# Camada de Kohonen

# Classificação por Similaridade

Classes agrupam elementos "similares" entre si



### versus classificadores arbitários

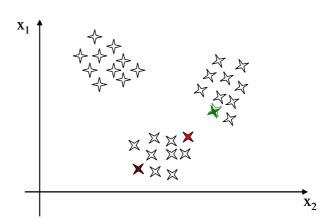
# 1 - Classificação por Similaridade - Critérios de pertinência

objetos físicos  $O_1 O_2$  objetos matemáticos  $\underline{\mathbf{x}}_1 \ \underline{\mathbf{x}}_2$  similaridade física similaridade matemática

$$\bigcirc_1 \approx \bigcirc_2 \qquad \underline{x}_1 \cong \underline{x}_2 \quad \text{ou} \quad |\underline{x}_1 - \underline{x}_2| <<$$

### Classificadores por similaridade

# Critérios de Pertinência à uma Classe



#### Critério básico

Dois elementos pertencem à mesma classe se estão próximos entre si

# Critério: k vizinhos mais próximos.

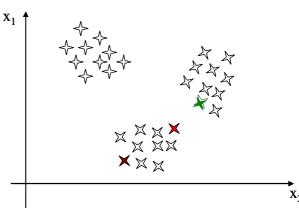
Considere os k elementos mais próximos da entrada que se pretende classificar. A entrada pertence à classe à qual pertencem a maioria dos k vizinhos.

### Critério simplificado: vizinho mais próximo.

Se k=1 uma entrada pertence a uma classe se seu vizinho mais próximo é um elemento desta classe

# pouco prático

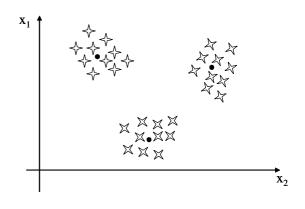
a descrição de uma classe exige a utilização de todos seus elementos (domínio da classe)



Padrão de classe:

Padrão  $\underline{w}_i$  da classe  $C_i$ 

$$\underline{w}_i = \underbrace{E}_{\forall \underline{x}_j \in C_i} \underline{x}_j = \frac{1}{N_i} \sum_{\forall \underline{x}_j \in C_i, j=1}^{N_i} \underline{x}_j$$



a **média** ou **baricentro** é escolhido como padrão da classe porque minimiza a **dissimilaridade**ou **erro de representação**da classe

Dispersão média intra classe ou erro (médio quadrático) de representação Da classe  $C_i$  para o seu padrão  $\vec{p}_j$ 

$$\sigma_{j}^{2} = E_{\forall x_{i} \in C_{j}} \|x_{i} - p_{j}\|^{2} = \frac{1}{n_{j}} \sum_{\forall x_{i} \in C_{j}} \|x_{i} - p_{j}\|^{2}$$

por componente k

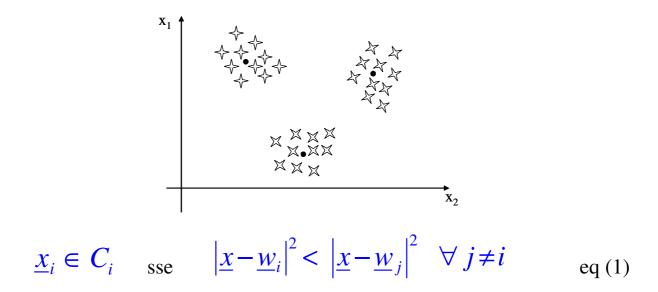
$$\sigma_{jk}^{2} = E_{\forall x_{i} \in C_{j}} (x_{ik} - p_{jk})^{2} = \frac{1}{n_{j}} \sum_{\forall x_{i} \in C_{j}} (x_{ik} - p_{jk})^{2}$$

após pequena álgebra 
$$\frac{\partial \sigma_{jk}^2}{\partial p_{jk}} = 0 \quad \Rightarrow \quad p_{jk} = m_{jk} \quad \Rightarrow \quad \vec{p}_j = \vec{m}_j$$

o padrão que minimiza a dispersão intra classe (ou erro de representação) de uma classe é o seu baricentro

