



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, π), rational numbers ($\frac{\text{integer}}{\text{integer}}$)

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

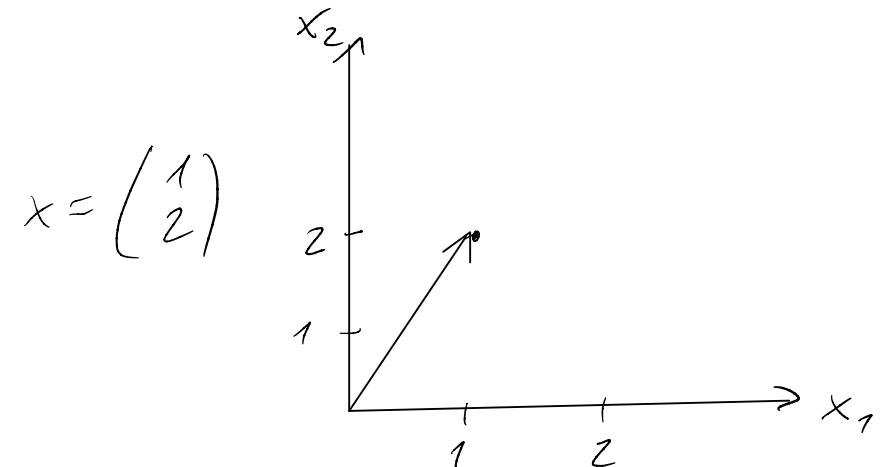
Vectors

- A one-dimensional array of numbers:

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

- 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:

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

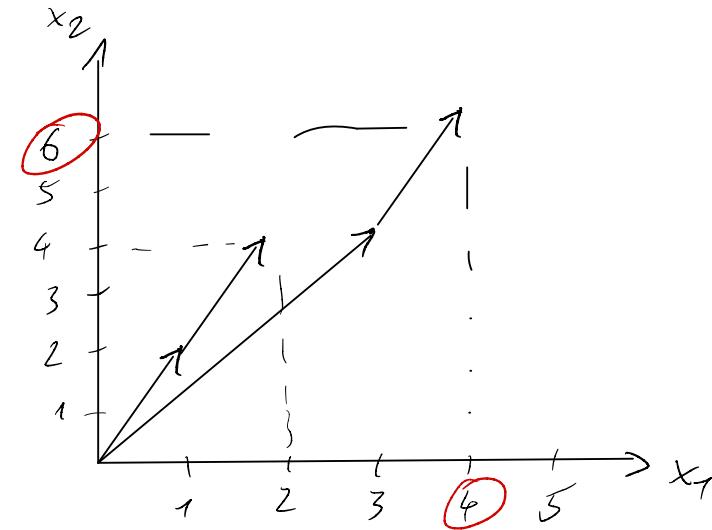
- Entries can be real numbers, binary, integer, ...
→ we usually denote a vector with the type of its entries and its size, e.g., $\mathbf{x} \in \mathbb{R}^n$

- 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

$n = 2$

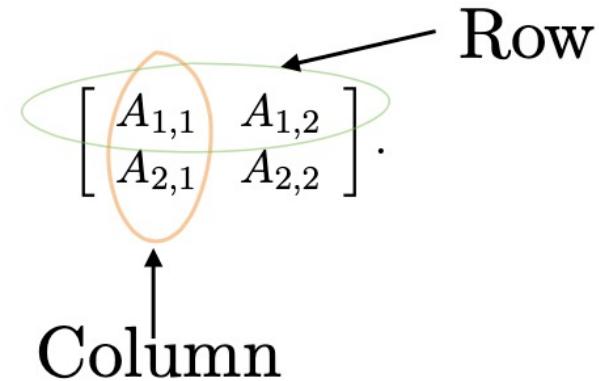
$m = 2$

- 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:

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

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

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

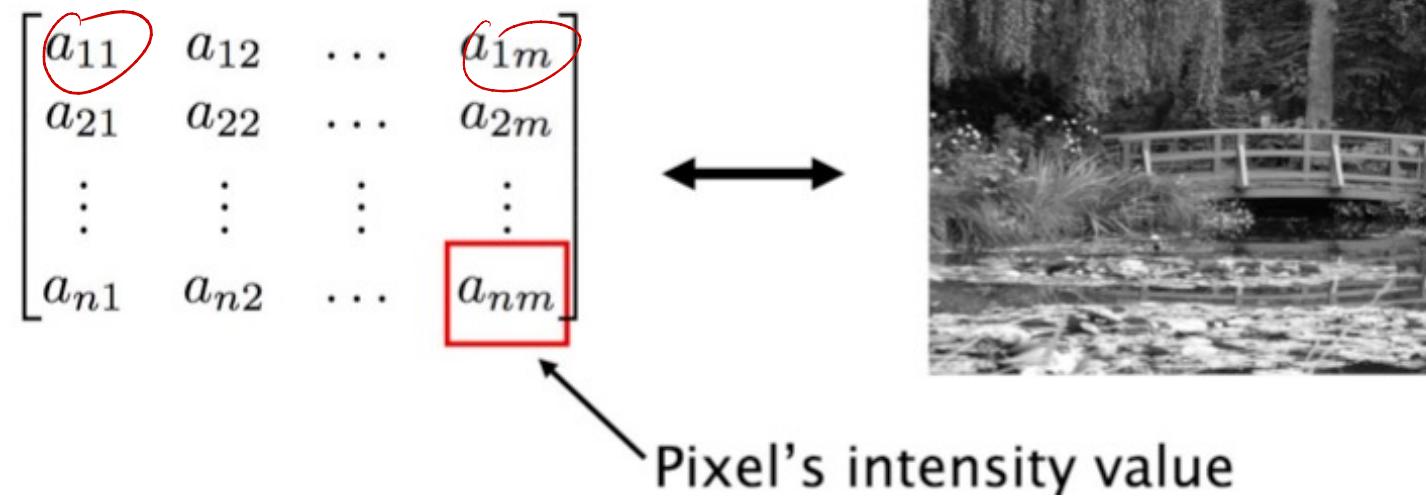
$$\begin{pmatrix} 1 & 2 \\ 3 & 4 \end{pmatrix} + \begin{pmatrix} 3 & 0 \\ 1 & 5 \end{pmatrix} \neq$$

