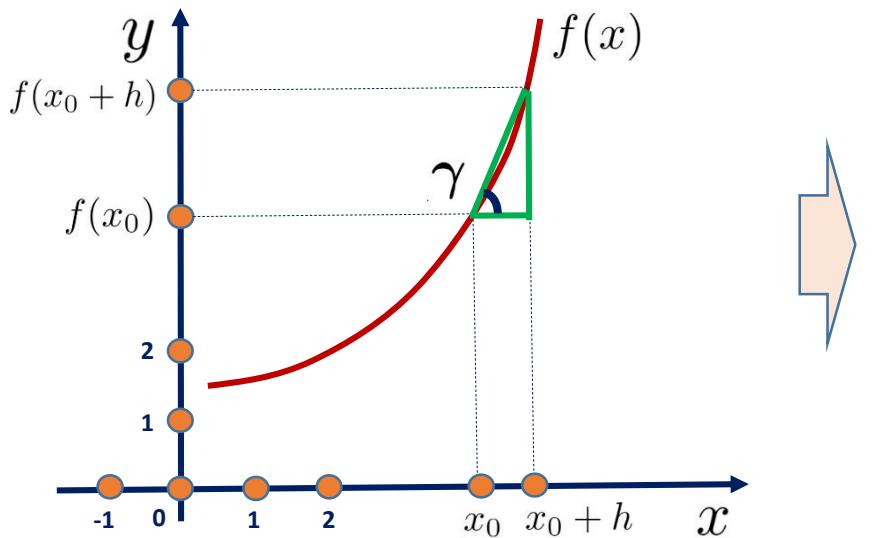
Let  $f : \mathbb{R} \rightarrow \mathbb{R}$  be a univariate function. The derivative of function  $f(x)$  at point  $x_0$  is the following finite limit( if it exists):

$$\lim_{h \rightarrow 0} \frac{f(x_0 + h) - f(x_0)}{h} = f'(x_0)$$

If function has continuous derivative – it is called smooth.

Geometric interpretation:

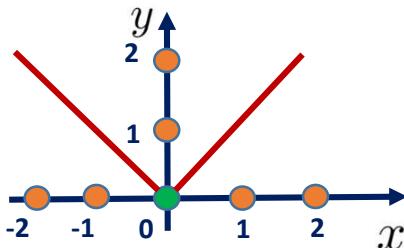
Ratio  $\frac{f(x_0 + h) - f(x_0)}{h} = \tan(\gamma)$



As  $h \rightarrow 0$  this ratio is equal to the slope of the tangent line at point  $x_0$

$$\lim_{h \rightarrow 0} \tan(\gamma) = f'(x_0)$$

Not all functions are differentiable:  $f(x) = |x|$  at  $x_0 = 0$



# Derivatives, Differentiation Rules

To compute derivative of a function we start with known derive expressions and the apply differentiation rules:

## Known Derivative Expressions:

$x' = 1$	$(x^2)' = 2x$	$(\csc x)' = -\cot x \csc x$	$(\arcsin x)' = \frac{1}{\sqrt{1-x^2}}$
$(x^n)' = nx^{n-1}$	$\left(\frac{1}{x}\right)' = -\frac{1}{x^2}$	$(\arccos x)' = -\frac{1}{\sqrt{1-x^2}}$	$(\arctan x)' = \frac{1}{1+x^2}$
$\left(\frac{1}{x^n}\right)' = -\frac{n}{x^{n+1}}$	$(\sqrt{x})' = \frac{1}{2\sqrt{x}}$	$(\text{arccot } x)' = -\frac{1}{1+x^2}$	$(\text{arcsec } x)' = \frac{1}{ x \sqrt{x^2-1}}$
$(\sqrt[m]{x})' = \frac{1}{m\sqrt[m]{x^{m-1}}}$	$(a^x)' = a^x \ln a$	$(\text{arccsc } x)' = -\frac{1}{ x \sqrt{x^2-1}}$	$(\sinh x)' = \cosh x$
$(e^x)' = e^x$	$(\log_a x)' = \frac{1}{x \ln a}$	$(\cosh x)' = \sinh x$	$(\tanh x)' = \operatorname{sech}^2 x$
$(\ln x)' = \frac{1}{x}$	$(\sin x)' = \cos x$	$(\coth x)' = -\operatorname{csch}^2 x$	$(\operatorname{sech} x)' = -\operatorname{sech} x \tanh x$
$(\cos x)' = -\sin x$	$(\tan x)' = \frac{1}{\cos^2 x} = \sec^2 x$	$(\operatorname{csch} x)' = -\operatorname{csch} x \coth x$	$(\operatorname{arcsinh} x)' = \frac{1}{\sqrt{x^2+1}}$
$(\cot x)' = -\frac{1}{\sin^2 x} = -\csc^2 x$	$(\sec x)' = \tan x \sec x$		

## Example:

$$(x^2)' = \lim_{h \rightarrow 0} \frac{(x+h)^2 - x^2}{h} = \lim_{h \rightarrow 0} \frac{x^2 + 2hx + h^2 - x^2}{h} = \lim_{h \rightarrow 0} \frac{2hx + h^2}{h} = \lim_{h \rightarrow 0} (2x + h) = 2x$$

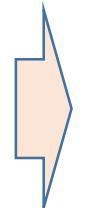
# Derivatives, Differentiation Rules

To compute derivative of a function we start with known derive expressions and the apply differentiation rules:

## Differentiation rules:

If functions  $f(x), g(x)$  are differentiable at point  $x$ , then:

- Product rule:  $(f(x)g(x))' = f'(x)g(x) + f(x)g'(x)$



$$(x^2 \cos x)' = (x^2)' \cos x + x^2(\cos x)' = 2x \cos x - x^2 \sin x$$

**Using:**  $(\cos x)' = -\sin x$  and  $(x^2)' = 2x$

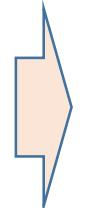
- Quotient rule:  $\left(\frac{f(x)}{g(x)}\right)' = \frac{f'(x)g(x) - f(x)g'(x)}{g^2(x)}$



$$\left(\frac{x^2}{\cos x}\right)' = \frac{(x^2)' \cos x - x^2(\cos x)'}{\cos^2 x} = \frac{2x \cos x + x^2 \sin x}{\cos^2 x}$$

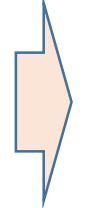
- Sum rule:  $(\alpha f(x) + \beta g(x))' = \alpha f'(x) + \beta g'(x)$

**for any scalars:**  $\alpha, \beta \in \mathbb{R}$



$$(4x^2 - 2 \cos x)' = 4(x^2)' - 2(\cos x)' = 8x + 2 \sin x$$

- Chain rule:  $(f[g(x)])' = f'[g(x)]g'(x)$



$$(\cos x^2)' = -\sin x^2 (x^2)' = -\sin x^2 (2x) = -2x \sin x^2$$

**Using:**  $f(z) = \cos z$  and  $g(x) = x^2$

## Examples:

# Partial Derivatives

Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  be a multivariate function. The partial derivative of function  $f(\mathbf{x})$  at point  $\mathbf{x}_0 = [[\mathbf{x}_0]_1, \dots, [\mathbf{x}_0]_n]^T$  are the following finite limits ( if they exist):

$$\frac{\partial f(\mathbf{x}_0)}{\partial x_1} = \lim_{h_1 \rightarrow 0} \frac{f([\mathbf{x}_0]_1 + h_1, [\mathbf{x}_0]_2, \dots, [\mathbf{x}_0]_n) - f([\mathbf{x}_0]_1, [\mathbf{x}_0]_2, \dots, [\mathbf{x}_0]_n)}{h_1}$$

$$\frac{\partial f(\mathbf{x}_0)}{\partial x_2} = \lim_{h_2 \rightarrow 0} \frac{f([\mathbf{x}_0]_1, [\mathbf{x}_0]_2 + h_2, \dots, [\mathbf{x}_0]_n) - f([\mathbf{x}_0]_1, [\mathbf{x}_0]_2, \dots, [\mathbf{x}_0]_n)}{h_2}$$

$$\vdots$$

$$\frac{\partial f(\mathbf{x}_0)}{\partial x_n} = \lim_{h_n \rightarrow 0} \frac{f([\mathbf{x}_0]_1, [\mathbf{x}_0]_2, \dots, [\mathbf{x}_0]_n + h_n) - f([\mathbf{x}_0]_1, [\mathbf{x}_0]_2, \dots, [\mathbf{x}_0]_n)}{h_n}$$



## Example:

For  $f(x_1, x_2, x_3) = x_1^2 + x_1 x_2 + \frac{x_3^3}{x_2}$  compute

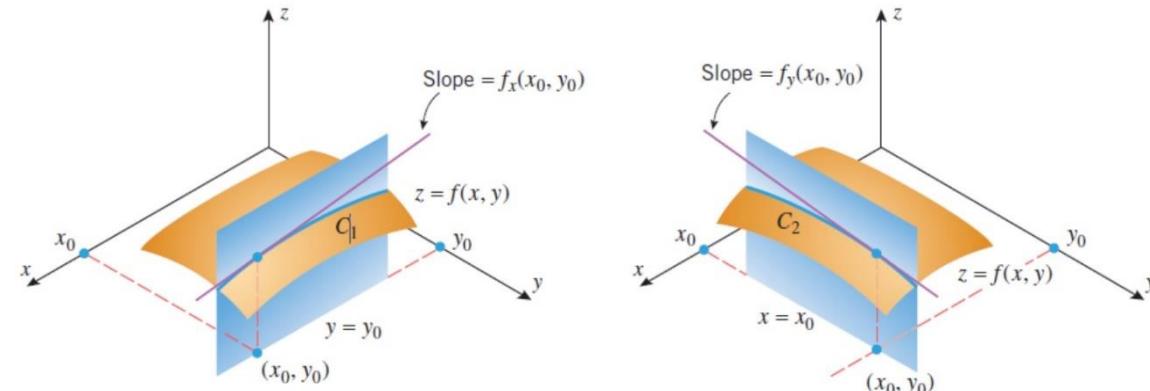
partial derivatives at point  $\mathbf{x}_0 = [1, 3, 5]^T$

$$\frac{\partial f(\mathbf{x})}{\partial x_1} = 2x_1 + x_2 \implies \frac{\partial f(\mathbf{x}_0)}{\partial x_1} = 2[\mathbf{x}_0]_1 + [\mathbf{x}_0]_2 = 2 \cdot 1 + 3 = 5$$

$$\frac{\partial f(\mathbf{x})}{\partial x_2} = x_1 - \frac{x_3^3}{x_2^2} \implies \frac{\partial f(\mathbf{x}_0)}{\partial x_2} = [\mathbf{x}_0]_1 - \frac{[\mathbf{x}_0]_3^3}{[\mathbf{x}_0]_2^2} = 1 - \frac{5^3}{3^2} = -\frac{116}{9}$$

$$\frac{\partial f(\mathbf{x})}{\partial x_3} = \frac{3x_3^2}{x_2} \implies \frac{\partial f(\mathbf{x}_0)}{\partial x_3} = \frac{3[\mathbf{x}_0]_3^2}{[\mathbf{x}_0]_2} = 3 \cdot \frac{5^2}{3} = 25$$

## Geometric interpretation:

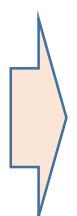


Partial derivatives define slopes of the tangent line to the curve defined by the intersection of surface of the function and the hyperplane

