# DEEP LEARNING

$x^{(i)} \in \mathbb{R}^d$  TRAIN DATA

$a^{(0)} = x$

$\left.\begin{array}{l}\end{array}\right]$ TRAIN DATA

for $l$ in $1, \ldots, L$

$\qquad z^{[l]} = W^{[l]} a^{[l-1]} + b^{[l]}$

$\qquad a^{[l]} = g(z^{[l]})$

MODEL WITH

L LAYERS

$\hat{y} = a^{[L]}$

$\mathcal{L} = \text{LOSS}(y, \hat{y})$

PREDICTION

AND LOSS

IN ORDER TO TRAIN THE MODEL, THE APPROACH IS TO MAXIMIZE THE LIKELIHOOD $\mathcal{L}$ (MINIMIZE THE $-\log \mathcal{L}$ )

## EX (*)

$x \in \mathbb{R}^d$ , $y \in \{0, 1\}$ (CLASSIFICATION PROBLEM)

$\mathcal{L} = -\left[ y \log \hat{y} + (1-y) \log 1-y \right]$
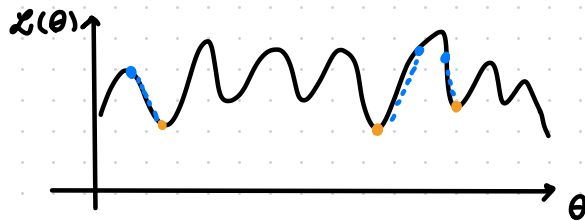
$$\boxed{\hat{y} = \text{Model}_\theta(x)}$$

# $\mathcal{L}$ OPTIMIZATION : BACKPROPAGATION

<u>NOTE</u> THE PARAMETERS ARE UPDATED $\forall_\ell$, BUT $\mathcal{L}$ IS APPLIED ONLY AT L (OUTPUT LAYER)

WE NEED TO COMPUTE $\nabla\mathcal{L}$ IN ORDER TO MAX THE LIKELIHOOD AND THEN APPLY GRADIENT DESCENT $\theta = \theta - \alpha \nabla\mathcal{L}$. IN NEURALNETS, THE OPTIMIZATION ALG. IS BACKPROP

<u>NOTE</u> DEEP LEARNING MODELS ARE <u>NOT</u> CONVEX! HOW DO WE FIND $\theta$?



- CONVERGE
- RAND INIT

LOCAL MINIMA? GLOBAL?
LOCAL MOST LIKELY!

LET'S DERIVE BACKPROP STARTING FROM THE EXAMPLE (*)

# ALGORITHM

→ INIT

$$\begin{cases} W^{[\ell]} \sim N\left(\vec{0}, \sqrt{\dfrac{2}{n^{[\ell]} + n^{[\ell-1]}}}\right) \\[2em] \vec{b} = 0 \end{cases}$$

RANDOM INIT OF W

BIAS INIT AT $\varnothing$

WE DEFINE $\Theta$ AS

$$\Theta = \left\{ W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]}, W^{[3]}, b^{[3]} \right\} \quad (3 \text{ LAYERS})$$

$$x^{(i)} \in \mathbb{R}^{d_0}$$

$$W^{[1]} \in \mathbb{R}^{d_1 \times d_0}$$

$$b^{[1]} \in \mathbb{R}^{d_1}$$

$$z^{[1]} = W^{[1]} \overbrace{a^{[0]}}^{x^{(i)}} + b^{[1]} \in \mathbb{R}^{d_1} \quad (\text{LOGIT})$$

$$a^{[1]} = g(z^{[1]}) \in \mathbb{R}^{d_1}$$

$$W^{[2]} = \mathbb{R}^{d_2 \times d_1}$$

MODEL
PARAMS

WE DEFINE $\dfrac{\partial \mathcal{L}}{\partial \Theta}$ AS THE FOLLOWING TENSOR (CONTAINS THE INFOS REGARDING DERIVATIV
OF $\mathcal{L}$ W.R.T. EACH WEIGHT & BIAS)

$$\frac{\partial \mathcal{L}}{\partial \Theta} = \begin{bmatrix} \dfrac{\partial \mathcal{L}}{\partial W^{[1]}} , & \dfrac{\partial \mathcal{L}}{\partial b^{[1]}} \\[2mm] \dfrac{\partial \mathcal{L}}{\partial W^{[2]}} , & \dfrac{\partial \mathcal{L}}{\partial b^{[2]}} \\[2mm] \dfrac{\partial \mathcal{L}}{\partial W^{[3]}} , & \dfrac{\partial \mathcal{L}}{\partial b^{[3]}} \end{bmatrix}$$

