**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 discuss the implementation of 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
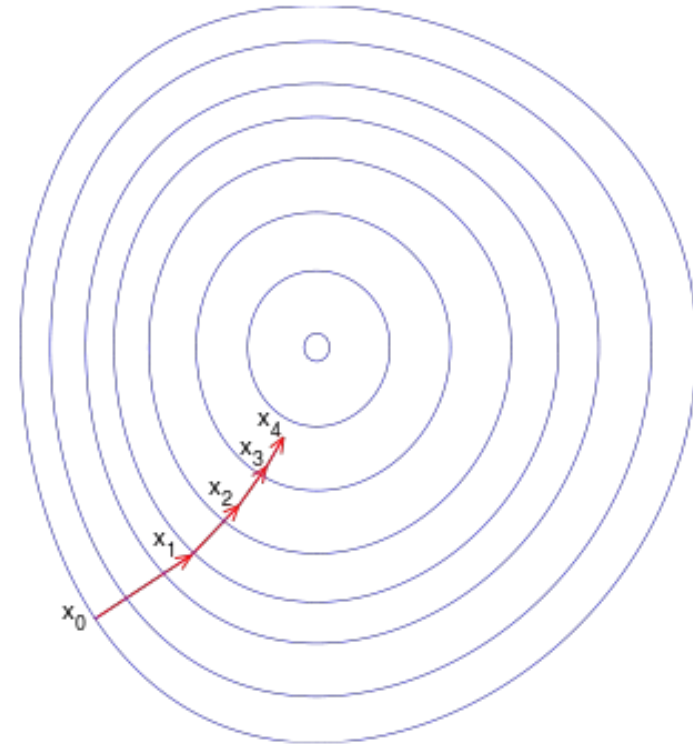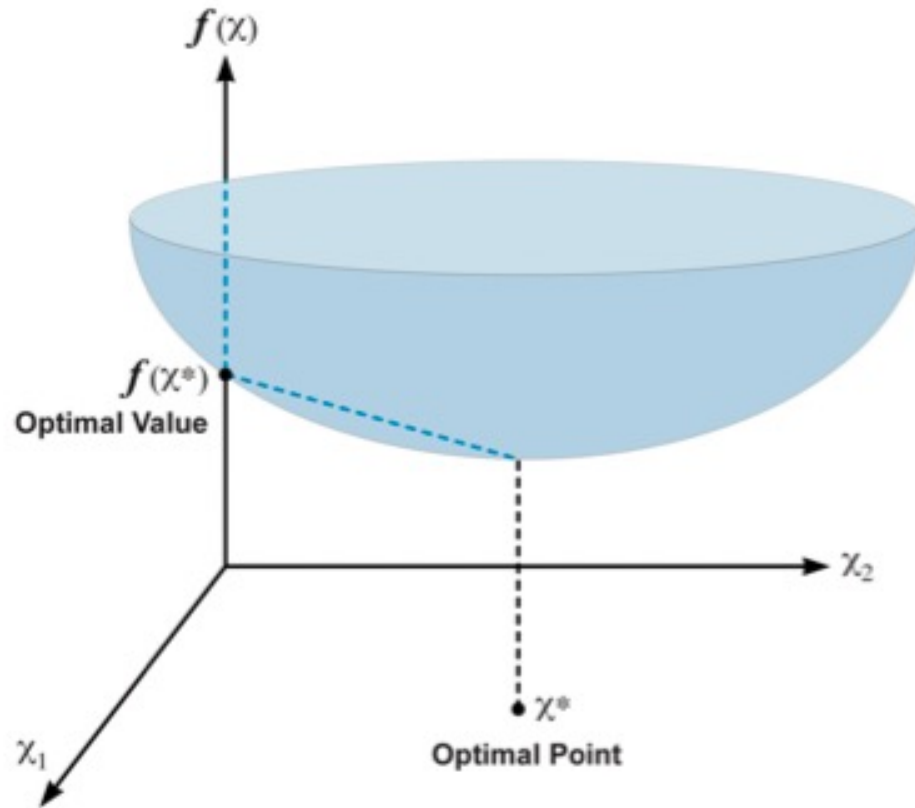- In the tutorial, we will implement more complex networks