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

Source: Goodfellow

A typical use of matrices in deep learning

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



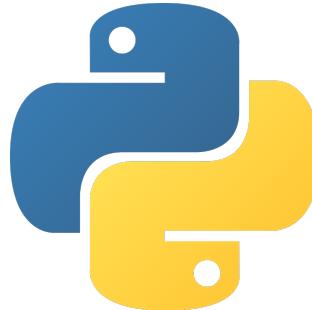
Source: Ivanovic

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

- 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:

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

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

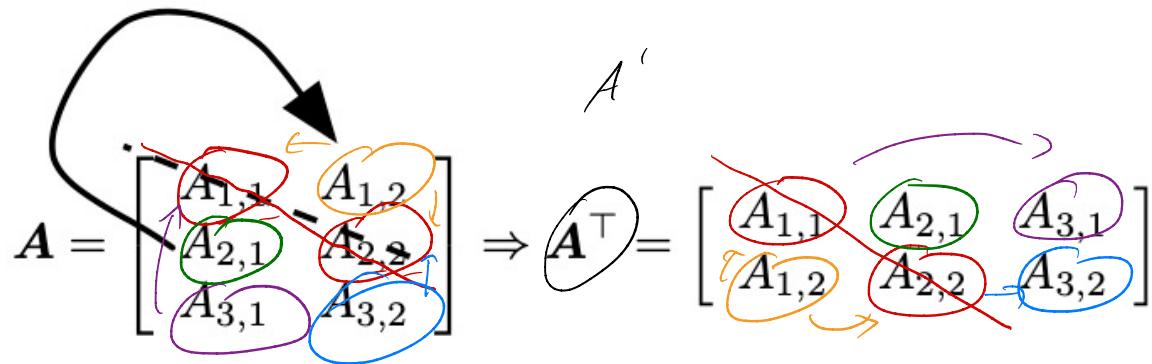
- Note that
 - $Tr(ABC) = Tr(CAB) = Tr(BCA)$
 - $Tr(A + B) = Tr(A) + Tr(B)$

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

Matrix transpose

- 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:

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

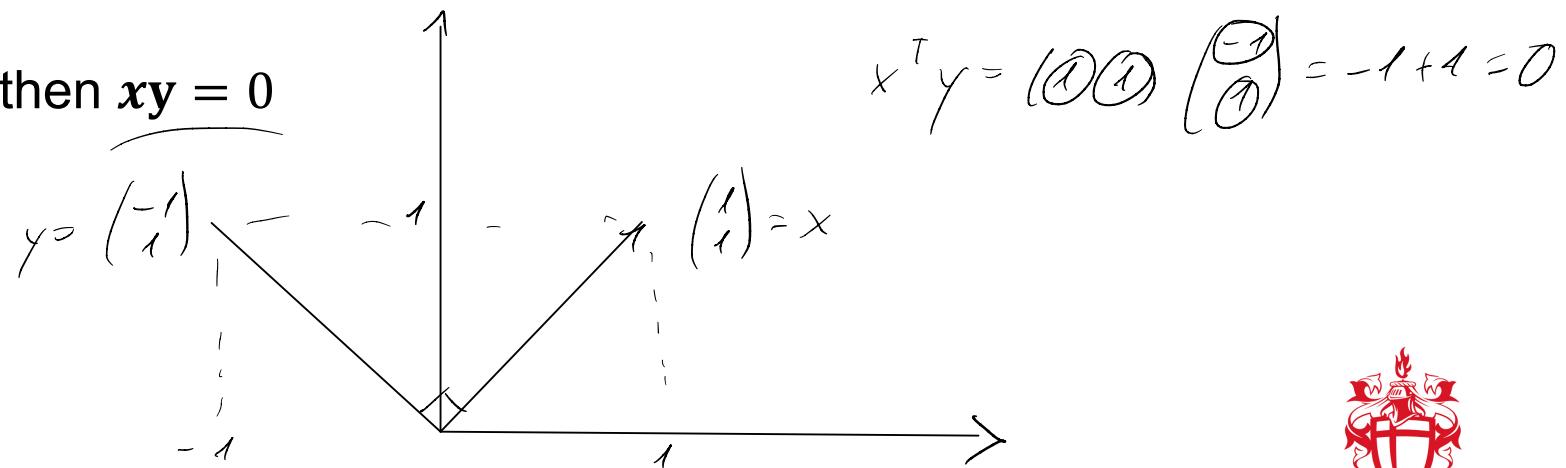
$$x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

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

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

- The inner product is a scalar!

- Note: if x and y are orthogonal, then $xy = 0$



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



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 of length n with entries x_1, x_2, \dots, x_n . Vector y is shown as a row of length m with entries y_1, y_2, \dots, y_m . An arrow labeled \rightsquigarrow points from the row y to the resulting matrix. The resulting matrix is a $n \times m$ matrix where each element $x_i y_j$ is highlighted with a red oval. The columns of the matrix are grouped by green ovals, and the rows are grouped by purple ovals.

- The outer product is a matrix!

Matrix-vector product

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

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

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

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

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

$$\begin{aligned} &= \begin{bmatrix} A_{1,:}\underline{y} \\ A_{2,:}\underline{y} \\ \vdots \\ A_{n,:}\underline{y} \end{bmatrix} \\ &= A_{:,1}\underline{y}_1 + A_{:,2}\underline{y}_2 + \dots + A_{:,m}\underline{y}_m \end{aligned}$$

