



Applied Deep Learning

Dr. Philippe Blaettchen
Bayes Business School (formerly Cass)

www.bayes.city.ac.uk

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 how to implement concepts in Python
- Building our very own logistic regression algorithm (the simplest neural network)



Why are we even doing this?

- Straight up: this will be the most tedious class for most of you
- However, deep learning, at the very core, functions with fundamental linear algebra operations and gradient descent algorithms (calculus!)
- Hence, we need to learn the gist of these concepts in order to:
 - Understand what is happening “behind the scenes” of neural networks
 - Build intuition about how models can be adjusted and improved – even many of the out-of-the-box tools to tune your networks don’t make sense if you don’t know what a gradient descent algorithm is
 - Create confidence in handling ultra-high dimensional data with millions of observations without having to rely on visual inspection
 - ...





Linear algebra – definitions

Scalars, Vectors, Matrices, and Tensors

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



Scalars, Vectors, Matrices, and Tensors

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

- Vector: One-dimensional array of scalars

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

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

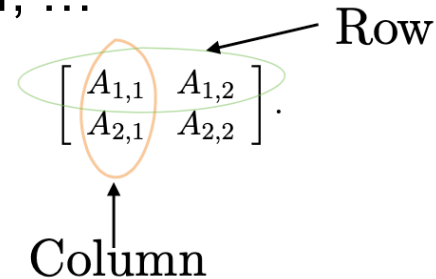
Scalars, Vectors, Matrices, and Tensors

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

- Vector: One-dimensional array of scalars

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

- Entries can be real numbers, binary, integer, ...
- Matrix: Two-dimensional array of numbers



A diagram of a 2x2 matrix $\begin{bmatrix} A_{1,1} & A_{1,2} \\ A_{2,1} & A_{2,2} \end{bmatrix}$. A green oval encircles the top row, with an arrow pointing to it from the label "Row". An orange oval encircles the first column, with an arrow pointing to it from the label "Column".

Scalars, Vectors, Matrices, and Tensors

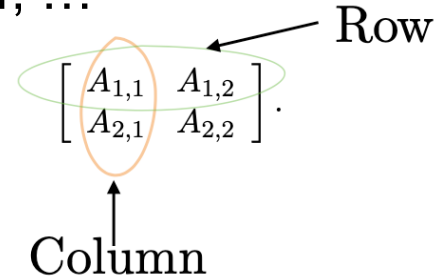
- Scalar: a single number
 - Integers (-1,0,1,2,...), real numbers (0.319375, 1.17, π), rational numbers $\left(\frac{integer}{integer}\right)$

- Vector: One-dimensional array of scalars

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

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

- Matrix: Two-dimensional array of numbers



A diagram of a 2x2 matrix $\begin{bmatrix} A_{1,1} & A_{1,2} \\ A_{2,1} & A_{2,2} \end{bmatrix}$. A green oval encircles the top row, with an arrow pointing to it from the label "Row". An orange oval encircles the first column, with an arrow pointing to it from the label "Column".

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

A typical use of matrices in deep learning

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

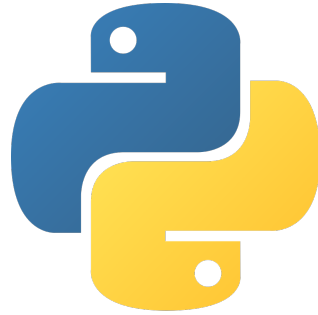


Pixel's intensity value

Source: Ivanovic

In Python

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






Linear algebra – typical operations

Matrix transpose

- Essentially a mirror image across the main diagonal


$$\mathbf{A} = \begin{bmatrix} A_{1,1} & A_{1,2} \\ A_{2,1} & A_{2,2} \\ A_{3,1} & A_{3,2} \end{bmatrix} \Rightarrow \mathbf{A}^{\top} = \begin{bmatrix} A_{1,1} & A_{2,1} & A_{3,1} \\ A_{1,2} & A_{2,2} & A_{3,2} \end{bmatrix}$$

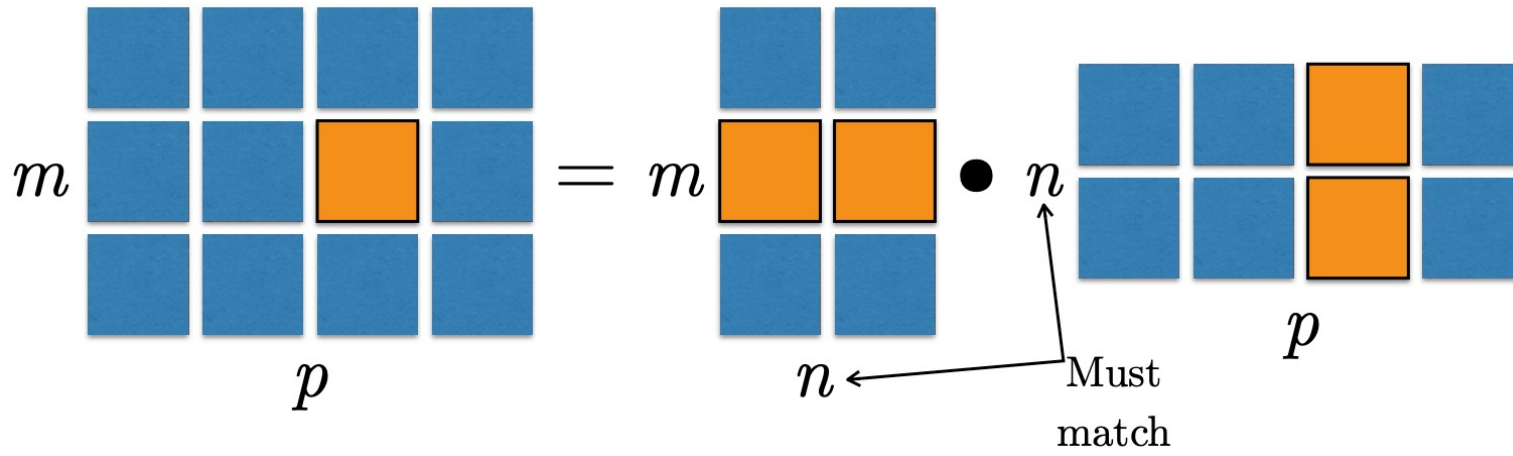
$$(\mathbf{A}^{\top})_{i,j} = A_{j,i}.$$

