**Applied Deep Learning**

Dr. Philippe Blaettchen
Bayes Business School (formerly Cass)

**Goals:** Creating neural networks – from logistic regression to feed-forward networks
- Use what we have learned about linear algebra and calculus to create a logistic regression algorithm from scratch
- Understand how what we learned generalizes to neural networks in general

**How will we do this?**
- We start with a quick recap on our discussion about logistic regression so far
- We then implement a logistic regression algorithm with numpy only
- Next, we visualize more general neural networks with the TensorFlow playground, before defining the concepts relevant for running our own networks

BAYES
BUSINESS SCHOOL
CITY, UNIVERSITY OF LONDON

# Taking a step back – logistic regression

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

- We model this probability, using the sigmoid function:

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

$$a_0 + a_1 x_1^{(i)} + a_2 x_2^{(i)} \cdots + a_m x_m^{(i)}$$

$$b + w_1 x_1 + w_2 x_2 + \cdots + w_m x_m$$

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

- We model this probability, using the sigmoid function:

$$\boxed{\boldsymbol{x} \quad \boldsymbol{w} \quad + \boldsymbol{b}}$$

$$(1, m) \quad (m, 1) \qquad (1, 1)$$



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

# 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 $\hat{y}^{(i)} = \dfrac{1}{1+e^{-\left(x^{(i)}w+b\right)}}$ :

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

$L^{(i)}$

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

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

BAYES
BUSINESS SCHOOL
CITY, UNIVERSITY OF LONDON

- 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
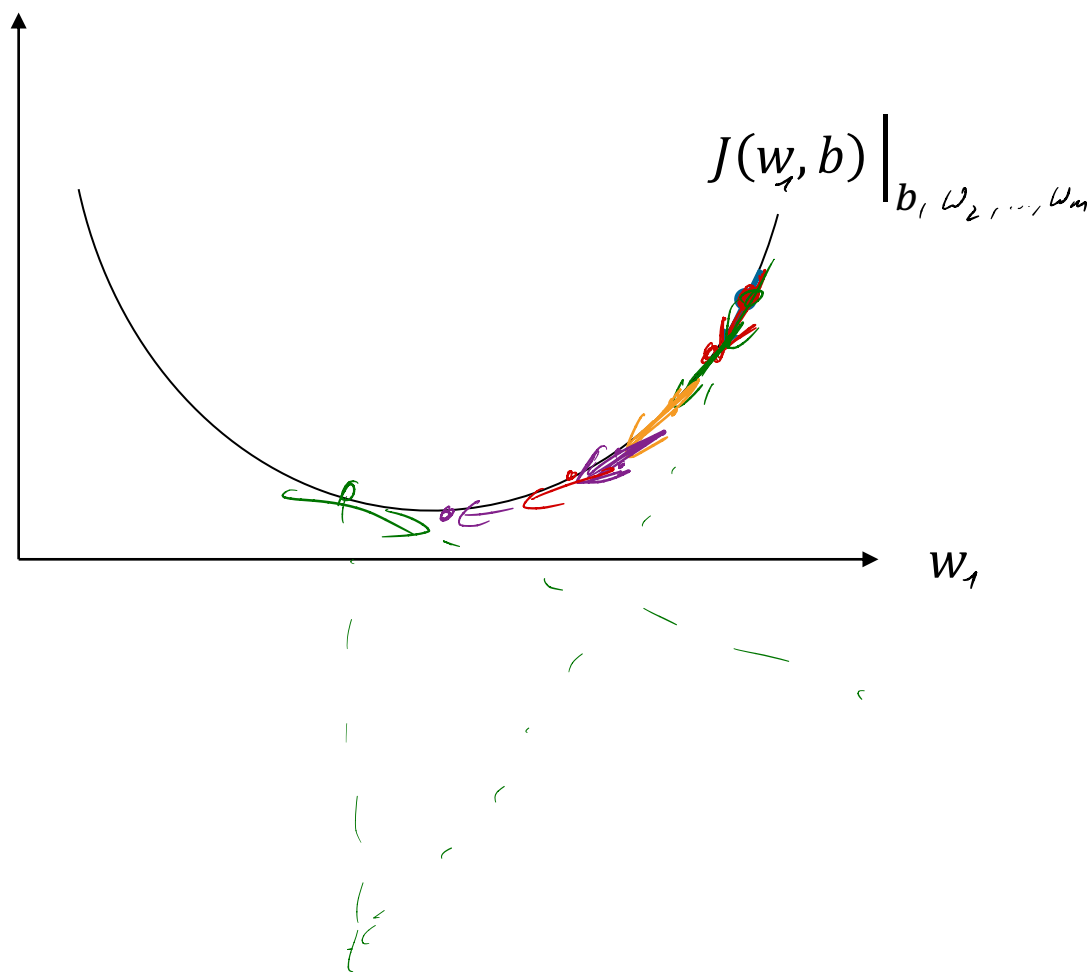- What we do is to minimize the "cost-function" $J(w, b)$, where $\hat{y}^{(i)} = \dfrac{1}{1+e^{-\left(x^{(i)}w+b\right)}}$ :

