



PECE Programa de
Educação Continuada

Escola Politécnica da USP

Introdução a Redes Neurais

Marlon Sproesser Mathias

Aula 2 – Parte 1

Perceptrons

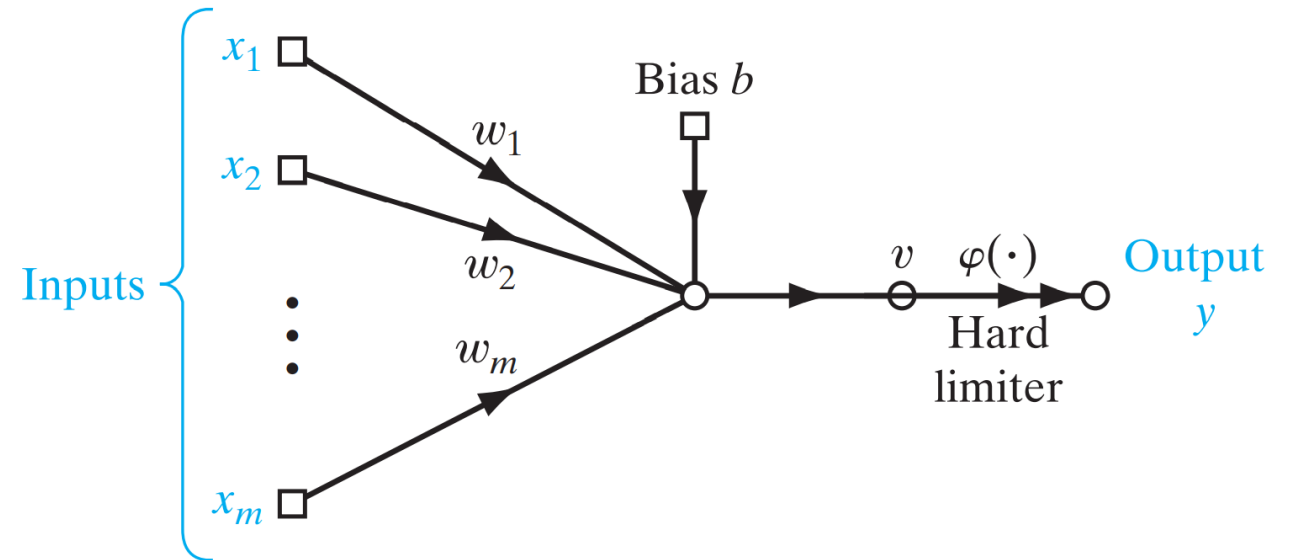
Classificação
binária

Gradientes
descendentes

Trabalho 1

Perceptrons

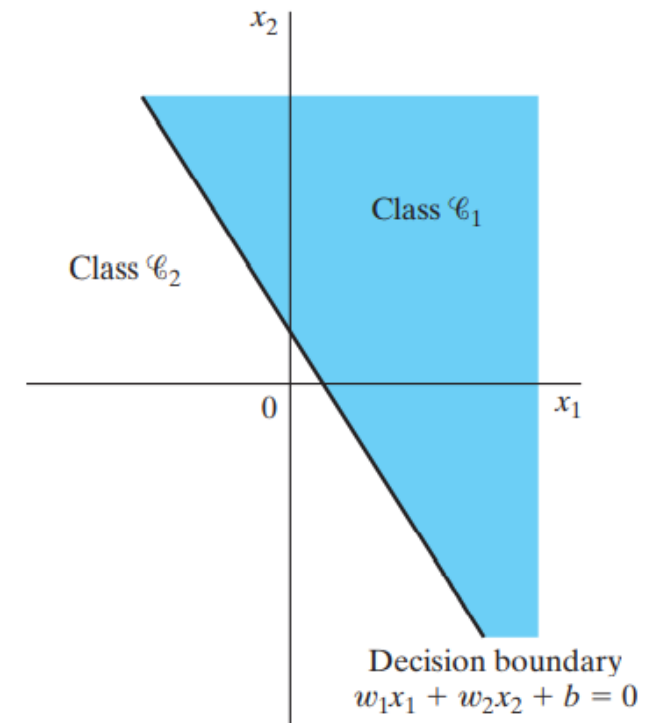
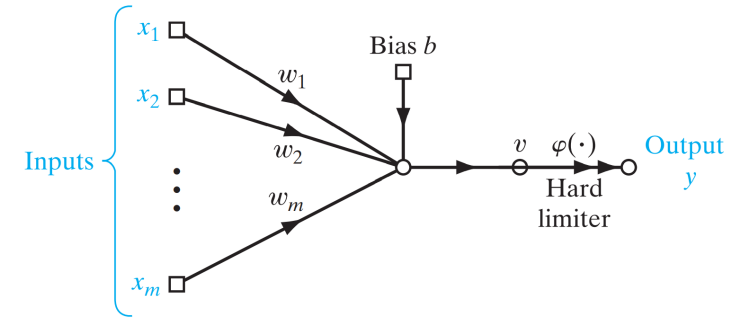
- Unidade de classificação
- Perceptron de Rosenblatt (1958)
- Cria uma hipersuperfície dentro do espaço das entradas



Haykin, S. Neural Networks and Learning Machines - 2009

Perceptrons

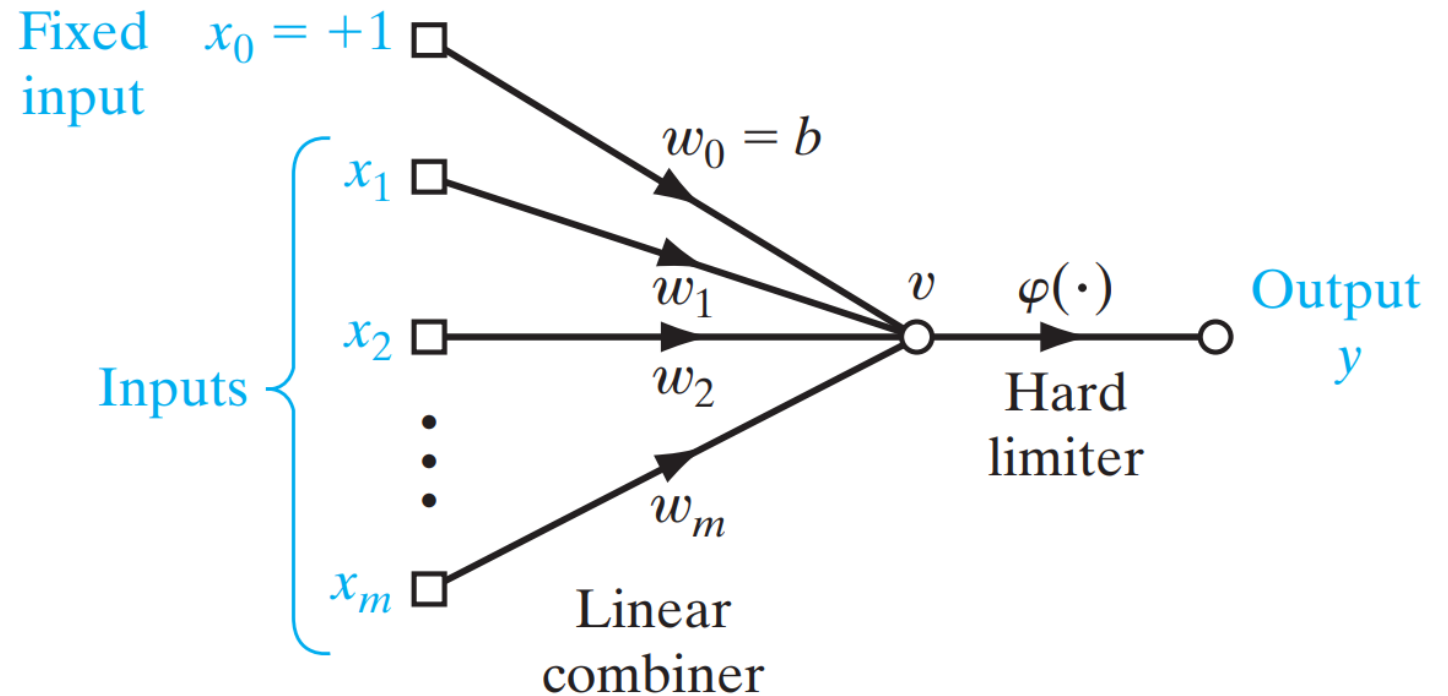
- Unidade de classificação
- Perceptron de Rosenblatt (1958)
- Cria uma hipersuperfície dentro do espaço das entradas
- Na prática, separa as saídas em duas classes



