

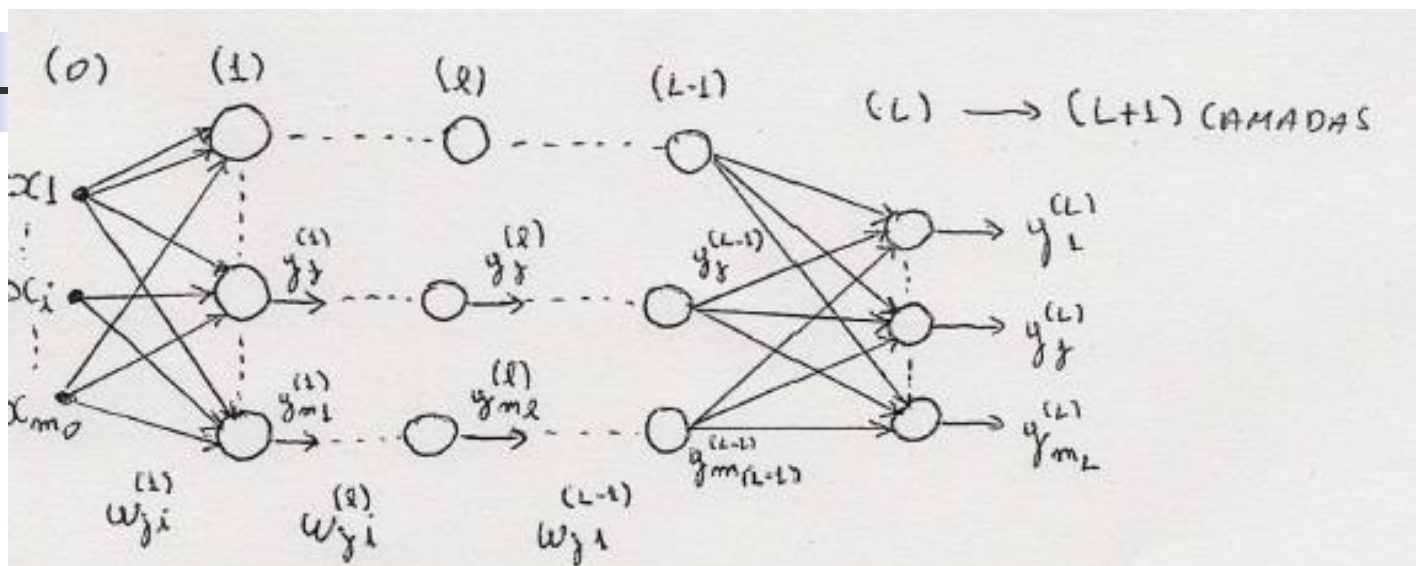


# Redes Multicamadas

---

Algoritmo de Retropropagação do  
Erro  
(*Backpropagation*)

# Notação



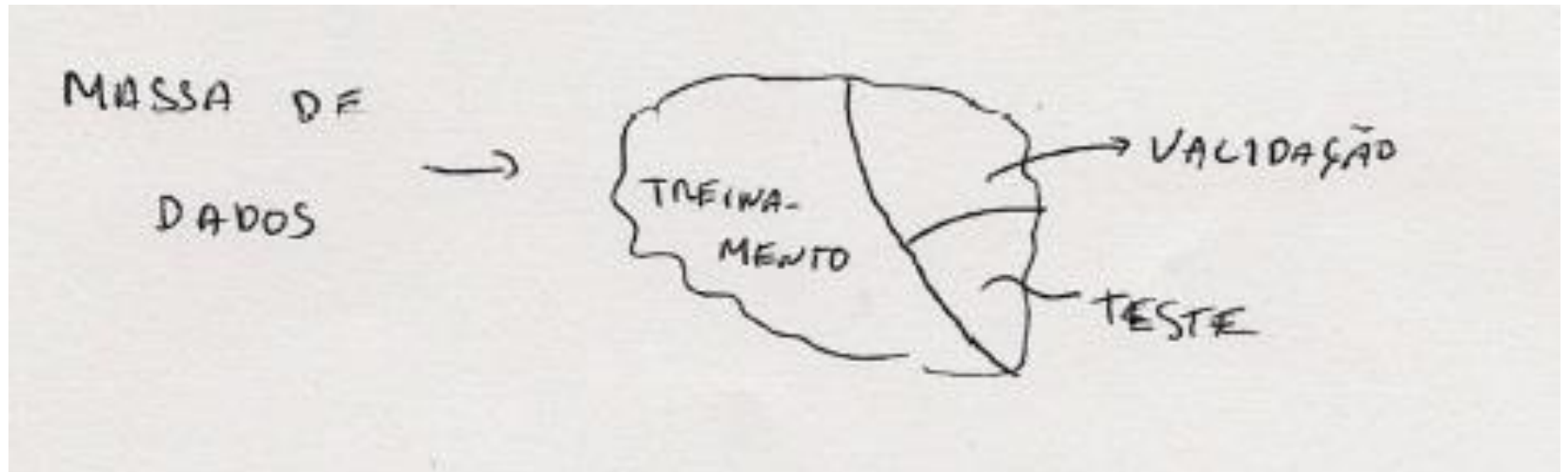
## NOTAÇÃO

$m_l = m_l^o$  - DE NEURÔNIOS NA CAMADA  $l$  (EXCLUINDO O "BIAS")

