

# APRENDIZAJE AUTOMÁTICO

Hasta aquí, hemos considerado el problema de memoria asociativa a través de redes recurrentes, en el cual hemos construido de inicio reglas, como la de Hopfield o sus modificaciones. Estas reglas "NO SE APRENDIERON", se "CONSTRUYERON".

A partir de hoy analizaremos un nuevo paradigma, y usaremos una nueva arquitectura.

*Nota: el modelo de Hopfield también se puede "aprender", aunque no vemos aquí este algoritmo.*

Como niños, no es trivial encontrar una regla para un problema tan complejo como el problema de memoria asociativa. La idea en el caso de "aprendizaje" es encontrar un algoritmo "adecuado" para resolver un problema dado.

El modelo de Hopfield es un caso de sistema "auto-asociativo", que identifique muchas configuraciones con una única configuración. Aquí el input y el output tiene la misma forma y el mismo tamaño. Pero hay muchísimos problemas en los cuales el input es de un tipo y el output de otro.

En general, más allá de las redes neuronales pero también incluyendo las, podemos decir que hay 3 tipos de algoritmos de aprendizaje. En el caso de redes, se trata de encontrar los acoplamientos  $w_{ij}$  en forma automática

## APRENDIZAJE SUPERVISADO

En este caso disponemos de datos de entrada que han sido ETIQUETADOS previamente, y nos permiten entonces encontrar el conjunto de las sinapsis que resuelven bien el problema para los ejemplos etiquetados, a los cuales llamamos "CONJUNTO DE ENTRENAMIENTO"

Ejemplo: clasificar como electrónico entrante como spam o no spam.

Algunos algoritmos desarrollados para aprendizaje supervisados son:

- K-Nearest Neighbors
- Regression lineal
- Regression logística
- Support Vector Machine
- Clasificados bayesianos
- Árboles de decisión
- Random forest
- Redes Neuronales
- Redes Neuronales profundas.

## APRENDIZAJE NO SUPERVISADO

En este caso no disponemos de datos previamente etiquetados y el algoritmo debe clasificar o procesar la información por sí solo.

Este algoritmo se usa generalmente para clasificar datos en pocas categorías.

Algunos algoritmos desarrollados para aprendizaje no supervisado son

- K-means
- Análisis de la componente principal
- Detección de anomalías.
- Redes neuronales

## APRENDIZAJE POR REFUERZO

El algoritmo es un "agente autónomo" que explora un "espacio" desconocido y decide que acciones realizar mediante prueba y error.

Recibirá premios y castigos. Cree la mejor estrategia posible para obtener la mayor recompensa. Esta estrategia definirá acciones a tomar ante cada situación que se enfrente.

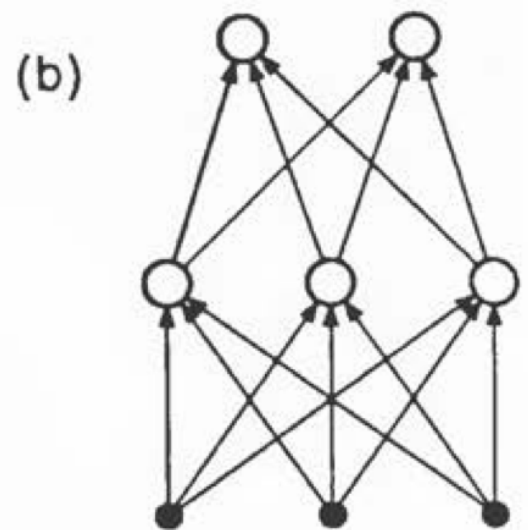
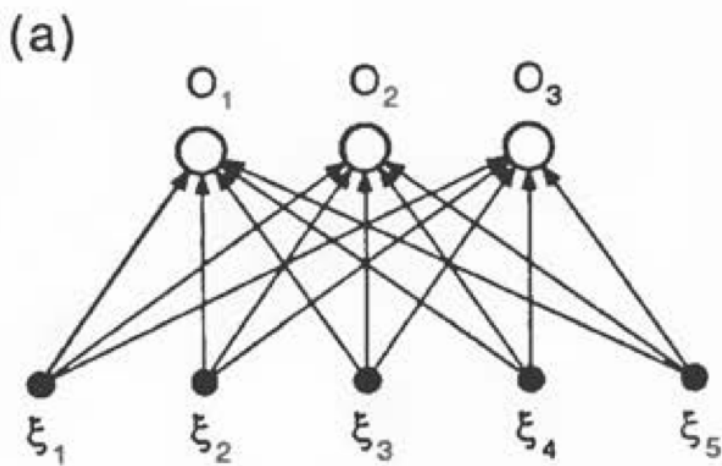
- MDP: Markov Decision Process



# REDES FEED-FORWARD

## PERCEPTRON

1958: Frank Rosenblatt (Cornell Aeronautical Laboratory)



los neuronas se organizan en capas o comedas. la informacion fluye de la entrada a la salida.

la salida es una capa

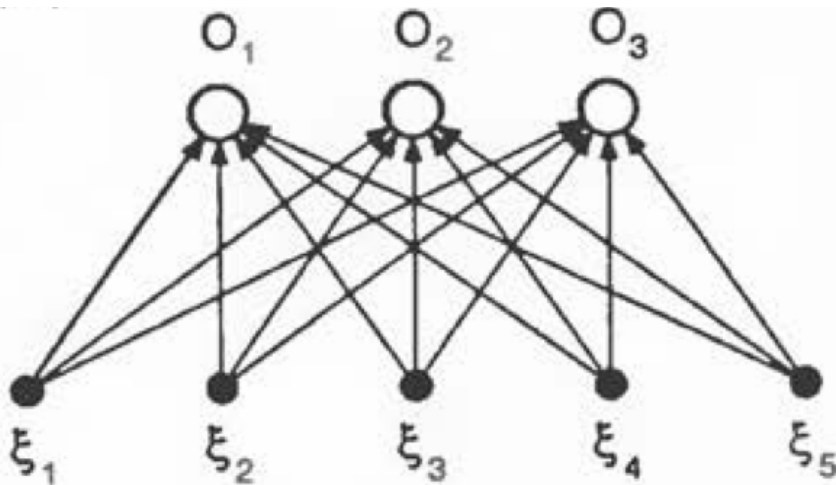
la entrada **NO ES** una capa.

En el ejemplo (a) hay 1 capa.

En el ejemplo (5) hay 2 capas.

- No se admiten conexiones entre neuronas de la misma capa
- No se admiten conexiones de una capa a otra previa
- Solo se admiten conexiones entre una capa y la siguiente.
- las conexiones son asimétricas.

## EL PERCEPTOR DE 1 CAPA



El vector

$$\vec{\xi} = (\xi_1, \xi_2, \xi_3, \dots, \xi_N)$$

representa la entrada. las variables  $\xi_i$  pueden ser discretos o continuos.

El vector

$$\vec{O} = (O_1, O_2, \dots, O_M)$$

representa la salida. Tenemos

$N \times M$  acoplamiento

$$W_{ij} \quad i=1, 2, \dots, M \quad j=1, 2, \dots, N$$

Tenemos que encontrar un conjunto  $\{W_{ij}\}$  tal que ante una entrada dada

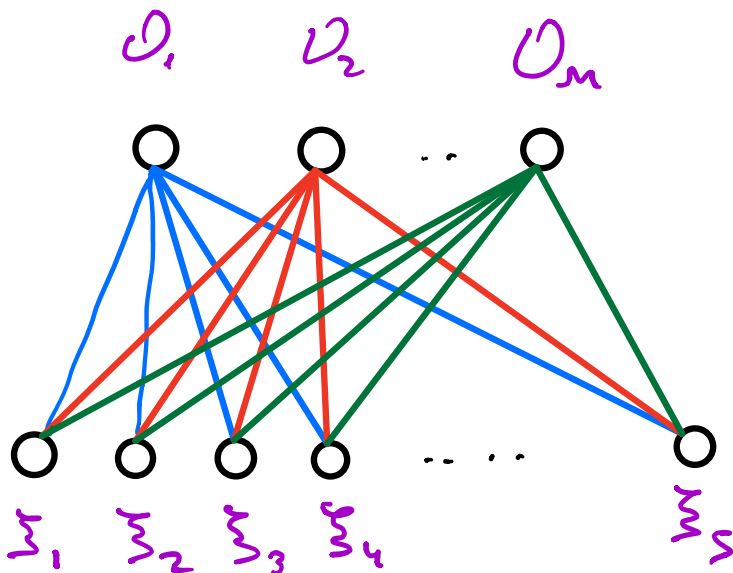
$$\vec{\xi}^M = (\xi_1^M, \xi_2^M, \dots, \xi_N^M)$$