# Gradient

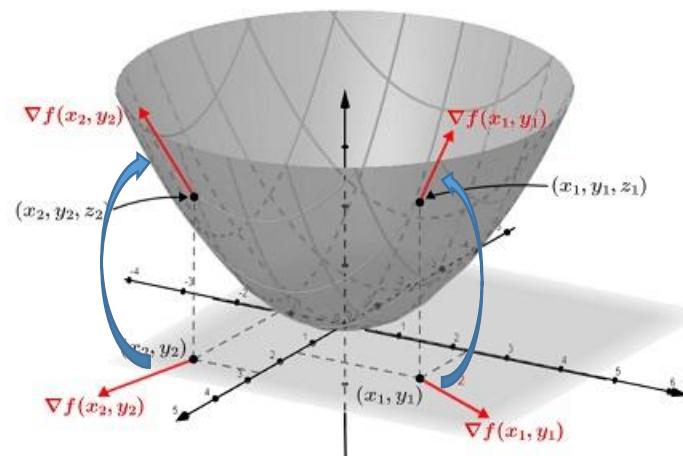
Gradient of function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  at point  $\mathbf{x}_0 = [[\mathbf{x}_0]_1, \dots, [\mathbf{x}_0]_n]^T$  is a vector of all partial derivatives:

$$\frac{df(\mathbf{x}_0)}{d\mathbf{x}} = \left[ \frac{\partial f(\mathbf{x}_0)}{\partial x_1} \quad \frac{\partial f(\mathbf{x}_0)}{\partial x_2} \quad \dots \quad \frac{\partial f(\mathbf{x}_0)}{\partial x_n} \right]^T \in \mathbb{R}^n$$

**Example:** For  $f(x_1, x_2, x_3) = x_1^2 + x_1x_2 + \frac{x_3^3}{x_2}$  compute the gradient at point  $\mathbf{x}_0 = [1, 3, 5]^T$  

$$\frac{df(\mathbf{x}_0)}{d\mathbf{x}} = \begin{bmatrix} \frac{\partial f(\mathbf{x}_0)}{\partial x_1} \\ \frac{\partial f(\mathbf{x}_0)}{\partial x_2} \\ \frac{\partial f(\mathbf{x}_0)}{\partial x_3} \end{bmatrix} = \begin{bmatrix} 5 \\ -\frac{116}{9} \\ 25 \end{bmatrix}$$

## Geometric interpretation:



Consider paraboloid  $z = f(x, y) = x^2 + y^2$ , then gradient vectors  $\nabla f(x_1, y_1)$  and  $\nabla f(x_2, y_2)$  drawn in the  $xy$ -plane have their initial point placed at  $(x_1, y_1)$  and  $(x_2, y_2)$  respectively.

Note that the gradient vectors are calculated in component form but are translated to their respective input points on the  $xy$ -plane to better show the direction associated with those points.

If we fold the  $xy$ -plane so that  $(x, y) \rightarrow (x, y, x^2 + y^2)$ , then each corresponding gradient vector can also be mapped onto the paraboloid.

The gradient vectors mapped to  $(x_1, y_1, z_1)$  and  $(x_2, y_2, z_2)$  show the direction of fastest increase.

# Gradient (ML Application)

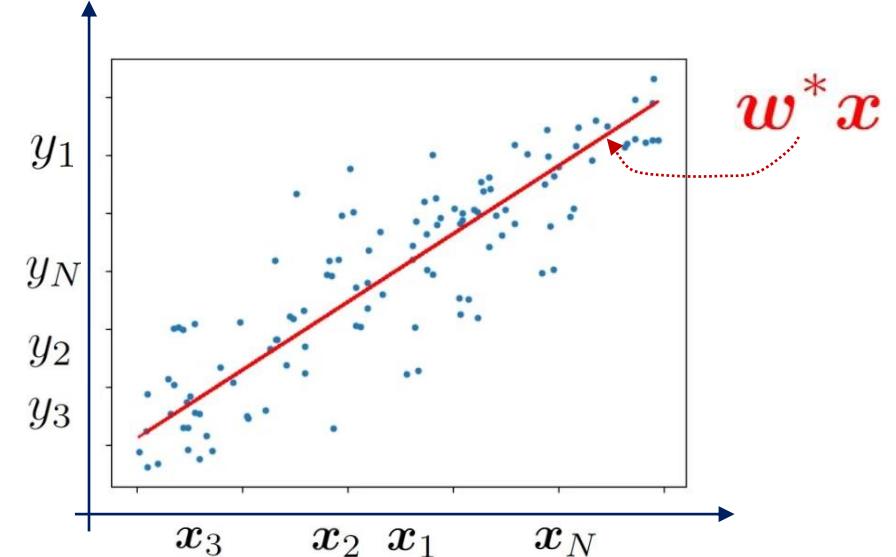
In linear regression task, we are given a data-set  $\mathcal{D} = \{\mathbf{x}_i, y_i\}_{i=1}^N$ .

Our goal is to find ML model  $\mathbf{w}^\top \mathbf{x}$  that fits this data.



To find the optimal model parameter we aim optimization problem:

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \mathcal{L}_{\text{loss}}(\mathbf{w}) = \sum_{i=1}^N (\mathbf{w}^\top \mathbf{x}_i - y_i)^2$$



At optimal point gradient vanishes:  $\nabla_{\mathbf{w}} \mathcal{L}_{\text{loss}}(\mathbf{w}^*) = \mathbf{0}$

Let us compute this gradient:

$$\begin{aligned} \nabla_{\mathbf{w}} \mathcal{L}_{\text{loss}}(\mathbf{w}) &= \begin{bmatrix} \frac{\partial \mathcal{L}_{\text{loss}}(\mathbf{w})}{\partial w_1} \\ \frac{\partial \mathcal{L}_{\text{loss}}(\mathbf{w})}{\partial w_2} \\ \vdots \\ \frac{\partial \mathcal{L}_{\text{loss}}(\mathbf{w})}{\partial w_n} \end{bmatrix} \stackrel{\text{fix arbitrary } c}{=} 2 \sum_{i=1}^N \mathbf{x}_{\text{data}}^\top \mathbf{x}_{\text{data}} \mathbf{w} - 2 \sum_{i=1}^N \mathbf{x}_{\text{data}}^\top \mathbf{y}_{\text{data}} \\ &= 2 \sum_{i=1}^N \mathbf{x}_{\text{data}}^\top \mathbf{x}_{\text{data}} \mathbf{w} - \sum_{i=1}^N \mathbf{x}_{\text{data}}^\top \mathbf{y}_{\text{data}} \sum_{i=1}^N y_i \mathbf{x}_i \cdots + w_n [\mathbf{x}_i]_n \end{aligned}$$

using sum rule  
↓ re-group terms  
↓

$$\nabla_{\mathbf{w}} \mathcal{L}_{\text{loss}}(\mathbf{w}^*) = \mathbf{0} \quad \boxed{X_{\text{data}}^\top X_{\text{data}} \mathbf{w}^* = X_{\text{data}}^\top \mathbf{y}_{\text{data}}}$$

Denote:  
 $X_{\text{data}} = \begin{bmatrix} \mathbf{x}_1^\top & \mathbf{x}_2^\top & \cdots & \mathbf{x}_N^\top \end{bmatrix}^\top$ ,  $\mathbf{y}_{\text{data}} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}$

$$y_{\text{data}} = \sum_{i=1}^N y_i \mathbf{x}_i = 2 \sum_{i=1}^N (c_i - y_i) [\mathbf{x}_i]_j$$

$$\sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^\top = X_{\text{data}}^\top X_{\text{data}}$$

$$\sum_{i=1}^N y_i \mathbf{x}_i = X_{\text{data}}^\top \mathbf{y}_{\text{data}}$$

Why?



# Hessian

Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  be a multivariate function. Each partial derivative  $\frac{\partial f(\mathbf{x})}{\partial x_j}$  is itself multivariate function:  $\frac{\partial f(\mathbf{x})}{\partial x_j} : \mathbb{R}^n \rightarrow \mathbb{R}$

Hence, if it is differentiable, we can study partial derivatives of  $\frac{\partial f(\mathbf{x})}{\partial x_j}$  itself:

$$\nabla_{\mathbf{x}} \frac{\partial f(\mathbf{x})}{\partial x_j} = \begin{bmatrix} \frac{\partial^2 f(\mathbf{x})}{\partial x_j \partial x_1} \\ \frac{\partial^2 f(\mathbf{x})}{\partial x_j \partial x_2} \\ \vdots \\ \frac{\partial^2 f(\mathbf{x})}{\partial x_j^2} \\ \vdots \\ \frac{\partial^2 f(\mathbf{x})}{\partial x_j \partial x_n} \end{bmatrix} \quad \text{combining into one matrix:}$$

$$\nabla_{\mathbf{x}}^2 f(\mathbf{x}) = \begin{bmatrix} \frac{\partial^2 f(\mathbf{x})}{\partial x_1^2} & \frac{\partial^2 f(\mathbf{x})}{\partial x_2 \partial x_1} & \frac{\partial^2 f(\mathbf{x})}{\partial x_3 \partial x_1} & \cdots & \frac{\partial^2 f(\mathbf{x})}{\partial x_n \partial x_1} \\ \frac{\partial^2 f(\mathbf{x})}{\partial x_1 \partial x_2} & \frac{\partial^2 f(\mathbf{x})}{\partial x_2^2} & \frac{\partial^2 f(\mathbf{x})}{\partial x_3 \partial x_2} & \cdots & \frac{\partial^2 f(\mathbf{x})}{\partial x_n \partial x_2} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f(\mathbf{x})}{\partial x_1 \partial x_n} & \frac{\partial^2 f(\mathbf{x})}{\partial x_2 \partial x_n} & \frac{\partial^2 f(\mathbf{x})}{\partial x_3 \partial x_n} & \cdots & \frac{\partial^2 f(\mathbf{x})}{\partial x_n^2} \end{bmatrix} \in \mathbb{R}^{n \times n}$$

Example: For  $f(x_1, x_2, x_3) = x_1^2 + x_1x_2 + \frac{x_3^3}{x_2}$  compute Hessian matrix at point  $x_0 = [1, 3, 5]^T$

Method of columns:

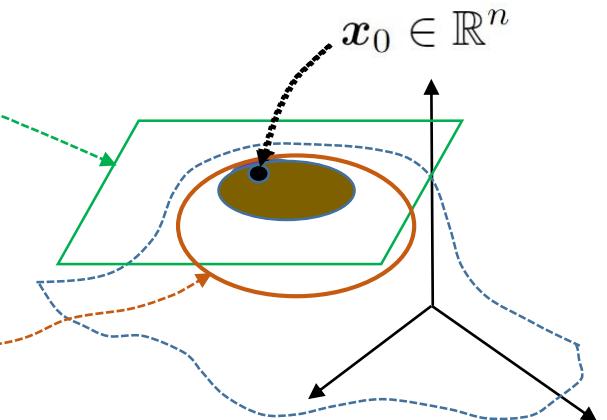
$$\nabla_{\mathbf{x}}^2 f(\mathbf{x}) = \begin{bmatrix} 2 & 1 & 0 \\ 1 & \frac{2x_3^3}{x_2^3} & -\frac{3x_3^2}{x_2^2} \\ 0 & -\frac{3x_3^2}{x_2^2} & -\frac{6x_3}{x_2} \end{bmatrix} \begin{bmatrix} \frac{\partial}{\partial x_1} \left( \frac{3x_3^2}{x_2} \right) \\ \frac{\partial}{\partial x_2} \left( \frac{3x_3^2}{x_2} \right) \\ \frac{\partial}{\partial x_3} \left( \frac{3x_3^2}{x_2} \right) \end{bmatrix} \nabla_{\mathbf{x}}^2 f(\mathbf{x}_0) = \begin{bmatrix} 2 & 1 & 0 \\ 1 & \frac{2[x_0]_3^3}{[x_0]_2^3} & -\frac{3[x_0]_3^2}{[x_0]_2^2} \\ 0 & -\frac{3[x_0]_3^2}{[x_0]_2^2} & -\frac{6[x_0]_3}{[x_0]_2} \end{bmatrix} = \begin{bmatrix} 2 & 1 & 0 \\ 1 & \frac{2 \cdot 125}{3^3} & -\frac{3 \cdot 5^2}{3^2} \\ 0 & -\frac{3 \cdot 5^2}{3^2} & -\frac{6 \cdot 5}{3} \end{bmatrix} = \begin{bmatrix} 2 & 1 & 0 \\ 1 & \frac{250}{27} & -\frac{25}{3} \\ 0 & -\frac{25}{3} & -10 \end{bmatrix}$$