$w_{ji}^{(l)}(m)$  = PESO SINÁPTICO CONECTANDO A SAÍDA DO NEURÔNIO " $i$ " DA CAMADA " $(l-1)$ " AO NEURÔNIO " $j$ " DA CAMADA " $l$ ".

↓  
ENTRADA  
↓  
NEURÔNIO

# Conjunto de dados



# Algoritmo- Erros

- ERRO NO NEURÔNIO DE SAÍDA  $j$  NA ITERAÇÃO  $m$  (APRESENTAÇÃO DO  $m$ -ÉSIMO EXEMPLO DE TREINAMENTO):

$$l_j^{(L)}(m) = d_j(m) - y_j^{(L)}(m), \quad j \in L$$

- ERRO INSTANTÂNEO TOTAL DAS SAÍDAS :

$$E(m) = \frac{1}{2} \sum_{j=1}^{m_L} \left[ l_j^{(L)}(m) \right]^2$$

- ERRO MÉDIO QUADRÁTICO :

$$E_{\text{avr}} = \frac{1}{N} \sum_{m=1}^N E(m)$$

$N$  = NÚMERO TOTAL DE PADRÕES CONTIDOS NO CONS. DE TREINAMENTO.

$E_{\text{avr}}$  É UMA MEDIDA DO DESEMPENHO DO APRENDIZADO  
↳ FUNÇÃO CUSTO.

# Algoritmo- objetivo

OBJETIVO: SINTONIZAR  $w_{ji}^{(l)}$  PARA OBTEN  $\min \{E_{ov}\}$

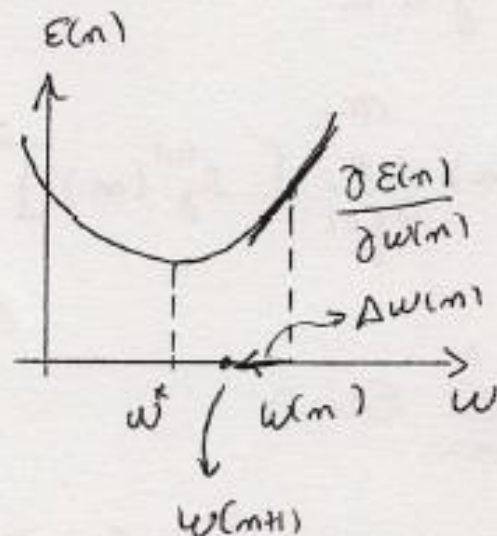
HIPÓTESE: VAMOS CONSIDERAR INICIALMENTE QUE OS PESOS  $w_{ji}^{(l)}$  SERÃO ATUALIZADOS ( $\Delta w_{ji}^{(l)}$ ) A CADA APRESENTAÇÃO DE UM PADRÃO DE TREINAMENTO (MODO "PADRÃO A PADRÃO" OU SEQUÊNCIAL), ATÉ O TÉRMINO DE UMA "ÉPOCA", ISTO É, ATÉ A APRESENTAÇÃO COMPLETA DE TODOS OS PADRÕES DO CONJ. DE TREINAMENTO.

$$w_{ji}^{(l)}(m+1) = w_{ji}^{(l)}(m) + \Delta w_{ji}^{(l)}(m)$$

$\Delta w_{ji}^{(l)}(m) \rightarrow$  ATUALIZADOS DE MANEIRA A MINIMIZAR  $E(m)$  p/ CADA PADRÃO  $\Rightarrow$  MINIMIZAR  $E_{ov}$ .

