

*Até o momento da "primeira" backpropagation*

$$w := w - \alpha \frac{\partial J(w, b)}{\partial w}$$

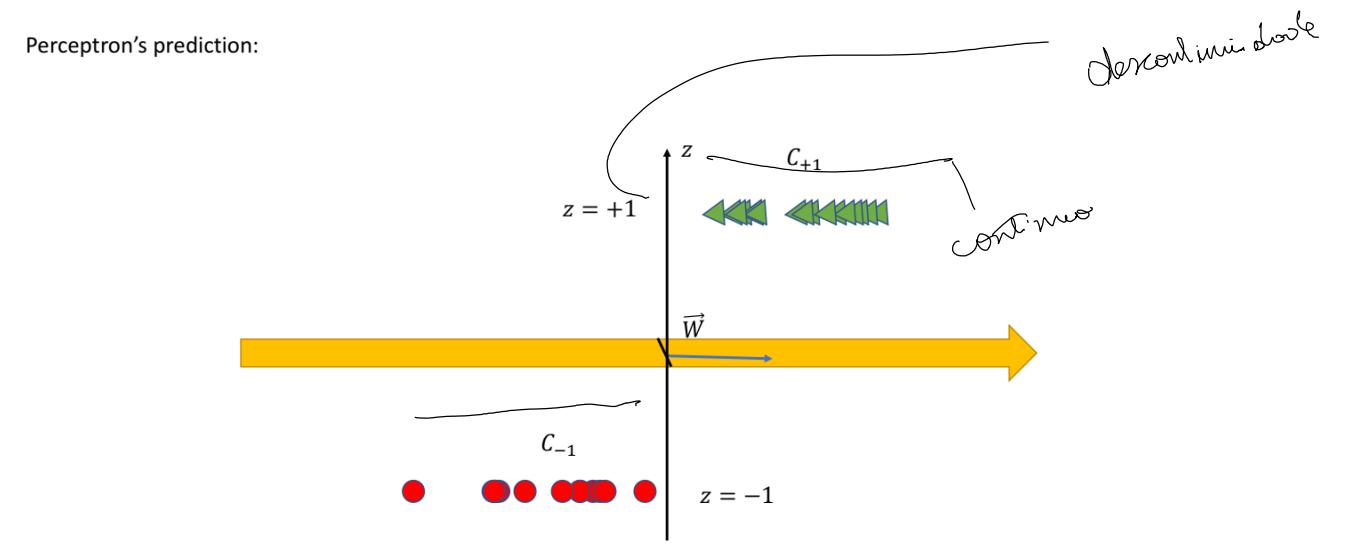
*nome da variável - comando*

$$b := b - \alpha \frac{\partial J(w, b)}{\partial b}$$

*notação de derivadas parciais*

$$\frac{\partial J(w, b)}{\partial w}$$

$$J = \frac{1}{m} \sum_{i=1}^m \log(\sigma(w^T x^{(i)}) + b)$$



→ Vídeo do Andrew Ng

Regras:  
 $z = w^T x + b$   
 $\hat{y} = a = \delta(z)$

Computational Graph

$$z = w_1 x_1 + w_2 x_2 + b \rightarrow \hat{y} = \delta(z) \rightarrow J(a, y)$$

$w_1, w_2, b$

Precisamos modificar os pesos ( $w_1, w_2$  e  $b$ ) para minimizar o loss

• 4 - Back propagation

↳ Aqui vai calculando os derivados parciais, de modo que para cada parâmetro é feito o cálculo de todos os derivados  $\frac{\partial J}{\partial w_1}, \frac{\partial J}{\partial w_2}, \frac{\partial J}{\partial b}$

Summary of gradient descent

$$\begin{aligned} \frac{\partial J}{\partial w} &= \frac{\partial J}{\partial z} \cdot \frac{\partial z}{\partial w} \\ \frac{\partial J}{\partial w} &= \frac{\partial J}{\partial z} \cdot \frac{\partial z}{\partial w} \\ \frac{\partial J}{\partial w} &= \frac{\partial J}{\partial z} \cdot \frac{\partial z}{\partial w} \\ \frac{\partial J}{\partial w} &= \frac{\partial J}{\partial z} \cdot \frac{\partial z}{\partial w} \\ \frac{\partial J}{\partial w} &= \frac{\partial J}{\partial z} \cdot \frac{\partial z}{\partial w} \end{aligned}$$

• 5 → Implementar o update rule

↳ Atualizamos os valores de "w" e "b", de modo a usar o aprendizado

↳ A regra base é:  $w = \theta - \alpha \frac{\partial J}{\partial \theta}$  sendo  $\theta$  = parâmetro

$\alpha$  = learning rate

Ex:  $w_1 \rightarrow w_1 = w_1 - \alpha \cdot \frac{\partial J}{\partial w_1}$  vêm dos passos anteriores

• 6 - backward

→ Backprop

(Queremos derivada voltando de trás para definir a variação de parâmetros  
 (comparar com  $\frac{\partial J(a, y)}{\partial a} \rightarrow \frac{\partial J(a, y)}{\partial a} = \frac{-y + \hat{y}}{1 - \hat{y}}$ )

"voltando", temos que  $\frac{\partial J}{\partial z} = \frac{\partial J(a, y)}{\partial z} = a - y$

Depois computa "w" e "b"  $\rightarrow \frac{\partial J}{\partial w_1} = x_1 \cdot \frac{\partial J}{\partial z}$   
 $\frac{\partial J}{\partial w_2} = x_2 \cdot \frac{\partial J}{\partial z}$   
 $\frac{\partial J}{\partial b} = \frac{\partial J}{\partial z}$

→ O ponto de atenção é: Se pegarmos um valor  $v$  e mudarmos, o que não vai acontecer com os subsequentes?

mas só contam "v", o que acontece com "j"?

$\frac{\partial J}{\partial v} \rightarrow$  Se mudarmos "v", o que acontece com "j"?

$J = 3v$

→ Back propagation entra si, porque você vai voltando no grupo para os outros mutantes e output

$J = 3v$

Se eu mudar  $v_1$ , vai mudar  $v$ , que vai mudar  $J$

$v \rightarrow v \rightarrow f \rightarrow \text{chain rule das derivadas}$

→ v foco é adiar a derivada da final output

→ O que é  $J$ ? Cost function

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m \log(\sigma(w^T x^{(i)} + b))$$

• Logistic Regression

$z = w^T x + b$   
 $\hat{y} = \Phi(z) = \sigma(z) = \frac{1}{1 + e^{-z}}$   
 $\hat{y} = \Phi(z) = \sigma(z)$  Sigmoid function  
 $\sigma(z) = \frac{1}{1 + e^{-z}}$