# Hessian (Geometric Interpretation). Taylor Series

Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  be a multivariate function, at consider the reference point  $x_0 \in \mathbb{R}^n$

Tangent plane  $= f(x_0) + \nabla_x^\top f(x_0)(x - x_0)$

Locally around point  $x_0 \in \mathbb{R}^n$  tangent plane gives linear approximation of the function:



Tangent second order surface  $= f(x_0) + \nabla_x^\top f(x_0)(x - x_0) + \frac{1}{2}(x - x_0)^\top \nabla_x^2 f(x_0)(x - x_0)$

Locally around point  $x_0 \in \mathbb{R}^n$  tangent second order surface gives quadratic approximation of the function:

## TAYLOR SERIES:

Similarly, adding higher order derivatives we can approximate "good" functions at any point:

$$f(x) = f(x_0) + \nabla_x^\top f(x_0)(x - x_0) + \frac{1}{2}(x - x_0)^\top \nabla_x^2 f(x_0)(x - x_0) + \dots + \frac{D^k f(x_0)}{k!}(x - x_0)^k + \dots$$

High-order  
Differentiation  
operator

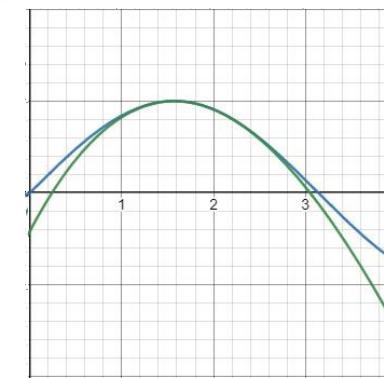
How it works: linear approximation      quadratic approximation

$$f(x) = \sin x$$

$$x_0 = 1.8$$

$$g_3(x) = \sin 1.8 + (\sin 1.8)'(x - 1.8) + \frac{1}{2}(x - 1.8)^2(\sin 1.8)'' + \frac{1}{3!}(x - 1.8)^3(\sin 1.8)'''$$

High order approximation  
seventh order  
approximation



Lets have 10 min break



# Gradient Extensions (I)

Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$  be a multivariate vector function:

$$f(\mathbf{x}) = \begin{bmatrix} f_1(\mathbf{x}) \\ f_2(\mathbf{x}) \\ \vdots \\ f_m(\mathbf{x}) \end{bmatrix}$$

each component of  $f_j(\mathbf{x})$  is a multivariate function itself:  $f_j : \mathbb{R}^n \rightarrow \mathbb{R}$

$\in \mathbb{R}^{m \times n}$

$$\left[ \begin{array}{cccc} \frac{\partial f_1(\mathbf{x})}{\partial x_1} & \frac{\partial f_1(\mathbf{x})}{\partial x_2} & \cdots & \frac{\partial f_1(\mathbf{x})}{\partial x_n} \\ \frac{\partial f_2(\mathbf{x})}{\partial x_1} & \frac{\partial f_2(\mathbf{x})}{\partial x_2} & \cdots & \frac{\partial f_2(\mathbf{x})}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_m(\mathbf{x})}{\partial x_1} & \frac{\partial f_m(\mathbf{x})}{\partial x_2} & \cdots & \frac{\partial f_m(\mathbf{x})}{\partial x_n} \end{array} \right]$$

**Jacobian**

We can define the notion of gradient for  $f(\mathbf{x})$ :  $\nabla_{\mathbf{x}} f(\mathbf{x}) =$

**Example:**  $f(x_1, x_2, x_3) = \begin{bmatrix} x_1^2 + x_2 x_3 \\ x_3 - x_1^2 x_2 \end{bmatrix}$   $\mathbf{x}_0 = [2, -1, 0]^\top$

$$\nabla_{\mathbf{x}} f(\mathbf{x}) = \begin{bmatrix} \frac{\partial}{\partial x_1}(x_1^2 + x_2 x_3) & \frac{\partial}{\partial x_2}(x_1^2 + x_2 x_3) & \frac{\partial}{\partial x_3}(x_1^2 + x_2 x_3) \\ \frac{\partial}{\partial x_1}(x_3 - x_1^2 x_2) & \frac{\partial}{\partial x_2}(x_3 - x_1^2 x_2) & \frac{\partial}{\partial x_3}(x_3 - x_1^2 x_2) \end{bmatrix} = \begin{bmatrix} 2x_1 & x_3 & x_2 \\ -2x_1 x_2 & -x_1^2 & 1 \end{bmatrix}$$



$$\nabla_{\mathbf{x}} f(\mathbf{x}_0) = \begin{bmatrix} 2 \cdot 2 & 0 & -1 \\ -2 \cdot (2 \cdot (-1)) & -(2)^2 & 1 \end{bmatrix} = \begin{bmatrix} 4 & 0 & -1 \\ 4 & -4 & 1 \end{bmatrix}$$

# Gradient Extensions (II)

Let  $f : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}$  be a function of a matrix  $\mathbf{X} \in \mathbb{R}^{m \times n}$ :

$$\mathbf{X} = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1n} \\ x_{21} & x_{22} & \dots & x_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ x_{m1} & x_{m2} & \dots & x_{mn} \end{bmatrix}$$

The gradient of  $f(\mathbf{X})$   
with respect to matrix :

$$\nabla_{\mathbf{X}} f(\mathbf{X}) = \begin{bmatrix} \frac{\partial f(\mathbf{X})}{\partial x_{11}} & \frac{\partial f(\mathbf{X})}{\partial x_{12}} & \dots & \frac{\partial f(\mathbf{X})}{\partial x_{1n}} \\ \frac{\partial f(\mathbf{X})}{\partial x_{21}} & \frac{\partial f(\mathbf{X})}{\partial x_{22}} & \dots & \frac{\partial f(\mathbf{X})}{\partial x_{2n}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f(\mathbf{X})}{\partial x_{m1}} & \frac{\partial f(\mathbf{X})}{\partial x_{m2}} & \dots & \frac{\partial f(\mathbf{X})}{\partial x_{mn}} \end{bmatrix} \in \mathbb{R}^{m \times n}$$

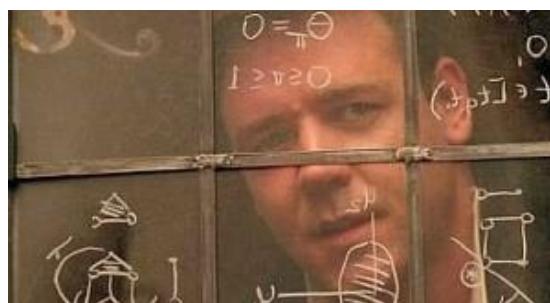
Example:  $f(\mathbf{X}) = \text{Tr}(\mathbf{X}^T \mathbf{X}), \quad \mathbf{X}_0 = \begin{bmatrix} 1 & 0 & -1 \\ 2 & -4 & 5 \end{bmatrix}$

Remember function carefully:

$$\nabla_{\mathbf{X}} f(\mathbf{X}) = \begin{bmatrix} 2x_{11} & 2x_{12} & \dots & 2x_{1n} \\ 2x_{21} & 2x_{22} & \dots & 2x_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ 2x_{m1} & 2x_{m2} & \dots & 2x_{mn} \end{bmatrix} = 2\mathbf{X} \quad \nabla_{\mathbf{X}} f(\mathbf{X}_0) = 2\mathbf{X}_0 = \begin{bmatrix} 2 & 0 & -2 \\ 4 & -8 & 10 \end{bmatrix}$$

Exercise:  $\mathbf{X}^T \quad \mathbf{X}$

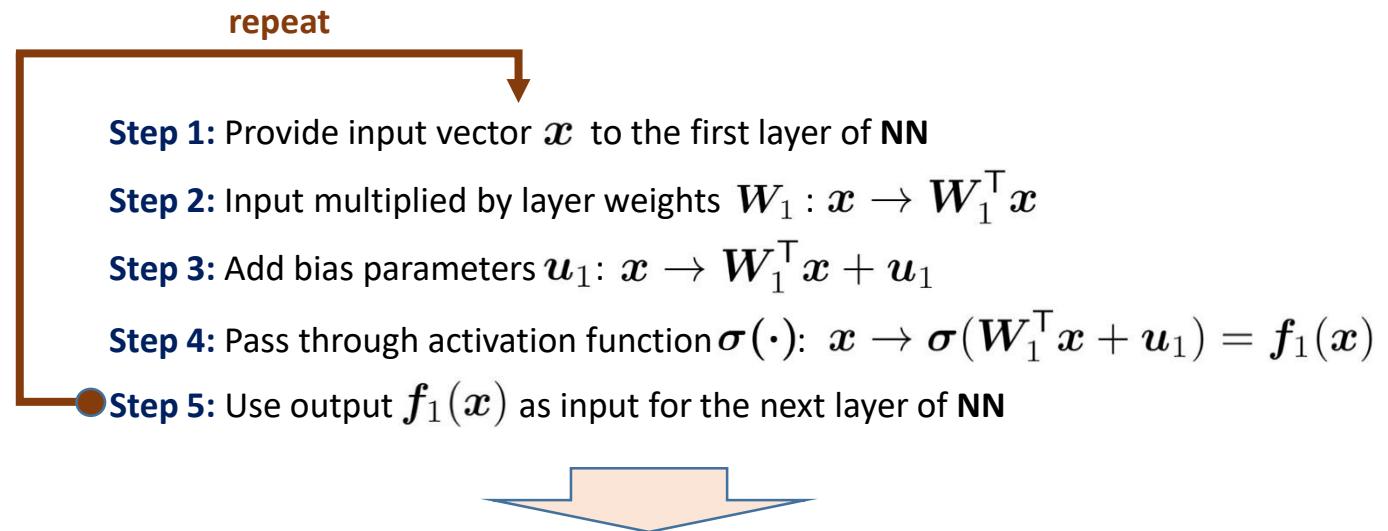
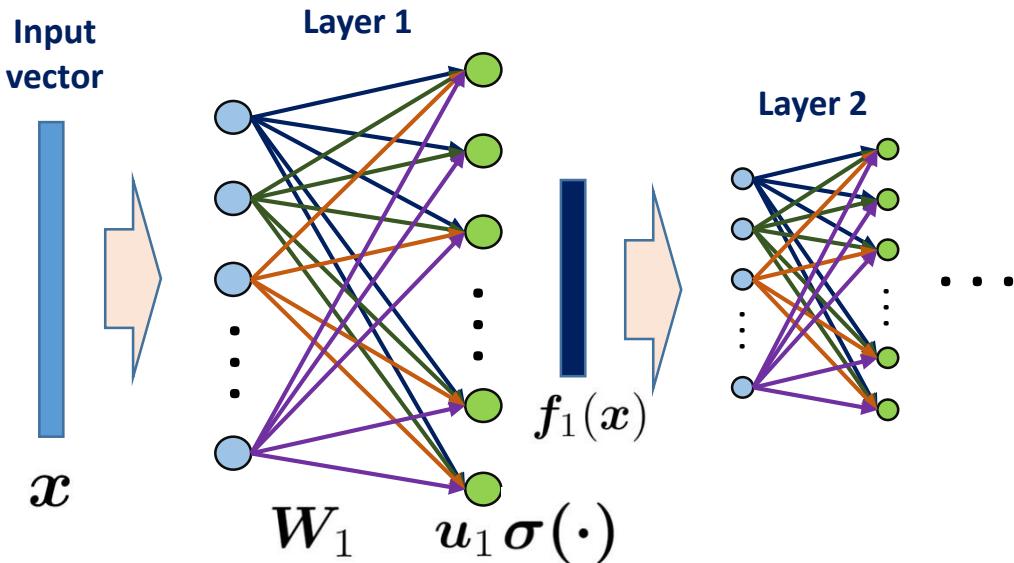
Show:  $\nabla_{\mathbf{X}} \mathbf{a}^T \mathbf{X} \mathbf{b} = \mathbf{a} \mathbf{b}^T, \quad \mathbf{a} \in \mathbb{R}^m, \mathbf{b} \in \mathbb{R}^n$



# ML Application (Backpropagation I)

For training deep neural network models, the backpropagation algorithm is an efficient way to compute the gradient of an error function with respect to the parameters of the model.