Matrix multiplication

$$C = \underline{A} \underline{B}.$$

$$C_{i,j} = \sum_k A_{i,k} B_{k,j}.$$

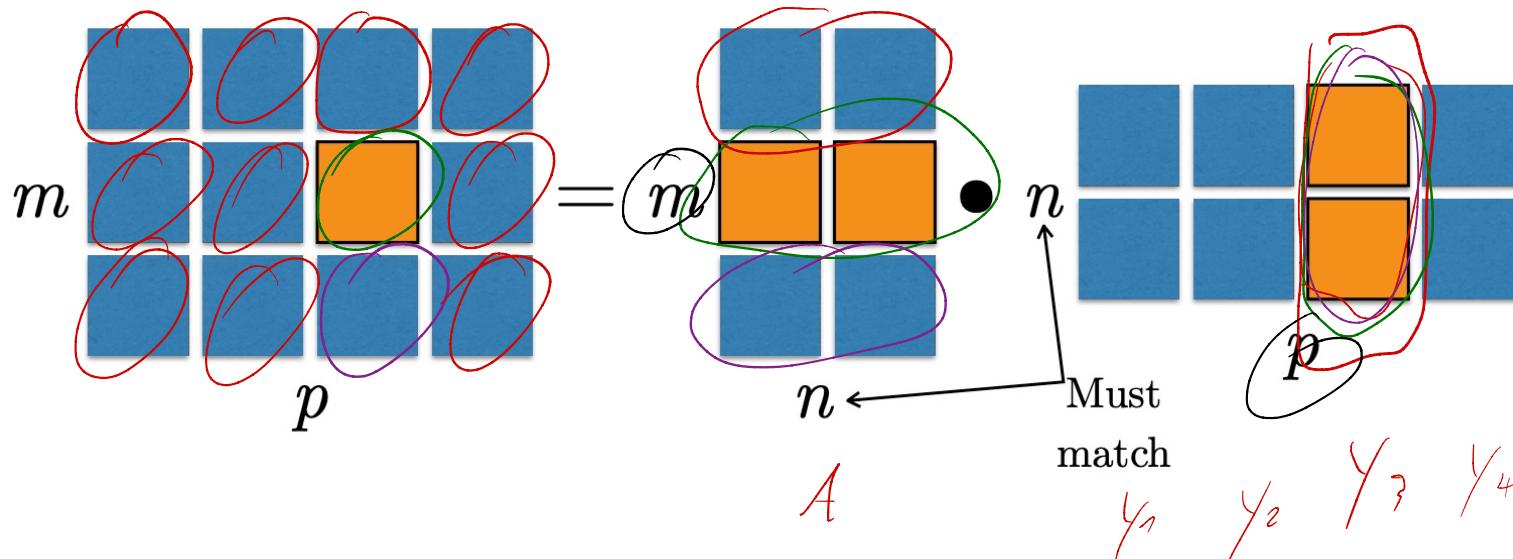
$$A = \begin{pmatrix} 1 & 2 \\ 3 & 4 \\ 5 & 6 \end{pmatrix} \quad \begin{pmatrix} 5 & 11 & \dots & - \\ 11 & 25 & - & \dots \\ \dots & \dots & \dots & \dots \end{pmatrix} = C$$

rows from A

columns from B

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

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



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)

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

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

- We say that vectors $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m \in \mathbb{R}^n$ are linearly **dependent** if any of the vectors can be written as a linear combination of the others, e.g., if, for some scalars $\alpha_1, \alpha_2, \dots, \alpha_m \in \mathbb{R}$,

$$\mathbf{x}_m = \sum_{j=1}^{m-1} \alpha_j \mathbf{x}_j$$

Linear independence

- We say that vectors $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m \in \mathbb{R}^n$ are linearly **dependent** if any of the vectors can be written as a linear combination of the others, e.g., if, for some scalars $\alpha_1, \alpha_2, \dots, \alpha_m \in \mathbb{R}$,

$$\mathbf{x}_m = \sum_{j=1}^{m-1} \alpha_j \mathbf{x}_j$$

- If this is **not** the case, we say that the vectors are linearly **independent**
- Going back to our matrix \mathbf{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 $\text{rank}(\mathbf{A}) = \text{rank}(\mathbf{A}^T)$
- \mathbf{A} is invertible if it is of “full rank” (that is, $\text{rank}(\mathbf{A}) = m = n$)

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:

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

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

...

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

- This system can have

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

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

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

- 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

- 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 $\|y\| = 1$

Norms defined for matrices

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

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

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

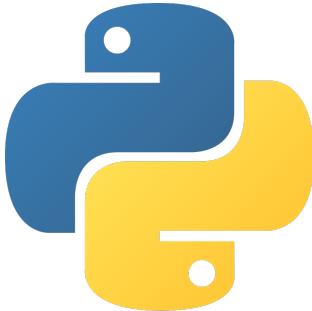
Linear algebra – singular value decomposition

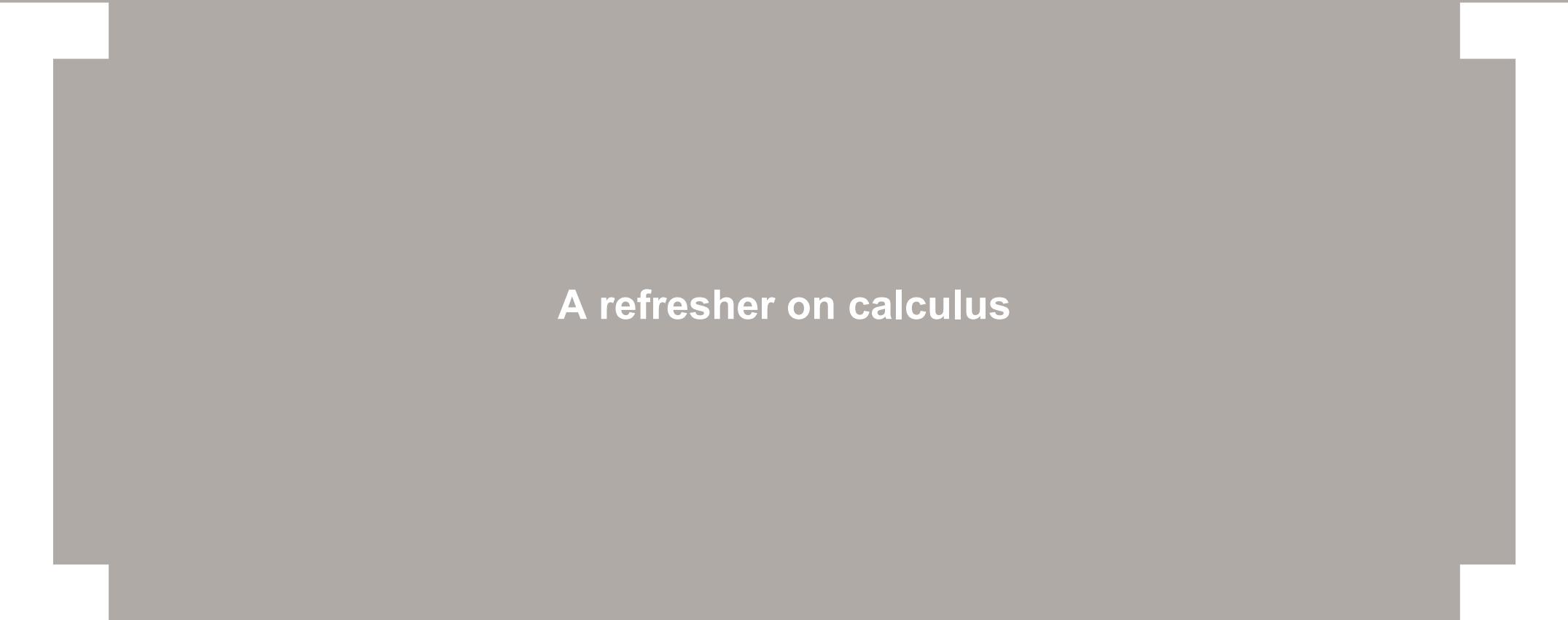
Eigenvectors, eigenvalues, and singular values