no lleve a una salida determinada

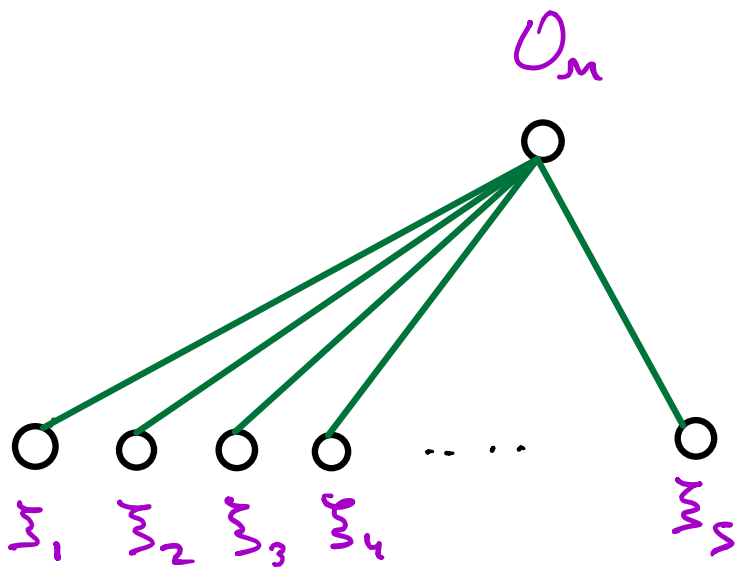
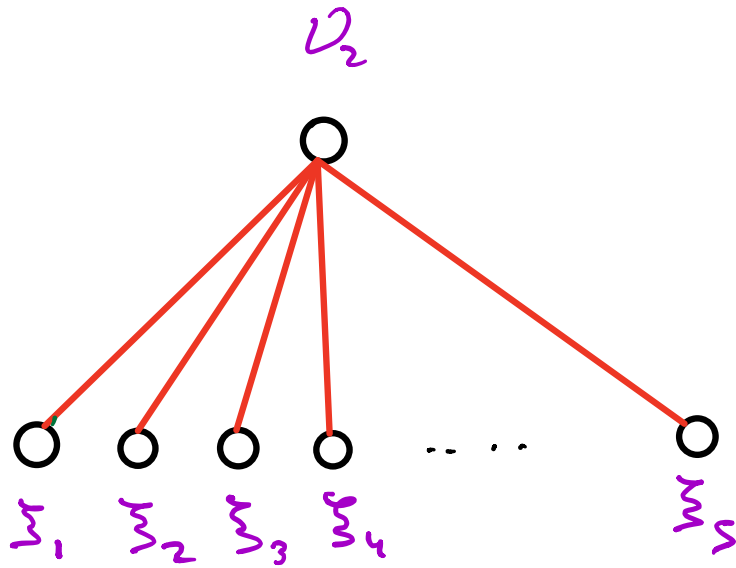
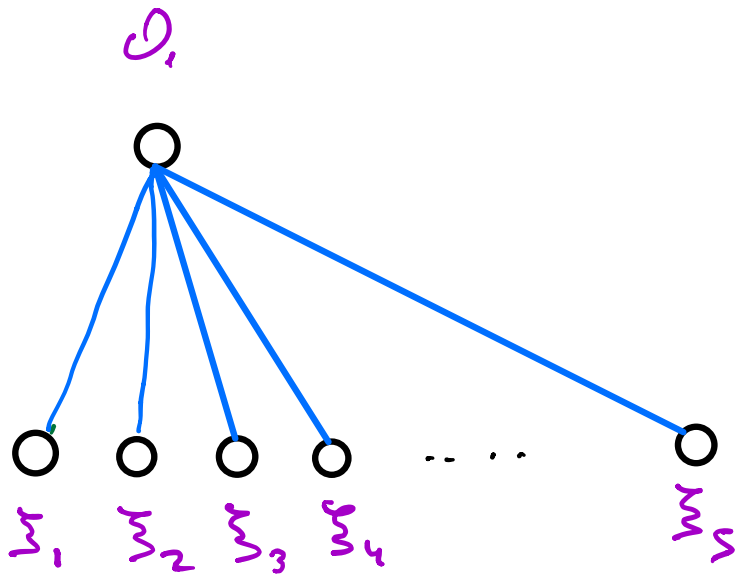
$$\vec{O}^M = (O_1^M, O_2^M, \dots, O_n^M)$$

## EL PERCEPTRON SIMPLE

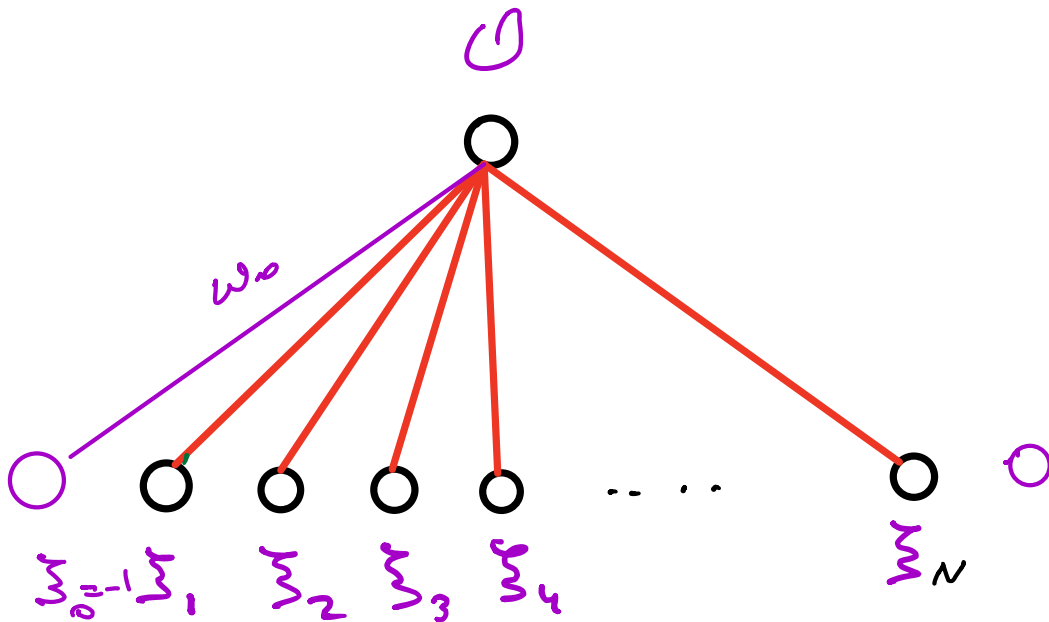
Podemos referir un perceptron de 1 capa en  $M$  perceptrones con 1 salida cada uno



Noten que no se mezclan los  $M$  perceptrones simples.



Podemos entonces analizar un perceptron simple por vez



Suprimiremos el indice del output.

$$O = g(h) = g\left(\sum_{k=1}^N w_k \sum_k\right)$$

donde  $g$  es la función de activación y puede ser cualquiera de las que ya usamos, o alguna nueva.

Hemos omitido los umbrales. ¿Porque?

$$0 = g\left(\sum_{k=1}^N w_k \xi_k - \theta\right)$$
$$= g\left(\sum_{k=0}^N w_k \xi_k\right)$$

si  $w_0 = \theta$  y  $\xi_0 = -1$  fijo.

O sea, el umbral se puede pensar como una variable de entrada fija en  $\xi_0 = -1$  con  $w_0 = \theta$ . De esta forma trata a todos de igual manera.

Supongamos que tenemos un conjunto de entradas etiquetadas; tal que si entre  $\sum_{k=1}^N$  sabemos que debe

don

$$O = \sum^M$$

$$\sum^M \rightarrow \sum^M$$

$$O^M = g(h^M) = g\left(\sum_{k=0}^N w_k \sum_k^M\right)$$

Todo dependerá sobre del Tipo de funciones de salida, o sea, de cual es la función de activación  $g(h)$

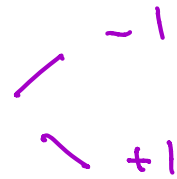
hetero-associación:

le salida es de naturaleza diferente al input



# EL CASO BINARIO

$$g(h) = \text{signo}(h)$$



Suponemos que  $\sum_i^N$  también es binario

$$\sum_i^N = \pm 1$$

Suponemos que de los  $2^N$  posibles inputs conseguimos  $p$  resultados

$$\vec{\sum}^N \longrightarrow \mathcal{J}^N = \pm 1$$

$$(w_1, w_2, \dots, w_N) = \vec{w}$$

$$\sum_k w_k \sum_k^N = \vec{w} \cdot \vec{\sum}^N$$

$$\mathcal{O}^N = \text{signo}(h^N) = \mathcal{J}^N$$

$$\text{signo}(h^M) = \text{signo}\left(\sum_k w_k \xi_k^M\right) = y^M$$

$$\text{signo}(\bar{w} \cdot \bar{\xi}^M) = y^M$$

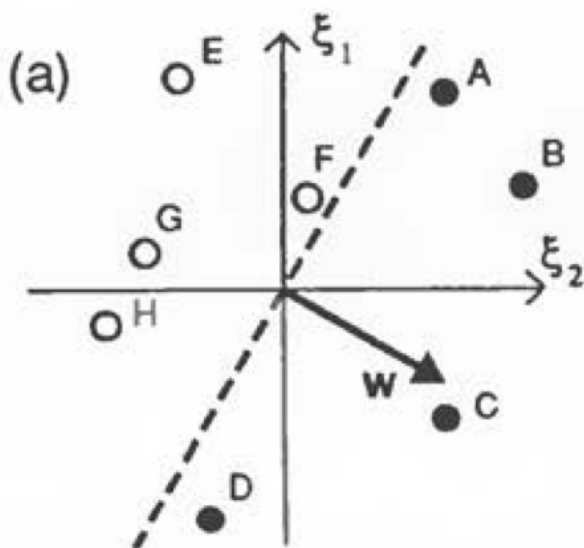
lo deseado

Esto debe valer para los  $p$  valores de  $y$ .

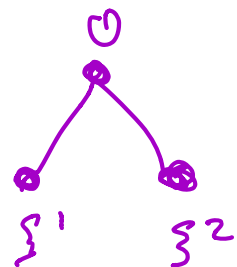
$\bar{w} \cdot \bar{\xi}^M$  es la proyección de  $\bar{\xi}^M$  sobre  $\bar{w}$

$$\bar{w} \cdot \bar{\xi}^M = |\bar{w}| |\bar{\xi}^M| \cos(\alpha)$$

donde  $\alpha$  es el ángulo entre  $\bar{w}$  y  $\bar{\xi}^M$ .



Debe existir un plano (en este caso una recta) que separe los  $\bullet$  de los  $\circ$ .