# Algoritmo - método

MÉTODO USADO: GRADIENTE DESCENDENTE



$$\Delta w(m) = -\eta \frac{\partial E(m)}{\partial w(m)}$$

$$\Delta w(m) = -\eta \nabla_w E(m)$$

ENTÃO:

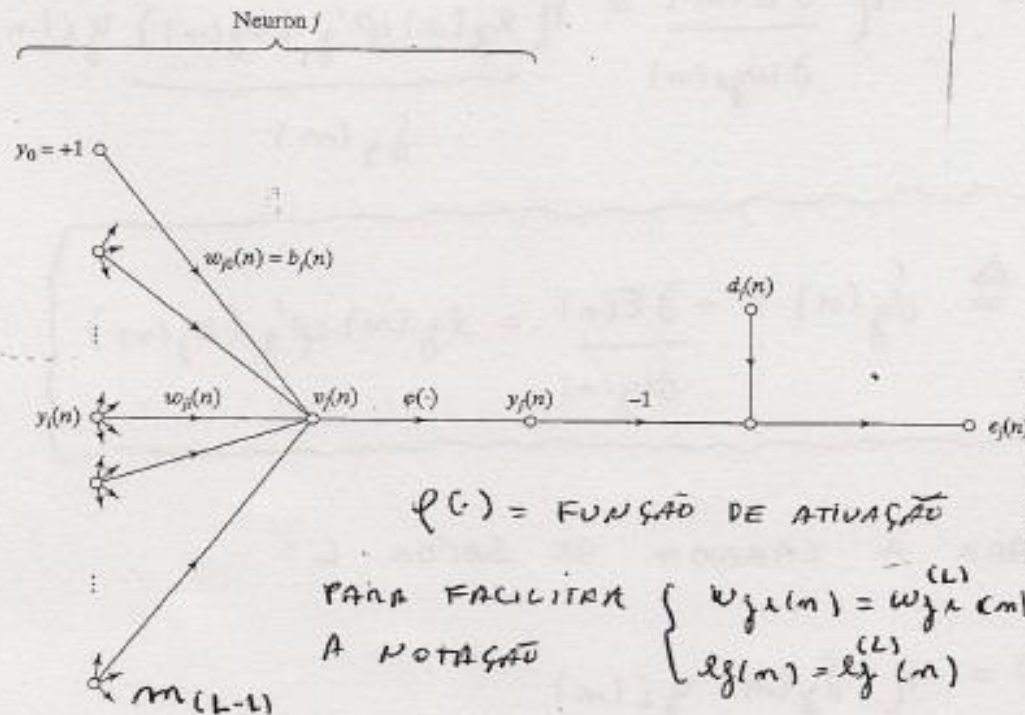
$$\Delta w_{ji}^{(l)}(m) = -\eta \frac{\partial E(m)}{\partial w_{ji}^{(l)}(m)}$$



# Caso1: neurônio camada saída

Caso 1

É um neurônio na camada de saída  $L$ .



$$v_j(n) = \sum_{i=0}^{m_{(L-1)}} w_{ji}(n) y_i(n) \quad , \quad y_j(n) = \phi_j(v_j(n))$$

# Caso1 (cont.)

REGRAS DA CADEIA:

$$\frac{\partial E(n)}{\partial w_{ji}(n)} = \frac{\partial E(n)}{\partial x_j(n)} \cdot \frac{\partial x_j(n)}{\partial y_j(n)} \cdot \frac{\partial y_j(n)}{\partial v_j(n)} \cdot \frac{\partial v_j(n)}{\partial w_{ji}(n)}$$

$\downarrow \quad \quad \downarrow \quad \quad \downarrow \quad \quad \downarrow$   
 $l_j(n) \quad -1 \quad \varphi'_j(v_j(n)) \quad y_i(n)$

$$\boxed{\frac{\partial E(n)}{\partial w_{ji}(n)} = -l_j(n) \cdot \varphi'_j(v_j(n)) y_i(n)}$$

NOTA QUE:

$$\frac{\partial E(n)}{\partial w_{ji}(n)} = \overbrace{\frac{\partial E(n)}{\partial v_j(n)}}^{\delta_j(n)} \cdot \frac{\partial v_j(n)}{\partial w_{ji}(n)} = -\delta_j(n) \cdot y_i(n)$$



$$\Delta w_{ji}(n) = -\eta \frac{\partial E(n)}{\partial w_{ji}(n)} = \eta \underbrace{x_j(n) \varphi'_j(v_j(n))}_{\delta_j(n)} y_i(n)$$

