

# Inteligencia Artificial

## Redes Neuronales

Dr. Ing. Martín G. Marchetta  
[martin.marchetta@ingenieria.uncuyo.edu.ar](mailto:martin.marchetta@ingenieria.uncuyo.edu.ar)



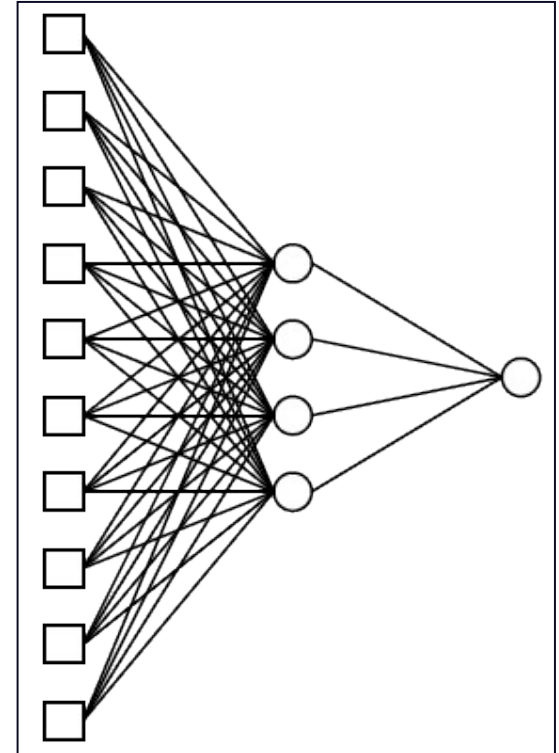
**UNCUYO**  
UNIVERSIDAD  
NACIONAL DE CUYO



- **Una neurona (biológica)**
  - Es una célula que recoge, procesa y emite señales eléctricas
  - Hipótesis: capacidad de procesamiento del cerebro se basa en redes de neuronas
- **Redes Neuronales Artificiales (ANN)**
  - También conocidos como: Artificial Neural Systems, Conexionismo, Computación Neuronal, etc.
  - Imitan la estructura del sistema nervioso
  - Enfoque bottom-up: comportamiento emergente a partir de un sistema “de bajo nivel”

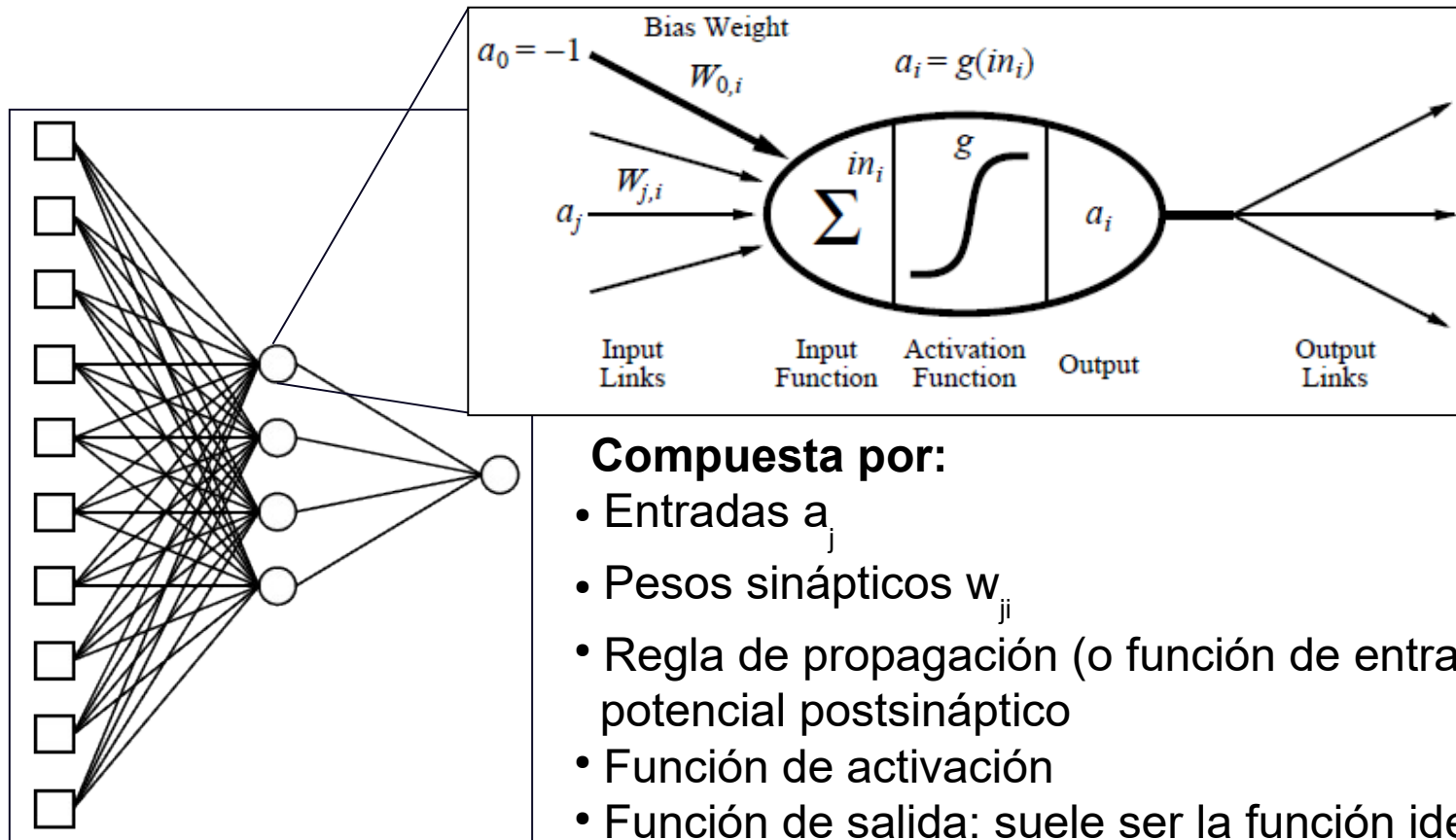
- Definición formal de Red Neuronal Artificial [Martin del Brio y Sanz, 2002]

- Grafo dirigido
- A cada nodo  $i$  se le asigna una variable de estado  $a_i$
- A cada conexión  $j,i$  se le asocia un peso  $w_{ji}$
- A cada nodo  $i$  se le asocia un umbral o sesgo  $\theta_i$
- Para cada nodo  $i$  se define una función  $g(x_i, w_{ji}, \theta_i)$  que proporciona la salida del nodo (función de activación)



- Los nodos son las neuronas artificiales
- Las aristas son las “conexiones sinápticas” que propagan la salida de un nodo a las neuronas conectadas

- Estructura general de una neurona artificial (McCulloch – Pitts, 1943)



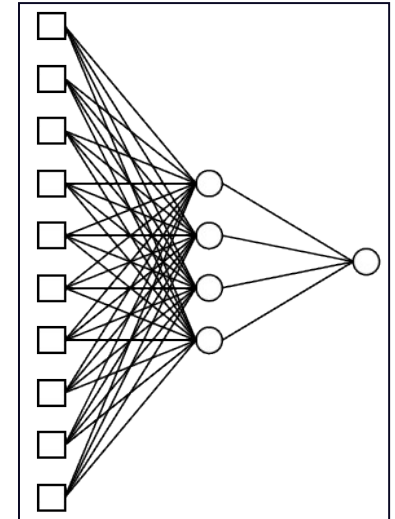
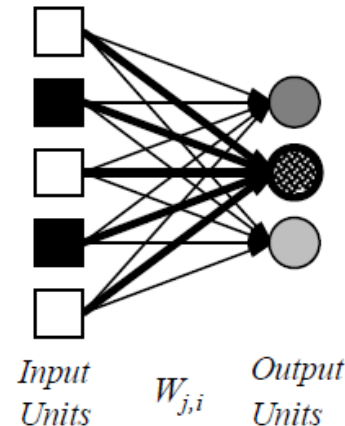
## Compuesta por:

- Entradas  $a_j$
- Pesos sinápticos  $w_{ji}$
- Regla de propagación (o función de entrada): potencial postsináptico
- Función de activación
- Función de salida: suele ser la función identidad  $f(x) = x$
- Salidas

- Clasificación y arquitecturas de ANN

- Por el estilo de aprendizaje

- Supervisados
    - No supervisados
    - Por refuerzo
    - Híbridos



- Por la arquitectura

- Monocapa vs. Multicapa
    - Unidireccionales vs. Realimentados

- Modos de operación de una ANN
  - Modo de aprendizaje (entrenamiento)
    - Aprendizaje: ajuste de los parámetros de la red
    - El tipo de aprendizaje define cómo se ajustan estos parámetros
    - Existen 2 niveles de aprendizaje
      - Ajuste de los valores de los parámetros (este esquema es el más común). Ej: ajuste de pesos sinápticos
      - Ajuste de la topología de la red (cantidad de nodos, conexiones, etc.)
    - La mayoría de los algoritmos de aprendizaje se basan en métodos numéricos iterativos para minimizar el error

- Modos de operación de una ANN
  - Modo de ejecución
    - También llamado “recuerdo” (recall)
    - Es la operación de la ANN “en producción”
    - Generalmente en este modo de ejecución el componente de aprendizaje “se desconecta”
    - La red propaga las señales de entrada a través de su topología y produce una salida

# Redes Neuronales Supervisadas



- Algunas ANN con feed-forward y aprendizaje supervisado
  - Asociador lineal con aprendizaje Hebbiano/Regla Pseudoinversa
  - Perceptrón simple
  - Adaline (ADaptive LInear NE, rebautizada ADaptive LInear Element)
  - Perceptrón Multicapa (MLP, Multi-Layer Perceptron)

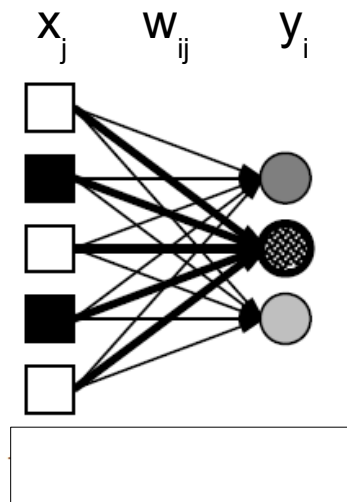
- Asociador lineal

- Objetivo: asociar  $p$  pares  $\{(\mathbf{x}^\mu, \mathbf{t}^\mu) / 1 \leq \mu \leq p\}$

- Regla de propagación 
$$y_i = \sum_{j=0}^n W_{ij} x_j$$

- Función de activación = Función de salida

Identidad  $\rightarrow f(x) = x$



- Asociador lineal (cont.)

- Aprendizaje Hebbiano

$$, \varepsilon \in (0, 1] \quad \Delta w_{ij} = \varepsilon y_i x_j$$

$\varepsilon$  se denomina **ritmo de aprendizaje**

- Tomando  $\varepsilon = 1$ , la **regla de Hebb** para el asociador lineal queda

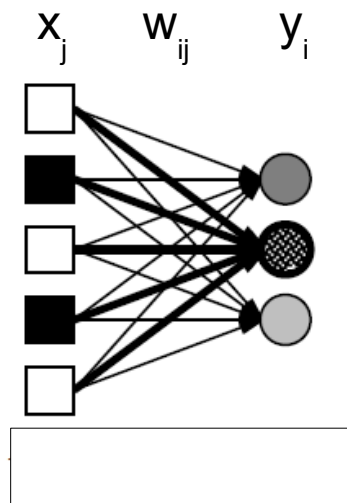
$$\Delta w_{ij}^u = t_i^u x_j^u \Rightarrow w_{ij}^{new} = w_{ij}^{old} + \Delta w_{ij}^u$$

- Si se parte de pesos nulos, el estado final para un peso  $w_{ij}$  es

$$w_{ij}^{final} = 0 + \Delta w_{ij}^1 + \dots + \Delta w_{ij}^p = t_i^1 x_j^1 + \dots + t_i^p x_j^p$$

Matricialmente:  $W = t^1 x^{1T} + t^2 x^{2T} + \dots + t^p x^{pT}$

$$y_i = \sum_{j=0}^n W_{ij} x_j$$



- Asociador lineal (cont.)