- We call a matrix symmetric if $\mathbf{A} = \mathbf{A}^{\top}$

Matrix multiplication

$$C = AB.$$

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



Source: Goodfellow

In Python



BAYES
BUSINESS SCHOOL
CITY, UNIVERSITY OF LONDON



Linear algebra – norms

- Functions f that measures the “length” of a vector x
- Such functions need to fulfill four conditions:
 - f needs to return non-negative values only (there is no negative length!)
 - The only vector that has length zero should be the 0-vector
 - The length “scales”: For all $x \in \mathbb{R}^n, \alpha \in \mathbb{R}, f(\alpha x) = |\alpha|f(x)$
 - The triangle inequality needs to hold: For all $x, y \in \mathbb{R}^n, f(x + y) \leq f(x) + f(y)$



Some commonly used norms

- L^p norm:
 - $\|\mathbf{x}\|_p = (\sum_i |x_i|^p)^{\frac{1}{p}}$
 - Most commonly used norm: $\|\mathbf{x}\|_2 = \sqrt{\sum_i x_i^2}$ (L^2 norm or “Euclidian” norm)
 - Quite common as well: $\|\mathbf{x}\|_1 = \sum_i |x_i|$ (L^1 norm)
- An extreme case: the max-norm $\|\mathbf{x}\|_\infty = \max_i |x_i|$



Norms defined for matrices

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

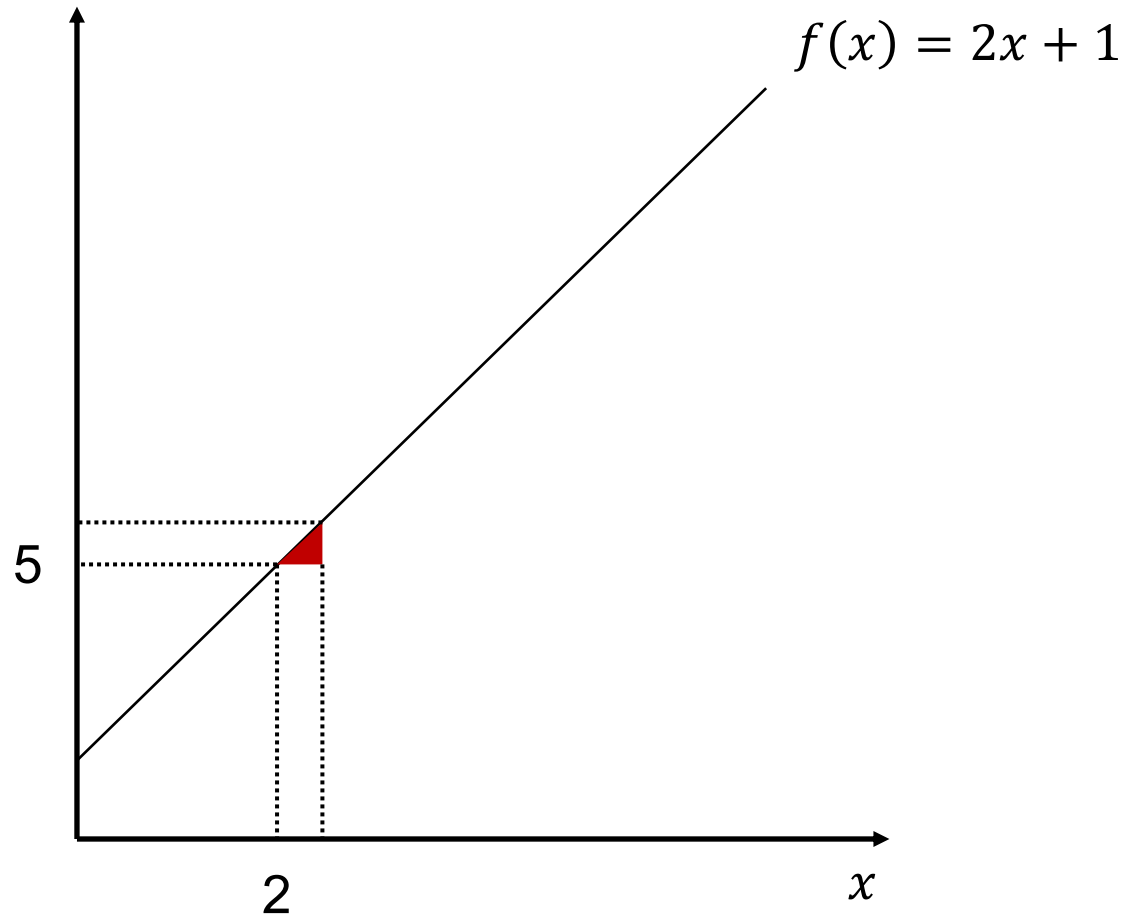
$$\|A\|_F = \sqrt{\sum_{i=1}^n \sum_{j=1}^m A_{i,j}^2}$$