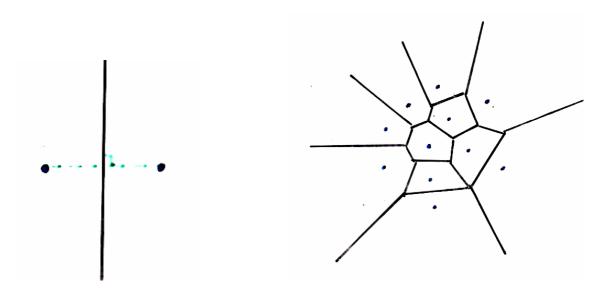
# 1.1 - Critério de pertinência 1 :

### Padrão mais similar à entrada



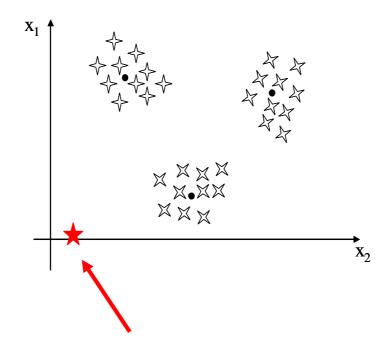
porque minimiza a dispersão intra-classes total

# Separadores para o critério 1

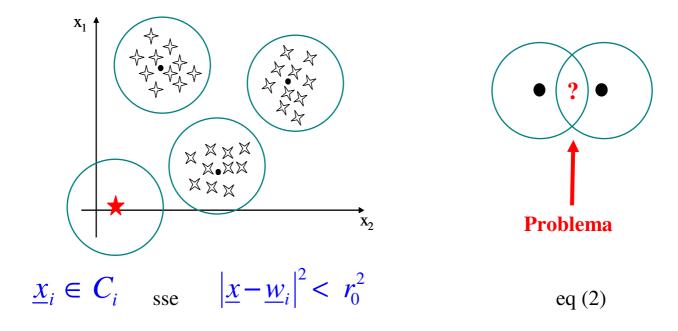


Tecelagem (diagrama) de Voronoi

# **Problema:**



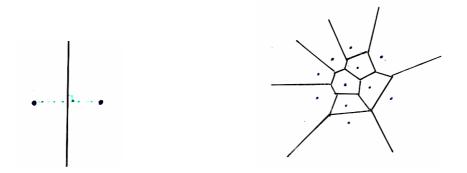
# 1.2 - Critério de pertinêcia 2 - Padrão que satisfaz uma similaridade minima



r<sub>0</sub> – raio de similaridade

# **Separadores**

# Critério 1 - padrão mais similar (mais próximo) a entrada Eq(1)

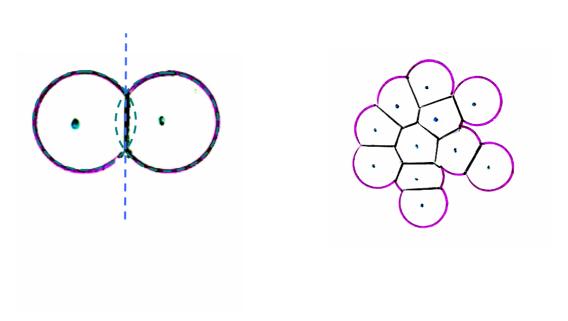


# Critério 2 - mínima similaridade exigida Eq (2)



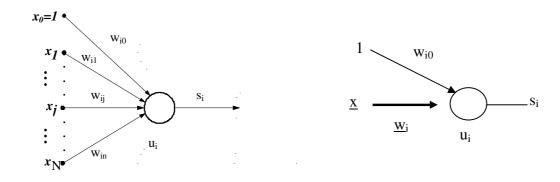
# 1.3 - Critérios $\underline{1}$ e $\underline{2}$ , Eq(1) + Eq(2)