- 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 λ)
- Singular values: non-negative square roots of eigenvalues of $A^T A$. Say σ_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

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

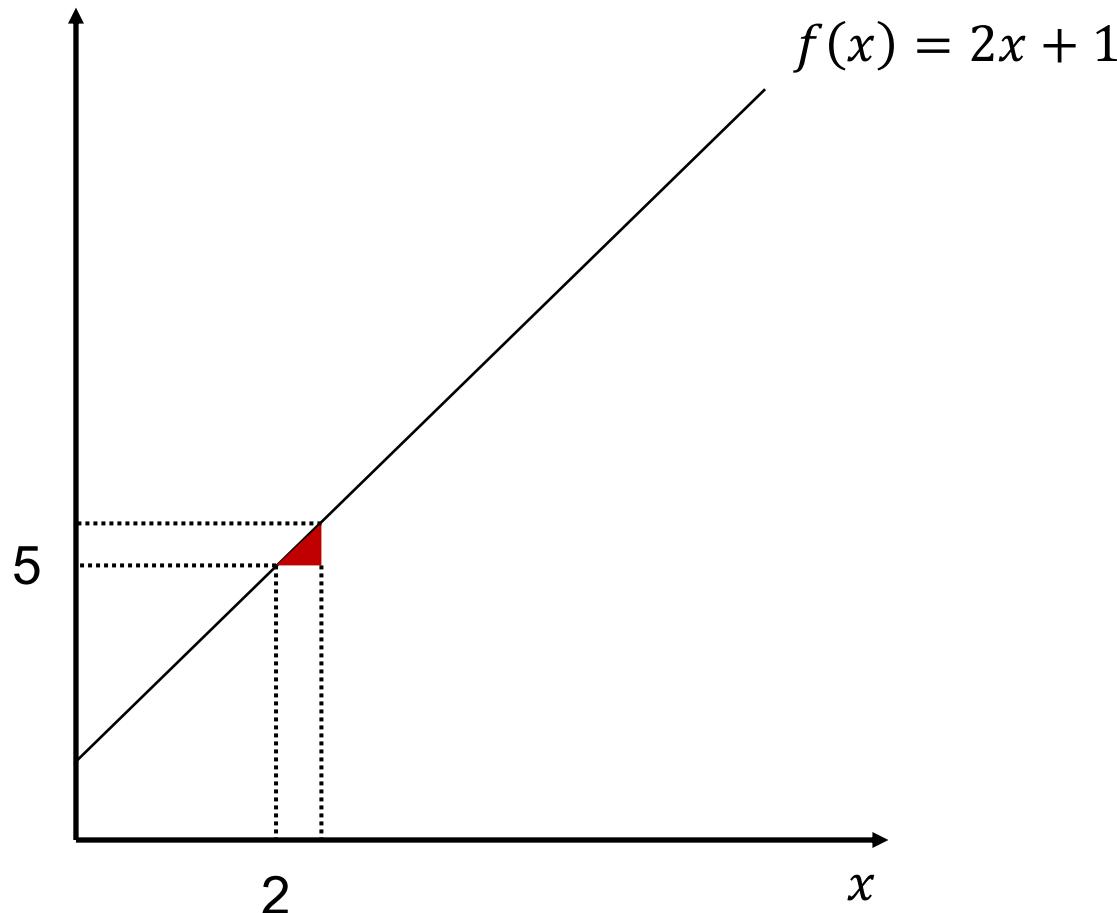
Matrix multiplication and co in Python





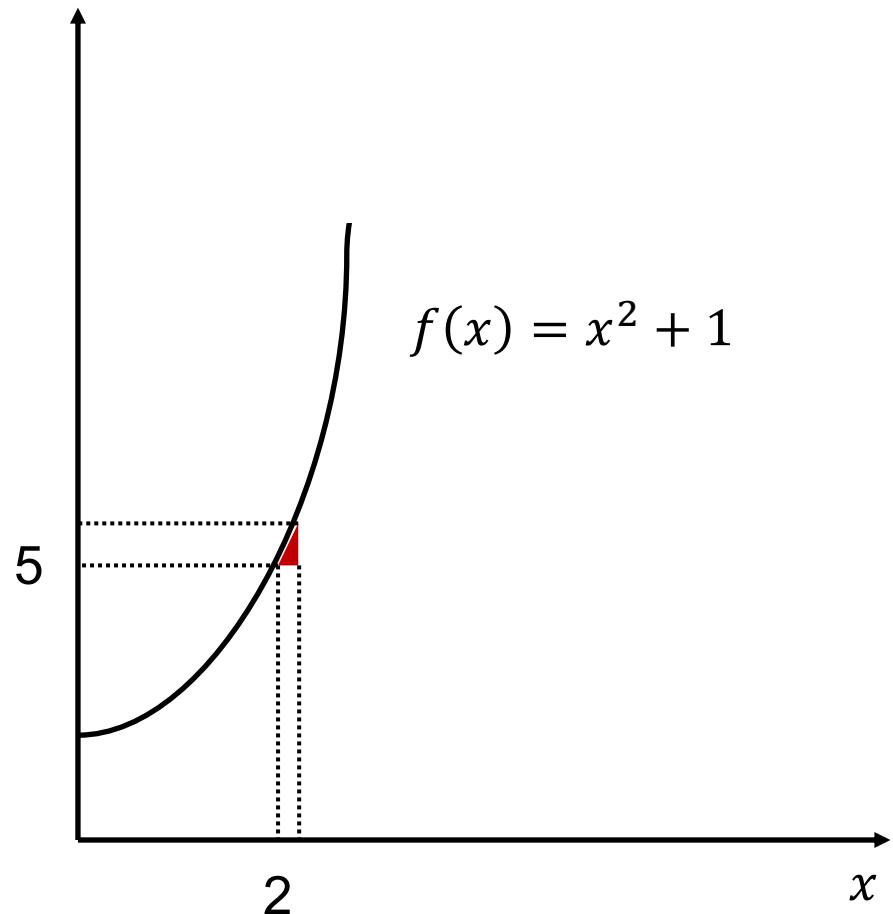
A refresher on calculus

What is a derivative?



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

What is a derivative?



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

A few important derivatives

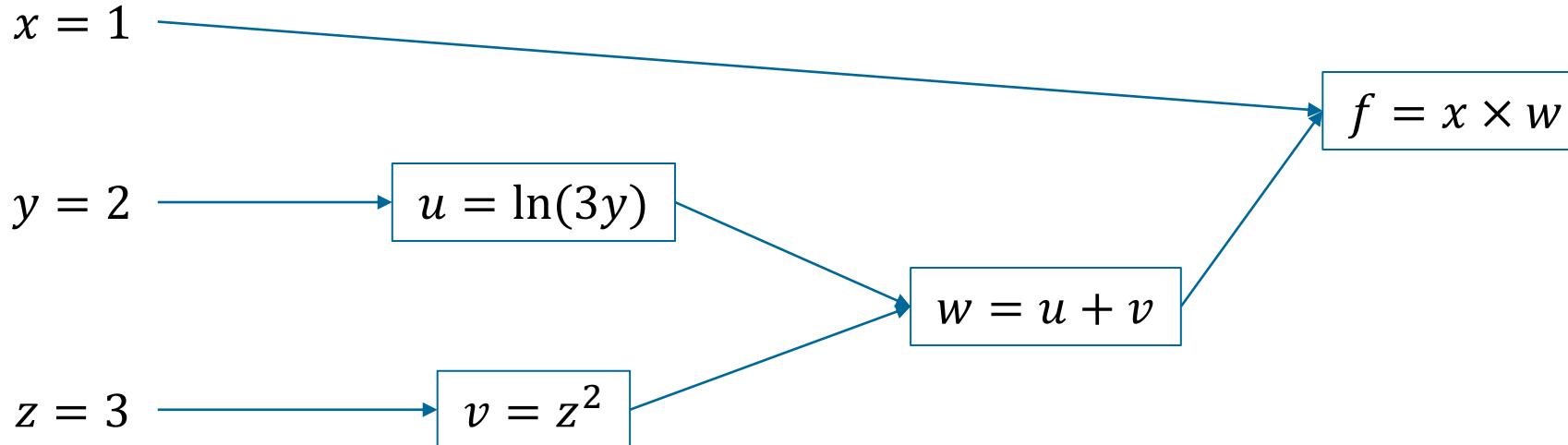
$f(x)$	$f'(x) = \frac{df(x)}{dx} = \frac{d}{dx} f(x)$
1	
x	
x^2	
x^3	
\sqrt{x}	
$\ln(x)$	
e^x	

