



# Applied Deep Learning

Dr. Philippe Blaettchen  
Bayes Business School (formerly Cass)

## Learning objectives of today

**Goals:** Understand the key concepts of linear algebra and calculus relevant to deep learning

- Basic definitions
- Taking derivatives
- Typical operations on vectors and matrices

**How will we do this?**

- Pick up where the videos left off
- Not a comprehensive review, but focusing on the most relevant concepts for understanding deep learning
- Introduction on how to implement concepts in Python

## **Linear algebra – definitions**

## Scalars

- Scalar: a single number
- Integers (-1,0,1,2,...), real numbers (0.319375, 1.17,  $\pi$ ), rational numbers ( $\frac{\text{integer}}{\text{integer}}$ )

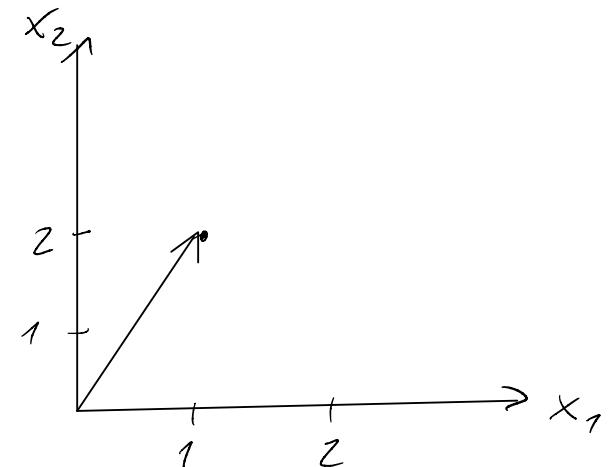
$$\alpha, t \quad \alpha \in \mathbb{R}, \quad t \in \mathbb{Z}$$

## Vectors

- A one-dimensional array of numbers:

$$\textcircled{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

$$x = \begin{pmatrix} 1 \\ 2 \end{pmatrix}$$



- Entries can be real numbers, binary, integer, ...

→ we usually denote a vector with the type of its entries and its size, e.g.,  $\mathbf{x} \in \underline{\mathbb{R}^n}$

## Vectors

- A one-dimensional array of numbers:

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

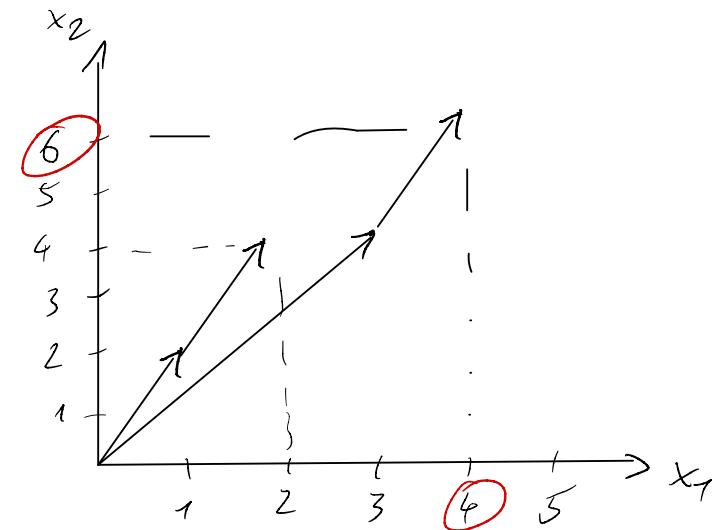
- Entries can be real numbers, binary, integer, ...  
→ we usually denote a vector with the type of its entries and its size, e.g.,  $\mathbf{x} \in \mathbb{R}^n$

- We can perform elementary algebra on vectors:

$$\bullet \quad \underline{\mathbf{x}} + \underline{\mathbf{y}} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} + \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} x_1 + y_1 \\ x_2 + y_2 \\ \vdots \\ x_n + y_n \end{bmatrix}$$

(check dimensions!)

$$\bullet \quad \cancel{\alpha \mathbf{x}} = \alpha \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} \alpha x_1 \\ \alpha x_2 \\ \vdots \\ \alpha x_n \end{bmatrix}$$



$$\mathbf{x} = \begin{pmatrix} 1 \\ 2 \end{pmatrix} \quad \mathbf{y} = \begin{pmatrix} 3 \\ 4 \end{pmatrix} \quad \mathbf{x} + \mathbf{y} = \begin{pmatrix} 1+3 \\ 2+4 \end{pmatrix} = \begin{pmatrix} 4 \\ 6 \end{pmatrix}$$

$$\alpha = 2, \quad \mathbf{x} = \begin{pmatrix} 1 \\ 2 \end{pmatrix} \quad \alpha \mathbf{x} = \begin{pmatrix} 2 \\ 4 \end{pmatrix}$$

## Matrices

- A two-dimensional array of numbers

$$A = \begin{bmatrix} A_{1,1} & A_{1,2} \\ A_{2,1} & A_{2,2} \end{bmatrix}.$$

Row

Column

$n = 2$

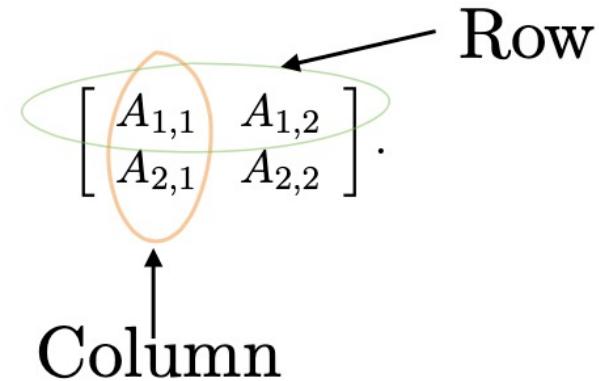
$m = 2$

- We again denote it with its type and shape:  $A \in \mathbb{R}^{n \times m}$ , where  $n$  is the number of rows and  $m$  is the number of columns

Source: Goodfellow

# Matrices

- A two-dimensional array of numbers



- We again denote it with its type and shape:  $A \in \mathbb{R}^{n \times m}$ , where  $n$  is the number of rows and  $m$  is the number of columns

- We can perform elementary algebra on vectors:

$$A + B = \underbrace{\begin{bmatrix} a_{1,1} & a_{1,2} \\ a_{2,1} & a_{2,2} \end{bmatrix}}_{\text{matrices}} + \underbrace{\begin{bmatrix} b_{1,1} & b_{1,2} \\ b_{2,1} & b_{2,2} \end{bmatrix}}_{\text{vectors}} = \begin{bmatrix} a_{1,1} + b_{1,1} & a_{1,2} + b_{1,2} \\ a_{2,1} + b_{2,1} & a_{2,2} + b_{2,2} \end{bmatrix}$$

$$A = \begin{pmatrix} 1 & 2 \\ 3 & 4 \end{pmatrix}, \quad B = \begin{pmatrix} 3 & 0 \\ 1 & 5 \end{pmatrix}$$

$$A + B = \begin{pmatrix} 1+3 & 2+0 \\ 3+1 & 4+5 \end{pmatrix} = \begin{pmatrix} 4 & 2 \\ 4 & 9 \end{pmatrix} \quad (\text{check dimensions!})$$

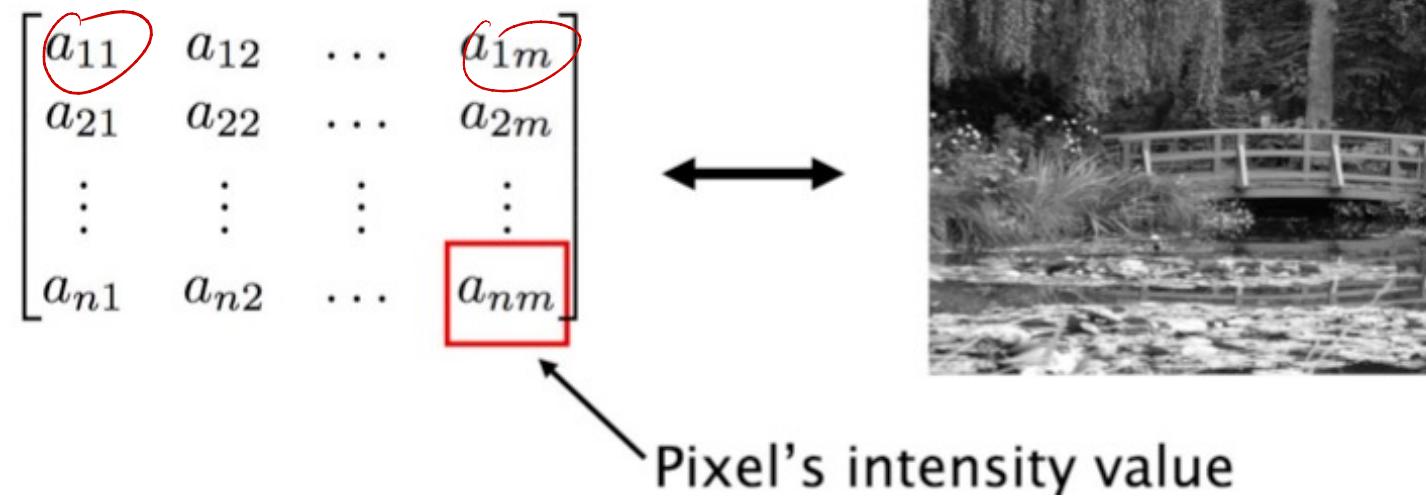
$$\begin{pmatrix} 1 & 2 \\ 3 & 4 \end{pmatrix} + \begin{pmatrix} 3 & 0 \\ 1 & 5 \end{pmatrix} \neq$$

$$\alpha A = \alpha \begin{bmatrix} a_{1,1} & a_{1,2} \\ a_{2,1} & a_{2,2} \end{bmatrix} = \begin{bmatrix} \alpha a_{1,1} & \alpha a_{1,2} \\ \alpha a_{2,1} & \alpha a_{2,2} \end{bmatrix}$$

Source: Goodfellow

## A typical use of matrices in deep learning

$$X = \begin{pmatrix} & & & \\ & 0 & \cdots & 0 \\ & \vdots & & \\ & 0 & & \end{pmatrix}$$



Source: Ivanovic

## Tensors

- Any array of numbers
- Could have zero dimensions (scalar), one dimension (vector), two dimensions (matrix)
- Could also have three or more dimensions

## In Python

*numpy*

- We generally use ~~Python~~ to work with vectors, matrices, and sometimes tensors
- Later, we'll also see the TensorFlow-specific implementation of tensors