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

$$w_1 \quad \leadsto \quad \frac{\partial J}{\partial w_1}$$

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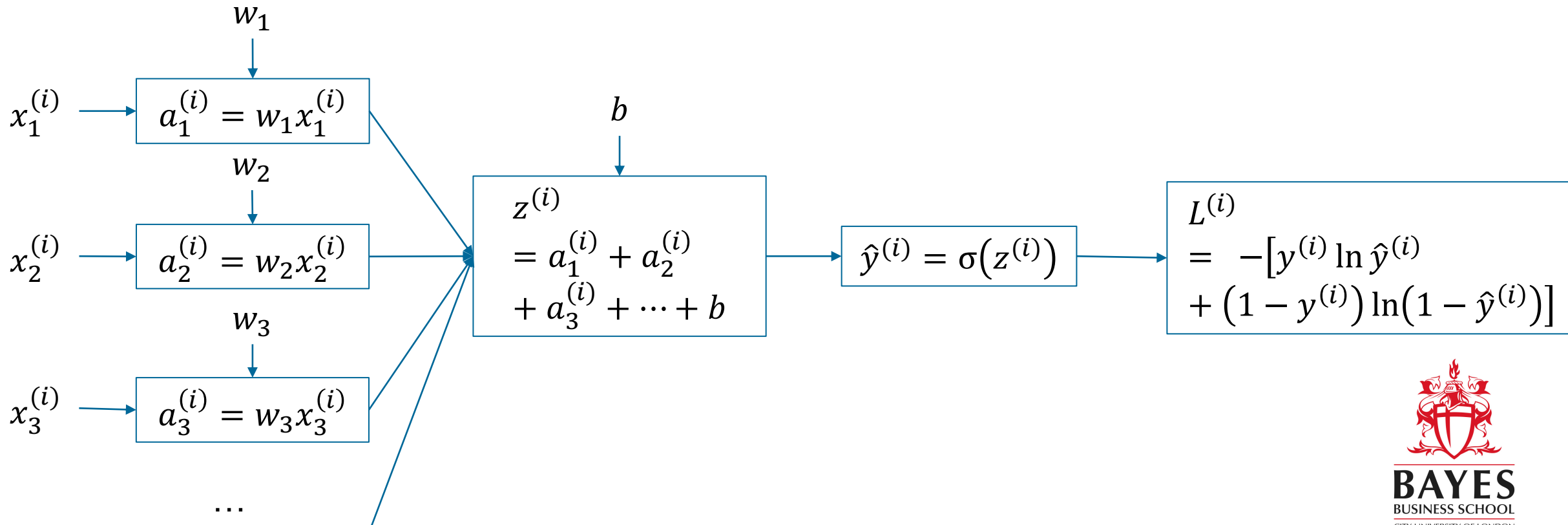
$$w_1 := w_1 - \alpha \frac{\partial J}{\partial w_1}$$

$$J(w_1, b)\Big|_{b, w_2, \ldots, w_m}$$

$$w_1$$

1. Decide a "learning rate" $\alpha$
2. Start with some $\boldsymbol{w}$ and $b$ and compute $J(\boldsymbol{w}, b)$
3. Until $J$ "doesn't change" anymore:

   - Let $w_1 := w_1 - \alpha \dfrac{\partial J(\boldsymbol{w}, b)}{\partial w_1}$

   - Let $w_2 := w_2 - \alpha \dfrac{\partial J(\boldsymbol{w}, b)}{\partial w_2}$

   - …

   - Let $w_m := w_m - \alpha \dfrac{\partial J(\boldsymbol{w}, b)}{\partial w_m}$

   - Let $b := b - \alpha \dfrac{\partial J(\boldsymbol{w}, b)}{\partial b}$

   - Recompute $J(\boldsymbol{w}, b)$

4. Enjoy the fruits of your labor: you have fit a logistic regression model manually!

**BAYES**
BUSINESS SCHOOL
CITY, UNIVERSITY OF LONDON

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

$w_1$

$x_1^{(i)} \longrightarrow \boxed{a_1^{(i)} = w_1 x_1^{(i)}}$

$w_2$

$b$

$x_2^{(i)} \longrightarrow \boxed{a_2^{(i)} = w_2 x_2^{(i)}}$

$\boxed{\begin{array}{l} z^{(i)} \\ = a_1^{(i)} + a_2^{(i)} \\ + a_3^{(i)} + \cdots + b \end{array}} \longrightarrow \boxed{\hat{y}^{(i)} = \sigma\left(z^{(i)}\right)} \longrightarrow \boxed{\begin{array}{l} L^{(i)} \\ = -\big[y^{(i)} \ln \hat{y}^{(i)} \\ + \left(1 - y^{(i)}\right) \ln\left(1 - \hat{y}^{(i)}\right)\big] \end{array}}$

