

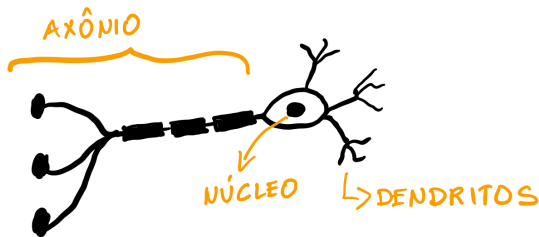
Redes Neurais com Perceptron Multicamadas

Advanced Institute for Artificial Intelligence – AI2

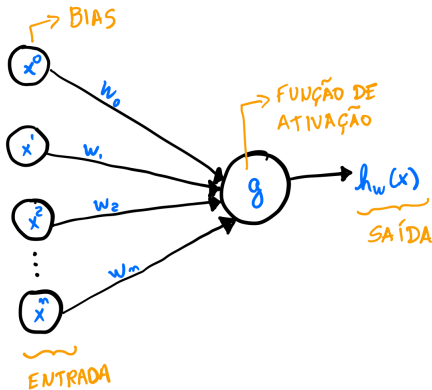
<https://advancedinstitute.ai>

Introdução

Redes Neurais foram amplamente utilizadas nos anos 80 e 90 e visam imitar o funcionamento do cérebro humano. Sua popularidade caiu no final dos anos 90, mas voltaram à cena com novas abordagens baseadas em aprendizado em profundidade. Mas como elas funcionam? Vamos dar uma olhada, primeiramente, na estrutura de um neurônio.



Matematicamente falando, podemos representar um neurônio da seguinte forma (Modelo de McCulloch-Pitts):



Alguns exemplos de funções de ativação:

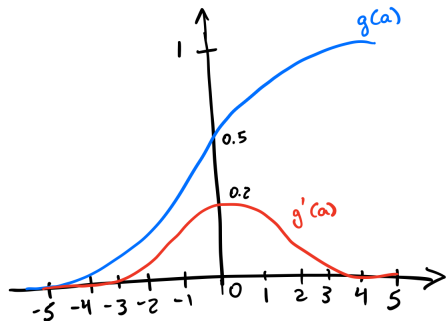
- Função logística (sigmoide): $g(a) = \frac{1}{1 + e^{-a}}$, tal que $g(a) \in [0, 1]$.
- Função limiar (degrau):

$$g_{\theta}(a) = \begin{cases} 1 & \text{caso } w^T x \geq \theta \\ 0 & \text{caso contrário.} \end{cases}$$

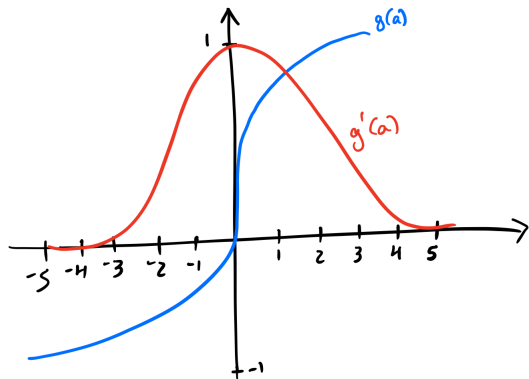
tal que $g(a) \in \{0, 1\}$.

- Função tangente hiperbólica: $g(a) = 2\sigma(2a) - 1$, tal que $g(a) \in [-1, 1]$ e $\sigma(x)$ corresponde à função logística.

É desejável que uma função de ativação seja **derivável**. Como escolhê-la? Depende de sua aplicação.



Suponha $g(a) = \frac{1}{1 + e^{-a}}$. Temos que $g'(a) = g(a)(1 - g(a))$. Note que $g'(a)$ satura quando $a > 5$ ou $a < -5$. Ademais, $g'(a) < 1$, $\forall a$. Isso significa que, para redes com muitas camadas, o gradiente tende a **desaparecer** durante o treinamento.