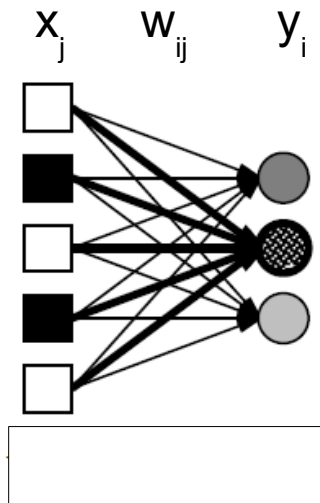
- Aprendizaje Hebbiano (cont.)

- Si los vectores de entrada  $\{\mathbf{x}^1, \mathbf{x}^2, \dots, \mathbf{x}^p\}$  son ortonormales (ortogonales y longitud 1), la red entrenada responderá así:

$$\begin{aligned} W \mathbf{x}^\mu &= (\mathbf{t}^1 \mathbf{x}^{1T} + \dots + \mathbf{t}^p \mathbf{x}^{pT}) \mathbf{x}^\mu = \\ &= \mathbf{t}^1 (\mathbf{x}^{1T} \mathbf{x}^\mu) + \dots + \mathbf{t}^p (\mathbf{x}^{pT} \mathbf{x}^\mu) = \mathbf{t}^\mu \end{aligned}$$

- Problema: condiciones muy restrictivas lo hacen poco útil (entradas ortonormales)

$$y_i = \sum_{j=0}^n W_{ij} x_j$$



- Perceptrón Simple

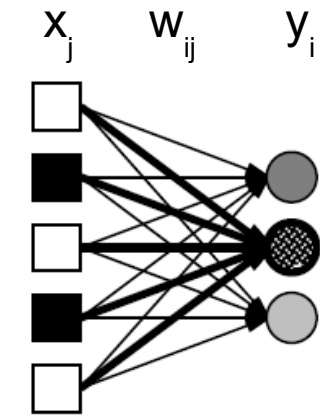
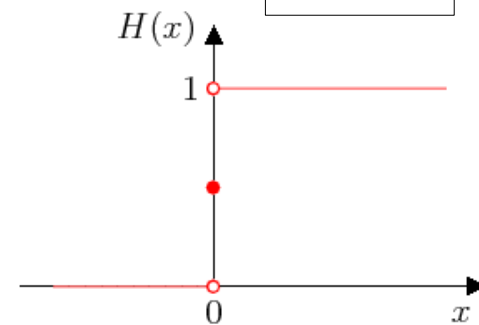
- Red monocapa, como el asociador lineal
- Regla de propagación

$$\sum_{j=0}^n w_{ij} x_j - \theta_i$$

- Función de activación: Heavyside
- Función de salida: Identidad
- Valor de salida (***n*** entradas y ***m*** salidas)

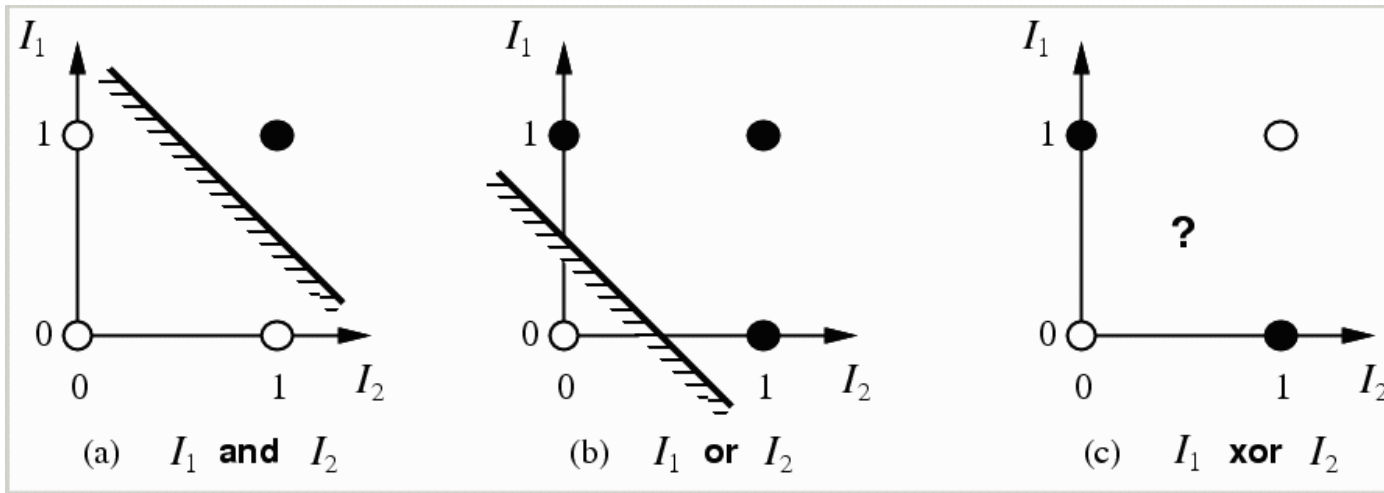
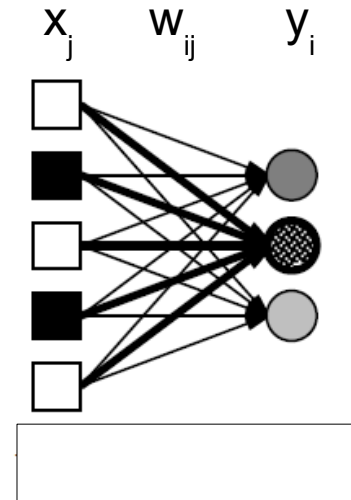
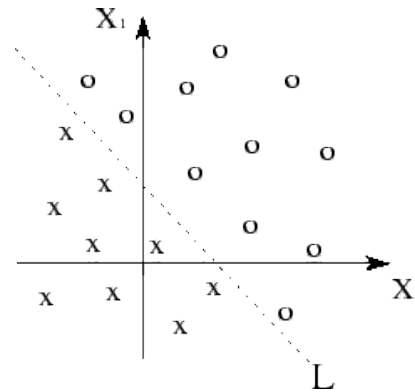
$$y_i = H\left(\sum_{j=0}^n w_{ij} x_j - \theta_i\right), \forall i, 1 \leq i \leq m, \forall j, 1 \leq j \leq n$$

$$H(x) = \begin{cases} 1 & \text{si } \sum_{j=0}^n w_{ij} x_j - \theta_i \geq 0 \\ 0 & \text{si } \sum_{j=0}^n w_{ij} x_j - \theta_i < 0 \end{cases}$$



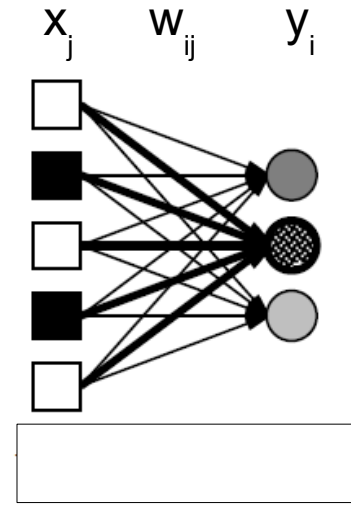
- Perceptrón Simple

- Problema: solo discrimina patrones ***linealmente separables***



- Perceptrón Simple

- Algoritmo de aprendizaje
  - Aprendizaje por corrección de errores
  - Ajuste en base a la diferencia entre salida esperada y salida obtenida
  - Regla de aprendizaje



$$\Delta w_{ij}^u(t) = \varepsilon (t_i^u - y_i^u) x_j^u$$

- El error que comete la red se va reduciendo hasta proporcionar la salida deseada (si se cumplen las condiciones necesarias de separación lineal)

- ADALINE (ADaptive Linear Element)

- Similar al perceptrón, pero la función de activación es la identidad
- Salida (***n*** neuronas de entrada y ***m*** de salida)

$$y_i(t) = \sum_{j=1}^n w_{ij} x_j - \theta_i, \quad \forall i \ 1 \leq i \leq m$$

- Limitación: al tener neuronas de respuesta lineal, sólo podrá separar patrones ***linealmente independientes***
- La mayor ventaja respecto de las anteriores es el algoritmo de aprendizaje, basado en la minimización del error cuadrático medio (***LMS, Least Mean Squares***)

- ADALINE (cont.)
  - Algoritmo de aprendizaje
    - Descenso por el gradiente
    - La función a optimizar es una función de costo  $E(W)$ , donde  $W$  son los pesos sinápticos
    - Actualización de pesos:

$$W(t+1) = W(t) - \varepsilon \nabla E(W)$$

- Compromiso en selección de  $\varepsilon$ : aprendizaje lento vs. oscilación cerca de min/max



- ADALINE (cont.)
  - Algoritmo de aprendizaje (cont.)
    - Regla **LMS** o Regla **Widrow-Hoff**

$$E[w_{ij}] = \frac{1}{2} \sum_{\mu=1}^p \sum_{i=1}^m (t_i^{\mu} - y_i^{\mu})^2, \text{ con } y_i^{\mu} = \sum_{j=1}^n w_{ij} x_j^{\mu} - \theta_i^{\mu}$$

$$\begin{aligned} \frac{\partial E[w_{ij}]}{\partial w_{ij}} &= (1/2) 2 \sum_{\mu=1}^p (t_i^{\mu} - y_i^{\mu}) \left( -\frac{dy_i^{\mu}}{dw_{ij}} \right) = \\ &= -(1/2) 2 \sum_{\mu=1}^p (t_i^{\mu} - y_i^{\mu}) \frac{dy_i^{\mu}}{dw_{ij}} = - \sum_{\mu=1}^p (t_i^{\mu} - y_i^{\mu}) x_j^{\mu} \end{aligned}$$

Por lo tanto

$$\Delta w_{ij} = -\varepsilon \frac{\partial E[w_{ij}]}{\partial w_{ij}} = \varepsilon \sum_{\mu=1}^p (t_i^{\mu} - y_i^{\mu}) x_j^{\mu}$$

- ADALINE (cont.)