$w_3$

$x_3^{(i)} \longrightarrow \boxed{a_3^{(i)} = w_3 x_3^{(i)}}$

$\ldots$

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

- We have that $\frac{\partial J(\boldsymbol{w}, b)}{\partial w_j} = \frac{1}{n}\sum_{i=1}^{n}\frac{\partial L^{(i)}}{\partial w_j}$
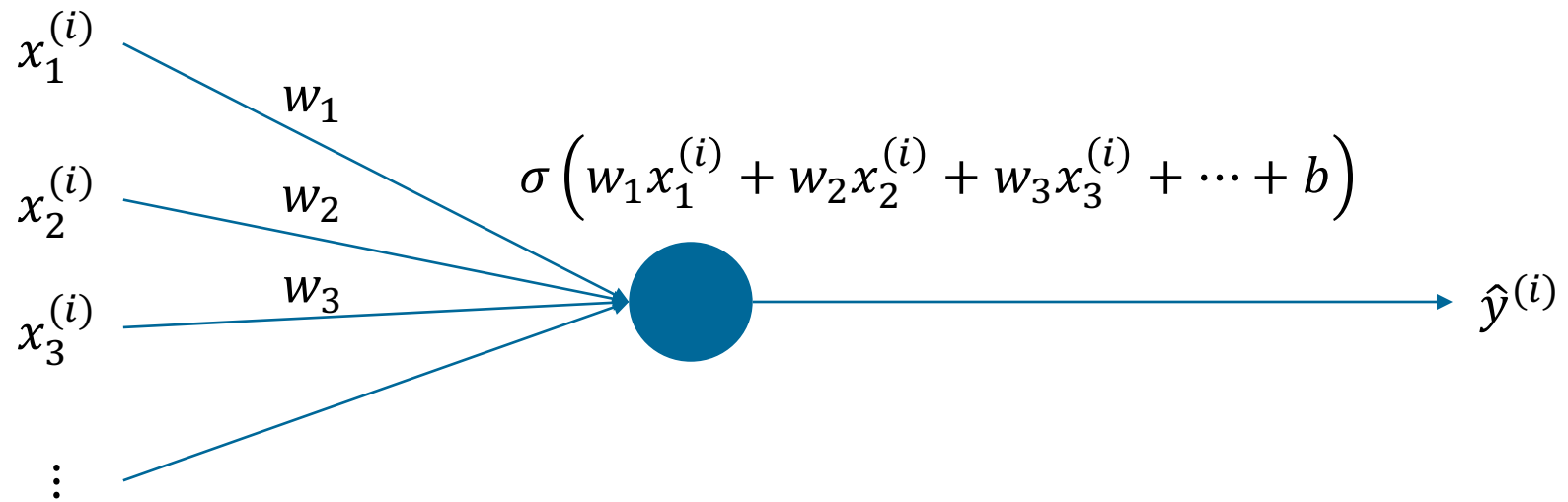
**Visualizing a logistic regression – our first neural network**

# Schema of a logistic regression

$$x_1^{(i)}$$

$$w_1$$

$$x_2^{(i)}$$

$$w_2$$

$$w_3$$

$$x_3^{(i)}$$

$$\sigma\left(w_1 x_1^{(i)} + w_2 x_2^{(i)} + w_3 x_3^{(i)} + \cdots + b\right)$$

$$\hat{y}^{(i)}$$

$$x_1^{(i)}$$

$$x_2^{(i)}$$

$$x_3^{(i)}$$

$$\vdots$$

$$\hat{y}^{(i)}$$

Source: Czarnecki

BAYES
BUSINESS SCHOOL
CITY, UNIVERSITY OF LONDON

Open https://playground.tensorflow.org/

1.  A simple case of binary classification:
    *   Change to the pattern on the lower left
    *   Set the level of "Noise" to 50
    *   Set "Ratio of training to test data" to 50%
    *   Set up the neural network: 1 hidden layer, 1 neuron, then press play
    *   Answer the following questions:
        *   Did the training eventually find a model that seems to capture the pattern in the data?
        *   How would you describe the pattern the model captured?
        *   Record the "Training loss" and "Test loss"
    *   How do your answers change when you select the pattern at the top right? What about setting the noise to 0?