## **Linear algebra – typical operations**

## Trace of a matrix

- Summing up all the entries on the diagonal of the matrix:

$$Tr(A) = \sum_{i=1}^n A_{i,i}$$

$$A = \begin{pmatrix} 1 & 2 & 5 \\ 4 & 3 & 0 \\ 2 & 1 & 0 \end{pmatrix}$$
$$Tr(A) = 1 + 3 + 0 = 4$$

- Note that
  - $Tr(ABC) = Tr(CAB) = Tr(BCA)$
  - $Tr(A + B) = Tr(A) + Tr(B)$

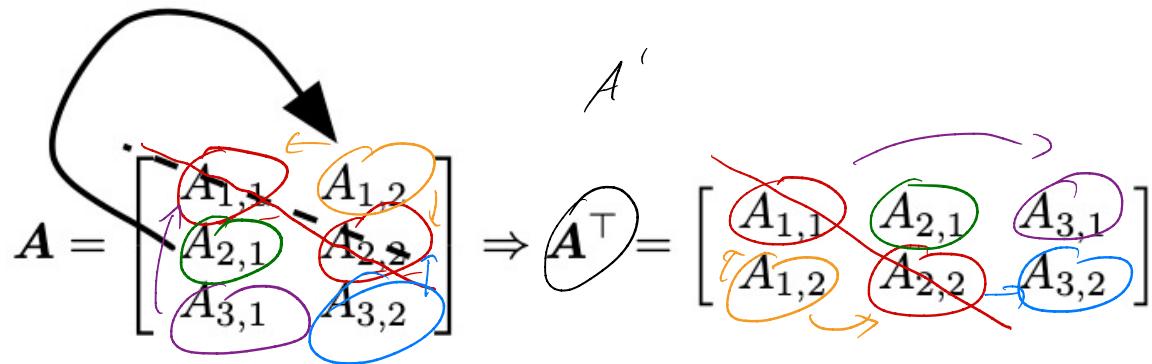
$$A = \begin{pmatrix} 0 & 2 \\ 3 & 0 \end{pmatrix} \quad B = \begin{pmatrix} 0 & 0 \\ 2 & 0 \end{pmatrix} \quad A + B = \begin{pmatrix} 0 & 2 \\ 5 & 0 \end{pmatrix}$$
$$Tr(A) = 5 \quad Tr(B) = 2 \quad Tr(A + B) = 7$$

## Matrix transpose

$$A = \begin{pmatrix} 1 & 2 \\ 3 & 4 \end{pmatrix} \quad A^T = \begin{pmatrix} 1 & 3 \\ 2 & 4 \end{pmatrix}$$

- Essentially a mirror image across the main diagonal

$A, T$



$$(A^T)_{i,j} = A_{j,i}.$$

- Note that:

- $(\underline{A^T})^T = A$
- $(\underline{AB})^T = \underline{B^T} \underline{A^T}$
- $(\underline{A + B})^T = \underline{A^T} + \underline{B^T}$

$$A = \begin{pmatrix} 1 & 3 \\ 3 & 2 \end{pmatrix} \quad A^T = \begin{pmatrix} 1 & 3 \\ 3 & 2 \end{pmatrix}$$

- We call a matrix symmetric if  $\underline{A} = \underline{A^T}$

Source: Goodfellow

## Inner product of vectors (also known as dot product)

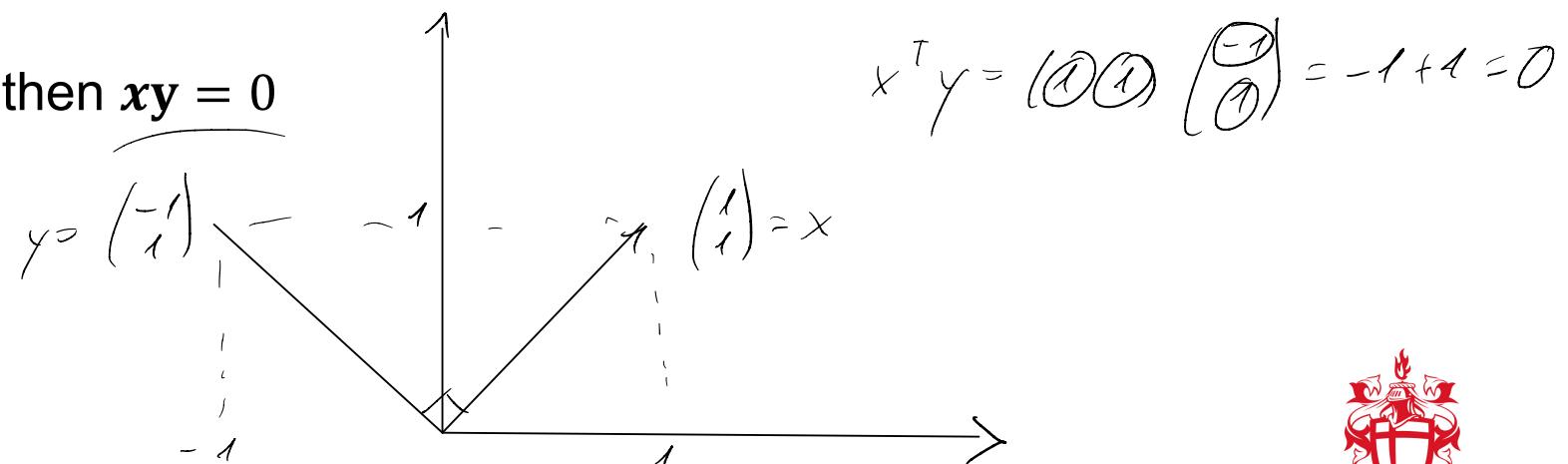
- Say  $x, y \in \mathbb{R}^n$ , then:

$$x^T y = [x_1 \ x_2 \ \dots \ x_n] \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \sum_{i=1}^n x_i y_i$$

$x = \begin{pmatrix} 1 \\ 2 \\ 3 \end{pmatrix} \quad y = \begin{pmatrix} 2 \\ 0 \\ 1 \end{pmatrix}$

$$x^T y = (1 \ 2 \ 3) \begin{pmatrix} 2 \\ 0 \\ 1 \end{pmatrix} = 2 + 0 + 3 = 5$$

- The inner product is a scalar!
- Note: if  $x$  and  $y$  are orthogonal, then  $xy = 0$



## Outer product of vectors

- Say  $x \in \mathbb{R}^n$ ,  $y \in \mathbb{R}^m$  then:

$$x y^T = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \begin{bmatrix} y_1 & y_2 & \cdots & y_m \end{bmatrix}^T = \begin{bmatrix} x_1 y_1 & x_1 y_2 & \cdots & x_1 y_m \\ x_2 y_1 & x_2 y_2 & \cdots & x_2 y_m \\ \vdots & \vdots & \ddots & \vdots \\ x_n y_1 & x_n y_2 & \cdots & x_n y_m \end{bmatrix}$$

The diagram illustrates the outer product  $x y^T$ . On the left, vector  $x$  is shown as a column with entries  $x_1, x_2, \dots, x_n$ , with a bracket labeled  $n$  indicating its dimension. To its right is a row vector  $y$  with entries  $y_1, y_2, \dots, y_m$ , with a bracket labeled  $m$  indicating its dimension. An arrow points from the row vector to the resulting matrix. The resulting matrix has  $n$  columns and  $m$  rows. Each entry in the matrix is circled in red and labeled  $x_i y_j$ , where  $i$  corresponds to the column index and  $j$  corresponds to the row index. Green ovals group the entries by column, and purple ovals group them by row.

- The outer product is a matrix!

$$\left( \quad \right)$$

A hand-drawn diagram of a matrix enclosed in parentheses. A vertical line on the left and a horizontal line at the top define the matrix boundaries. Red arrows point from the labels  $n$  and  $m$  to the respective dimensions of the matrix. The matrix itself is empty, represented by four small red dots at the corners of the parentheses.

## Matrix-vector product

$$A_{0,0} \quad d_{:,}$$

- Say  $\underline{y} \in \mathbb{R}^m$ ,  $A \in \mathbb{R}^{n \times m}$  then:

$$\underline{x} = A\underline{y} = \begin{bmatrix} a_{1,1} & a_{1,2} & \dots & a_{1,m} \\ a_{2,1} & a_{2,2} & \dots & a_{2,m} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n,1} & a_{n,2} & \dots & a_{n,m} \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix} = \begin{bmatrix} a_{1,1}y_1 + a_{1,2}y_2 + \dots + a_{1,m}y_m \\ a_{2,1}y_1 + a_{2,2}y_2 + \dots + a_{2,m}y_m \\ \vdots \\ a_{n,1}y_1 + a_{n,2}y_2 + \dots + a_{n,m}y_m \end{bmatrix}$$

$$A = \begin{pmatrix} 1 & 2 \\ 3 & 4 \end{pmatrix}, \quad Y = \begin{pmatrix} 1 \\ 2 \end{pmatrix}$$

$$AY = \begin{pmatrix} 5 \\ 11 \end{pmatrix}$$

$$= \begin{bmatrix} A_{1,:}\underline{y} \\ A_{2,:}\underline{y} \\ \vdots \\ A_{n,:}\underline{y} \end{bmatrix} = A_{:,1}y_1 + A_{:,2}y_2 + \dots + A_{:,m}y_m$$