# padrão mais similar & mínima similaridade atingida



#### 2 - Redes Neurais

### Neurônio para redes feedforward



$$u_{i} = \sum_{j=1}^{N} w_{ij} x_{j} + w_{i0} = \underline{w}_{i}^{t} \underline{x} + w_{i0}$$

$$\widetilde{y}_{k} = \varphi_{k}(u_{k}) = \begin{cases} u_{k} & neurônio \ linear \\ tgh(u_{k}) & neurônio \ tgh \end{cases}$$

### Neurônio como comparador de similaridade

de uma entrada  $\underline{x}$  com dois padrões  $\underline{w}_i$  e  $\underline{w}_j$ 

$$d_{i}^{2} = |\underline{x} - \underline{w}_{i}|^{2} = |\underline{x}|^{2} + |\underline{w}|^{2} - 2\underline{w}^{t}\underline{x}$$

$$u_{i} = \underline{w}^{t}\underline{x} + w_{i0}$$

$$fazendo \quad w_{i0} = -\frac{1}{2}|\underline{w}_{i}|^{2}$$

$$u_{i} = \frac{1}{2}(|\underline{x}|^{2} - d_{i}^{2})$$

$$u_i > u_j \iff d_i^2 < d_j^2$$

# e u é um comparador de similaridade

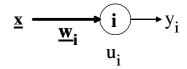
y<sub>i</sub> – depende dos outros neurônios

#### Neurônio como medidor de similaridade

entre uma entrada x e um padrão w<sub>i</sub>

#### uma outra definição

$$u_i = -\left\|\underline{x} - \underline{w}_i\right\|^2 = -d_i^2 \le 0$$



 $u_i$  - medida de similaridade entre  $\underline{x}$  e  $\underline{w}_i$ 

 $u_i = 0$  distância nula = máxima similaridade

y<sub>i</sub> – depende dos outros neurônios

### 3- Camada de Kohonen

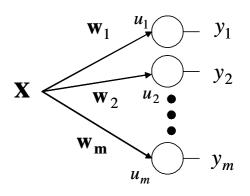
### **Template Matching**

$$u_i = -d_i^2$$

maior  $u_i =$ 

menor distância =

maior similaridade



 $u_i$  é uma medida de similaridade entre  $\underline{x}$  e  $\underline{w}_i$ 

#### Winner-takes-all

$$y_i = 1$$
 sse  $u_i > u_j \quad \forall j \neq i$   
 $y_i = 0$  caso contrário

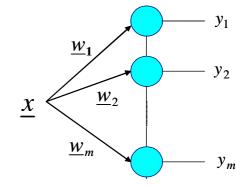
# Camada de Kohonen

[versão unidimensional, simplificada (D=0)]

# Classe C<sub>i</sub>

Padrão w<sub>i</sub>

Indicador y<sub>i</sub>

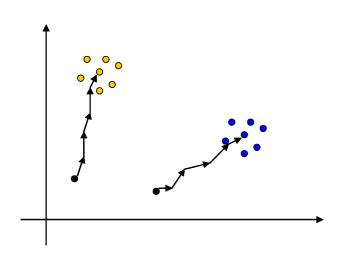


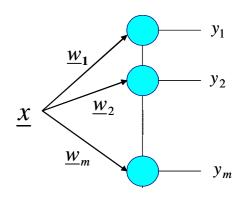
 $Se \quad y_i = 1 \quad ent \tilde{a}o \quad \underline{x} \quad \in \quad C_i$ 

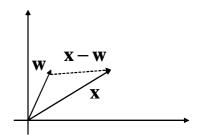
pelo critério 1 (padrão mais similar a entrada).

# Camada de Kohonen

# Treinamento Supervisionado





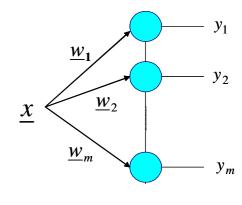


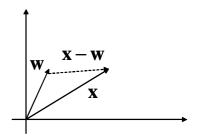
# 4 - Treinamento da Camada de Kohonen

$$\underline{x}(n) \in C_i \quad \Rightarrow \quad y_i = 1$$
$$y_j = 0 \ \forall \ j \neq i$$