Consider feedforward computation of NN:



The computation of  $K$ -layer NN with weights  $\{W_1, W_2, \dots, W_K\}$  and biases  $\{u_1, u_2, \dots, u_K\}$  allows iterative form:

$$f_0 = x$$

$$f_i = \sigma (W_{i-1}^T f_{i-1} + u_{i-1}), \quad i = 1, \dots, K$$

Group unknown parameter associated with each layer of NN:  $\Theta_1 = \{W_1, u_1\}, \Theta_2 = \{W_2, u_2\}, \dots, \Theta_K = \{W_K, u_K\}$

or more compactly:

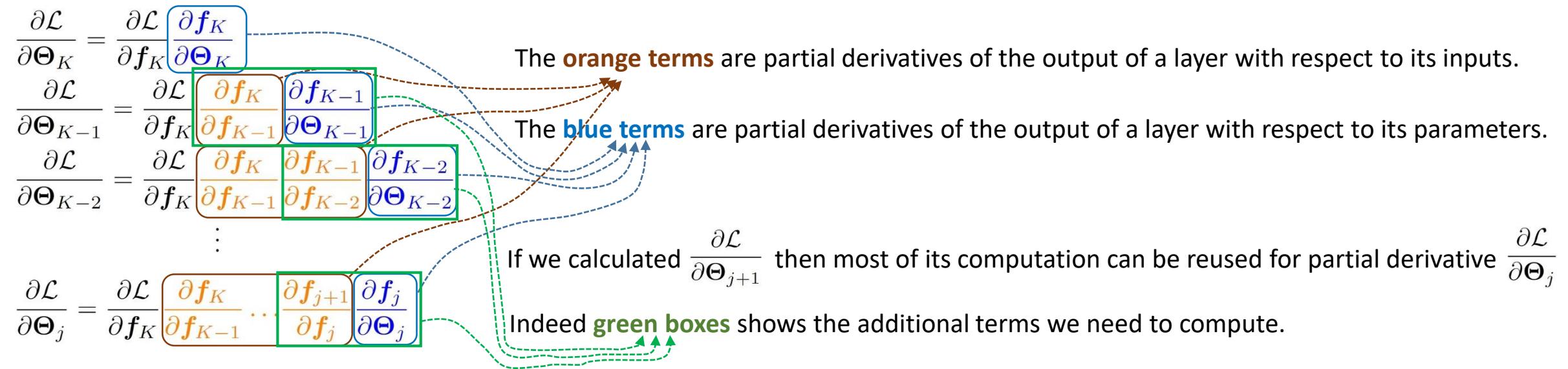
$$\Theta = \{\Theta_1, \Theta_2, \dots, \Theta_K\}$$

# ML Application (Backpropagation II)

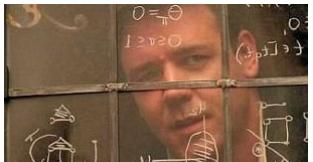
Training **NN** constitutes in finding parameters  $\Theta$  which best fits observed data  $\mathcal{D} = \{\mathbf{X}_{\text{data}}, \mathbf{y}_{\text{data}}\}$ .

The fitness measure is defined via loss function:  $\mathcal{L}(\Theta) = \|\mathbf{y}_{\text{data}} - \mathbf{f}_K(\Theta, \mathbf{X}_{\text{data}})\|^2$

To obtain the gradients with respect to the parameter  $\Theta$ , we compute the partial derivatives of  $\mathcal{L}(\Theta)$  with respect to the parameters  $\Theta_j = \{\mathbf{W}_j, \mathbf{u}_j\}$  of each layer  $j = 1, \dots, K$ . The chain rule allows us to determine the partial derivatives as:



## Exercise:

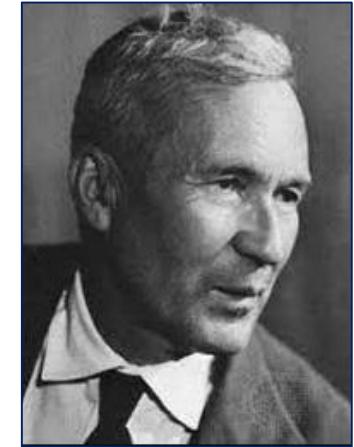


Using derivatives of vector function,  
calculate:  $\frac{\partial \mathbf{f}_{j+1}}{\partial \mathbf{f}_j}$  and  $\frac{\partial \mathbf{f}_j}{\partial \Theta_j}$

$$f_0 = \mathbf{x}$$
$$f_i = \sigma (\mathbf{W}_{i-1}^\top \mathbf{f}_{i-1} + \mathbf{u}_{i-1}), \quad i = 1, \dots, K$$

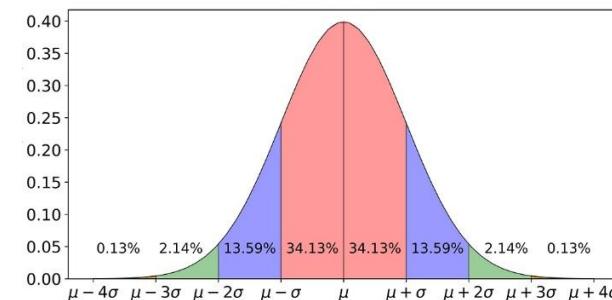
# Probability Theory

- Probability and Random Variables.
- Discrete and Continuous Probabilities.
- Sum Rule, Sum Product, Bayes Theorem.
- Summary Statistics, Independence.
- Gaussian Distribution.



Andrey Nikolaevich Kolmogorov  
(1903-1987)

$$P(A|B) = \frac{P(B|A) P(A)}{P(B)}$$



Source: Google Images

# Probability Space

The probability space models a real-world process (referred to as an experiment) with random outcomes. It consists of three components

## THE SAMPLE SPACE $\Omega$

The sample space is the set of all possible elementary outcomes of the experiment.



### Example:

Coin toss:  $\Omega = \{ \begin{array}{c} \text{Head} \\ \text{Tail} \end{array} \}$

Roll a dice:  $\Omega = \{ \begin{array}{c} \text{1 dot} \\ \text{2 dots} \\ \text{3 dots} \\ \text{4 dots} \\ \text{5 dots} \\ \text{6 dots} \end{array} \}$

## THE EVENT SPACE $\mathcal{A}$

The event space is the space of potential results of the experiment.

Each event  $A \in \mathcal{A}$  consists of elementary outcomes, i.e.  $A \subseteq \Omega$

The set of all events is a power set of  $\Omega$ , and denotes:  $\mathcal{A} = 2^\Omega$



Event of getting TAIL in a coin-flip:  $A = \text{Tail}$

Event of getting even number in a roll dice:  $B = \{ \text{2 dots}, \text{4 dots}, \text{6 dots} \}$

## THE PROBABILITY $\mathbb{P}$

With each event  $A \in \mathcal{A}$ , we associate a number  $\mathbb{P}(A) \in [0, 1]$  that measures the probability or degree of belief that the event will occur.  $\mathbb{P}(A)$  is called the probability of  $A$



Probability of getting TAIL in a coin-flip:  $\mathbb{P}(A) = \frac{1}{2}$

Probability of getting even number in a roll dice:  $\mathbb{P}(B) = \frac{1}{2}$

**Remark:** chance of ... = the probability of ...

The chance of ... is 33% = The probability of ... is 0.33

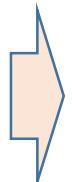
# Random Variable

In machine learning, we often don't refer to the probability space, but instead refer to probabilities on quantities of interest denoted by  $\mathcal{T}$

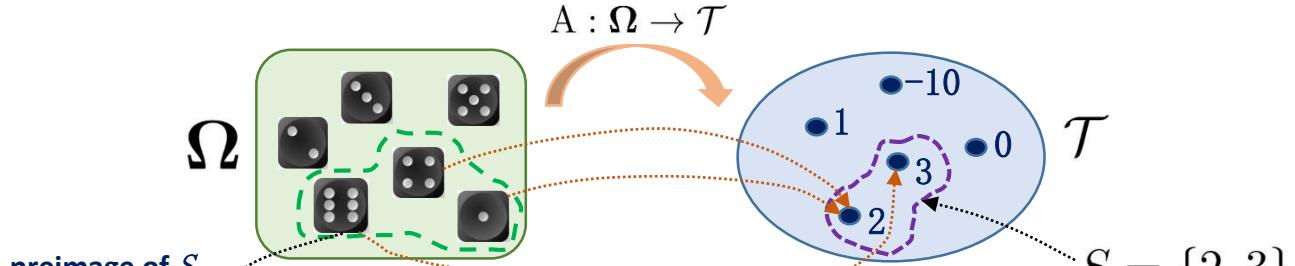
A random variable  $A : \Omega \rightarrow \mathcal{T}$  is a mapping that takes an element  $\omega \in \Omega$  and returns a an element  $x \in \mathcal{T}$ .

## Example:

Random variable defined on a set of roll dices:


$$\begin{array}{lll} A(\bullet) = 2 & A(\circ\bullet) = 1 & A(\circ\circ) = -10 \\ A(\circ\circ\bullet) = 2 & A(\circ\circ\circ) = 0 & A(\circ\circ\circ\bullet) = 3 \end{array}$$

In this case:  $\mathcal{T} = \{2, 3, -10, 0, 1\}$



$$\begin{aligned} \mathbb{P}_A(S) &= \mathbb{P}(A^{-1}(S)) = \mathbb{P}(\bullet\bullet\bullet, \bullet\bullet\circ, \circ\bullet\bullet) = \mathbb{P}(\bullet\bullet\bullet) + \mathbb{P}(\bullet\bullet\circ) + \mathbb{P}(\circ\bullet\bullet) = \\ &= \frac{1}{6} + \frac{1}{6} + \frac{1}{6} = \frac{1}{2} \end{aligned}$$

Random variable  $A$  defines probability distribution  $\mathbb{P}_A(S) \in [0, 1]$  on the target space  $\mathcal{T}$ .

# Discrete Probabilities

Consider the case when we can enumerate all elements in the target set  $\mathcal{T}$ .

In other words,  $\mathcal{T} = \{x_1, x_2, \dots, x_k, \dots\}$

Random variable  $A$  defines a probability distribution on  $\mathcal{T}$ , such that for any  $x_i \in \mathcal{T}$  its probability  $\mathbb{P}_A(x_i) = \mathbb{P}(A = x_i) = p_i$ . Hence, discrete random variable is characterized by a probability mass function  $\mathbb{P}_A() = \{p_1, p_2, \dots\}$ , such that  $p_i \geq 0$  and  $\sum_i p_i = 1$

$A = \text{Bernoulli}(p)$

Consider  $\mathcal{T} = \{0, 1\}$  and :

$$\mathbb{P}_A(1) = \mathbb{P}(A = 1) = p, \quad \mathbb{P}_A(0) = \mathbb{P}(A = 0) = 1 - p$$

$A = \text{Poisson}(\lambda)$

Consider  $\mathcal{T} = \{0, 1, 2, \dots\}$  and :

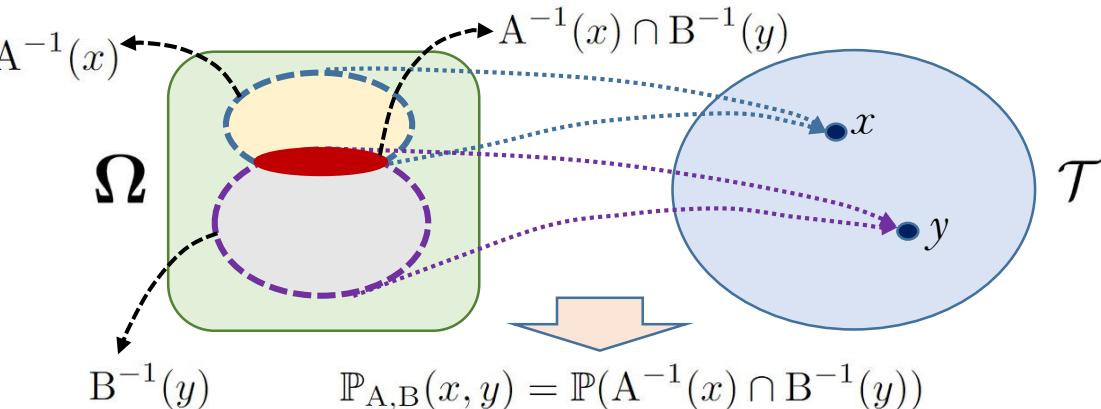
$$\mathbb{P}_A(k) = \mathbb{P}(A = k) = \frac{\lambda^k e^{-\lambda}}{k!} \quad \text{with } \lambda \in (0, +\infty)$$

## JOINT PROBABILITY FOR RANDOM VARIABLES:

Discrete random variables:  $A : \Omega \rightarrow \mathcal{T}$  and  $B : \Omega \rightarrow \mathcal{T}$ , and  $x, y \in \mathcal{T}$ .