- Dado que las neuronas son lineales, la función **E(W)** es cuadrática en los pesos, permitiendo alcanzar siempre un mínimo global (es un paraboloide)

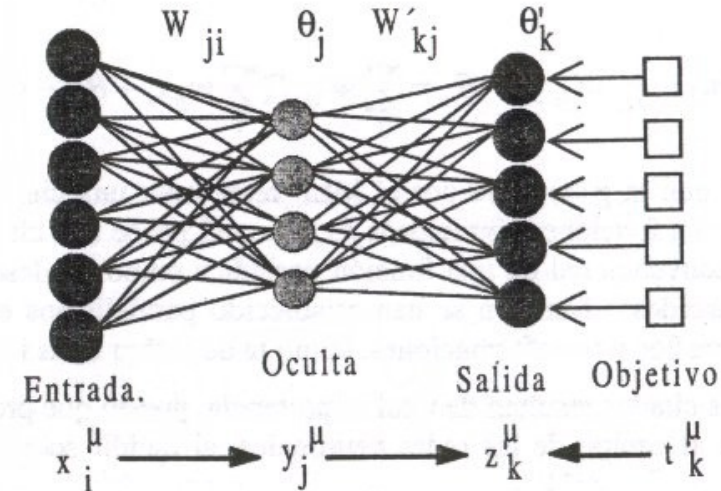
$$E[w_{ij}] = \frac{1}{2} \sum_{\mu=1}^p \sum_{i=1}^n (t_i^{\mu} - y_i^{\mu})^2$$

- Existe una versión multicapa de la ADALINE llamada MADALINE (Many ADALINEs) para superar las limitaciones como red monocapa

- MLP (Multi-Layer Perceptron)

- Es una red multicapa
- Regla de propagación  
(usualmente igual en todas las capas)

$$\sum_{i=1}^n w_{ji} x_i - \theta_j$$

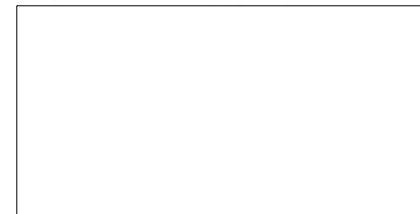
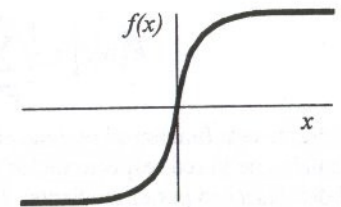


- Función de activación: suele ser lineal en la capa de salida, y sigmoide en la capa oculta

Ej:

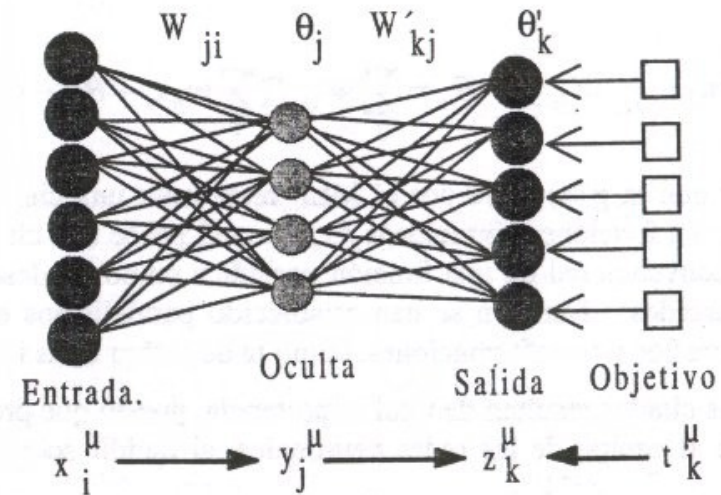
$$f(x) = \frac{1}{1 + e^{-x}}$$

$$f(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}} = \tanh(x)$$



- MLP (Multi-Layer Perceptron)
  - Salida de la red

$$z_k = g \left( \sum_j w'_{kj} y_j - \theta'_k \right) =$$
$$= g \left( \sum_j w'_{kj} f \left( \sum_i w_{ji} x_i - \theta_j \right) - \theta'_k \right)$$



donde  $g()$  es la función de activación de las neuronas de salida, y  $f()$  es la función de activación de las neuronas de la capa oculta

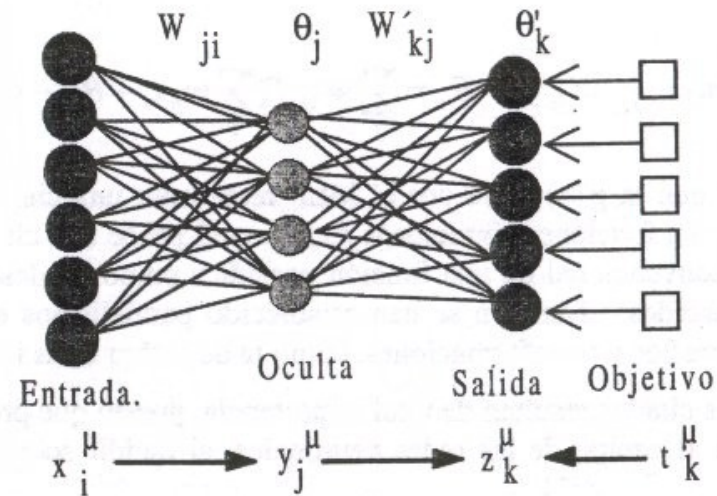
- MLP (Multi-Layer Perceptron)
  - Algoritmo de aprendizaje: BP (backpropagation)

$$z_k^u = g \left( \sum_j w'_{kj} \underbrace{f \left( \sum_i w_{ji} x_i^u - \theta_j \right)}_{y_j^u} - \theta'_k \right)$$

$$E(w_{ji}, w'_{kj}) = \frac{1}{2} \sum_u \sum_k \left[ t_k^u - g \left( \sum_j w'_{kj} y_j^u - \theta'_k \right) \right]^2$$

$$\delta w'_{kj} = -\varepsilon \frac{\partial E}{\partial w'_{kj}}$$

$$\delta w_{ji} = -\varepsilon \frac{\partial E}{\partial w_{ji}}$$



- MLP (Multi-Layer Perceptron)

- Algoritmo de aprendizaje: BP (backpropagation)  
Derivando respecto a los pesos y aplicando la regla de la cadena se tiene

$$\delta w'_{kj} = -\varepsilon \frac{\partial E}{\partial w'_{kj}} = \varepsilon \sum_{\mu} \Delta'^{\mu}_k y^{\mu}_j$$

$$\Delta'^{\mu}_k = [t^{\mu}_k - g(h'^{\mu}_k)] \frac{\partial g(h'^{\mu}_k)}{\partial h'^{\mu}_k}$$

$$h'^{\mu}_k = \sum w'_{kj} y^{\mu}_j - \theta'_k$$

$$\delta w_{ji} = -\varepsilon \frac{\partial E}{\partial w_{ji}} = \varepsilon \sum_{\mu} \Delta^{\mu}_j x^{\mu}_i$$

$$\Delta^{\mu}_j = \left[ \sum_k \Delta'^{\mu}_k w'_{kj} \right] \frac{\partial f(h^{\mu}_j)}{\partial h^{\mu}_j}$$

$$h^{\mu}_j = \sum_i w_{ji} x^{\mu}_i - \theta_j$$

- MLP (Multi-Layer Perceptron)
  - Cantidad de ejemplos de entrenamiento
    - Para una red con  $w$  pesos sinápticos se requiere aproximadamente una cantidad de patrones  $p=w/\epsilon$  para obtener un error del orden de  $\epsilon$ .

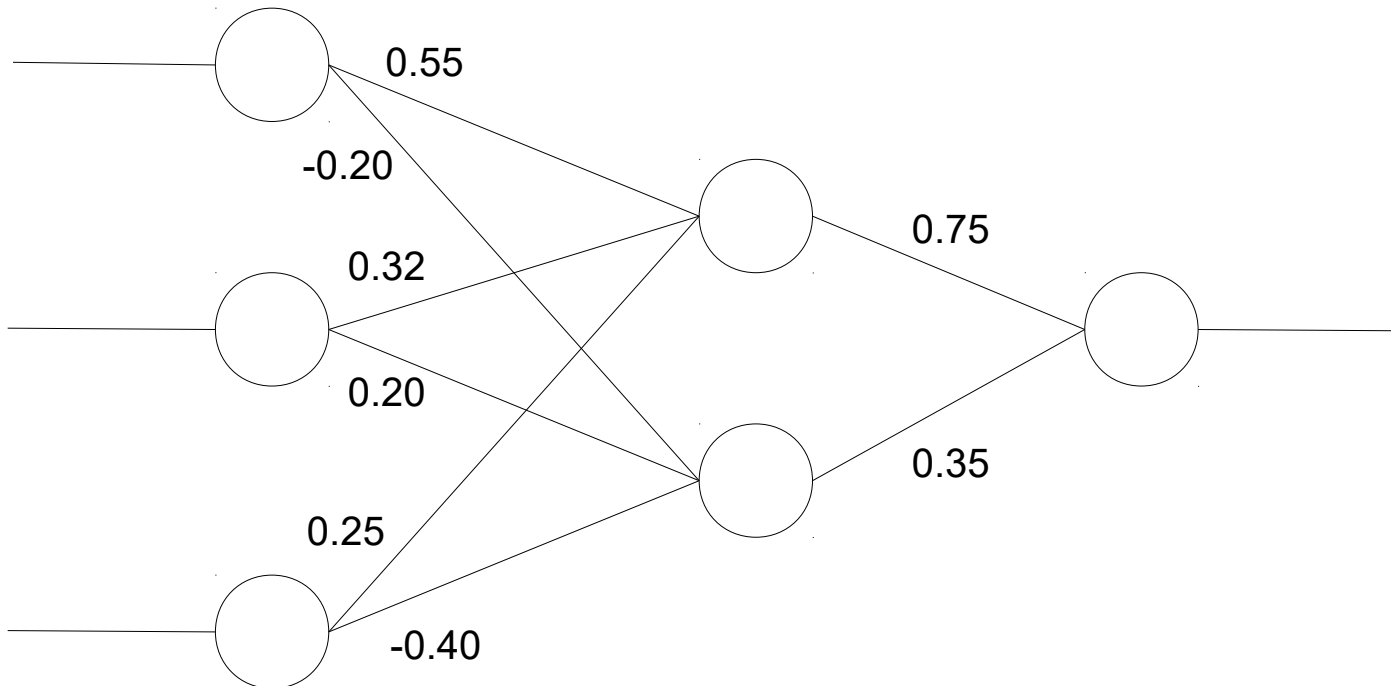
- MLP: Un ejemplo

- **Arquitectura:**

- 3 entradas, 1 salida, 1 capa oculta

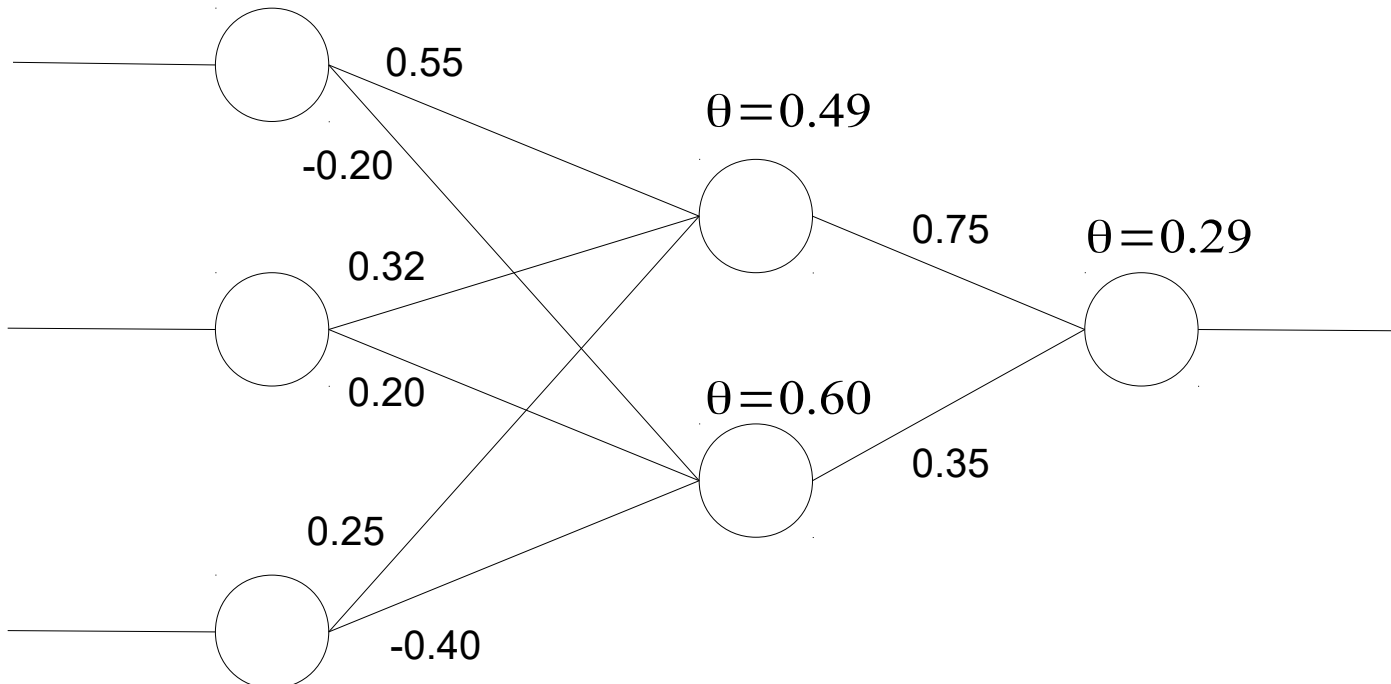
- Función de activación capa oculta y de salida:  $f(x) = g(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}} = \tanh(x)$