nos continuamos tendo que calcular o dot product e offset

Sigmoid tem como equação

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

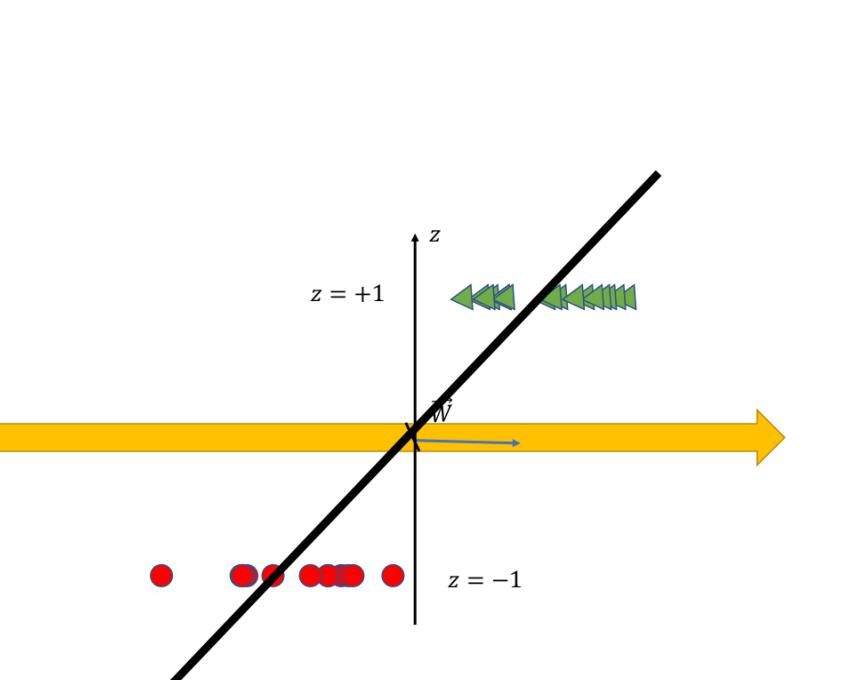
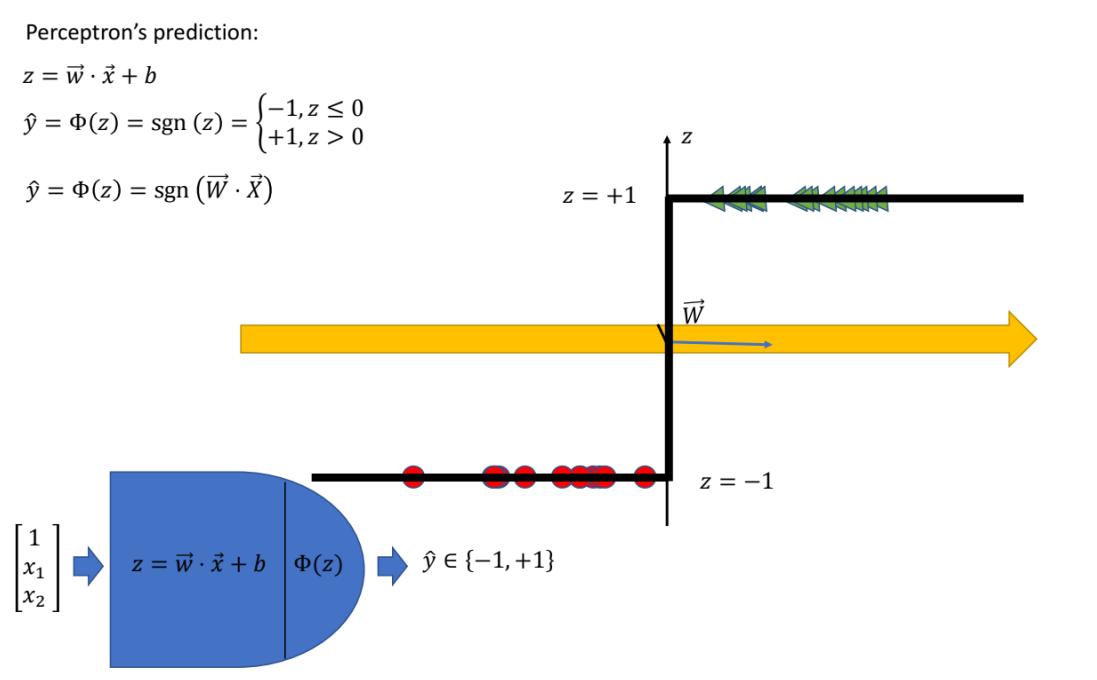
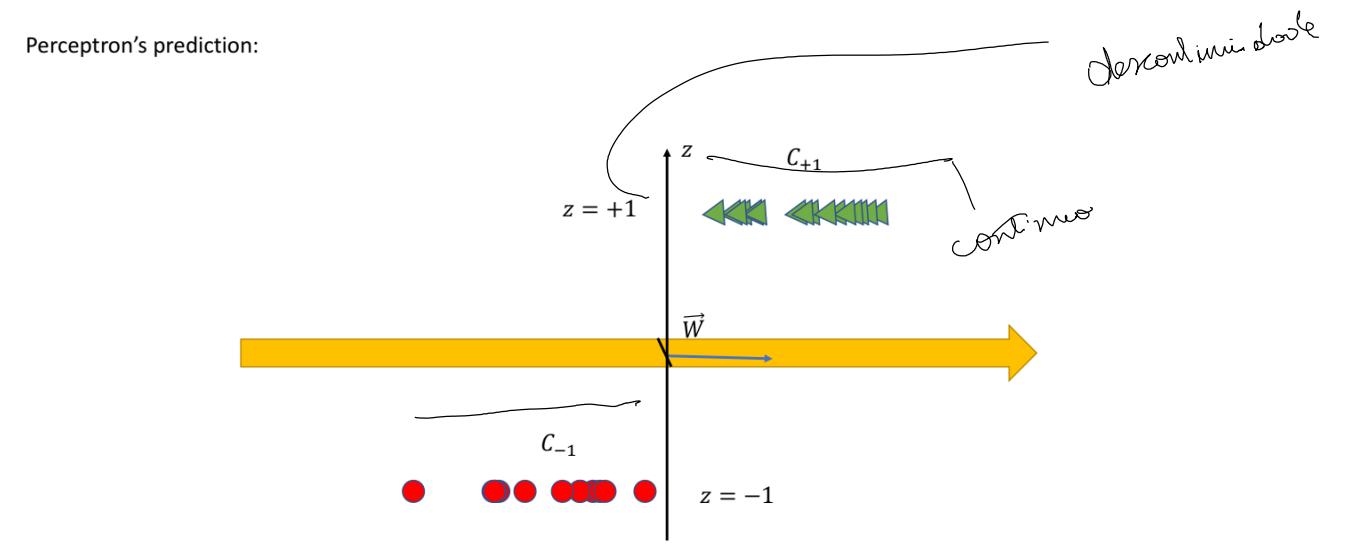
↳ função de probabilidade

• Até o momento da "primeira" backpropagation

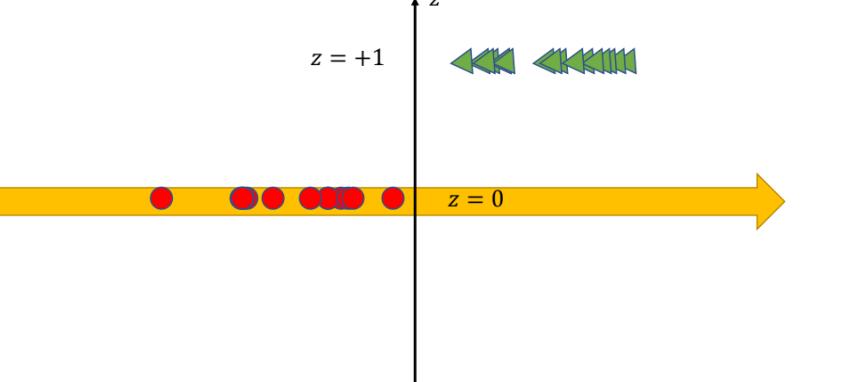
$$w := w - \alpha \frac{\partial J(w, b)}{\partial w}$$

"dwi"

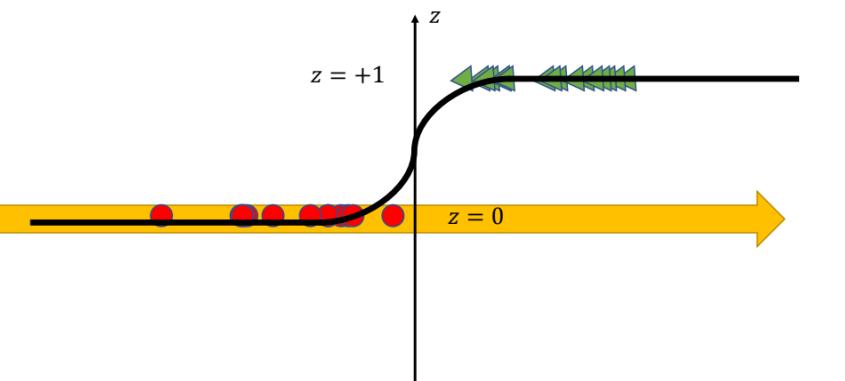
$$b := b - \alpha \frac{\partial J(w, b)}{\partial b}$$