Joint Probability of  $A = x$  and  $B = y$  is defined as:

$\mathbb{P}_{A,B}(x, y) = \mathbb{P}(\omega \in \Omega, \text{ such that } A(\omega) = x \text{ and } B(\omega) = y)$



$$A(\bullet) = 2 \quad A(\square) = 1$$

$$A(\blacksquare) = 2 \quad A(\blacksquare\square) = -10 \quad \text{and}$$

$$A(\blacksquare\square\square) = 0 \quad A(\blacksquare\square\square\square) = 3$$

$$B(\bullet) = 1 \quad B(\square) = 2 \quad x = 2$$

$$B(\blacksquare) = 2 \quad B(\blacksquare\square) = 3, \quad y = 1$$

$$B(\blacksquare\square\square) = 0 \quad B(\blacksquare\square\square\square) = 3$$

$\downarrow$   $A^{-1}(2) = \{\bullet, \square\}$  and  $B^{-1}(1) = \{\bullet\}$

$\downarrow$   $A^{-1}(2) \cap B^{-1}(1) = \{\bullet\}$  Hence:  $\mathbb{P}_{A,B}(2, 1) = \mathbb{P}(\bullet) = \frac{1}{6}$

# Discrete Probabilities

Consider the case when we can enumerate all elements in the target set  $\mathcal{T}$ .

In other words,  $\mathcal{T} = \{x_1, x_2, \dots, x_k, \dots\}$

Random variable  $A$  defines a probability distribution on  $\mathcal{T}$ , such that

for any  $x_i \in \mathcal{T}$  its probability  $\mathbb{P}_A(x_i) = \mathbb{P}(A = x_i) = p_i$ .

Hence, discrete random variable is characterized by a probability

mass function  $\mathbb{P}_A() = \{p_1, p_2, \dots\}$ , such that  $p_i \geq 0$  and  $\sum_i p_i = 1$



## Examples:

$$A = \text{Bernoulli}(p)$$

Consider  $\mathcal{T} = \{0, 1\}$  and :

$$\mathbb{P}_A(1) = \mathbb{P}(A = 1) = p, \quad \mathbb{P}_A(0) = \mathbb{P}(A = 0) = 1 - p$$

$$A = \text{Poisson}(\lambda)$$

Consider  $\mathcal{T} = \{0, 1, 2, \dots\}$  and :

$$\mathbb{P}_A(k) = \mathbb{P}(A = k) = \frac{\lambda^k e^{-\lambda}}{k!} \quad \text{with } \lambda \in (0, +\infty)$$

## CONDITIONAL PROBABILITY FOR RANDOM VARIABLES:

Discrete random variables:  $A : \Omega \rightarrow \mathcal{T}$  and  $B : \Omega \rightarrow \mathcal{T}$ , and  $x, y \in \mathcal{T}$ .

Conditional Probability of  $B = y$  condition on  $A = x$  is defined as:

$$\mathbb{P}_{B|A}(y, x) = \mathbb{P}(\omega \in A^{-1}(x), \text{ such that } B(\omega) = y)$$



$$\mathbb{P}_{A,B}(x, y) = \mathbb{P}(\omega \in \Omega, \text{ such that } A(\omega) = x \text{ and } B(\omega) = y)$$

Notice, in contrast with joint probability, here we are looking for elementary outcomes giving  $B(\omega) = y$  only in the preimage  $A^{-1}(x)$ , rather than the whole set  $\Omega$ .

$$A(\bullet) = 2 \quad A(\circ) = 1$$

$$A(\bullet\bullet) = 2 \quad A(\circ\circ) = -10 \quad \text{and}$$

$$A(\bullet\circ) = 0 \quad A(\circ\bullet) = 3$$

$$B(\bullet) = 1 \quad B(\circ) = 2 \quad x = 2$$

$$B(\bullet\bullet) = 2 \quad B(\circ\circ) = 3, \quad y = 1$$

$$B(\bullet\circ) = 0 \quad B(\circ\bullet) = 3$$

$$A^{-1}(2) = \{\bullet\bullet, \circ\circ\} \quad \text{and} \quad B^{-1}(1) = \{\bullet\}$$

Hence:

$$\mathbb{P}_{B|A}(1, 2) = \mathbb{P}(\bullet \text{ out of } \bullet\bullet, \circ\circ) = \frac{1}{2}$$

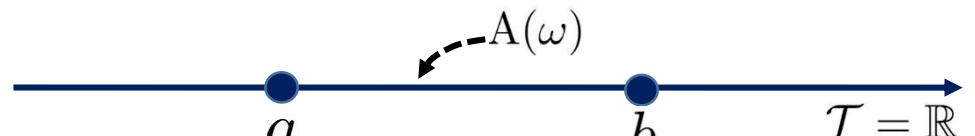
# Continuous Probabilities

Consider now when the target set is not countable, for example  $\mathcal{T} = \mathbb{R}$ .

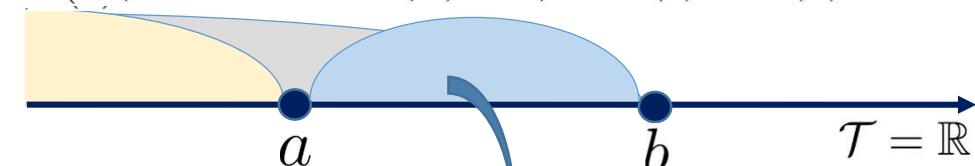
We interested in the probability that a random variable  $A$  is in the interval, i.e.  $\mathbb{P}(\omega \in \Omega : a \leq A(\omega) < b)$  for  $a, b \in \mathbb{R}$ .

To compute it, we introduce cumulative distribution function (cdf):

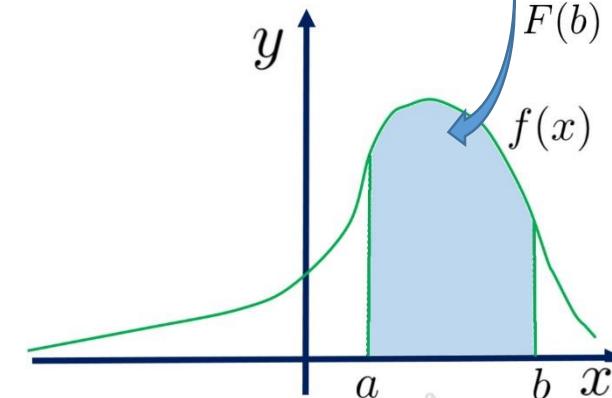
$$F(x) = \mathbb{P}(\omega \in \Omega : A(\omega) \leq x)$$



$$F(\mathbb{P}(\omega \in \Omega : a \leq A(\omega) < b)) = F(b) - F(a)$$



$$F(b) - F(a) = \int_a^b f(x)dx$$

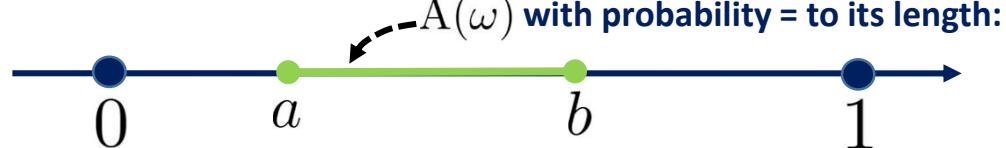


A function  $f : \mathbb{R} \rightarrow \mathbb{R}$  is called probability density function (pdf) if:

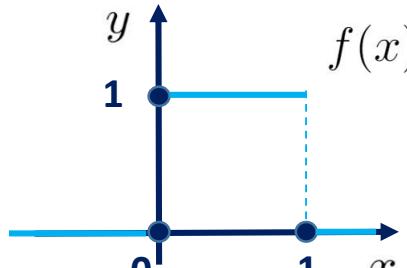
- Non-negative:  $f(x) \geq 0$  for any  $x \in \mathbb{R}$
- Integrable:  $\int_{\mathbb{R}} f(x)dx = 1$



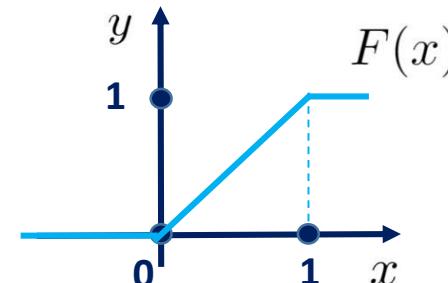
Example:  $A = \text{Uniform } ([0,1])$



Probability density function:



Cumulative distribution function:



For any interval  $(a, b) \subseteq [0, 1]$  we have:  $\mathbb{P}(A \in (a, b)) = b - a$

# Sum Rule

Unify notation for discrete and continuous random variables

DISCRETE RANDOM VARIABLES:		CONTINUOUS RANDOM VARIABLES:	
$\mathbb{P}_A(x)$	$=$	$p(x)$	$= p_A(x)$
$\mathbb{P}_{A,B}(x, y)$	$=$	$p(x, y)$	$= p_{A,B}(x, y)$
$\mathbb{P}_{A B}(x y)$	$=$	$p(x y)$	$= p_{A B}(x y)$