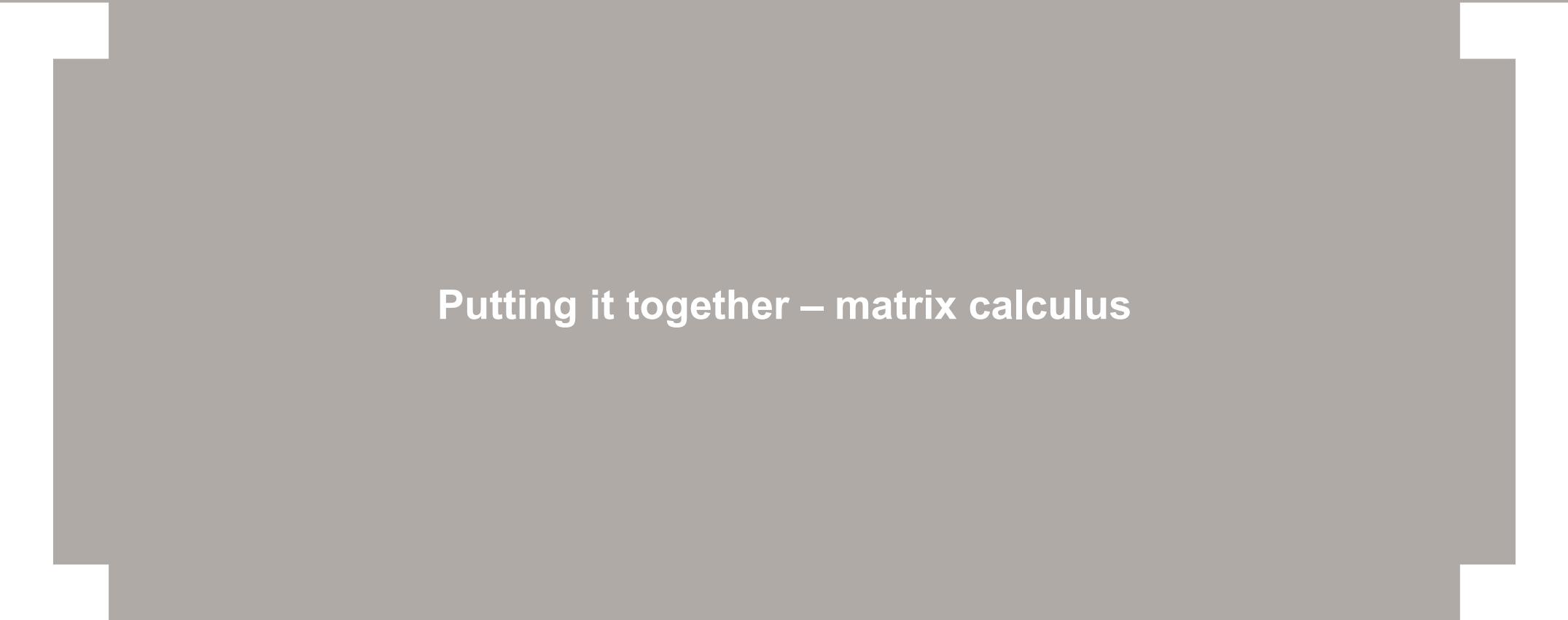
Some rules about handling derivatives

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

Using a computation graph for derivatives

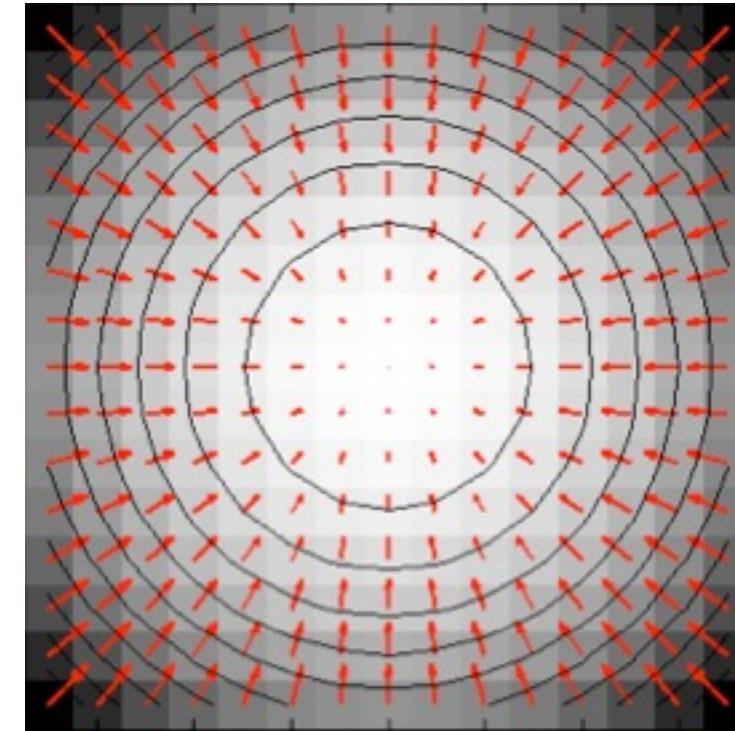
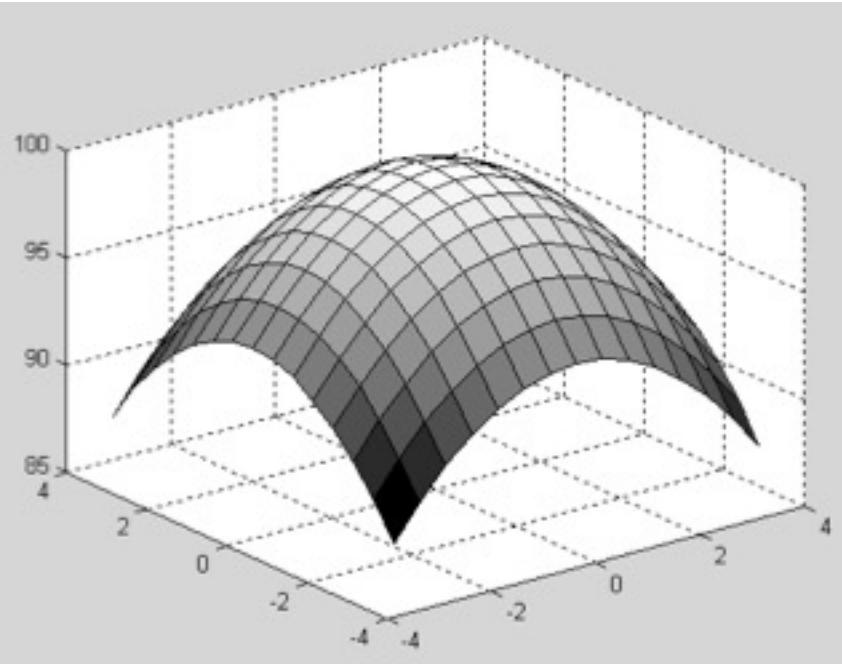
- We can use the “chain rule” to simplify the search for derivatives
- Say, we want to compute the derivatives of $f = x(\ln(3y) + z^2)$ to x, y, z at $x = 1, y = 2, z = 3$





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 $\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$
- The gradient (with respect to \mathbf{x}) is the vector pointing in the direction of fastest increase:

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

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 $\mathbf{A} = \begin{bmatrix} a_{1,1} & a_{1,2} & \dots & a_{1,m} \\ a_{2,1} & a_{2,2} & \dots & a_{2,m} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n,1} & a_{n,2} & \dots & a_{n,m} \end{bmatrix}$
- The gradient (with respect to \mathbf{A}) is now also a matrix $\nabla_{\mathbf{A}} f(\mathbf{A}) = \begin{bmatrix} \frac{\partial f(\mathbf{A})}{\partial a_{1,1}} & \frac{\partial f(\mathbf{A})}{\partial a_{1,2}} & \dots & \frac{\partial f(\mathbf{A})}{\partial a_{1,m}} \\ \frac{\partial f(\mathbf{A})}{\partial a_{2,1}} & \frac{\partial f(\mathbf{A})}{\partial a_{2,2}} & \dots & \frac{\partial f(\mathbf{A})}{\partial a_{2,m}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f(\mathbf{A})}{\partial a_{n,1}} & \frac{\partial f(\mathbf{A})}{\partial a_{n,2}} & \dots & \frac{\partial f(\mathbf{A})}{\partial a_{n,m}} \end{bmatrix}$