2. A shallow neural network:
- Stick with the pattern at the top right, a noise of 0 and a ratio of 50%
- Now use 3 neurons for your hidden layer
- Answer the three questions from before:
  - Did the training eventually find a model that seems to capture the pattern in the data?
  - How would you describe the pattern the model captured?
  - Record the "Training loss" and "Test loss"
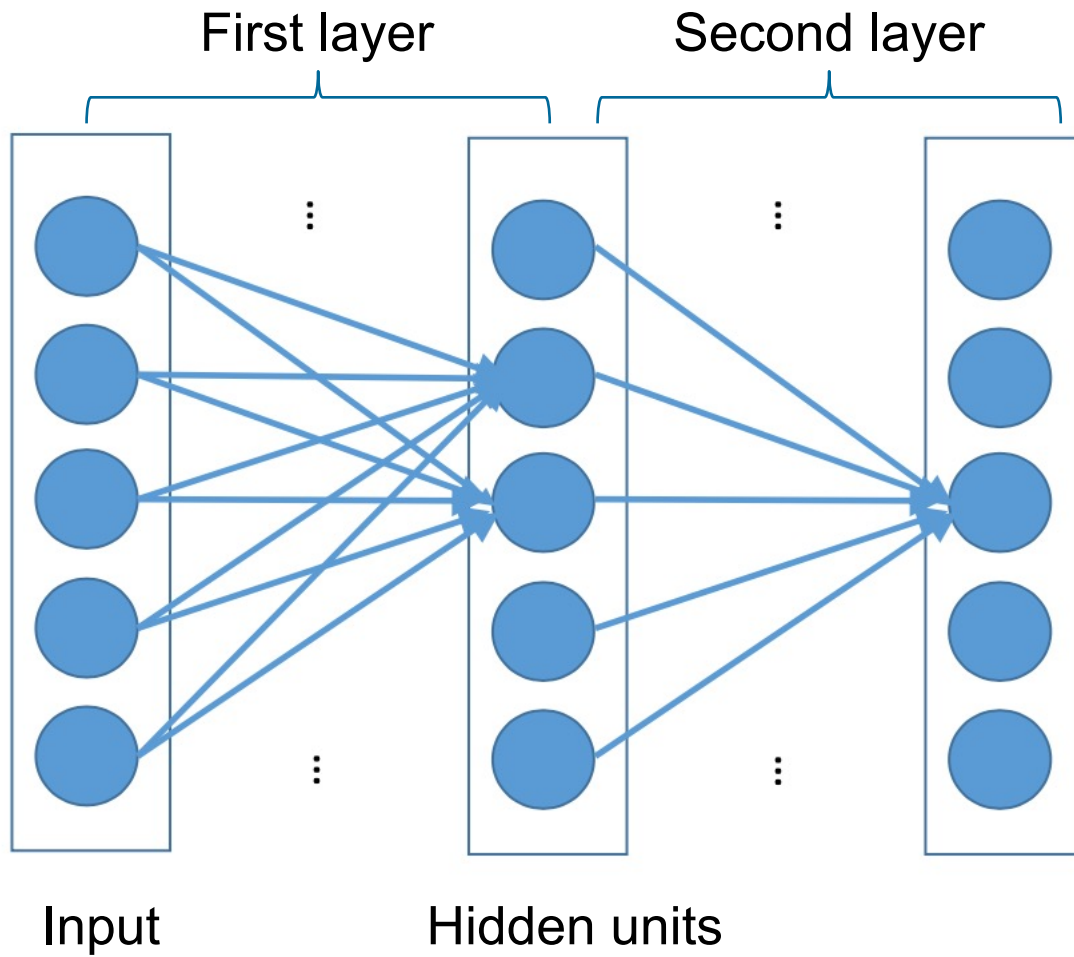- How do your answers change when you use 6 neurons instead?

3. A deep neural network:
- Use a second hidden layer, with 3 neurons each (and the other setups from 2.)
- How do your answers change now?

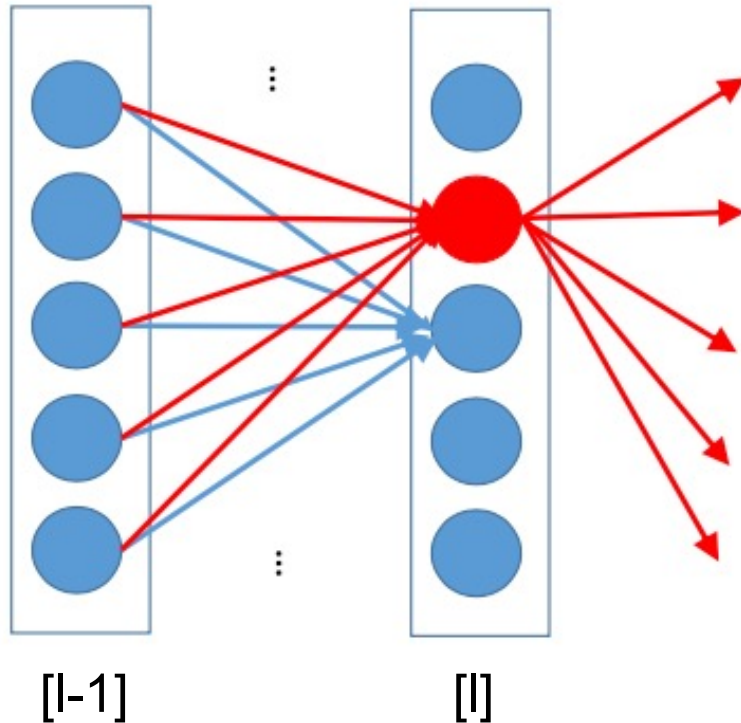Key components of a neural network

**Components**

First layer · Second layer · Output layer (L)