Probability density functions

## SUM RULE:

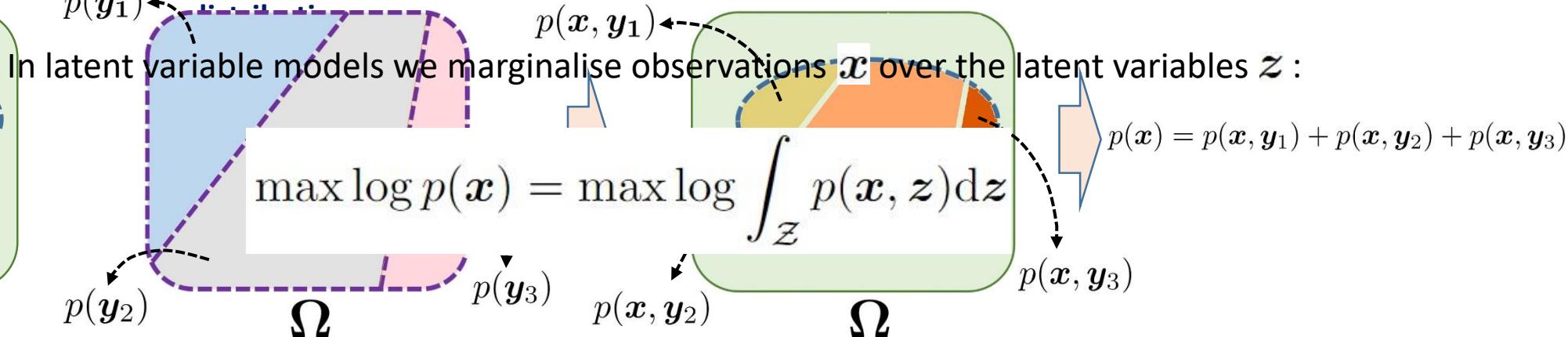
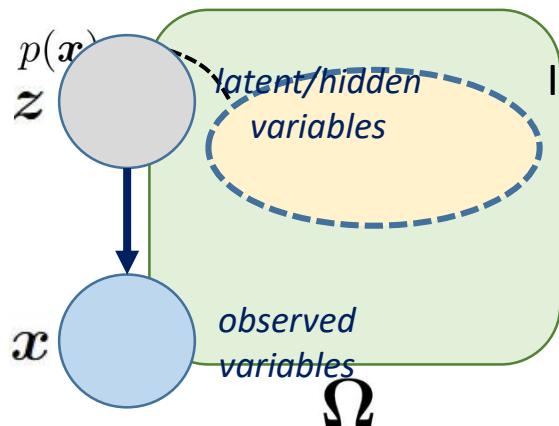
Consider random variables:

$$A : \Omega \rightarrow \mathcal{T}_A \text{ and } B : \Omega \rightarrow \mathcal{T}_B$$

$$p(x) = \begin{cases} \sum_{y \in \mathcal{T}_B} p(x, y) & \text{if } \mathcal{T}_B \text{ is discrete} \\ \int_{\mathcal{T}_B} p(x, y) dy & \text{if } \mathcal{T}_B \text{ is continuous} \end{cases}$$

sum out (or integrate out) the set of states  $y$  of the random variable  $B$ . This operation is called marginalisation

## Metric application:



# Product Rule

Unify notation for discrete  
and continuous random variables

**DISCRETE RANDOM VARIABLES:**

Probability mass functions

$$\begin{cases} \mathbb{P}_A(x) = p(x) \\ \mathbb{P}_{A,B}(x,y) = p(x,y) \\ \mathbb{P}_{A|B}(x|y) = p(x|y) \end{cases}$$

**CONTINUOUS RANDOM VARIABLES:**

$$\begin{cases} p(x) = p_A(x) \\ p(x,y) = p_{A,B}(x,y) \\ p(x|y) = p_{A|B}(x|y) \end{cases}$$

Probability density functions

## PRODUCT RULE:

Consider random variables:

$$A : \Omega \rightarrow \mathcal{T}_A \text{ and } B : \Omega \rightarrow \mathcal{T}_B$$

$$p(x, y) = p(y|x)p(x)$$

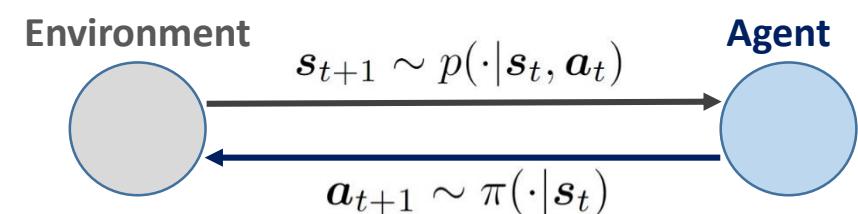
Every joint distribution of two random variables can be factorized of two other distributions. The two factors are the marginal distribution of the first random variable  $p(x)$ , and the conditional distribution of the second random variable given the first  $p(y|x)$

## ML application:

In Reinforcement Learning the probability of state action trajectory:

$\tau = (s_0, a_0, s_1, a_1, \dots, a_{H-1}, s_H)$  is factorised using of product rule:

$$p(\tau) = p(s_0, a_0, s_1, a_1, \dots, a_{H-1}, s_H) = p(s_0) \prod_{t=1}^H p(s_t | a_{t-1}, s_{t-1}) \pi(a_{t-1} | s_{t-1})$$



# Bayes Theorem

In machine learning and Bayesian statistics, we are often interested in making inferences of unobserved random variable  $\mathbf{x}$ , given that we have observed other random variable  $\mathbf{y}$ , related to  $\mathbf{x}$



$p(\mathbf{x})$  - our initial guess about hidden variable  $\mathbf{x}$

$p(\mathbf{y}|\mathbf{x})$  - relationship between  $\mathbf{x}$  and observed  $\mathbf{y}$



Can we draw any conclusions about  $\mathbf{x}$  given our observations  $\mathbf{y}$ ?



Thomas Bayes  
(1701-1761)

YES – using Bayes Formula

$p(\mathbf{x})$  - prior distribution over hidden variable  $\mathbf{x}$

$p(\mathbf{y}|\mathbf{x})$  - likelihood, describing conditional probability of  $\mathbf{y}$  given  $\mathbf{x}$ .

$p(\mathbf{x}|\mathbf{y})$  - posterior distribution, which is the quantity of interest, as it reflects conditional probability of  $\mathbf{x}$  given  $\mathbf{y}$ .

$p(\mathbf{y})$  - marginal likelihood/evidence “weighted average” of likelihood  $p(\mathbf{y}|\mathbf{x})$ , with weights defined by prior  $p(\mathbf{x})$

- using sum rule
- and then product rule

$$p(\mathbf{y}) = \sum_{\mathbf{x}} p(\mathbf{y}, \mathbf{x}) = \sum_{\mathbf{x}} p(\mathbf{y}|\mathbf{x})p(\mathbf{x})$$

Using product rule:

$$p(\mathbf{x}, \mathbf{y}) = p(\mathbf{y}|\mathbf{x})p(\mathbf{x})$$

$$p(\mathbf{x}, \mathbf{y}) = p(\mathbf{y})p(\mathbf{x}|\mathbf{y})$$

$$p(\mathbf{x}|\mathbf{y}) = \frac{p(\mathbf{y}|\mathbf{x})p(\mathbf{x})}{p(\mathbf{y})}$$

Similarly:

# Summary Statistics (Expected Value)

Weighted average expression over distribution is a special case of statistics. A statistic of a random variable is a deterministic function of that random variable. They provide a general view on the behavior of a random variable.

Two most popular statistics of random variable are mean and variance.

## MEAN (MATHEMATICAL EXPECTATION, EXPECTED VALUE, FIRST MOMENT)

### CASE I (Discrete):

Let  $A : \Omega \rightarrow \mathcal{T}$  be a discrete random variable with the target set  $\mathcal{T} = \{x_1, x_2, \dots, x_k, \dots\}$  and the probability mass function  $p(A = x_i) = p_i$  for  $i = 1, 2, \dots$ . The expected value of  $A$  is a sum:

$$\mathbb{E}[A] = \sum_i x_i p_i$$



Consider a fair dice (all sides are equally probable  $= \frac{1}{6}$ ) and random variable  $A$ :

$$\begin{aligned} A(\bullet) &= 2 & A(\circ) &= 1 \\ A(\square) &= 2 & A(\diamond) &= -10 \\ A(\blacksquare) &= 0 & A(\blackdiamond) &= 3 \end{aligned}$$

### Example:

$$\mathbb{E}[A] = 2 \cdot \frac{1}{6} + 1 \cdot \frac{1}{6} + (-10) \cdot \frac{1}{6} + 2 \cdot \frac{1}{6} + 0 \cdot \frac{1}{6} + 3 \cdot \frac{1}{6} = \boxed{-\frac{1}{3}}$$

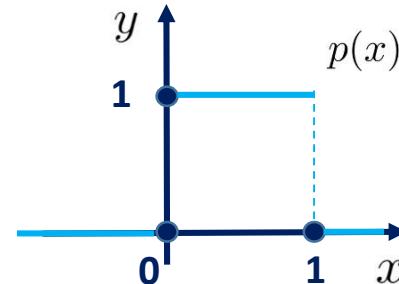
### CASE II (Continuous):

Let  $A : \Omega \rightarrow \mathcal{T}$  be a continuous random variable with the target set  $\mathcal{T} = \mathbb{R}^d$  and the probability density function  $p(x)$ . The mean value of  $A$  is given by an integral (if it exists):

$$\mathbb{E}[A] = \int_{\mathcal{T}} x p(x) dx$$



Consider  $A = \text{Uniform } ([0,1])$  with pdf  $p(x)$ :



$$\begin{aligned} \mathbb{E}[A] &= \int_{-\infty}^{+\infty} x p(x) dx = \\ &= \int_0^1 x dx = \frac{x^2}{2} \Big|_0^1 = \frac{1}{2} \end{aligned}$$

# Summary Statistics (Variance)

Weighted average expression over distribution is a special case of statistics. A statistic of a random variable is a deterministic function of that random variable. They provide a general view on the behavior of a random variable.

Two most popular statistics of random variable are mean and variance.

## VARIANCE (DISPERSION, SECOND MOMENT)

### CASE I (Discrete):

Let  $A : \Omega \rightarrow \mathcal{T}$  be a discrete random variable with the target set  $\mathcal{T} = \{x_1, x_2, \dots, x_k, \dots\}$  and the probability mass function  $p(A = x_i) = p_i$  for  $i = 1, 2, \dots$ . The variance of  $A$  is a sum:

$$\mathbb{V}[A] = \mathbb{E}[(A - \mathbb{E}[A])^2] = \sum_i p_i(x_i - \mathbb{E}[A])^2$$

### CASE II (Continuous):

Let  $A : \Omega \rightarrow \mathcal{T}$  be a continuous random variable with the target set  $\mathcal{T} = \mathbb{R}^d$  and the probability density function  $p(\mathbf{x})$ . Then, the variance of  $A$  is given by an integral (if it exist):

$$\mathbb{V}[A] = \int_{\mathcal{T}} \|\mathbf{x} - \mathbb{E}[A]\|_2^2 p(\mathbf{x}) d\mathbf{x}$$

### Example:

Consider a fair dice (all sides are equally probable =  $\frac{1}{6}$ ) and random variable  $A$ :

We know:  $\mathbb{E}[A] = \frac{1}{3}$

$$\mathbb{V}[A] = 2 \cdot \frac{1}{6} \left[2 + \frac{1}{3}\right]^2 + \frac{1}{6} \left[1 + \frac{1}{3}\right]^2 + \frac{1}{6} \left[3 + \frac{1}{3}\right]^2 + \frac{1}{6} \left[0 + \frac{1}{3}\right]^2 + \frac{1}{6} \left[-10 + \frac{1}{3}\right]^2 = \frac{176}{9}$$

Consider  $A = \text{Uniform } ([0,1])$ :

We know:  $\mathbb{E}[A] = \frac{1}{2}$

$$\mathbb{V}[A] = \int_{-\infty}^{+\infty} \left[x - \frac{1}{2}\right]^2 p(x) dx = \int_0^1 \left[x - \frac{1}{2}\right]^2 dx = \left[\frac{x^3}{3} - \frac{x^2}{2} + \frac{x}{4}\right]_0^1 = \frac{1}{12}$$

# Summary Statistics (ML Application)

Weighted average expression over distribution is a special case of statistics. A statistic of a random variable is a deterministic function of that random variable. They provide a general view on the behavior of a random variable.