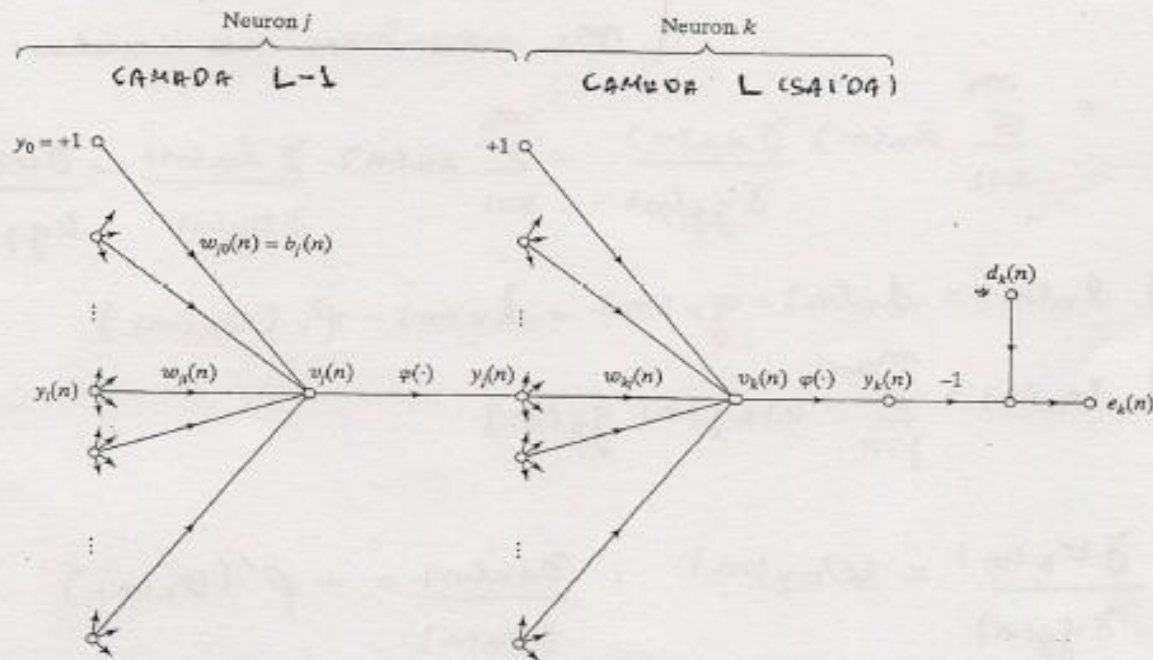
ENTÃO, PARA A CAMADA DE SAÍDA L:

$$\left\{ \begin{array}{l} \Delta w_{ji}(n) = \eta \delta_j(n) y_i(n) \\ \delta_j(n) = e_j(n) \varphi'(w_j(n)) \end{array} \right.$$

# Caso2: neurônio camada escondida

CASO 2: NEURÔNIO  $j$  ESTÁ EM UMA UNIDADE ESCONDIDA,

- QUAL É O SINAL DE ERRO PARA UMA UNIDADE ESCONDIDA?
- R: RECUSÃO



SIMPLIFICANDO A NOTAÇÃO: 
$$\begin{cases} w_{kj}(n) = w_{ji}^{(L)}(n) \\ w_{jk}(n) = w_{kj}^{(L-1)}(n) \end{cases}$$

## Caso 2: gradiente camada escondida

GRADIENTE LOCAL  $\delta_j(m)$  PARA O NEURÔNIO  $j$  (CAMADA ESCONDIRA)

$$\delta_j(m) = - \frac{\partial E(m)}{\partial v_j(m)} = - \frac{\partial E(m)}{\partial y_j(m)} \cdot \frac{\partial y_j(m)}{\partial v_j(m)}$$

$$\delta_j(m) = - \underbrace{\frac{\partial E(m)}{\partial y_j(m)}}_{\rightarrow ?} \cdot \varphi'_j(v_j(m)) \quad (1)$$

# Caso 2: cálculo gradiente

• CALCULANDO  $\partial E(n) / \partial y_j(n)$

$$E(n) = \frac{1}{2} \sum_{k=1}^{m_L} l_k^2(n) \quad \rightarrow \left\{ \begin{array}{l} \text{neurônio } k \text{ está na camada} \\ \text{de saída} \\ m_L \text{ neurônios na saída} \end{array} \right.$$

$$\frac{\partial E(n)}{\partial y_j(n)} = \sum_{k=1}^{m_L} l_k(n) \frac{\partial l_k(n)}{\partial y_j(n)} = \sum_{k=1}^{m_L} l_k(n) \cdot \frac{\partial l_k(n)}{\partial v_k(n)} \cdot \frac{\partial v_k(n)}{\partial y_j(n)}$$

$$\text{COMO } \left\{ \begin{array}{l} l_k(n) = d_k(n) - g_k(n) = d_k(n) - \varphi_k(v_k(n)) \\ v_k(n) = \sum_{j=0}^{m_{(L-1)}} w_{kj}(n) y_j(n) \end{array} \right.$$

$$\text{ENTÃO: } \frac{\partial v_k(n)}{\partial y_j(n)} = w_{kj}(n) \quad , \quad \frac{\partial l_k(n)}{\partial v_k(n)} = -\varphi'(v_k(n))$$



## Caso2: gradiente (cont.)

Logo:

$$\frac{\partial E(n)}{\partial y_j(n)} = - \sum_{k=1}^{m_L} \delta_k(n) \varphi'_k(u_k(n)) \cdot w_{kj}(n)$$

$$\frac{\partial E(n)}{\partial y_j(n)} = - \sum_{k=1}^{m_L} \underbrace{\delta_k(n)}_{\substack{\text{gradiente local da camada} \\ \text{L (SAIDA) NEURÔNIO K}}} \cdot w_{kj}(n)$$