## Matrix multiplication

$$A \cdot B = C$$

$$B \cdot A$$

$$(n \times p) \quad (m \times n) \quad C_{i,j} = \sum_k A_{i,k} B_{k,j}.$$

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

$$B \in \mathbb{R}^{n \times p}$$

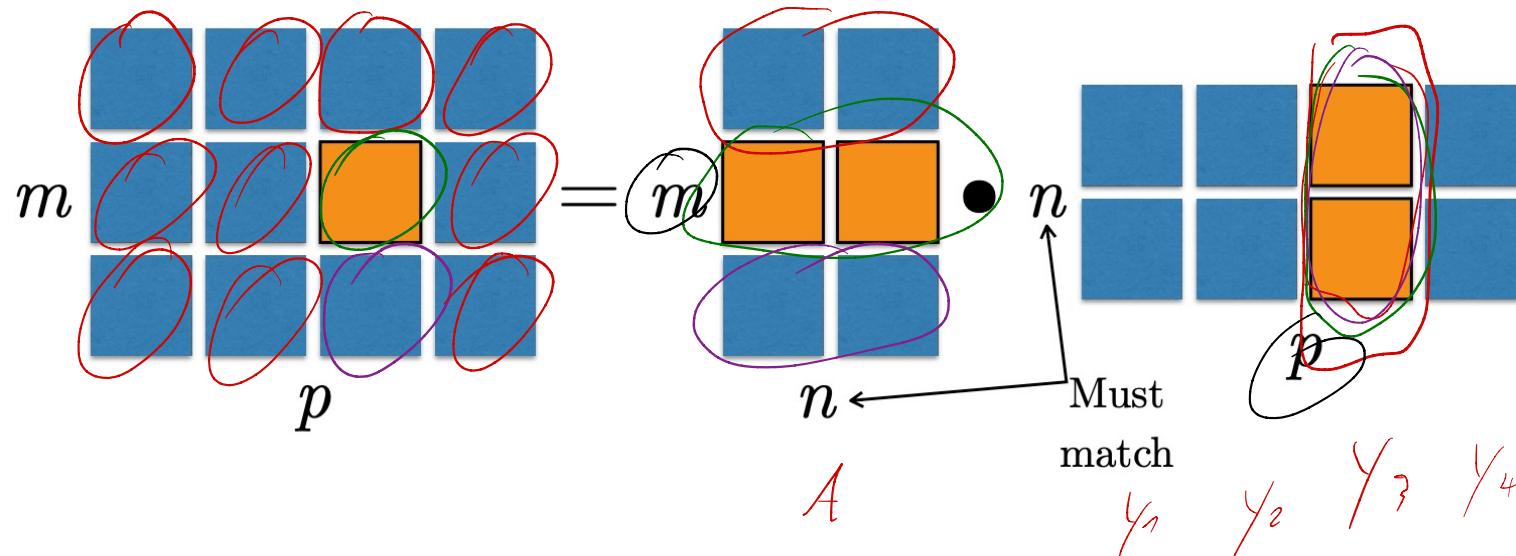
$$C \in \mathbb{R}^{m \times p}$$

$$C = \underline{\underline{AB}}.$$

$$A = \begin{pmatrix} 1 & 2 \\ 3 & 4 \\ 5 & 6 \end{pmatrix} \quad \begin{pmatrix} 5 & 11 & \dots & - \\ 11 & 25 & \dots & - \\ \dots & \dots & \dots & \dots \end{pmatrix} = C$$

rows from A  
columns from B

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



$$C \in \mathbb{R}^{m \times p}$$

$$\mathbb{R}^{3 \times 4}$$

Source: Goodfellow

## A few important properties of matrix multiplication

- Associative:  $(\underline{AB})C = A(\underline{BC})$
- Distributive:  $\underline{A}(\underline{B + C}) = \underline{AB} + \underline{AC}$
- Generally, not commutative:  $\underline{AB} \neq \underline{BA}$   
(also, just because  $AB$  exists, it doesn't mean that  $BA$  exists)

$$\begin{aligned} A &\in \mathbb{R}^{n \times n} & n \neq p \\ B &\in \mathbb{R}^{m \times p} \\ \Rightarrow AB &\in \mathbb{R}^{n \times p} \\ BA &\cancel{\in \mathbb{R}^m} \end{aligned}$$

$$3 \times 4 = 4 \times 3$$

## Matrix inversion

- A matrix  $A^{-1}$  such that  $A^{-1}A = I_n$

- Here,  $I_n$  is the “identity matrix” of rank  $n$ . For example,  $I_3 =$

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

- Why identity matrix? Because  $I_n x = x$  for all  $x \in \mathbb{R}^n$  and  $A I_n = A = I_n A$

## Matrix inversion

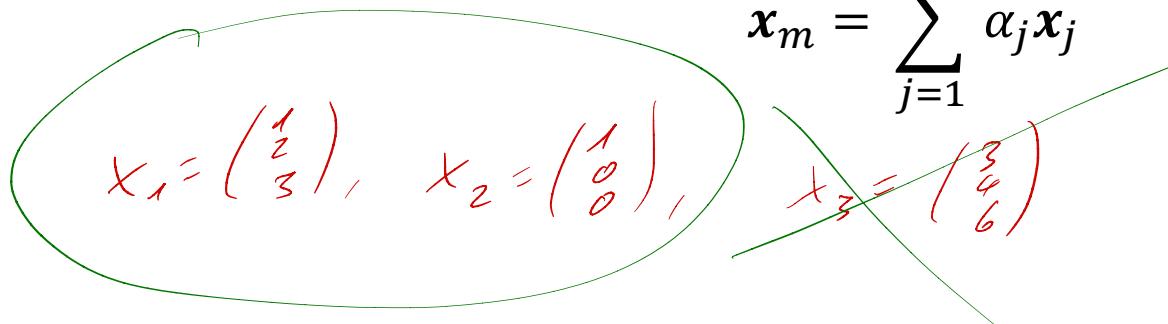
- A matrix  $A^{-1}$  such that  $A^{-1}A = I_n$
- Here,  $I_n$  is the “identity matrix” of rank  $n$ . For example,  $I_3 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$ 
  - Why identity matrix? Because  $I_n x = x$  for all  $x \in \mathbb{R}^n$  and  $A I_n = A = I_n A$
  - Note that  $A^{-1}$  only exists if the number of rows = the number of columns and both rows and columns are **linearly independent**
  - A few properties (if the inverses exist):
    - $(A^{-1})^{-1} = A$
    - $(AB)^{-1} = B^{-1}A^{-1}$
    - $(A^{-1})^T = (A^T)^{-1} = "A^{-T}"$
    - If  $A^T A = I$ , (or  $A^T = A^{-1}$ ) we say that  $A$  is an orthogonal matrix

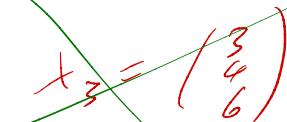
## Linear independence

$$x_1 = \begin{pmatrix} x_{1,1} \\ \vdots \\ x_{n,1} \end{pmatrix} \dots x_m = \begin{pmatrix} x_{1,m} \\ \vdots \\ x_{n,m} \end{pmatrix}$$

- We say that vectors  $x_1, x_2, \dots, x_m \in \mathbb{R}^n$  are linearly **dependent** if any of the vectors can be written as a linear combination of the others, e.g., if, for some scalars  $\alpha_1, \alpha_2, \dots, \alpha_m \in \mathbb{R}$ ,

$$x_m = \sum_{j=1}^{m-1} \alpha_j x_j$$


$$x_1 = \begin{pmatrix} 1 \\ 2 \\ 3 \end{pmatrix}, \quad x_2 = \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}$$

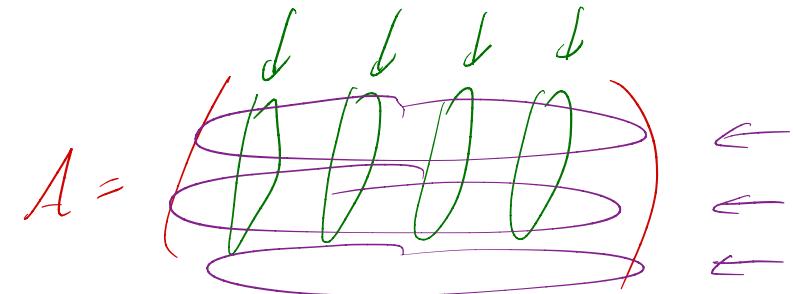

$$x_3 = \begin{pmatrix} 3 \\ 4 \\ 6 \end{pmatrix}$$

$$x_3 = 2x_1 + 1x_2$$

## Linear independence

- We say that vectors  $x_1, x_2, \dots, x_m \in \mathbb{R}^n$  are linearly **dependent** if any of the vectors can be written as a linear combination of the others, e.g., if, for some scalars  $\alpha_1, \alpha_2, \dots, \alpha_m \in \mathbb{R}$ ,

$$x_m = \sum_{j=1}^{m-1} \alpha_j x_j$$



- If this is **not** the case, we say that the vectors are linearly **independent**
- Going back to our matrix  $A$ , we say that
  - the column rank is the largest number of columns that are linearly independent
  - the row rank is the largest number of rows that are linearly independent
- It turns out, row rank = column rank = “rank”, so also  $\text{rank}(A) = \text{rank}(A^T)$
- $A$  is **invertible** if it is of “full rank” (that is,  $\text{rank}(A) = m = n$ )

## Using matrix inversions to solve systems of equations

- Say you are interested in finding the vector  $x$  that solves  $Ax = b$
- This can also be written as a system of  $m$  equations:

$$\underline{A_{1,:}x = b_1}$$

$$\underline{A_{2,:}x = b_2}$$

...

$$\underline{A_{m,:}x = b_m}$$

- This system can have
  - No solution
  - Many solutions
  - Exactly one solution (if  $A$  is invertible):  $Ax = b$

$$A^{-1}Ax = \underline{A^{-1}b}$$

$$I_n x = A^{-1}b$$

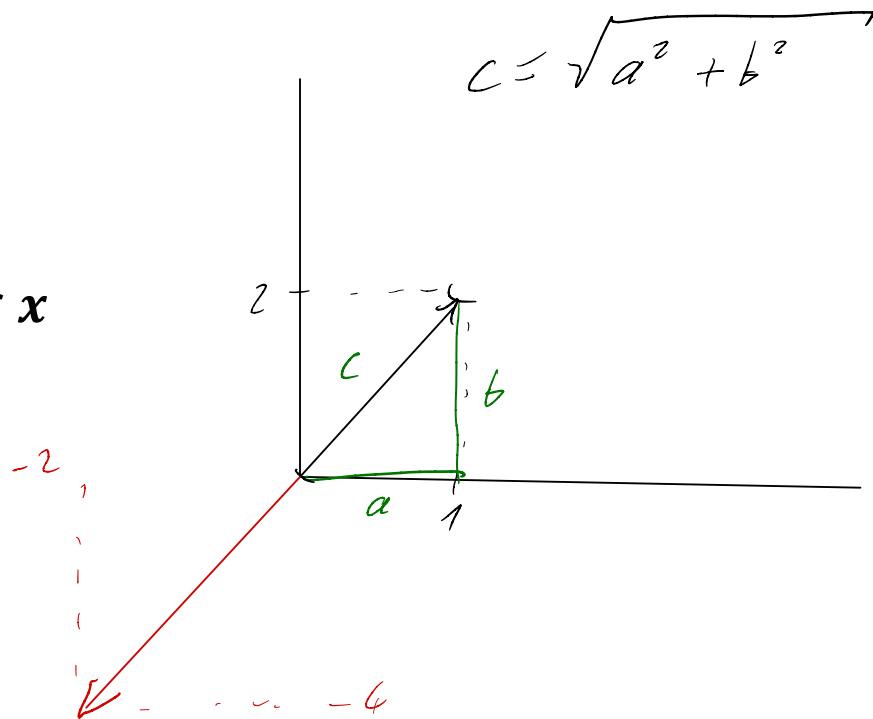
- Note: this way of solving the system is numerically unstable