$$\underline{w}_{i}(n+1) = \underline{w}_{i}(n) + \alpha[\underline{x}(n) - \underline{w}_{i}(n)]$$
$$= (1 - \alpha)\underline{w}_{i}(n) + \alpha\underline{x}(n)$$

$$\underline{w}_{j}(n+1) = \underline{w}_{j}(n) \quad \forall \ j \neq i$$





# Onde estabiliza, qual o valor final de $\underline{\mathbf{w}}$ ?

### Evolução do w de uma classe

Passo n:  $\underline{\mathbf{x}}(\mathbf{n})$   $\underline{\mathbf{w}}(\mathbf{n})$ 

$$x(n) \in C_i$$

$$\begin{cases} \underline{\mathbf{w}}_{i}(n+1) = (1-\alpha)\underline{\mathbf{w}}_{i}(n) + \alpha\underline{\mathbf{x}}(n) \\ \\ \underline{\mathbf{w}}_{i}(n+1) = \underline{\mathbf{w}}_{i}(n) \end{cases} \forall j \neq i$$

$$\underline{\mathbf{w}}(0)$$

$$\underline{\mathbf{w}}(1) = (1-\alpha) \underline{\mathbf{w}}(0) + \alpha \underline{\mathbf{x}}(1)$$

$$\underline{\mathbf{w}}(2) = (1-\alpha) \underline{\mathbf{w}}(1) + \alpha \underline{\mathbf{x}}(2)$$

$$= (1-\alpha)[(1-\alpha) \underline{\mathbf{w}}(0) + \alpha \underline{\mathbf{x}}(1)] + \alpha \underline{\mathbf{x}}(2)$$

$$= (1-\alpha)^2 \underline{\mathbf{w}}(1) + \alpha(1-\alpha) \underline{\mathbf{x}}(1) + \alpha \underline{\mathbf{x}}(2)$$

$$\underline{\mathbf{w}}(3) = (1-\alpha) \, \underline{\mathbf{w}}(2) + \alpha \, \underline{\mathbf{x}}(3)$$

$$= (1-\alpha)[(1-\alpha)^2 \, \underline{\mathbf{w}}(0) + \alpha(1-\alpha) \, \underline{\mathbf{x}}(1) + \alpha \, \underline{\mathbf{x}}(2)] + \alpha \, \underline{\mathbf{x}}(3)$$

$$= (1-\alpha)^3 \, \underline{\mathbf{w}}(0) + \alpha(1-\alpha)^2 \, \underline{\mathbf{x}}(1) + \alpha(1-\alpha) \, \underline{\mathbf{x}}(2) + \alpha \, \mathbf{x}(3)$$

:

$$\mathbf{w}(\mathbf{n}) = (1-\alpha)^{\mathbf{n}} \, \underline{\mathbf{w}}(0) + \alpha \left[ (1-\alpha)^{\mathbf{n}-1} \, \underline{\mathbf{x}}(1) + (1-\alpha)^{\mathbf{n}-2} \, \underline{\mathbf{x}}(2) + \dots + \underline{\mathbf{x}}_{\mathbf{n}} \right]$$
$$= (1-\alpha)^{\mathbf{n}} \, \underline{\mathbf{w}}(0) + \alpha \, \sum_{i=1}^{n} \, (1-\alpha)^{\mathbf{n}-i} \, \underline{\mathbf{x}}(i)$$

algumas aproximações:

$$0 < \alpha < 1 \qquad 0 < (1-\alpha) < 1 \qquad n >> 1 \qquad (1-\alpha)^{n} \to 0$$

$$1 + (1-\alpha) + (1-\alpha)^{2} + (1-\alpha)^{3} + \dots = \sum_{i=0}^{\infty} (1-\alpha)^{i} = \frac{1}{1-(1-\alpha)} = \frac{1}{\alpha}$$

$$n >> 1 \qquad \alpha \approx \frac{1}{\sum_{i=0}^{n-1} (1-\alpha)^{i}}$$

$$\underline{\mathbf{w}}(\mathbf{n}) = (1-\alpha)^{n} \underline{\mathbf{w}}(0) + \alpha \sum_{i=1}^{n} (1-\alpha)^{n-i} \underline{\mathbf{x}}(\mathbf{i})$$

$$\cong \frac{\sum_{i=1}^{n} (1-\alpha)^{n-i} \overline{\mathbf{x}}(\mathbf{i})}{\sum_{i=1}^{n-1} (1-\alpha)^{i}} = \frac{\sum_{i=0}^{n-1} (1-\alpha)^{i} \overline{\mathbf{x}}(n-i)}{\sum_{i=1}^{n-1} (1-\alpha)^{i}}$$

