

Redes Neurais Feedforward

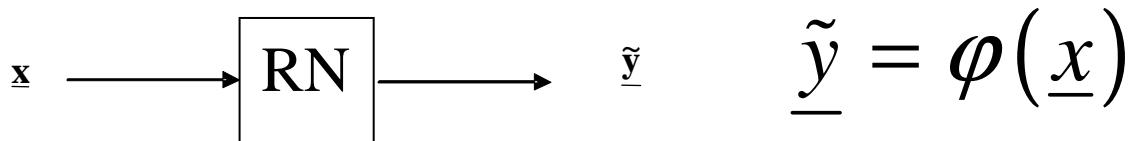
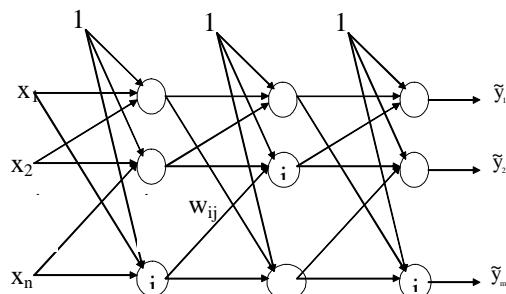
Aprendizado (Treinamento)

Backpropagation

Aproximador Universal

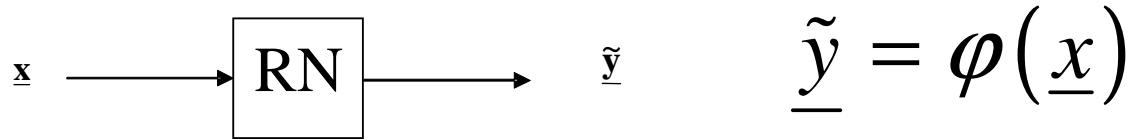
Redes Neurais *FeedForward*

Redes “*Backpropagation*”



Mapeador não linear

Aproximador Universal !

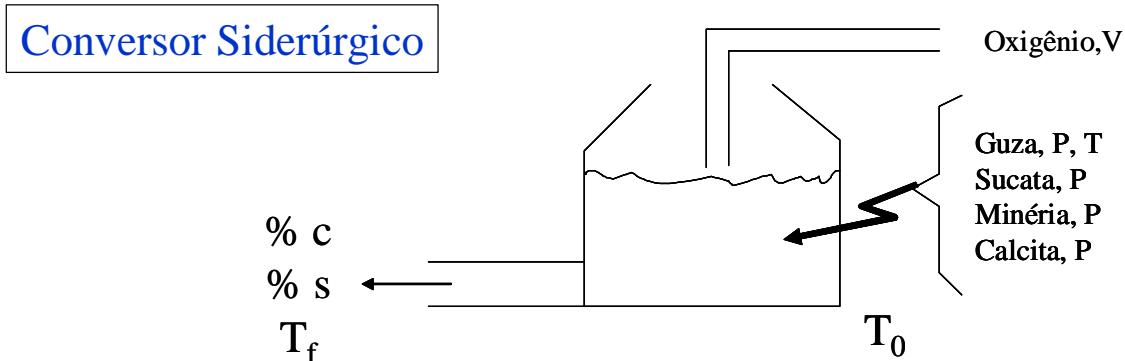


Para que serve ?

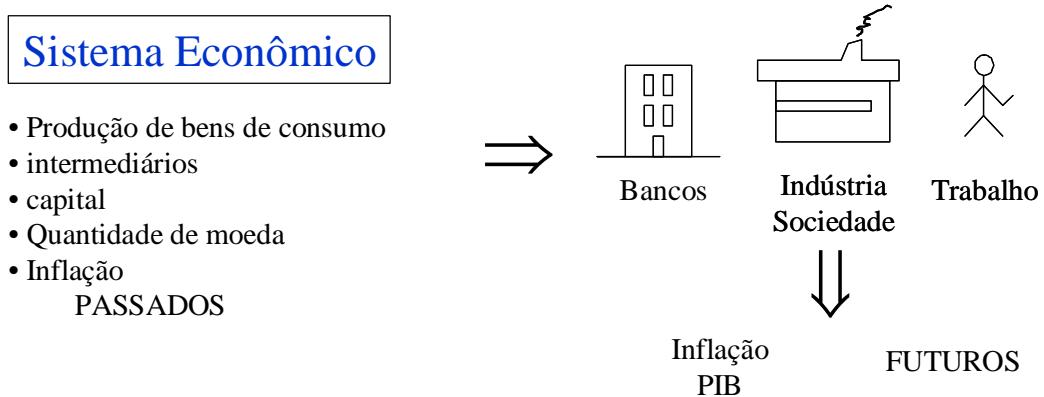
Aplicações:

Simuladores;
Controladores;
Classificadores;
Reconhecimento de padrões;
Filtragem não linear;
etc...

Simulação de Sistemas



Simulação de Sistemas



Simulação de Sistemas

Sistema Real, Planta

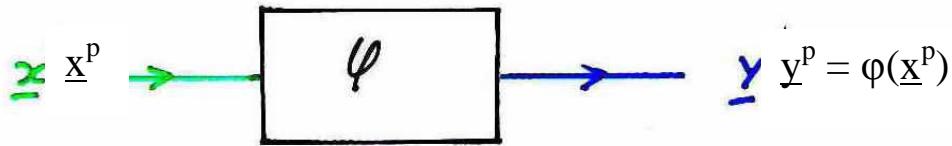
$$(\underline{\mathbf{x}}^p, \underline{\mathbf{y}}^p) \quad p = 1, 2, \dots, P$$



Simulador



Sistema à emular



Conhecimento do Sistema

Fenomenológico – equações – modelo caixa branca - NÃO

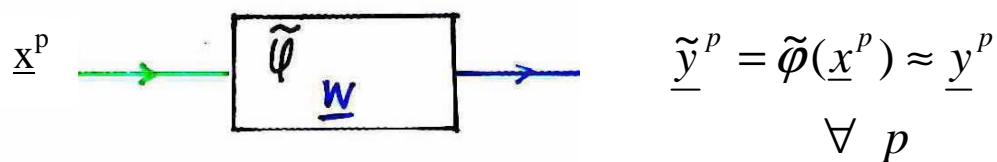
Funcional – relação entrada/saída – modelo caixa preta

Pares entrada-saída

$$(\underline{x}^p, \underline{y}^p) \quad p = 1, \dots P$$

Modelo Matemático Fenomenológico / Numérico -

Simulador



Modelo Fenomenológico – Caixa branca

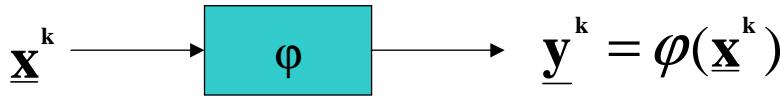
Modelo Numérico Geral (e.g. rede neural) – Caixa Preta

Ajuste dos parâmetros w ?