Haykin, S. Neural Networks and Learning Machines - 2009

Perceptrons

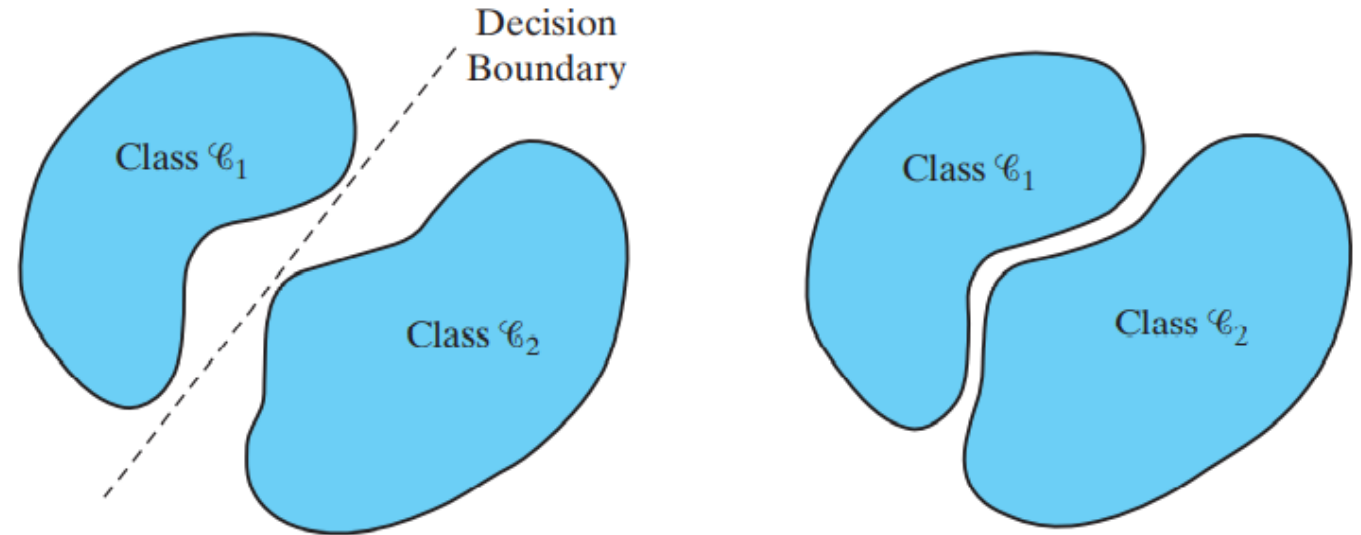
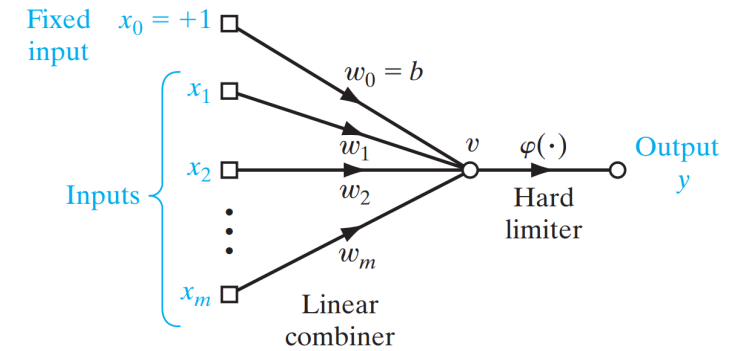
- Quais variáveis escolhemos?



Haykin, S. Neural Networks and Learning Machines - 2009

Perceptrons

- Quais variáveis escolhemos?
- Limitado a separações mais simples → Apenas lineares

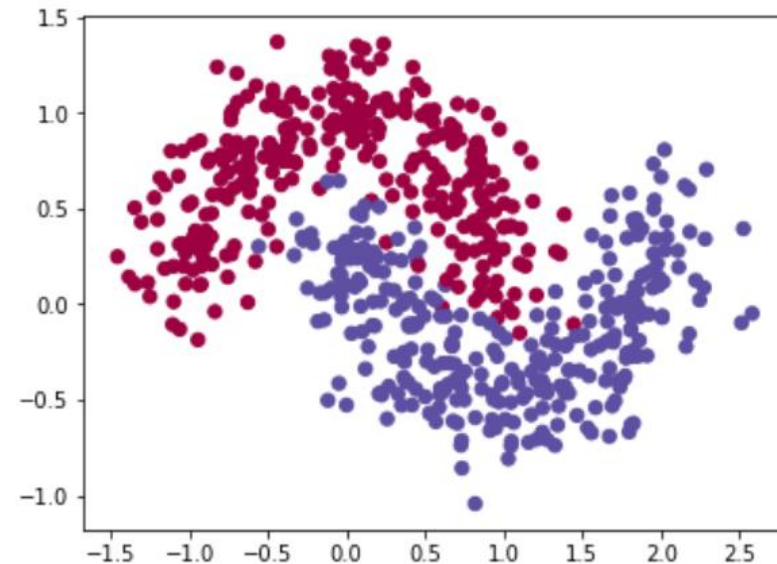


Haykin, S. Neural Networks and Learning Machines - 2009

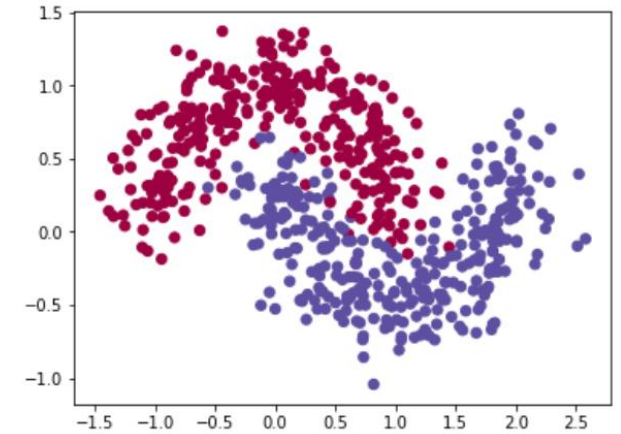
Classificação binária

- Separar dados em dois grupos
- Treinamento supervisionado → Já temos alguns exemplos

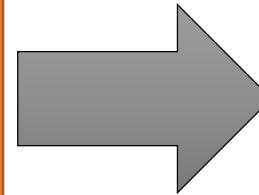
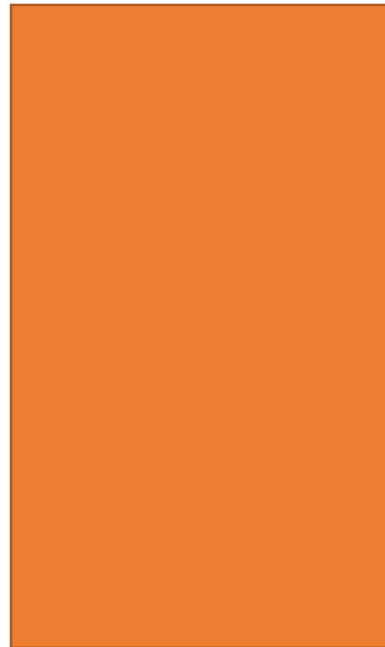
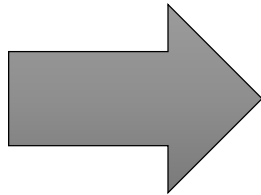
- Dados pertencem a dois grupos distintos
- Cada ponto tem coordenadas (x_1, x_2)
- Como separamos os grupos?
- Como treinamos uma RNA para separar os grupos?