#### Média das entradas que pertencem a classe

#### ponderada geometricamente pelo tempo!

$$\underline{\underline{w}}_{n} \cong \frac{\sum_{i=1}^{n} (1-\alpha)^{n-i} \vec{x}(i)}{\sum_{i=1}^{n} (1-\alpha)^{n-i}} = \left(\frac{1}{\sum_{i=0}^{n-1} (1-\alpha)^{i}}\right) \sum_{i=0}^{n-1} (1-\alpha)^{i} \vec{x}(n-i)$$

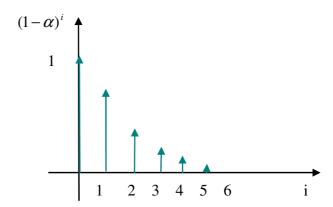
instante atual = n

 $\underline{\mathbf{x}}(\mathbf{n}-\mathbf{i}) = \mathbf{e}\mathbf{n}\mathbf{t}\mathbf{r}\mathbf{a}\mathbf{d}\mathbf{a}$ 

de i intervalos de tempo

 $(1-\alpha)^{i}$  = ponderador

para a entrada  $\underline{x}(n-\hat{i})$ 



### Note que se

$$\alpha \cong 1$$

e a estatística de  $\vec{x}$  for invariante no tempo

$$\vec{w}_n \cong \left(\frac{1}{\sum_{i=0}^{n-1} (1-\alpha)^i}\right) \sum_{i=0}^{n-1} (1-\alpha)^i \vec{x} (n-i) \cong E[\vec{x}]$$

$$\vec{w}_n \cong E[\vec{x}]$$

### 4.1 – Tempo de medida

Fim (prático) da soma ponderada (tempo de medida)

Atraso zero

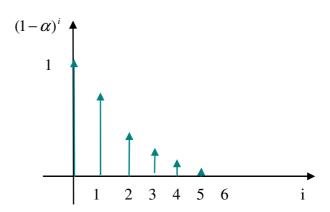
Ponderador = 
$$(1-\alpha)^0 = 1$$

Último atraso à ser considerado:

Ponderador =  $.02 = (1-\alpha)^{i}$ 

$$i = \frac{\ln .02}{\ln (1 - \alpha)} \approx \frac{-4}{-\alpha} |_{\alpha \to 0}$$

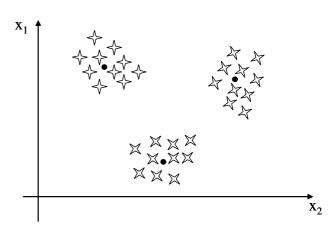
$$i \approx \frac{4}{\alpha}$$



# 4.2 - Erro na determinação do baricentro:

Se em uma classificação por similaridade cada elemento  $\underline{x}$  de uma classe  $C_k$  pode ser representado pelo padrão  $\underline{w}_k$  da classe adicionado de um vetor de ruído  $\underline{r}$  com média nula.

$$\underline{\mathbf{x}} = \underline{\mathbf{w}}_{\mathbf{k}} + \underline{\mathbf{r}}$$



Cada componente  $x_j$  de  $\underline{x}$  é então representada pela componente correspondente de  $\underline{w}_k$ ,  $w_j$ , adicionada de um ruído  $r_j$  de média nula e variância  ${\sigma_{xj}}^2$ 

$$x_j = w_j + r_j$$

Com que precisão componente  $w_j$  é calculada pela camada de Kohonen, qual a sua variância  $\sigma_{wj}^2$ ? A componente j de  $\underline{w}$ ,  $w_j$ , em um instante n >> 1 é dada por

$$w_{j} \cong \left(\frac{1}{\sum_{i=0}^{n} (1-\alpha)^{i}}\right) \sum_{i=0}^{n} (1-\alpha)^{i} x_{j} (n-i)$$