• significa output deseado  $\zeta^M = +1$

○ significa output no deseado  $\zeta^M = -1$

Outra forma:

$$\text{signo}(\bar{w} \cdot \bar{\zeta}^M) = \zeta^M$$

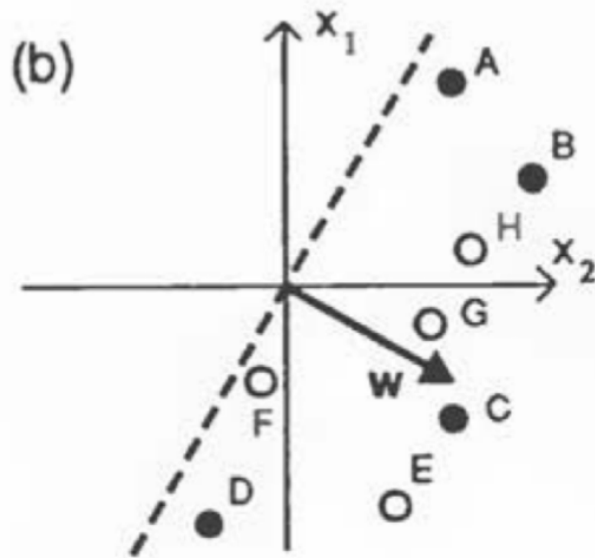
$$\zeta^M \cdot \text{signo}(\bar{w} \cdot \bar{\zeta}^M) = \zeta^M \zeta^M = +1$$

$$\text{signo}(\zeta^M (\bar{w} \cdot \bar{\zeta}^M)) = +1$$

$$\text{signo}(\bar{w} \cdot (\zeta^M \bar{\zeta}^M)) = +1$$

$$\text{signo}(\bar{w} \cdot \bar{x}^M) = +1$$

$$\bar{w} \cdot \bar{x}^M > 0 \quad (\text{deseado})$$



¿Qué pasa si no se pueden separar?  
 En este caso no se puede resolver.

Un problema de clasificación binaria es linealmente separable si existe un plano en el espacio  $\sum$  que separe los casos  $y^{\mu} = +1$  de los casos  $y^{\mu} = -1$ .

# SEPARABILIDAD LINEAL

Miremos que pasa si incluimos explícitamente el umbral

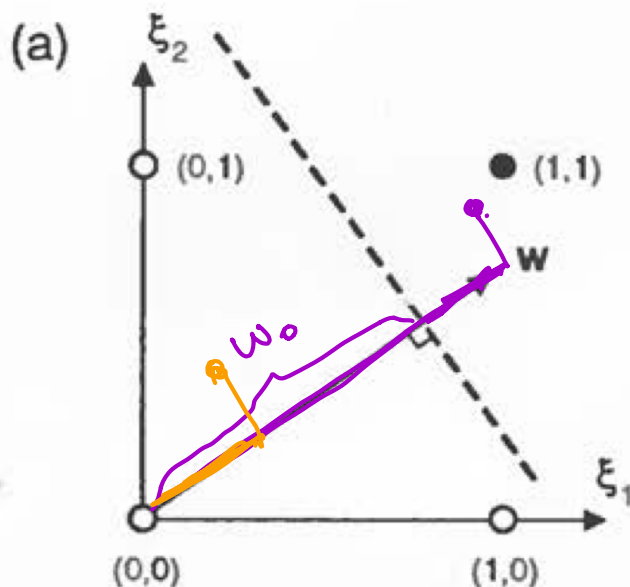
$$O = \text{sign}(\bar{w} \cdot \bar{\xi} - w_0)$$

Las regiones con  $Y^H = +1$  y  $Y^H = -1$  deben separarse por un plano de dimensión  $N-1$ .

El plano se separa por

$$\bar{w} \cdot \bar{\xi} = w_0$$

$$\text{Signo}(\bar{w} \cdot \bar{\xi} - w_0)$$



# EJEMPLOS

LA FUNCIÓN "Y" O "AND".

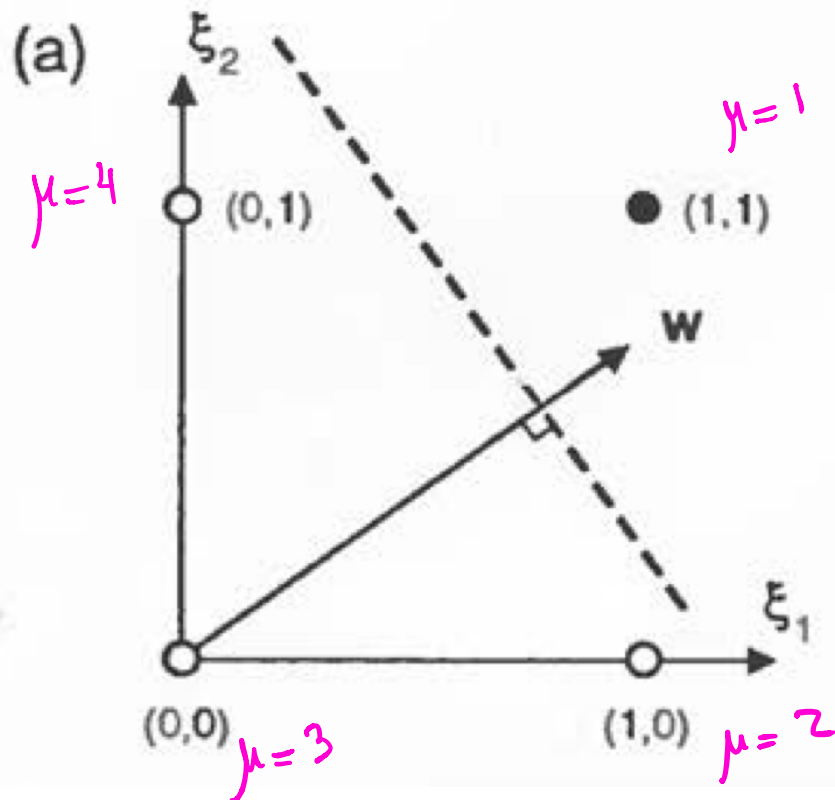
$\mu$	$\xi_1$	$\xi_2$	$y$
-------	---------	---------	-----

1	1	1	+1
---	---	---	----

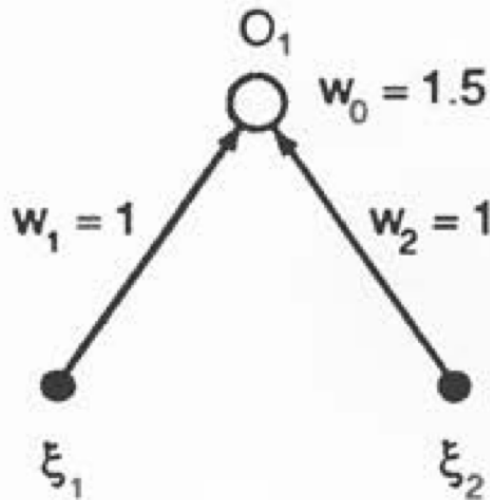
2	1	0	-1
---	---	---	----

3	0	1	-1
---	---	---	----

4	0	0	-1
---	---	---	----



(b)



$$\mu = 1 \quad \xi_1 = 1 \quad \xi_2 = 1 \quad \zeta^1 = +1$$

$$\begin{aligned} \zeta^1 &= \text{signo} (w_1 \xi_1 + w_2 \xi_2 - 1.5) \\ &= \text{signo} (1 \cdot 1 + 1 \cdot 1 - 1.5) = \text{signo} (2 - 1.5) \\ &= \text{signo} (0.5) = +1 = \zeta^1 \end{aligned}$$

$$\mu = 2 \quad \xi_1^2 = 1 \quad \xi_2^2 = 0 \quad \zeta^2 = -1$$

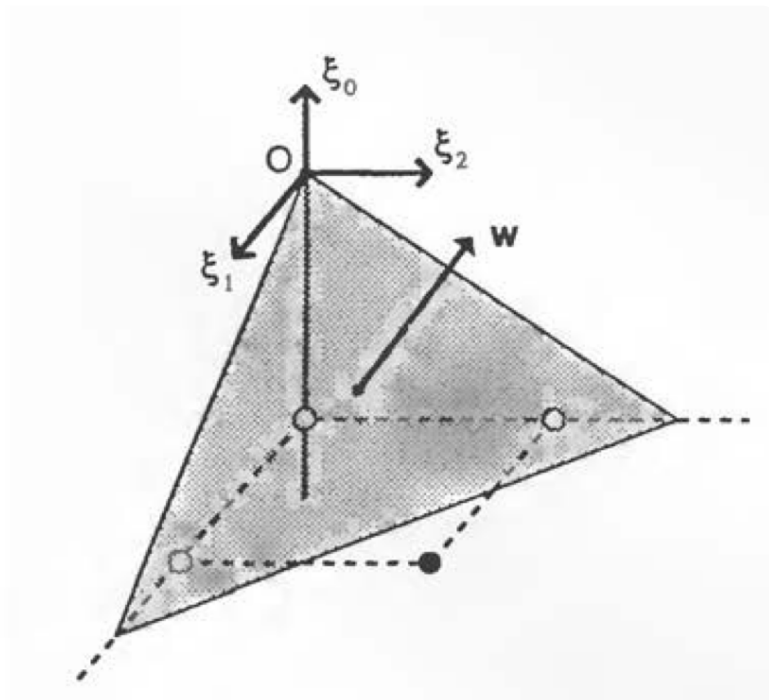
$$\begin{aligned} \zeta^2 &= \text{signo} (1 \cdot 1 + 1 \cdot 0 - 1.5) = \text{signo} (1 - 1.5) \\ &= \text{signo} (-0.5) = -1 = \zeta^2 \end{aligned}$$

$$\mu = 3 \quad \sum_{i=1}^3 \xi_i = 0 \quad \sum_{i=1}^3 \xi_i^2 = -1 \quad \zeta^3 = -1$$