Source: Goodfellow

# Linear algebra – norms

## Norms

- Functions  $f$  that measures the “length” of a vector  $x$
- Such functions need to fulfill four conditions:
  - $f(x) \geq 0$
  - $f(x) = 0 \Leftrightarrow x = 0$
  - For all  $x \in \mathbb{R}^n, \alpha \in \mathbb{R}, f(\alpha x) = |\alpha|f(x)$
  - For all  $x, y \in \mathbb{R}^n, f(x + y) \leq f(x) + f(y)$  (triangle inequality)



## Some commonly used norms

$$\|x\|_2 = \left( \sum_i |x_i|^2 \right)^{\frac{1}{2}} = \sqrt{\sum_i |x_i|^2}$$

- $L^p$  norm:

- $\|x\|_p = (\sum_i |x_i|^p)^{\frac{1}{p}}$

- Most commonly used norm:  $\|x\|_2 = \sqrt{\sum_i x_i^2}$  ( $L^2$  norm or “Euclidian” norm)
- Quite common as well:  $\|x\|_1 = \sum_i |x_i|$  ( $L^1$  norm)
- An extreme case: the max-norm  $\|x\|_\infty = \max_i |x_i|$
- For a given norm, we call the vector  $x$  with  $\|x\| = 1$  the “unit vector”
  - Note: Let  $y = \frac{x}{\|x\|}$ , then  $\|y\| = 1$

$$x = \begin{pmatrix} 1 \\ 2 \end{pmatrix}, \|x\|_2 = \sqrt{1^2 + 2^2} = \sqrt{5}$$

$$y = \begin{pmatrix} 1/\sqrt{5} \\ 2/\sqrt{5} \end{pmatrix} \Rightarrow \|y\|_2 = \frac{1}{\sqrt{5}} \|x\|_2 = \frac{\sqrt{5}}{\sqrt{5}}$$

## Norms defined for matrices

$$\cancel{\|A\|_2} \quad \|A\|_F$$

- There are many matrix norms, but we will only need one: the Frobenius norm

$$\|A\|_F = \sqrt{\sum_{i=1}^n \sum_{j=1}^m A_{i,j}^2} = \sqrt{Tr(A^T A)}$$

- Note the similarity with the  $L^2$  norm for vectors

# **Linear algebra – singular value decomposition**

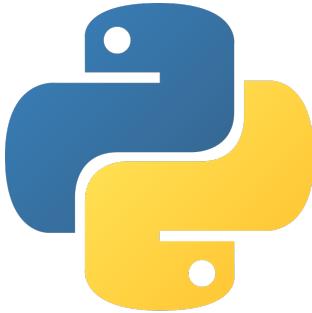
## Eigenvectors, eigenvalues, and singular values

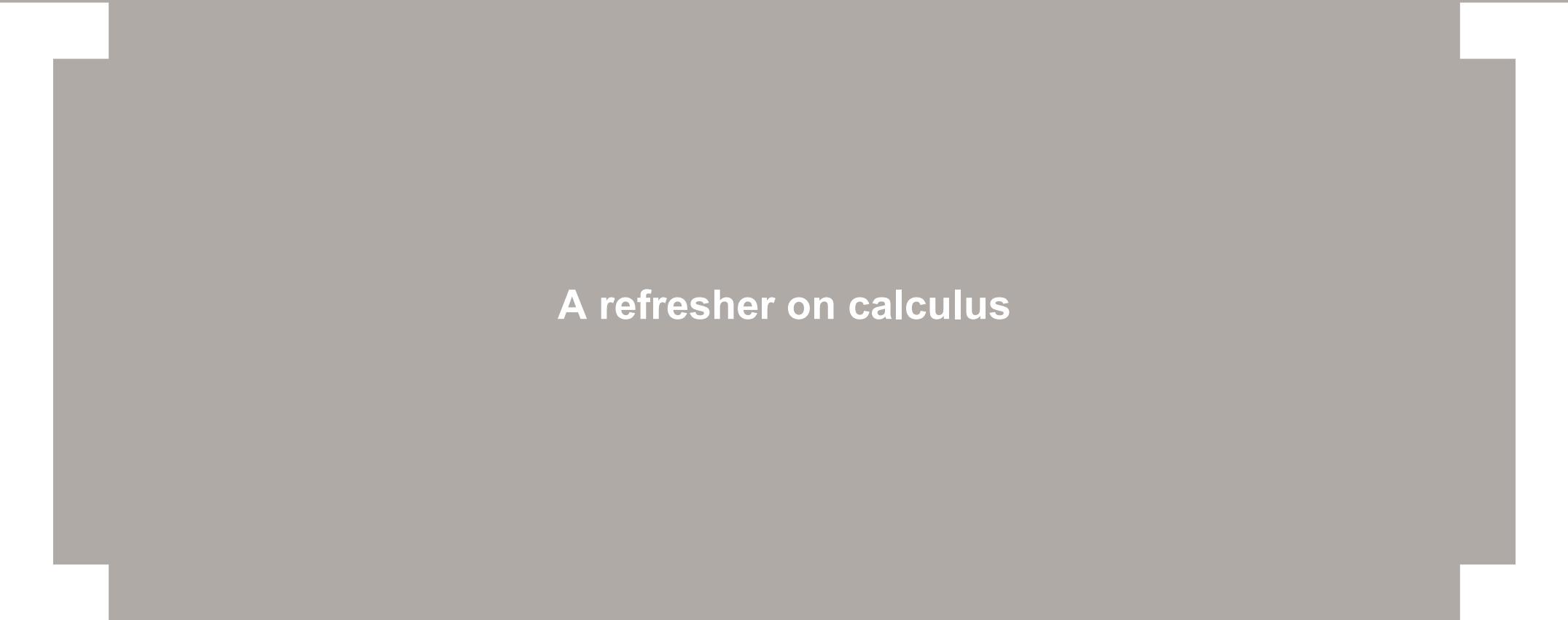
$$\underbrace{\begin{pmatrix} 4 & 2 \\ 2 & 1 \end{pmatrix}}_A \underbrace{\begin{pmatrix} 2 \\ 1 \end{pmatrix}}_v = \begin{pmatrix} 10 \\ 5 \end{pmatrix} = \underbrace{5}_{\lambda} \underbrace{\begin{pmatrix} 2 \\ 1 \end{pmatrix}}_v$$

- If  $A\mathbf{v} = \lambda\mathbf{v}$ , where  $\mathbf{v}$  is a non-zero vector, then we say  $\mathbf{v}$  is an eigenvector of square-matrix  $A$  (with eigenvalue  $\lambda$ )
- Singular values: non-negative square roots of eigenvalues of  $A^T A$ . Say  $\sigma_i$ ,  $i = 1, \dots, m$
- Singular value decomposition: If  $A \in \mathbb{R}^{n \times m}$ , there exists orthogonal matrices  $U \in \mathbb{R}^{n \times n}$  and  $V \in \mathbb{R}^{m \times m}$  such that

$$U^{-1} A V = \Sigma = \begin{bmatrix} \sigma_1 & 0 & 0 & 0 \\ 0 & \sigma_1 & 0 & 0 \\ 0 & 0 & \ddots & 0 \\ 0 & 0 & 0 & \sigma_m \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

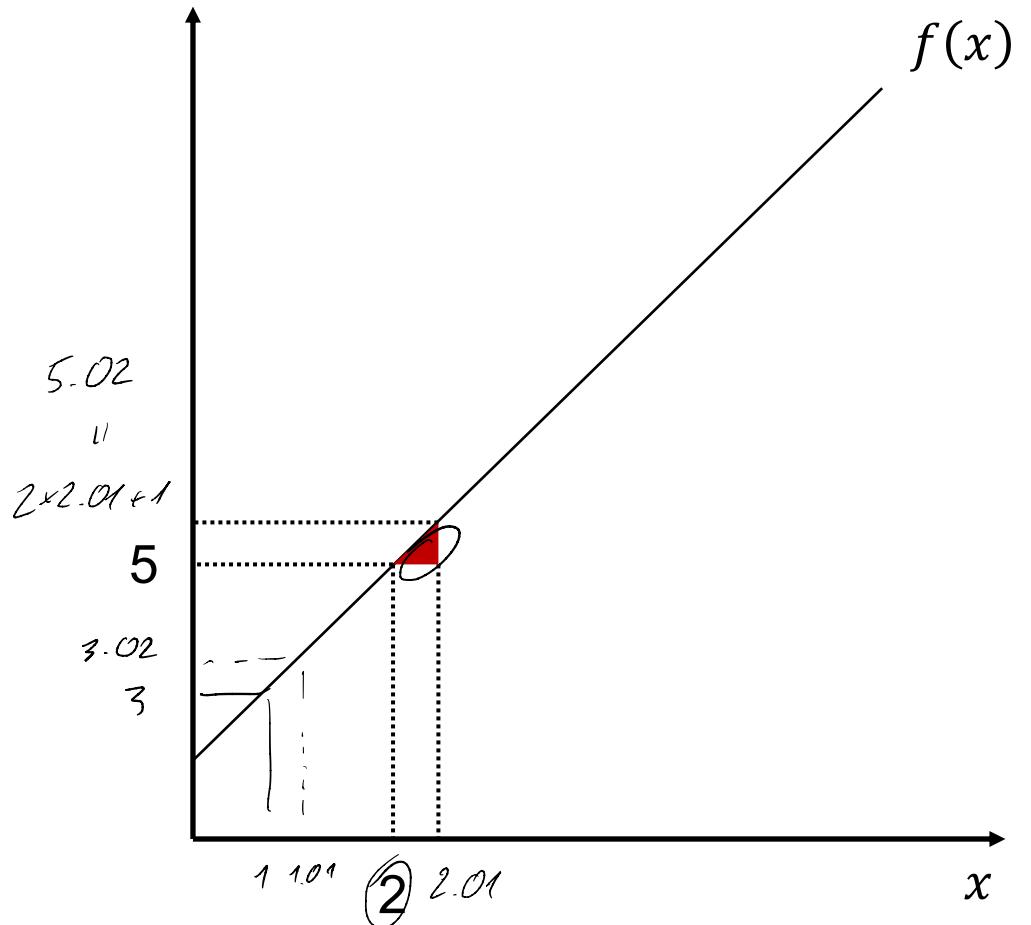
# Matrix multiplication and co in Python