Sendo uma soma ponderada sua variância será dada por

$$\sigma_{wj}^2 \cong \left(\frac{1}{\sum_{i=0}^{\infty} (1-\alpha)^i}\right)^2 \sum_{i=0}^{\infty} (1-\alpha)^{2i} \sigma_{xj}^2 \qquad ou \qquad \sigma_{wj}^2 \cong \frac{\alpha}{2} \sigma_{xj}^2$$

### 4. 3 – Compromisso entre Erro vs. Tempo de Medida

Quanto menor α mais precisamente o baricentro será determinado. Mas, como esperado, maior o tempo (número de passos) necessário para calculá-lo

$$\sigma_{wj}^2 \cong \frac{\alpha}{2} \sigma_{xj}^2$$
 # passos =  $i \cong \frac{4}{\alpha}$ 

#### **Exemplo:**

$$\sigma_{\rm x}=.05$$
 (5%) 
$$\sigma_{\rm w}=.01$$
 (1%) requerido 
$$\alpha=\frac{2\sigma_{\rm wj}^2}{\sigma_{\rm xi}^2}=.08 \qquad \# passos=\frac{4}{\alpha}=50 \ passos$$

Mas se:

$$\sigma_{\rm w} = .001 \ (.1\%) \ {\rm requerido}$$
  $\alpha = \frac{2\sigma_{\rm wj}^2}{\sigma_{\rm xj}^2} = .0008 \ \# \ passos = \frac{4}{\alpha} = 5.000 \ passos !!$ 

# 4.4 Considerações:

#### 4.4.1. Fim do treinamento?



$$E[\Delta \underline{w}] = \underline{0}$$



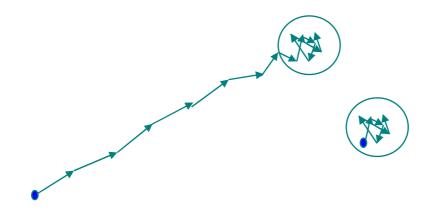
$$\Delta \vec{w} = \alpha \left( \vec{x} - \vec{w} \right)$$

$$E(\Delta \vec{w}) = E[\alpha(\vec{x} - \vec{w})] = \alpha[E(\vec{x}) - \vec{w}] = 0$$

$$\vec{w} = E(\vec{x})$$

### 4.4.2. Valores iniciais das sinapses

Irrelevantes! Mas escolher o valor inicial como o primeiro elemento da classe acelera o treinamento, porque a sinapse já começa dentro da classe.



#### 4.4.3. Passo de treinamento

reduzir ao longo do tempo

$$\alpha(n) = \alpha_0 e^{-\frac{n}{N_0}}$$

$$\alpha(0) = \alpha_0$$

$$\alpha(n+1) = k \alpha(n)$$

$$\alpha(0) = \alpha_0$$
  $\alpha(n+1) = k \alpha(n)$  onde  $k = 1 - \frac{1}{N_0}$ 

o processo acaba ( $\alpha(n) \ll \alpha_0$ ) para  $n > 4N_0$ 

e as sinapses pouco ativadas (classes pouco populosas)? usar α diferenciado por sinapse

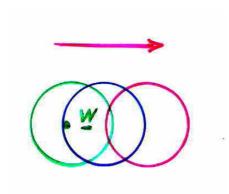
$$\alpha_{w_i}(0) = \alpha_0$$

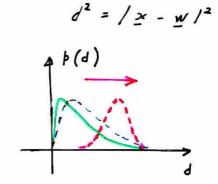
$$\alpha_{w_i}(0) = \alpha_0$$
  $\alpha_{w_i}(n_j + 1) = k \alpha_{w_i}(n_j)$ 

onde  $n_i$  é o número de vezes que a sinapse  $\underline{w}_i$  foi treinada.

# 4.4.4 Treinamento dinâmico, adaptativo

Uma classe varia de posição no tempo. Como saber?