$$\begin{aligned} \mathcal{O} &= \text{signo}(0 \times 1 + 1 \cdot 1 - 1 \cdot 5) = \text{signo}(-0.5) \\ &= -1 = \zeta^3 \end{aligned}$$

$$\mu = 4 \quad \sum_{i=1}^4 \xi_i = 0 \quad \sum_{i=1}^4 \xi_i^2 = 0 \quad \zeta^4 = -1$$

$$\begin{aligned} \mathcal{O} &= \text{signo}(0 \times 1 + 0 \times 1 - 1 \cdot 5) = \text{signo}(-1.5) \\ &= -1 = \zeta^4 \end{aligned}$$



$$\bar{\omega} = (\omega_0, \omega_1, \omega_2)$$



# La función XOR (OR exclusivo)

$x$	$\xi_1$	$\xi_2$	$y$
-----	---------	---------	-----

1	1	1	-1
---	---	---	----

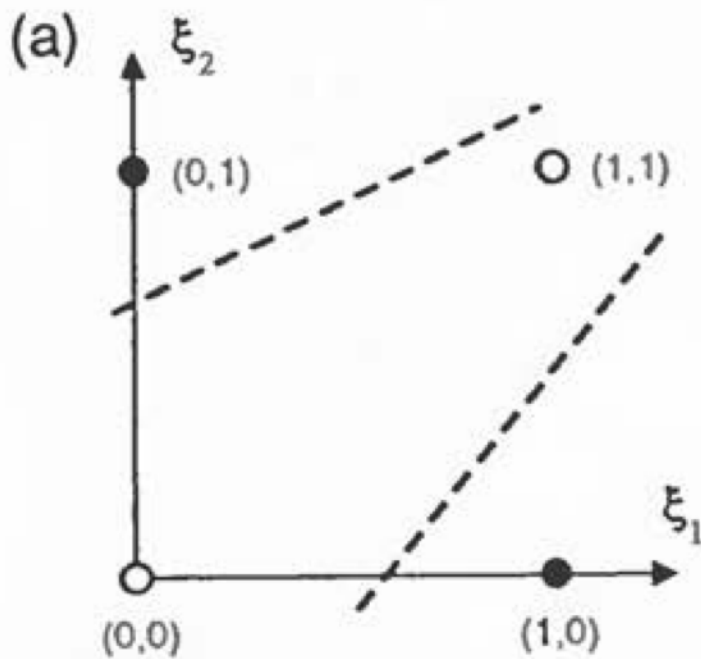
2	1	0	+1
---	---	---	----

3	0	1	+1
---	---	---	----

4	0	0	-1
---	---	---	----

$(\xi^1 \ \xi^2)$

$(w^1 \ w^2)$



Minsky y Papert (1969)

; No es linealmente separable!

Esto frenó el estudio de las redes neuronales.

El sistema a resolver es:

$$W_1 + W_2 - W_0 < 0 \Rightarrow W_1 + W_2 < W_0$$

(i)

$$-W_1 - W_2 - W_0 < 0 \Rightarrow -W_1 - W_2 < W_0$$

(ii)

$$W_1 - W_2 - W_0 > 0 \Rightarrow W_1 - W_2 > W_0$$

(iii)

$$-W_1 + W_2 - W_0 > 0 \Rightarrow -W_1 + W_2 > W_0$$

(iv)

De (i) y (iv)  $W_1 < 0$

De (ii) y (iii)  $W_1 > 0$

Esto muestra que no es posible encontrar  $W_1$ ,  $W_2$  y  $W_0$  que resuelvan el problema.

# UN ALGORITMO DE APRENDIZAJE

Comenzamos con valores aleatorios

$W_k$  con  $k=1, 2, \dots, N$ .

$$W_{ik}^{\text{nuevo}} = W_{ik}^{\text{viejo}} + \Delta W_k$$

$$\Delta W_{ik} = \begin{cases} 2\eta \sum_i^M \sum_k^M & \text{si } \sum_i^M \neq O_i^M \\ 0 & \text{si } \sum_i^M = O_i^M \end{cases}$$

$a_i$

$$\Delta W_{ik} = \eta (1 - \sum_i^M O_i^M) \sum_i^M \sum_i^M$$

$a_i$

$$\Delta W_{ik} = \eta (\sum_i^M - O_i^M) \sum_i^M$$

$\eta$  : learning rate.

Um proceso más "robusto" es:

$$\sum_i^M h_i^M \equiv \sum_i^M \sum_k w_{ik} \sum_k^M > N\alpha$$

$$\Delta w_{ik} = \eta \Theta(N\alpha - \sum_i^M h_i^M) \sum_i^M \sum_k^M$$

$\Theta$  : función escalón.

Este es la regla de aprendizaje del  
perceptron simple (Rosenblatt, 1962)  
para una única salida ( $M=1$ )

$$\vec{w} \cdot \vec{x}^M > N\alpha$$

Planes de problemas geométricamente

$$\Delta \bar{w} = \eta \Theta (N_k - \bar{w} \cdot \bar{x}^M) \bar{x}^M$$

0 sea  $\Delta \bar{w}$  apunta hacia  $\bar{x}^M$  si

$$N_k < \bar{w} \cdot \bar{x}^M = |\bar{w}| |\bar{x}^M| \cos(\alpha)$$

$\frac{N_k}{|\bar{w}|} < |\bar{x}^M| \cos(\alpha)$  : proyección de  $\bar{x}^M$  sobre  $\bar{w}$ .

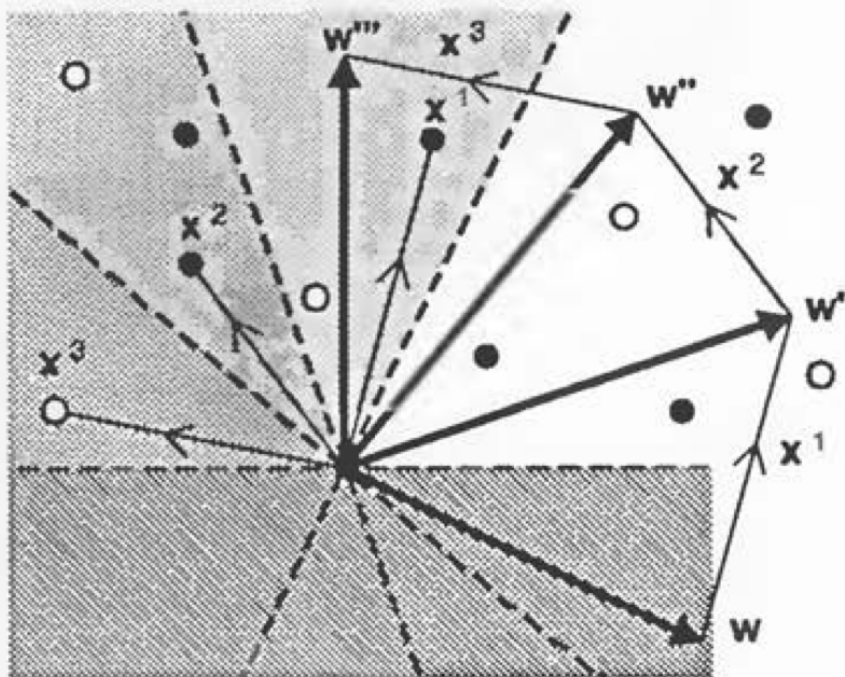
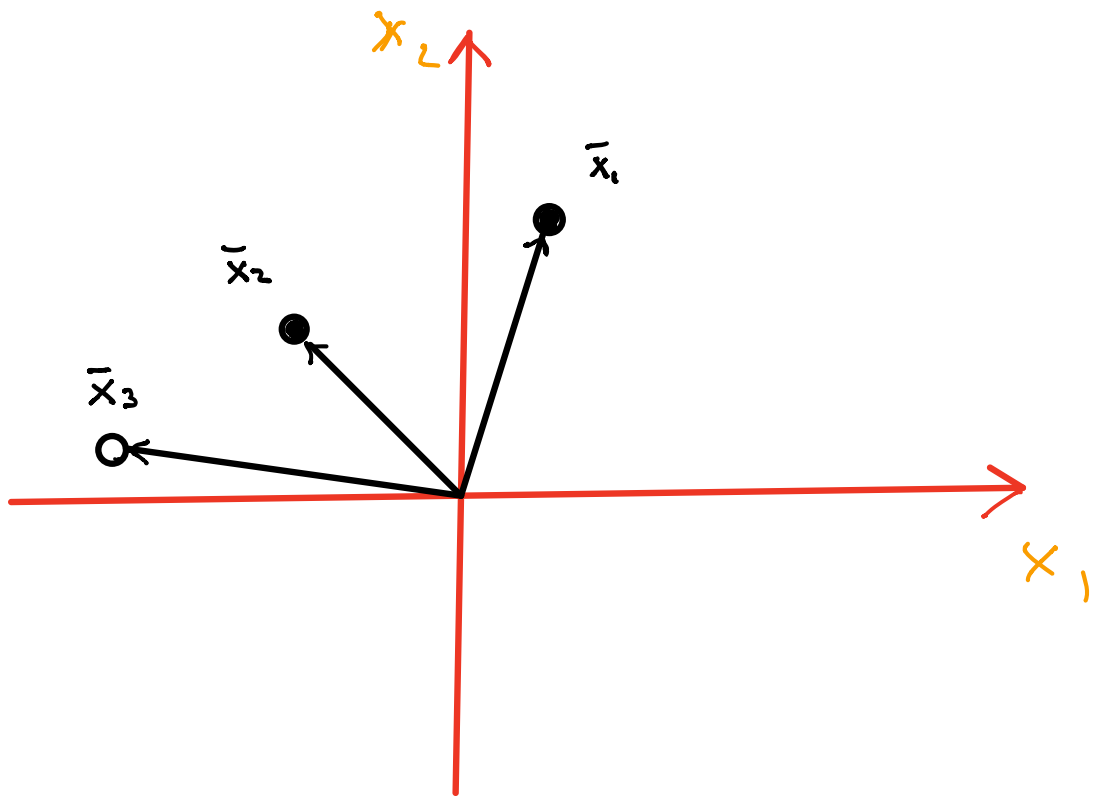


FIGURE 5.6 How the weight vector evolves during training for  $\kappa = 0$ ,  $\eta = 1$ . Successive values of the weight vector are shown by  $w$ ,  $w'$ ,  $w''$ , and  $w'''$ . The darker and darker shading shows the "bad" region where  $w \cdot x < 0$  for the successive  $w$  vectors. Each  $w$  is found from the previous one (e.g.,  $w'$  from  $w$ ) by adding an  $x^M$  from the current bad region.