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



A refresher on calculus

What is a derivative?

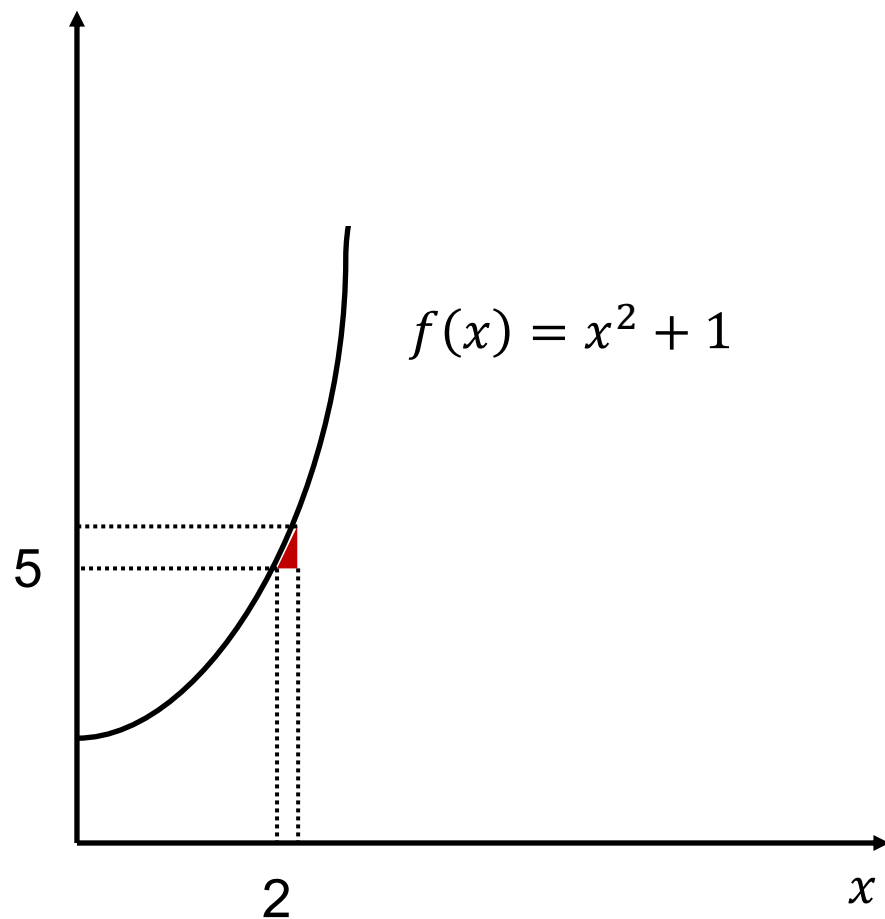


Slope (derivative) of $f(x)$ at 2 is $\frac{\text{"height"}}{\text{"width"}}$



BAYES
BUSINESS SCHOOL
CITY, UNIVERSITY OF LONDON

What is a derivative?



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

A few important derivatives

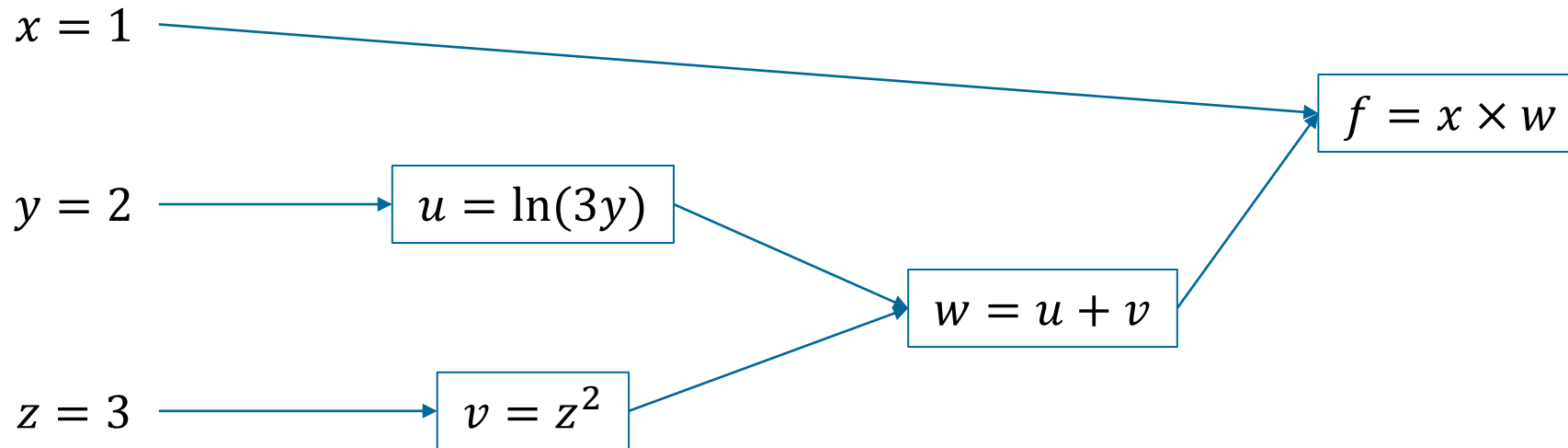
$f(x)$	$f'(x) = \frac{df(x)}{dx} = \frac{d}{dx}f(x)$
1	
x	
x^2	
x^3	
\sqrt{x}	
$\ln(x)$	
e^x	