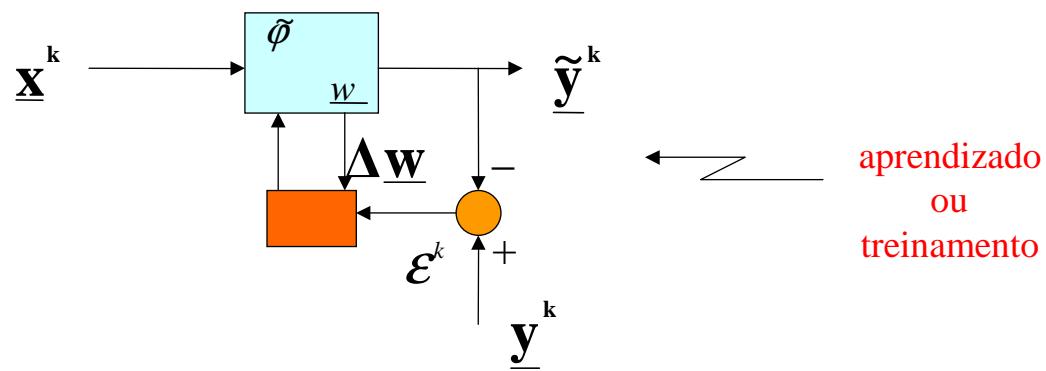
Simulação de Sistemas

Sistema Real, Planta

$$(\underline{\mathbf{x}}^k, \underline{\mathbf{y}}^k) \quad k=1,2,\dots,P$$



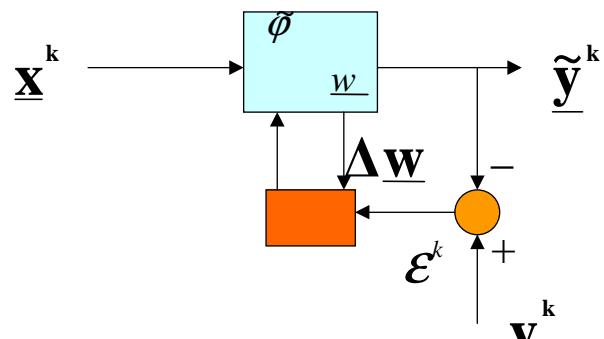
Simulador



Aprendizado

(ou treinamento, ou ajuste de parâmetros)

como um processo de otimização



Aprendizado =
minimização do erro na saída

Erro a minimizar:

$$\varepsilon^{2^k} = \left\| \underline{\mathbf{y}}^k - \tilde{\underline{\mathbf{y}}}^k \right\|^2 = \sum_{l=1}^m (\underline{\mathbf{y}}_l^k - \tilde{\underline{\mathbf{y}}}_l^k)^2$$

Erro médio quadrático:

$$F_0 = E(\varepsilon^{2^k}) = \frac{1}{P} \sum_{k=1}^P \varepsilon^{2^k}$$

$$F_0 = E(\varepsilon^{2^k}) = F_0(\underline{\mathbf{w}}) \geq 0$$

Treinamento, Aprendizado: Minimizar o erro na saída

Função objetivo à minimizar:

$$F_0 = E(\varepsilon^{2^k}) = F_0(\underline{\mathbf{w}}) \geq 0$$

Problema:

Como encontrar o vetor $\underline{\mathbf{w}}_0$ que minimize $F_0(\underline{\mathbf{w}})$?

$$F_0(\underline{\mathbf{w}}_0) \leq F_0(\underline{\mathbf{w}}) \quad \forall \underline{\mathbf{w}} \neq \underline{\mathbf{w}}_0$$

Gradiente:

$$F_0(w_1, w_2, \dots) = F(\underline{w})$$

$$\underline{w} = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \end{bmatrix} \quad \nabla_{\underline{w}} F_0 = \underline{\nabla} = \begin{bmatrix} \frac{\partial F_0}{\partial w_1} \\ \frac{\partial F_0}{\partial w_2} \\ \vdots \end{bmatrix}$$

Propriedades do Gradiente:

1. $\underline{w} \rightarrow \underline{w} + \Delta \underline{w} \Rightarrow F_0 \rightarrow F_0 + \Delta F_0$

se $\Delta \underline{w} = \alpha \underline{\nabla} \Rightarrow \Delta F_0$ é máximo

então se $\Delta \underline{w} = -\alpha \underline{\nabla} \Rightarrow \Delta F_0$ é mínimo

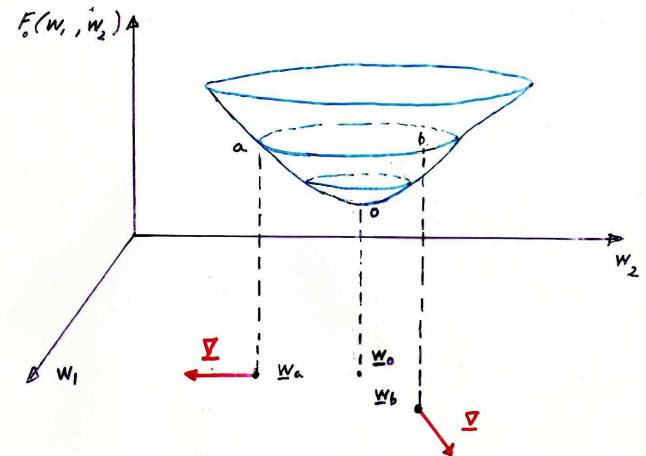
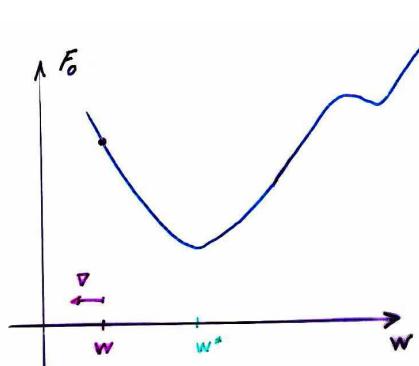
i.e. $\Delta F_0 < 0$ e $|\Delta F_0|$ é máximo

2. $\underline{\nabla}|_{\underline{w}_0} = \underline{0}$

Gradiente descendente

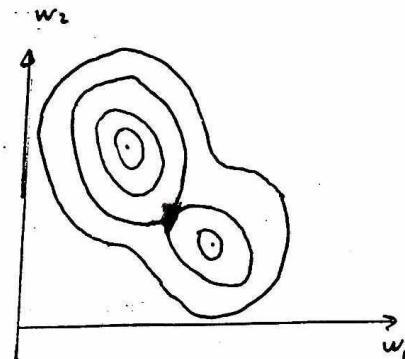
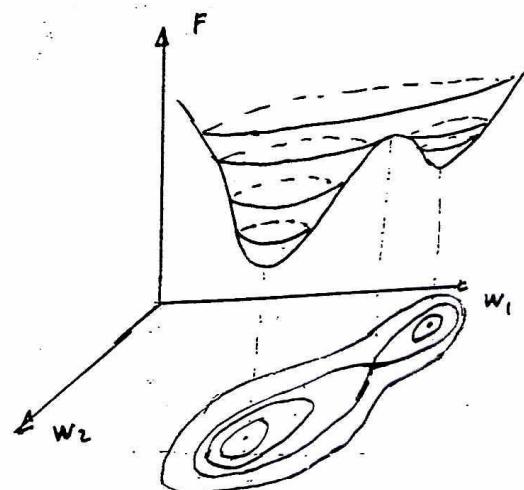
ou descida por gradiente

$$\underline{w} \longrightarrow \Delta \underline{w} = -\alpha \nabla(\underline{w}) \longrightarrow \underline{w} = \underline{w} + \Delta \underline{w}$$

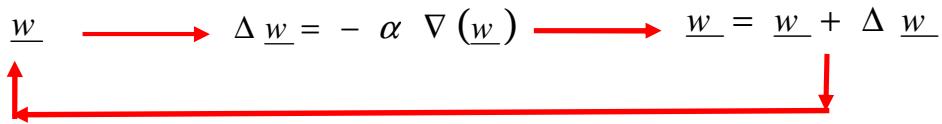


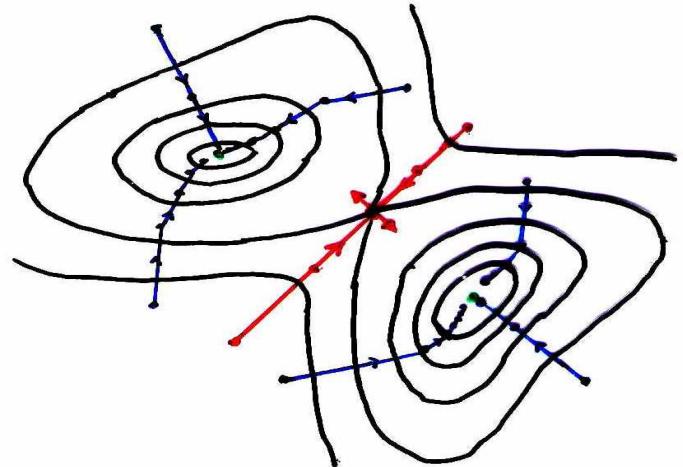
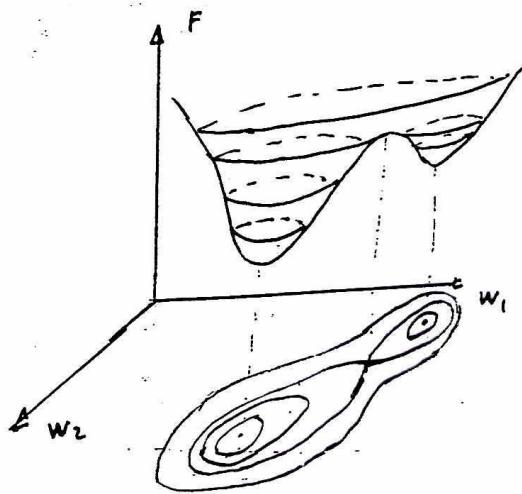
Visualização

por Curvas de Nível



Gradiente descendente ou descida por gradiente

$$\underline{w} \longrightarrow \Delta \underline{w} = -\alpha \nabla (\underline{w}) \longrightarrow \underline{w} = \underline{w} + \Delta \underline{w}$$




sela instável mínimos locais

Complexidade de cálculo:

$$\Delta w_{ij} = -\alpha \frac{\partial F}{\partial w_{ij}}$$

Fim do processo:

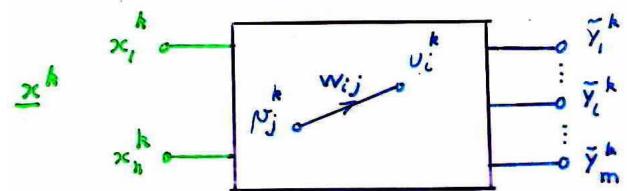
$$\nabla \underline{w}_{otimo} = \underline{0} \quad \rightarrow \quad \Delta \underline{w}_{otimo} = \underline{0}$$

Acréscimo a aplicar em cada sinapse:

$$F_0 = E(\varepsilon^{2^p}) = F_0(\underline{\mathbf{w}}) \quad \Delta w_{ij} = -\alpha \frac{\partial F_0(w_{ij})}{\partial w_{ij}}$$

Como calcular

$$\frac{\partial F_0(w_{ij})}{\partial w_{ij}} \quad ???$$



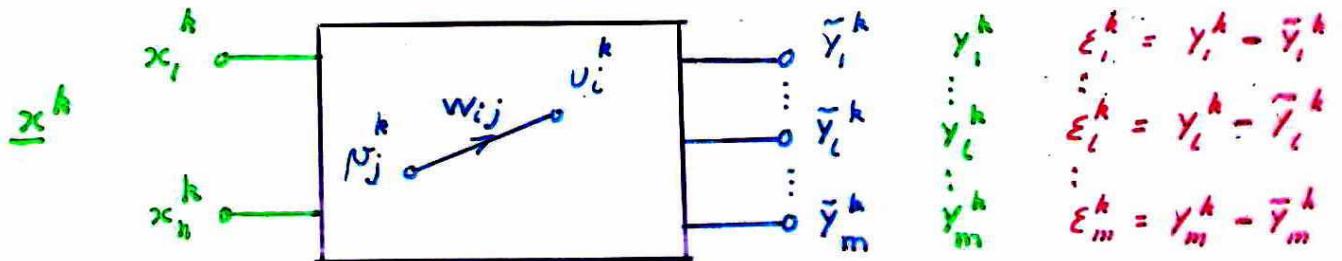
$$1 - F_0(w_{ij}) = E_k \left[(\varepsilon^k)^2 \right] = \frac{1}{K} \sum_{k=1}^K (\varepsilon^k)^2$$

$$\frac{\partial F_0(w_{ij})}{\partial w_{ij}} = \frac{\partial E_k [(\varepsilon^k)^2]}{\partial w_{ij}} = E_k \left[\frac{\partial (\varepsilon^k)^2}{\partial w_{ij}} \right]$$

Propriedade importante:

o gradiente do valor esperado do erro quadrático é igual ao valor esperado do gradiente do erro quadrático
cada par entrada-saída é tratado isoladamente.

2 – Cálculo de $\frac{\partial(\varepsilon^k)^2}{\partial w_{ij}}$ para cada par entrada – saída k

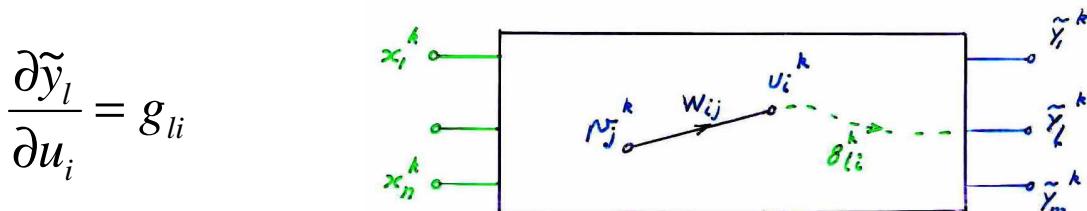


$$(\varepsilon)^2 = \sum_{l=1}^m (\varepsilon_l)^2 = \sum_{l=1}^m (y_l - \tilde{y}_l)^2$$

$$\frac{\partial(\varepsilon)^2}{\partial w_{ij}} \quad ???$$

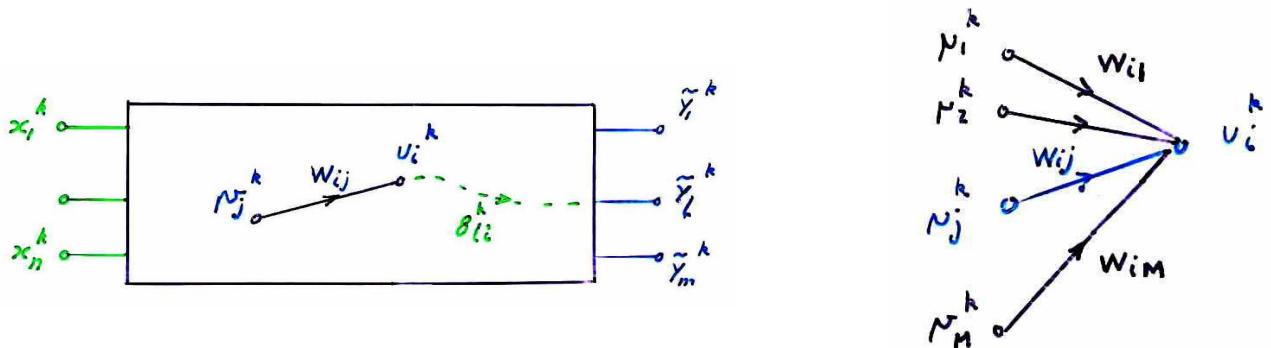
$$\frac{\partial(\varepsilon)^2}{\partial w_{ij}} = \frac{\partial}{\partial w_{ij}} \sum_{l=1}^m \varepsilon_l^2 = \sum_{l=1}^m \frac{\partial \varepsilon_l^2}{\partial w_{ij}} = 2 \sum_{l=1}^m \varepsilon_l \frac{\partial \varepsilon_l}{\partial w_{ij}} =$$

$$= 2 \sum_{l=1}^m \varepsilon_l \frac{\partial (y_l - \tilde{y}_l)}{\partial w_{ij}} = -2 \sum_{l=1}^m \varepsilon_l \frac{\partial \tilde{y}_l}{\partial w_{ij}} = -2 \sum_{l=1}^m \varepsilon_l \frac{\partial \tilde{y}_l}{\partial u_i} \frac{\partial u_i}{\partial w_{ij}}$$