- Derivadas:  $f'(x) = g'(x) = 1 - \tanh^2(x)$





- MLP: Un ejemplo
  - **Estado inicial de la red:** pesos aleatorios, pequeños, positivos y negativos

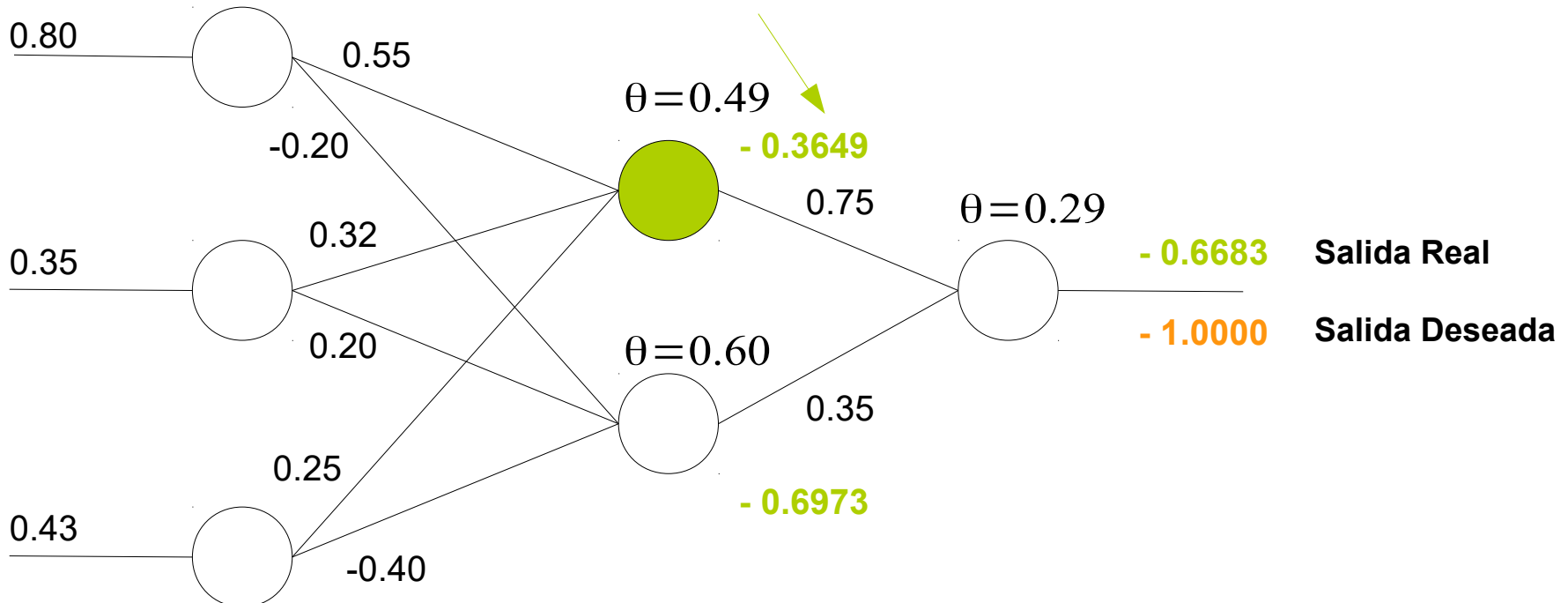


- MLP: Un ejemplo

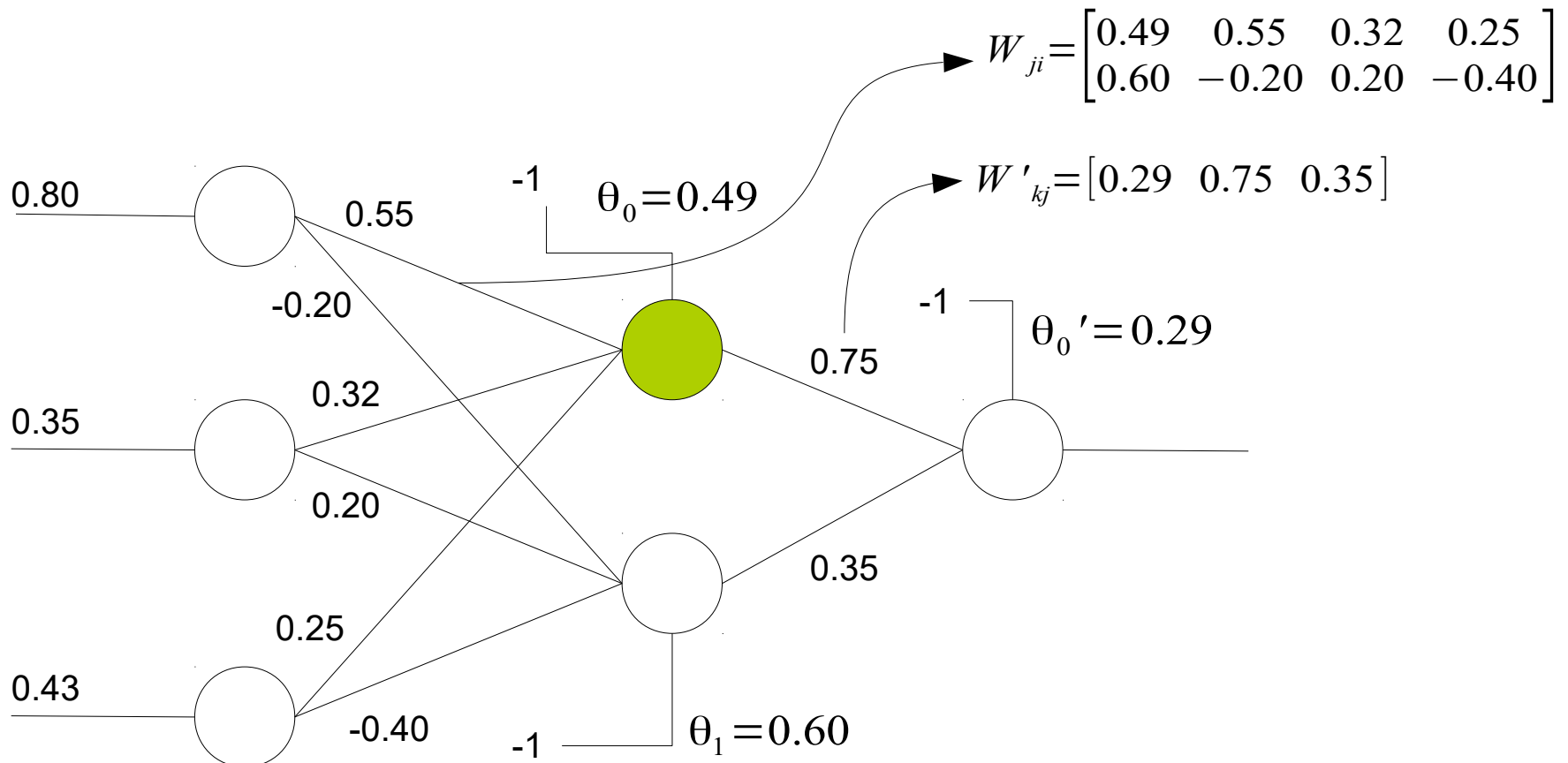
- Entrenamiento: (1) Cálculo de la salida de la red**

**Regla de Propagación:**  $\sum_{i=1}^n w_{ji} x_i - \theta_j = 0.55 * 0.80 + 0.32 * 0.35 + 0.25 * 0.43 - 0.49 = -0.3825$

**Función de Activación:**  $\tanh(-0.3825) = -0.3649$



- MLP: Un ejemplo
  - **Entrenamiento: (2) Backpropagation**



## • MLP: Un ejemplo

### • Entrenamiento: (2.a) Backpropagation para $w'_{kj}$

$$h'_0 = \sum_j w_{0j}' y_j - \theta_0' = [0,75 \times (-0,3649)] + [0,35 \times (-0,6973)] - 0,29 = -0,8077$$

$$g'(h'_k) = 1 - \tanh(h'_k) \quad \leftarrow$$

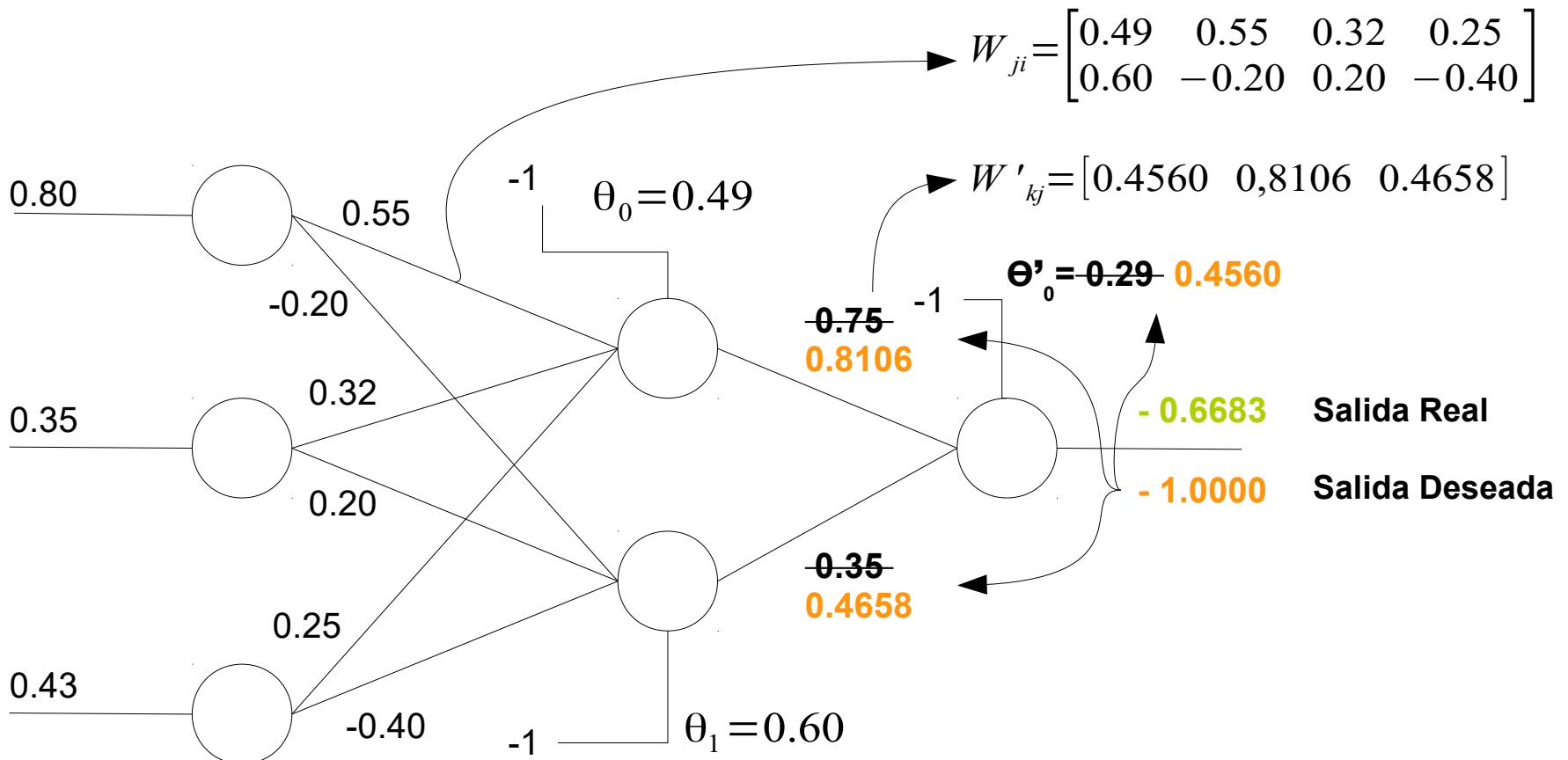
$$\Delta'_0 = [t_0 - g(h'_0)] \frac{\partial g(h'_0)}{\partial h'_0} = (-1.000 - (-0.6683)) (1 - \tanh(-0.8077)) = -0.5534$$

$$\delta w'_{0j} = -\varepsilon \frac{\partial E}{\partial w'_{0j}} = \varepsilon \Delta'_0 y_j \quad \left\{ \begin{array}{l} \delta w'_{00} = \varepsilon \Delta'_0 \theta'_0 = 0.30 (-0.5534) (-1.0000) = 0,1660 \\ \delta w'_{01} = \varepsilon \Delta'_0 y_0 = 0.30 (-0.5534) (-0.3649) = 0,0606 \\ \delta w'_{02} = \varepsilon \Delta'_0 y_1 = 0.30 (-0.5534) (-0.6973) = 0,1158 \\ \text{donde } \varepsilon = 0.30 \end{array} \right.$$

$$W'_{kj}(t+1) = W'_{kj}(t) + \delta w'_{kj} = [0.29 \quad 0.75 \quad 0.35] + [0.1660 \quad 0.0606 \quad 0.1158] = [0,4560 \quad 0,8106 \quad 0,4658]$$