Input · Hidden units

y

Source: Liang

$$\text{``}\boldsymbol{x}\text{''} = \boldsymbol{a}^{[l-1]} \qquad z = \boldsymbol{a}^{[\boldsymbol{l-1}]}\boldsymbol{w} + \boldsymbol{b} \qquad f(z)$$
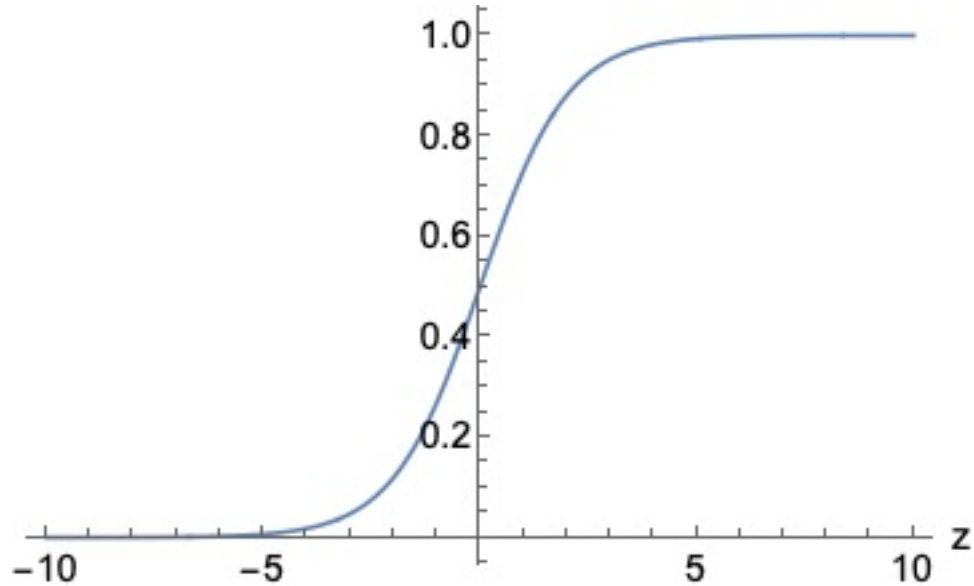


[l-1]                              [l]

- $f$ is what we call an "activation function"

- There are many activation functions, and new ones are invented all the time

- Many of these functions do just fine, or slightly better than existing ones

Source: Liang

BAYES
BUSINESS SCHOOL
CITY, UNIVERSITY OF LONDON

# Typical activation functions: logistic (sigmoid) function

### Logistic (sigmoid) function



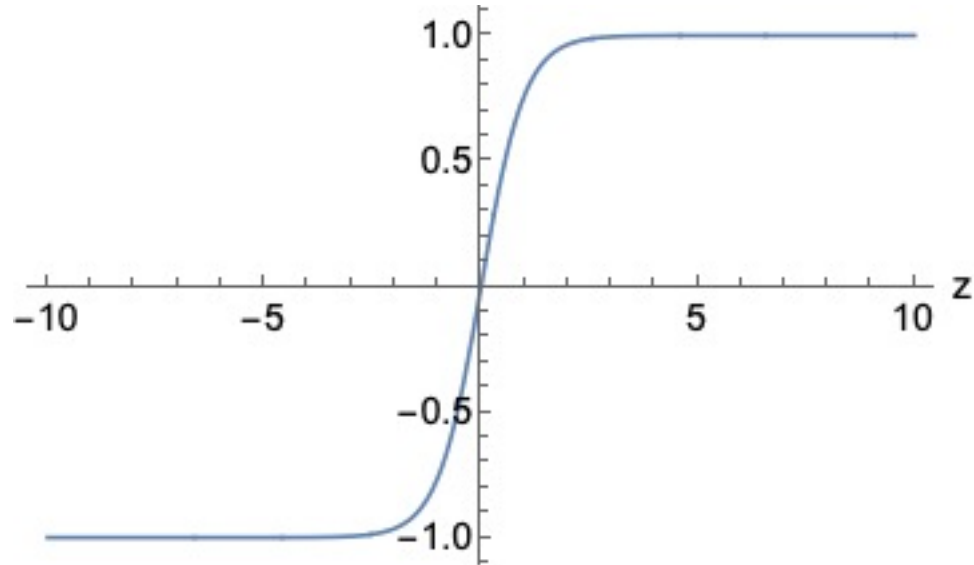$$f(z) = \sigma(z) = \frac{1}{1 + e^{-z}}$$