LET'S ANALYZE EACH MEMBER $\frac{\partial \mathcal{L}}{\partial W^{[i]}}$

<u>EX</u>

$$\frac{\partial \mathcal{L}}{\partial W^{[2]}} \longrightarrow \begin{cases} \mathcal{L} \in \mathbb{R} \\ W^{[2]} \in \mathbb{R}^{d_2 \times d_1} \\ \frac{\partial \mathcal{L}}{\partial W^{[2]}} \in \mathbb{R}^{d_2 \times d_1} \end{cases}$$

THE OPTIMIZATION OPERATION, IN GENERAL, LOOKS LIKE THIS (G.D.)

$$\Theta = \Theta - \alpha \frac{\partial \mathcal{L}}{\partial \Theta}$$

WE SPLIT THIS OPERATION INTO SUBPROBLEMS

$$\begin{cases} W^{[1]} = W^{[1]} - \alpha \frac{\partial \mathcal{L}}{\partial W^{[1]}} \\ \quad \vdots \\ b^{[3]} = b^{[3]} - \alpha \frac{\partial \mathcal{L}}{\partial b^{[3]}} \end{cases} \longrightarrow$$

$$W^{[2]} = \begin{bmatrix} W_{11}^{[2]} & \cdots & W_{1\,d_1} \\ \vdots & \ddots & \vdots \\ W_{d_2\,1}^{[2]} & \cdots & W_{d_2\,d_1} \end{bmatrix} \longrightarrow \frac{\partial \mathcal{L}}{\partial W^{[2]}} = \begin{bmatrix} \frac{\partial \mathcal{L}}{\partial W_{11}^{[2]}} & \cdots & \frac{\partial \mathcal{L}}{\partial W_{1\,d_1}} \\ \frac{\partial \mathcal{L}}{\partial W_{d_2\,1}^{[2]}} & \cdots & \frac{\partial \mathcal{L}}{\partial W_{d_2\,d_1}} \end{bmatrix}$$

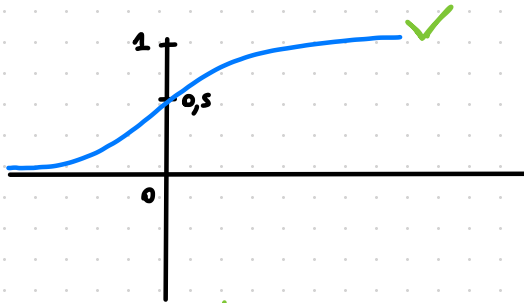WE DO THIS $\forall W, b$ AND GET A SET OF PARTIAL DERIVATIVES

DEPENDING ON THE NATURE OF $\hat{y}$, THE LAST NON-LINEARITY FUNCTION SHOULD BE CHOOSED ACCORDINGLY
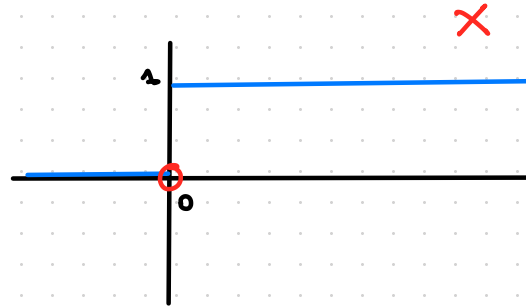
## Ex

$$y \in \{0, 1\}$$

$$\hat{y} \in \begin{cases} 0 & z < 0 \\ z & z \geq 0 \end{cases}$$

**NOTE** HAS TO BE SOME NON-LINEARITY DIFFERENTIABLE IN EACH POINT



$$\hat{y} = \frac{1}{1 + e^{(-z)}} \qquad \hat{y} = \mathbb{1}\left[z > 0\right]$$

BY APPLYING THE SIGMOID $\sigma(z)$ FUNCTION AS NON LINEARITY AND CHOOSE A BINARY CLASSIFICATION PROBLEM WE GET THE FOLLOWING RESULT

$$\begin{cases} \hat{y} = \dfrac{1}{1+e^{-z}} \\[2mm] \mathcal{L} = y\log\hat{y} + (1-y)\log(1-\hat{y}) \end{cases} \longrightarrow \boxed{\dfrac{\partial \mathcal{L}}{\partial z} = \hat{y} - y}$$