Suponha $g(a) = 2\sigma(2a) - 1$. Temos que $g'(a) = 1 - g^2(a)$. Muito embora saturações aconteçam, $g'(a)$ atinge valores maiores, chegando, inclusive, ao máximo de 1 quando $a = 0$.

Perceptron

O Perceptron foi formalmente proposto por McCulloch e Pitts na década de 40 com o intuito de modelar matematicamente o neurônio humano. Muito embora ele tenha sido a base para muitos algoritmos, seu poder discriminativo é limitado, pois consegue aprender apenas hiperplanos como função de decisão.

Definição do problema: seja um conjunto de dados $\mathcal{X} = \{(x_1, y_1), (x_2, y_2), \dots, (x_z, y_z)\}$ tal que $x_i \in \mathbb{R}^{n+1}$ corresponde ao dado de entrada e $y_i \in [-1, +1]$ denota o seu respectivo valor de saída. Temos, ainda, que \mathcal{X} pode ser **particionado** da seguinte forma: $\mathcal{X} = \mathcal{X}^1 \cup \mathcal{X}^2$, em que \mathcal{X}^1 e \mathcal{X}^2 denotam os conjuntos de dados de **treinamento** e **teste**, respectivamente. Nosso objetivo é, dado o conjunto de treinamento, aprender uma função $h : \mathbb{R}^{n+1} \rightarrow \{-1, +1\}$ que consiga atribuir a classe correta à uma determinada amostra.

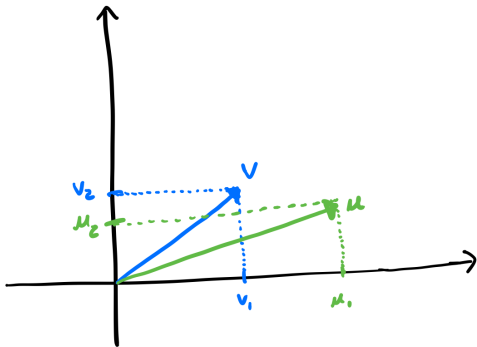
Vamos, agora, adaptar a função de ativação limiar para que possamos criar a nossa função hipótese:

$$h_{\mathbf{w}}(\mathbf{x}) = \begin{cases} +1 & \text{caso } \mathbf{w}^T \mathbf{x} \geq \theta \\ -1 & \text{caso contrário.} \end{cases} \quad (1)$$

Para simplificar a notação, é usual trazermos θ para o lado esquerdo da equação e atribuírmos $w_0 = \theta$. Novamente, consideraremos $x^0 = 1$. Assim, temos a função hipótese atualizada como segue:

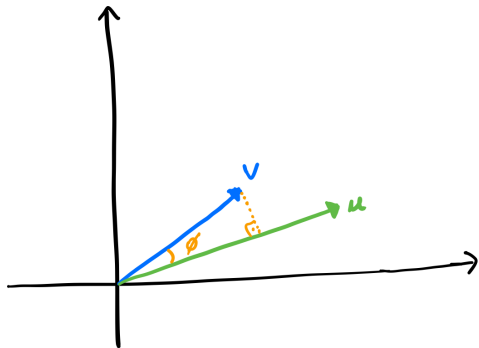
$$h_{\mathbf{w}}(\mathbf{x}) = \begin{cases} +1 & \text{caso } \mathbf{w}^T \mathbf{x} \geq 0 \\ -1 & \text{caso contrário.} \end{cases} \quad (2)$$

Vamos, agora, revisar alguns conceitos de Geometria Analítica. Suponha que tenhamos dois vetores $\mathbf{u} = [u_1 \ u_2]$ e $\mathbf{v} = [v_1 \ v_2]$. Geometricamente, podemos representá-los como segue:



Lembrando que $\|\mathbf{u}\| = \sqrt{u_1^2 + u_2^2}$ denota o comprimento (magnitude) de \mathbf{u} .