### Derivative



$$f'(z) = \sigma(z)\big(1 - \sigma(z)\big)$$

BAYES
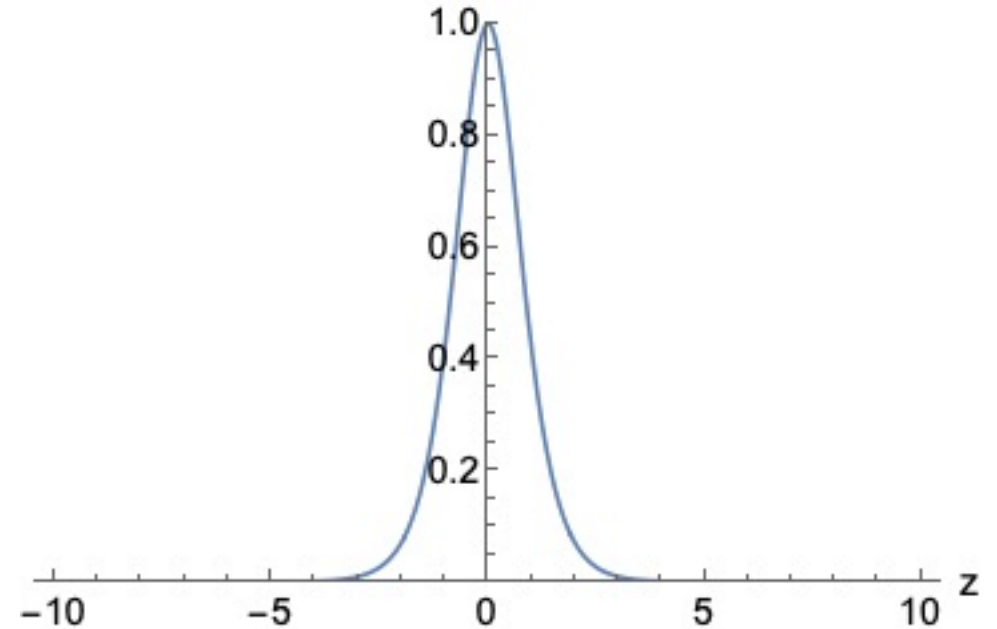BUSINESS SCHOOL
CITY, UNIVERSITY OF LONDON

# Typical activation functions: hyperbolic tangent

### Hyperbolic tangent



$$f(z) = \tanh(z) = \frac{e^{2z} - 1}{e^{2z} + 1}$$

### Derivative


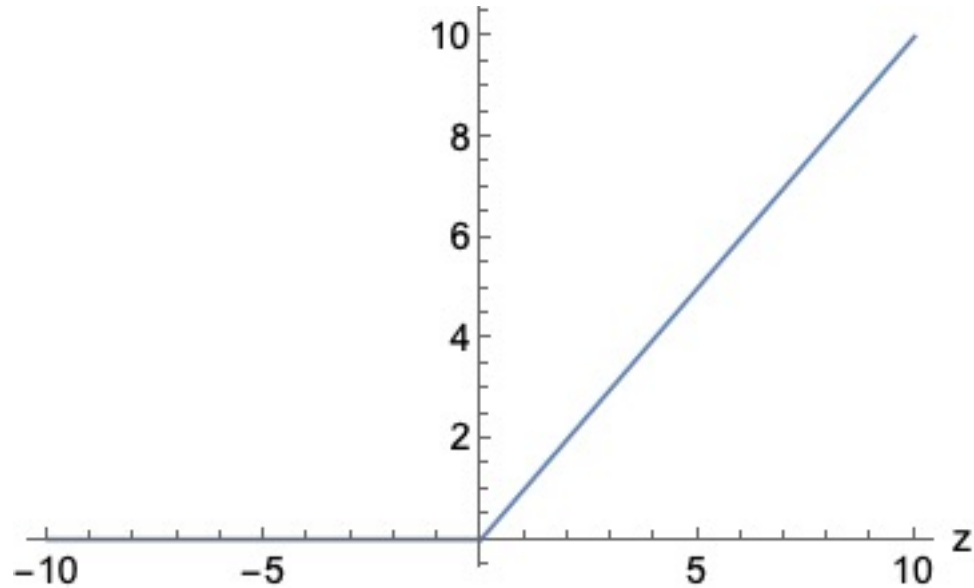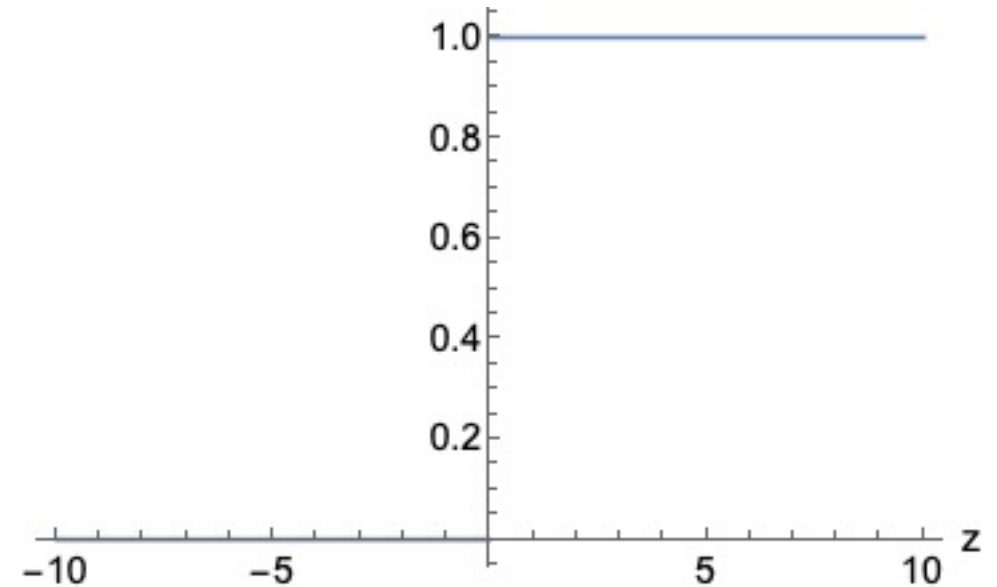
$$f'(z) = \operatorname{sech}(z)^2$$