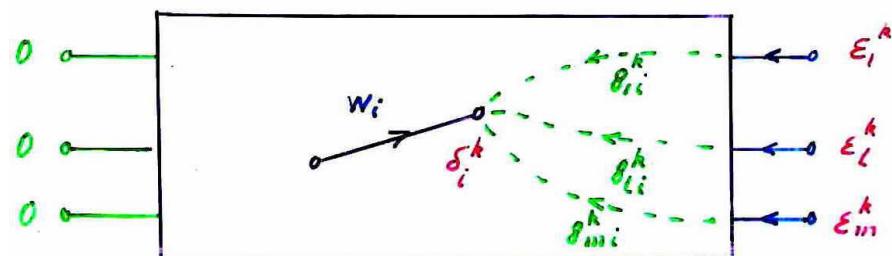
$$\frac{\partial(\varepsilon)^2}{\partial w_{ij}} = -2 \sum_{l=1}^m \varepsilon_l \frac{\partial \tilde{y}_l}{\partial u_i} \frac{\partial u_i}{\partial w_{ij}} = -2 \sum_{l=1}^m \varepsilon_l g_{li} v_j$$

$$\frac{\partial \tilde{y}_l}{\partial u_i} = g_{li} \quad u_i = \sum_{m=1}^M w_{mj} v_j \quad \frac{\partial u_i}{\partial w_{ij}} = v_j$$



$$\frac{\partial(\varepsilon)^2}{\partial w_{ij}} = -2 \sum_{l=1}^m \varepsilon_l g_{li} v_j = -2 v_j \sum_{l=1}^m \varepsilon_l g_{li} = -2 v_j \delta_i$$

$$\delta_i = \sum_{l=1}^m \varepsilon_l \frac{\partial \tilde{y}_l}{\partial u_i} = \sum_{l=1}^m \varepsilon_l g_{li}$$



δ_i é o erro retropropagado da saída até a extremidade da sinapse

3 - como:

$$\Delta w_{ij} = -\alpha \frac{\partial F_0(w_{ij})}{\partial w_{ij}}$$

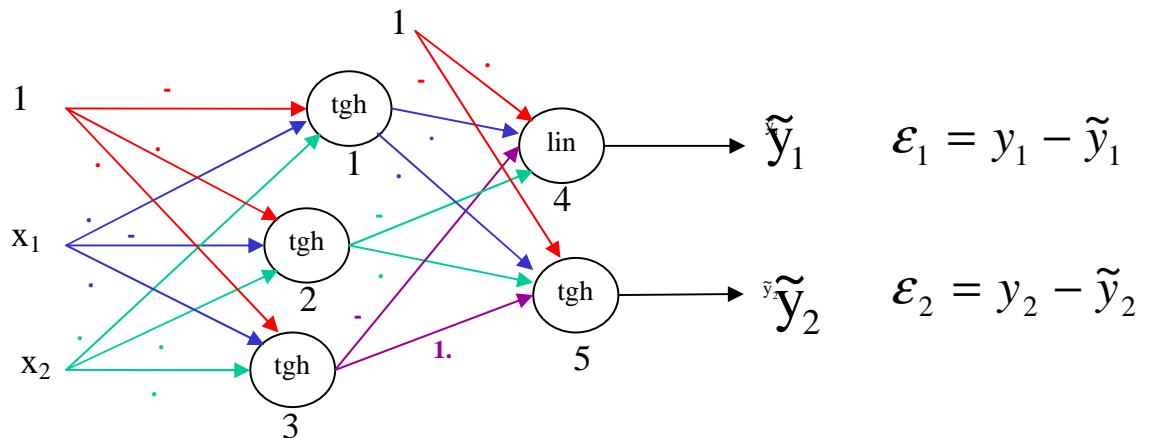
$$\frac{\partial F_0(w_{ij})}{\partial w_{ij}} = \frac{1}{P} \sum_{p=1}^P \frac{\partial}{\partial w_{ij}} (\epsilon^{2^p})$$

$$\frac{\partial (\epsilon)^2}{\partial w_{ij}} = -2 v_j \delta_i$$

$$\boxed{\Delta w_{ij} = 2 \alpha \frac{1}{P} \sum_{p=1}^P v_j \delta_i|_p}$$

Error Backpropagation -Princípio do Algoritmo:

1 - Rede Original:



Propagar o sinal na rede original e conservar o valor do sinal v_j na entrada de cada sinapse w_{ij} . Calcular o erro em cada saída.

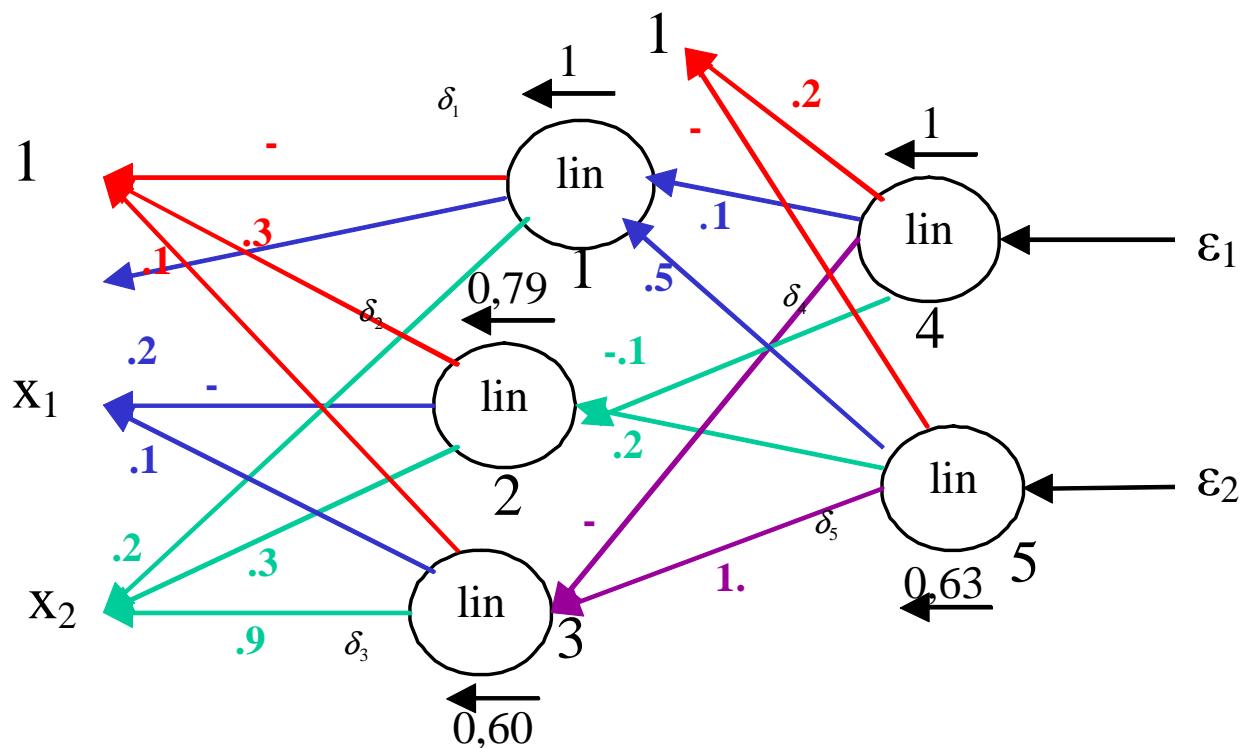
2 – Criar uma “Rede Associada” a partir da Rede Original