Classificação binária



Coordenadas
(x_1, x_2)



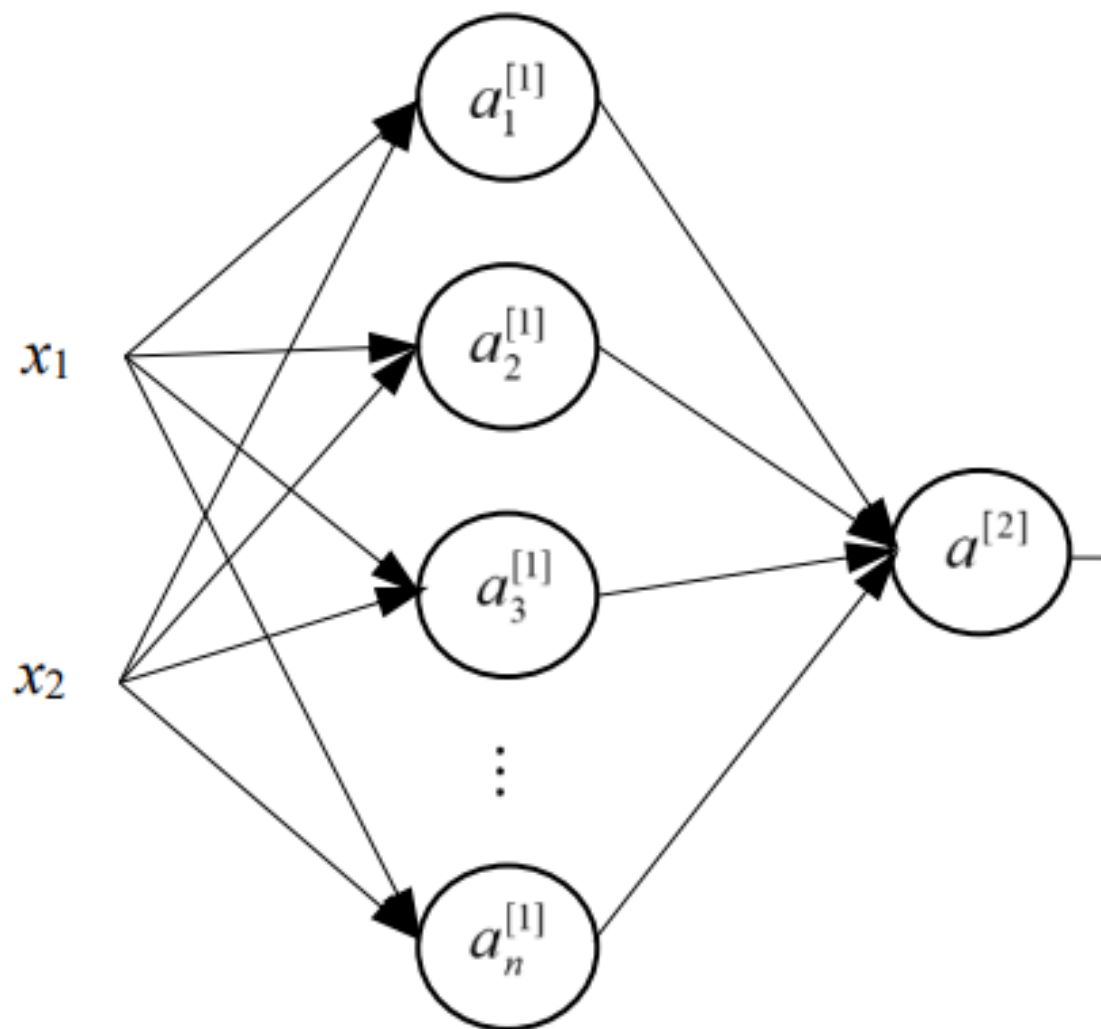
Classe
(y)

0 → Azul

1 → Vermelha

Estrutura da RNA

- Entrada \rightarrow Vetor $\mathbf{x}^{(i)} \in \mathbb{R}^2$
- Saída $\rightarrow \mathbf{y}^{(i)} \in \mathbb{R}$
- Camada intermediária
 - $n^{[1]}$ neurônios
 - Ativação de tangente hiperbólica
- Camada de saída
 - $n^{[2]} = 1$ neurônio
 - Ativação de sigmoide

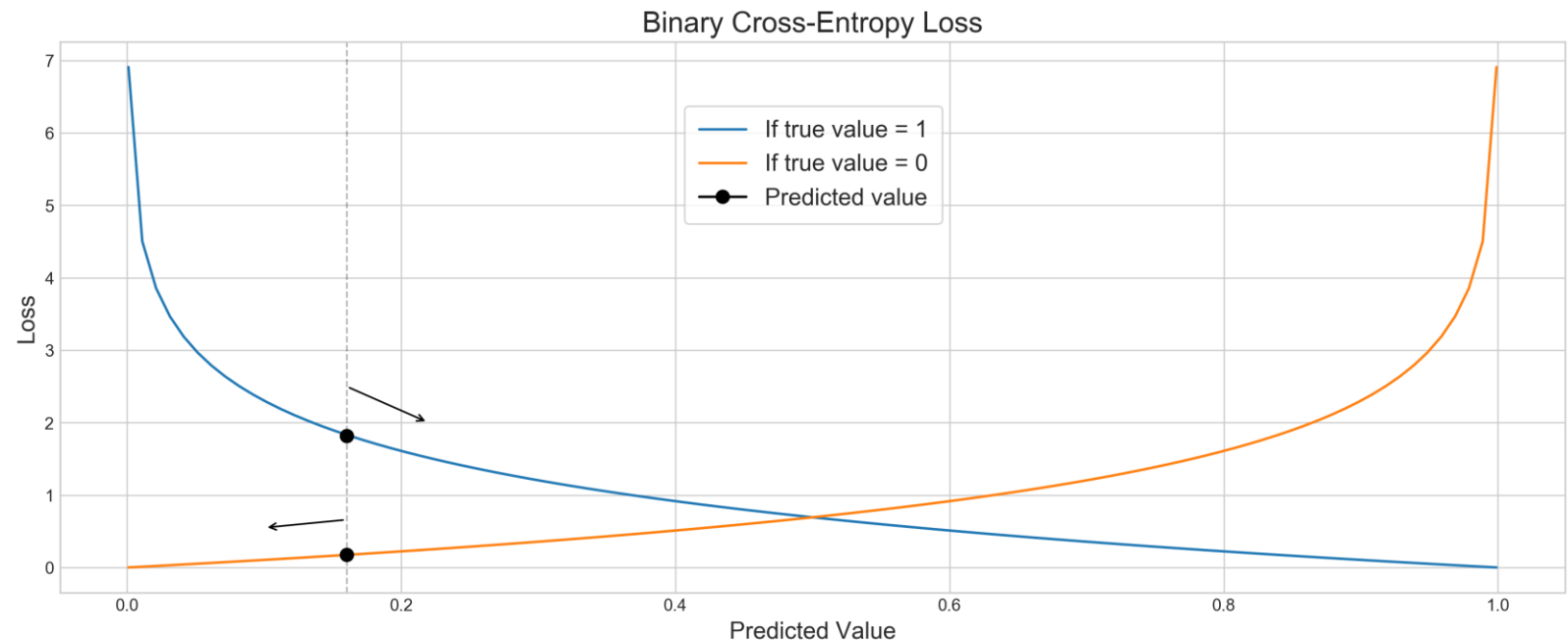


Eduardo Lobo Lustosa Cabral

Função de custo

- Como comparar a saída da RNA com a “verdade”
- $\hat{y}^{(i)}$ é a previsão da rede no ponto i
- $y^{(i)}$ é o valor de referência no ponto i
- Função de erro logística

$$L(\hat{y}^{(i)}, y^{(i)}) = -y^{(i)} \log \hat{y}^{(i)} - (1 - y^{(i)}) \log(1 - \hat{y}^{(i)})$$