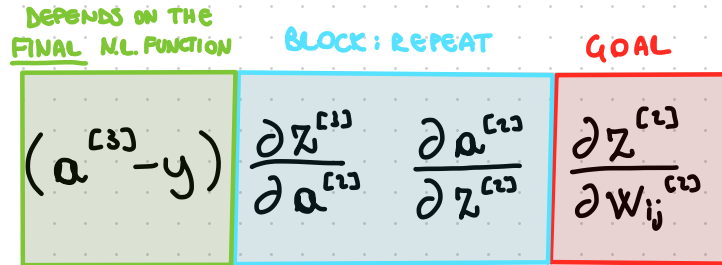WHICH IS EXTREMELY IMPORTANT IN THE DEFINITION PROCESS OF BACKPROP

FROM THE DEFINITIONS OF THE MODEL, WE CAN APPLY THE CHAIN RULE OF DERIVATION IN ORDER TO OBTAIN $\dfrac{\partial \mathcal{L}}{\partial W_{ij}^{(1)}}$. THE CHAIN RULE CONSISTS OF BREAKING DOWN DERIVATIVES INTO JACOBIANS

<u>EX</u>

$$\frac{\partial \mathcal{L}}{\partial W_{ij}^{(1)}} = \frac{\partial \mathcal{L}}{\partial \hat{y}}\,\frac{\partial \hat{y}}{\partial W_{ij}^{(1)}} =$$

$$= \frac{\partial \mathcal{L}}{\partial \hat{y}}\,\frac{\partial \hat{y}}{\partial W_{ij}^{(1)}} = \frac{\partial \mathcal{L}}{\partial a^{(2)}}\,\frac{\partial a^{(2)}}{\partial W_{ij}^{(1)}} =$$

$$= \frac{\partial \mathcal{L}}{\partial a^{(2)}}\,\frac{\partial a^{(2)}}{\partial W_{ij}^{(1)}} = \frac{\partial \mathcal{L}}{\partial a^{(2)}}\,\frac{\partial a^{(2)}}{\partial z^{(2)}}\,\frac{\partial z^{(2)}}{\partial W_{ij}^{(1)}}$$

$$\frac{\partial \mathcal{L}}{\partial a^{[3]}} \frac{\partial a^{[3]}}{\partial z^{[3]}} \frac{\partial z^{[3]}}{\partial W_{ij}^{[3]}} = \overbrace{\frac{\partial \mathcal{L}}{\partial a^{[3]}} \frac{\partial a^{[3]}}{\partial z^{[3]}}}^{\frac{\partial \mathcal{L}}{\partial z^{2}}} \frac{\partial z^{[3]}}{\partial a^{[2]}} \frac{\partial a^{[2]}}{\partial W_{ij}^{[2]}} =$$

$$= \underbrace{\frac{\partial \mathcal{L}}{\partial z^{[3]}}}_{(a^{[3]}-y)} \frac{\partial z^{[3]}}{\partial a^{[2]}} \frac{\partial a^{[2]}}{\partial W_{ij}^{[2]}} = \left( a^{[3]} - y \right) \frac{\partial z^{[3]}}{\partial a^{[2]}} \frac{\partial a^{[2]}}{\partial W_{ij}^{[2]}}$$

AND FINALLY OBTAINING OUR GOAL EXPRESSION

DEPENDS ON THE FINAL N.L. FUNCTION        BLOCK: REPEAT        GOAL

$$\left( a^{[3]} - y \right) \quad \frac{\partial z^{[3]}}{\partial a^{[2]}} \frac{\partial a^{[2]}}{\partial z^{[2]}} \quad \frac{\partial z^{[2]}}{\partial W_{ij}^{[2]}}$$

IN GENERAL

- GOAL: CALCULATE $\dfrac{\partial \mathcal{L}}{\partial W_{ij}^{[L]}}$

- BLOCKS: $\dfrac{\partial z^{[L]}}{\partial a^{[L-1]}} \dots \dfrac{\partial z^{[\ell+1]}}{\partial a^{[\ell]}} \dfrac{\partial a^{[\ell]}}{\partial z^{[\ell]}} \cdot {\color{red}\dfrac{\partial z^{[\ell]}}{\partial W_{ij}^{[\ell]}}}$

THE OBJECTIVE IS TO REACH THE $z$ OF THE $\longrightarrow$ GO BACKWARDS! CORRESPONDING LAYER (OF THE W,b W.R.T. WE ARE DERIVATING)