mesma arquitetura, mas

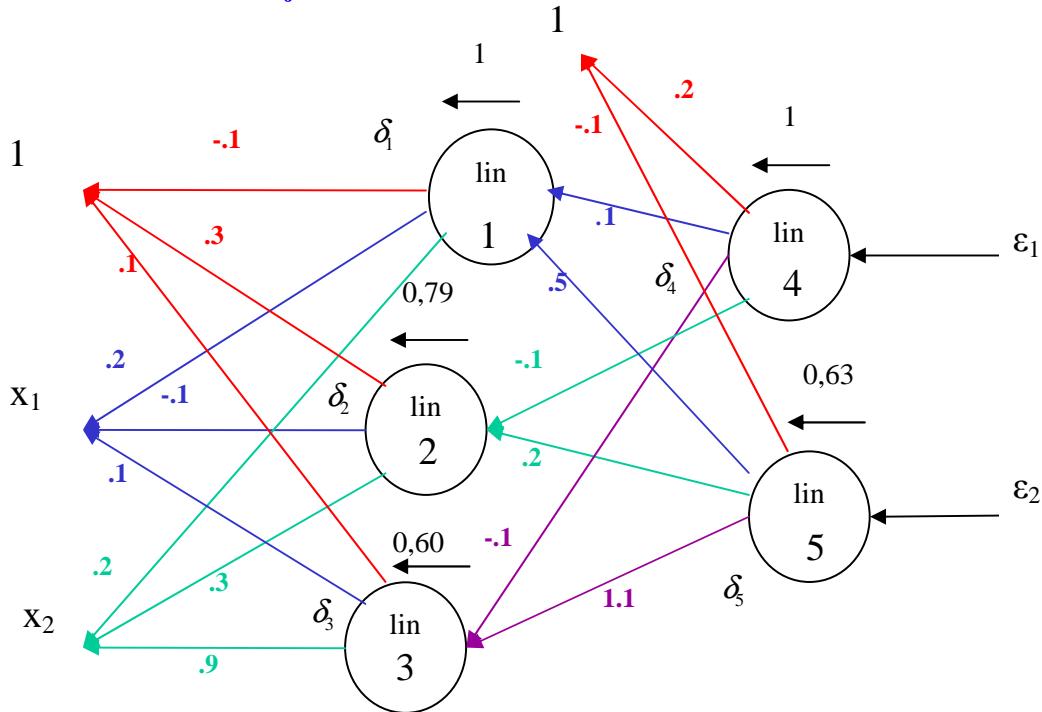
2.1 - “Linearizar” (os neurônios) da rede original

Rede Original	\rightarrow	Rede Associada
Neurônio não linear $v_0 = \tanh(u_0)$	\rightarrow	Neurônio linear $v = (1-v_0^2) u$
Neurônio linear $v_0 = u_0$	\rightarrow	Neurônio linear $v = u$

2.2- Inverter todos os sentidos de transmissão (dos neurônios e sinapses)



3 – (Retro)propagar o erro ε_i em cada saída através da rede “associada” e conservar o valor do sinal δ_i de erro retropropagado que chega na saída (original) de cada sinapse w_{ij} .



4 – O acréscimo para cada sinapse w_{ij} para o par p entrada-saída em operação é dado por:

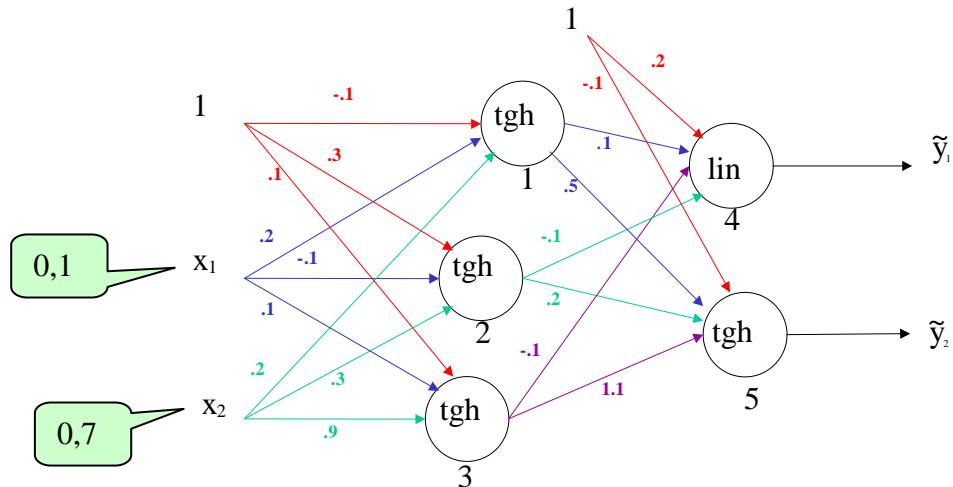
$$\Delta_p w_{ij} = 2 \alpha v_j \delta_i$$

5 – O acréscimo a ser aplicado na sinapse w_{ij} é o valor esperado dos acréscimos calculados para todos os pares entrada-saída.

$$\Delta w_{ij} = E_p (\Delta_p w_{ij}) = 2 \alpha E_p (v_j \delta_i)$$

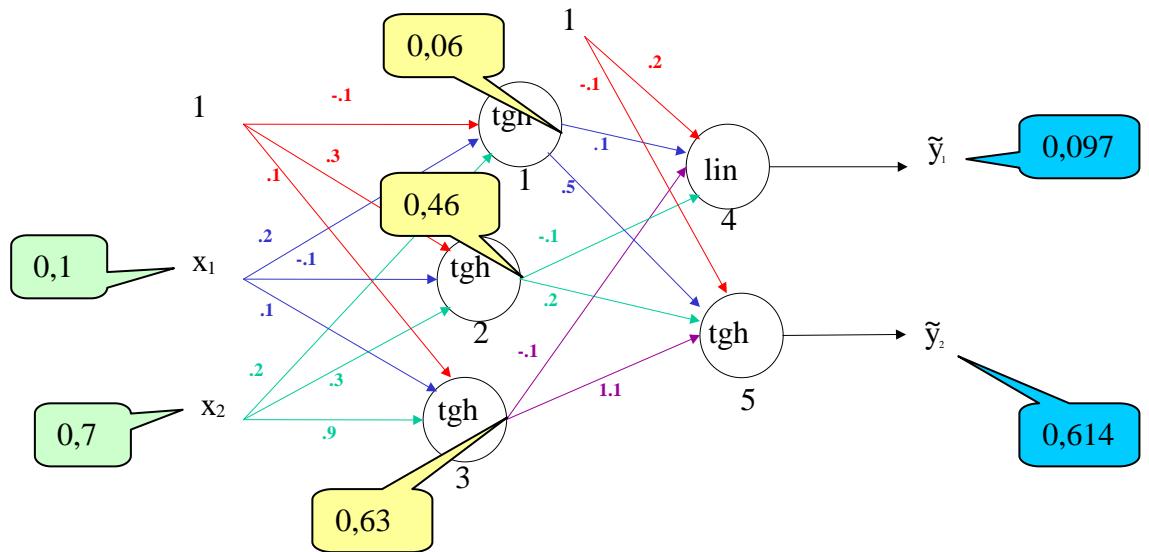
Processo em Batelada – Batch

Exemplo anterior:



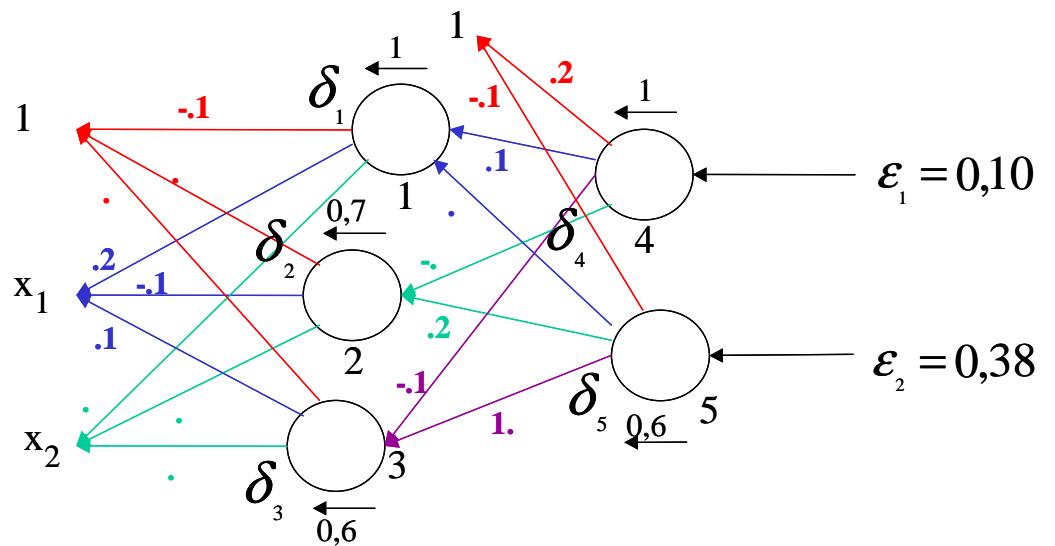
$$\underline{x} = \begin{bmatrix} 0,1 \\ 0,7 \end{bmatrix} \quad \tilde{\underline{y}} = ? \quad \underline{y} = \begin{bmatrix} 0,2 \\ 0,1 \end{bmatrix}$$