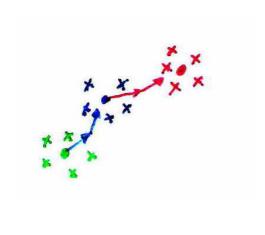


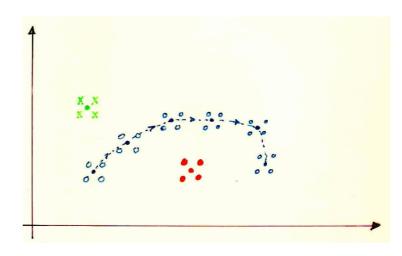


O valor médio de  $u_i = -d_i^2$  diminui

(o valor da distância média (ou a da moda) aumenta)

Como corrigir ? Ligar o treinamento até que d<sup>2</sup> i volte a seus valores normais.





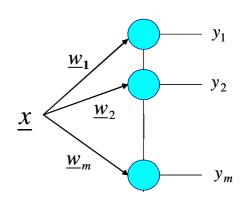
Se a variação do baricentro for lenta e o treinamento estiver ligado o padrão segue o baricentro

### **5 – LVQ – Learning Vector Quantization**

$$\vec{x}(n) \in C_i$$

e

neurônio vencedor  $y_j$ 



$$j = i$$
 ou  $j \neq i$ 

se 
$$j = i$$
  $\vec{w}_j(n+1) = \vec{w}_j(n) + \alpha \left[ \vec{x}(n) - \vec{w}_j(n) \right]$ 

se a entrada pertence à classe vencedora, aproxima a sinapse vencedora da entrada e conseqüentemente do centro da classe

se 
$$j \neq i$$
  $\vec{w}_j(n+1) = \vec{w}_j(n) - \alpha \left[ \vec{x}(n) - \vec{w}_j(n) \right]$ 

se a entrada não pertence à classe vencedora, afasta a sinapse vencedora da entrada e, conseqüentemente, do centro  $\underline{\mathbf{w}}_i$  da classe a que a entrada pertence. Isto aumenta a chance da classe certa vencer em uma próxima competição.

Note que se as sinapses tiverem sido inicializadas com um elemento da classe a probabilidade do segundo caso (  $j \neq i$  ) ocorrer e é muito pequena. Isto também ocorre após a fase inicial do treinamento, quando as sinapdes já foram arrastadas para dentro de suas respectivas classes.

#### Obs:

1 - O passo de aprendizagem deve decrescer com o tempo

$$\alpha = \alpha(n) = \alpha_0 \exp\left(-\frac{n}{\tau_{\alpha}}\right)$$
  $\alpha_0 \approx 0.1$   $\tau_{\alpha} \approx 1000$ 

**2** - O equilíbrio é atingido quando  $E[\Delta \vec{w}_i] = \vec{0}$  , o que ocorre quando

$$\vec{w}_i = \vec{m}_i = \sum_{\forall \vec{x} \in C_i} (\vec{x})$$

**3** - Este processo corresponde a minimizar a função objetivo

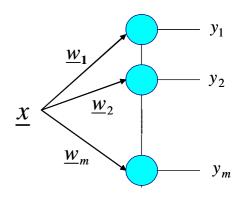
$$F = \sum_{\forall C_i} F(\vec{w}_i) \qquad F(\vec{w}_i) = \sum_{\forall \vec{x} \in C_i} |\vec{x} - \vec{w}_i|^2$$

isto é, minimizar a dispersão intra-classe total

### 5a – Um método alternativo (???)

#### LVQ - Learning Vector Quantization

 $\vec{x}(n) \in C_i$  aproxima o padrão vencedor de  $\vec{x}(n)$  e afasta os demais padrões de  $\vec{x}(n)$ 



$$\vec{x}(n) \in C_i \qquad \begin{cases} y_i = 1 \\ y_j = 0 \quad \forall \ j \neq i \end{cases}$$

$$\vec{w}_i(n+1) = \vec{w}_i(n) + \alpha \left[ \vec{x}(n) - \vec{w}_i(n) \right]$$

$$\vec{w}_j(n+1) = \vec{w}_j(n) - \alpha \left[ \vec{x}(n) - \vec{w}_j(n) \right] \quad \forall \ j \neq i$$

ou, usando uma fórmula única

$$\vec{x}(n) \in C_i \qquad \begin{cases} y_i = 1 \\ y_j = 0 \quad \forall \ j \neq i \end{cases}$$

$$\vec{w}_{i}(n+1) = \vec{w}_{i}(n) + (2y_{i} - 1)\alpha \left[\vec{x}(n) - \vec{w}_{i}(n)\right] \qquad \forall j$$