BAYES
BUSINESS SCHOOL
CITY, UNIVERSITY OF LONDON

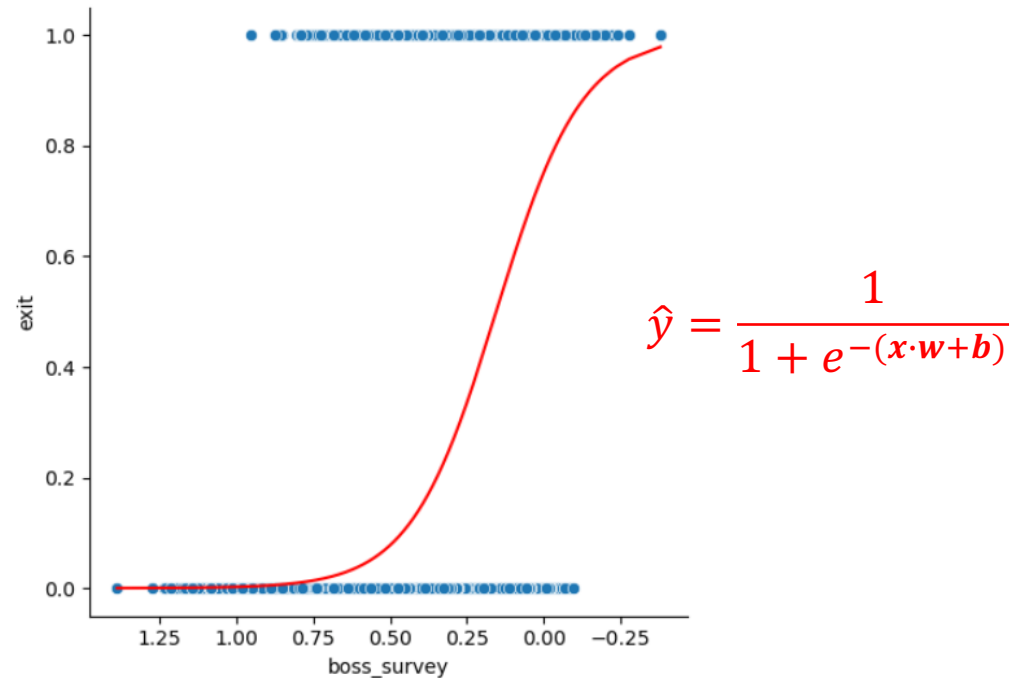# Recap – logistic regression

# Gradient descent – the idea

- 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 $y^{(i)} = 1$: $\hat{y}^{(i)} = P\big(y^{(i)} = 1 \big| \boldsymbol{x}^{(i)}\big)$

- We model this probability, using the sigmoid function:

- 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:



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

- Remember that $\boldsymbol{w} \in \mathbb{R}^m$ and $b \in \mathbb{R}$
- To get to the "right" model, we optimize our parameters $\boldsymbol{w}, b$ 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" $J(\boldsymbol{w}, b)$, where $\hat{y}^{(i)} = \dfrac{1}{1+e^{-\left(\boldsymbol{x}^{(i)}\boldsymbol{w}+\boldsymbol{b}\right)}}$ :