- MLP: Un ejemplo

- Entrenamiento: (2.a) Backpropagation para  $w'_{kj}$



## • MLP: Un ejemplo

### • **Entrenamiento:** (2.b) Backpropagation para $w_{ji}$

Primera neurona de la capa oculta ( $j=0$ )

$$h_0 = \sum_i w_{0i} x_i - \theta_0 = 0.55 \times 0.80 + 0.32 \times 0.35 + 0.25 \times 0.43 - 0.49 = 0.1695$$

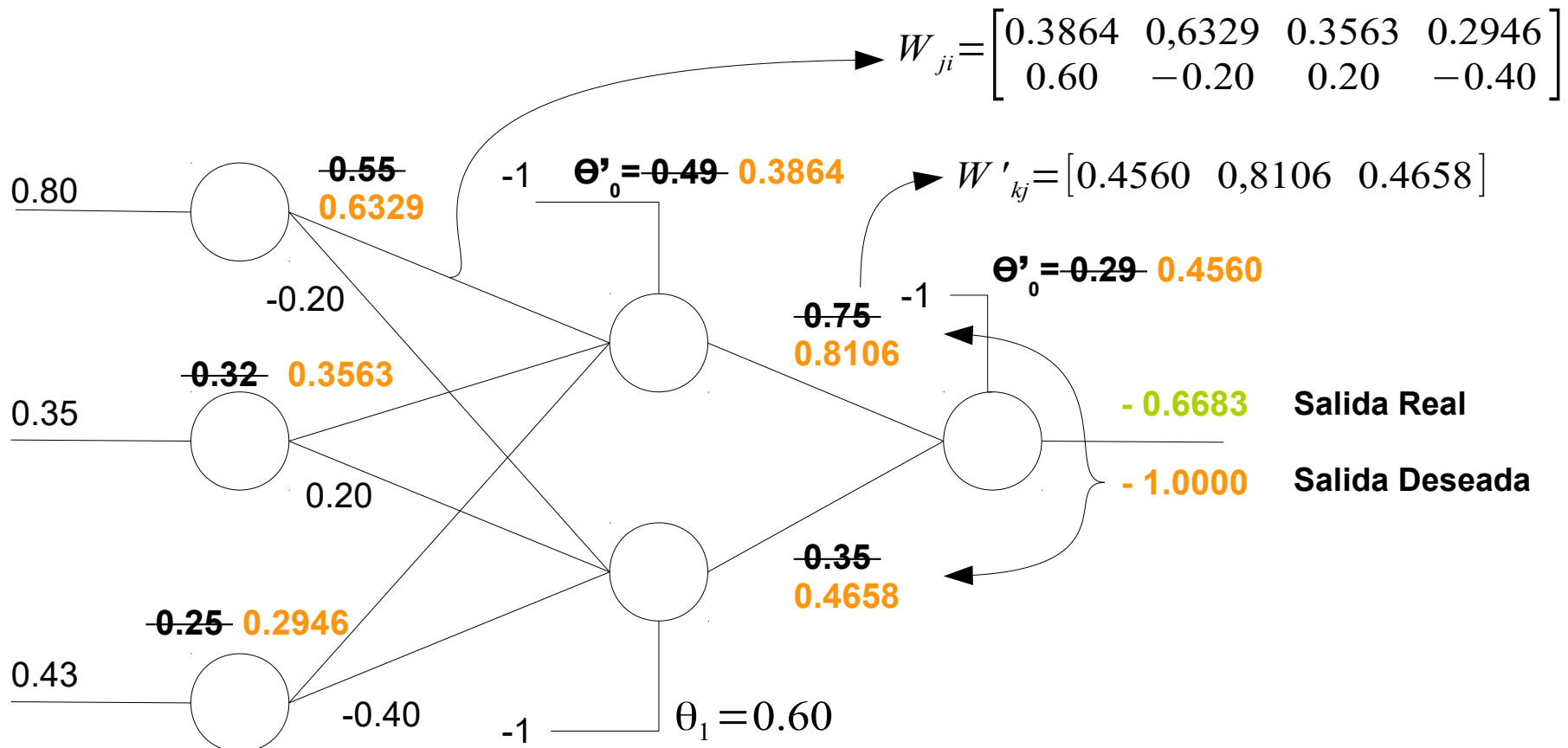
$$\Delta_0 = \left[ \sum_k \Delta'_k w'_{k0} \right] \frac{\partial f(h_0)}{\partial h_0} = \Delta'_0 w'_{00} f'(h_0) = \overset{\Delta'_0}{(-0.5534)} \overset{w'_{00}}{0.75} \overset{f'(h'_0)}{(1 - \tanh(0.1695))} = -0.3454$$

$$\delta w_{0i} = -\varepsilon \frac{\partial E}{\partial w_{0i}} = \varepsilon \Delta_0 x_i \quad \left\{ \begin{array}{l} \delta w_{00} = \varepsilon \Delta_0 \theta_0 = 0.30 (0.3454) (-1) = -0.1036 \\ \delta w_{01} = \varepsilon \Delta_0 x_0 = 0.30 (0.3454) (0.80) = 0.0829 \\ \delta w_{02} = \varepsilon \Delta_0 x_1 = 0.30 (0.3454) (0.35) = 0.0363 \\ \delta w_{03} = \varepsilon \Delta_0 x_2 = 0.30 (0.3454) (0.43) = 0.0446 \\ \text{donde } \varepsilon = 0.30 \end{array} \right.$$

$$W_{0i}(t+1) = W_{0i}(t) + \delta w_{0i} = [0.49 \ 0.55 \ 0.32 \ 0.25] + [-0.1036 \ 0.0829 \ 0.0363 \ 0.0446] = [0.3864 \ 0.6329 \ 0.3563 \ 0.2946]$$

- MLP: Un ejemplo

- Entrenamiento: (2.a) Backpropagation para  $w'_{kj}$



## • MLP: Un ejemplo

### • **Entrenamiento:** (2.b) Backpropagation para $w_{ji}$

Segunda neurona de la capa oculta ( $j=1$ )

$$h_1 = \sum_i w_{1i} x_i - \theta_1 = -0.20 \times 0.80 + 0.20 \times 0.35 + 0.25 \times -0.40 - 0.60 = -0.79$$

$$\Delta_1 = \left[ \sum_k \Delta'_k w'_{k1} \right] \frac{\partial f(h_1)}{\partial h_1} = \Delta'_0 w'_{01} f'(h_1) = \overset{\Delta'_0}{(-0.5534)} \overset{w'_{01}}{0.35} \overset{f'(h'_1)}{(1 - \tanh(-0.79))} = -0.3212$$

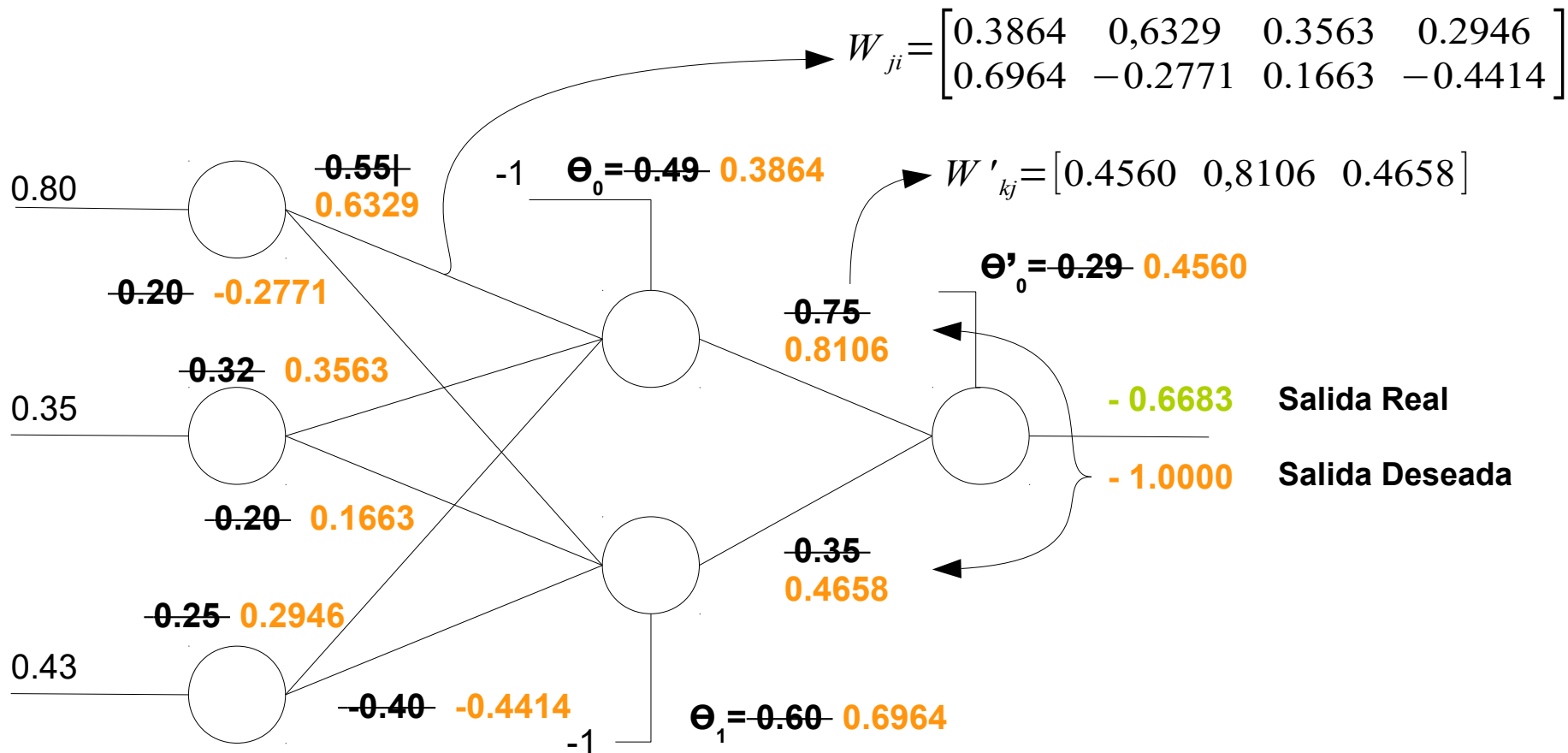
$$\delta w_{1i} = -\varepsilon \frac{\partial E}{\partial w_{1i}} = \varepsilon \Delta_1 x_i \quad \left\{ \begin{array}{l} \delta w'_{10} = \varepsilon \Delta_1 \theta_1 = 0.30 (-0.3212) (-1.0) = 0.0964 \\ \delta w'_{11} = \varepsilon \Delta_1 x_0 = 0.30 (-0.3212) (0.80) = -0.0771 \\ \delta w'_{12} = \varepsilon \Delta_1 x_1 = 0.30 (-0.3212) (0.35) = -0.0337 \\ \delta w'_{13} = \varepsilon \Delta_1 x_2 = 0.30 (-0.3212) (0.43) = -0.0414 \\ \text{donde } \varepsilon = 0.30 \end{array} \right.$$

$$W_{1i}(t+1) = W_{1i}(t) + \delta w_{1i} = [0.60 \quad -0.20 \quad 0.20 \quad -0.40] + [0.0964 \quad -0.0771 \quad -0.0337 \quad -0.0414] = [0.6964 \quad -0.2771 \quad 0.1663 \quad -0.4414]$$