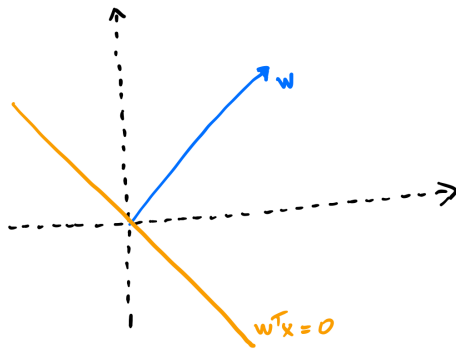
Temos, também, a definição de **projeção escalar** entre dois vetores, dada como segue:



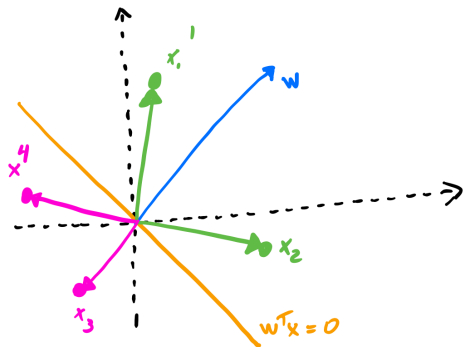
Podemos, então, representar a projeção escalar entre dois vetores da seguinte forma: $w^T x = \|w\| \cdot \cos \phi$. Lembrando que temos as seguintes situações:

- $\cos \phi > 0$ quando $0 < \phi < 90^\circ$.
- $\cos \phi < 0$ quando $90^\circ < \phi < 270^\circ$.
- $\cos \phi = 0$ quando $\phi = 90^\circ$ ou $\phi = 270^\circ$

Voltando ao nosso problema inicial, temos que a equação $w^T x = 0$ define um hiperplano ortogonal ao vetor de pesos w deslocado de $-\theta$ (assumindo que $w_0 = -\theta$ e $x_0 = 1$). Vamos supor que $\theta = 0$, ou seja, o hiperplano passa pela origem.



Como aprender o conjunto de pesos w ? A ideia intuitiva é **mover** o vetor de pesos de tal forma que as amostras fiquem posicionadas corretamente no espaço de características. No exemplo abaixo, temos um conjunto de dados $\mathcal{X}^1 = \{(\mathbf{x}_1, +1), (\mathbf{x}_2, +1), (\mathbf{x}_3, -1), (\mathbf{x}_4, -1)\}$, de tal forma que o hiperplano $w^T \mathbf{x} = 0$ já se encontra posicionado corretamente.



A regra de atualização dos pesos é dada pela seguinte fórmula:

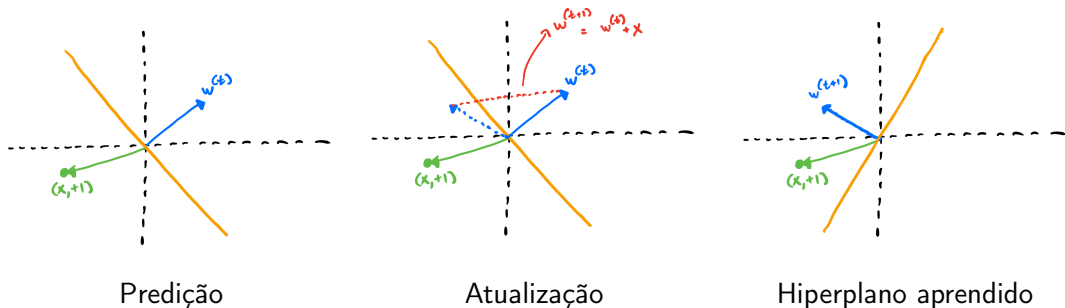
$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} + \alpha(y_i - h_{\mathbf{w}^{(t)}}(\mathbf{x}_i))\mathbf{x}_i, \quad (3)$$

tal que $\forall i = 1, 2, \dots, m$.