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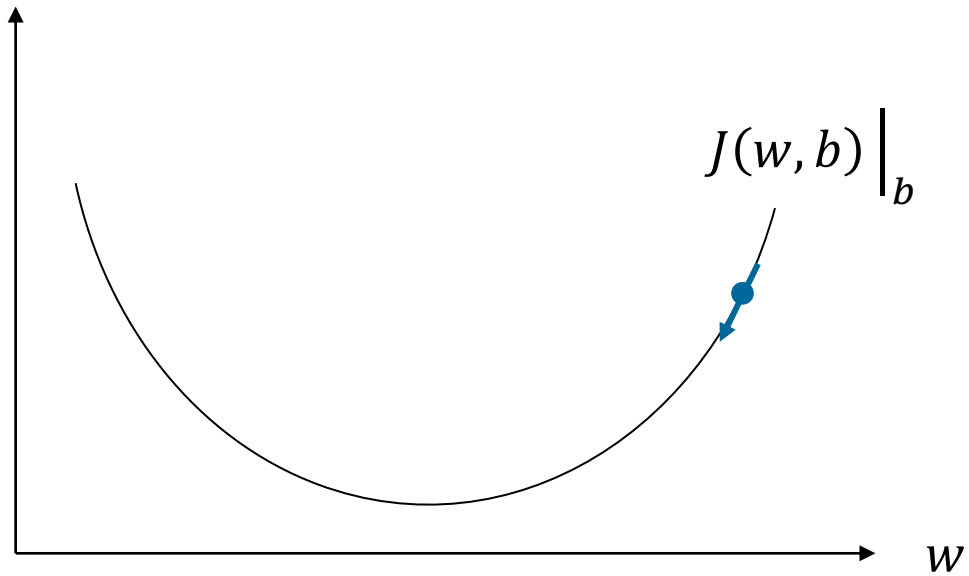
- Remember that $\boldsymbol{w} \in \mathbb{R}^m$ and $b \in \mathbb{R}$
- To get to the "right" model, we optimize our parameters $\boldsymbol{w}, b$ so that the $\hat{y}^{(i)}$s are "as close as possible" to the $y^i$s
- What we do is to minimize the "cost-function" $J(\boldsymbol{w}, b)$, where $\hat{y}^{(i)}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

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

$$J(w, b)\Big|_b$$

$$w$$

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

$w_1$

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

$w_2$

$b$

$x_2^{(i)} \longrightarrow \boxed{a_2^{(i)} = w_2 x_2^{(i)}} \longrightarrow \boxed{\begin{array}{l} z^{(i)} \\ = a_1^{(i)} + a_2^{(i)} \\ + a_3^{(i)} + \cdots + b \end{array}} \longrightarrow \boxed{\hat{y}^{(i)} = \sigma\left(z^{(i)}\right)} \longrightarrow \boxed{\begin{array}{l} L^{(i)} \\ = -\left[y^{(i)} \ln \hat{y}^{(i)} \right. \\ \left. + \left(1 - y^{(i)}\right) \ln\left(1 - \hat{y}^{(i)}\right)\right] \end{array}}$

$w_3$

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

$\cdots$

BAYES
BUSINESS SCHOOL
CITY, UNIVERSITY OF LONDON

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

From logistic regression to neural network
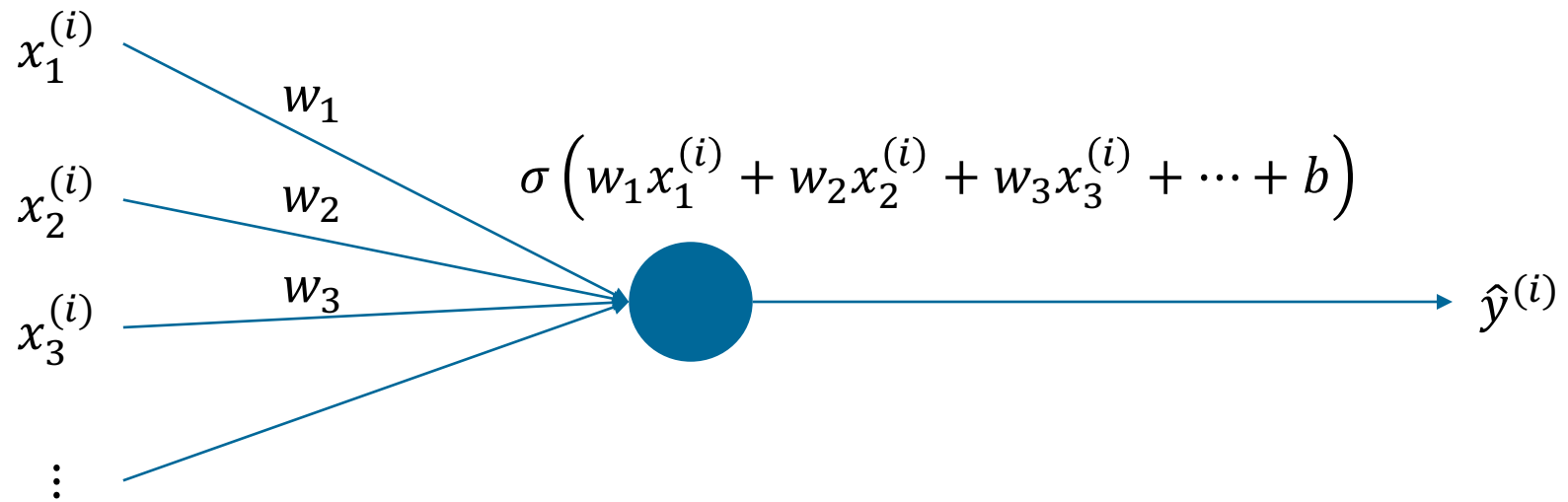
# Schema of a logistic regression