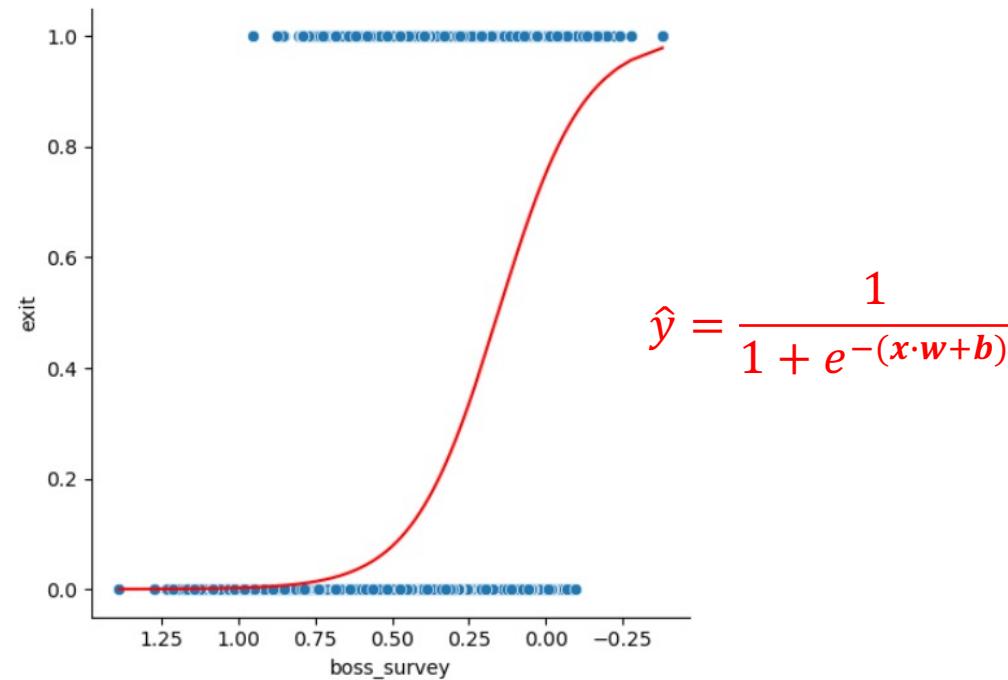
Taking a step back – logistic regression

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

- We are given values $(x^{(i)}, y^{(i)})$, 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:

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:



The optimization part

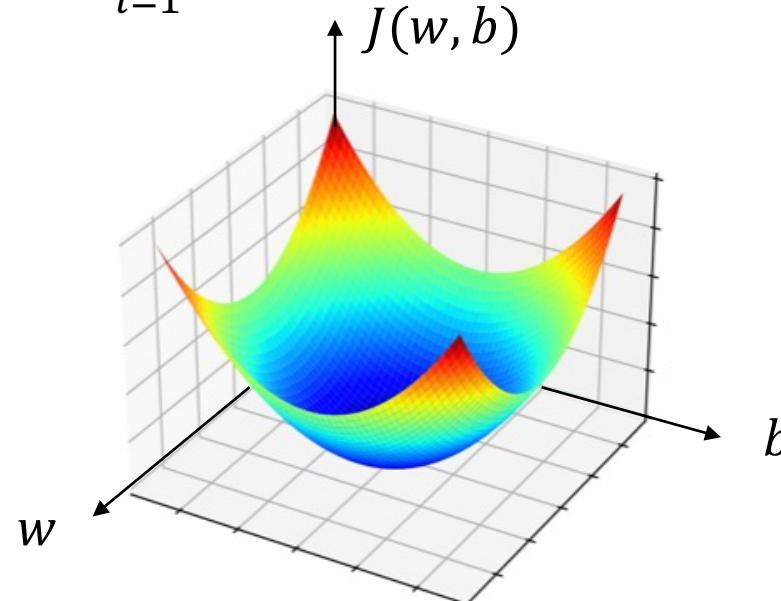
- Remember that $\mathbf{w} \in \mathbb{R}^m$ and $b \in \mathbb{R}$
- To get to the “right” model, we optimize our parameters \mathbf{w}, b so that the $\hat{y}^{(i)}$ s are “as close as possible” to the y^i s
- What we do is to minimize the “cost-function” $J(\mathbf{w}, b)$, where $\hat{y}^{(i)} =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