COMO DE (1):

$$\delta_j(n) = - \frac{\partial E(n)}{\partial y_j(n)} \cdot \varphi'_j(u_j(n))$$

$$\delta_j(n) = \underbrace{\varphi'_j(u_j(n))}_{\substack{\downarrow \\ \text{DEPENDE SOMENTE} \\ \text{DO NEURÔNIO } j}} \cdot \underbrace{\sum_{k=1}^{m_L} \delta_k(n) w_{kj}(n)}_{\substack{\downarrow \\ \text{depende dos neurônios da} \\ \text{camada imediatamente à direita}}} \quad , j = \text{neurônio da camada L-L}$$

# Atualização pesos camada escondida

PODEMOS ESCRIVER:

$$\left\{ \begin{array}{l} \Delta w_{ji}(n) = \eta \delta_j(n) y_i(n) \quad (2) \\ \text{com} \\ \delta_j(n) = \varphi'_j(v_j(n)) \sum_{k=1}^{m_L} \delta_k(n) w_{kj}(n) \quad \text{p/ } j \text{ E CAMADA ESCONDIDA (L-1)} \\ \delta_j(n) = e_j(n) \varphi'(v_j(n)) \quad \text{p/ } j \text{ E CAMADA DE SAÍDA L} \end{array} \right.$$





# Retropropagação

SE O NEURÔNIO  $j$  ESTÁ EM UMA CAMADA ESCONDIDA  $l$ ,  $\delta_j^{(l)}$  PODE SER OBTIDO PELA SOMA DOS GRADIENTES LOCAIS DA CAMADA SEGUINTE  $(l+1)$ , DE MANEIRA RECURSIVA, CAMADA POR CAMADA, INICIANDO DA CAMADA DE SAÍDA  $L$  E CALCULANDO-SE ENTÃO A MUDANÇA NOS PESOS DADA POR (2) EM CADA CAMADA.

# Retropopagação com momento

MOMENTO ( MELHORA A ESTABILIDADE DO ALGORITMO )

$$\Delta w_{ji}(n) = \alpha \Delta w_{ji}(n-1) + \eta \delta_j(n) y_i(n)$$

↓  
FATOR DE  
MOMENTO

$$\begin{array}{l} 0 < \eta < 1 \\ 0 \leq \alpha < 1 \end{array}$$

↓  
TAXA DE APRENDIZAGEM

ou

$$\Delta w_{ji}(n) = \alpha \Delta w_{ji}(n-1) + \underbrace{\eta^1 (1-\alpha)}_{\eta} \delta_j(n) y_i(n)$$

# Resumo Algoritmo

MODO "PADRÃO A PADRÃO", SEQUENCIAL, OU ESTOCASTICO

Neste caso os pesos são atualizados após a apresentação de cada exemplo de treinamento

①- Inicializar  $w_{ji}^{(l)}$ ,  $m=0$ ,  $n, \alpha$

②- "RANDOMIZAR" A ORDEM DO CONJUNTO DE TREINAMENTO  $(x, d)$

# Resumo Algoritmo (cont.)

③ - PARA CADA PADRÃO DO CONJ. DE TREINAMENTO  $X$ :  
( $X = [x_1 \ x_2 \ \dots \ x_{m_0}]^T$ )

③.1 COMPUTAÇÃO DIRETA:

$$\text{CALCULAR } v_j^{(l)}(n) = \sum_{i=0}^{m_l} w_{ji}^{(l)}(n) y_i^{(l-1)}(n)$$

(começando da entrada e seguindo em direção à saída)

$$\text{OBS: } y_j^0(n) = x_j(n) \quad , \quad y_0^{(l)}(n) = +1$$

③.2 DETERMINAR O ERRO NA SAÍDA:  $\delta_j^{(l)}(n) = d_j^{(l)}(n) - y_j^{(l)}(n)$

③.3 RETROPROPAGAÇÃO: calcular o ~~pro~~ gradiente local

$$\delta_j^{(l)}(n) = \begin{cases} \delta_j^{(l)}(n) \varphi_j'(v_j^{(l)}(n)) & , \text{ se } j \in L \\ \varphi_j'(v_j^{(l)}(n)) \sum_{k=1}^{m_{l+1}} \delta_k^{(l+1)}(n) w_{kj}^{(l+1)}(n) & , \text{ se } j \notin L \end{cases}$$

③.4 AJUSTAR OS PESOS A CADA PADRÃO APRESENTADO:

$$w_{ji}^{(l)}(n+1) = w_{ji}^{(l)}(n) + \alpha [\Delta w_{ji}^{(l)}(n-1)] + \eta \delta_j^{(l)}(n) y_i^{(l-1)}(n)$$

$$\text{com } \Delta w_{ji}^{(l)}(n-1) = w_{ji}^{(l)}(n) - w_{ji}^{(l)}(n-1)$$

④ SE NÃO ATINGIR CRITÉRIO DE PARADA VÁ PARA ②,  
CASO CONTRÁRIO FIM.