**A refresher on calculus**

# What is a derivative?



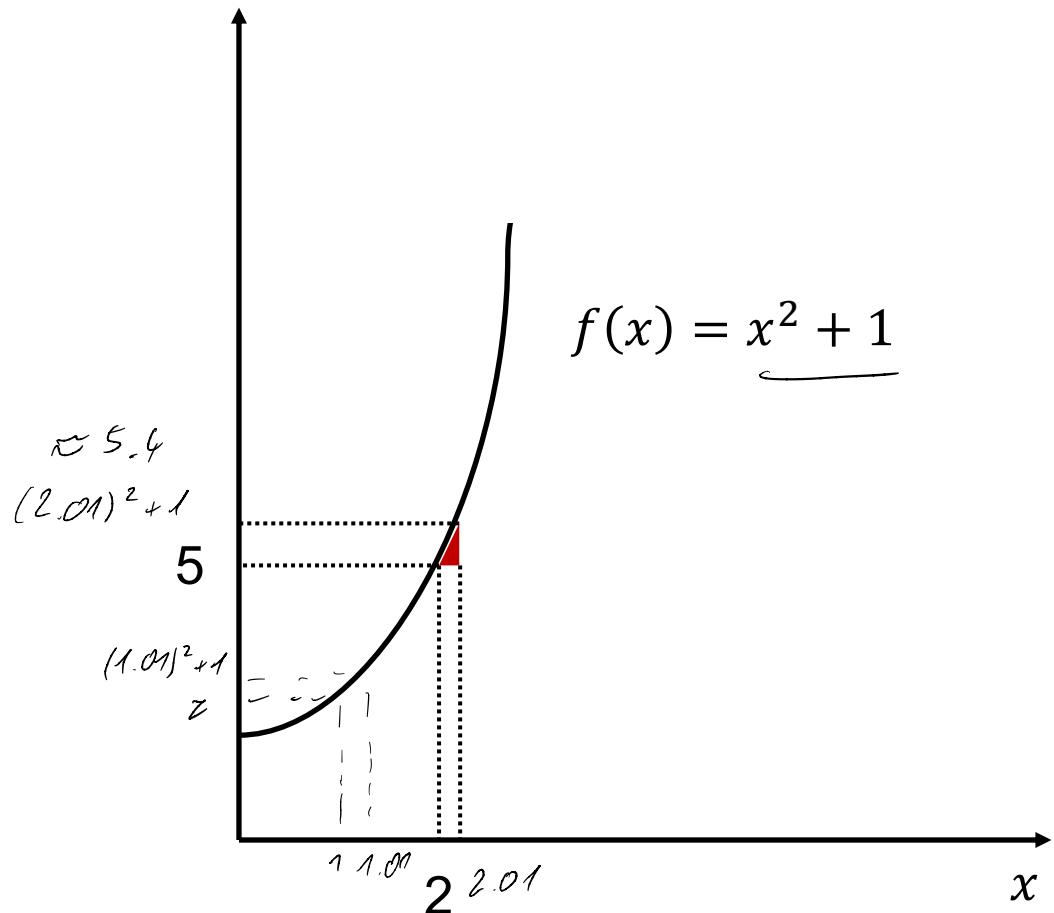
Slope (derivative) of  $f(x)$  at 2 is  $\frac{\text{"height"}}{\text{"width"}}$

$$\frac{d f(x)}{dx} = f'(x) \text{ at } 2 \approx \frac{5.02 - 5}{2.01 - 2} = \frac{0.02}{0.01} = 2$$

$$\underline{\Delta y} \approx \frac{3.02 - 3}{1.01 - 1} = 2$$



# What is a derivative?



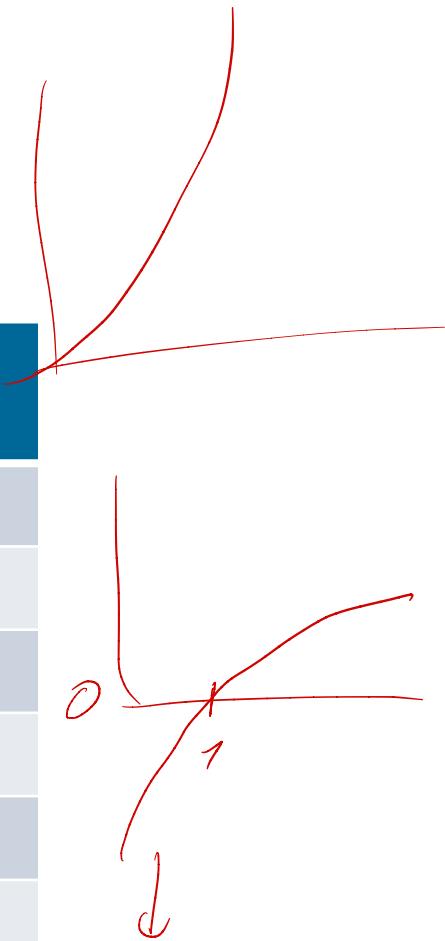
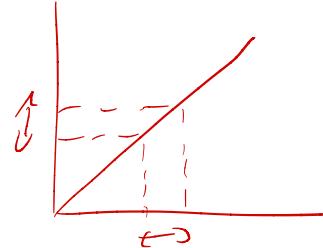
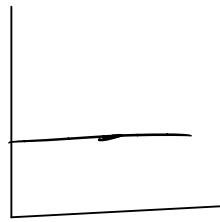
Slope (derivative) of  $f(x)$  at 2 is  $\frac{\text{"height"}}{\text{"width"}}$

$$\frac{\cancel{x+1}}{2} \quad \frac{(2.01)^2 + 1 - 5}{2.01 - 2} \approx \frac{0.4}{0.1} = 4$$

$$\cancel{x+1} \quad \frac{0.2}{0.1} = 2$$

$$f'(x) = 2x$$

## A few important derivatives



$f(x)$	$f'(x) = \frac{df(x)}{dx} = \frac{d}{dx} f(x)$
1	0
$x$	1
$x^2$	$2x$
$x^3 - 1$	$3x^2$
$\sqrt{x} = x^{\frac{1}{2}}$	$\frac{1}{2}x^{-\frac{1}{2}}$
$\ln(x)$	$\frac{1}{x}$
$e^x$	$e^x$



## Some rules about handling derivatives

$$f'(x) = \frac{df(x)}{dx} \quad \frac{df}{dx} = \frac{df}{dg} \frac{dg}{dx}$$

$$h(x) = f(g(x)) = f(x^2)$$

$$f(y) = e^y$$

$h(x)$	$h'(x)$	Example
$c f(x)$	$\underline{c} \underline{f'(x)}$	$h(x) = 18 x^k \quad h'(x) = 18 k x^{k-1}$
$f(x) + g(x)$	$\underline{f'(x)} + \underline{g'(x)}$	$h(x) = \ln(x) - x^2 + 5 \quad h'(x) = \frac{1}{x} + (-2x)$
$f(x)g(x)$	$\underline{f'(x)} \underline{g(x)} + \underline{f(x)} \underline{g'(x)}$	$h(x) = \frac{2e^x x}{c f(x) g(x)} \quad h'(x) = 2[e^x x + \underline{e^x}]$
$\frac{1}{f(x)}$	$-\frac{f'(x)}{f(x)^2}$	$h(x) = \frac{1}{\ln(x)} \quad h'(x) = -\frac{1/x}{(\ln x)^2}$
$\frac{f(x)}{g(x)}$	$\frac{f'(x)g(x) - f(x)g'(x)}{g(x)^2}$	$h(x) = \frac{\ln(x)}{x^2} \quad h'(x) = \frac{\frac{1}{x} x^2 - \log(x) 2x}{x^4}$
$f(g(x))$	$\underline{f'(g(x))} \underline{g'(x)}$	$h(x) = e^{\frac{x^2}{f(x)}} \quad h'(x) = e^{\frac{x^2}{f(x)}} 2x$

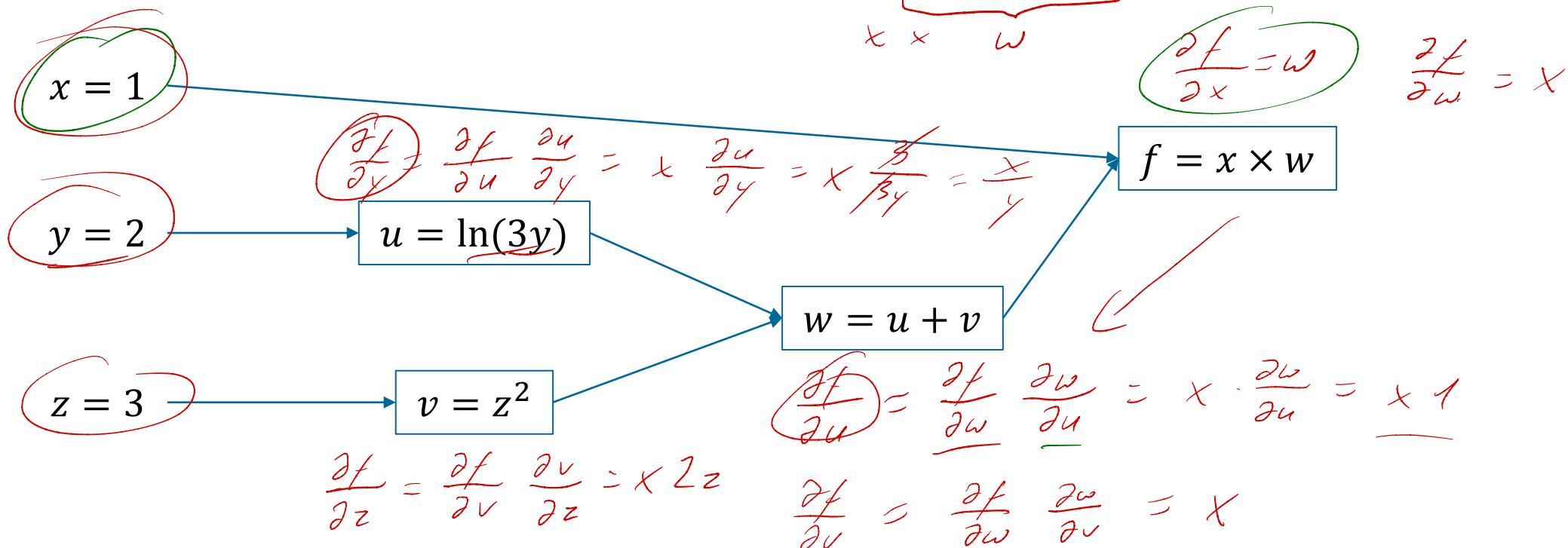
## Using a computation graph for derivatives