Mas como ele funciona na prática? Suponha uma amostra $x \in \mathcal{X}$ tal que $y = +1$ e $h_w(x) = -1$. Para um valor de $\alpha = 0.5$, temos que:

$$\begin{aligned} w^{(t+1)} &= w^{(t)} + \alpha(1 - (-1))x \\ &= w^{(t)} + 0.5(2)x = w^{(t)} + x. \end{aligned} \quad (4)$$

Desta forma, o vetor de pesos w será rotacionado de modo que x seja positivo.



De maneira análoga, caso o rótulo de x seja negativo, ou seja, $y = -1$, temos que:

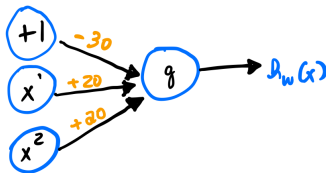
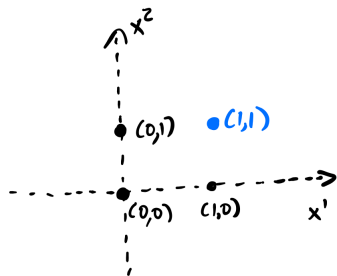
$$\begin{aligned} \mathbf{w}^{(t+1)} &= \mathbf{w}^{(t)} + \alpha(-1 - (+1))\mathbf{x} \\ &= \mathbf{w}^{(t)} + 0.5(-2)\mathbf{x} = \mathbf{w}^{(t)} - \mathbf{x}. \end{aligned} \tag{5}$$

Assim, o vetor de pesos será rotacionado (por meio das projeções) para o outro lado. Para um conjunto de dados que seja **linearmente separável**, foi provado matematicamente que o algoritmo do Perceptron tem sua **convergência garantida**.

Como funciona o seu algoritmo? Vamos lá:

- ➊ Atribua pesos aleatórios para w .
- ➋ Inicialize α .
- ➌ $t = 0$.
- ➍ Para cada amostra $x_i \in \mathcal{X}^1$, faça:
 - ➊ $w^{(t+1)} = w^{(t)} + \alpha(y_i - h_{w^{(t)}}(x_i))x_i$
- ➎ Repita o passo 4 até algum critério de convergência for estabelecido.

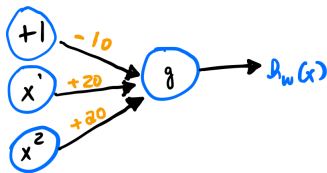
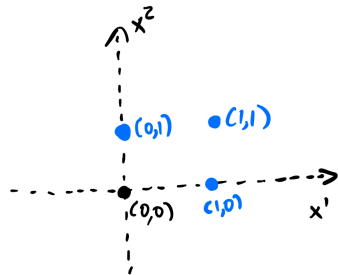
Vejamos alguns exemplos de funcionamento do Perceptron. Considere o problema de resolver a seguinte equação lógica: $y = x^1 \text{ AND } x^2$. Para duas entradas, temos 2^2 possibilidades de amostras, isto é, o nosso conjunto de dados é composto pelos seguintes elementos: $\mathcal{X} = \{([0 \ 0], 0), ([0 \ 1], 0), ([1 \ 0], 0), ([1 \ 1], 1)\}$. Conseguimos encontrar o hiperplano separador?



Nossa função hipótese é dada por $h_w(x) = g(-30 + 20x^1 + 20x^2)$. Utilizando g como uma função logística, temos que:

x^1	x^2	$h_w(x)$
0	0	$g(-30) \approx 0$
0	1	$g(-10) \approx 0$
1	0	$g(-10) \approx 0$
1	1	$g(10) \approx 1$

Considere, agora, o problema de resolver a seguinte equação lógica: $y = x^1 \text{ OR } x^2$. Para duas entradas, temos 2^2 possibilidades de amostras, isto é, o nosso conjunto de dados é composto pelos seguintes elementos: $\mathcal{X} = \{([0 \ 0], 0), ([0 \ 1], 1), ([1 \ 0], 1), ([1 \ 1], 1)\}$. Conseguimos encontrar o hiperplano separador?