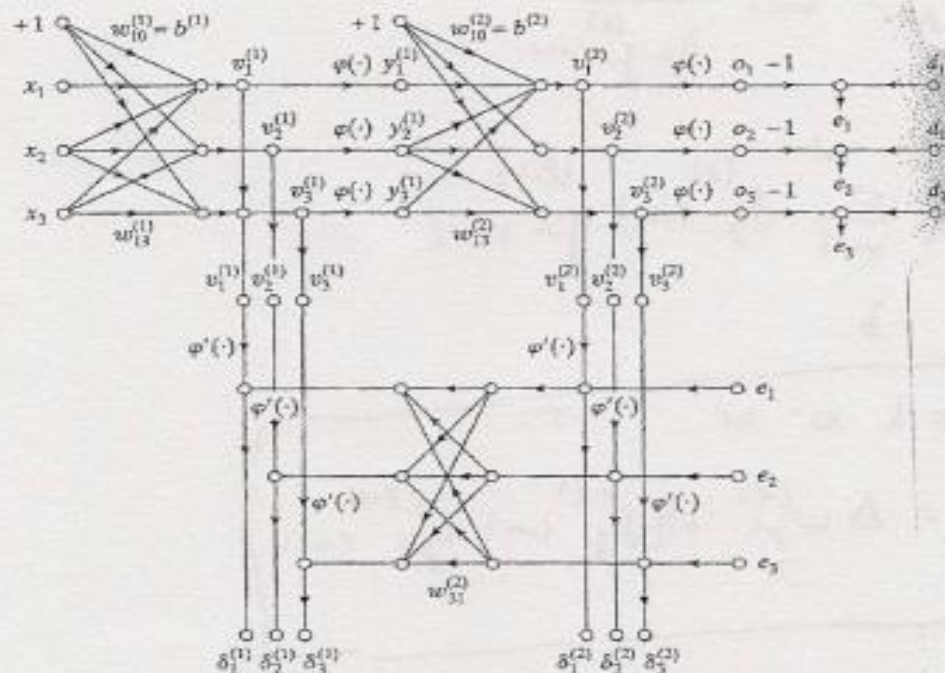


# Os dois passos do algoritmo

## DIAGRAMA DE FLUXO DE SINAL

DOIS PASSOS:

- COMPUTAÇÃO DIRETA
- RETROPROPAGAÇÃO DOS ERROS



# Modo Batelada (*Batch*)

- A conexão dos pesos é realizada após a apresentação de todos os  $N$  padrões de treinamento.

$$\Delta w_{ji}^{(l)} = -\eta' \frac{\partial E_{av}}{\partial w_{ji}^{(l)}} =$$

$$\Delta w_{ji}^{(l)} = \left( \frac{\eta'}{N} \right) \sum_{n=1}^N \frac{\partial E_{lm}}{\partial w_{ji}^{(l)}} =$$

$$\Delta w_{ji}^{(l)} = \eta \sum_{n=1}^N \delta_j^{(l)} y_{ji}^{(l-1)}$$

↓

PARA  $m = 1 \text{ a } N$

$$\Delta w_{ji}^{(l)} = \Delta w_{ji}^{(l)} + \eta \delta_j^{(l)} y_{ji}^{(l-1)}$$



# Algoritmo modo batelada

$p$  = número da época;

$\Delta W_{pji}^{(l)}(p)$  = variação do peso no fim da época  $p$ ;

$W_{ji}^{(l)}(p)$  = peso na época  $p$ .

- ① INICIALIZAÇÃO :  $p=0$  ,  $\Delta W_{pji}^{(l)}(p)=0$  ,  $\eta, \alpha$   
 $W_{ji}^{(l)}(p)$  = valor aleatório  $[-1,1]$

# Algoritmo modo batelada (cont.)

- ② Para cada apresentação  $p$  de uma nova época
- ②.1  $p = p + 1$  ; incremento número de épocas
- ②.2  $\Delta w_{ji}^{(l)} = 0$  ; zero acumulador de variações dos pesos na época