Rectified Linear Unit (ReLU)

"Derivative"



$$f(z) = \max\{0, z\}$$

$$f'(z) = \begin{cases} 0, \; if \; z < 0 \\ 1, \; if \; z > 0 \end{cases}$$

BAYES
BUSINESS SCHOOL
CITY, UNIVERSITY OF LONDON

Leaky ReLU



$$f(z) = \max\{0.1, z\}$$

"Derivative"



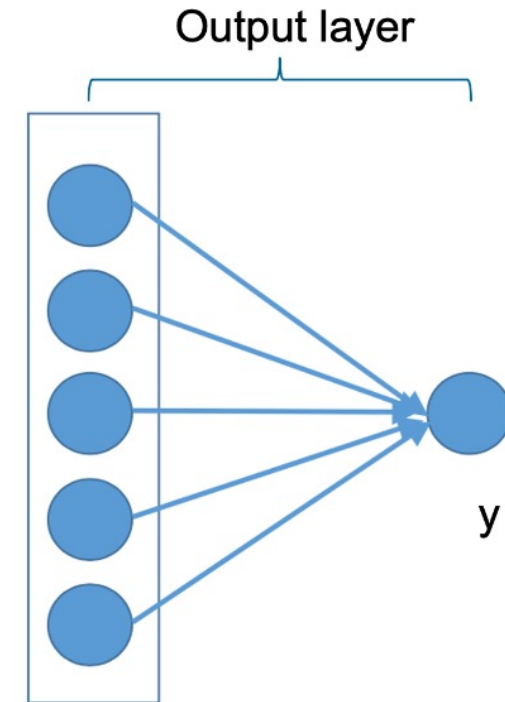$$f'(z) = \begin{cases} 0.1, if\ z < 0 \\ 1, if\ z > 0 \end{cases}$$

# Binary classification
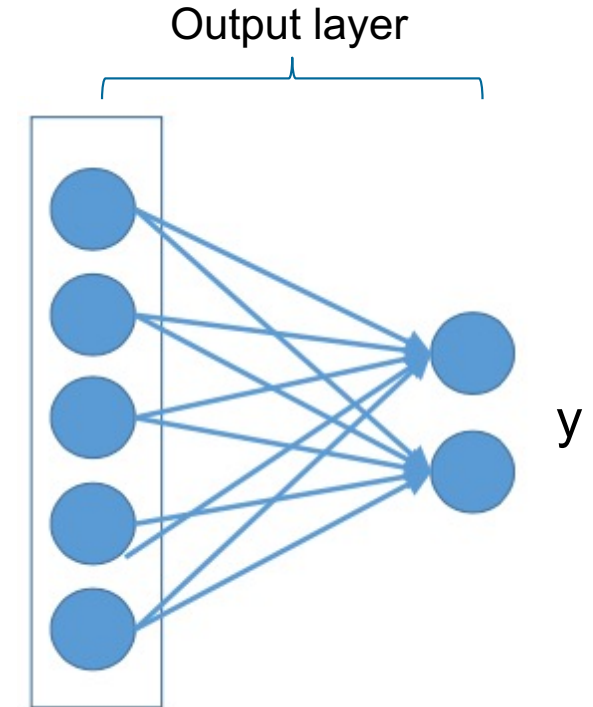
- Input: $\boldsymbol{a}^{[L-1](i)}$

- As usual, we make a linear transformation:
$z^{[L](i)} = \boldsymbol{a}^{[L-1](i)}\boldsymbol{w}^{[L]} + b^{[L]}$

- We then use the logistic sigmoid function
$\hat{y}^{(i)} = f\left(z^{[L](i)}\right) = \sigma(z^{[L](i)}) = \frac{1}{1+e^{-z^{[L](i)}}}$

- We can interpret the output as the probability
of $y^{(i)} = 1$

Output layer



Source: Liang

BAYES
BUSINESS SCHOOL
CITY, UNIVERSITY OF LONDON

# Multi-class classification

- We again make a linear transformation with a matrix of weights: $z^{[L](i)} = a^{[L-1](i)}W^{[L]} + b^{[L]}$