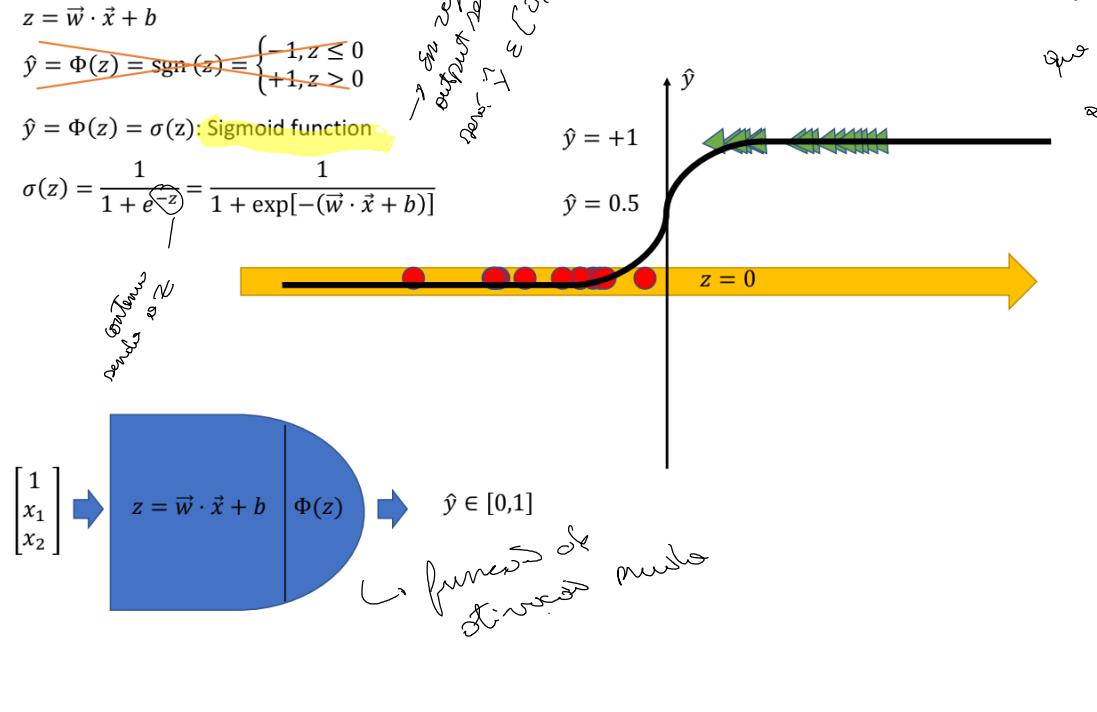
Probabilistic view



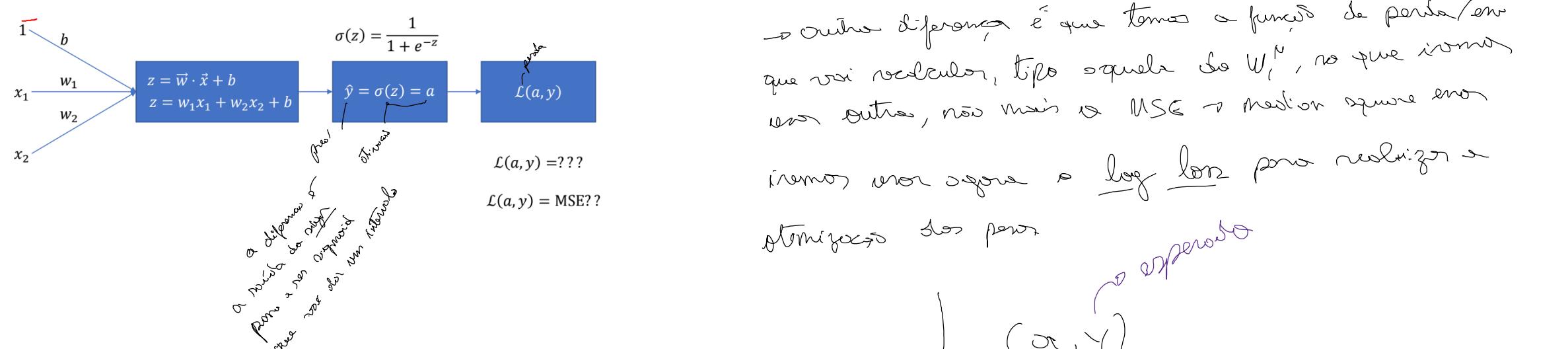
Logistic Regression!



Logistic Regression!



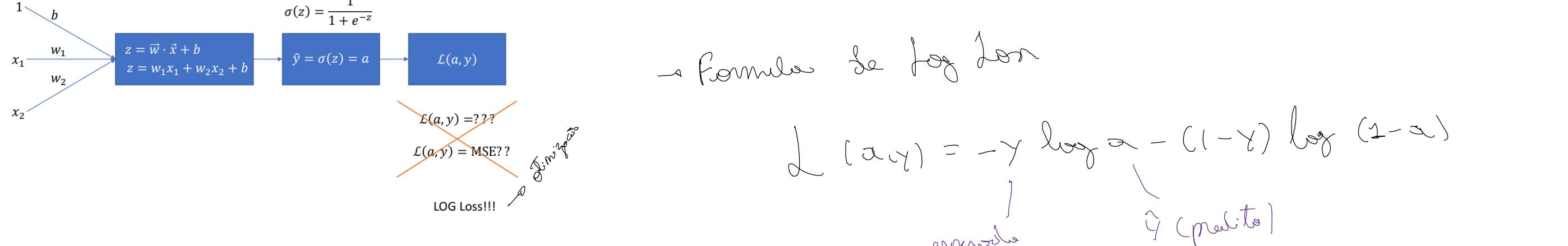
## Training



o Outra diferença é que temos a função de perda/métrica que vai calcular, tipo quando da  $W^T$ , se queremos usar outra, nós mudamos para median square error ou vamos usar alguma a log loss pra regularizar os pesos