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

$x_1^{(i)}$

$w_1$

$x_2^{(i)}$

$w_2$

$x_3^{(i)}$

$w_3$

$\vdots$

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

$x_1^{(i)}$

$x_2^{(i)}$

$x_3^{(i)}$

$\vdots$

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

Source: Czarnecki

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"
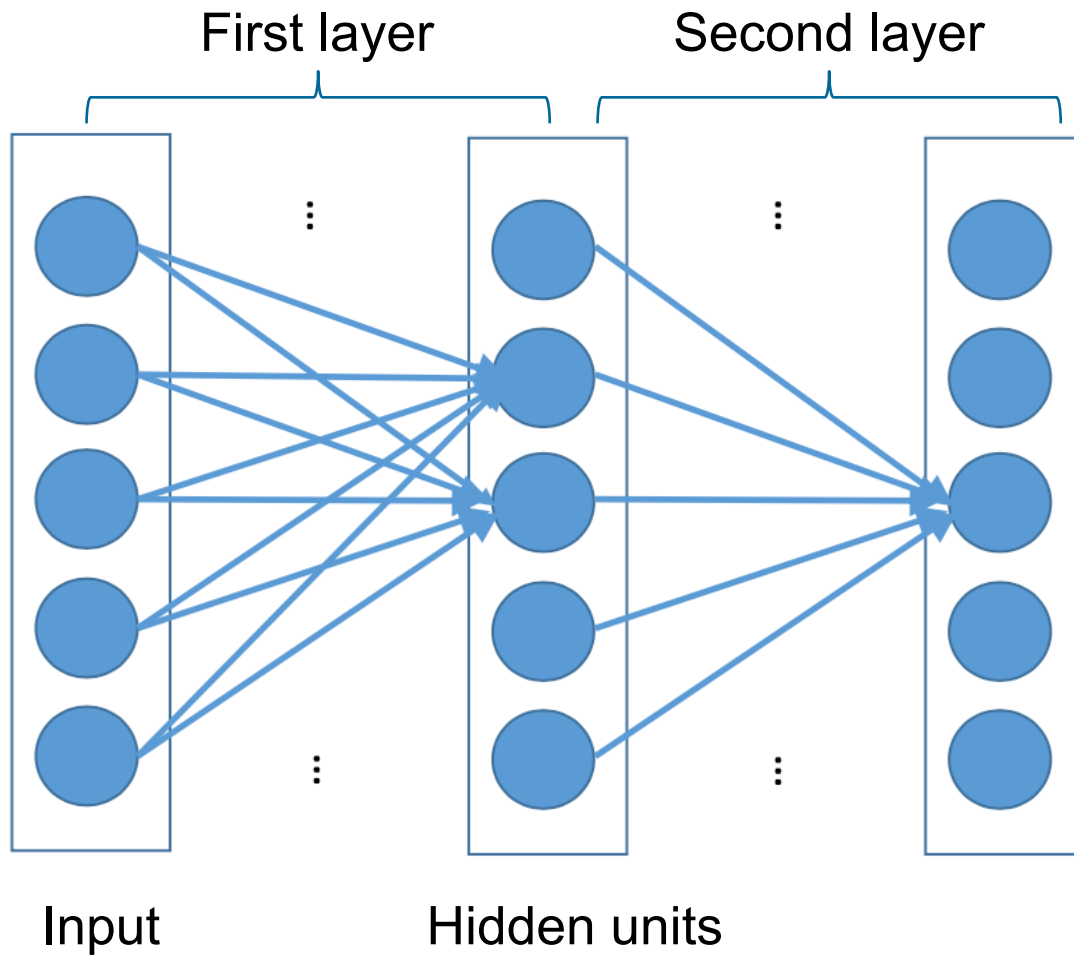- How do your answers change when you use 6 neurons instead?