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

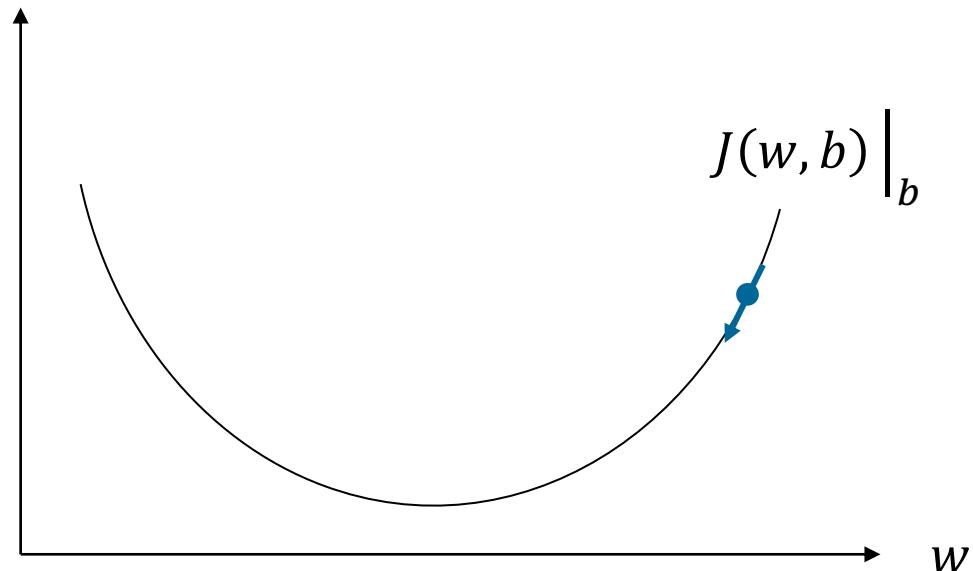
The optimization part

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

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



Solving the optimization problem through gradient descent

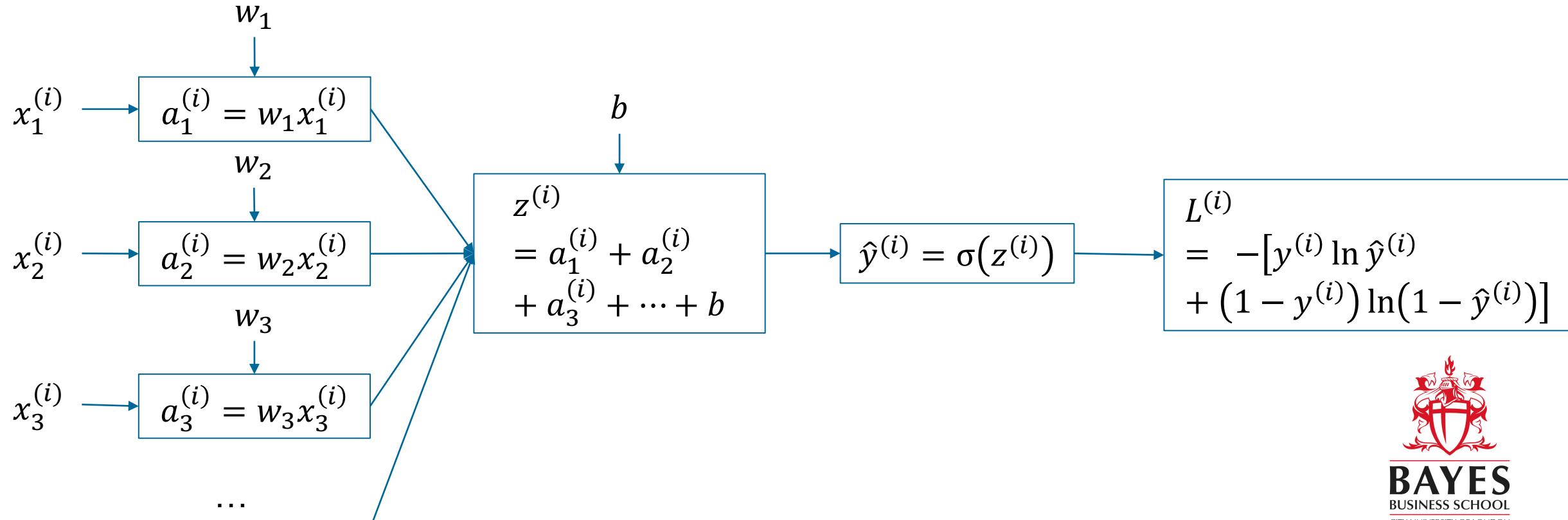


Our first optimization algorithm