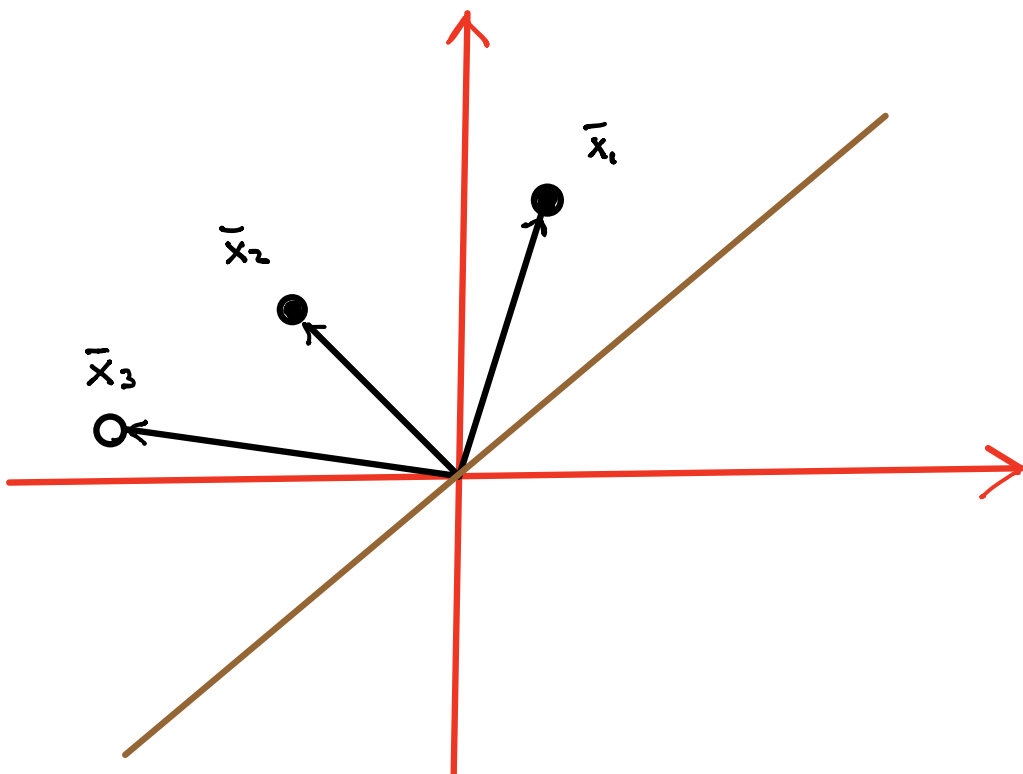
Ejemplo con  $\kappa = 0$

$$D(\bar{\omega}) = \frac{1}{|\bar{\omega}|} \min_{\mu} (\bar{\omega} \cdot \bar{x}^{\mu})$$

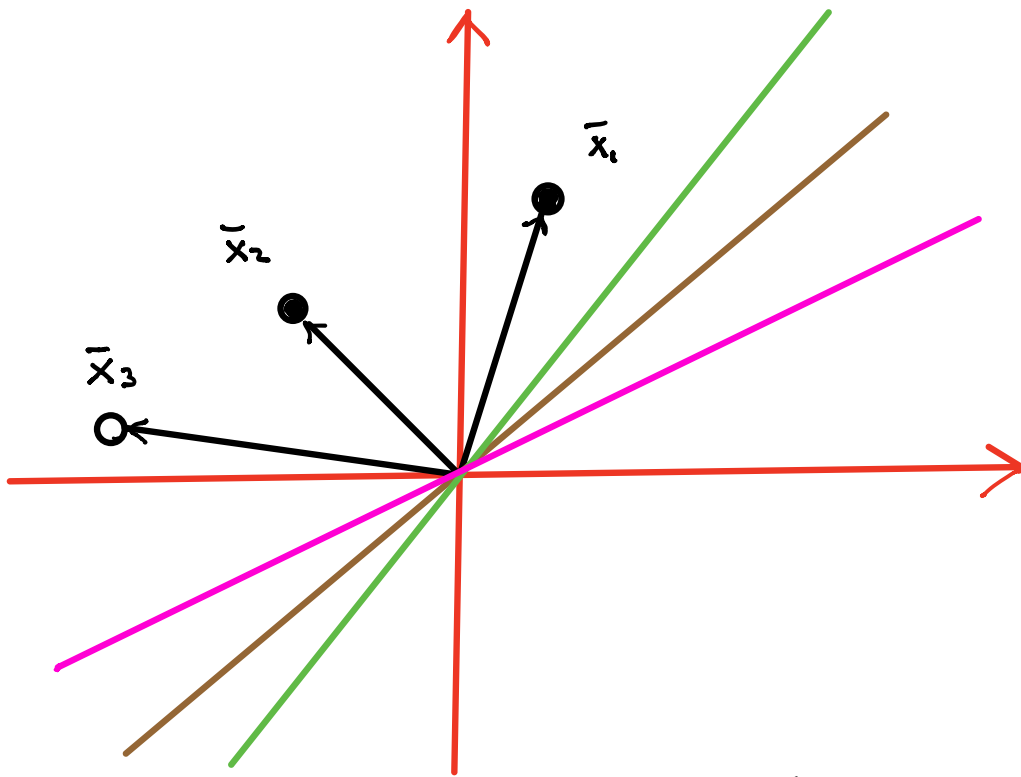
$$D_{\max} \equiv \max_{\bar{\omega}} D(\bar{\omega})$$



Vemos que el problema admite solución pues todos los puntos de entrensamiento quedan en un semi plano. Como  $x=0$ , la superficie que refiere en 2 semi-planos debe pasar por el origen.



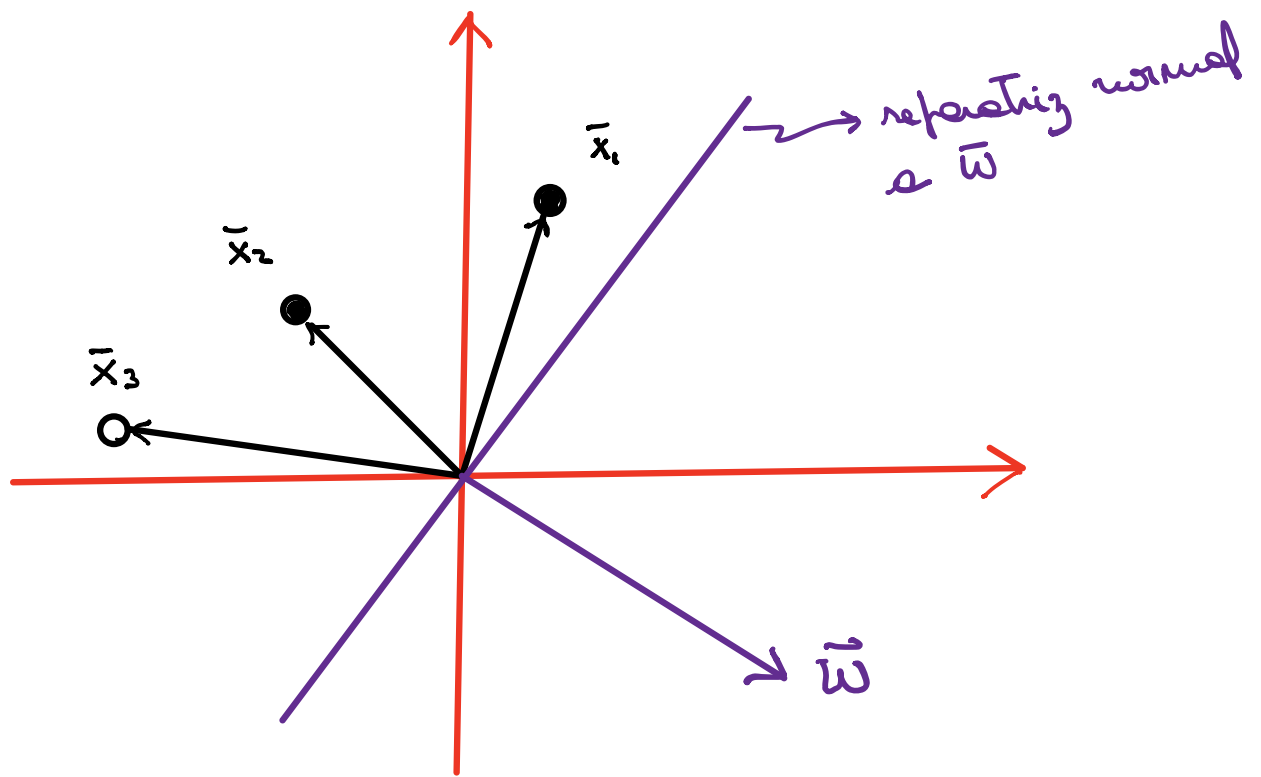
La línea marrón es una posible separatrix, pero hay muchas otras posibles. Este es un problema fácil.