Mean and Variance are used almost everywhere in ML:

## VARIOUS ML APPLICATIONS OF MEAN/VARIANCE:

- **In Supervised Learning:**

Training loss function for ML models is usually represented as empirical mean:

$$\mathcal{L}_{\text{training}}(\Theta) = \mathbb{E}_{(x,y) \sim \mathcal{D}_{\text{tr}}} [(y - \text{NN}_\Theta(x))^2] = \frac{1}{N} \sum_{i=1}^N (y_i - \text{NN}_\Theta(x_i))^2$$



Here:

$\mathcal{D}_{\text{tr}} = \{\mathbf{x}_i, y_i\}_{i=1}^N$  — training data-set



$\text{NN}_\Theta(\mathbf{x})$  — ML model  
(Neural Network)

- **In Bayesian Optimisation:**

UCB acquisition function for Bayesian Optimisation has the following form:

$$\alpha_{\text{UCB}}(\mathbf{x}) = \mu_{\text{post}}(\mathbf{x}; \mathcal{D}_{\text{obs}}) + \beta \sigma_{\text{post}}(\mathbf{x}; \mathcal{D}_{\text{obs}})$$

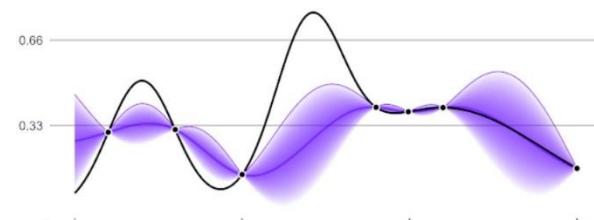


Here:

$\mu_{\text{post}}(\mathbf{x}; \mathcal{D}_{\text{obs}})$  — is posterior mean of GP model

$\sigma_{\text{post}}^2(\mathbf{x}; \mathcal{D}_{\text{obs}})$  — is posterior variance of GP model

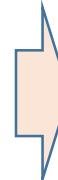
$\mathcal{D}_{\text{obs}} = \{\mathbf{x}_i, y_i\}_{i=1}^M$  — is collection of black-box observations.



# Independence

Two random variables A, B are statistically independent if and only if:

$$p(\mathbf{x}, \mathbf{y}) = p(\mathbf{x})p(\mathbf{y})$$



Intuitively, A, B are independent if the value of one of them does not add any additional information about the value of another.



## ML application:

In machine learning, we often consider problems that can be modeled as independent and identically distributed(i.i.d.) random variables. For example, in supervised learning, training data-set is assumed to consist of i.i.d. samples. In Reinforcement learning, for a fixed policy the sampled trajectories are i.i.d.

## SOME PROPERTIES OF RANDOM VARIABLES:

- If A, B are independent, then:  $p(\mathbf{y}|\mathbf{x}) = p(\mathbf{y})$  and  $p(\mathbf{x}|\mathbf{y}) = p(\mathbf{x})$
- If A, B are independent, then:  $\mathbb{V}[A + B] = \mathbb{V}[A] + \mathbb{V}[B]$
- If A, B are independent, then:  $\mathbb{E}[AB] = \mathbb{E}[A]\mathbb{E}[B]$  and  $\mathbb{E}\left[\frac{A}{B}\right] = \frac{\mathbb{E}[A]}{\mathbb{E}[B]}$
- If A, B are independent, then:  $\mathbb{E}[f(A)g(B)] = \mathbb{E}[f(A)]\mathbb{E}[g(B)]$ , where  $f(\cdot), g(\cdot)$  – are integrable.
- If A, B are random variables (might be dependent), then:  $\mathbb{E}[\alpha A + \beta B] = \alpha\mathbb{E}[A] + \beta\mathbb{E}[B]$ , where  $\alpha, \beta \in \mathbb{R}$

## Examples:

Flipping a coin several times, then random variables associated with each such trial are independent.

Roll a dice (random variable A) and then sample a point  $B \sim \text{Uniform } ([0,1])$ . Then, A, B are independent

Roll a dice (random variable A) and then sample a point  $B \sim \text{Uniform } ([0,1])$ . Consider random variable

$$C = \begin{cases} 2B & \text{if } A = \{ \text{dice faces} \} \\ \frac{1}{2}B & \text{otherwise} \end{cases}$$

Then, A, C are not independent.

**Linearity of Expectation**

# Gaussian Random Variable

The Gaussian (normal) distribution is the most well-studied probability distribution for continuous-valued random variables. It has application in many fields of ML, including diffusion models, normalizing flows, Kalman filters, Gaussian processes ,etc.



Carl Friedrich Gauss  
(1777-1855)

The Gaussian random variable  $X$  can be 1-dimensional (i.e.  $X \in \mathbb{R}$ ) or  $d$ -dimensional (i.e.  $X \in \mathbb{R}^d$ ) and can be fully determined by its first and second moments:

## 1-DIMENSIONAL CASE:

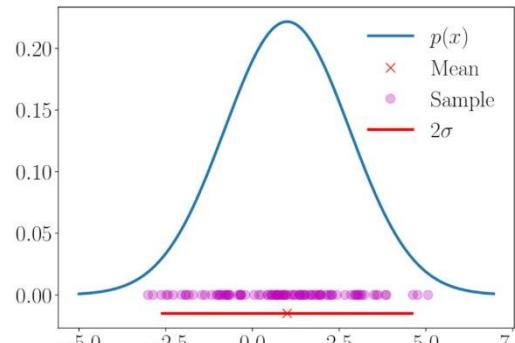
Probability density function:

$$p(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

$\mu \in \mathbb{R}$  - is expected value of  $X$

$\sigma^2 \in \mathbb{R}$  - is variance of  $X$

Examples\*:



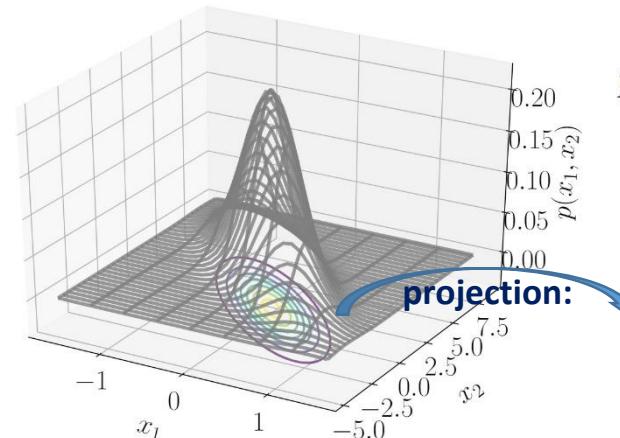
$$p(x) = \mathcal{N}(1, 1.69)$$

## $d$ -DIMENSIONAL CASE:

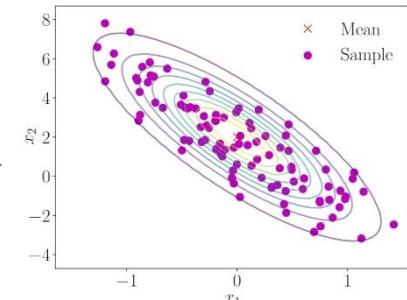
$$p(x|\mu, \Sigma) = (2\pi)^{-\frac{d}{2}} |\Sigma|^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(x-\mu)^\top \Sigma^{-1} (x-\mu)\right)$$

$\mu \in \mathbb{R}^d$  - is expected value of  $X$

$\Sigma \in \mathbb{R}^{d \times d}$  - is covariance matrix:  
 $\Sigma = \mathbb{E}[(x-\mu)^\top (x-\mu)]$



$$p(x_1, x_2) = \mathcal{N}\left(\begin{bmatrix} 1 \\ 2 \end{bmatrix}, \begin{bmatrix} 0.3 & -1 \\ -1 & 5 \end{bmatrix}\right)$$



# Properties of Gaussian

The Gaussian (normal) distribution is the most well-studied probability distribution for continuous-valued random variables. It has application in many fields of ML, including diffusion models, normalizing flows, Kalman filters, Gaussian processes ,etc.



Carl Friedrich Gauss  
(1777-1855)

## 1-DIMENSIONAL CASE:

$$p(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

## $d$ -DIMENSIONAL CASE:

$$p(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = (2\pi)^{-\frac{d}{2}} |\boldsymbol{\Sigma}|^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x}-\boldsymbol{\mu})\right)$$

### Sum of independent Gaussians:

Let  $X = \mathcal{N}(\mu_1, \sigma_1^2)$  and  $Y = \mathcal{N}(\mu_2, \sigma_2^2)$  be two independent Gaussian random variables. Then, for any  $\alpha, \beta \in \mathbb{R}$ :

$$\alpha X + \beta Y = \mathcal{N}(\alpha\mu_1 + \beta\mu_2, \alpha^2\sigma_1^2 + \beta^2\sigma_2^2)$$

$d$ -dimensional version for  $X = \mathcal{N}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1) \in \mathbb{R}^d$  and  $Y = \mathcal{N}(\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2) \in \mathbb{R}^d$ :

$$\alpha X + \beta Y = \mathcal{N}(\alpha\boldsymbol{\mu}_1 + \beta\boldsymbol{\mu}_2, \alpha^2\boldsymbol{\Sigma}_1 + \beta^2\boldsymbol{\Sigma}_2)$$

### Linear transformation of Gaussian:

Let  $X = \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \in \mathbb{R}^d$  be a multivariate Gaussian, and  $\mathbf{A} \in \mathbb{R}^{d \times p}$  be an arbitrary matrix. Consider random variable

$$Y = \mathbf{AX} \in \mathbb{R}^p - \text{which is linear transformation of } X \text{ by } \mathbf{A}. \text{ Then: } Y = \mathcal{N}(\mathbf{A}\boldsymbol{\mu}, \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^\top)$$

### Mixture of univariate normal densities:

Let  $p_1(x|\mu_1, \sigma_1^2)$  and  $p_2(x|\mu_2, \sigma_2^2)$  be two Gaussian densities and  $\alpha \in (0, 1)$ . Consider density:

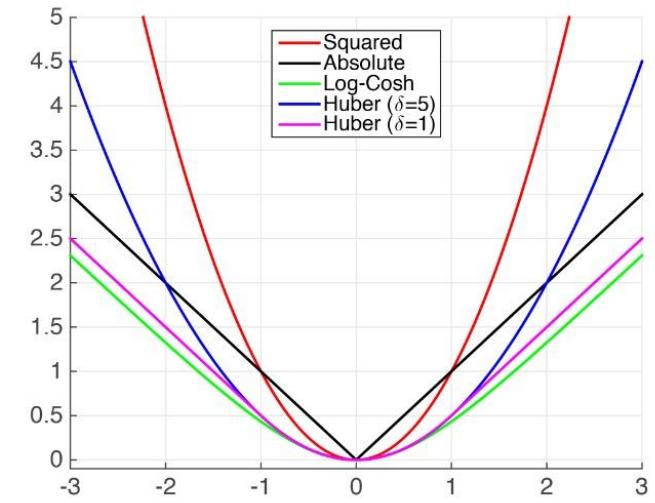
$$q(x) = \alpha p_1(x|\mu_1, \sigma_1^2) + (1 - \alpha)p_2(x|\mu_2, \sigma_2^2)$$

$$\mathbb{E}[Q] = \alpha\mu_1 + (1 - \alpha)\mu_2$$

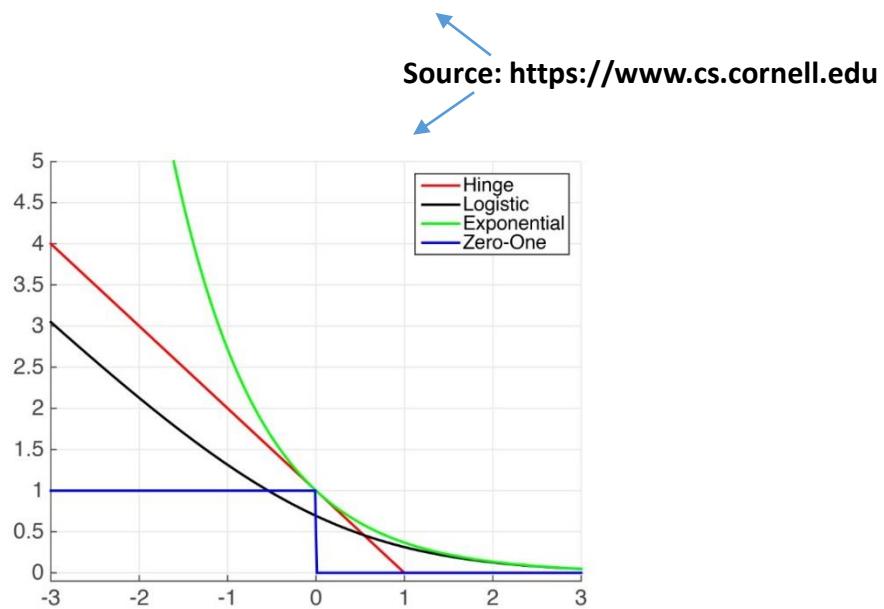
$$\mathbb{V}[Q] = [\alpha\sigma_1^2 + (1 - \alpha)\sigma_2^2] + ([\alpha\mu_1^2 + (1 - \alpha)\mu_2^2] - [\alpha\mu_1 + (1 - \alpha)\mu_2]^2)$$

# Loss Functions in ML

- Regression Loss Function



- Classification Loss Function



Source: <https://www.cs.cornell.edu>

# Regression Loss Function (MLE)

REGRESSION TASK SETUP:

$f(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}$  - unknown function we want to model.