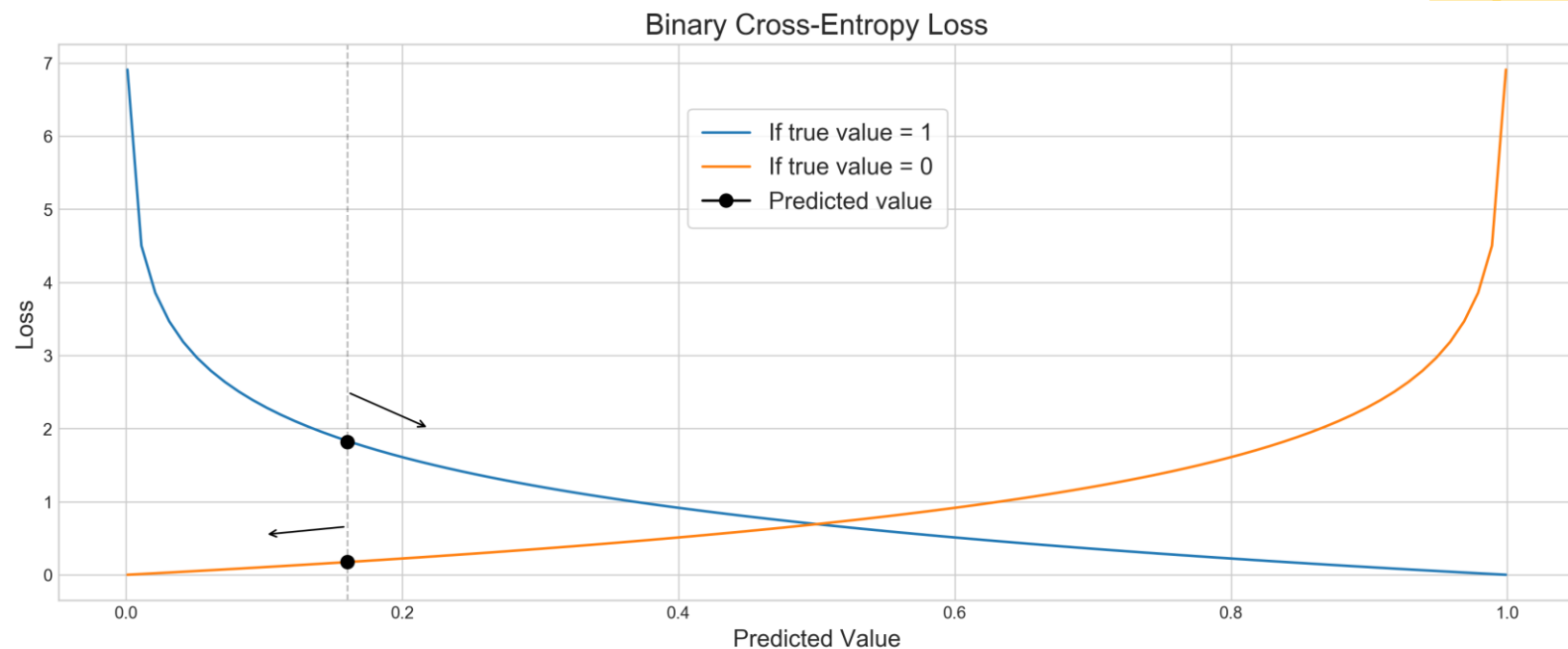


<https://towardsdatascience.com/logistic-regression-from-scratch-69db4f587e17>

Função de custo

- Como comparar a saída da RNA com a “verdade”
- $\hat{y}^{(i)}$ é a previsão da rede no ponto i
- $y^{(i)}$ é o valor de referência no ponto i
- Função de erro logística

$$L(\hat{y}^{(i)}, y^{(i)}) = -y^{(i)} \log \hat{y}^{(i)} - (1 - y^{(i)}) \log(1 - \hat{y}^{(i)})$$



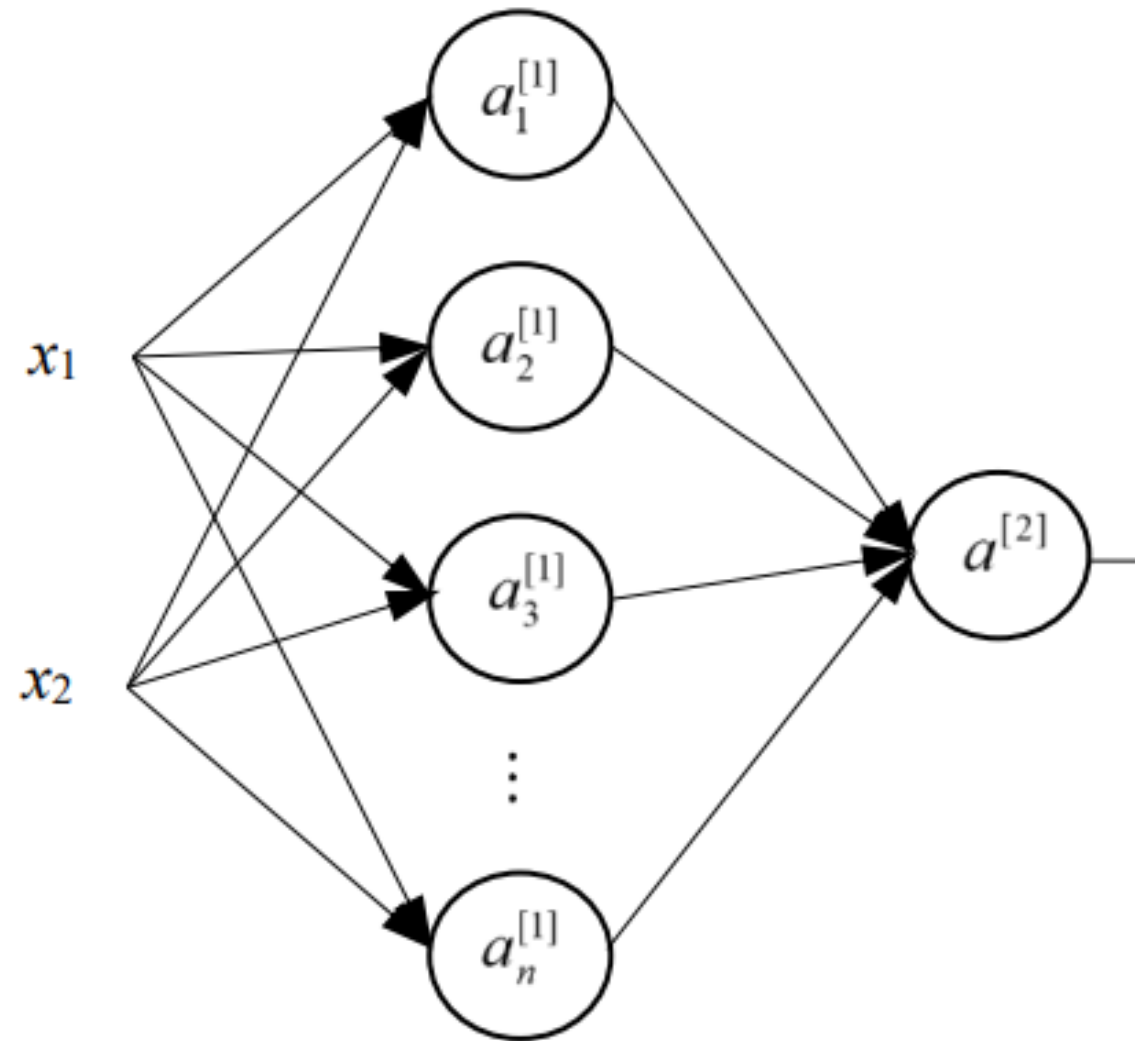
<https://towardsdatascience.com/logistic-regression-from-scratch-69db4f587e17>