$\int (a, y)$

$\hat{y}$ , resultados da regressão



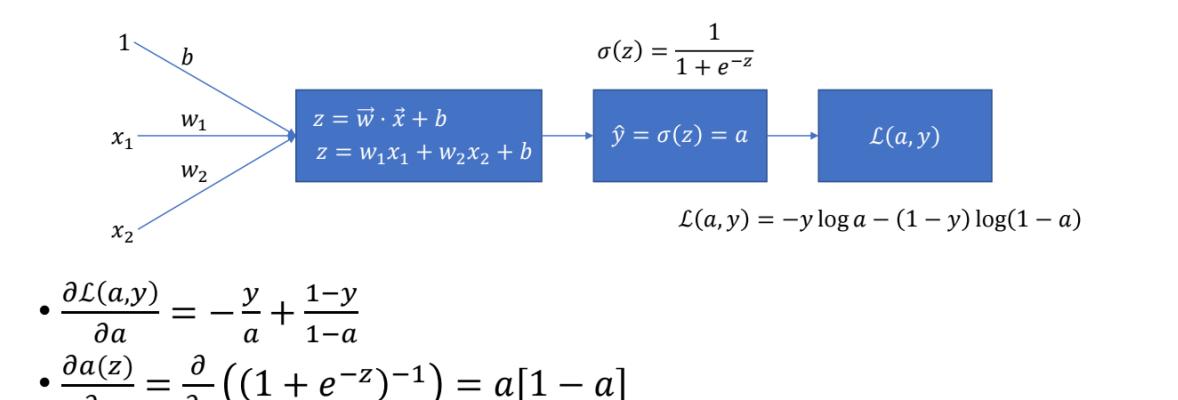
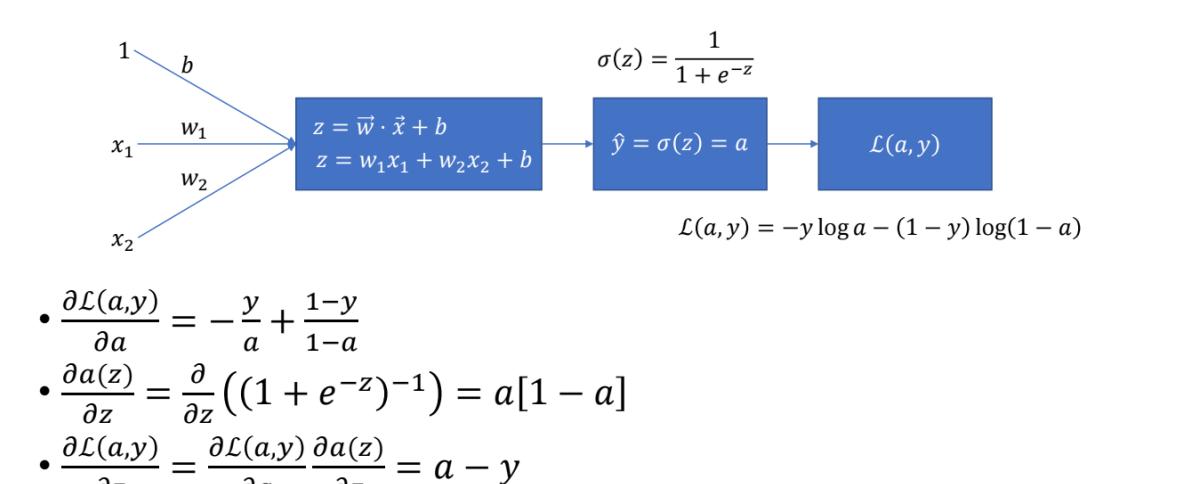
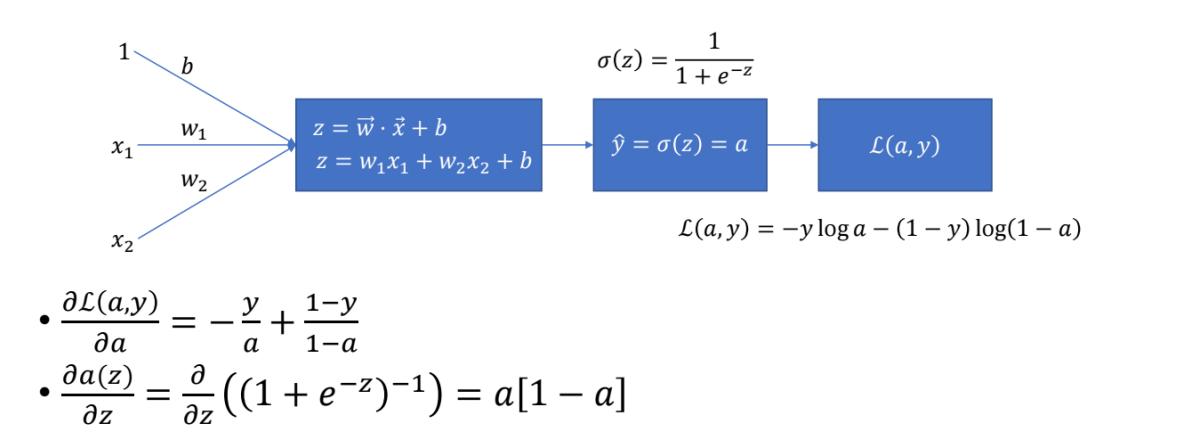
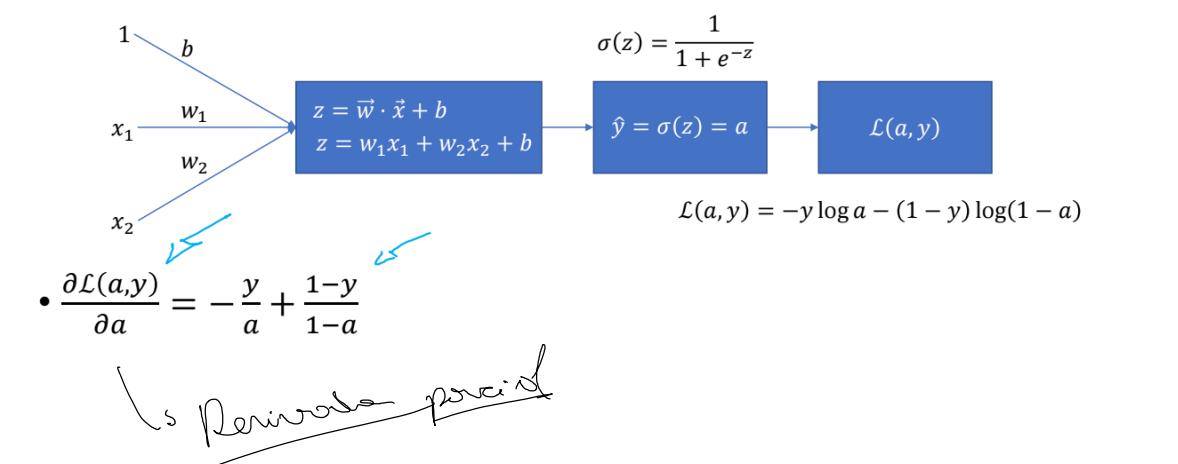
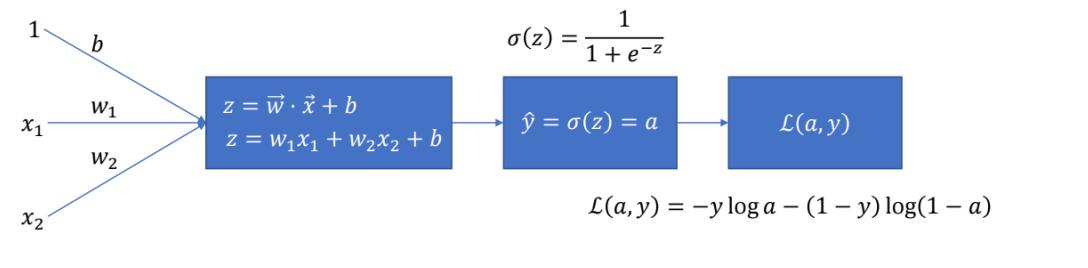
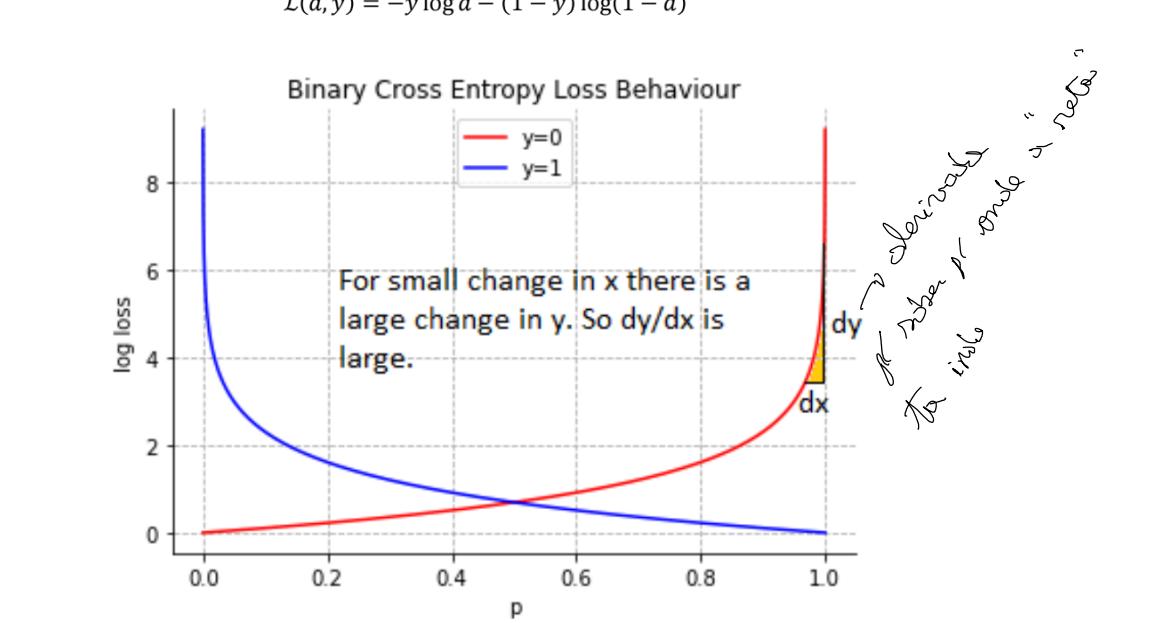
as fórmulas da log loss