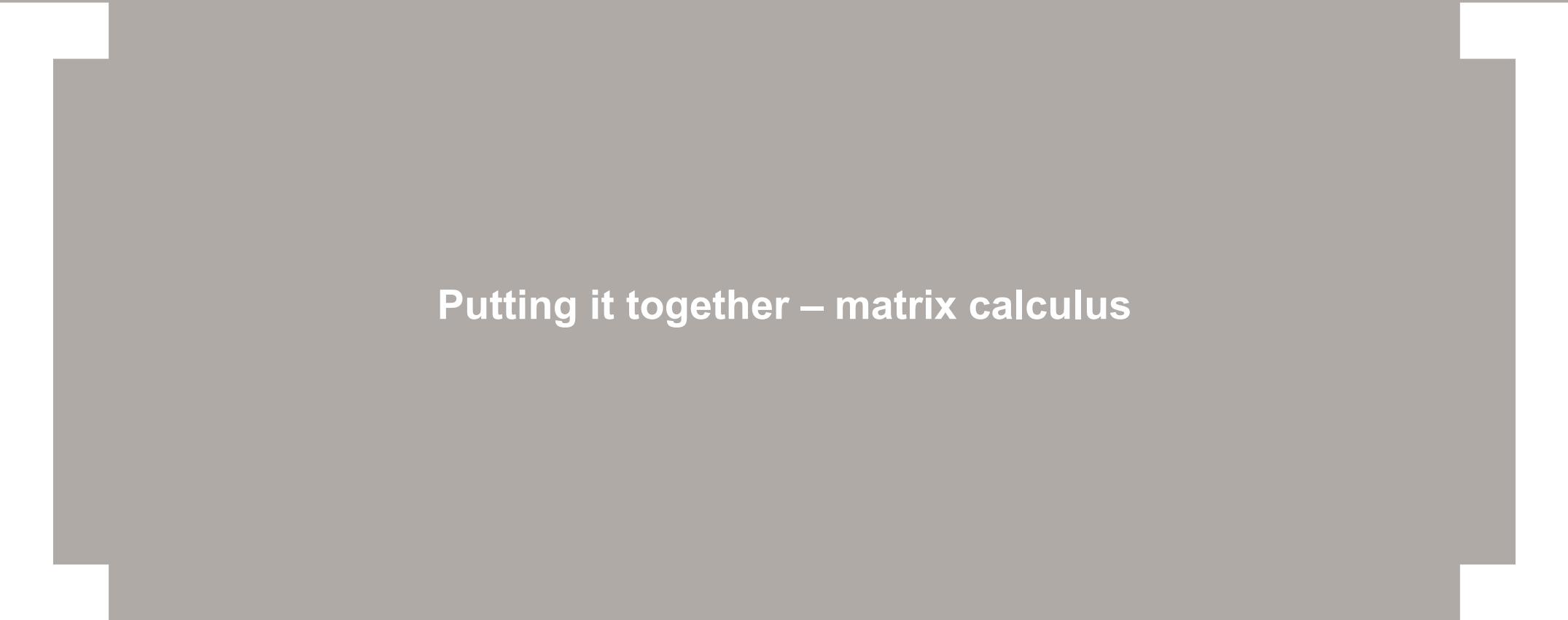
$$u = \ln(6)$$

$$v = 9 = 3^2$$

$$\omega = \underline{\ln(6) + 9}$$

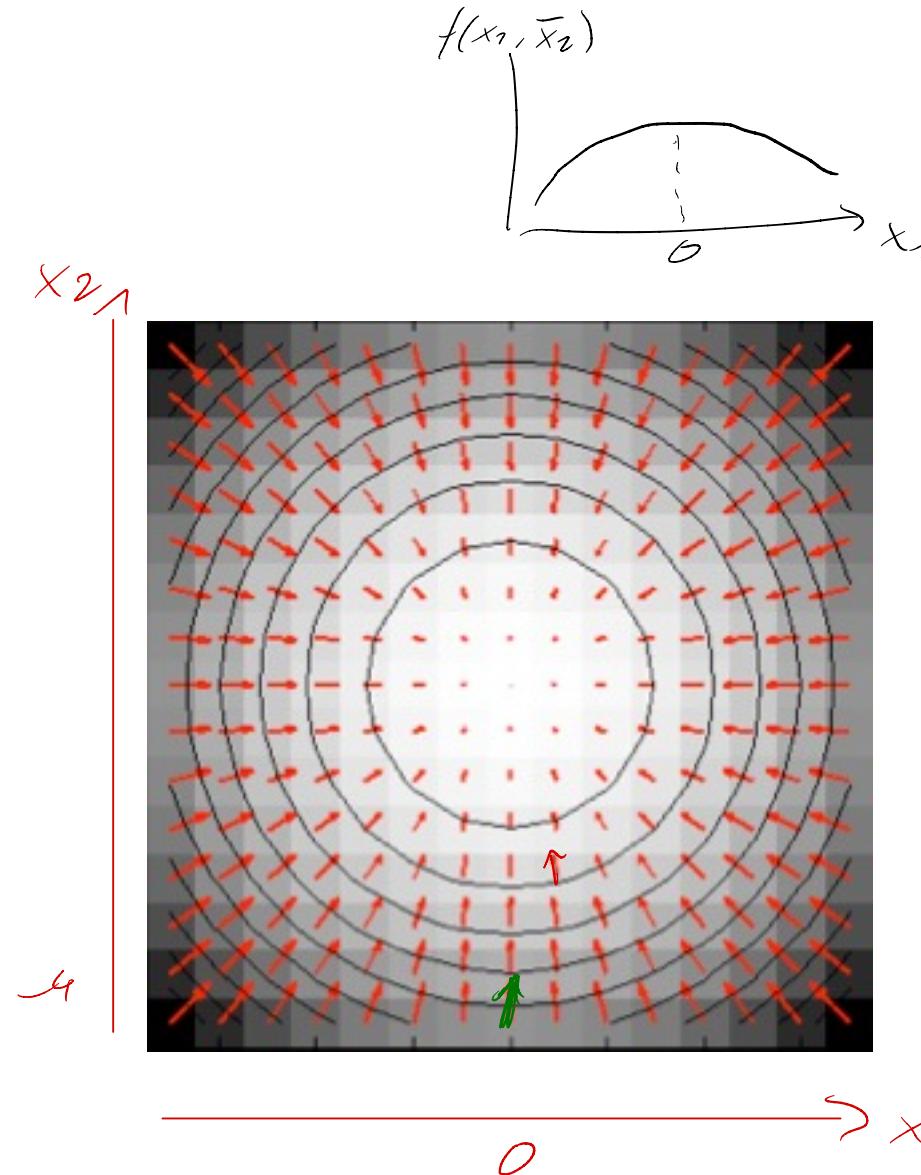
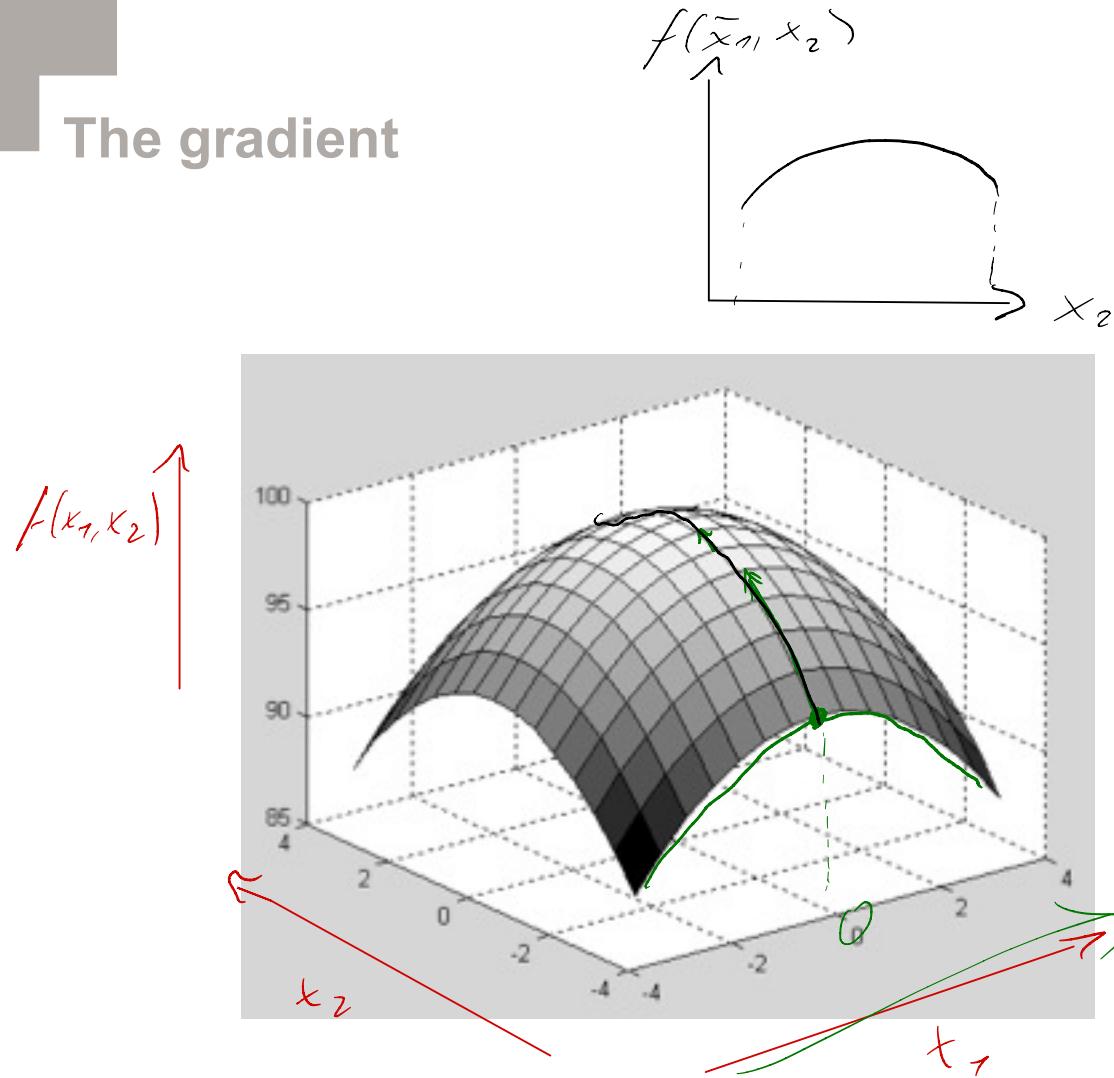
- We can use the “chain rule” to simplify the search for derivatives
- Say, we want to compute the derivatives of  $f = x(\ln(3y) + z^2)$  to  $x, y, z$  at  $x = 1, y = 2, z = 3$