Vemos que las líneas verde y rosa también son solución. Recuerden que  $\vec{w}$  debe ser normal a la separatrix.

Ahora elegimos un vector  $\vec{w}$  inicial arbitrario y lo graficamos en el plano  $\bar{x} = (x_1, x_2)$ .





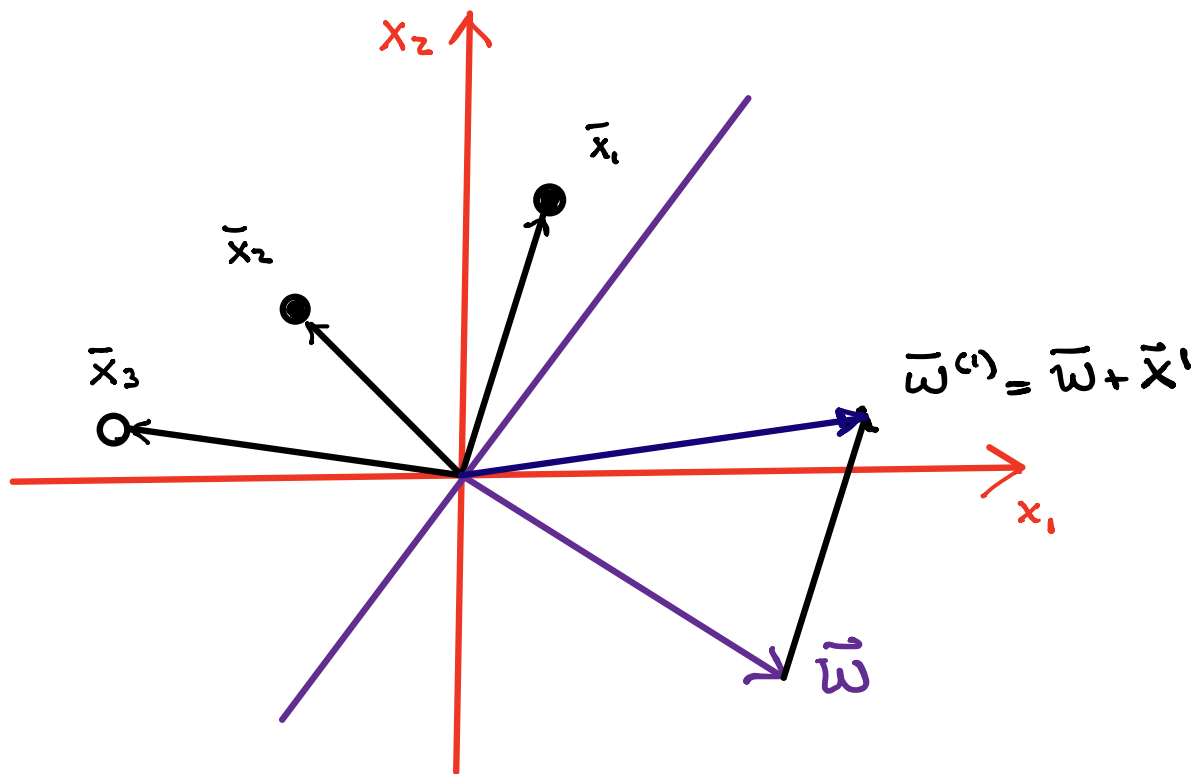
Vemos que este  $\bar{w}$  no cumple la condición de tener proyección positiva con todos los  $\bar{x}^i$  (en este caso  $i=1, 2, 3$ ). Este  $\bar{w}$  en particular es HORRIBLE, pues no tiene proyección positiva con ningún ejemplo 😞

Ahora aplico la regla de aprendizaje con el ejemplo 1. Como la proyección con  $\bar{x}^1$  es negativa

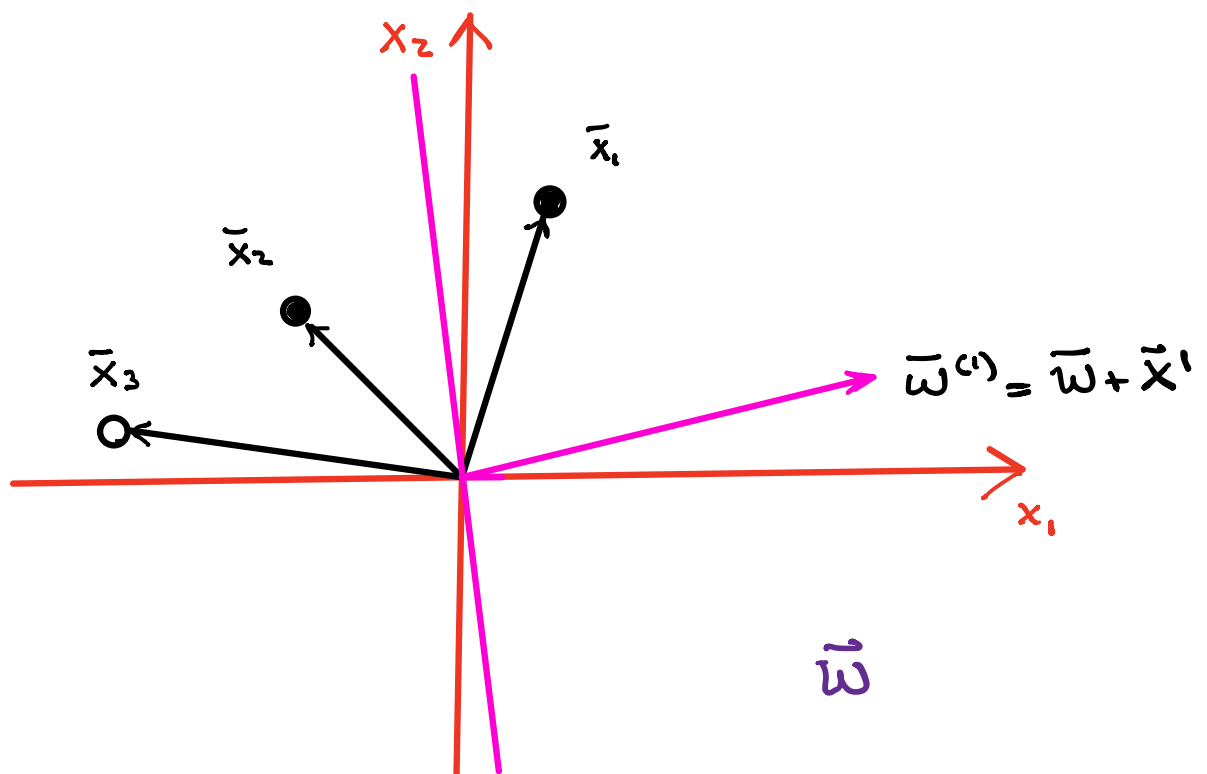
$$\bar{w} \longmapsto \bar{w} + \eta \bar{x}^1$$

Si hago  $\eta=1$

$$\bar{w} \longmapsto \bar{w} + \bar{x}^1 = \bar{w}^{(1)}$$



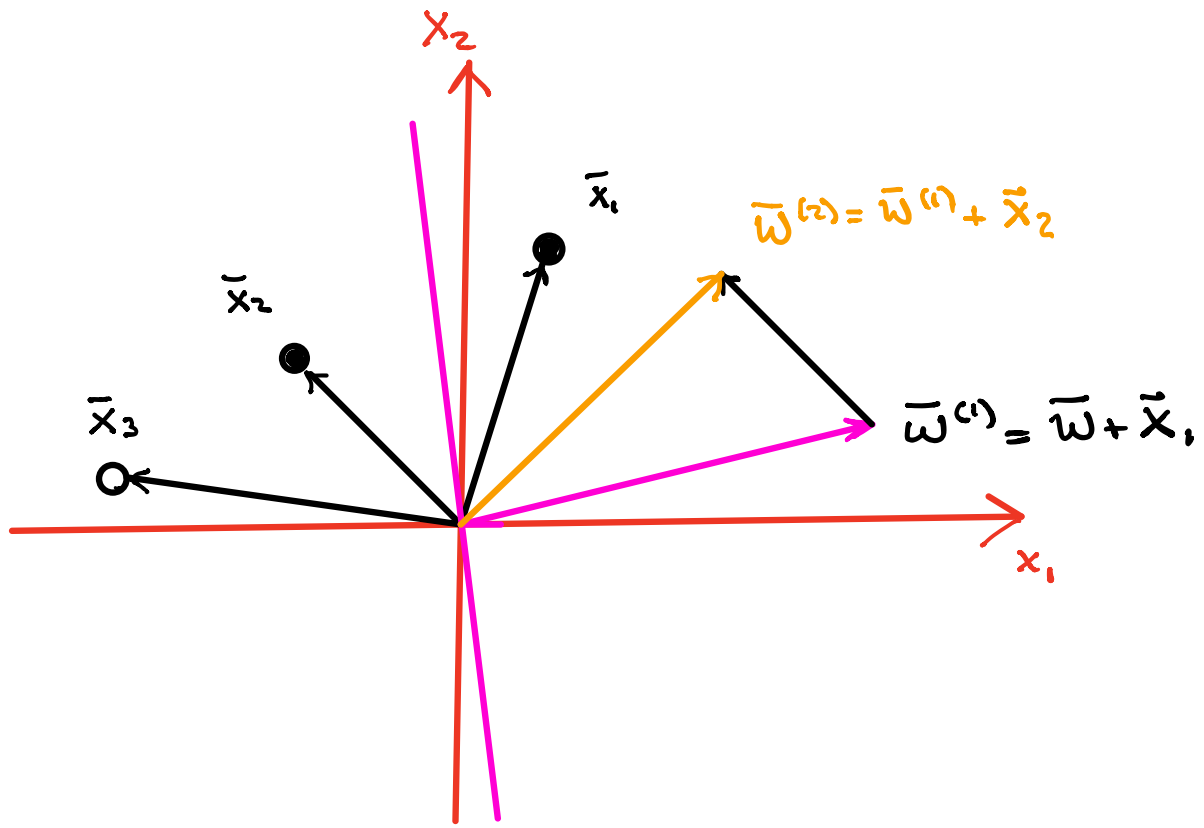
Buscamos el  $\vec{w}$  inicial y su recta normal para  
 aclarar