- Note that $z^{[L](i)} = \left( z_1^{[L](i)} \quad z_2^{[L](i)} \quad \cdots \quad z_K^{[L](i)} \right)$

- We then use the softmax function on each of the outputs:
$$\hat{y}_k^{(i)} = f\left(z^{[L](i)}\right) = \frac{e^{-z_k^{[L](i)}}}{\sum_{k=1}^{K} e^{-z_k^{[L](i)}}}$$

- This implies that $\hat{y}_k^{(i)} \in (0,1)$ and $\sum_{k=1}^{K} \hat{y}_k^{(i)} = 1$

- Hence, we can interpret $\hat{y}_k^{(i)}$ as the probability that $y^{(i)} = k$ ("belongs to class $k$")

Output layer



y

Source: Liang

- Recall from logistic regression:

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

- Generally, to learn parameters $\boldsymbol{\theta}$, we define the cross-entropy

$$J(\boldsymbol{\theta}) = \frac{1}{n}\sum_{i=1}^{n} L^{(i)} = -\frac{1}{n}\sum_{i=1}^{n}\ln p_{\boldsymbol{\theta}}\left(y^{(i)}\big|\boldsymbol{x}^{(i)}\right)$$

- Also known as "maximum likelihood estimator"

- Mean square error tends to perform poorly, especially when we have activation functions with $e^z$

1. Decide a "learning rate" $\alpha$
2. Start with some parameters $\boldsymbol{\theta}$ and compute $J(\boldsymbol{\theta})$      (forward propagation)
3. Until $J$ "doesn't change" anymore:
   - Let $\boldsymbol{\theta} := \boldsymbol{\theta} - \alpha \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$      (back-propagation)
   - Recompute $J(\boldsymbol{\theta})$      (forward propagation)
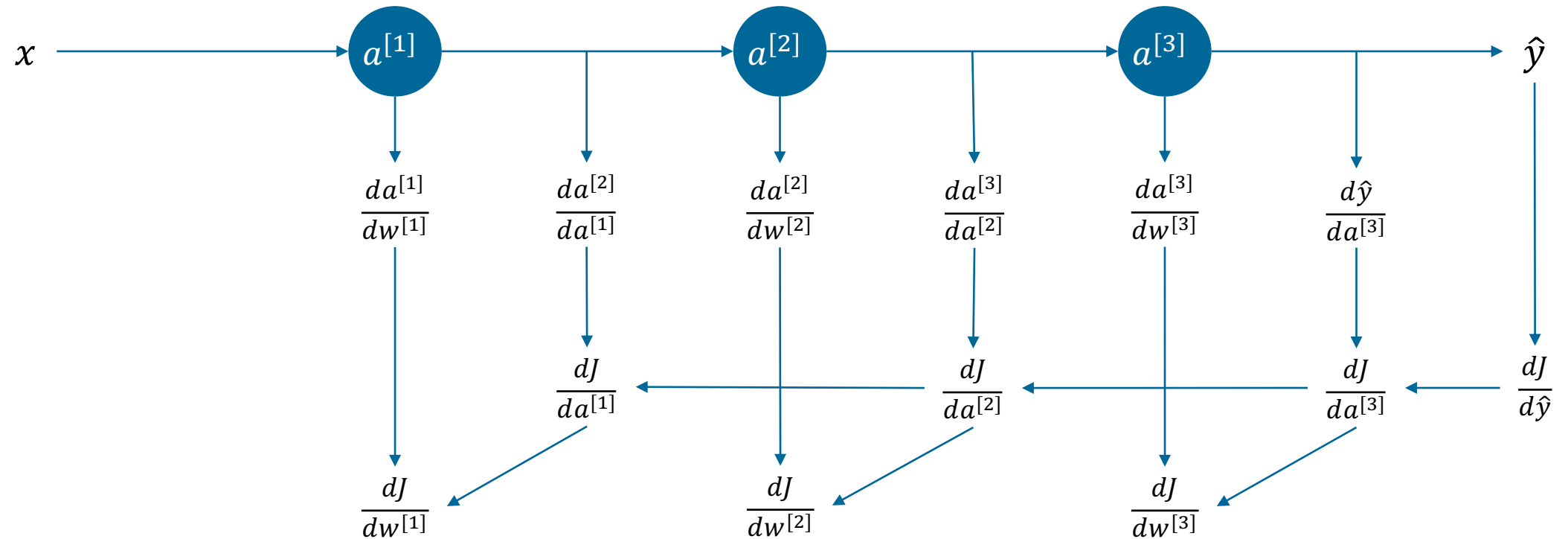
- Initialize weights to small random values ( e.g., *np.random.randn(**shape of W**) * 0.01* )

- Bias terms can be initialized randomly, but can also just be initialized to zero ( e.g., *np.zeros(**shape of b**)* )

See you next week!