Aplicando a todos os pontos de treino:

$$J(\mathbf{W}, \mathbf{B}) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log \hat{y}^{(i)} + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)})]$$

Propagação para frente

- $z_1^{[1]} = \mathbf{W}_{1,1}^{[1]}x_1 + \mathbf{W}_{1,2}^{[1]}x_2 + b_1^{[1]}$
- Ou, de forma matricial:
- $\mathbf{z}^{[1]} = \mathbf{W}^{[1]}\mathbf{x} + \mathbf{b}^{[1]}$
- $\mathbf{a}^{[1]} = \mathbf{g}^{[1]}(\mathbf{z}^{[1]}) \leftarrow$ Tangente hiperbólica
- $\mathbf{z}^{[2]} = \mathbf{W}^{[2]}\mathbf{a}^{[1]} + \mathbf{b}^{[2]}$
- $\mathbf{a}^{[2]} = \mathbf{g}^{[2]}(\mathbf{z}^{[2]}) \leftarrow$ Sigmoid
- $\hat{y} = \mathbf{a}^{[2]}$



Propagação para frente

$$\mathbf{z}^{[1]} = \mathbf{W}^{[1]}\mathbf{x} + \mathbf{b}^{[1]}$$

$$\mathbf{a}^{[1]} = \mathbf{g}^{[1]}(\mathbf{z}^{[1]})$$

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

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

$$\hat{y} = \mathbf{a}^{[2]}$$

Variável	Linhas	Colunas
x		
$\mathbf{z}^{[1]}$		
$\mathbf{W}^{[1]}$		
$\mathbf{b}^{[1]}$		
$\mathbf{a}^{[1]}$		
$\mathbf{z}^{[2]}$		
$\mathbf{W}^{[2]}$		
$\mathbf{b}^{[2]}$		
$\mathbf{a}^{[2]}$		

Propagação para frente

$$\mathbf{z}^{[1]} = \mathbf{W}^{[1]}\mathbf{x} + \mathbf{b}^{[1]}$$

$$\mathbf{a}^{[1]} = \mathbf{g}^{[1]}(\mathbf{z}^{[1]})$$

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

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

$$\hat{y} = \mathbf{a}^{[2]}$$

Variável	Linhas	Colunas
x	2	n_s
$W^{[1]}$	n_h	2
$b^{[1]}$	n_h	1
$z^{[1]}$	n_h	n_s
$a^{[1]}$	n_h	n_s
$W^{[2]}$	1	n_h
$b^{[2]}$	1	1
$z^{[2]}$	1	n_s
$a^{[2]}$	1	n_s

Como treinar a RNA?

Devemos encontrar os valores de W e b que minimizam o erro J

$$\mathbf{z}^{[1]} = \mathbf{W}^{[1]} \mathbf{x} + \mathbf{b}^{[1]}$$

$$\mathbf{a}^{[1]} = \mathbf{g}^{[1]}(\mathbf{z}^{[1]})$$

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

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

$$\hat{y} = \mathbf{a}^{[2]}$$

$$J(\mathbf{W}, \mathbf{B}) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log \hat{y}^{(i)} + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)})]$$

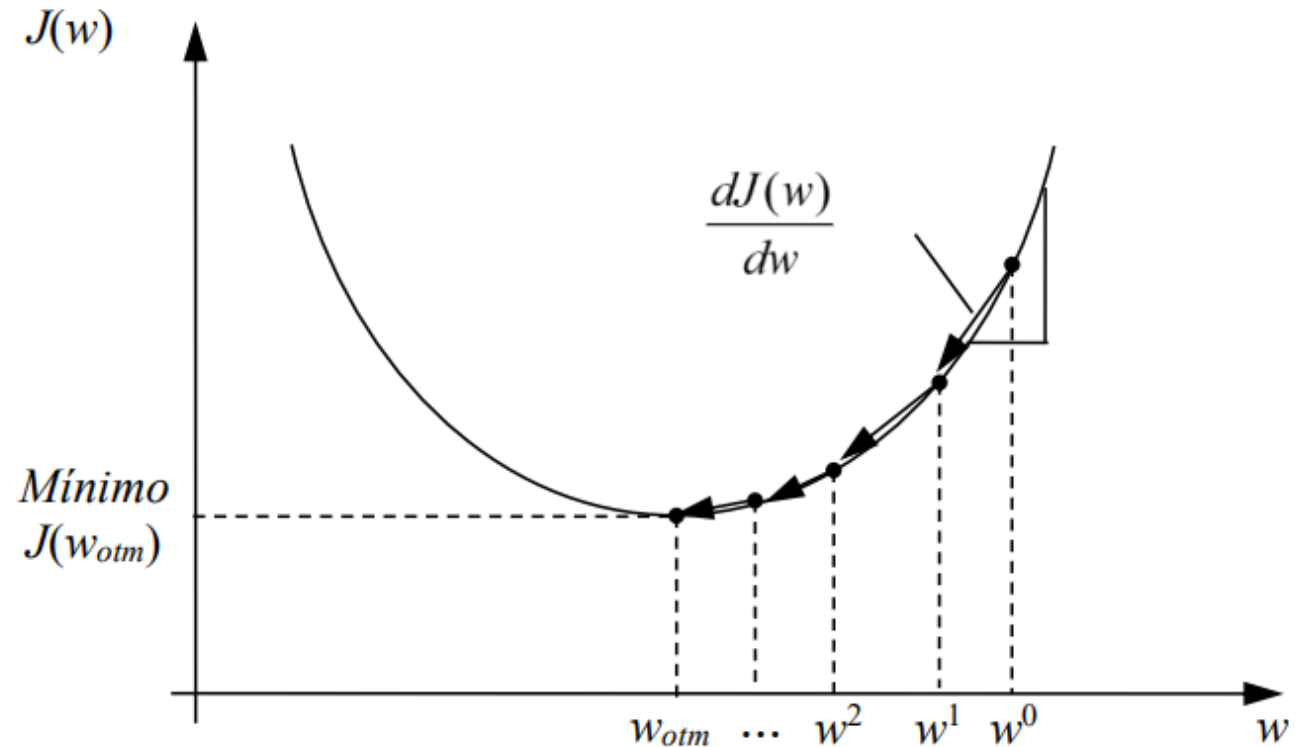
Variável	Linhas	Colunas
x	2	n_s
$W^{[1]}$	n_h	2
$b^{[1]}$	n_h	1
$z^{[1]}$	n_h	n_s
$a^{[1]}$	n_h	n_s
$W^{[2]}$	1	n_h
$b^{[2]}$	1	1
$z^{[2]}$	1	n_s
$a^{[2]}$	1	n_s