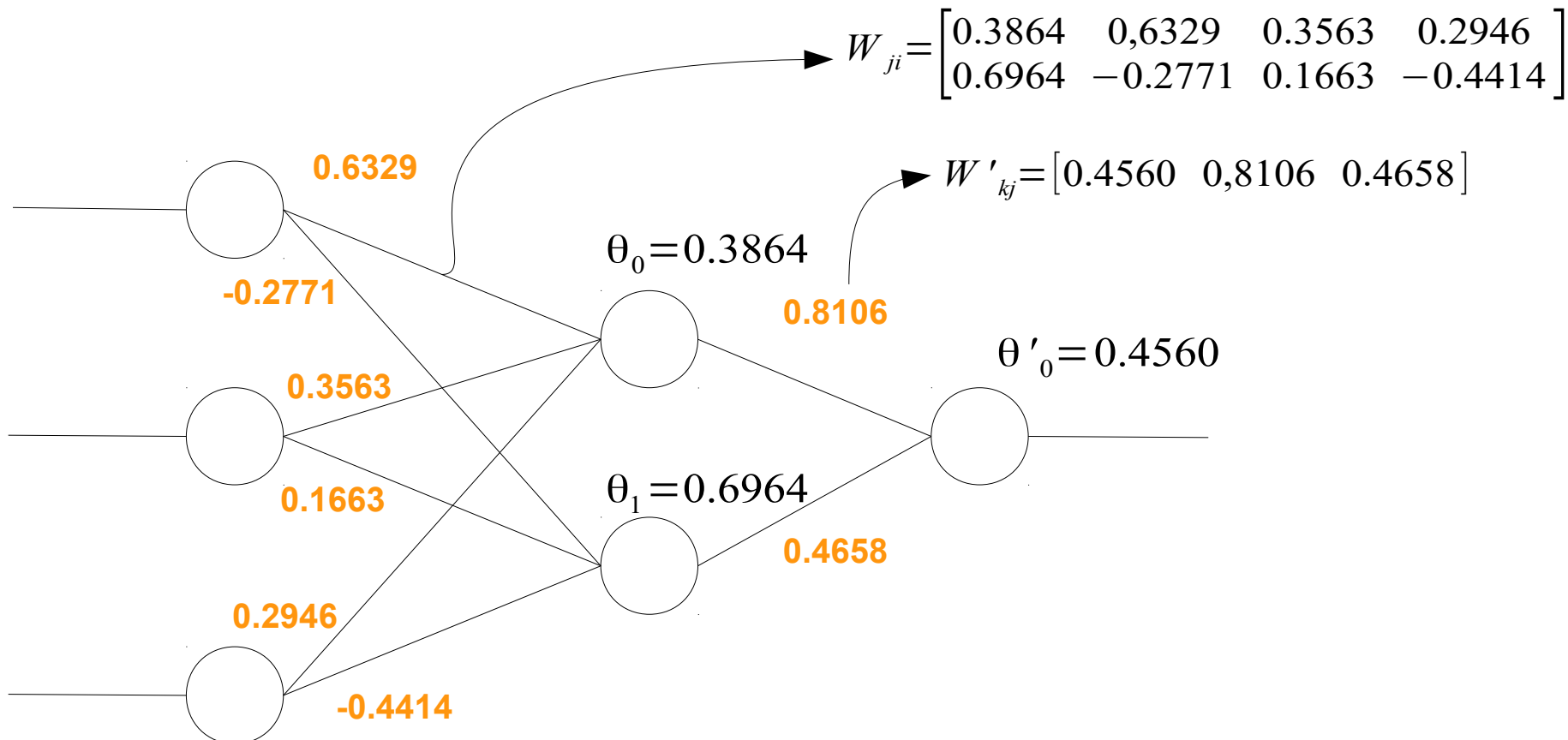


- MLP: Un ejemplo

- Entrenamiento: (2.b) Backpropagation para  $w_{ji}$



- MLP: Un ejemplo
  - Estado de la red luego de entrenarla con 1 ejemplo



- Aspectos prácticos de la aplicación de un ANS
  - Definir claramente el problema a resolver (tipo de entradas, tipo de salidas que se desea obtener, etc)
  - Revisión bibliográfica
  - Elegir adecuadamente la cantidad de entradas y salidas (muy pocas limitará la capacidad de generalización; demasiadas podría generar asociaciones inválidas entre entradas y salidas) → Dilema Sesgo vs. Varianza
  - Elección de conjuntos de aprendizaje y test
    - El set de entrenamiento debe cubrir razonablemente el fenómeno que se desea modelar
    - Se recomienda validación cruzada con parada temprana para evitar overfitting
    - Los ejemplos de aprendizaje no deben incluirse en el conjunto de test

- Aspectos prácticos de la aplicación de un ANS
  - Preprocesamiento de las entradas
    - Se aconseja: buena distribución, entradas con rangos parecidos y dentro del intervalo de trabajo de la función de activación de la capa oculta
    - Las entradas pueden escalarse a un intervalo  $[0, 1]$  o  $[-1, 1]$ .  
Enfoques
      - Por patrones (por filas): se normaliza el valor de los componentes del vector de entrada al intervalo  $[0, 1]$  linealmente en base al máximo y mínimo del patrón (**solo aplicable cuando todas las entradas tienen rangos de valores parecidos**)
      - Global por columnas: se normalizan los valores de cada entrada, de acuerdo a los máximos y mínimos de esa entrada a lo largo de todos los patrones

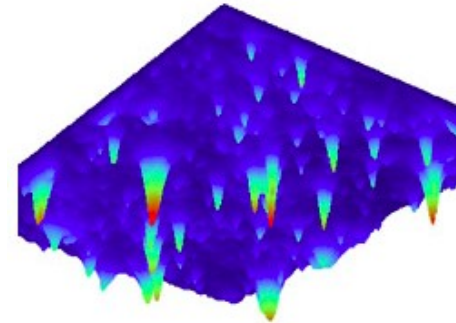
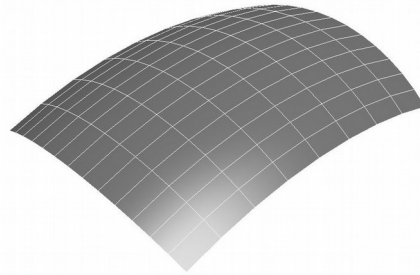
- Aspectos prácticos de la aplicación de un ANS
  - Preprocesamiento de las entradas
    - Un procesamiento que suele dar buenos resultados:
      - Estandarizar: transformar las entradas para que tengan media 0 y varianza 1.  
Cálculo: por cada componente del vector de entrada, calcular:

$$x' = \frac{x - \bar{x}}{\sigma}$$

- Luego de estandarizar, escalar al intervalo [0,1] o [-1,1] mediante una transformación lineal

- Aspectos prácticos de la aplicación de un ANS
  - Entrenamiento
    - Inicialización de pesos: normalmente aleatorios, pequeños, positivos y negativos (ej:  $[-0.3, 0.3]$ ,  $[-1.1]$ , etc)
    - Ritmo de aprendizaje lo más grande posible que no cause oscilaciones en el algoritmo. Se recomienda un ritmo de aprendizaje variable (se verá más adelante)
    - Neuronas ocultas: idealmente, la menor cantidad posible que produzca una buena respuesta (navaja de Ockham). Se recomienda probar con pocas (1 o 2 capas ocultas, con unas pocas neuronas cada una), y experimentar agregando paulatinamente neuronas a cada capa

- Velocidad de convergencia del algoritmo BP
  - Se ve afectada por algunos problemas
    - En superficies con poca pendiente, el avance es lento
    - En superficies con mucha pendiente, hay riesgo de oscilación entre picos



- En zonas donde las derivadas de las funciones de activación son pequeñas, los ajustes de los pesos son pequeños (el tamaño del ajuste depende del valor de las derivadas)

- Mejoras para la velocidad de convergencia y estabilidad del algoritmo BP: actualización de pesos
  - Métodos de primer orden: siguen usando solo las derivadas primeras. Ej: Momentum, SAB, SuperSAB, RPROP
  - Métodos de segundo orden: usan información de las derivadas segundas. Ej: Levenberg-Marquardt
  - Ninguno de ellos domina a los otros
    - En algunos casos alguno funciona mejor que otros, y en otros casos ocurre lo contrario
    - Se debe experimentar para determinar cuál conviene en cada problema



- Actualización de pesos – Métodos de 1er orden
  - Momentum
    - Añadir al cálculo de variación de pesos un **término de inercia**, proporcional al incremento en la iteración anterior
    - Utiliza los mismos parámetros para todos los pesos
    - Tiende a acelerar las variaciones en zonas con poca pendiente, y a hacerlas más suaves en zonas escarpadas (reduciendo/evitando oscilaciones)

$$\delta w'_{kj}(t) = -\varepsilon \frac{\partial E}{\partial w'_{kj}}(t) + \alpha \delta w'_{kj}(t-1)$$

$$\delta w_{ji}(t) = -\varepsilon \frac{\partial E}{\partial w_{ji}}(t) + \alpha \delta w_{ji}(t-1)$$

$$\alpha \in [0, 1]$$

- Actualización de pesos – Métodos de 1er orden
  - SuperSAB (Super Self-Adapting Backpropagation)
    - Agrega inercia y ritmo de aprendizaje independiente para cada peso
    - El ajuste de pesos y el cambio en el ritmo de aprendizaje dependen del signo del gradiente en iteraciones continuas
    - El ajuste de pesos depende además de la magnitud del gradiente

$$\Delta w_{ij}(t) = \begin{cases} -\varepsilon_{ij}(t) \frac{\partial E}{\partial w_{ij}}(t) + \alpha \Delta w_{ij}(t-1) - \text{decay } w_{ij}(t), & \text{si } \frac{\partial E}{\partial w_{ij}}(t) \frac{\partial E}{\partial w_{ij}}(t-1) \geq 0 \\ 0 & \text{en otro caso (es decir, } w_{ij}(t) = w_{ij}(t-1) \text{)} \\ \text{donde decay es un parámetro del modelo} \end{cases}$$
$$\varepsilon_{ij}(t) = \begin{cases} \varepsilon_{ij}(t-1)u, & \text{si } \frac{\partial E}{\partial w_{ij}}(t) \frac{\partial E}{\partial w_{ij}}(t-1) \geq 0 \quad \text{y} \quad \varepsilon_{ij}(t-1) \leq \varepsilon_{\max} \\ \varepsilon_{ij}(t-1)d, & \text{en otro caso} \\ \text{donde } u \in [1.1, 1.3] \quad \text{y} \quad d \in [0.7, 0.9] \end{cases}$$

- Actualización de pesos – Métodos de 1er orden
  - RPROP (Resilient Backpropagation) – Sin weight Backtracking
    - Sólo utiliza el signo del gradiente para el tamaño del ajuste de los pesos
    - Es independiente de la magnitud del gradiente

$$\Delta_{ij}(t) = \begin{cases} \eta^+ \Delta_{ij}(t-1), & \text{si } \frac{\partial E}{\partial w_{ij}}(t-1) \frac{\partial E}{\partial w_{ij}}(t) > 0 \\ \eta^- \Delta_{ij}(t-1), & \text{si } \frac{\partial E}{\partial w_{ij}}(t-1) \frac{\partial E}{\partial w_{ij}}(t) < 0 \\ \Delta_{ij}(t-1) & \text{en otro caso} \end{cases}$$

donde  $0 < \eta^- < 1 < \eta^+$  (originalmente 0.5 y 1.2 respectivamente)

$$\delta w_{ij}(t) = \begin{cases} -\Delta_{ij}(t) & \text{si } \frac{\partial E}{\partial w_{ij}}(t) > 0 \\ +\Delta_{ij}(t) & \text{si } \frac{\partial E}{\partial w_{ij}}(t) < 0 \\ 0 & \text{en otro caso} \end{cases}$$

