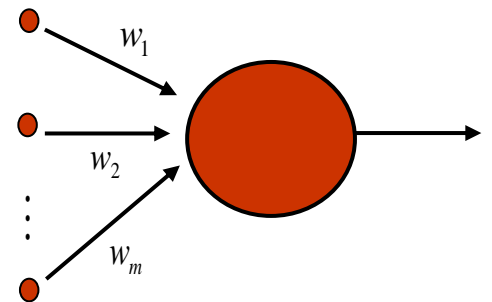




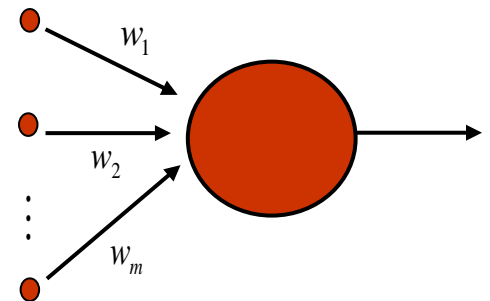
# ***ALGUNOS MODELOS DE UNA NEURONA***

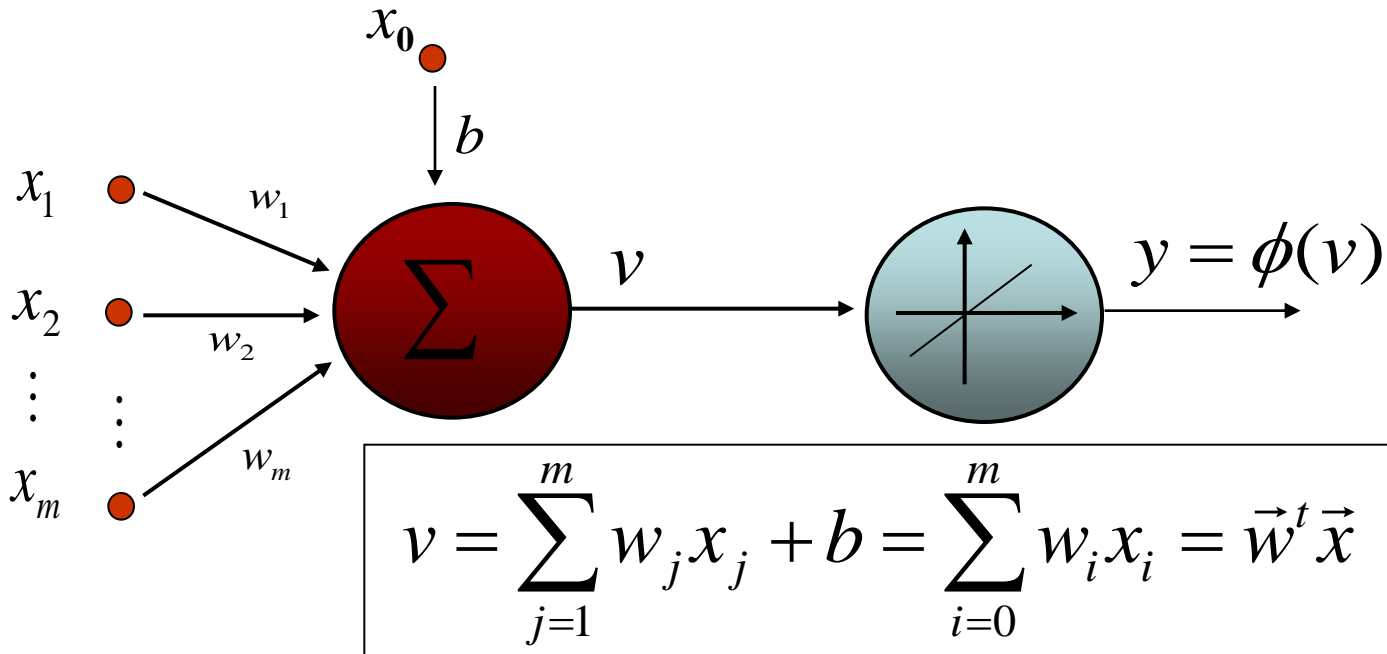




# ADALINE (ADAptive Linear Element)

Widrow y Hoff (1960)





con

$$\vec{w} = [b \quad w_1 \quad w_2 \quad \dots \quad w_m]^t$$
$$\vec{x} = [1 \quad x_1 \quad x_2 \quad \dots \quad x_m]^t$$
$$y = v$$