Gradiente descendente

- Problema de otimização
- Minimizar uma função
- Queremos achar W que minimiza J

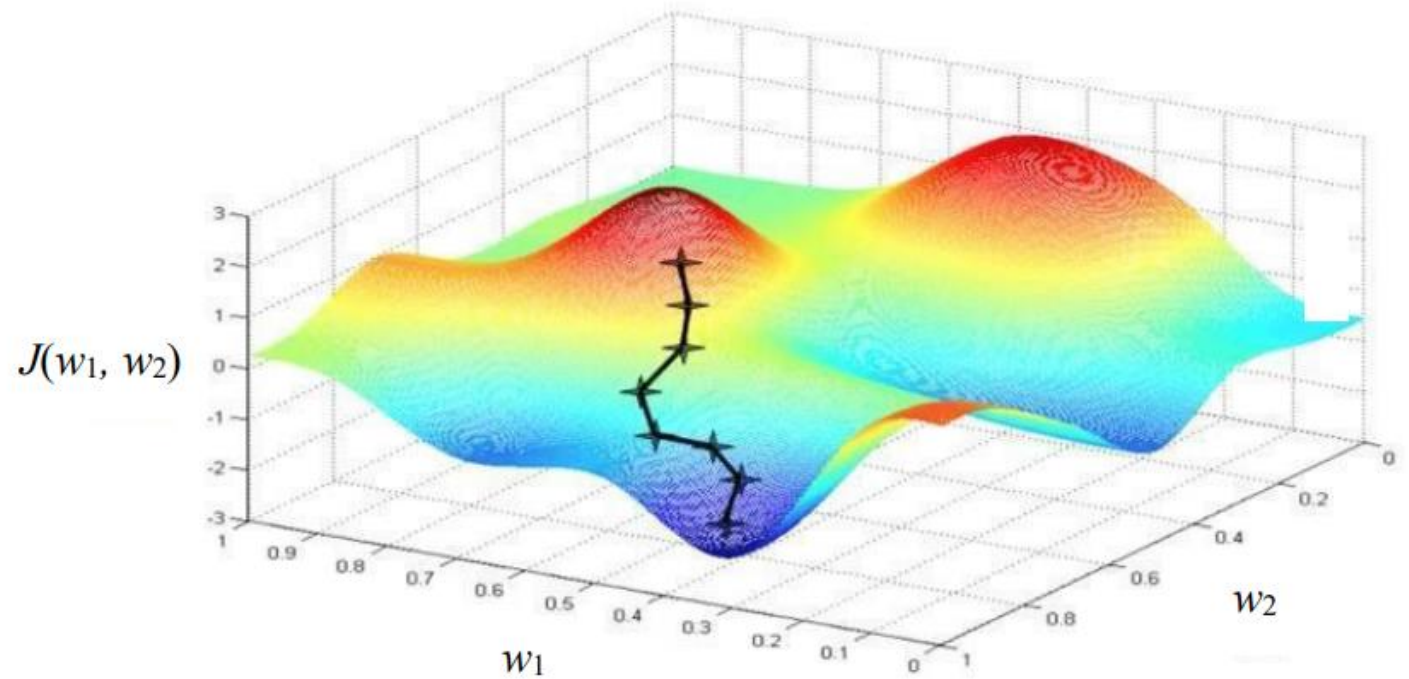
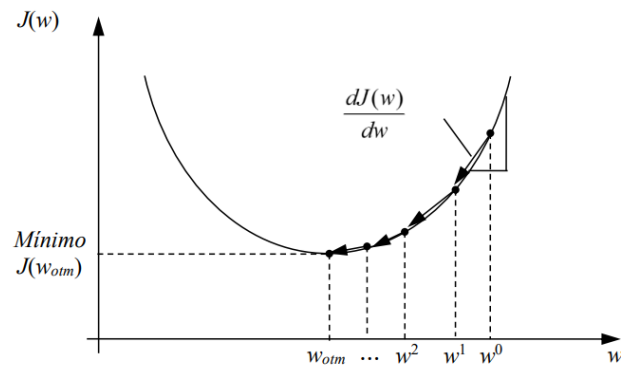
$$\frac{\partial J}{\partial W}$$

Qual o formato dessa derivada?



Gradiente descendente

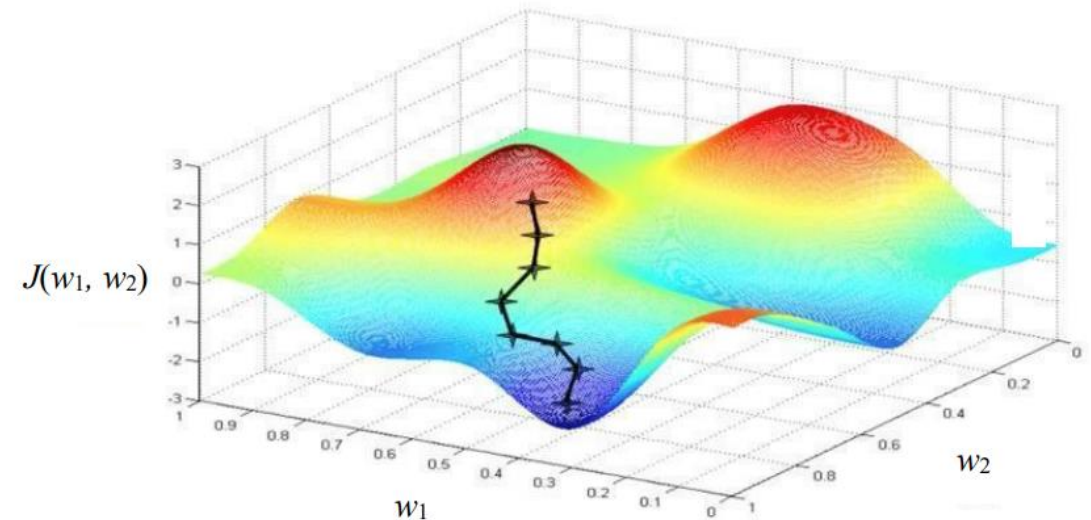
- Problema de otimização
- Minimizar uma função



<https://github.com/RNogales94/ISI-Tensorflow>

Gradiente descendente

- Precisamos do gradiente do custo em relação a cada um dos pesos
- Aqui entra o *backpropagation*



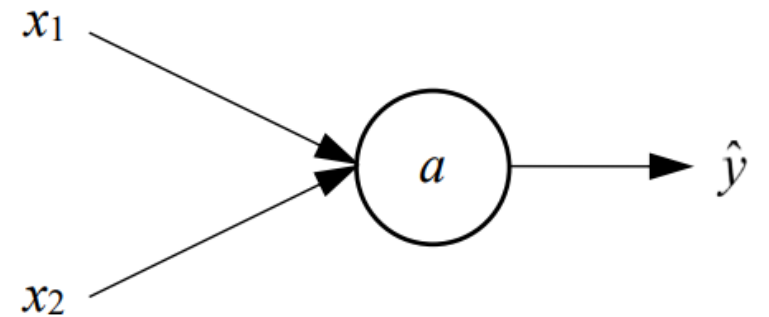
<https://github.com/RNogales94/ISI-Tensorflow>

$$w_{k,j}^{[l]} = w_{k,j}^{[l]} - \alpha \frac{\partial J(\mathbf{W}, \mathbf{B})}{\partial w_{k,j}^{[l]}}, \text{ para } l = 1, \dots, L, k = 1, \dots, n^{[l]}, j = 1, \dots, n^{[l-1]}$$

$$b_k^{[l]} = b_k^{[l]} - \alpha \frac{\partial J(\mathbf{W}, \mathbf{B})}{\partial b_k^{[l]}}, \text{ para } l = 1, \dots, L, k = 1, \dots, n^{[l]}$$