$\Phi_{\Theta}(\mathbf{x}) : \mathbb{R}^d \times \mathbb{R}^n \rightarrow \mathbb{R}$  - our ML model (typically NN) to approximate  $f(\mathbf{x})$ .

$\mathcal{D}_{\text{tr}} = \{\mathbf{x}_i, y_i\}_{i=1}^m$  - observed data-set, where  $\mathbf{x}_i \in \mathbb{R}^n$  input variables, and  $y_i = f(\mathbf{x}_i) + \epsilon_i$  are noisy observations.

GOAL: Find parameter  $\underline{\Theta^*} \in \mathbb{R}^d$  such that  $\Phi_{\Theta^*}(\mathbf{x}) \approx f(\mathbf{x})$



We assume that  $y(\mathbf{x}) = \underline{\Phi_{\Theta}(\mathbf{x})} + \mathcal{N}(0, \sigma_{\text{noise}}^2) \Rightarrow y(\mathbf{x}) \sim \mathcal{N}(\Phi_{\Theta}(\mathbf{x}), \sigma_{\text{noise}}^2)$

Hence, probability density of points sampled independently defines the likelihood function:

$$p(y_1, \dots, y_m | \Phi_{\Theta}(\mathbf{x}_1), \dots, \Phi_{\Theta}(\mathbf{x}_m)) = \prod_{i=1}^m p(y_i | \Phi_{\Theta}(\mathbf{x}_i), \sigma_{\text{noise}}^2) \stackrel{\text{i.i.d. samples}}{=} \prod_{i=1}^m \frac{1}{\sigma_{\text{noise}} \sqrt{2\pi}} \exp\left(-\frac{(y_i - \Phi_{\Theta}(\mathbf{x}_i))^2}{2\sigma_{\text{noise}}^2}\right)$$



$$\Theta^* = \arg \min_{\Theta} \sum_{i=1}^m (y_i - \Phi_{\Theta}(\mathbf{x}_i))^2 = \arg \max_{\Theta} \log \left[ \prod_{i=1}^m \frac{1}{\sigma_{\text{noise}} \sqrt{2\pi}} \exp\left(-\frac{(y_i - \Phi_{\Theta}(\mathbf{x}_i))^2}{2\sigma_{\text{noise}}^2}\right) \right] = \arg \max_{\Theta} \sum_{i=1}^m \frac{(y_i - \Phi_{\Theta}(\mathbf{x}_i))^2}{2\sigma_{\text{noise}}^2}$$

If we replace  $(y_i - \Phi_{\Theta}(\mathbf{x}_i))^2 \Rightarrow |y_i - \Phi_{\Theta}(\mathbf{x}_i)|$  we get mean absolute error (MAE) loss:  $\Theta^* = \arg \min_{\Theta} \sum_{i=1}^m |y_i - \Phi_{\Theta}(\mathbf{x}_i)|$

# Regression Loss Function (MAP)

## REGRESSION TASK SETUP:

$f(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}$  - unknown function we want to model.

$\Phi_{\Theta}(\mathbf{x}) : \mathbb{R}^d \times \mathbb{R}^n \rightarrow \mathbb{R}$  - our ML model (typically NN) to approximate  $f(\mathbf{x})$ .

$\mathcal{D}_{\text{tr}} = \{\mathbf{x}_i, y_i\}_{i=1}^m$  - observed data-set, where  $\mathbf{x}_i \in \mathbb{R}^n$  input variables, and  $y_i = f(\mathbf{x}_i) + \epsilon_i$  are noisy observations.

$$\epsilon_i \sim \mathcal{N}(0, \sigma_{\text{noise}}^2)$$

Previously, we found parameter value  $\Theta^*$  maximising the probability of observing data points  $\{y_1, y_2, \dots, y_m\}$

Now, let's treat  $\Theta \in \mathbb{R}^d$  as random variable itself and for a given data  $\{y_1, y_2, \dots, y_m\}$  aim to find most probable value of model parameter  $\tilde{\Theta}$ :

Let assume prior distribution  $\Theta \sim p_{\text{prior}} = \mathcal{N}(\mathbf{0}, \tau^2 \mathbf{I})$ . Then we need to maximise posterior distribution:  $p(\Theta | y_1, \dots, y_m)$

Using Bayes Formula:  $p(\Theta | y_1, \dots, y_m) = \frac{p(y_1, \dots, y_m | \Phi_{\Theta}(x_1), \dots, \Phi_{\Theta}(x_m)) p_{\text{prior}}(\Theta)}{p_{\text{prior}}(\Theta)}$  gives:

$$\begin{aligned} \tilde{\Theta} &= \arg \max_{\Theta} \sum_{i=1}^m \left[ -\frac{(y_i - \Phi_{\Theta}(x_i))^2}{2\sigma_{\text{noise}}^2} \right] - \frac{\|\Theta\|_2^2}{2\tau^2} = \arg \min_{\Theta} \sum_{i=1}^m \frac{(y_i - \Phi_{\Theta}(x_i))^2}{2\sigma_{\text{noise}}^2} + \frac{\sigma_{\text{noise}}^2}{\tau^2} \|\Theta\|_2^2 \\ &= \arg \max_{\Theta} \log [p(y_1, \dots, y_m | \Phi_{\Theta}(x_1), \dots, \Phi_{\Theta}(x_m)) p_{\text{prior}}(\Theta)] = \end{aligned}$$

Step 1: This step does not distribute over the model

In other words, maximising a posterior probability we arrive at L-2 regularized mean squared error (MSE) loss function.

If we replace prior from Gaussian to Laplace, i.e.  $p_{\text{prior}} = \text{Laplace}(0, \beta)$  then we arrive to L-1 regularised mean squared loss function:

$$= \arg \max_{\Theta} \left[ \tilde{\Theta} = \arg \min_{\Theta} \sum_{i=1}^m \left( \frac{(y_i - \Phi_{\Theta}(x_i))^2}{2\sigma_{\text{noise}}^2} + \frac{2\sigma_{\text{noise}}^2}{\beta} \|\Theta\|_1^2 \right) \right] =$$

# Classification Loss Function

## CLASSIFICATION TASK SETUP:

$\mathcal{D}_{\text{tr}} = \{\mathbf{x}_i, y_i\}_{i=1}^m$  - observed data-set, where  $\mathbf{x}_i \in \mathbb{R}^n$  input variables, and each  $y_i \in \{0, 1\}$ .

Classifier  $\Theta(\mathbf{x}) = \begin{bmatrix} \Phi_{\Theta}(\mathbf{x}) \\ 1 - \Phi_{\Theta}(\mathbf{x}) \end{bmatrix}$  - our ML model classifier (typically NN), outputting two dimensional stochastic vector.  
where  $\Phi_{\Theta}(\mathbf{x}) : \mathbb{R}^d \times \mathbb{R}^n \rightarrow [0, 1]$  defines probability that  $y(\mathbf{x}) = 1$

**GOAL:** Find parameter  $\Theta^* \in \mathbb{R}^d$  such that classifier maximizes the proper labelling of input points.

It means,  $\Phi_{\Theta}(\mathbf{x}) : \mathbb{R}^d \times \mathbb{R}^n \rightarrow [0, 1]$  should be close to one when input point  $\mathbf{x}$  belong to the first class ( $y(\mathbf{x}) = 1$ ), and it should be close to zero when input point  $\mathbf{x}$  belongs to the second class ( $y(\mathbf{x}) = 0$ ).

To combine these requirements, consider the following likelihood expression:

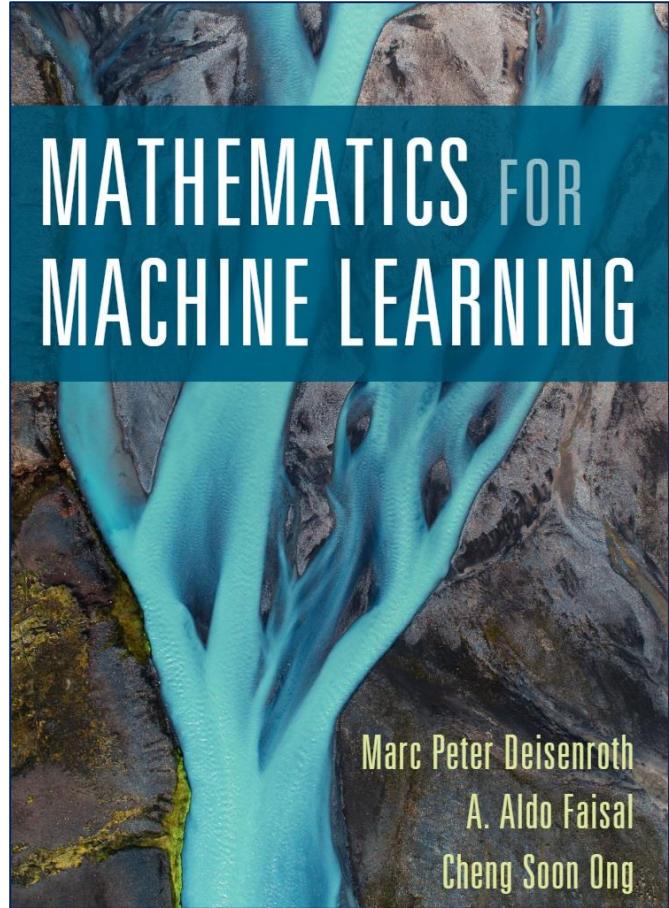
$$p(y(\mathbf{x})|\Theta) = \Phi_{\Theta}(\mathbf{x})^{y(\mathbf{x})}(1 - \Phi_{\Theta}(\mathbf{x}))^{1-y(\mathbf{x})} \quad \Rightarrow \quad \text{Notice, each output } y(\mathbf{x}), \text{ it returns the expression we need to maximization:}$$
$$p(y(\mathbf{x})|\Theta) = \begin{cases} \Phi_{\Theta}(\mathbf{x}) & \text{if } y(\mathbf{x}) = 1 \\ 1 - \Phi_{\Theta}(\mathbf{x}) & \text{if } y(\mathbf{x}) = 0 \end{cases}$$

Aiming maximisation of the likelihood:

$$\Theta^* = \arg \max_{\Theta} p(y_1, \dots, y_m | \Theta) \stackrel{\text{Step 1}}{=} \arg \max_{\Theta} \prod_{i=1}^m p(y_i | \Theta) \stackrel{\text{Step 2}}{=} \arg \max_{\Theta} \prod_{i=1}^m \Phi_{\Theta}(\mathbf{x}_i)^{y_i} (1 - \Phi_{\Theta}(\mathbf{x}_i))^{1-y_i} =$$
$$\stackrel{\text{Step 3}}{=} \arg \max_{\Theta} \log \left[ \prod_{i=1}^m \Phi_{\Theta}(\mathbf{x}_i)^{y_i} (1 - \Phi_{\Theta}(\mathbf{x}_i))^{1-y_i} \right] = \arg \max_{\Theta} \sum_{i=1}^m y_i \log [\Phi_{\Theta}(\mathbf{x}_i)] + (1 - y_i) \log [1 - \Phi_{\Theta}(\mathbf{x}_i)]$$

**Step 2: Distinguishing for the likelihood**

We arrive at cross entropy loss function.



Thank you