$$\int (a, y) = -y \log a - (1-y) \log (1-a)$$

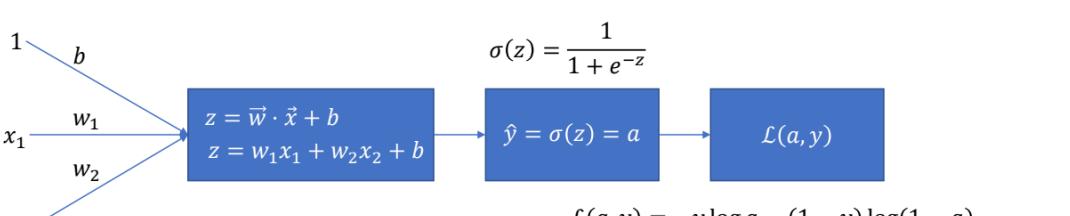
explosão

$\hat{y}$  (predito)

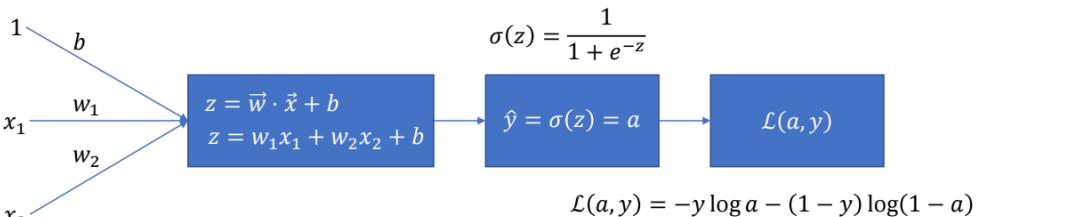
Log loss



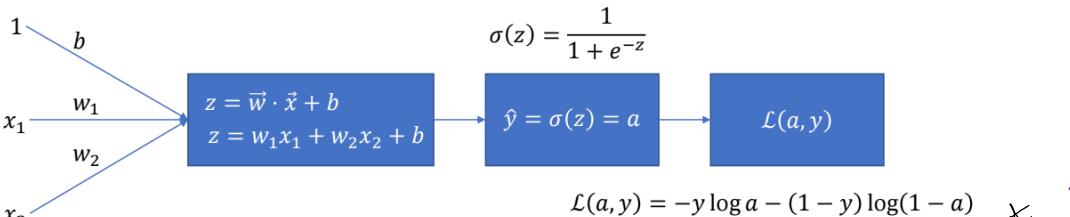
- $\frac{\partial \mathcal{L}(a,y)}{\partial z} = \frac{\partial \mathcal{L}(a,y)}{\partial a} \frac{\partial a}{\partial z} = a - y$
- $\frac{\partial \mathcal{L}}{\partial w_1} = \frac{\partial \mathcal{L}}{\partial z} \frac{\partial z}{\partial w_1} = (a - y)x_1$



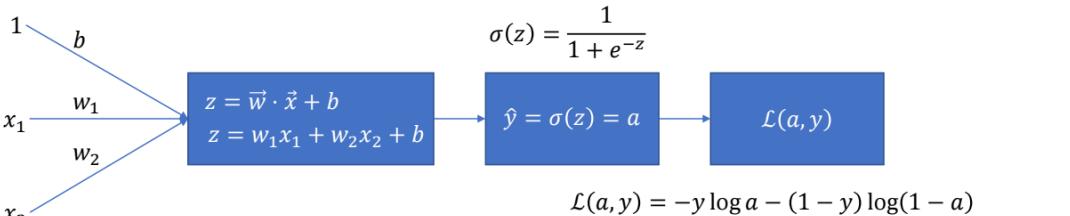
- $\frac{\partial \mathcal{L}(a,y)}{\partial a} = -\frac{y}{a} + \frac{1-y}{1-a}$
- $\frac{\partial a}{\partial z} = \frac{\partial}{\partial z} ((1+e^{-z})^{-1}) = a[1-a]$
- $\frac{\partial \mathcal{L}(a,y)}{\partial z} = \frac{\partial \mathcal{L}(a,y)}{\partial a} \frac{\partial a}{\partial z} = a - y$
- $\frac{\partial \mathcal{L}}{\partial w_1} = \frac{\partial \mathcal{L}}{\partial z} \frac{\partial z}{\partial w_1} = (a - y)x_1$
- $\frac{\partial \mathcal{L}}{\partial w_2} = \frac{\partial \mathcal{L}}{\partial z} \frac{\partial z}{\partial w_2} = (a - y)x_2$



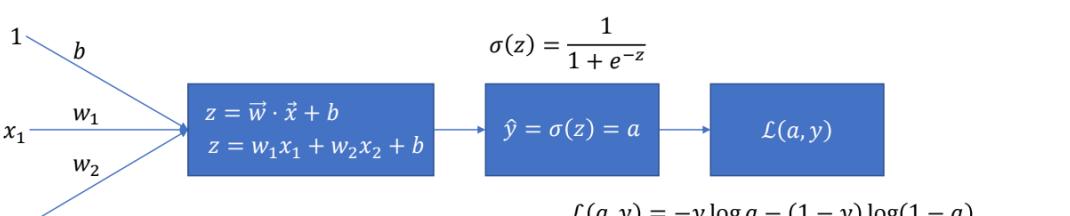
- $\frac{\partial \mathcal{L}(a,y)}{\partial a} = -\frac{y}{a} + \frac{1-y}{1-a}$
- $\frac{\partial a}{\partial z} = \frac{\partial}{\partial z} ((1+e^{-z})^{-1}) = a[1-a]$
- $\frac{\partial \mathcal{L}(a,y)}{\partial z} = \frac{\partial \mathcal{L}(a,y)}{\partial a} \frac{\partial a}{\partial z} = a - y$
- $\frac{\partial \mathcal{L}}{\partial w_1} = \frac{\partial \mathcal{L}}{\partial z} \frac{\partial z}{\partial w_1} = (a - y)x_1$
- $\frac{\partial \mathcal{L}}{\partial w_2} = \frac{\partial \mathcal{L}}{\partial z} \frac{\partial z}{\partial w_2} = (a - y)x_2$
- $\frac{\partial \mathcal{L}}{\partial b} = \frac{\partial \mathcal{L}}{\partial z} \frac{\partial z}{\partial b} = (a - y)$



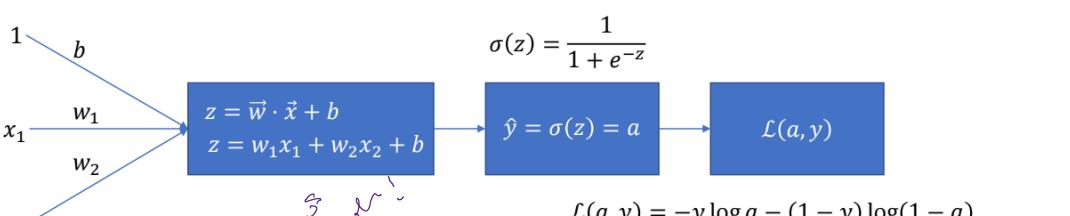
- $\frac{\partial \mathcal{L}(a,y)}{\partial a} = -\frac{y}{a} + \frac{1-y}{1-a}$
- $\frac{\partial a}{\partial z} = \frac{\partial}{\partial z} ((1+e^{-z})^{-1}) = a[1-a]$
- $\frac{\partial \mathcal{L}(a,y)}{\partial z} = \frac{\partial \mathcal{L}(a,y)}{\partial a} \frac{\partial a}{\partial z} = a - y$
- $\frac{\partial \mathcal{L}}{\partial w_1} = \frac{\partial \mathcal{L}}{\partial z} \frac{\partial z}{\partial w_1} = (a - y)x_1$
- $\frac{\partial \mathcal{L}}{\partial w_2} = \frac{\partial \mathcal{L}}{\partial z} \frac{\partial z}{\partial w_2} = (a - y)x_2$
- $\frac{\partial \mathcal{L}}{\partial b} = \frac{\partial \mathcal{L}}{\partial z} \frac{\partial z}{\partial b} = (a - y)$