**Putting it together – matrix calculus**

## The gradient



Source: Collins

## What are we seeing here?

- Say, we have a function  $f$ , taking as input a vector  $x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$
- The gradient (with respect to  $x$ ) is the vector pointing in the direction of fastest increase:

$$x = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$$

$$\nabla_x f(x) = \begin{pmatrix} \frac{\partial f(x)}{\partial x_1} \\ \frac{\partial f(x)}{\partial x_2} \end{pmatrix}$$

$$\nabla_x f(x) = \begin{bmatrix} \frac{\partial f(x)}{\partial x_1} \\ \frac{\partial f(x)}{\partial x_2} \\ \vdots \\ \frac{\partial f(x)}{\partial x_n} \end{bmatrix}$$

## What are we seeing here?

- Say, we have a function  $f$ , taking as input a vector  $\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$
- The gradient (with respect to  $\mathbf{x}$ ) is the vector pointing in the direction of fastest increase:
- Based on the previous ideas, note that:
  - $\nabla_{\mathbf{x}}(f(\mathbf{x}) + g(\mathbf{x})) = \nabla_{\mathbf{x}}f(\mathbf{x}) + \nabla_{\mathbf{x}}g(\mathbf{x})$
  - For all  $\alpha \in \mathbb{R}$ ,  $\nabla_{\mathbf{x}}(\alpha f(\mathbf{x})) = \alpha \nabla_{\mathbf{x}}f(\mathbf{x})$
  - If  $f(\mathbf{x}) = \sum_{i=1}^n b_i x_i = \mathbf{b}^T \mathbf{x}$ , then  $\nabla_{\mathbf{x}}f(\mathbf{x}) = \mathbf{b}$

$$\nabla_{\mathbf{x}}f(\mathbf{x}) = \begin{bmatrix} \frac{\partial f(\mathbf{x})}{\partial x_1} \\ \frac{\partial f(\mathbf{x})}{\partial x_2} \\ \vdots \\ \frac{\partial f(\mathbf{x})}{\partial x_n} \end{bmatrix}$$

## This naturally extends to functions that take as input a matrix

- Say, now we have a function  $f$ , taking as input a matrix  $\mathbf{A}$  =

$$f \left( \begin{pmatrix} 1 & 3 \\ 4 & 2 \end{pmatrix} \right) = 6$$

$$f(\mathbf{A}) = 5$$

- The gradient (with respect to  $\mathbf{A}$ ) is now also a matrix  $\nabla_{\mathbf{A}} f(\mathbf{A})$  =

$$\begin{bmatrix} a_{1,1} & a_{1,2} & \dots & a_{1,m} \\ a_{2,1} & a_{2,2} & \dots & a_{2,m} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n,1} & a_{n,2} & \dots & a_{n,m} \end{bmatrix}$$

$$\begin{bmatrix} \frac{\partial f(\mathbf{A})}{\partial a_{1,1}} & \frac{\partial f(\mathbf{A})}{\partial a_{1,2}} & \dots & \frac{\partial f(\mathbf{A})}{\partial a_{1,m}} \\ \frac{\partial f(\mathbf{A})}{\partial a_{2,1}} & \frac{\partial f(\mathbf{A})}{\partial a_{2,2}} & \dots & \frac{\partial f(\mathbf{A})}{\partial a_{2,m}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f(\mathbf{A})}{\partial a_{n,1}} & \frac{\partial f(\mathbf{A})}{\partial a_{n,2}} & \dots & \frac{\partial f(\mathbf{A})}{\partial a_{n,m}} \end{bmatrix}$$



**Taking a step back – logistic regression**

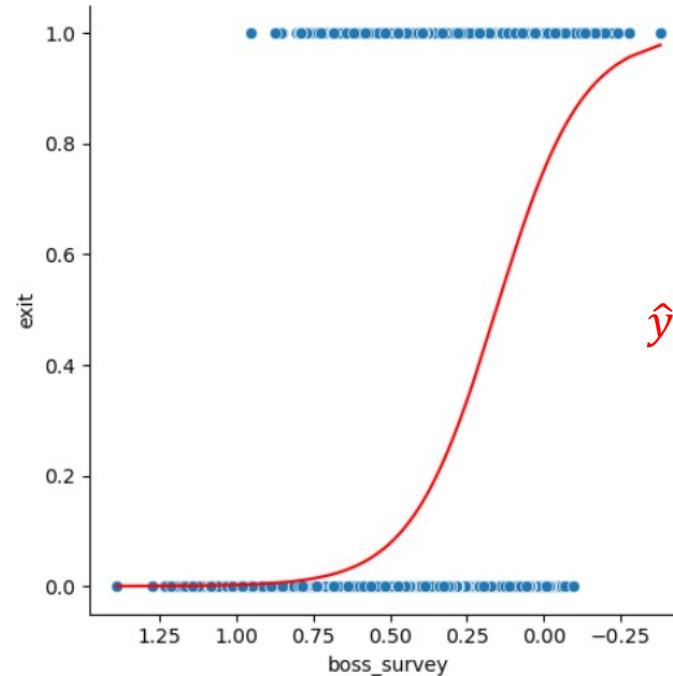
## What do we actually do when training a logistic regression model?

- We are given values  $(\mathbf{x}^{(i)}, y^{(i)})$ , where  $\mathbf{x}^{(i)} \in \mathbb{R}^m$  and  $y^{(i)} \in \{0,1\}$
- Our prediction  $\hat{y}^{(i)}$  should reflect the probability that  $y^{(i)} = 1$ :  $\hat{y}^{(i)} = P(y^{(i)} = 1 | \mathbf{x}^{(i)})$
- We model this probability, using the sigmoid function:

$$\mathbf{X} = \begin{pmatrix} 1 & \dots & \mathbf{x}^{(1)} \\ & \ddots & \vdots \\ 1 & \dots & \mathbf{x}^{(n)} \end{pmatrix}$$
$$\mathbf{Y} = \begin{pmatrix} y^{(1)} \\ \vdots \\ y^{(n)} \end{pmatrix}$$

## What do we actually do when training a logistic regression model?

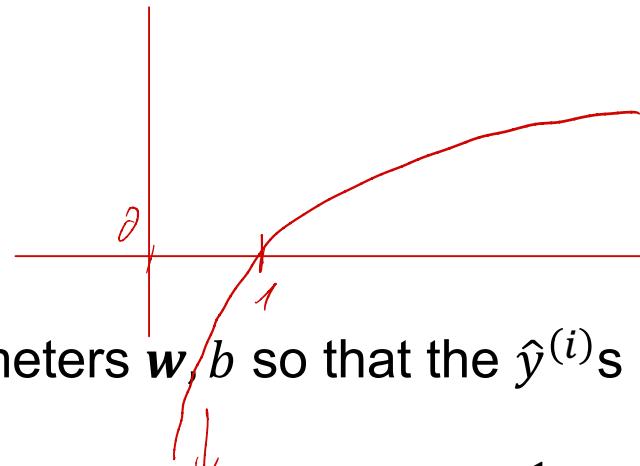
- We are given values  $(x^{(i)}, y^{(i)})$ , where  $x^{(i)} \in \mathbb{R}^m$  and  $y^{(i)} \in \{0,1\}$
- Our prediction  $\hat{y}^{(i)}$  should reflect the probability that  $y^{(i)} = 1$ :  $\hat{y}^{(i)} = P(y^{(i)} = 1 | x^{(i)})$
- We model this probability, using the sigmoid function:



$$\begin{aligned} & \alpha_0 + \alpha_1 x_1 + \alpha_2 x_2 + \dots + \alpha_m x_m \\ & b + w_1 x_1 + w_2 x_2 + \dots + w_m x_m \\ & \textcircled{x} \textcircled{w} + \textcircled{b} \\ & (1, m) \quad (m, 1) \quad (1, 1) \\ \hat{y} &= \frac{1}{1 + e^{-(x \cdot w + b)}} = \frac{e^{(x \cdot w + b)}}{1 + e^{(x \cdot w + b)}} \end{aligned}$$

## The optimization part

- Remember that  $w \in \mathbb{R}^m$  and  $b \in \mathbb{R}$
- To get to the “right” model, we optimize our parameters  $w, b$  so that the  $\hat{y}^{(i)}$ s are “as close as possible” to the  $y^i$ s
- What we do is to minimize the “cost-function”  $J(w, b)$ , where  $\hat{y}^{(i)} = \frac{1}{1+e^{-(x^{(i)}w+b)}}$ :



$$J(w, b) = -\frac{1}{n} \sum_{i=1}^n [y^{(i)} \ln \hat{y}^{(i)} + (1 - y^{(i)}) \ln(1 - \hat{y}^{(i)})]$$

$\hat{y}^{(i)} = 1 : \cancel{0} \left[ \ln \cancel{\hat{y}^{(i)}} + 0 \right] \rightarrow 0 : \hat{y}^{(i)} \rightarrow 1$