O passo de aprendizagem deve decrescer com o tempo

$$\alpha = \alpha(n) = \alpha_0 \exp\left(-\frac{n}{\tau_\alpha}\right)$$
  $\alpha_0 = 0.1$   $\tau_\alpha = 1000$ 

Este processo melhora as fronteiras de classificação (Haykin).

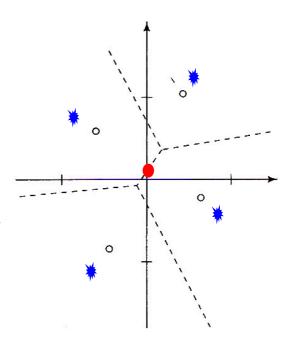
O equilíbrio é atingido quando  $E[\Delta \vec{w}_i] = \vec{0}$ , o que ocorre quando

$$\vec{w}_{i} = p_{i} \underbrace{E}_{\forall \vec{x} \in C_{i}} (\vec{x}) - (1 - p_{i}) \underbrace{E}_{\forall \vec{x} \notin C_{i}} (\vec{x}) =$$

$$= 2p_{i} \underbrace{E}_{\forall \vec{x} \in C_{i}} (\vec{x}) - \underbrace{E}_{\forall \vec{x}} (\vec{x}) =$$

$$= 2p_{i} \vec{m}_{i} - \vec{m}$$

onde  $p_i$  é a probabilidade de uma entrada  $\vec{x}$  pertencer à classe  $C_i$ ,  $\vec{m}_i$  é o baricentro de  $C_i$  e  $\vec{m}$  é o baricentro de todas as entradas, normalmente  $\vec{m} = \vec{0}$ . O efeito é deslocar os padrões das classes (\*) para longe do centro de todas as entradas (•) mas também dos baricentros (•) das mesmas.



Este processo corresponde a otimizar a função objetivo

$$F(\vec{w}_i) = \sum_{\forall \vec{x} \in C_i} |\vec{x} - \vec{w}_i|^2 - \sum_{\forall \vec{x} \notin C_i} |\vec{x} - \vec{w}_i|^2$$

isto é, minimizar a dispersão intra classe de  $C_i$  e maximizar a distância de  $w_i$  aos elementos das outras classes.

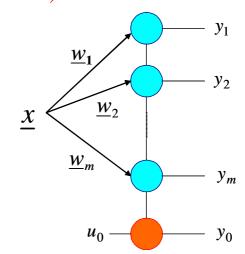
# 6 - E o segundo critério (similaridade mínima)?

#### Camada de Kohonen (aumentada)

$$u_0 = -r_0^2$$

Se  $y_i = 1$  entao

$$\underline{x} \ \in \ C_i$$



pelos critérios 1 (padrao mais similar a entrada) e

2 (satisfaz a similaridade minima exigida).

Se  $y_i = 1$  então

$$\underline{\mathbf{x}} \in \mathbf{C}_{\mathbf{i}}$$

pelos critérios

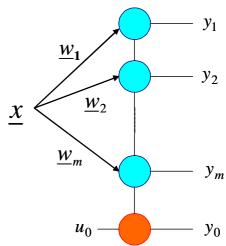
1 (centro de classe mais similar à entrada) e

2 (satisfaz à similaridade mínima)

Se 
$$y_0 = 1$$
 então

para nenhuma classe

 $\underline{\mathbf{x}} \notin \mathbf{C}_{\mathbf{i}} \forall i$ x não satisfaz ao critério 2



# 7 – HVQ – Hierarquical Vector Quantization