## Notación:

- $D$  es el conjunto de datos  $\{(x(k), d(k)), k = 1, \dots, L\}$ , con  $x(k)$  un vector

Queremos minimizar el error cuadrático medio

$$E(\vec{w}) = \frac{1}{2} \sum_{x(k) \in D} (d(k) - y(k))^2$$

**Esta función por definición es concava, diferenciable y tiene un mínimo.**



De la estadística sabemos que la solución de mínimos cuadrado es:

$$W = \underbrace{(X^t X)^{-1}}_{\text{Pseudo-inversa}} X^t d$$

**Pseudo-inversa**

con

$$X = \begin{pmatrix} x_1 & x_1 & \cdots & x_1 \\ x_2 & x_2 & \cdots & x_2 \\ \vdots & \vdots & \ddots & \vdots \\ x_N & x_N & \cdots & x_N \end{pmatrix}$$

**Esta forma tiene un costo computacional muy alto y puede ser prohibitiva cuando tenemos tamaños de datos muy grandes.**



Podemos tener como en el caso del perceptrón varios ADALINEs juntos y en este caso se le denomina MADALINE. Vamos a centrarnos en una sola neurona.

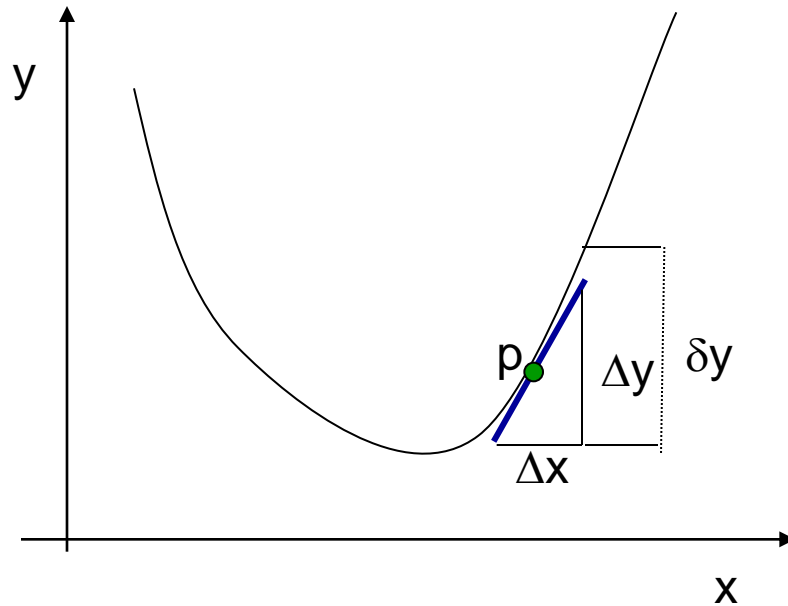
### El Entrenamiento

La regla de actualización de pesos es la **LMS (least mean square algorithm)** que se basa en el método de descenso de gradiente para minimizar la diferencia entre la salida deseada y la actual (entrenamiento supervisado).

El método de descenso de gradiente forma parte del conjunto de métodos de optimización sin restricciones.

En el caso del **perceptrón** se requería que el problema fuera **linealmente separable**. Aquí buscamos un clasificador lineal que sea óptimo en algún sentido (minimizando una función de costos).

La idea es producir una secuencia de pesos  $W(n)$  tal que en cada iteración la función de costos disminuya, esto se logra con **correcciones al vector en dirección opuesta al gradiente** (que resulta ser la de mayor descenso)



La pendiente de la función es el gradiente de la tangente a la curva en ese punto

$$P = \Delta y / \Delta x$$

en 1-D: Si inicio en  $a$ , para  $\gamma$  suficientemente pequeño,  $b = a + \gamma y'(a)$  es tal que

$$y(b) = y(a + \underbrace{\gamma y'(a)}_h) = y(a) + h y'(a) + O(h^2) = y(a) + \gamma (y'(a))^2 > y(a)$$

Hay que caminar en dirección opuesta!!