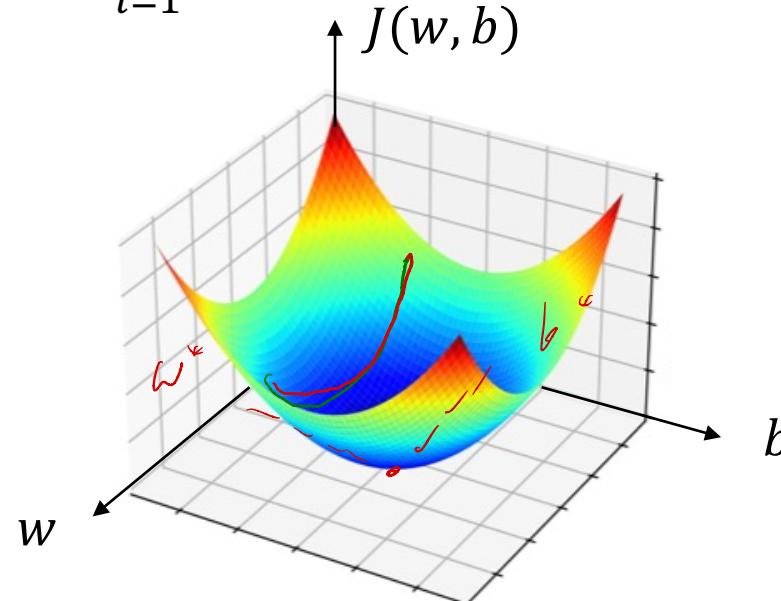
$\hat{y}^{(i)} = 0 : -[0 + \ln(1 - \cancel{\hat{y}^{(i)}})] \rightarrow 0 : \hat{y}^{(i)} \rightarrow 0$



## The optimization part

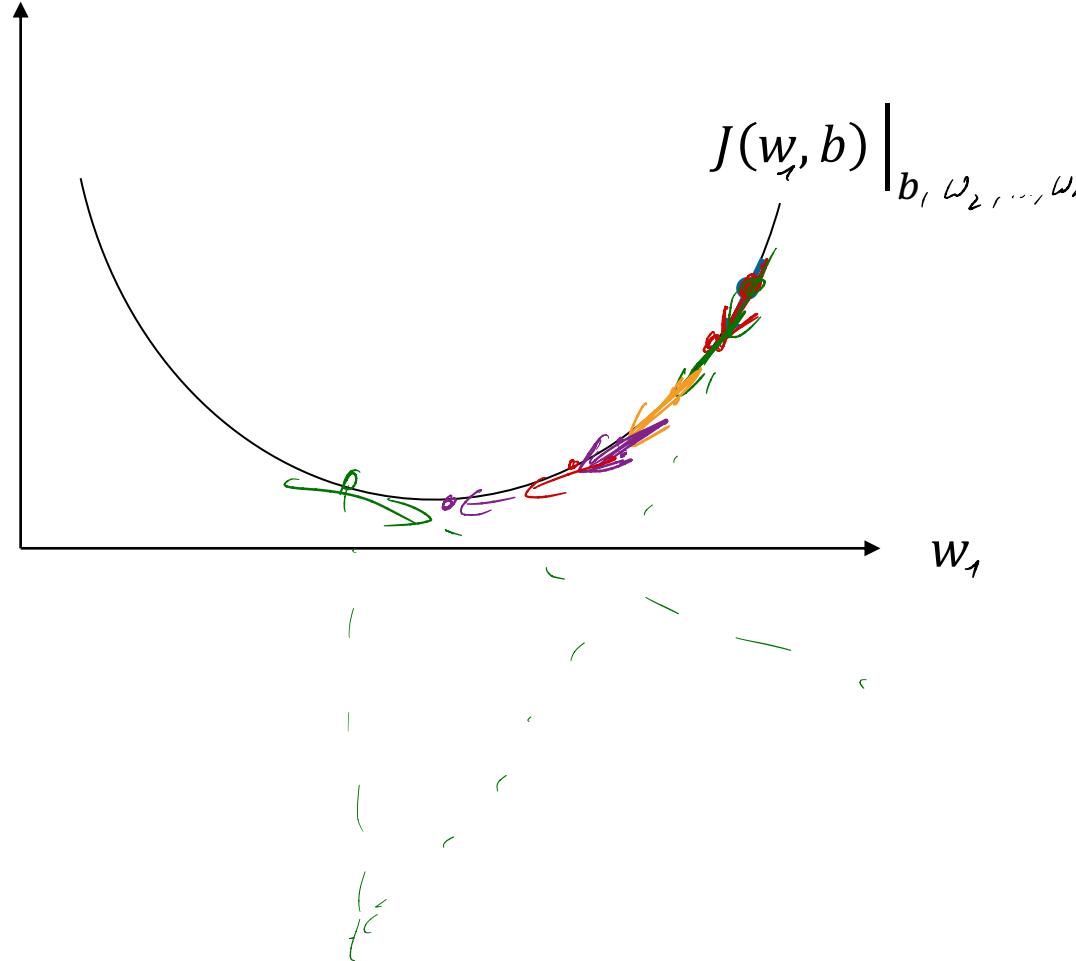
- Remember that  $w \in \mathbb{R}^m$  and  $b \in \mathbb{R}$
- To get to the “right” model, we optimize our parameters  $w, b$  so that the  $\hat{y}^{(i)}$ s are “as close as possible” to the  $y^i$ s
- What we do is to minimize the “cost-function”  $J(w, b)$ , where  $\hat{y}^{(i)} = \frac{1}{1+e^{-(x^{(i)}w+b)}}$ :

$$J(w, b) = -\frac{1}{n} \sum_{i=1}^n [y^{(i)} \ln \hat{y}^{(i)} + (1 - y^{(i)}) \ln(1 - \hat{y}^{(i)})]$$



$$\hat{y}^{(i)} = \frac{1}{1 + e^{-w^*x^{(i)} + b^*}}$$

## Solving the optimization problem through gradient descent



$$w_1 \rightsquigarrow \frac{\partial J}{\partial w_1}$$

$$(w_1) := w_1 - \alpha \frac{\partial J}{\partial w_1}$$

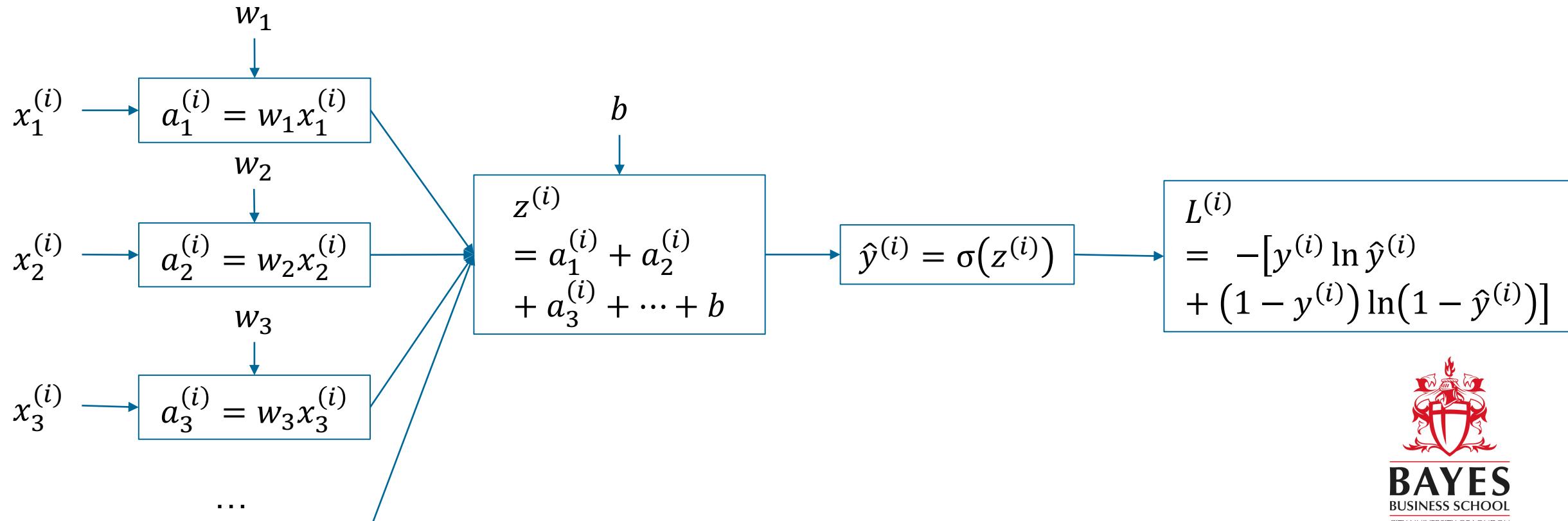
$$\rightarrow w_1 := w_1 - \alpha \frac{\partial J}{\partial w_1} +$$

## Our first optimization algorithm

1. Decide a “learning rate”  $\alpha$
2. Start with some  $w$  and  $b$  and compute  $J(w, b)$
3. Until  $J$  “doesn’t change” anymore:
  - Let  $w_1 := w_1 - \alpha \frac{\partial J(w, b)}{\partial w_1}$
  - Let  $w_2 := w_2 - \alpha \frac{\partial J(w, b)}{\partial w_2}$
  - ...
  - Let  $w_m := w_m - \alpha \frac{\partial J(w, b)}{\partial w_m}$
  - Let  $b := b - \alpha \frac{\partial J(w, b)}{\partial b}$
  - Recompute  $J(w, b)$
4. Enjoy the fruits of your labor: you have fit a logistic regression model manually!

## Wait a second, how do we find all those derivatives?

- We can use again the computation graph!
- Recall that  $\hat{y}^{(i)} = \frac{1}{1+e^{-(x^{(i)}w+b)}} = \sigma(x^{(i)}w + b)$



## As the same parameters influence all examples, we have to consider one final step

- Recall that  $J(\mathbf{w}, b) = -\frac{1}{n} \sum_{i=1}^n [y^{(i)} \ln \hat{y}^{(i)} + (1 - y^{(i)}) \ln(1 - \hat{y}^{(i)})] = \frac{1}{n} \sum_{i=1}^n L^{(i)}$
- We have that  $\frac{\partial J(\mathbf{w}, b)}{\partial w_j} = \frac{1}{n} \sum_{i=1}^n \frac{\partial L^{(i)}}{\partial w_j}$

We can now implement a logistic regression





See you next week!

## Sources

- Bhaskhar, 2021, Linear Algebra: <https://cs229.stanford.edu/notes2021fall/section1notes-linear-algebra-review.pdf>
- Collins, 2012, Intensity Surfaces and Gradients:  
[http://www.cse.psu.edu/~rtc12/CSE486/lecture02\\_6pp.pdf](http://www.cse.psu.edu/~rtc12/CSE486/lecture02_6pp.pdf)
- Goodfellow, Bengio, Courville, 2016, The Deep Learning Book:  
<http://www.deeplearningbook.org>
- Kolter, Do, & Ma, 2020, Linear Algebra Review and Reference:  
<https://cs229.stanford.edu/summer2020/cs229-linalg.pdf>
- Ivanovic, 2017, Python Introduction and Linear Algebra Review:  
<https://web.stanford.edu/class/cs231a/section/section1.pdf>