Exemplo em uma RNA de um neurônio

- 2 entradas
- 1 saída
- m exemplos de treinamento
- 1 camada
- 1 neurônio



Eduardo Lobo Lustosa Cabral

Exemplo em uma RNA de um neurônio

Propagação para frente

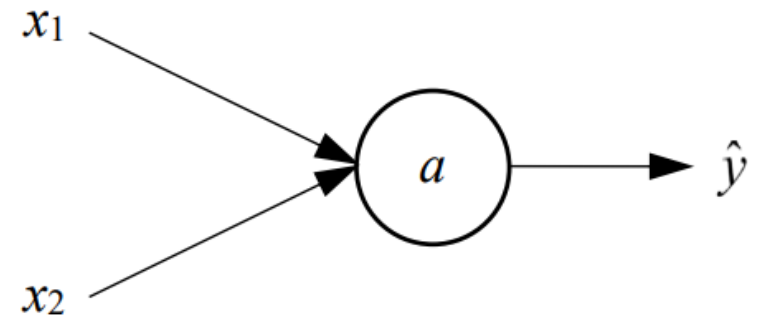
- $z = w_1x_1 + w_2x_2 + b$
- $\hat{y} = a = g(z)$

Parâmetros:

- w_1, w_2, b

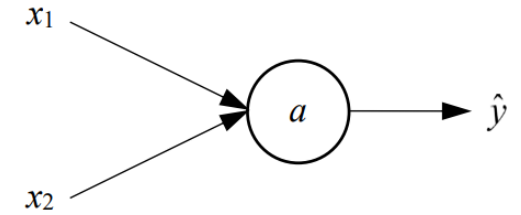
Custo:

- $J(w_1, w_2, b) = \frac{1}{m} \sum_{i=1}^m E(\hat{y}^{(i)}, y^{(i)})$



Eduardo Lobo Lustosa Cabral

Exemplo em uma RNA de um neurônio



Propagação para frente

- $z = w_1x_1 + w_2x_2 + b$
- $\hat{y} = a = g(z)$

Parâmetros:

- w_1, w_2, b

Custo:

- $J(w_1, w_2, b) = \frac{1}{m} \sum_{i=1}^m E(\hat{y}^{(i)}, y^{(i)})$

$\left[\frac{\partial E(a, y)}{\partial a} \right] =$ derivada da função de erro em relação à saída calculada ($\hat{y} = a$);

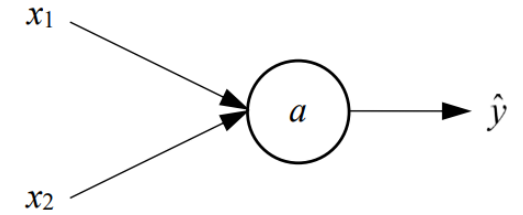
$\frac{\partial a}{\partial z} = \frac{d[g(z)]}{dz} \Rightarrow$ derivada da função de ativação g

$$\frac{\partial z}{\partial w_1} = x_1$$

$$\frac{\partial z}{\partial w_2} = x_2$$

$$\frac{\partial z}{\partial b} = 1$$

Exemplo em uma RNA de um neurônio



$$\left[\frac{\partial E(a, y)}{\partial a} \right] = \text{derivada da função de erro em relação à saída calculada } (\hat{y} = a);$$

$$\frac{\partial a}{\partial z} = \frac{d[g(z)]}{dz} \Rightarrow \text{derivada da função de ativação } g$$

$$\frac{\partial z}{\partial w_1} = x_1$$

$$\frac{\partial z}{\partial w_2} = x_2$$

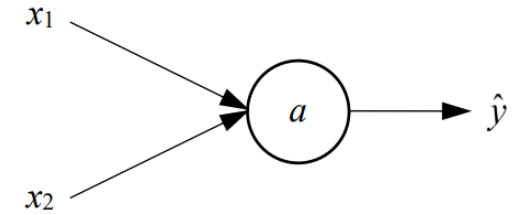
$$\frac{\partial z}{\partial b} = 1$$

$$\frac{\partial E(a, y)}{\partial w_1} = \left[\frac{\partial E(a, y)}{\partial a} \right] \left[\frac{d[g(z)]}{dz} \right] x_1$$

$$\frac{\partial E(a, y)}{\partial w_2} = \left[\frac{\partial E(a, y)}{\partial a} \right] \left[\frac{d[g(z)]}{dz} \right] x_2$$

$$\frac{\partial E(a, y)}{\partial b} = \left[\frac{\partial E(a, y)}{\partial a} \right] \left[\frac{d[g(z)]}{dz} \right]$$

Exemplo em uma RNA de um neurônio



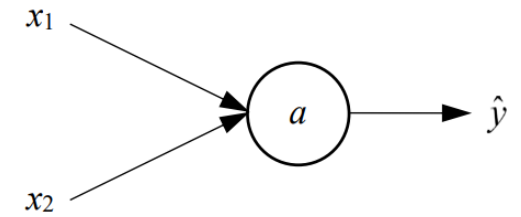
$$\boxed{\frac{\partial E(a, y)}{\partial w_1} = \left[\frac{\partial E(a, y)}{\partial a} \right] \left[\frac{d[g(z)]}{\partial z} \right] x_1}$$

Para função logística

$$\frac{\partial E(a, y)}{\partial a} = \begin{cases} y = 0 \rightarrow -1/(1 - a) \\ y = 1 \rightarrow -1/a \end{cases}$$

$$L(\hat{y}^{(i)}, y^{(i)}) = -y^{(i)} \log \hat{y}^{(i)} - (1 - y^{(i)}) \log(1 - \hat{y}^{(i)})$$

Exemplo em uma RNA de um neurônio

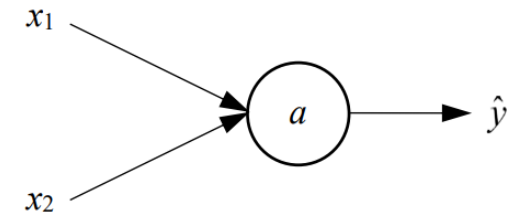


$$\boxed{\frac{\partial E(a, y)}{\partial w_1} = \left[\frac{\partial E(a, y)}{\partial a} \right] \left[\frac{d[g(z)]}{dz} \right] x_1}$$

Para tangente hiperbólica

$$\frac{dg(z)}{dz} = \frac{d}{dz} \left[\frac{e^z - e^{-z}}{e^z + e^{-z}} \right] = \frac{(e^z + e^{-z})^2 - (e^z - e^{-z})^2}{(e^z + e^{-z})^2} = \underbrace{\frac{(e^z + e^{-z})^2}{(e^z + e^{-z})^2}}_1 - \underbrace{\left[\frac{(e^z - e^{-z})}{(e^z + e^{-z})} \right]^2}_{\tanh(z)} = 1 - a^2$$

Exemplo em uma RNA de um neurônio

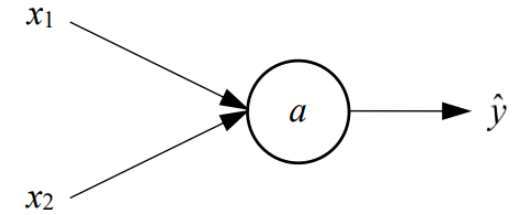


$$\boxed{\frac{\partial E(a, y)}{\partial w_1} = \left[\frac{\partial E(a, y)}{\partial a} \right] \left[\frac{d[g(z)]}{dz} \right] x_1}$$