Signal Feedforward



$$\underline{x} = \begin{bmatrix} 0,1 \\ 0,7 \end{bmatrix} \quad \underline{y} = \begin{bmatrix} 0,2 \\ 1 \end{bmatrix} \quad \tilde{\underline{y}} = \begin{bmatrix} 0,097 \\ 0,614 \end{bmatrix} \quad \underline{\epsilon} = \begin{bmatrix} 0,103 \\ 0,386 \end{bmatrix}$$

Rede associada e retropropagação do erro:



$$\delta_5 = (0,386)(0,63) = 0,24 \quad \delta_4 = (0,103)(1) = 0,103$$

$$\delta_3 = [(0,103)(-0,1) + (0,24)(1,1)](0,60) = 0,152$$

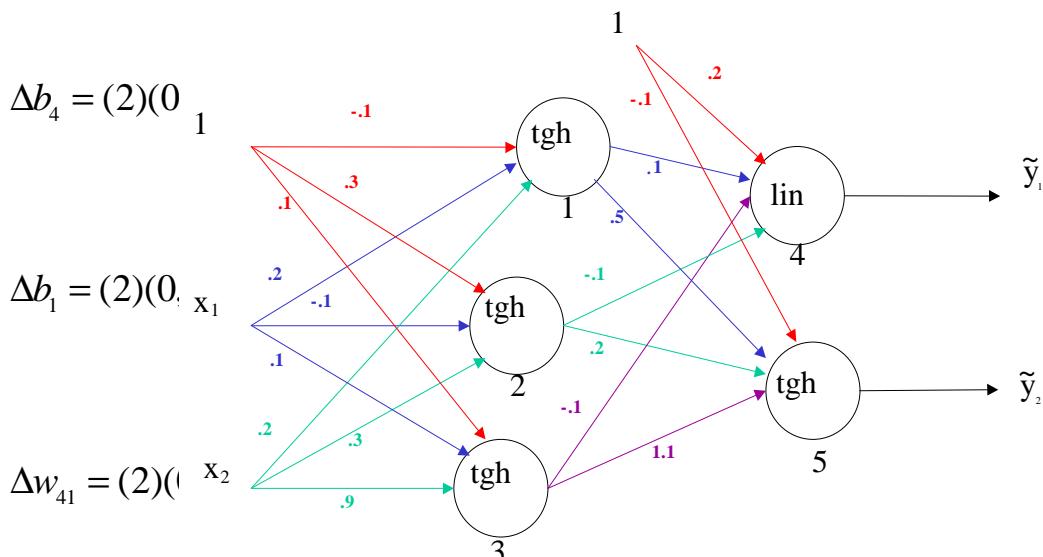
$$\delta_2 = [(0,103)(-0,1) + (0,24)(0,2)](0,79) = 0,030$$

$$\delta_1 = [(0,103)(0,1) + (0,24)(0,5)](1) = 0,130$$

Treinamento Regra Delta

Atualização dos valores das sinapses:

$$\Delta w_{ij} = 2\alpha v_j \delta_i$$



$$\Delta w_{3b} = (2)(0,1)(0,7)(0,152) = 0,021$$

$$w_{3b} \rightarrow 0,9 + 0,021 = 0,921$$

Algorítmos BP para cálculo do acréscimo na sinapse w_{ij} devido ao par p , $\Delta_p w_{ij}$

para dois casos específicos:

Redes com uma única camada

Definição das variáveis:

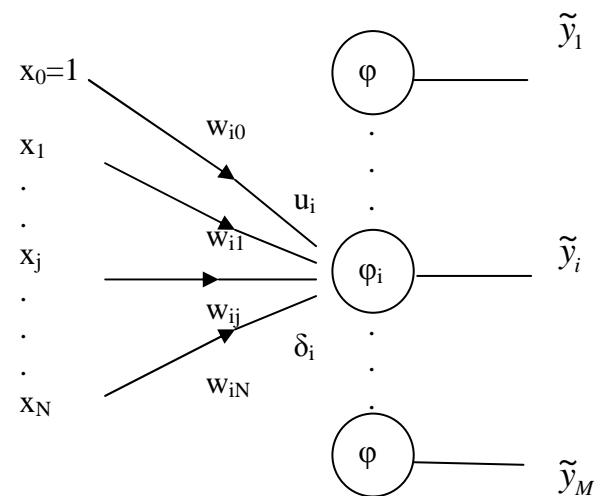
N entradas: x_1, x_2, \dots, x_N e x_0 (bias) = 1

w_{ij} - sinapse conectando a entrada j
ao neurônio i

M neurônios com função de ativação
 $\varphi(\cdot)$ linear ou tgh

u_i - excitação interna do neurônio i
 \tilde{y}_i - saída do neurônio i

Rede Original



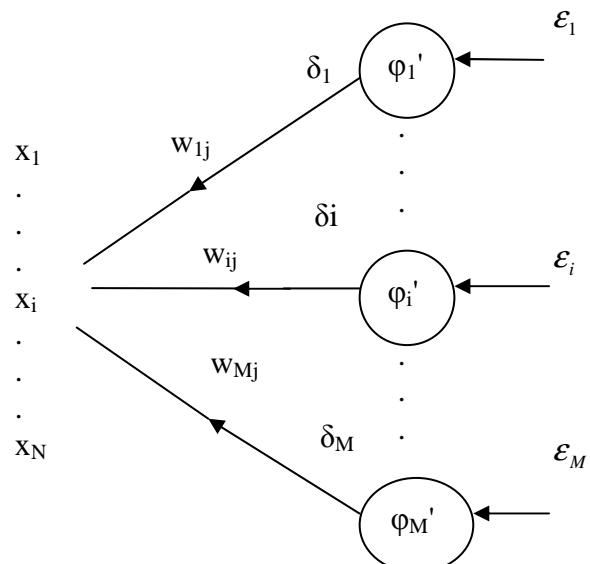
Definição das variáveis:

ε_i - erro na saída do neurônio i

δ_i - erro retropropagado até a entrada do neurônio i

φ'_i - valor numérico da derivada de φ_i no ponto de operação

Rede Associada



Algoritmo BP acréscimo Δw_{ij}^p - Rede 1 camada

Para o par entrada-saída p ($\underline{x}, \underline{y}$)

Signal feedforward: Propague o sinal para a frente, calcule a saída e o erro em cada saída:

$$u_i = \sum_{j=0}^N w_{ij} x_j$$

$$\tilde{y}_i = \varphi_i(u_i) = \begin{cases} u_i & \text{neurônio linear} \\ \tgh(u_i) & \text{neurônio tgh} \end{cases}$$

$$\varepsilon_i = y_i - \tilde{y}_i$$

Calcule o erro retropropagado até a entrada de cada neurônio:

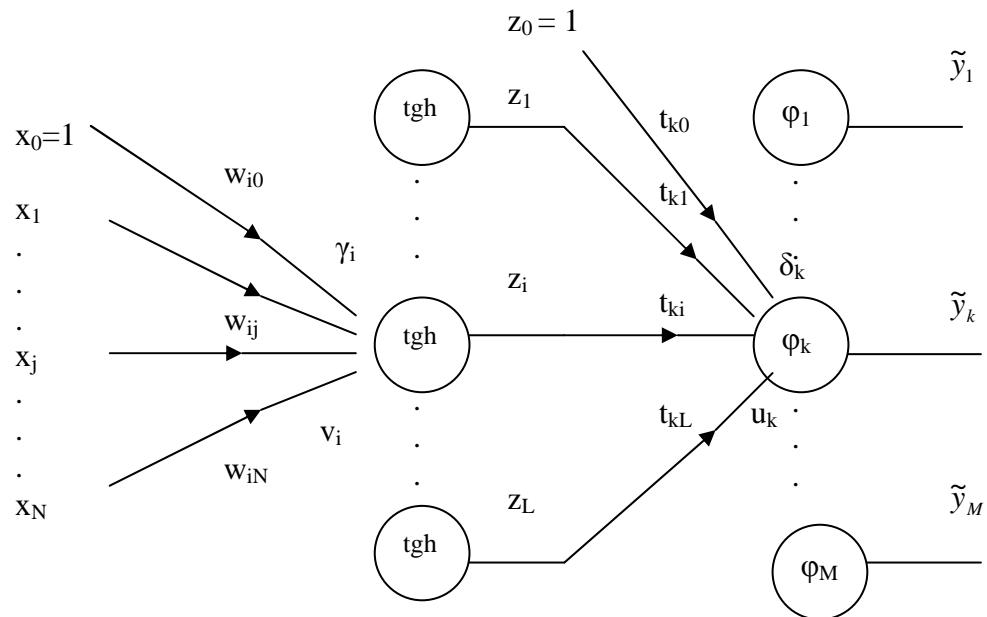
$$\delta_i = \begin{cases} \varepsilon_i & \text{neurônio linear} \\ (1 - \tilde{y}_i^2) \varepsilon_i & \text{neurônio tgh} \end{cases}$$