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

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

[l-1]          [l]

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

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

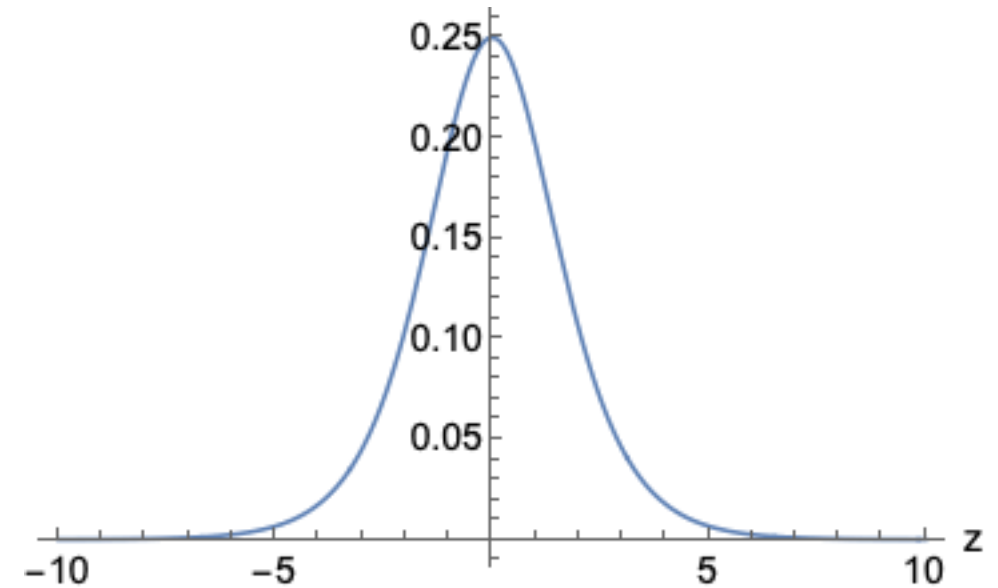- Many of these functions do just fine, or slightly better than existing ones

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



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

## Hyperbolic tangent



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

## Derivative



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

Rectified Linear Unit (ReLU)



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

"Derivative"



$$f'(z) = \begin{cases} 0, & if \ z < 0 \\ 1, & if \ z > 0 \end{cases}$$
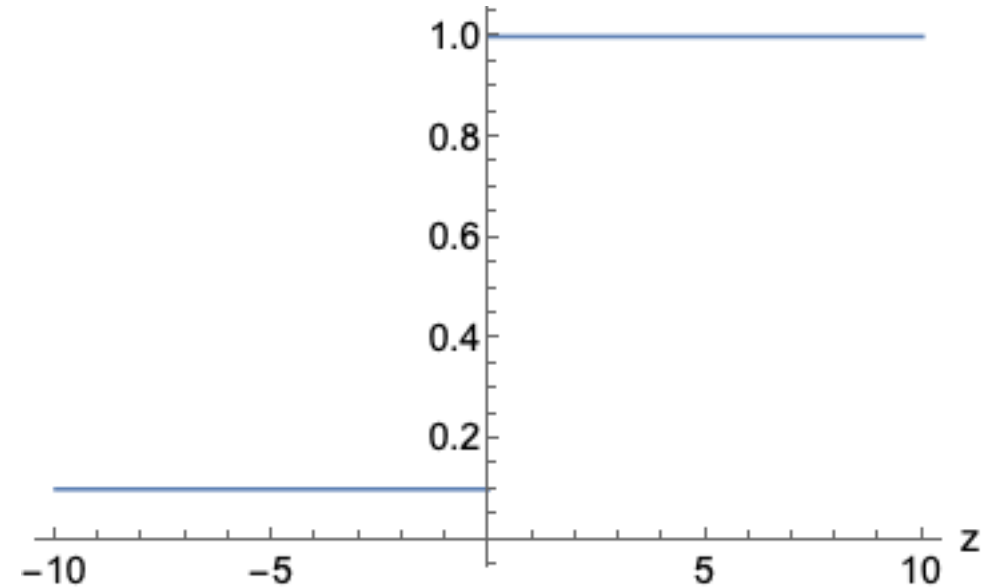
**BAYES**
BUSINESS SCHOOL
CITY, UNIVERSITY OF LONDON