- $\frac{\partial \mathcal{L}(a,y)}{\partial a} = -\frac{y}{a} + \frac{1-y}{1-a}$
- $\frac{\partial a}{\partial z} = \frac{\partial}{\partial z} ((1+e^{-z})^{-1}) = a[1-a]$
- $\frac{\partial \mathcal{L}(a,y)}{\partial z} = \frac{\partial \mathcal{L}(a,y)}{\partial a} \frac{\partial a}{\partial z} = a - y$
- $\frac{\partial \mathcal{L}}{\partial w_1} = \frac{\partial \mathcal{L}}{\partial z} \frac{\partial z}{\partial w_1} = (a - y)x_1$
- $\frac{\partial \mathcal{L}}{\partial w_2} = \frac{\partial \mathcal{L}}{\partial z} \frac{\partial z}{\partial w_2} = (a - y)x_2$
- $\frac{\partial \mathcal{L}}{\partial b} = \frac{\partial \mathcal{L}}{\partial z} \frac{\partial z}{\partial b} = (a - y)$

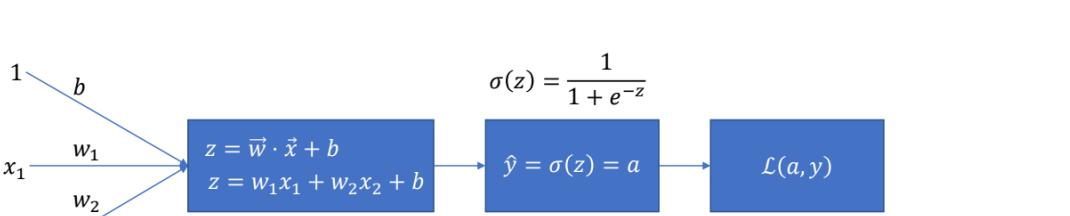


- $\frac{\partial \mathcal{L}(a,y)}{\partial a} = -\frac{y}{a} + \frac{1-y}{1-a}$
- $\frac{\partial a}{\partial z} = \frac{\partial}{\partial z} ((1+e^{-z})^{-1}) = a[1-a]$
- $\frac{\partial \mathcal{L}(a,y)}{\partial z} = \frac{\partial \mathcal{L}(a,y)}{\partial a} \frac{\partial a}{\partial z} = a - y$
- $\frac{\partial \mathcal{L}}{\partial w_1} = \frac{\partial \mathcal{L}}{\partial z} \frac{\partial z}{\partial w_1} = (a - y)x_1$
- $\frac{\partial \mathcal{L}}{\partial w_2} = \frac{\partial \mathcal{L}}{\partial z} \frac{\partial z}{\partial w_2} = (a - y)x_2$
- $\frac{\partial \mathcal{L}}{\partial b} = \frac{\partial \mathcal{L}}{\partial z} \frac{\partial z}{\partial b} = (a - y)$

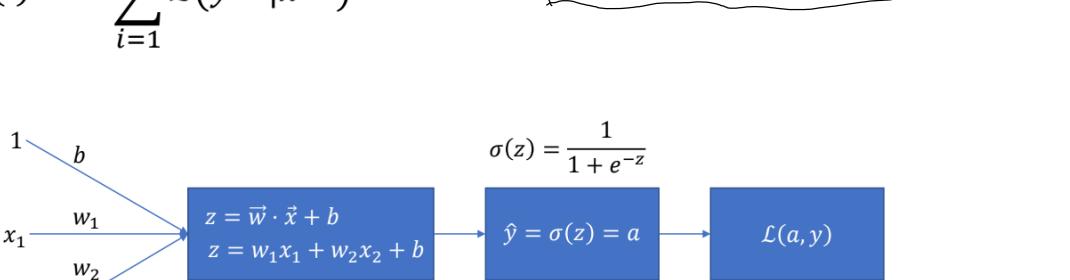


- Cost of  $m$  samples:  
 $P(\cdot) = \prod_{i=1}^m P(y^{(i)}|\vec{x}^{(i)})$  (Maximum Likelihood Estimation)

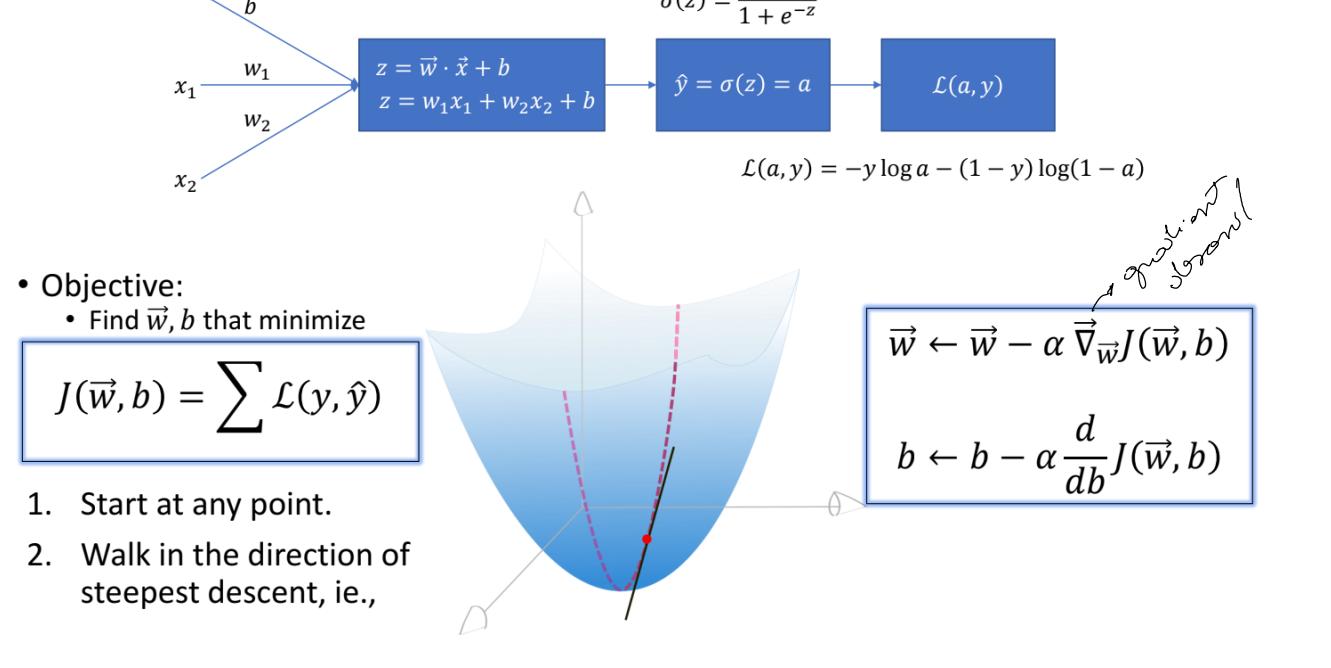
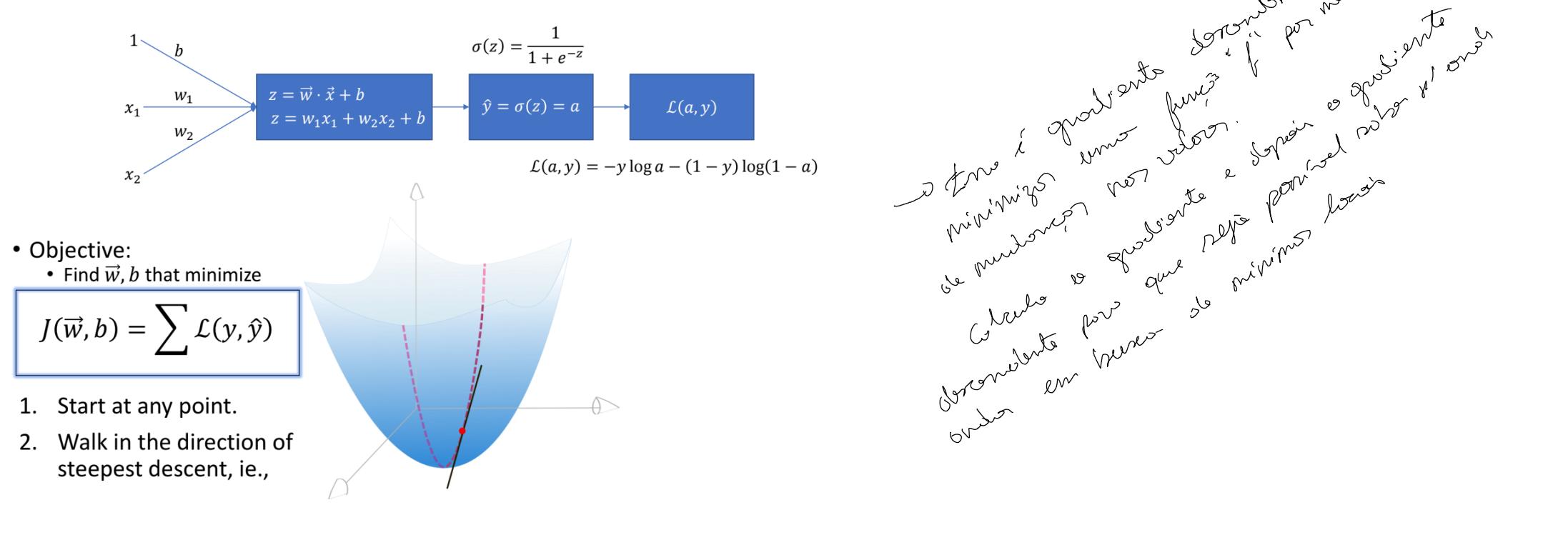
The MLE estimates the set of parameters  $\vec{\theta} = (\vec{w}, b)$  so that they are as close as possible to the real set of parameters  $\vec{\theta}$  that generated such data samples by minimizing the loss  $L(y, \hat{y})$ .