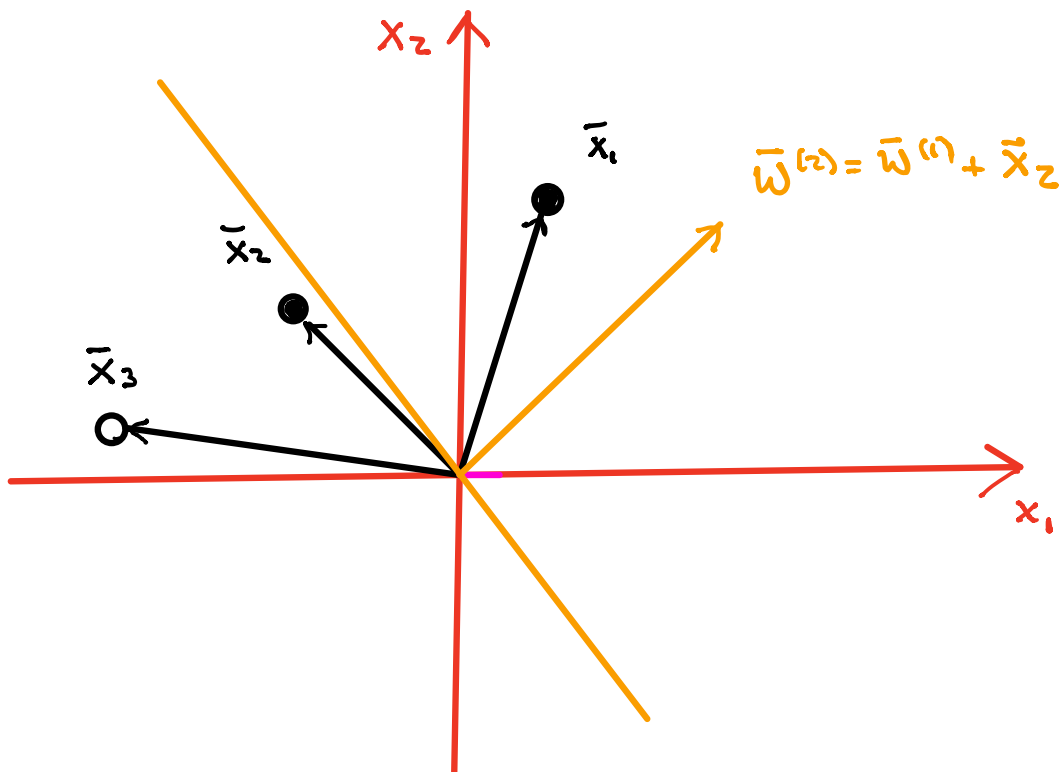
Noten que ahora el nuevo  $\vec{w}^{(1)}$  tiene el menor  
 proyección positiva con el ejemplo  $\mu=1$

Ahora repetimos lo mismo para  $k=2$ .

$$\bar{w}^{(2)} = \bar{w}^{(1)} + \bar{x}_2$$



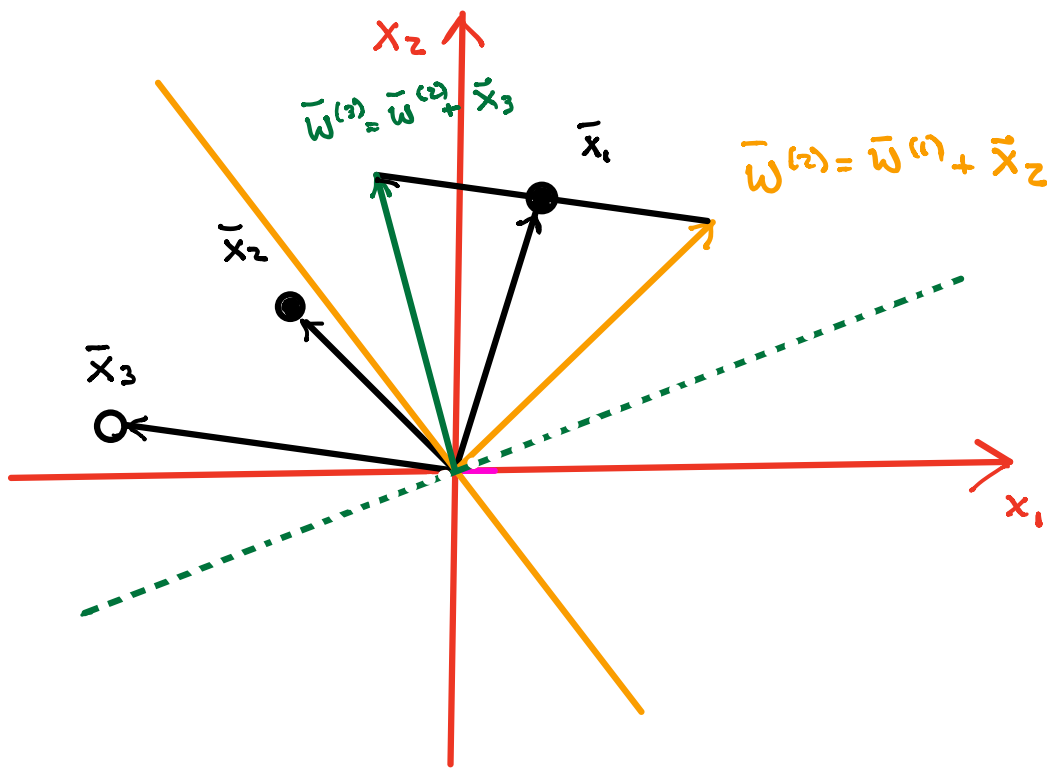
Repetimos lo mismo:



Queremos teniendo proyección positiva solo con  $\mu=1$  pero estamos más cerca de resolver  $\mu=2$  😊

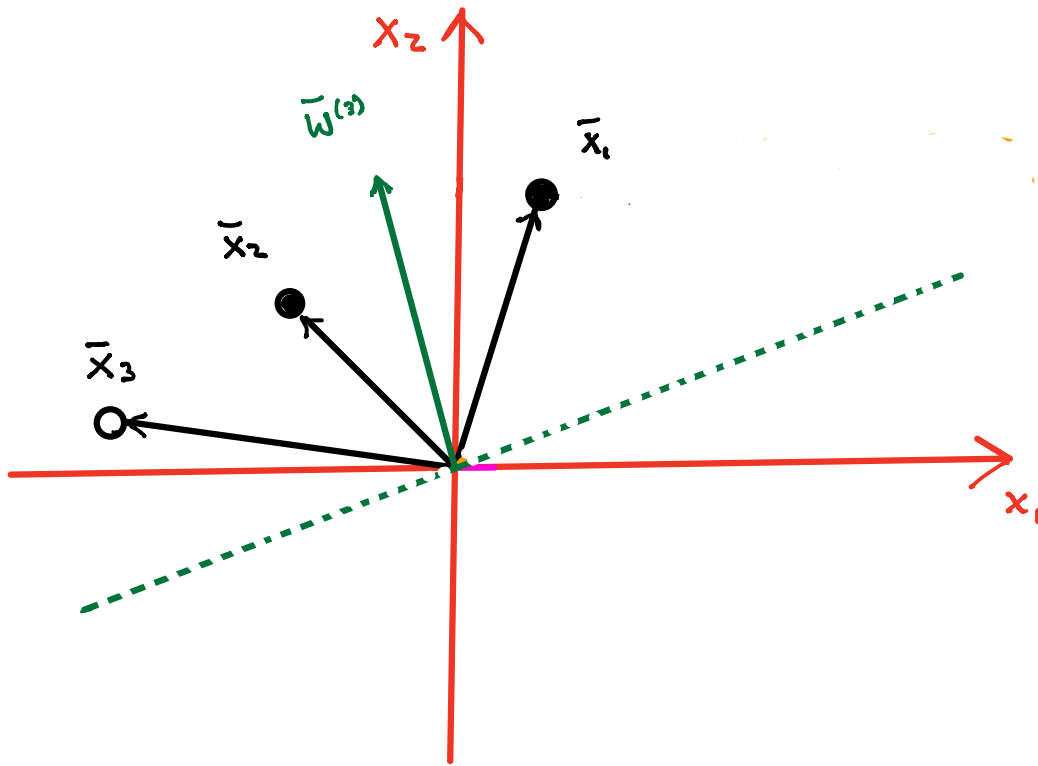
Repitamos para  $\mu=3$ .

$$\bar{w}^{(3)} = \bar{w}^{(2)} + \bar{x}^3$$



Ya resolvimos el problema en tres pasos, visitando una vez cada memoria. Noten que  $\bar{w}^{(3)}$  define una recta perpendicular a él, que pasa por el origen y que todos los  $\bar{x}_\mu$  están de un lado del plano (en un semiplano), y es el lado en el que vive  $\bar{w}^{(3)}$ .