# LET'S DO SOME ANALYSIS OF THE DIMENSIONS

$$a^{[2]} = g(z^{[2]})$$

$$\frac{\partial a^{[2]}}{\partial z^{[2]}} = \text{diag}\left(g'(z^{[2]})\right)$$

→ JACOBIAN MATRIX (DERIVATIVE APPLIED TO EACH ELEMENT)

→ IT IS A DIAGONAL MATRIX

→ $g$ IS APPLIED ELEMENT-WISE

$$\left(a^{[3]} - y\right) \quad \frac{\partial z^{[3]}}{\partial a^{[2]}} \quad \frac{\partial a^{[2]}}{\partial z^{[2]}} \quad \frac{\partial z^{[2]}}{\partial W_{ij}^{[2]}}$$

$\mathbb{R}$

$$z^{[3]} = W^{[3]} a^{[2]} + b^{[3]} \in \mathbb{R}^{d_3}$$

$$\frac{\partial z^{[3]}}{\partial a^{[2]}} = W^{[3]} \in \mathbb{R}^{d_3 \times d_2}$$

NOTE   WE CAN THINK IT IN THIS WAY

$$z^{[L]} \longrightarrow b^{[L]}$$
$$z^{[L]} \longrightarrow W^{[L]}$$

$$\cdots \leftarrow a^{[L-1]} \leftarrow z^{[L]} \leftarrow a^{[L]} \leftarrow \hat{y} \leftarrow \mathcal{L}$$

WE OBTAIN THE FOLLOWING FORM

$$z^{[2]} = W^{[2]} a^{[1]} + b^{[2]}$$

$$\begin{bmatrix} z_i^{[2]} \end{bmatrix} = \begin{bmatrix} \leftarrow W_{ij}^{[2]} \rightarrow \end{bmatrix} \begin{bmatrix} a_j^{[1]} \end{bmatrix} + \begin{bmatrix} b^{[2]} \end{bmatrix}$$

$$z_i^{[2]} = \sum_j W_{ij}^{[2]} \underbrace{a_j^{[1]}}_{\mathbb{R}^{d_1}} + b^{[2]}$$

$$\frac{\partial z^{[2]}}{\partial W_{ij}^{[2]}} = \boxed{a_j^{[i]} \cdot C_j} = a_j^{[1]} \cdot e_j$$

WHERE $e_j$ IS THE $j^{TH}$ BASIS VECTOR

$$a_j^{[i]} e_i = a_j^{[i]} \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ a_j^{[1]} \\ \vdots \\ 0 \end{bmatrix} \rightarrow j^{TH}$$

NOW, WE CAN FINALLY WORK OUR WAY TO A CLOSED FORM EXPRESSION

$$\frac{\partial \mathcal{L}}{\partial W_{ij}^{[1]}} = (a^{[3]} - y) \, W^{[3]} \, diag(g'(z^{[1]})) \; \frac{\partial z^{[1]}}{\partial W_{ij}^{[1]}}$$

$$= (a^{[3]} - y) \, W^{[3]} \underline{\odot} \; g'(z^{[1]}) \cdot a_j^{[1]} \, e_i$$

<span style="color:blue">HADAMART PRODUCT
(ELEMENT-WISE MATRIX
MULTIPLICATION)</span>

$$\frac{\partial \mathcal{L}}{\partial W_{ij}^{[1]}} = \left[ (a^{[3]} - y) \, W^{[3]} \odot g'(z^{[1]}) \right]_i \, a_j^{[1]}$$

IN GENERAL

$$\boxed{\frac{\partial \mathcal{L}}{\partial W_{ij}^{[\ell]}} = \left[ (a^{[L]} - y) \, W^{[\ell+1]} \odot g'(z^{[\ell]}) \right]_i \, a_j^{[\ell-1]^T}}$$

WE ARE EXTRACTING THE i-TH TERM AND MULTIPLYING IT BY $a_j$

ONCE WE HAVE THE $z$ OF THE LAST LAYER WE CAN WORK OUR WAY BACK TO OUR $z$ OF INTEREST.

<span style="color:red">**NOTE**</span> COMPUTATIONALLY, THE EXPRESSION WE FOUND IS MUCH MORE CONVENIENT THAN CALCULATING THE CHAIN BY MAT MULT ALL THE $L$ BLOCKS

$$(a^{[L]} - y) W^{[L]} \text{diag}(\ ) W^{[L-1]} \text{diag}(\ ) W^{[\ ]} \text{diag}(\ ) \cdots \color{red}{\times} \quad \widehat{\frown}$$