## Leaky ReLU



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

## "Derivative"
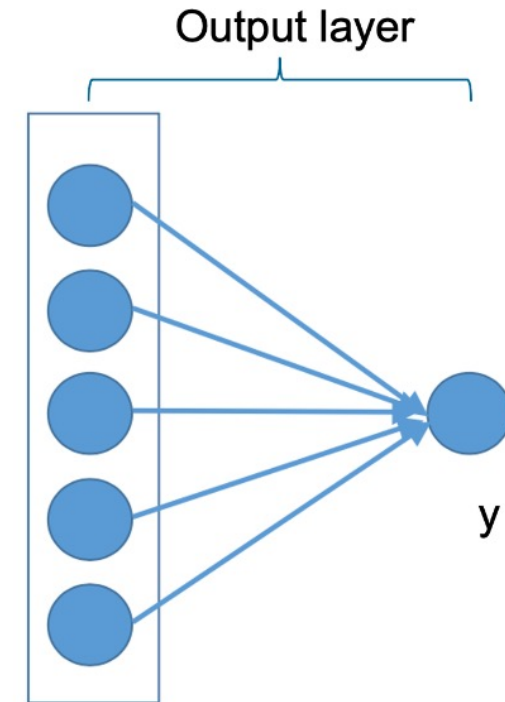


$$f'(z) = \begin{cases} 0.1, & if\ z < 0 \\ 1, & if\ z > 0 \end{cases}$$

# Binary classification

- Input: $a^{[L-1](i)}$

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

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

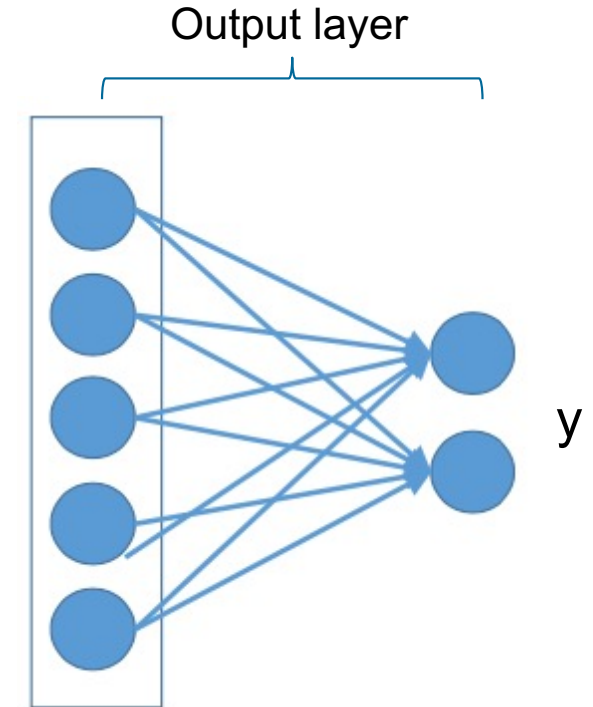- We can interpret the output as the probability
of $y^{(i)} = 1$

Output layer

y

# Multi-class classification

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

- Note that $\mathbf{z}^{[L](i)} = \begin{pmatrix} z_1^{[L](i)} & z_2^{[L](i)} & \cdots & z_K^{[L](i)} \end{pmatrix}$

- We then use the softmax function on each of the outputs:

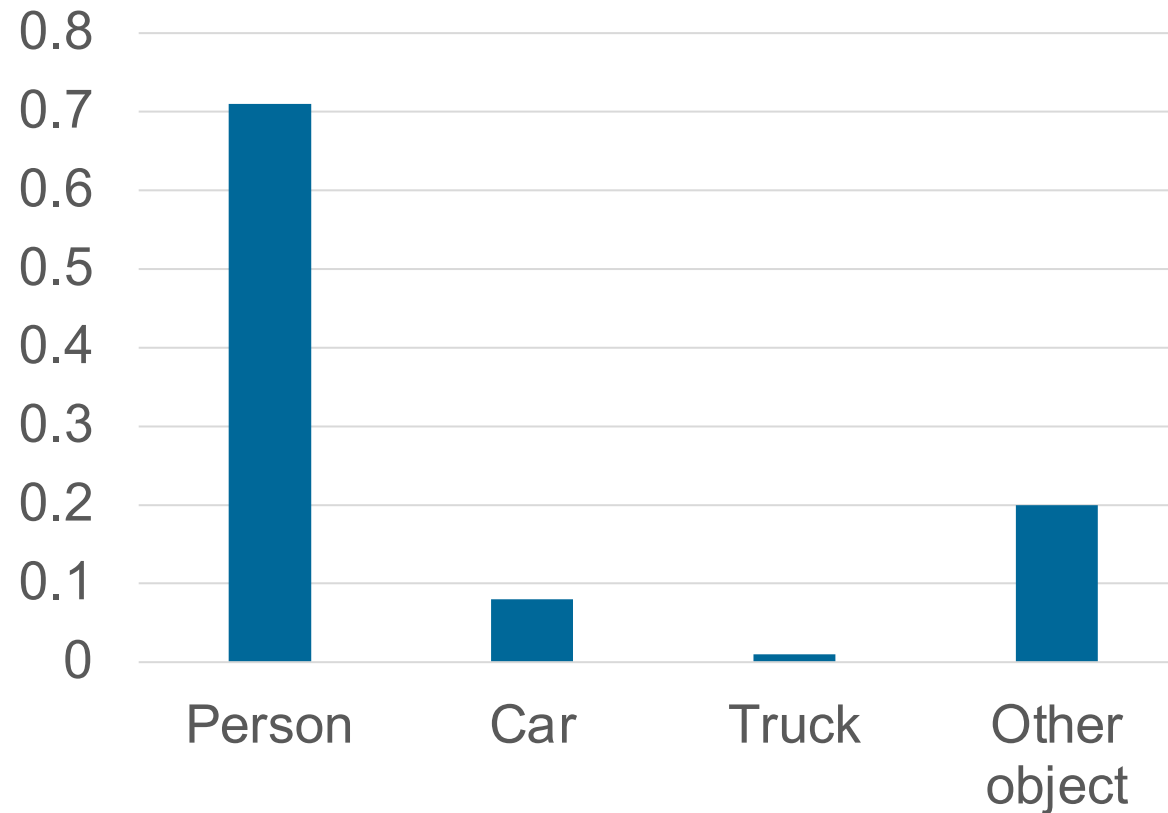$$\hat{y}_k^{(i)} = f\big(\mathbf{z}^{[L](i)}\big) = \frac{e^{-z_k^{[L](i)}}}{\sum_{k=1}^K e^{-z_k^{[L](i)}}}$$

- 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

BAYES
BUSINESS SCHOOL
CITY, UNIVERSITY OF LONDON

- E.g., when performing object recognition, we might represent our prediction $\widehat{\boldsymbol{y}}^{(i)}$ as

- 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 $e^z$

# Learning with gradient descent

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)