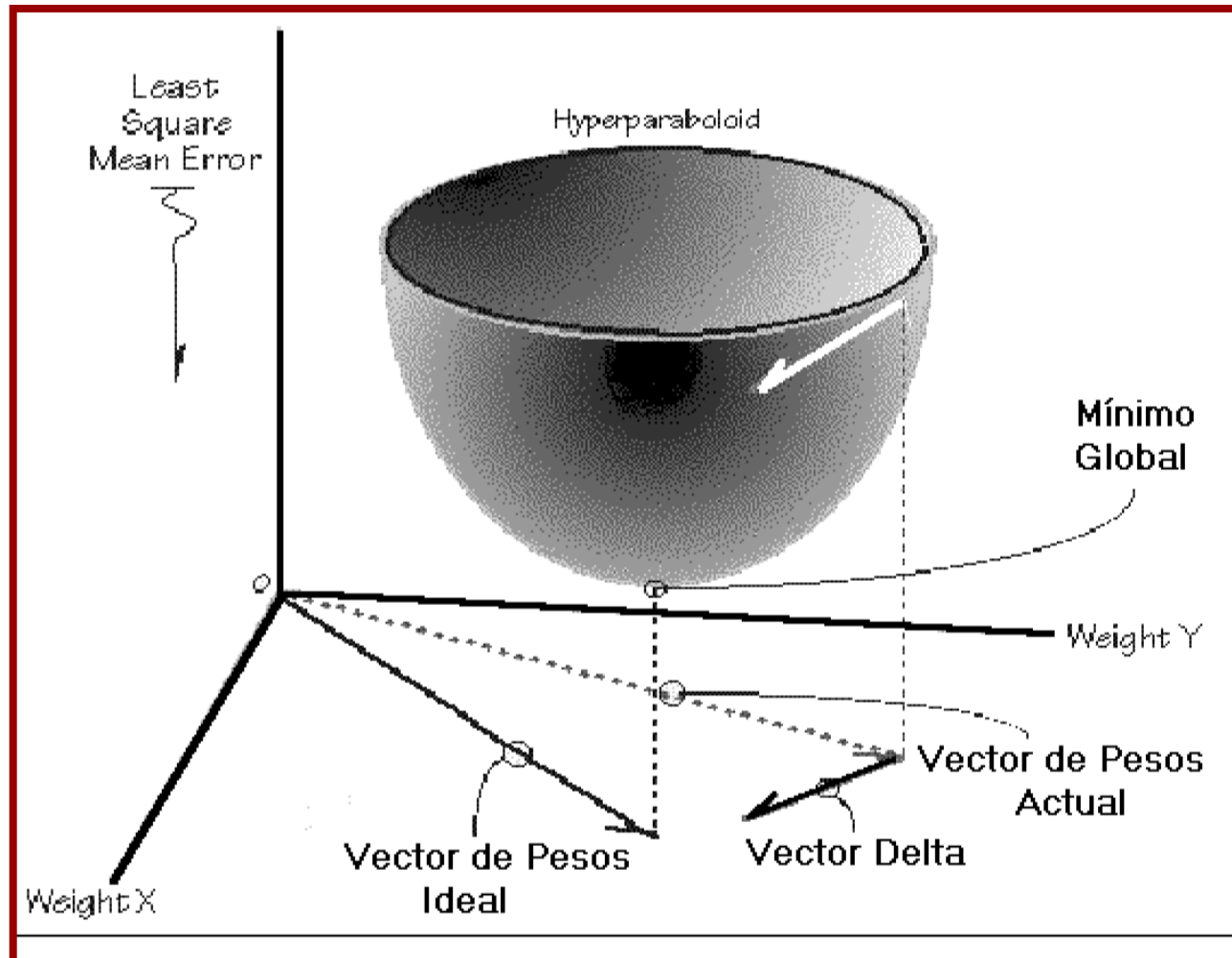


**Dicho de otra manera: En 1D para encontrar un mínimo local ,**

- Si  $f'(x) > 0$ , me muevo hacia la izquierda,**
- Si  $f'(x) < 0$ , me muevo hacia la derecha**
- Si  $f'(x) = 0$ , me quedo quieta**