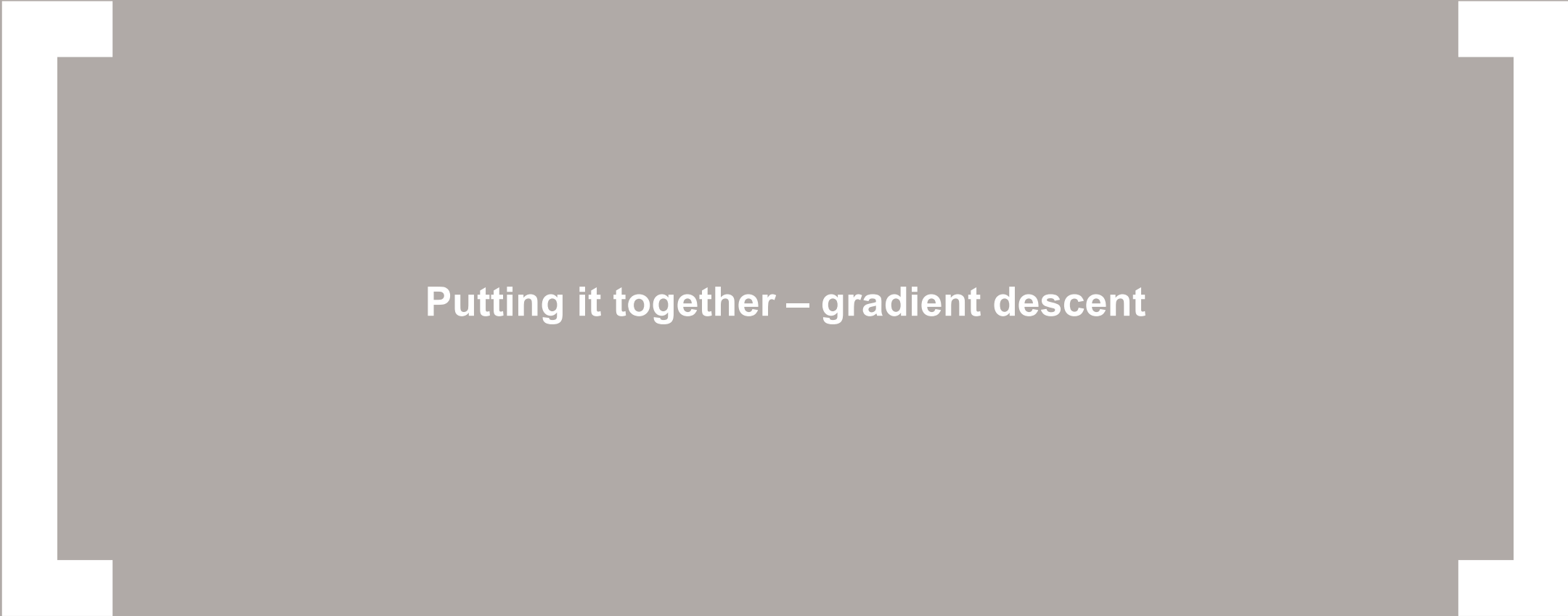
Some rules about handling derivatives

$h(x)$	$h'(x)$	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

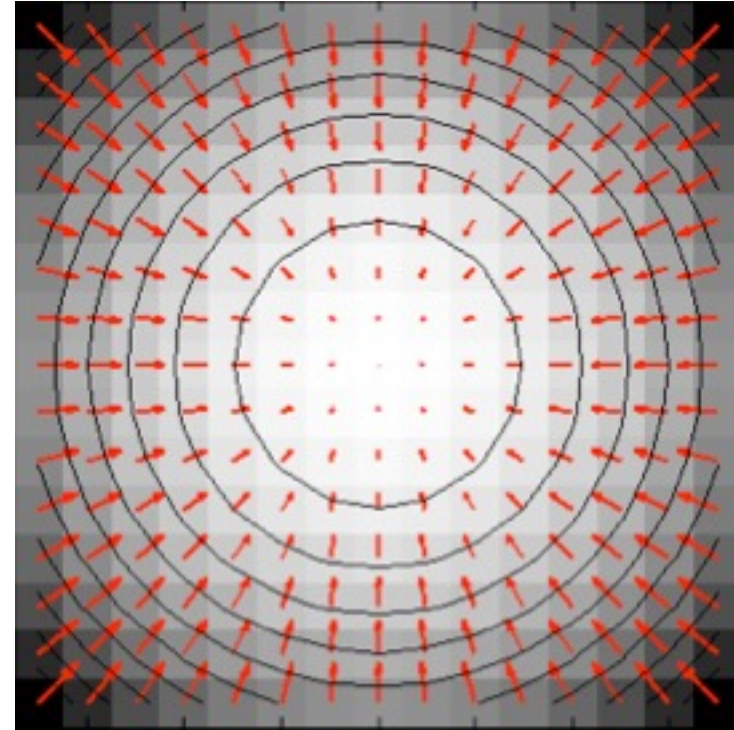
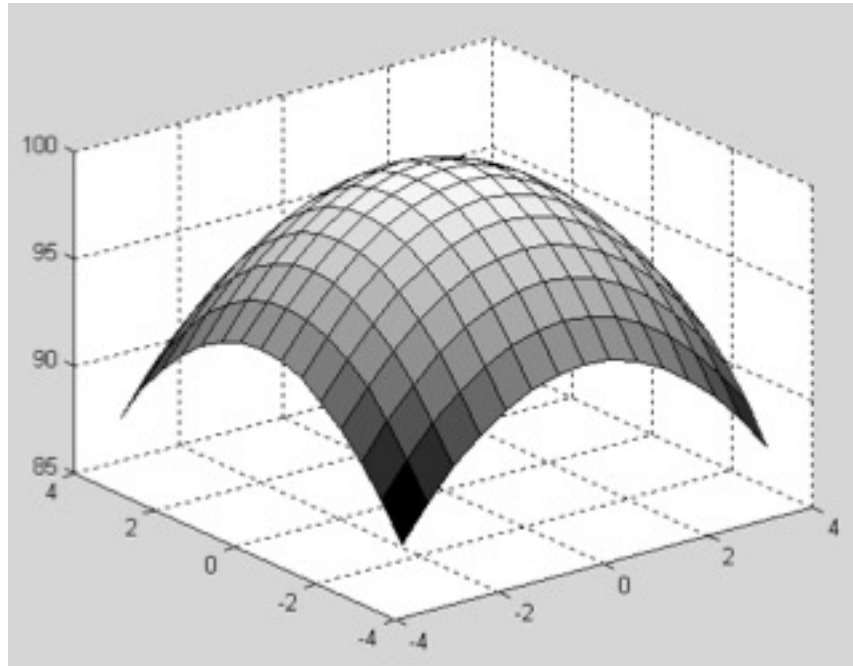
- We can use the “chain rule” to simplify the search for derivatives
- Say, we want to compute the derivatives of $f = x(\ln(3y) + z^2)$ to x, y, z at $x = 1, y = 2, z = 3$





Putting it together – gradient descent

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:

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



This naturally extends to functions that take as input a matrix

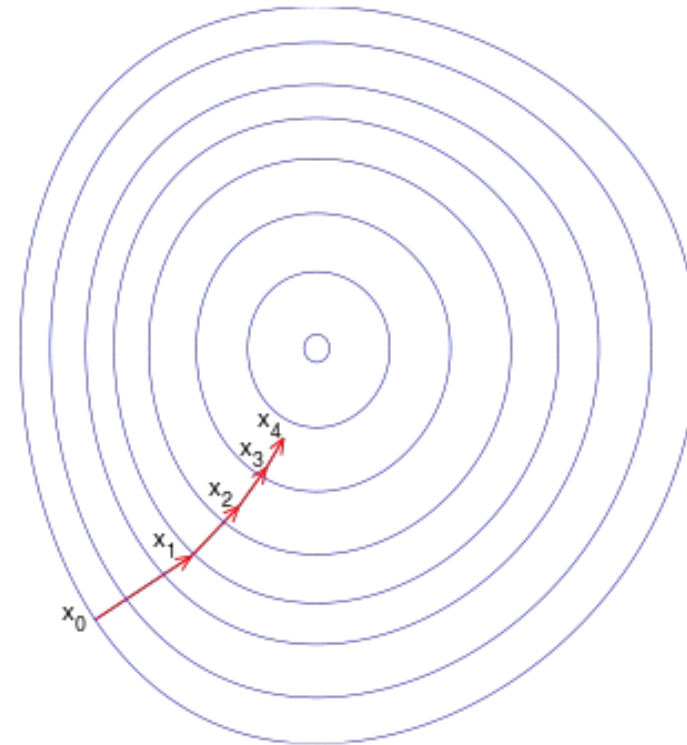
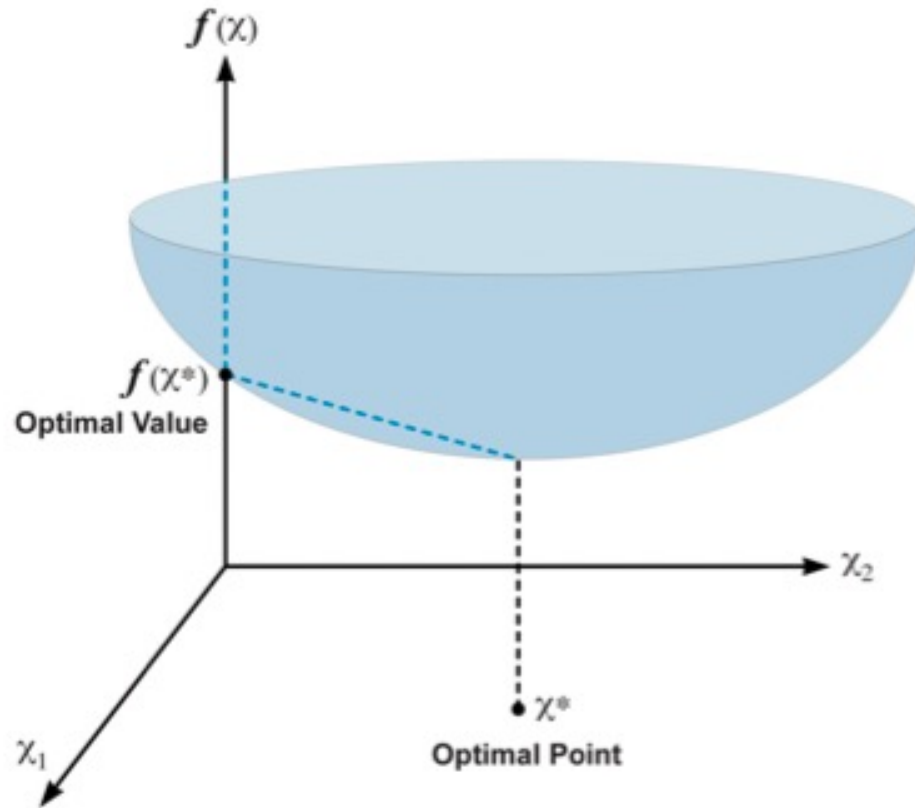
- Say, now we have a function f , taking as input a matrix $A = \begin{bmatrix} a_{1,1} & a_{1,2} & \dots & a_{1,m} \\ a_{2,1} & a_{2,2} & \dots & a_{2,m} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n,1} & a_{n,2} & \dots & a_{n,m} \end{bmatrix}$
- The gradient (with respect to A) is now also a matrix $\nabla_A f(A) = \begin{bmatrix} \frac{\partial f(A)}{\partial a_{1,1}} & \frac{\partial f(A)}{\partial a_{1,2}} & \dots & \frac{\partial f(A)}{\partial a_{1,m}} \\ \frac{\partial f(A)}{\partial a_{2,1}} & \frac{\partial f(A)}{\partial a_{2,2}} & \dots & \frac{\partial