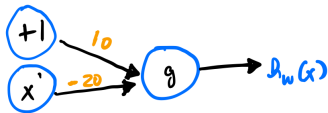


Nossa função hipótese é dada por $h_w(x) = g(-10 + 20x^1 + 20x^2)$. Utilizando g como uma função logística, temos que:

x^1	x^2	$h_w(x)$
0	0	$g(-10) \approx 0$
0	1	$g(10) \approx 1$
1	0	$g(10) \approx 1$
1	1	$g(30) \approx 1$

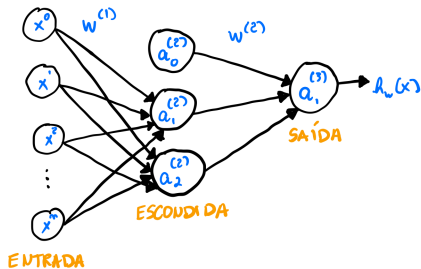
Considere, agora, o problema de resolver a seguinte equação lógica: $y = \text{NOT } x^1$. Agora temos apenas uma entrada, ou seja, x^1 . Assim, o nosso conjunto de dados é composto pelos seguintes elementos: $\mathcal{X} = \{(0, 1), (1, 0)\}$. Conseguimos encontrar o hiperplano separador?

Nossa função hipótese é dada por $h_w(x) = g(10 - 20x^1)$. Utilizando g como uma função logística, temos que:



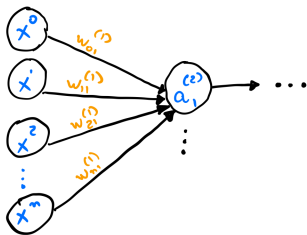
x^1	$h_w(x)$
0	$g(10) \approx 1$
1	$g(-10) \approx 0$

O que seria então uma Rede Neural do tipo Perceptron Multicamadas (*Multilayer Perceptron* - MLP)? Basicamente, é um grupo de neurônios que, combinados, permitem o aprendizado de um número maior de funções de decisão.



Na ilustração acima, temos que $a_i^{(j)}$ denota o neurônio i da camada j , e $\mathbf{W}^{(l)}$ é a matriz de pesos que conecta as camadas l e $l + 1$. Essa arquitetura é geralmente representada por $n:3:1$.

Considerando a rede neural anterior, vamos analisar uma situação mais específica.

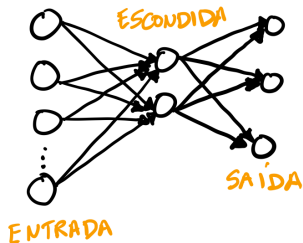


$$a_1^{(2)} = g(\underbrace{w_{01}^{(1)}x^0 + w_{11}^{(1)}x^1 + w_{21}^{(1)}x^2 + \dots + w_{n1}^{(1)}x^n}_{\sum_{i=0}^n w_{i1}^{(1)}x^i = b_1^{(2)}})$$

A função de decisão final da rede neural é dada pela seguinte formulação:

$$h_{\mathbf{x}}(\mathbf{w}) = a_1^{(3)} = g \left(w_{01}^{(2)}a_0^{(2)} + w_{11}^{(2)}a_1^{(2)} + w_{21}^{(2)}a_2^{(2)} \right). \quad (6)$$

E no caso de problemas com múltiplas classes? Neste caso, para um problema com c classes, nossa camada de saída necessita de c neurônios.