- ②.3 Para cada padrão  $n$  do conj. de treinamento

- ②.3.1 Computação direta

Calcular  $v_j^{(l)}(n)$  obtendo a saída da rede

- ②.3.2 Determinar  $\delta_j^{(l)}$  e calcular

$$\Delta w_{ji}^{(l)} = \Delta w_{ji}^{(l)} + \eta \delta_j^{(l)} y_i^{(l-1)}$$

- ③ ATUALIZA OS PESOS A CADA ÉPOCA APRESENTADA

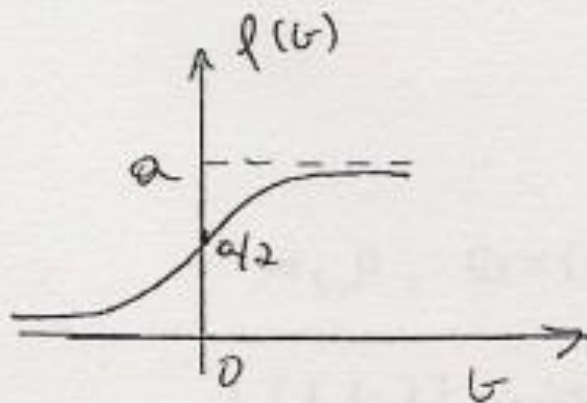
$$w_{ji}^{(l)}(p+1) = w_{ji}^{(l)}(p) + \Delta w_{ji}^{(l)} + \alpha \Delta w_{epji}^{(l)}(p-1)$$

$$\Delta w_{epji}^{(l)}(p) = w_{ji}^{(l)}(p+1) - w_{ji}^{(l)}(p)$$

- ④ SE NÃO ATINGIR CRITÉRIO DE PARADA VÁ PARA ②, CASO CONTRÁRIO FIM.

# Funções de Ativação

## 1) SIGMOIDAL LOGÍSTICA



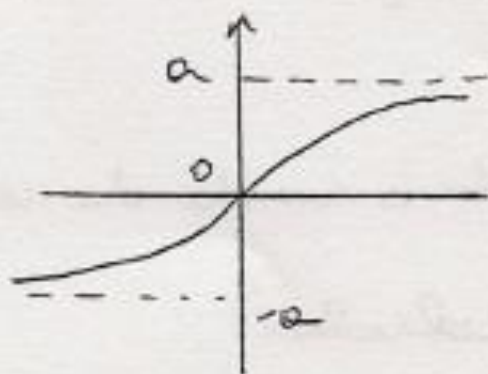
$$\varphi(v) = \frac{a}{1 + e^{-bv}}$$

$$\varphi'(v) = \frac{b}{a} \varphi(v) [a - \varphi(v)]$$

$$\varphi'_j(v_j(m)) = \frac{b}{a} \varphi_j(m) [a - \varphi_j(m)]$$

## Funções de Ativação (cont.)

2) SIGMOIDAL TANGENTE



$$\varphi(v) = a \tanh(bv)$$

$$\varphi(v) = a \left[ \frac{1 - e^{-bv}}{1 + e^{-bv}} \right]$$

$$\varphi'(v) = \frac{b}{a} [a^2 - \varphi^2(v)]$$

$$\varphi'_j(v_j(m)) = \frac{b}{a} [a^2 - \varphi_j^2(m)]$$

# Critérios de parada

•  $p$  = número de épocas

① VECTOR GRADIENTE ( $\|\cdot\| \rightarrow$  norma euclidiana)

$$\|\nabla_w E(w(p))\| \leq \theta$$

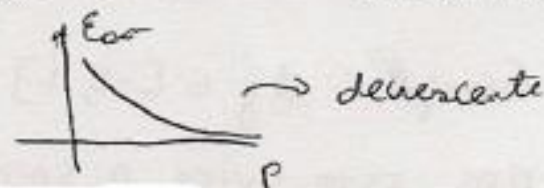
$$\|\Delta w_{ji}^{(e)}(p)\| \leq \theta$$

} PARA MODO BATELADA

② VALOR DO ERRO MÉDIO QUADRÁTICO

$$E_{\text{err}}(w(p)) \leq \varepsilon$$

(PODE LEVAR A "OVERFITTING")





# Critérios de parada (cont.)

③

VARIACÃO DO ERRO QUADRÁTICO MÉDIO

• VARIACÃO ABSOLUTA

$$|\Delta E_{\text{er}}(w)| \leq \lambda$$

$$\Rightarrow E_{\text{er}}(w(p)) - E_{\text{er}}(w(p-1)) \geq -\lambda$$

• VARIACÃO RELATIVA

$$\frac{|\Delta E_{\text{er}}(w(p))|}{E_{\text{er}}(w(p))} \leq \lambda \Rightarrow \frac{E_{\text{er}}(w(p)) - E_{\text{er}}(w(p-1))}{E_{\text{er}}(w(p-1))} \geq -\lambda$$

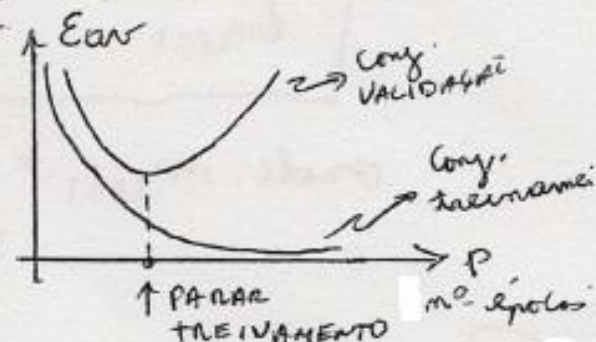
④

PARADA ANTECIPADA (MELHOR!)