Calcule o acréscimo em cada sinapse devido ao par p:

$$\Delta w_{ij}^p = 2\alpha x_j \delta_i$$

fim do algoritmo

Redes com duas camadas - Rede Original



Definição das variáveis:

N entradas: x_1, x_2, \dots, x_N e x_0 (bias) = 1

L neurônios na primeira camada com função de ativação tgh

M neurônios na camada de saída com função de ativação $\varphi(\cdot)$ linear ou tgh

w_{ij} - sinapse conectando a entrada j ao neurônio i da camada intermediária

t_{ki} - sinapse conectando a saída z_i do neurônio i da camada intermediária com a entrada do neurônio k da camada de saída.

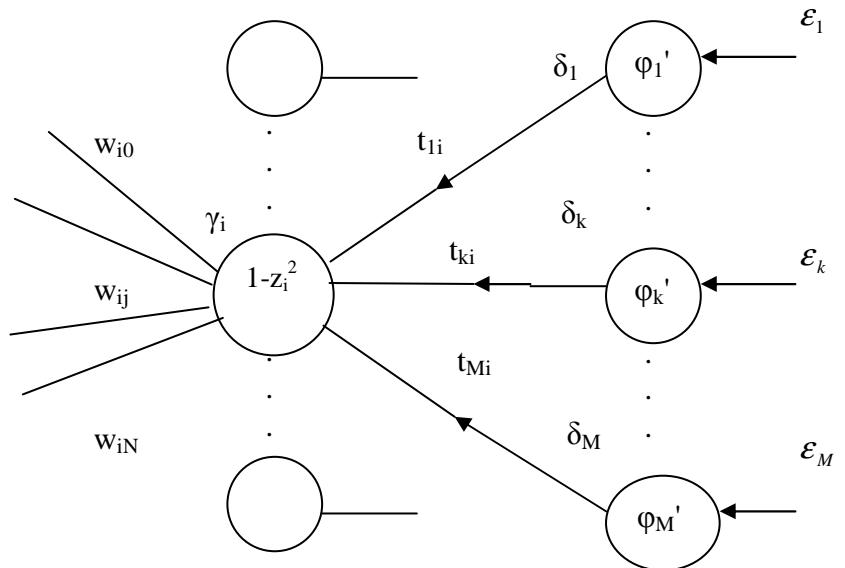
v_i - excitação interna do neurônio i da primeira camada

z_i - saída do neurônio i da camada intermediária. z_0 (bias) = 1

u_k - excitação interna do neurônio k da camada intermediária

\tilde{y}_k - saída do neurônio k da camada de saída

Rede Associada



Definição das variáveis:

ϵ_k - erro na saída do neurônio da camada de saída

δ_k - erro retropropagado até a entrada do neurônio k da camada de saída

φ_i' - valor numérico da derivada de φ_i no ponto de operação

Algoritmo BP acréscimos Δw_{ij}^p e Δt_{ki}^p - Rede com 2 camadas

Para o par entrada-saída p ($\underline{x}, \underline{y}$)

Signal feedforward: Propague o sinal para a frente, calcule a saída e o erro em cada saída:

Para todos os neurônios da primeira camada $i = 1, \dots, L$

$$\begin{aligned} v_i &= \sum_{j=0}^N w_{ij} x_j & x_0 &= 1 \\ z_i &= \text{tgh}(v_i) & z_0 &= 1 \end{aligned}$$

Para todos os neurônios da segunda camada $k = 1, \dots, M$

$$\begin{aligned} u_k &= \sum_{i=0}^L t_{ki} z_i & z_0 &= 1 \\ \tilde{y}_k &= \varphi_k(u_k) = \begin{cases} u_k & \text{neurônio linear} \\ \text{tgh}(u_k) & \text{neurônio tgh} \end{cases} \\ \epsilon_k &= y_k - \tilde{y}_k \end{aligned}$$

Retropropagação do erro

Para todo neurônio da segunda camada, $k = 1, \dots, M$:

$$\delta_k = \begin{cases} \varepsilon_k & \text{neurônio linear} \\ (1 - \tilde{y}_k^2) \varepsilon_k & \text{neurônio tgh} \end{cases}$$

Para todo neurônio da primeira camada, $i = 1, \dots, L$

$$\gamma_i = (1 - z_i^2) \sum_{k=1}^M t_{ki} \delta_k$$

Acréscimo nas sinapses da primeira camada devido ao par p:

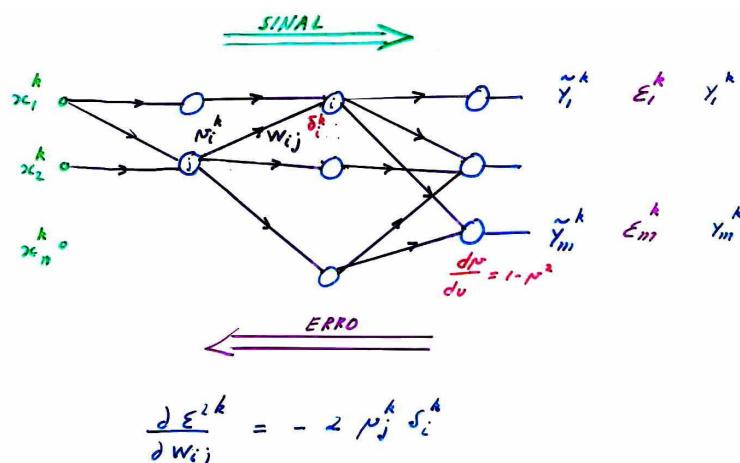
$$\Delta w_{ij}^p = 2 \alpha x_j \gamma_i \quad \forall j = 0, 1, \dots, N \quad e \quad i = 1, 2, \dots, L$$

Acréscimo nas sinapses da segunda camada devido ao par p:

$$\Delta t_{ki}^p = 2 \alpha z_i \delta_k \quad \forall i = 0, 1, \dots, L \quad e \quad k = 1, 2, \dots, M$$

fim do algoritmo

Error Backpropagation - Resumo:



$$\frac{\partial \varepsilon^k}{\partial w_{ij}} = -2 \mu_j^k \delta_i^k$$

$$\nabla E = \left[\frac{\partial E}{\partial w_{ij}} \right] = -2 \cdot E(\mu_j \delta_i) \left[\right]$$