- Cost of  $m$  samples:  
 $\log P(\cdot) = \log \prod_{i=1}^m P(y^{(i)}|\vec{x}^{(i)})$
- $\log P(\cdot) = \sum_{i=1}^m \log P(y^{(i)}|\vec{x}^{(i)})$
- $\log P(\cdot) = -\sum_{i=1}^m \mathcal{L}(y^{(i)}|\vec{x}^{(i)})$



- Objective:  
 $J(\vec{w}, b) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(y^{(i)}, \hat{y}^{(i)})$



## Notation

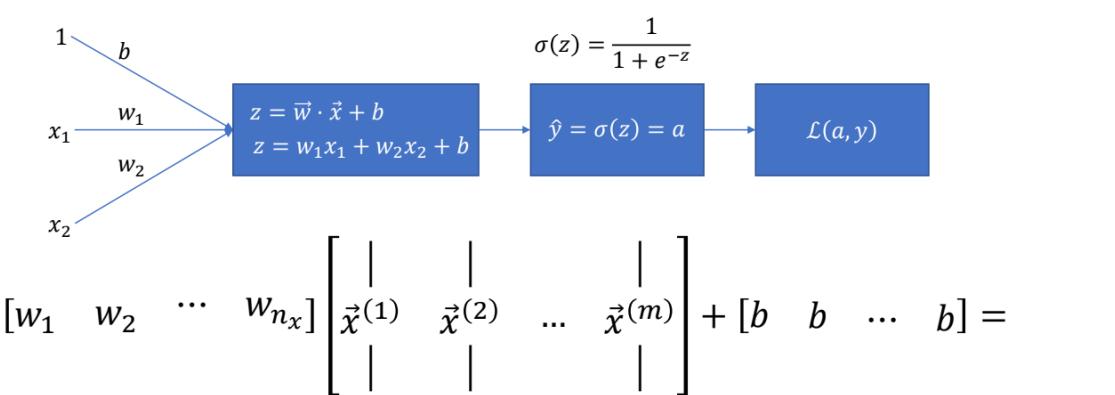
- Single Training Sample:  $(\vec{x}, y)$ , where:
  - $\vec{x} \in \mathbb{R}^{n_x}$ .
  - $y \in \{0,1\}$ .

- Single Training Sample:  $(\vec{x}, y)$ , where:
  - $\vec{x} \in \mathbb{R}^{n_x}$ .
  - $y \in \{0,1\}$ .
- $m$  training samples:  $\{(\vec{x}^{(1)}, y^{(1)}), (\vec{x}^{(2)}, y^{(2)}), \dots, (\vec{x}^{(m)}, y^{(m)})\} = \{\vec{x}^{(i)}, y^{(i)}\}_{i=1}^m$ , where:
  - $m_{train}$  = Number of training samples.
  - $m_{test}$  = Number of test samples.

- Single Training Sample:  $(\vec{x}, y)$ , where:
  - $\vec{x} \in \mathbb{R}^{n_x}$ .
  - $y \in \{0,1\}$ .
- $m$  training samples:  $\{(\vec{x}^{(1)}, y^{(1)}), (\vec{x}^{(2)}, y^{(2)}), \dots, (\vec{x}^{(m)}, y^{(m)})\} = \{\vec{x}^{(i)}, y^{(i)}\}_{i=1}^m$ , where:
  - $m_{train}$  = Number of training samples.
  - $m_{test}$  = Number of test samples.
- $X = \begin{bmatrix} | & | & | \\ \vec{x}^{(1)} & \vec{x}^{(2)} & \dots & \vec{x}^{(m)} \\ | & | & | \end{bmatrix}_{n_x \times m}$ ,  $Y = \begin{bmatrix} y^{(1)} & y^{(2)} & \dots & y^{(m)} \end{bmatrix}$
- $X \in \mathbb{R}^{(n_x \times m)}$ , i.e.,  $X.shape = (n_x, m)$ .
- $Y \in \mathbb{R}^{(1 \times m)}$ , i.e.,  $Y.shape = (1, m)$ .

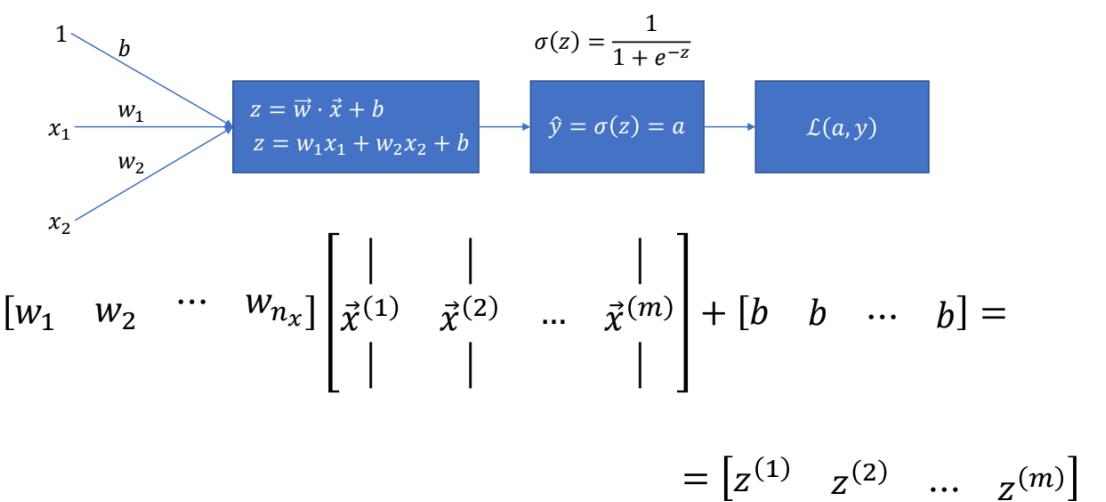
## Gradient Descent on Logistic Regression