**En espacios de dimensión  $>1$ , nos movemos en la dirección opuesta al gradiente.**







El gradiente se calcula como:

$$\nabla E(\vec{w}) = \left[ \frac{\partial E}{\partial w_0}, \frac{\partial E}{\partial w_1}, \dots, \frac{\partial E}{\partial w_m} \right]$$

En el algoritmo del descenso del gradiente tomamos un paso en dirección de máximo descenso que es menos el gradiente

$$\Delta \vec{w} = -\eta \nabla E(\vec{w})$$

En cada componente

$$\Delta w_i = -\eta \frac{\partial E}{\partial w_i}$$



$$\begin{aligned}\frac{\partial E}{\partial w_i} &= \frac{\partial}{\partial w_i} \left[ \frac{1}{2} \sum_{x(k) \in D} (d(k) - y(k))^2 \right] \\ &= \frac{1}{2} \sum_{x(k) \in D} \frac{\partial}{\partial w_i} (d(k) - y(k))^2 \\ &= \frac{1}{2} \sum_{x(k) \in D} 2(d(k) - y(k)) \frac{\partial}{\partial w_i} (d(k) - y(k)) \\ &= \sum_{x(k) \in D} (d(k) - y(k)) \frac{\partial}{\partial w_i} (d(k) - \vec{w}x(k)) \\ &= \sum_{x(k) \in D} -(d(k) - y(k))x_i(k)\end{aligned}$$



Por lo tanto la regla de actualización es

$$\Delta w_i = \eta \sum_{x(k) \in D} (d(k) - y(k)) x_i(k)$$

El método del gradiente garantiza convergencia hacia 0 cuando los pesos sinápticos aproximan la relación de entrada y salida. La convergencia es asintótica, mientras que en el caso del perceptrón sabemos que converge en un número finito de pasos.

**OJO!!!!**

Hay dos tipos de actualizaciones distintas que inducen a tipos de convergencia diferentes: **Actualización en lote** y **actualización continua**



### Actualización por lotes (Batch update)

Repetir hasta satisfacer la condición de parada

- Buscar el gradiente del error  $\nabla E_D(\vec{w})$
- $\vec{w} \leftarrow \vec{w} - \eta \nabla E_D(\vec{w})$

### Actualización continua (on-line update)

Repetir hasta satisfacer la condición de parada

- Para cada par  $(x,d)$  en  $D$ 
  - \*\* Calcule  $\nabla E_d(\vec{w}) = (d(k) - y(k))x(k)$
  - \*\*  $\vec{w} \leftarrow \vec{w} - \eta \nabla E_d(\vec{w})$



En actualización continua, usamos como función de error el **error cometido en la n-ésima iteración**, mientras que en el descenso de gradiente usamos la **suma de los errores**.

### Algoritmo

- 1) Iniciar los pesos en valores aleatorios pequeños
- 2) Poner  $E=0$
- 3) Para cada  $w_i$ , inicializar  $\Delta w_i=0$
- 4) Para cada  $(x,d)$  en  $D$ 
  - $\delta = d - y(x)$
  - $E = E + \delta^2$
  - $\Delta w_i = \Delta w_i + \eta \delta x_i$ , para  $i=1, \dots, m$
- 5) Para cada  $w_i$  actualizar  $w_i = w_i + \Delta w_i$
- 6) Si no se verifica la condición de parada, continuar al paso 2)



**Este algoritmo actualiza los pesos al final de cada época en vez de continuamente con cada presentación de los estímulos, la convergencia está garantizada por la convergencia de descenso de gradiente.**