"EARLY STOPPING"



EVITA "OVERFITTING"







# Técnicas para melhorar o algoritmo de retropropagação – Função de ativação

---

## 1) FUNÇÃO DE ATIVAÇÃO

Usar sigmoideal tangente (leva a um aprendizado mais rápido) com  $a = 1,7159$ ,  $b = 2/3$

## Técnicas para melhorar o algoritmo de retropropagação – Valor das saídas

2) VALOR DAS SAÍDAS (QDO NÃO USANDO NEURÔNIO LINEAR NA SAÍDA).

SUPONDO FUNÇÃO NÃO LINEAR COM SATURAÇÃO EM  $[-a, +a]$

$$\begin{cases} +a \Rightarrow dy = a - \epsilon \\ -a \Rightarrow dy = -a + \epsilon \end{cases}$$

Exemplo:  $dy \in [-1, 1]$  quando  $a = 1,715$

OBS: ISTO EVITA A SATURAÇÃO DOS NEURÔNIOS DAS CAMADAS ESCONDIDAS, O QUE ACABARIA UMA DIMINUIÇÃO NA VELOCIDADE DO APRENDIZADO.

## Técnicas para melhorar o algoritmo de retropropagação – Inicialização dos pesos

### 3) INICIALIZAÇÃO DOS PESOS

Para aprendizado uniforme e rápido, escolhemos aleatoriamente o valor inicial dos pesos de uma distribuição uniforme com:

$$\frac{-1}{\sqrt{m_{(l-1)}}} \leq w_{ji}^l(0) \leq \frac{1}{\sqrt{m_{(l-1)}}}$$

onde  $m_{(l-1)}$  = n.º de neurónios da ~~me~~ camada anterior à camada  $l$

## Técnicas para melhorar o algoritmo de retropropagação – Normalizar entradas

### ④ NORMALIZAR OS DADOS DE ENTRADA

Cada componente do vetor de entrada deve ser normalizado para média nula e variância unitária

$$X = [x_1 - x_1 \dots x_m]^T$$

$$x_{i, \text{norm}} = \frac{x_i - \mu_i}{\sigma_i}$$

OBS: É interessante que os  $x_i$  sejam descorrelacionados para acelerar a aprendizagem.

## Técnicas para melhorar o algoritmo de retropropagação – Taxa de aprendizagem

### ⑤ TAXA DE APRENDIZAGEM

- $\eta = 0,1$  COMO PRIMEIRA TENTATIVA
- Diminua-se  $\eta$  se a função custo  $E$  oscilar ou oscila durante a aprendizagem.
- $\uparrow \eta$  se a aprendizagem parece inadvertidamente ter parado.
- $\eta$  pode variar durante a aprendizagem.





# Técnicas para melhorar o algoritmo de retropropagação – momento

⑥

MOMENTO

$\alpha \approx 0,9$  TÍPICO.

⑦

TREINAMENTO SEQUENCIAL, DATELADA OU "ON LINE"

ON LINE  $\Rightarrow$  CONT. TREINAMENTO GRANDE

DATELADA É TÍPICAMENTE MAIS VAGAROSO QUE  
SEQUENCIAL





# Conjunto de dados

---

- Técnicas de pré-processamento são frequentemente utilizadas para tornar o conjunto de dados mais adequados para uso de algoritmos de aprendizagem:
  - **Eliminação manual de atributos:** quando um atributo de entrada não contribui para a estimativa do atributo alvo (saída) ele é considerado irrelevante (ex: valor igual para todos os padrões/objetos);
  - **Integração de dados:** dados a serem utilizados em diferentes conjuntos;
  - **Amostragem de dados:** a amostragem de dados deve ser representativa dos conjunto de dados original;
  - **Balanceamento de dados:** evitar que o número de objetos de uma dada classe seja muito diferente das demais
  - **Limpeza de dados:** dados inconsistentes, dados redundantes, dados com ruídos;
  - **Redução de dimensionalidade:** agregação (PCA), seleção;
  - **Transformação de dados:** ex: normalizaçãotransformação de atributo simbólicos.