Podemos utilizar a metodologia de codificação *one-hot* para representar cada neurônio de saída, em que $h_w(x) \in \mathbb{R}^3$:

$$h_w^1(x) \approx \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$$

Classe 1

$$h_w^2(x) \approx \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$$

Classe 2

$$h_w^3(x) \approx \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$$

Classe 3

O mesmo acontece com o rótulo y de cada amostra, que passa a ser agora um vetor $y \in \mathbb{R}^3$.

Existem diversos algoritmos de treinamento para Redes Neurais MLP, em que o mais conhecido é chamado de **retropropagação**, do inglês *backpropagation*. Ele possui dois passos: (i) propagação para frente (*forward propagation*) e (ii) propagação para trás (*backpropagation*). Antes de estudarmos o seu funcionamento, vamos dar uma olhada na função de custo:

$$J(\mathbf{w}) = \frac{1}{m} \left[\sum_{i=1}^m \sum_{k=1}^c -y_i^k \log(h_{\mathbf{w}}^k(\mathbf{x}_i)) - (1 - y_i^k) \log(1 - h_{\mathbf{w}}^k(\mathbf{x}_i)) \right]. \quad (7)$$

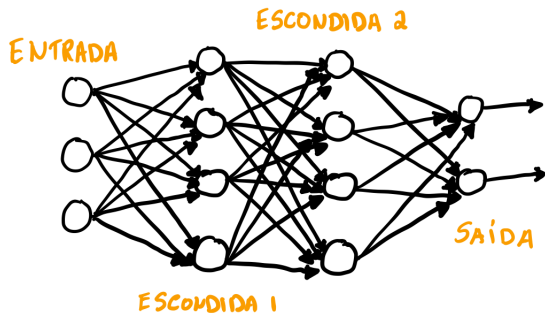
Fazendo uma analogia, a formulação acima contempla uma rede neural com c regressores logísticos caso tenhamos uma função de ativação logística na camada de saída.

Podemos utilizar, novamente, o gradiente descendente para otimizar a função de custo $J(\mathbf{w})$. No entanto, note que o problema passa a ser mais complexo, pois precisamos calcular as derivadas parciais com relação à todos os pesos da rede, ou seja:

$$\frac{\partial J(\mathbf{w})}{\partial w_{ij}^{(l)}}, \quad (8)$$

em que $l = 1, 2, \dots, L - 1$ tal que L denota a quantidade de camadas da rede neural.

Seja uma rede do tipo 3:4:4:2 ilustrada abaixo. Suponha que tenhamos apenas uma amostra no conjunto de treinamento. Neste caso, o passo *forward* é dado pelas seguintes etapas:



$$\textcircled{1} \mathbf{a}^{(1)} \leftarrow \mathbf{x}$$

$$\textcircled{2} \mathbf{b}^{(2)} \leftarrow \left(\mathbf{W}^{(1)} \right)^T \mathbf{a}^{(1)}$$

$$\textcircled{3} \mathbf{a}^{(2)} \leftarrow g \left(\mathbf{b}^{(2)} \right)$$

$$\textcircled{4} \mathbf{b}^{(3)} \leftarrow \left(\mathbf{W}^{(2)} \right)^T \mathbf{a}^{(2)}$$

$$\textcircled{5} \mathbf{a}^{(3)} \leftarrow g \left(\mathbf{b}^{(3)} \right)$$

$$\textcircled{6} \mathbf{b}^{(4)} \leftarrow \left(\mathbf{W}^{(3)} \right)^T \mathbf{a}^{(3)}$$

$$\textcircled{7} h_w(\mathbf{x}) = \mathbf{a}^{(4)} \leftarrow g \left(\mathbf{b}^{(4)} \right)$$

Para o passo de *backpropagation*, temos a definição de uma nova variável $\delta_j^{(l)}$ que denota um erro parcial acumulado no neurônio j da camada l . Este erro deve ser calculado de forma diferente para os neurônios de saída e das camadas escondidas, como segue:

- Camada de saída ($l = 4$):

$$\begin{aligned}\delta_j^{(4)} &= a_j^{(4)} - y^j \\ &= h_w^j(\mathbf{x}) - y^j\end{aligned}\tag{9}$$

Em notação vetorial, temos que $\boldsymbol{\delta}^{(4)} = \mathbf{a}^{(4)} - \mathbf{y} = \mathbf{h}_w(\mathbf{x}) - \mathbf{y}$.