**Parámetros necesarios**

$$\Delta_0 \quad \Delta_{max}$$

**Parámetro opcional**

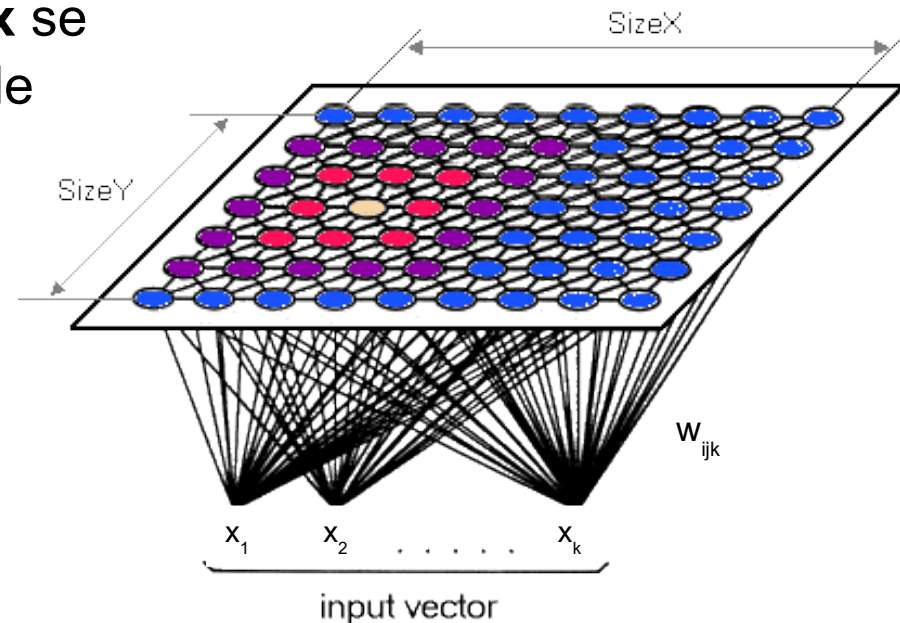
$$\Delta_{min}$$

# Redes Neuronales No Supervisadas

- Redes neuronales No supervisadas
  - Característica principal: durante el entrenamiento no se asocia la salida deseada con los patrones presentados
  - Estas redes se **autoorganizan** para agrupar patrones con ciertas similitudes
  - Algunas aplicaciones:
    - Análisis de similitud/novedad de patrones
    - Análisis de componentes principales
    - Clustering/Clasificación
    - Memoria asociativa
    - Codificación
    - Mapas de rasgos
    - Análisis exploratorio

- Redes neuronales No supervisadas
  - Suelen ser monocapa
  - Pueden combinarse con redes supervisadas, constituyendo redes híbridas (usualmente estas redes incluyen capas supervisadas y capas no supervisadas)
  - Hay 2 tipos
    - Redes no supervisadas Hebbianas: múltiples neuronas pueden quedar activadas
    - Redes no supervisadas Competitivas: sólo una neurona (o un pequeño grupo) pueden quedar activadas
      - También se las denomina WTA: Winner Takes All
  - En los modelos competitivos, la/s neurona/s ganadora/s se premian reforzando sus conexiones sinápticas

- SOFM (Self-Organizing Feature Maps)
  - 2 capas: una de entrada y otra de procesamiento/salida
  - La capa de entrada tiene un vector de  $k$  neuronas, cada una representando un elemento del vector de entrada  $\mathbf{x}$
  - La capa de procesamiento usualmente es un rectángulo con  $i \times j$  neuronas
  - Cada elemento  $x_k$  de la entrada  $\mathbf{x}$  se vincula con todas las neuronas de la capa de salida con un peso sináptico  $w_{ijk}$



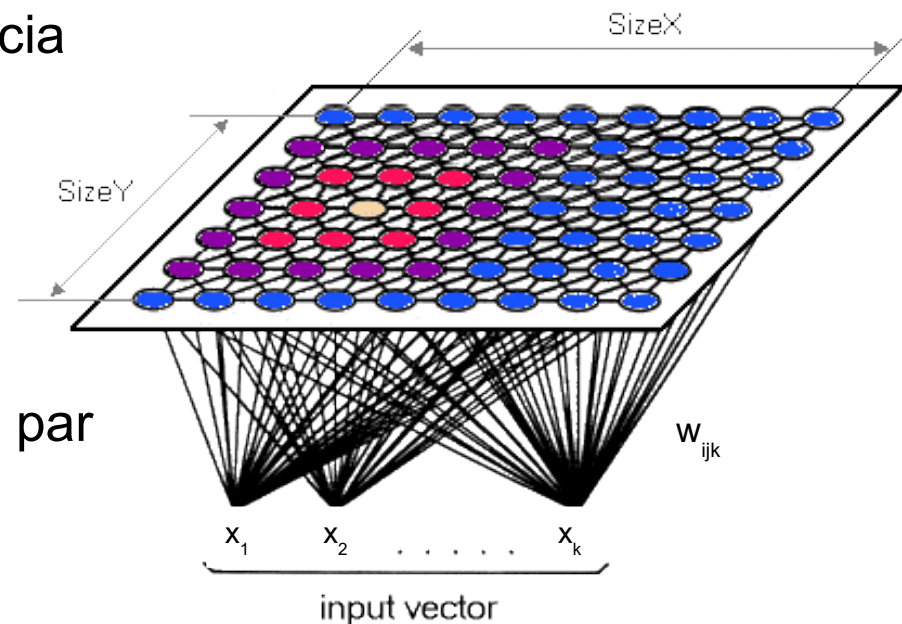
- SOFM (Self-Organizing Feature Maps)

- **Fase de ejecución**

- Los pesos permanecen fijos (aprendizaje se “desenchufa”)
    - Cada neurona **(i, j)** calcula la **similitud** del vector de entrada **x** con su vector de pesos sinápticos **w<sub>ij</sub>**
    - Existen diversas **funciones de similitud**
    - Gana la neurona cuya distancia sea menor

$$d(\mathbf{w}_g, \mathbf{x}) = \min_{ij} \{d(\mathbf{w}_{ij}, \mathbf{x})\}$$

- La neurona ganadora **g** es el par ij que minimiza la distancia

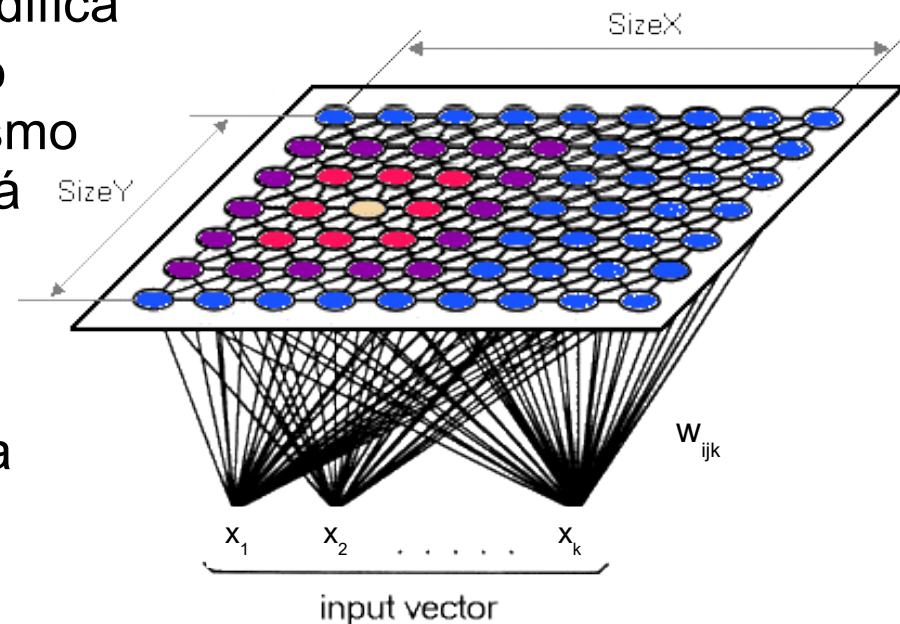




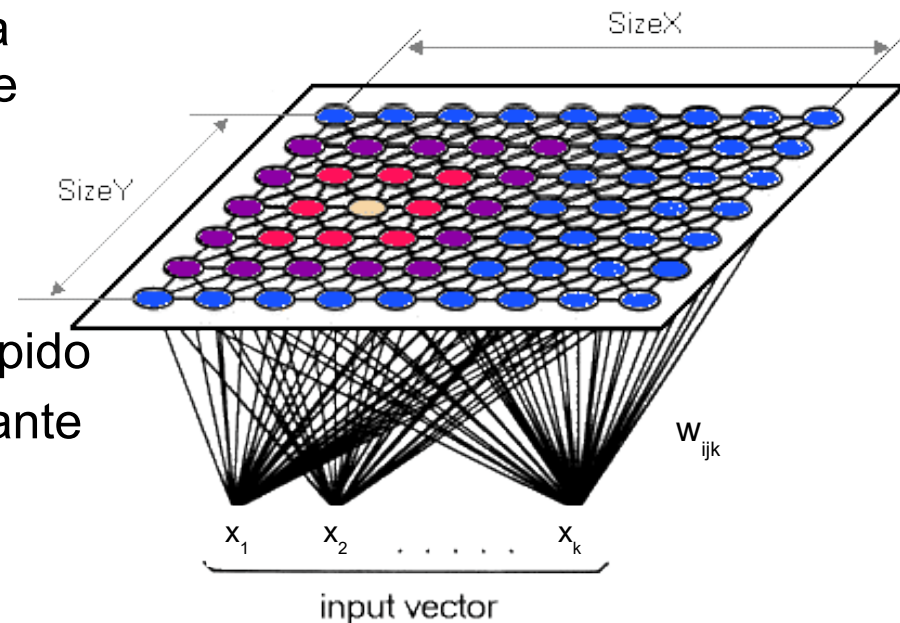
- SOFM (Self-Organizing Feature Maps)

- **Fase de aprendizaje**

- Cada neurona se “especializa” en un grupo o familia de patrones
    - Para el aprendizaje, se le presentan a la red múltiples patrones en secuencia
    - Por cada patrón, se calcula la neurona ganadora, y el vector de pesos  $w_{ij}$  de la misma se modifica para que se parezca un poco más a la entrada (ante el mismo patrón, la neurona se activará con más intensidad en iteraciones subsiguientes)
    - Esto representa un modelo competitivo sencillo (una sola neurona ganadora)



- SOFM (Self-Organizing Feature Maps)
  - **Fase de aprendizaje**
    - El modelo de Kohonen incorpora una **función de vecindad** que indica las neuronas vecinas cuyos pesos también se actualizarán
    - Patrones similares en el espacio sensorial tendrán representaciones cercanas en la superficie de la red
    - La función de vecindad premia a las neuronas vecinas a la ganadora (interacciones laterales) e inhibe a las alejadas
    - El uso de esta función de vecindad tiene 2 ventajas:
      - Ritmo de convergencia + rápido
      - El sistema es más robusto ante variaciones en los valores iniciales



- SOFM (Self-Organizing Feature Maps)
  - **Procedimiento general de aprendizaje**
    - Inicialización de pesos  $\mathbf{w}_{ijk}$ . Pueden ser nulos, aleatorios (pequeños) o con un valor determinado
    - En cada iteración  $t$ 
      - Presentación del patrón  $\mathbf{x}(t)$  tomado al azar del conjunto de entrenamiento
      - Cada neurona calcula la similitud entre su vector de pesos  $\mathbf{w}_{ij}$  y la entrada  $\mathbf{x}$ . **Ejemplo:** distancia euclídea

$$d^2(\mathbf{w}_{ij}, \mathbf{x}) = \sum_{k=1}^n (w_{ijk} - x_k)^2$$

- Determinación de la **neurona ganadora**

- SOFM (Self-Organizing Feature Maps)
  - **Procedimiento general de aprendizaje (cont.)**

- En cada iteración  $t$  (cont.)