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

- 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**)* )
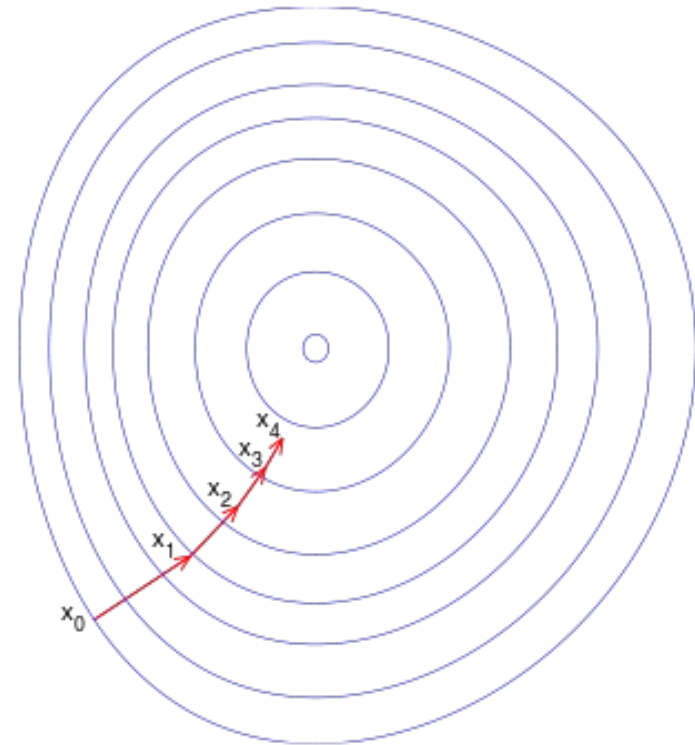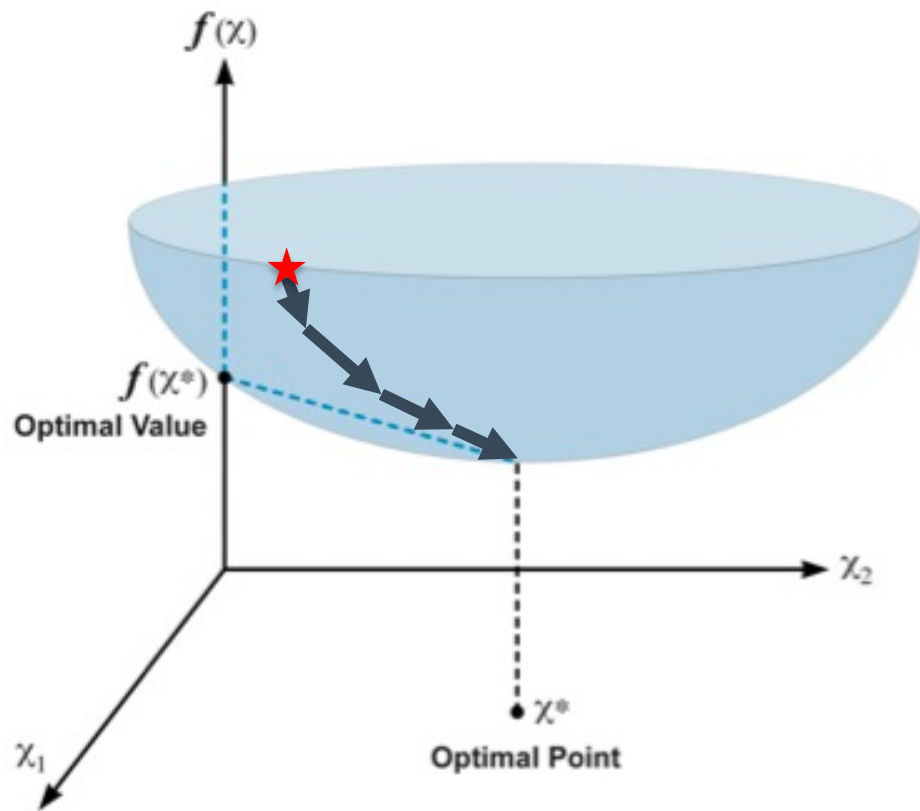
$x \longrightarrow \quad a^{[1]} \longrightarrow a^{[2]} \longrightarrow a^{[3]} \longrightarrow \hat{y}$

$w^{[1]} \qquad\qquad w^{[2]} \qquad\qquad w^{[3]}$

→ No matter which approach, when the function is convex, we will find a global minimum

# Pathological non-convex cases



(gradient *ascend*)

Source: Wikipedia

See you next week!

# Sources

- 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
- Liang, 2016, Introduction to Deep Learning: https://www.cs.princeton.edu/courses/archive/spring16/cos495/
- Wikipedia, n.d., Gradient ascent: https://en.wikipedia.org/wiki/File:Gradient_ascent_(surface).png

BAYES
BUSINESS SCHOOL
CITY, UNIVERSITY OF LONDON