### Forward Propagation for $m$ examples

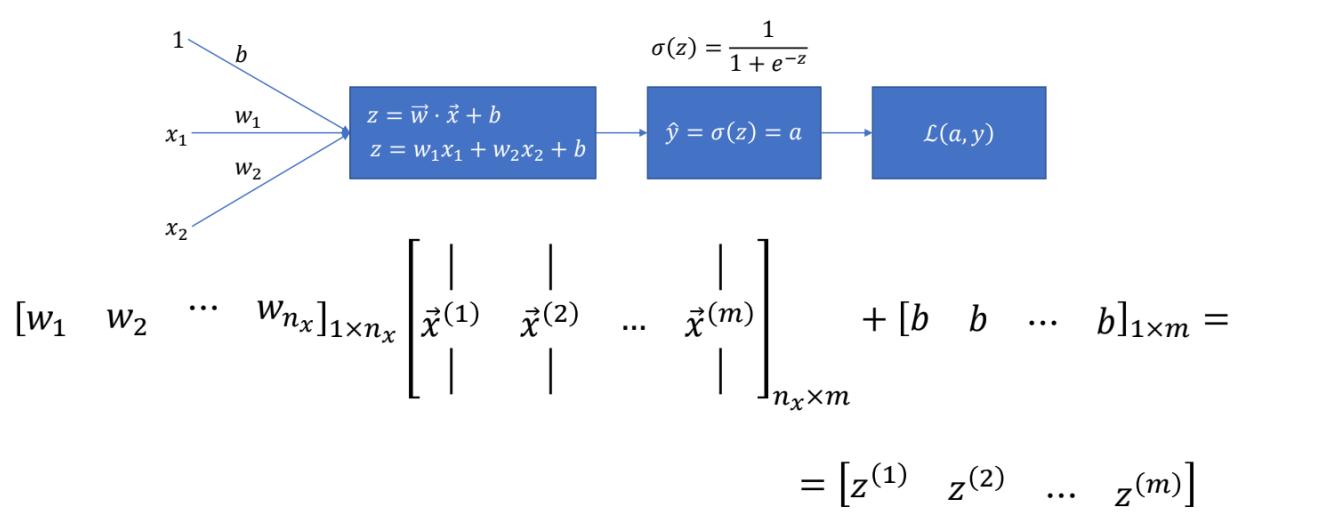


### Forward Propagation for $m$ examples

## Forward Propagation for 1 example



Forward Propagation for  $m$  examples



Forward Propagation for  $m$  examples

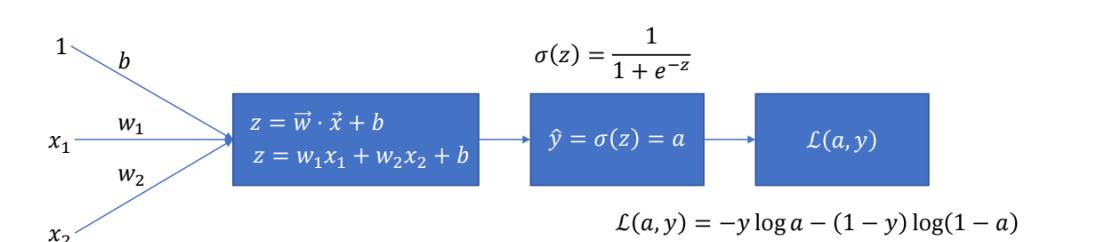
$$[w_1 \ w_2 \ \dots \ w_{n_x}]_{1 \times n_x} \begin{bmatrix} | & | & | \\ \hat{x}^{(1)} & \hat{x}^{(2)} & \dots & \hat{x}^{(m)} \\ | & | & | \end{bmatrix}_{n_x \times m} + [b \ b \ \dots \ b]_{1 \times m} = \\ = Z = [z^{(1)} \ z^{(2)} \ \dots \ z^{(m)}]$$

Forward Propagation for  $m$  examples

$$[w_1 \ w_2 \ \dots \ w_{n_x}]_{1 \times n_x} \begin{bmatrix} | & | & | \\ \hat{x}^{(1)} & \hat{x}^{(2)} & \dots & \hat{x}^{(m)} \\ | & | & | \end{bmatrix}_{n_x \times m} + [b \ b \ \dots \ b]_{1 \times m} = \\ = Z = [z^{(1)} \ z^{(2)} \ \dots \ z^{(m)}] \\ \hat{y}(Z) = [a^{(1)} \ a^{(2)} \ \dots \ a^{(m)}] \\ A = [y^{(1)} \ y^{(2)} \ \dots \ y^{(m)}]$$

Forward Propagation for  $m$  examples

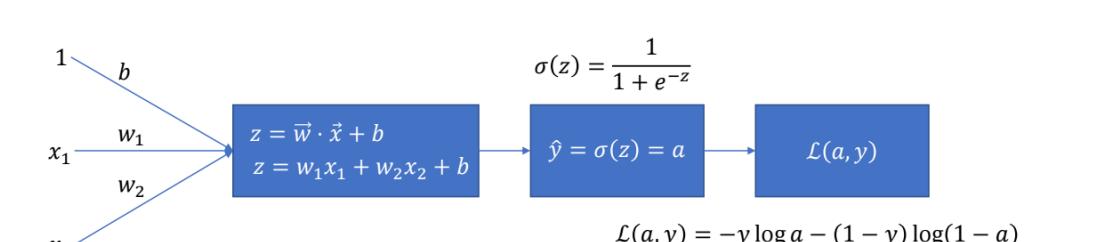
$$[w_1 \ w_2 \ \dots \ w_{n_x}]_{1 \times n_x} \begin{bmatrix} | & | & | \\ \hat{x}^{(1)} & \hat{x}^{(2)} & \dots & \hat{x}^{(m)} \\ | & | & | \end{bmatrix}_{n_x \times m} + [b \ b \ \dots \ b]_{1 \times m} = \\ = Z = [z^{(1)} \ z^{(2)} \ \dots \ z^{(m)}] \\ \frac{\partial \mathcal{L}(a, y)}{\partial z} = \frac{\partial \mathcal{L}(a, y)}{\partial a} \frac{\partial a}{\partial z} = a - y \\ \left. \begin{aligned} A\sigma(Z) &= [a^{(1)} \ a^{(2)} \ \dots \ a^{(m)}] \\ A &= [y^{(1)} \ y^{(2)} \ \dots \ y^{(m)}] \end{aligned} \right\} \\ dZ = A - Y = [a^{(1)} - y^{(1)} \ a^{(2)} - y^{(2)} \ \dots \ a^{(m)} - y^{(m)}]$$



Recall the gradient descent for 1 sample:

- $\frac{\partial \mathcal{L}}{\partial w_1} = \frac{\partial \mathcal{L}}{\partial a} \frac{\partial a}{\partial w_1} = (\mathbf{a} - \mathbf{y})x_1$
- $\frac{\partial \mathcal{L}}{\partial w_2} = \frac{\partial \mathcal{L}}{\partial a} \frac{\partial a}{\partial w_2} = (\mathbf{a} - \mathbf{y})x_2$
- $\frac{\partial \mathcal{L}}{\partial b} = \frac{\partial \mathcal{L}}{\partial a} \frac{\partial a}{\partial b} = (\mathbf{a} - \mathbf{y})$

This is the loss related to 1 sample!



Gradient descent for  $m$  samples:

$$\mathcal{J}(\vec{w}, b) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(y^{(i)}, \hat{y}^{(i)}) \\ \frac{\partial}{\partial w_k} \mathcal{J}(\vec{w}, b) = \frac{1}{m} \sum_{i=1}^m \frac{\partial}{\partial w_k} \mathcal{L}(y^{(i)}, \hat{y}^{(i)})$$

$$\left[ w_k \leftarrow w_k - \eta \frac{\partial}{\partial w_k} \mathcal{J}(\vec{w}, b) \right]_{k=1}^{n_x} \\ b \leftarrow b - \eta \frac{\partial}{\partial b} \mathcal{J}(\vec{w}, b)$$

Thank you