- Camadas escondidas $l = \{2, 3\}$:
 - $\delta^{(3)} = \left(\mathbf{W}^{(3)}\right)^T \delta^{(4)} \cdot * g' \left(\mathbf{b}^{(3)}\right)$
 - $\delta^{(2)} = \left(\mathbf{W}^{(2)}\right)^T \delta^{(3)} \cdot * g' \left(\mathbf{b}^{(2)}\right)$

Na prática, temos a seguinte formulação para propagação dos erros nas camadas intermediárias:

$$\delta^{(l)} = \left(\mathbf{W}^{(l)}\right)^T \delta^{(l+1)} \cdot * g' \left(\mathbf{b}^{(l)}\right). \quad (10)$$

O nome **retropropagação** deve-se ao fato de o algoritmo “propagar para trás” o erro estimado em cada camada. As derivadas parciais podem ser calculadas como segue:

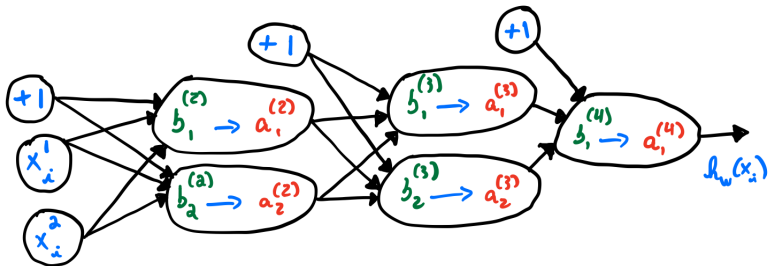
$$\frac{\partial J(\mathbf{w})}{\partial w_{ij}^{(l)}} = a_i^{(l)} \delta_j^{(l+1)}. \quad (11)$$

Note, então, que as derivadas parciais são calculadas com relação à todos os pesos da rede neural.

Segue algoritmo de retropropagação para treinamento de uma rede neural MLP.

- ❶ Atribua pesos aleatórios para $w_{ij}^{(l)} \quad \forall l, i, j$
- ❷ Execute os passos abaixo ate o critério de parada ter sido estabelecido: Laço das épocas
 - ❶ $\Delta_{ij}^{(l)} = 0 \quad \forall l, i, j$ Variável utilizada para armazenar $\frac{\partial J(\mathbf{w})}{\partial w_{ij}^{(l)}}$
 - ❷ Para cada amostra $\mathbf{x}_i \in \mathcal{X}^1$, faça:
 - ❶ Execute o passo de *forward propagation* para calcular $\mathbf{a}^l, \quad l = 2, 3, \dots, L$
 - ❷ Execute o passo de *backward propagation* para calcular:
 - ❶ $\delta^l, \quad l = L, L - 1, \dots, 2$ Erro em cada neurônio
 - ❷ $\Delta_{ij}^{(l)} = \Delta_{ij}^{(l)} + a_i^{(l)} \delta_j^{(l+1)}$ Derivadas parciais são acumuladas
 - ❸ $D_{ij}^{(l)} = \frac{1}{m} \Delta_{ij}^{(l)} \quad \forall l, i, j$ Calcula o gradiente médio
 - ❹ $w_{ij}^{(l)} = w_{ij}^{(l)} - \alpha D_{ij}^{(l)} \quad \forall l, i, j$ Atualiza pesos com gradiente descendente
 - ❺ Avalie a função de custo $J(\mathbf{w})$

Resumindo: passo *forward*.



Resumindo: passo *backward*.