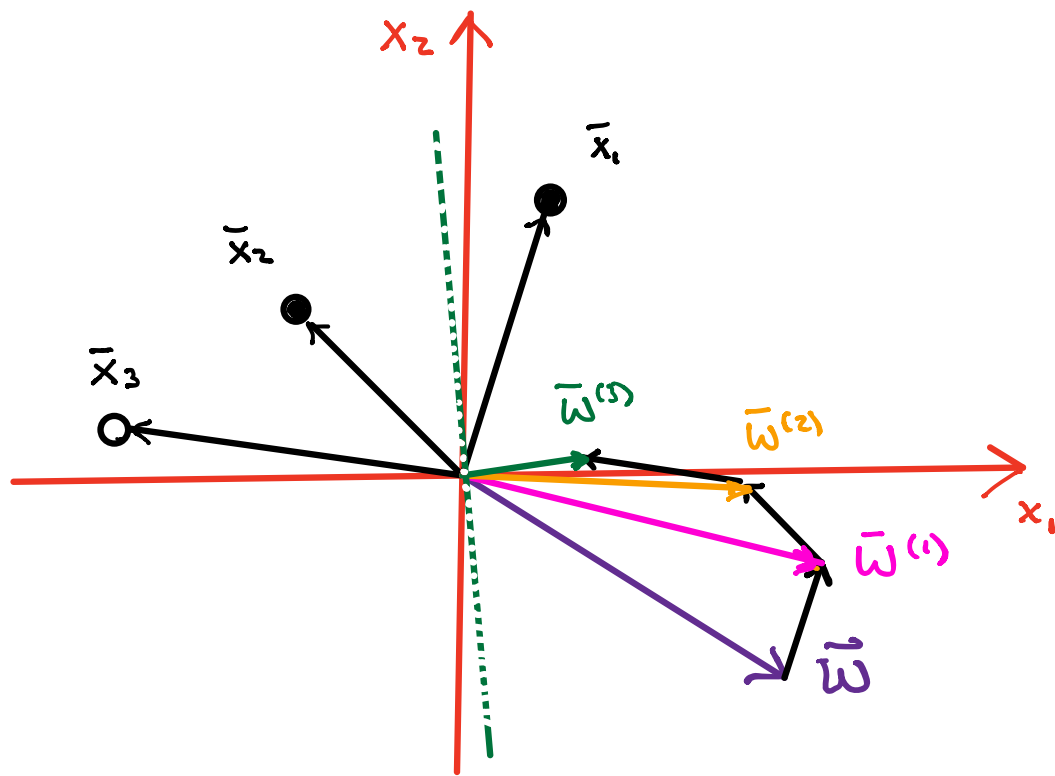
Esto quiere decir que  $\bar{w}^{(3)}$  tiene proyección positiva con los tres vectores  $\bar{x}_1$ ,  $\bar{x}_2$  y  $\bar{x}_3$



Este ejemplo es muy muy simple, pero los casos reales son mucho más complicados.

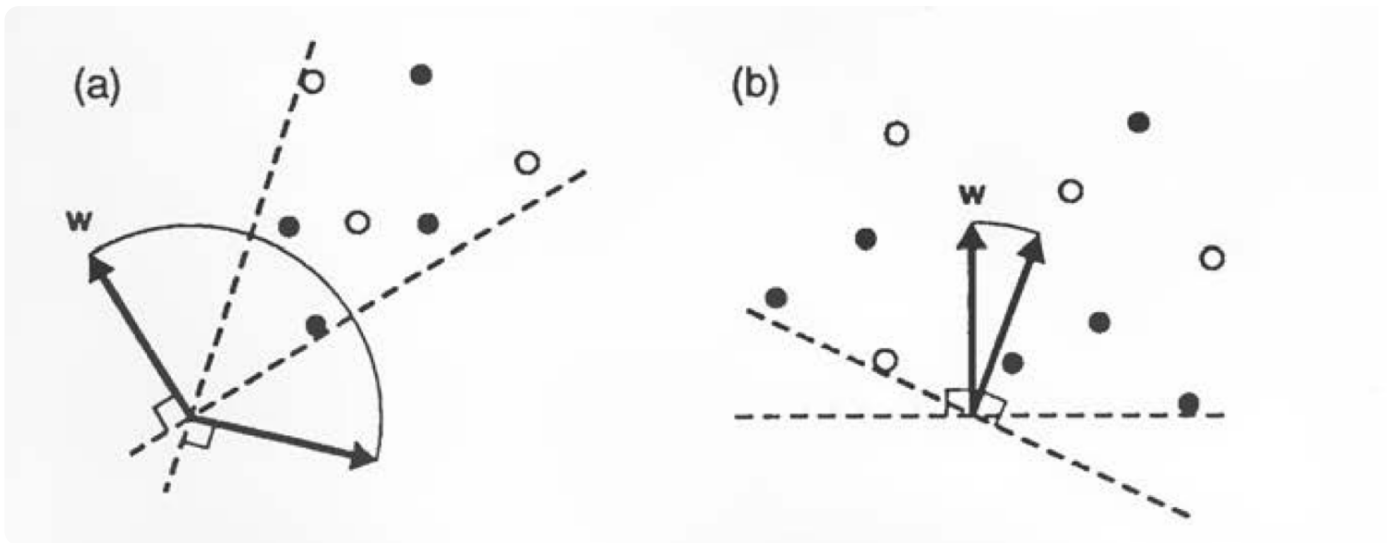
- con unce tenemos 2 variables de entrada sino muchos más.
- con unce tenemos 3 ejemplos etiquetados en el conjunto de entrenamiento, sino muchos más.
- con unce funciona usar  $\eta = 1$ , en lo cual todo se hace mucho más lento

Abajo nuestro como hubiere sido el proceso si hubiésemos tomado  $\eta = \frac{1}{2}$



Vemos que  $\bar{w}^{(3)}$  obtenido con  $\eta = \frac{1}{2}$  no resuelve pero vemos que si tenemos muchos ejemplos en  $\eta$  pequeños (i.e. casi pequeños) es mucho mejor. Nuestro algoritmo dice que ahora debemos volver a repetir el mismo procedimiento pero solo para los ejemplos que aún no fueron resueltos (en este caso  $\mu=2$ ,  $\mu=3$  pues  $\mu=1$  está resuelto con  $\bar{w}^{(3)}$ ).

Vemos que en general hay problemas más fáciles que otros



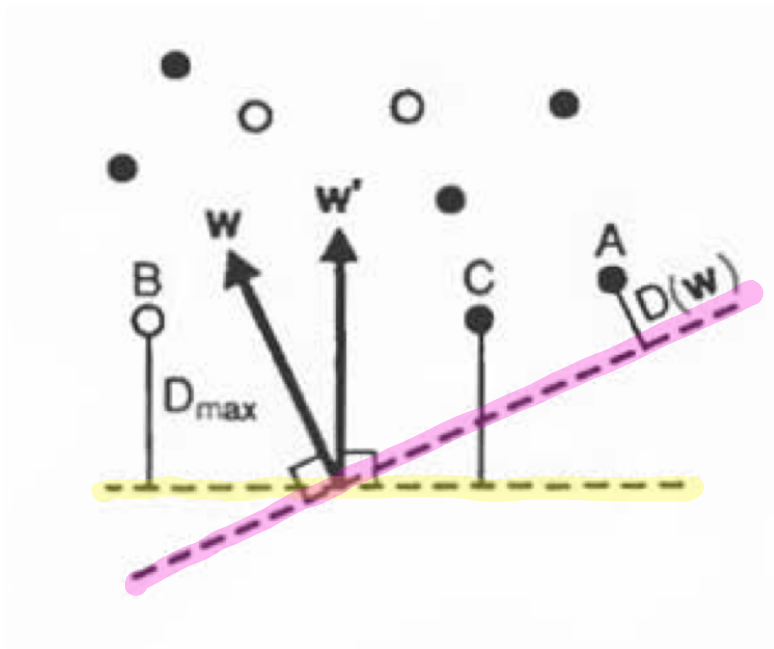
Discutamos por qué el ejemplo (a) es fácil y el (b) es difícil.

Notemos que no existe una única solución.

- si cambiamos el módulo del vector  $\bar{w}$  que resuelve el problema, el nuevo también lo resuelve.

(no podemos cambiar la dirección).

- además de cambiar el módulo tenemos cierto margen de variación para elegir  $\bar{w}$ .



Vemos que la separación sencilla es más robusta que la rosa, y por ende preferimos  $\bar{w}'$  a  $\bar{w}$ . Tomaremos como "OPTIMA" aquella que nos da el mayor  $D(\bar{w})$ . Recuerden que

$$D(\bar{w}) = \frac{1}{|\bar{w}|} \min_{\mu} (\bar{w} \cdot \bar{x}^{\mu})$$

$$D_{\max} \equiv \max_{\bar{w}} D(\bar{w})$$



Pero, ¿cómo sabemos que encontraremos una solución? Gracias a la matemática

## Teorema de convergencia de la regla del perceptron

El teorema tiene una formalización compleja para este curso. Básicamente establece que el algoritmo de la regla de aprendizaje converge a una solución en un número finito de pasos, suponiendo que tenemos un conjunto de datos de entrenamiento linealmente separables.

No lo demostramos aquí.

Si tenemos suerte, este  $w$  resolverá también casos que no estaban en el conjunto de entrenamiento. Esto significa que puede generalizar a partir de los ejemplos aprendidos.

## Resumen hasta aquí.

- Consideramos un perceptron de una capa de  $M$  neuronas y lo referenciamos en  $M$  perceptrons simples (con una única neurona de salida, o sea,  $M=1$ ).
- Suponemos que la función de activación es la función signo (o sea  $\sigma = \pm 1$ ) y no nos preocupamos por los entrenados.
- Vimos que el conjunto de datos de entrenamiento deben ser **LINEALMENTE SEPARABLE**.
- Si los datos son linealmente separables existe un algoritmo de aprendizaje que nos da una solución posible en un número finito de iteraciones.

Nosotros podemos fijar el valor de  $\eta$  y el orden en que visitamos los ejemplos.

Combinando  $M$  perceptrons simples podemos hacer un perceptron de  $N$  entradas y  $M$  salidas.