1. Decide a “learning rate” α
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 $\hat{y}^{(i)} = \frac{1}{1+e^{-(x^{(i)}w+b)}} = \sigma(x^{(i)}w + b)$

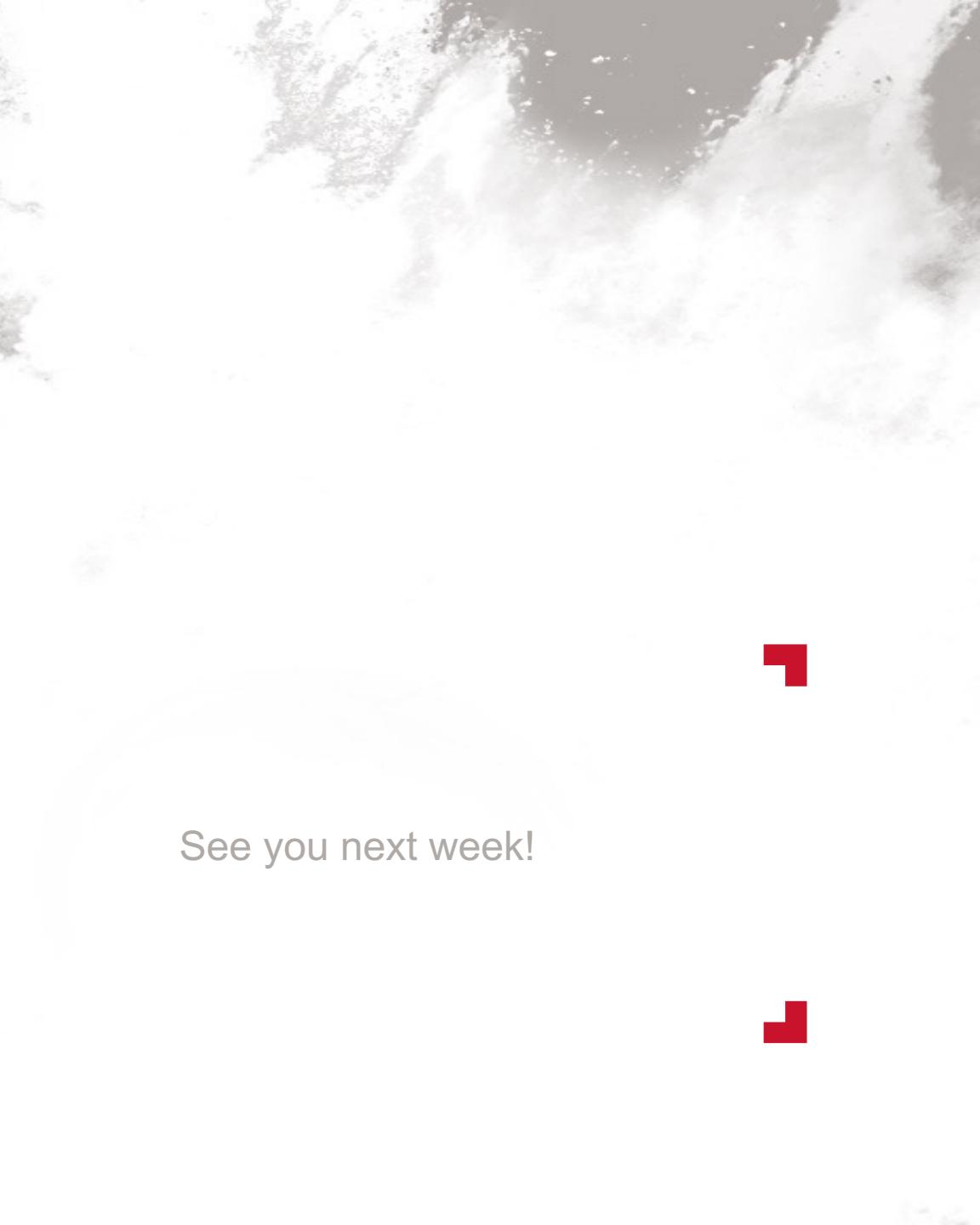


As the same parameters influence all examples, we have to consider one final step

- 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
- 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>