**Sin embargo el LMS que se presenta en el Haykin es la actualización continua**

**1) Iniciar los pesos en valores aleatorios pequeños**

**2) Para cada  $(x,d)$  en  $D$**

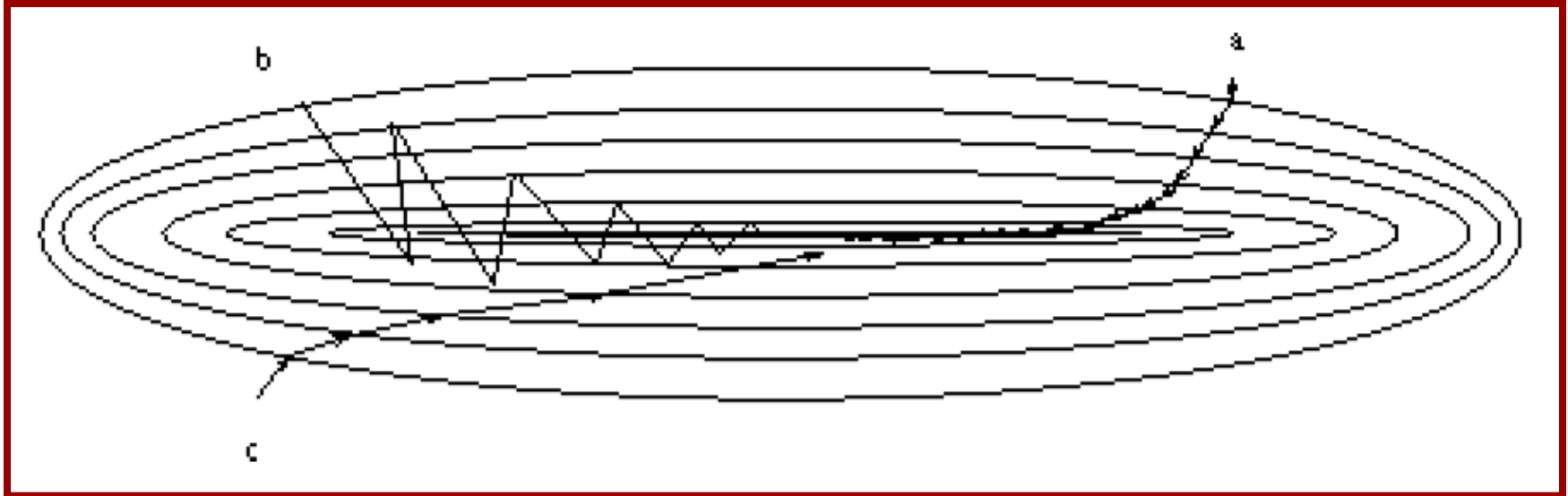
**-  $\delta = d - y(x)$**

**-  $w_i = w_i + \eta \delta x_i$ , para  $i=1 \dots m$**



propiedades	<b>On-line</b>	<b>batch</b>
convergencia	Los pesos siguen una trayectoria “aleatoria” (stochastic gradient descent). A medida que aumenta el número de iteraciones los pesos se comportan como un movimiento Browniano.	Converge al mínimo lentamente (depende de $\eta$ ). Si $\eta$ es suf. pequeño los pesos siguen una trayectoria suave.
aproximaciones	Aproxima la actualización por lotes (batch) cuando $\eta$ es pequeño	
evaluaciones		Mas eficiente





$\alpha$ )  $\eta$  pequeña

$\beta$ )  $\eta$  alta



**Se puede demostrar que LMS (continua) converge si  $0 < \eta < (2/\lambda_{\max})$  donde  $\lambda_{\max}$  es el mayor autovalor de la matriz de correlación para las entradas.**

**Un estimado esta dado por  $0 < \eta < (2/\text{tr}(R_x))$ , donde  $\text{tr}(R_x)$  es la suma cuadrada de los estímulos.**