Para sigmoide

$$\frac{d\sigma(z)}{dz} = \frac{d}{dz} \left[\frac{1}{1 + e^{-z}} \right] = \frac{e^{-z}}{(1 + e^{-z})^2} = \underbrace{\frac{1}{(1 + e^{-z})}}_a \underbrace{\left[1 - \frac{1}{(1 + e^{-z})} \right]}_{(1-a)} = a(1-a)$$

Exemplo em uma RNA de um neurônio



$$\frac{\partial E(a, y)}{\partial w_1} = \left[\frac{\partial E(a, y)}{\partial a} \right] \left[\frac{d[g(z)]}{\partial z} \right] x_1$$

$$\frac{\partial E(a, y)}{\partial w_2} = \left[\frac{\partial E(a, y)}{\partial a} \right] \left[\frac{d[g(z)]}{\partial z} \right] x_2$$

$$\frac{\partial E(a, y)}{\partial b} = \left[\frac{\partial E(a, y)}{\partial a} \right] \left[\frac{d[g(z)]}{\partial z} \right]$$

Lembrando que:

$$J(w_1, w_2, b) = \frac{1}{m} \sum_{i=1}^m E(\hat{y}^{(i)}, y^{(i)})$$

Podemos somar $\frac{\partial E}{\partial w_i}$ para cada ponto e obter $\frac{\partial J(w_1, w_2, b)}{\partial w_i}$

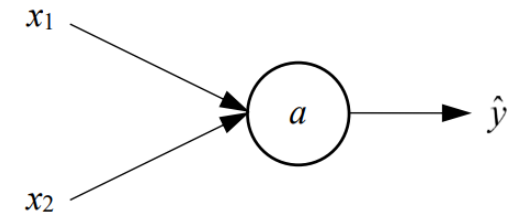
Exemplo em uma RNA de um neurônio

Lembrando que:

$$J(w_1, w_2, b) = \frac{1}{m} \sum_{i=1}^m E(\hat{y}^{(i)}, y^{(i)})$$

Podemos usar o $\frac{\partial E}{\partial w_i}$ de cada ponto e obter

$$\frac{\partial J(w_1, w_2, b)}{\partial w_i} = \frac{1}{m} \sum_{i=1}^m \frac{\partial E(\hat{y}^{(i)}, y^{(i)})}{\partial w_i}$$



Finalmente, para atualizar os pesos, usamos:

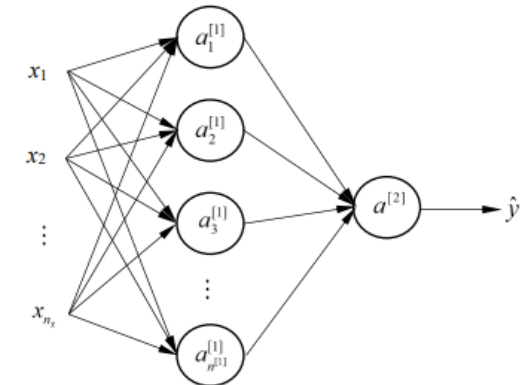
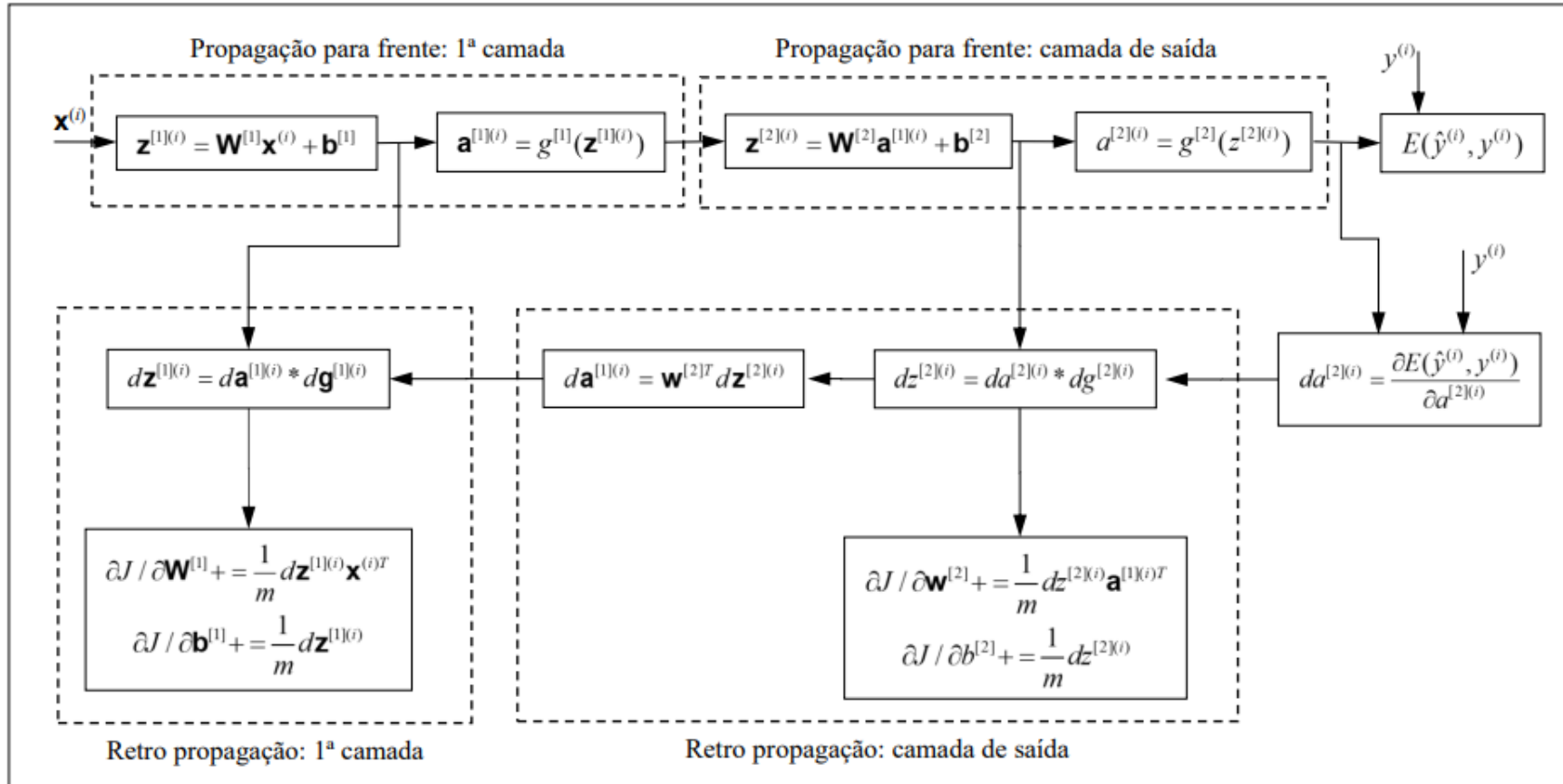
$$w_1 = w_1 - \alpha \frac{\partial J(w_1, w_2, b)}{\partial w_1},$$

$$w_2 = w_2 - \alpha \frac{\partial J(w_1, w_2, b)}{\partial w_2},$$

$$b = b - \alpha \frac{\partial J(w_1, w_2, b)}{\partial b}.$$

Onde α é a taxa de aprendizagem

E com mais camadas?



Eduardo Lobo Lustosa Cabral

Trabalho 1

- Classificação binária
- Implementado num notebook de Python
- Sem usar bibliotecas de ML
- Rede neural rasa