f(A)}{\partial a_{2,m}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f(A)}{\partial a_{n,1}} & \frac{\partial f(A)}{\partial a_{n,2}} & \dots & \frac{\partial f(A)}{\partial a_{n,m}} \end{bmatrix}$



Why do we need the gradient? A thought experiment:

- You wake up somewhere in the mountain
- Your goal is to reach the lowest point as quickly as possible
- You have no map, GPS, and can only see a few meters ahead because of trees and fog
- How do you do it?

Gradient descent – the idea



The algorithm

- Take small steps
- For each step, go downhill in the locally steepest descent direction
- Repeat until you are on a flat surface



Taking a step back – logistic regression

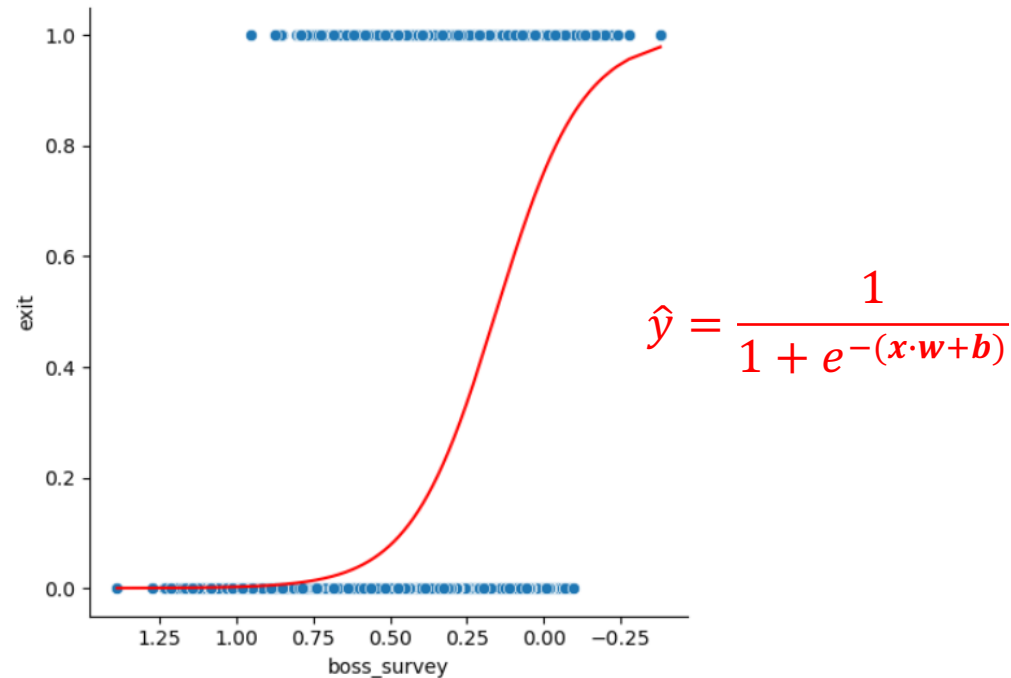
What do we actually do when training a logistic regression model?