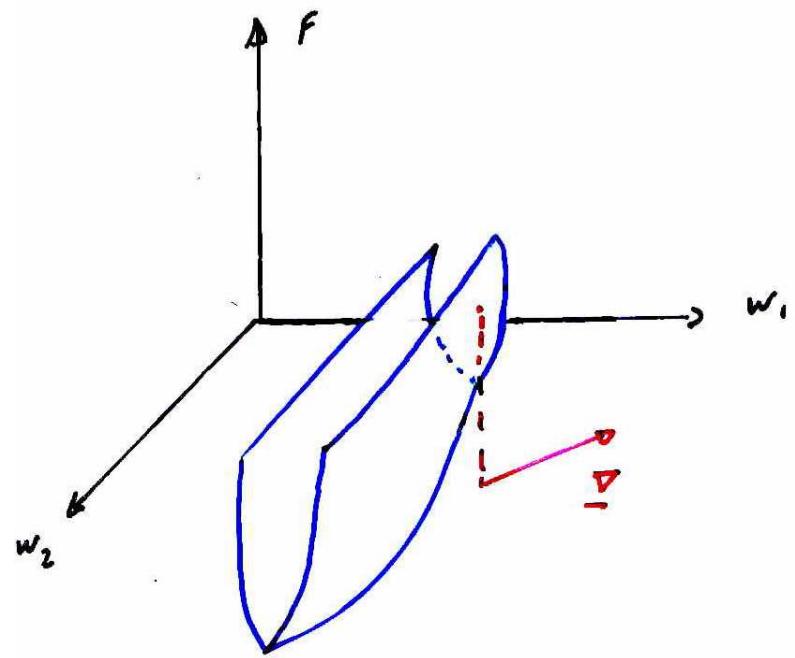
$$\Delta w_{ij} = 2 \alpha E(\mu_j \delta_i)$$

Método de Primeira Ordem Aprimorado

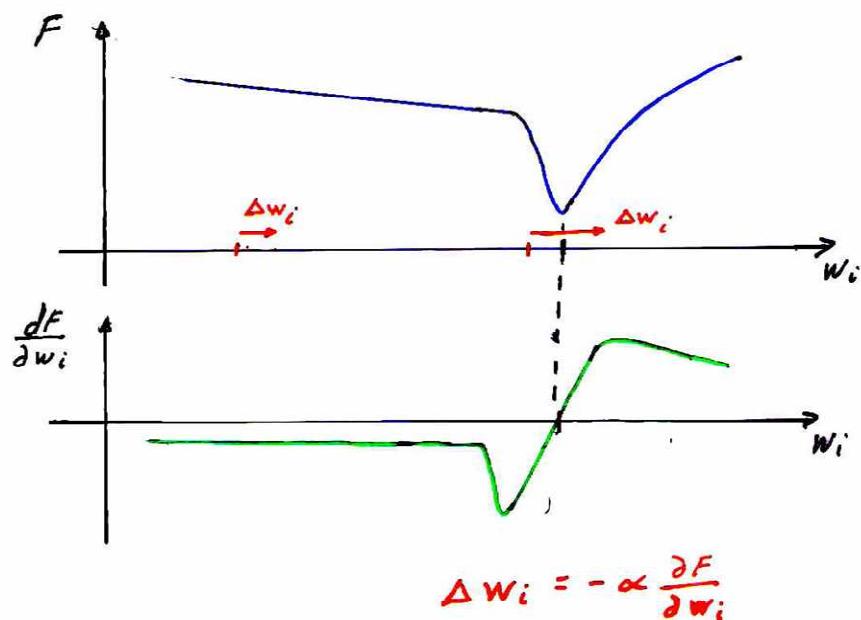
Resilient Backpropagation

Acelerando o processo

passo variável
por sinapse



Acelerando o processo – passo variável no tempo



Resilient Backpropagation (modificado)

(Silva e Almeida, Super SAB, etc.)

$$\alpha \text{ variável} \quad \left\{ \begin{array}{l} \text{por variável } w_i \\ \text{por passo de treinamento } n \end{array} \right.$$

$$\Delta w_i(n) = -\alpha_i(n) \frac{\partial F}{\partial w_i}(n)$$

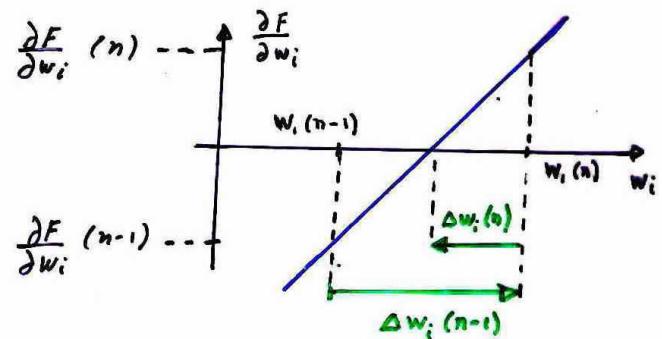
$$\alpha_i(n) = \begin{cases} a \alpha_i(n-1) & \text{se } \operatorname{sign} \frac{\partial F}{\partial w_i}(n) = \operatorname{sign} \frac{\partial F}{\partial w_i}(n-1) \\ b \alpha_i(n-1) & \text{se } \operatorname{sign} \frac{\partial F}{\partial w_i}(n) \neq \operatorname{sign} \frac{\partial F}{\partial w_i}(n-1) \end{cases}$$

$$\text{se } \alpha_i(n) > \alpha^* \text{ faça } \alpha = \alpha^*$$

$$\alpha_i(0) \approx .05 \quad \alpha^* \approx .3 \quad a \approx 1.05 \quad b \approx .9$$

Outra alternativa

$$b = \frac{\frac{\partial F}{\partial w_i}(n-1)}{\frac{\partial F}{\partial w_i}(n-1) - \frac{\partial F}{\partial w_i}(n)}$$



BP Resiliente - Algoritmo

Iniciar: $n = 1; \alpha_i(1) = .1; \alpha_{\max} \approx .5; a \approx 1.05; b \approx .9; \frac{\partial F}{\partial w_i}(0) = 0^+;$
 $w_i(1) \in [.2, -.2] \text{ randômicos};$

Até que o critério de parada esteja satisfeito

Passo de treinamento n

Cálculo das derivadas por sinapse

para cada par $(\underline{x}^p, \underline{y}^p) \quad p = 1, \dots, P$

$$\text{calcular} \quad \underline{\tilde{y}}^p = \phi(\underline{x}^p) \quad \epsilon_m^p = y_m^p - \tilde{y}_m^p \quad \text{e} \quad \frac{\partial \tilde{y}_m^p}{\partial w_i} \quad \forall m, i$$

$$\frac{\partial \epsilon^p}{\partial w_i} = -2 \sum_{m=1}^M \epsilon_m^p \frac{\partial y_m^p}{\partial w_i} \quad \forall i$$

outro par

$$\text{calcular} \quad \frac{\partial F}{\partial w_i}(n) = E_p \frac{\partial \epsilon^p}{\partial w_i}$$

Cálculo do acréscimo

para $i = 1, 2, \dots$

$$\alpha_i(n) = \begin{cases} a \alpha_i(n-1) & \text{se } \text{sign} \frac{\partial F}{\partial w_i}(n) = \text{sign} \frac{\partial F}{\partial w_i}(n-1) \\ & \text{mas se } \alpha_i(n) > \alpha_{\max} \text{ faça } \alpha_i(n) = \alpha_{\max} \\ \frac{\frac{\partial F}{\partial w_i}(n-1)}{\frac{\partial F}{\partial w_i}(n-1) - \frac{\partial F}{\partial w_i}(n)} \alpha_i(n-1) & \text{se } \text{sign} \frac{\partial F}{\partial w_i}(n) \neq \text{sign} \frac{\partial F}{\partial w_i}(n-1) \end{cases}$$

$$\Delta w_i(n) = -\alpha_i(n) \frac{\partial F}{\partial w_i}(n)$$

Atualização das sinapses

$$w_i(n+1) = w_i(n) + \Delta w_i(n)$$

$$n = n + 1$$

fim do algoritmo

Métodos de Segunda Ordem

necessitam das derivadas de segunda ordem e $\underline{d} = -\underline{H}^{-1} \nabla$
da inversão da Hessiana !

Métodos de 2^a Ordem

Newton

Newton Amortecido

Levemberg-Marquadt

Quasi Newton

BFGS

Gradiente Conjugado

Comentários finais

Primeira ordem

Passo fixo (“vanilla”) – convergência usualmente
muito lenta

BP Resiliente (Silva & Almeida, SuperSAB) –
bom compromisso entre estabilidade e rapidez na
convergência.

Segunda ordem – rápidos, mas mais complexos e com risco
de instabilidade na convergência

Levenberg – Marquardt (Newton amortecido)

BFGS (quase Newton)