- Actualización de los pesos de la neurona ganadora y vecinas. Ej:

$$w_{ijk}(t+1) = w_{ijk}(t) + \alpha(t) h\left(\left| \underbrace{\mathbf{i}}_{(i,j)} - \underbrace{\mathbf{g}}_{(m,n)} \right|, t\right) (x_k(t) - w_{ijk}(t))$$

donde

- $\alpha(t)$  es el ritmo de aprendizaje,
      - $h()$  es la función de vecindad que depende de la distancia de las neuronas  $\mathbf{i}$  y  $\mathbf{g}$  (ganadora), que vale 0 cuando  $\mathbf{i}$  no pertenece a la vecindad de  $\mathbf{g}$ , y un  $n^\circ$  positivo cuando sí es vecina.
      - $\alpha()$  y  $h()$  disminuyen monótonamente con  $t$
    - Luego se puede realizar una segunda etapa de **ajuste fino** del mapa, tomando  $\alpha()$  un valor pequeño constante y  $h() = 1$  solo para la ganadora
    - Entre 50 y 100 iteraciones por neurona suelen ser suficientes, pero idealmente deberían ser 500 por neurona

- SOFM (Self-Organizing Feature Maps)
  - **Procedimiento general de aprendizaje (*cont.*)**

- Ritmo de aprendizaje  $\alpha()$ 
  - Son habituales 2 opciones:
    - decreciente linealmente

$$\alpha(t) = \alpha_0 + (\alpha_f - \alpha_0) \frac{t}{t_\alpha}$$

- decreciente exponencialmente

$$\alpha(t) = \alpha_0 \left( \frac{\alpha_f}{\alpha_0} \right)^{t/t_\alpha}$$

donde  $\alpha_0$  es el ritmo inicial ( $|\alpha_0| < 1$ ),  $\alpha_f$  es el ritmo final ( $\sim 0.01$ ) y  $t_\alpha$  el máximo número de iteraciones para llegar a  $\alpha_f$

- SOFM (Self-Organizing Feature Maps)

- **Procedimiento general de aprendizaje (cont.)**

- Función de vecindad  $h(|\mathbf{i} - \mathbf{g}|, t)$

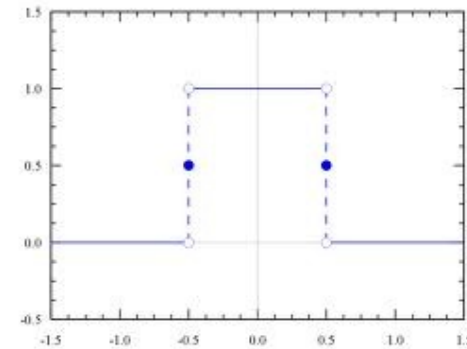
- Usualmente es simétrica y centrada en la ganadora  $\mathbf{g}$  por lo que la distancia de una neurona  $\mathbf{i}=(i_i, i_j)$  a  $\mathbf{g}=(g_i, g_j)$  generalmente se mide con la siguiente ecuación

$$|\mathbf{i} - \mathbf{g}| = \sqrt{(i_i - g_i)^2 + (i_j - g_j)^2}$$

- La función  $\mathbf{h}()$  puede tomar varias formas, que dependen de una función  $R(t)$  que determina el **radio** de “vecindad”

- Una función simple es la escalón

$$h(|\mathbf{i} - \mathbf{g}|, t) = \begin{cases} 0, & \text{si } |\mathbf{i} - \mathbf{g}| > R(t) \\ 1, & \text{si } |\mathbf{i} - \mathbf{g}| \leq R(t) \end{cases}$$

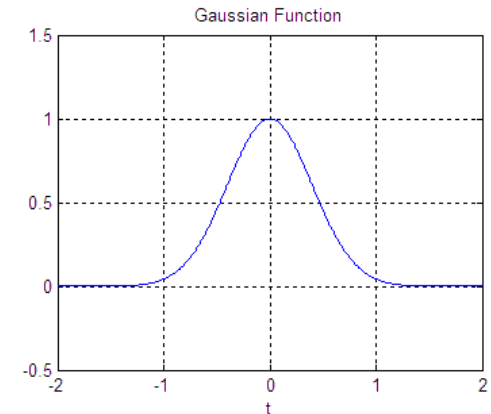


- SOFM (Self-Organizing Feature Maps)

- **Procedimiento general de aprendizaje**

- Función de vecindad  $h(|\mathbf{i} - \mathbf{g}|, t)$

- Otras funciones pueden tomar otras formas, como ser formas gaussianas



- También pueden adoptar formas cóncavas o convexas
      - La función de vecindad en sí misma no cambia con el tiempo, pero sí la función que define el radio  $R(t)$ , la cual generalmente disminuye monótonamente
        - Ej: disminución lineal

$$R(t) = R_0 + (R_f - R_0) \frac{t}{t_R}$$

donde  $t$  es la iteración, y  $t_R$  el número de iteraciones para alcanzar  $R_f$

- SOFM (Self-Organizing Feature Maps)
  - **Consideraciones prácticas**
    - Los SOFM tal como se han descripto anteriormente, pueden utilizarse para análisis de datos (exploración de datos, monitorización de procesos, etc) y mapping (cuantificación vectorial, reducción de dimensiones, etc)
    - Puede también adaptarse la técnica para resolver problemas de **clasificación**. Ej: utilizando el ajuste fino LVQ (*Learning Vector Quantization*)
    - Existen técnicas complementarias enfocadas en
      - Tamaño adaptativo de los mapas
      - Solución de la asimetría (“desventaja”) de neuronas en los bordes del mapa
      - Procesamiento de series de tiempo
      - Aplicación en redes híbridas: multi-capas con partes supervisadas + partes no supervisadas (Ej: red de contrapropagación, RBF, etc)



- SOFM (Self-Organizing Feature Maps)

- Métricas de similitud
  - Correlación o producto escalar

$$C_{ij} = \sum_{k=1}^n w_{ijk} x_k$$

- Conicida con el modelo de neurona estándar (suma ponderada)
- Desventaja: sensible al tamaño de los vectores (en general requiere vectores normalizados para funcionar bien)

- Coseno

$$\cos(\mathbf{w}_{ij}, \mathbf{x}) = \frac{\sum_{k=1}^n w_{ijk} x_k}{\|\mathbf{w}_{ij}\| \cdot \|\mathbf{x}\|}$$

- Se deriva de la correlación, y su ventaja es que es independiente del tamaño de los vectores

- SOFM (Self-Organizing Feature Maps)

- Métricas de similitud
  - Distancia euclídea

$$d^2(\mathbf{w}_{ij}, \mathbf{x}) = \sum_{k=1}^n (w_{ijk} - x_k)^2$$

- Distancia de Manhattan

$$d(\mathbf{w}_{ij}, \mathbf{x}) = \sum_{k=1}^n |w_{ijk} - x_k|$$

- Ventaja: muy simple computacionalmente, lo que la hace muy rápida (no requiere potenciación, radicación ni multiplicaciones)

- SOFM (Self-Organizing Feature Maps)

- Reglas de aprendizaje para las métricas de similitud más usadas
  - Aprendizaje para la distancia euclídea

$$\mathbf{w}_i(t+1) = \mathbf{w}_i(t) + \alpha(t) h(|\mathbf{i} - \mathbf{g}|, t) (x_k(t) - w_{ijk}(t))$$

- Aprendizaje para la distancia Manhattan

$$\mathbf{w}_i(t+1) = \mathbf{w}_i(t) + \alpha(t) h(|\mathbf{i} - \mathbf{g}|, t) \text{signo}(x_k(t) - w_{ijk}(t))$$

o bien

$$\mathbf{w}_i(t+1) = \mathbf{w}_i(t) + \Delta w_{ijk}(t)$$

$$\text{con } \Delta w_{ijk}(t) = \begin{cases} +\alpha(t) h(|\mathbf{i} - \mathbf{g}|, t) & \text{si } x_k(t) > w_{ijk}(t) \\ 0 & \text{si } x_k(t) = w_{ijk}(t) \\ -\alpha(t) h(|\mathbf{i} - \mathbf{g}|, t) & \text{si } x_k(t) < w_{ijk}(t) \end{cases}$$

- SOFM (Self-Organizing Feature Maps)
  - Reglas de aprendizaje
    - La regla de aprendizaje debe ser compatible con la métrica de similitud para que el mapa se desarrolle correctamente
    - Procedimiento sistemático para derivar reglas de aprendizaje
      - Sea  $d(\mathbf{w}_{ij}, \mathbf{x})$  una función de distancia genérica, derivable
      - La neurona ganadora  $\mathbf{g}$  cumple

$$\mathbf{g} = \arg \min_{ij} \{d(\mathbf{w}_{ij}, \mathbf{x})\}$$

- Se define una función de error que pretende representar la diferencia entre los pesos sinápticos y las entradas del espacio sensorial

- SOFM (Self-Organizing Feature Maps)

- Reglas de aprendizaje
  - Procedimiento sistemático para derivar reglas de aprendizaje (cont.)

- Una **aproximación** a tal función **global** de error es

$$E'(t) = \sum_i h(|\mathbf{i} - \mathbf{g}|, t) f(d(\mathbf{x}(t), \mathbf{w}_i(t)))$$

con **f()** una función, usualmente identidad, incluida por generalidad

- Aplicando el mecanismo de descenso por el gradiente, se pueden derivar actualizaciones de peso dependientes de f(), d() y h()

$$\mathbf{w}_i(t+1) = \mathbf{w}_i(t) - \lambda(t) \nabla_{\mathbf{w}_i} E'(t)$$

con  $\lambda(t)$  ritmo de aprendizaje

- SOFM (Self-Organizing Feature Maps)
  - **LVQ (*Learning Vector Quantization*)**
    - Permite ajuste fino de un mapa para aplicaciones de clasificación de patrones
    - Cada neurona o grupo de neuronas representa una **clase**, y cada una almacena un vector de pesos  $w_{ij}$  que representa el **prototipo** de una clase
    - Aprendizaje: para un patrón  $\mu$  dado
      - Si la neurona ganadora representa la clase del patrón, se la **premia** ajustando sus pesos **hacia** el patrón (similar al mapa de Kohonen tradicional)
      - Si la neurona ganadora no representa la clase del patrón, se la **castiga** ajustando sus pesos en el **sentido contrario** al del patrón

- SOFM (Self-Organizing Feature Maps)

- **LVQ (*Learning Vector Quantization*)**

- Aprendizaje: ajuste de pesos

Suponiendo un patrón de entrada  $\mathbf{x}^\mu$  que pertenece a una clase  $C_{x^\mu}$  y una neurona ganadora  $\mathbf{w}_c$  que representa una clase  $C_{w_c}$

- Si  $C_{w_c} = C_{x^\mu}$  (clasificación correcta) los pesos de la neurona ganadora se actualizan con

$$\mathbf{w}_c(t+1) = \mathbf{w}_c(t) + \alpha_t [\mathbf{x}^\mu - \mathbf{w}_c(t)]$$

- Si  $C_{w_c} \neq C_{x^\mu}$  (clasificación incorrecta) los pesos de la neurona ganadora se actualizan con

$$\mathbf{w}_c(t+1) = \mathbf{w}_c(t) - \alpha_t [\mathbf{x}^\mu - \mathbf{w}_c(t)]$$

- El resto de los pesos no se modifican
      - Para evitar overfitting se puede usar parada temprana
      - Una cantidad de iteraciones de entre 50 y 200 veces la cantidad de neuronas suele ser suficiente para el entrenamiento
      - Inicialización del mapa: (a) un ejemplo de la clase; (b) k-means; (c) entrenamiento SOFM tradicional (sin ajuste fino)