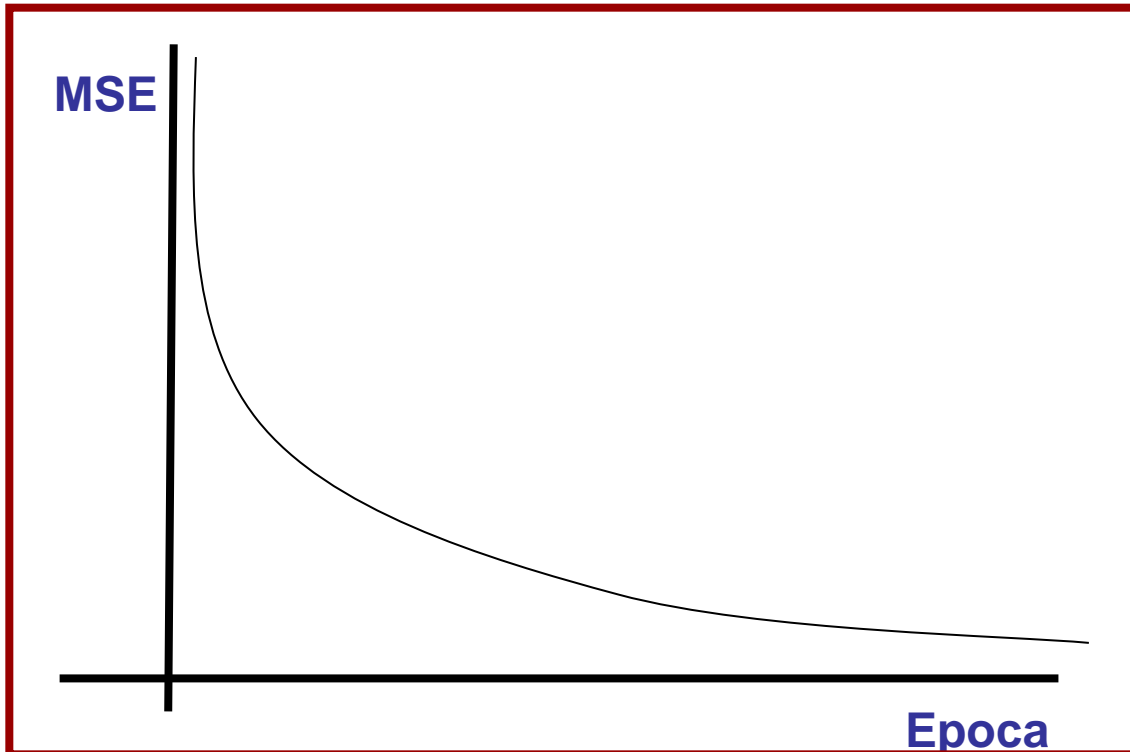
$$R_x = E[x^*x^t]$$

**A la actualización continua se le puede mejorar su convergencia si se utilizan parámetros adaptativos**

- $\eta = \eta_0$
- $\eta(n) = c/n$  (aproximación estocástica)
- $\eta(n) = \eta_0 / (1 + n/\tau)$



Si hemos elegido bien al parámetro de aprendizaje, en cada iteración vamos a acercarnos más hacia un mínimo (en nuestro caso global), por lo que el error en cada época debería ser menor que la anterior.



### FAQ:

Cómo es la curva de aprendizaje del perceptrón?

Porqué no es monótona si la regla de actualización es la misma?



**1) Clasificación.** Recordar que un ADALINE busca un hiperplano. En esta aplicación es el que separa. El error cometido no será 0, por lo que típicamente la condición de parada debe cambiar.

- 1) Diferencia entre errores en épocas sucesivas sea pequeña
- 2) Número máximo de épocas.

**2) Regresión lineal.** Se busca un modelo de la forma:

$$Y = a + b_1x_1 + b_2x_2 + b_3x_3 + \dots + b_nx_n$$

**3) Modelos autorregresivos.** Se busca un modelo de la forma:

$$Y = a + b_1x(n) + b_2x(n-1) + b_3x(n-2) + \dots + b_kx(n-k+1)$$

**4) Aproximación polinómica.** Se busca un modelo de la forma:

$$Y = a + b_1x + b_2x^2 + b_3x^3 + \dots + b_kx^k$$



1. **Definición de red neural, vinculación con cerebro humano.**
2. **Paradigmas del aprendizaje.**
3. **Definimos una red neural artificial .**
4. **Modelos de una neurona:**
  1. **Perceptrón**
    1. **Neurona no-lineal.**
    2. **Clasificación, data linealmente separables.**
    3. **Usa LMS (actualización continua) como regla de aprendizaje.**
    4. **Convergencia garantizada para cualquier valor de la tasa de aprendizaje.**
  2. **Adaline**
    1. **Dispositivo lineal.**
    2. **Usa LMS como regla de aprendizaje.**
    3. **Convergencia garantizada según estilo de iteración, sensible a cambios en la tasa de aprendizaje.**
5. **Vimos ejemplos de declaración de cada dispositivo en Matlab.**





**Hay problemas difíciles de resolver con una sola neurona, o con una capa de neuronas y hace falta recurrir a arquitecturas más complejas para resolverlas.**

**Vamos a ver redes artificiales con más de una capa de neuronas, y la forma de aprendizaje:**

**retropropagación (Backpropagation)**