- 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:



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:



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

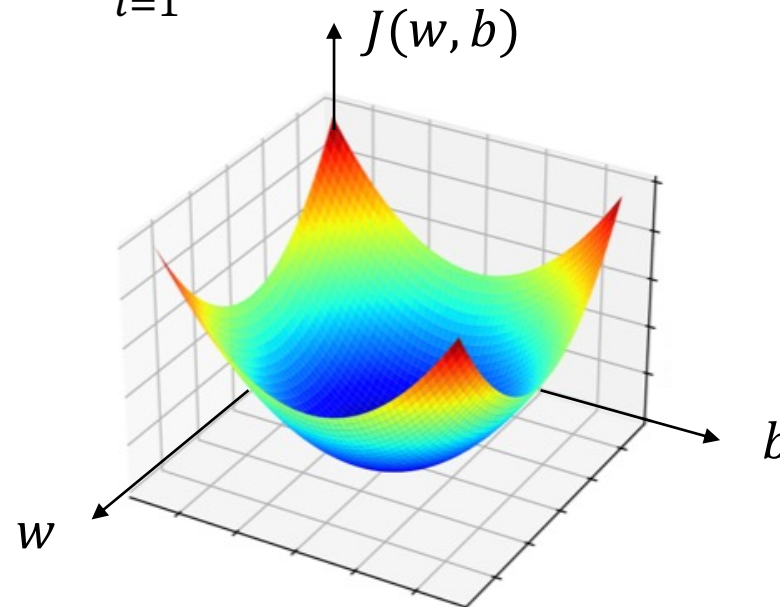
$$J(\mathbf{w}, b) = -\frac{1}{n} \sum_{i=1}^n [y^{(i)} \ln \hat{y}^{(i)} + (1 - y^{(i)}) \ln(1 - \hat{y}^{(i)})]$$



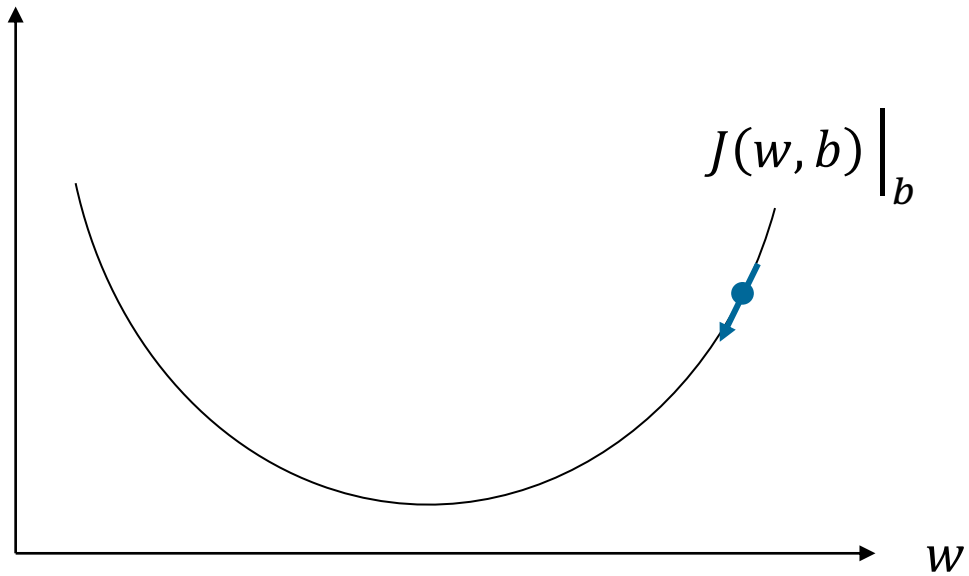
The optimization part

- Remember that $\mathbf{w} \in \mathbb{R}^m$ 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 $\hat{y}^{(i)} = \frac{1}{1+e^{-(x^{(i)}\mathbf{w}+b)}}$:

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



Solving the optimization problem through gradient descent



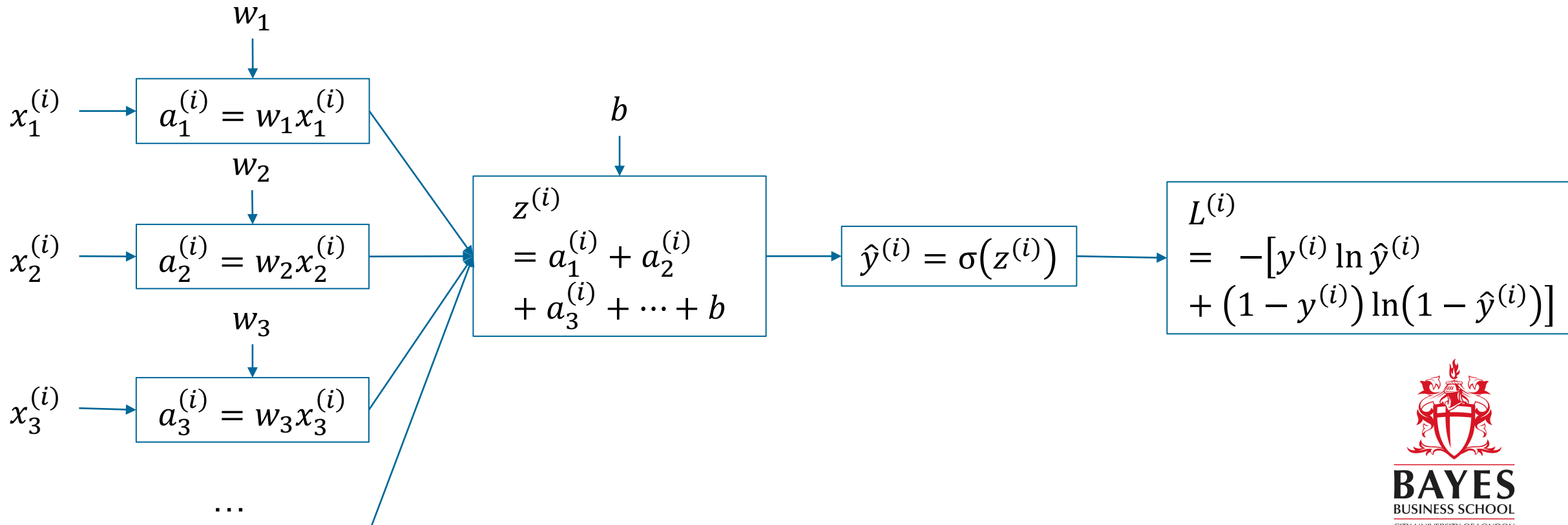
Our first optimization algorithm

1. Decide a “learning rate” α
2. Start with some \mathbf{w} and b and compute $J(\mathbf{w}, b)$
3. Until J “doesn’t change” anymore:
 - Let $w_1 := w_1 - \alpha \frac{\partial J(\mathbf{w}, b)}{\partial w_1}$
 - Let $w_2 := w_2 - \alpha \frac{\partial J(\mathbf{w}, b)}{\partial w_2}$
 - ...
 - Let $w_m := w_m - \alpha \frac{\partial J(\mathbf{w}, b)}{\partial w_m}$
 - Let $b := b - \alpha \frac{\partial J(\mathbf{w}, b)}{\partial b}$
 - Recompute $J(\mathbf{w}, b)$
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 $\hat{y}^{(i)} = \frac{1}{1+e^{-(x^{(i)}w+b)}} = \sigma(x^{(i)}w+b)$

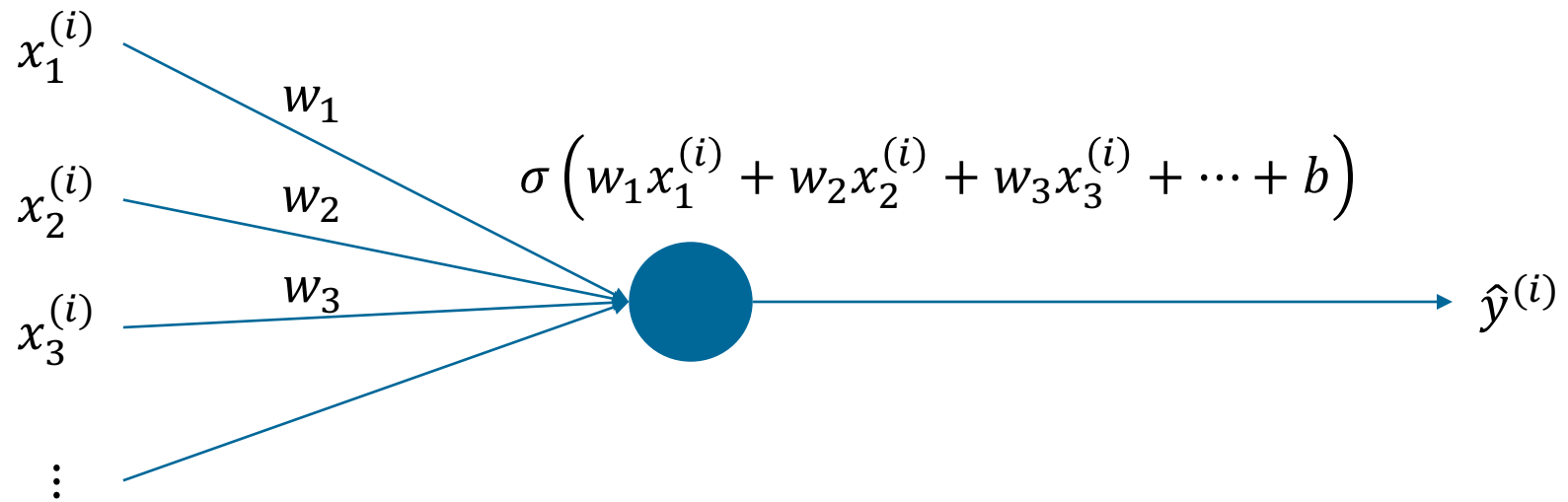


As the same parameters influence all examples, we have to consider one final step

- 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 $\frac{\partial J(\mathbf{w}, b)}{\partial w_j} = \frac{1}{n} \sum_{i=1}^n \frac{\partial L^{(i)}}{\partial w_j}$



Schema of a logistic regression



We can now implement a logistic regression



BAYES
BUSINESS SCHOOL
CITY, UNIVERSITY OF LONDON



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>

