



REDES NEURONALES 2024

Clase 15 parte 3

Lunes 7 de octubre 2024

FAMAF, UNIVERSIDAD NACIONAL DE CÓRDOBA

INSTITUTO DE FÍSICA ENRIQUE GAVIOLA (UNC-CONICET)

J12/10/23 c15p2

Hasta aquí vimos bajo cuáles condiciones existe una solución al problema de clasificación, pero no tenemos ninguna regla para construir un vector sináptico \bar{W} que resuelva el problema.

Y si esto lo debemos trabajar con un número grande de entradas, como sucede en la vida real y con un número muy grande de elementos en el conjunto de entrenamiento, el problema se torna inviable si no lo hacemos de forma algorítmica.

A continuación presentaremos el primer y más simple algoritmo de aprendizaje de este curso y de la historia de la Inteligencia Artificial.

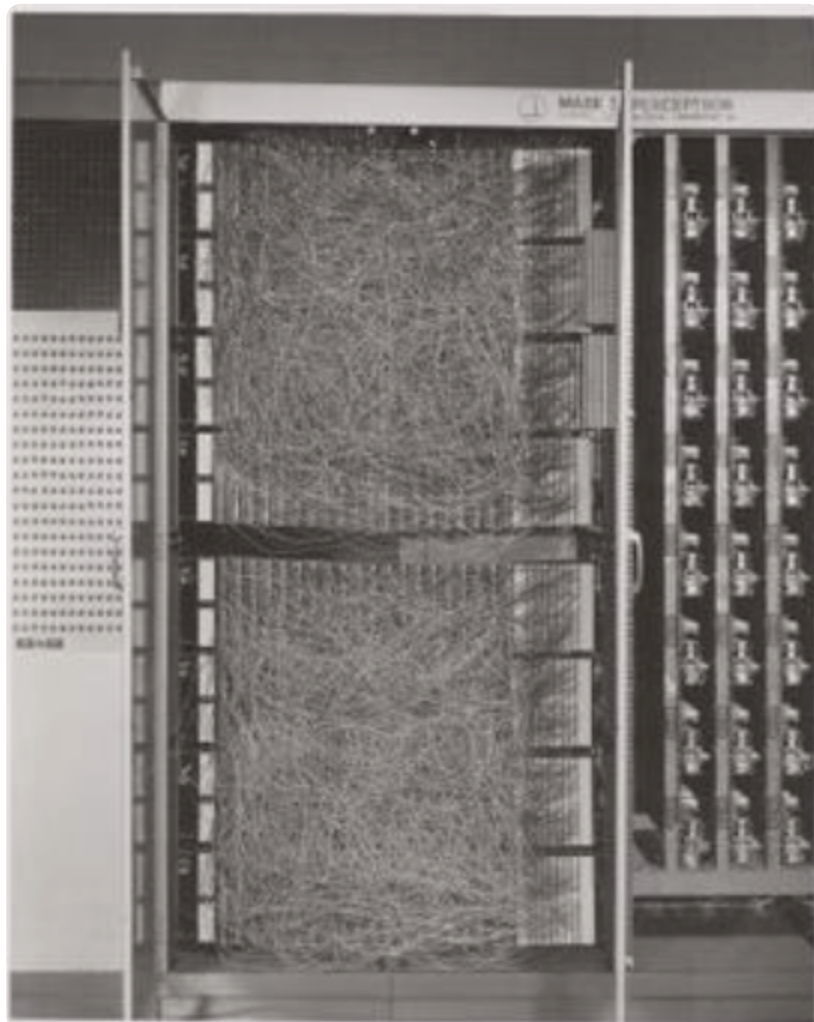
La idea de una neurona artificial matemática, recordemos, había sido introducida por *McCulloch and Pitts* en 1943, inspirados en los trabajos pioneros sobre el funcionamiento de las neuronas de *Santiago Ramón y Cajal* y *Charles Scott Sherrington*. Aquí aparece otra persona central de la historia de la inteligencia artificial llamada Frank Rosenblatt. Frank era un joven formado en ciencias que desde temprano se orientó al estudio de la psicología cognitiva. Fue contratado en el *Cornell Aeronautical Laboratory* en la ciudad de Buffalo, donde haría sus grandes contribuciones. Pero a nosotros nos interesa el desarrollo del primer algoritmo para encontrar automáticamente los pesos del modelo de McCulloch y Pitts. Este fue el punto inicial del Aprendizaje Automático.

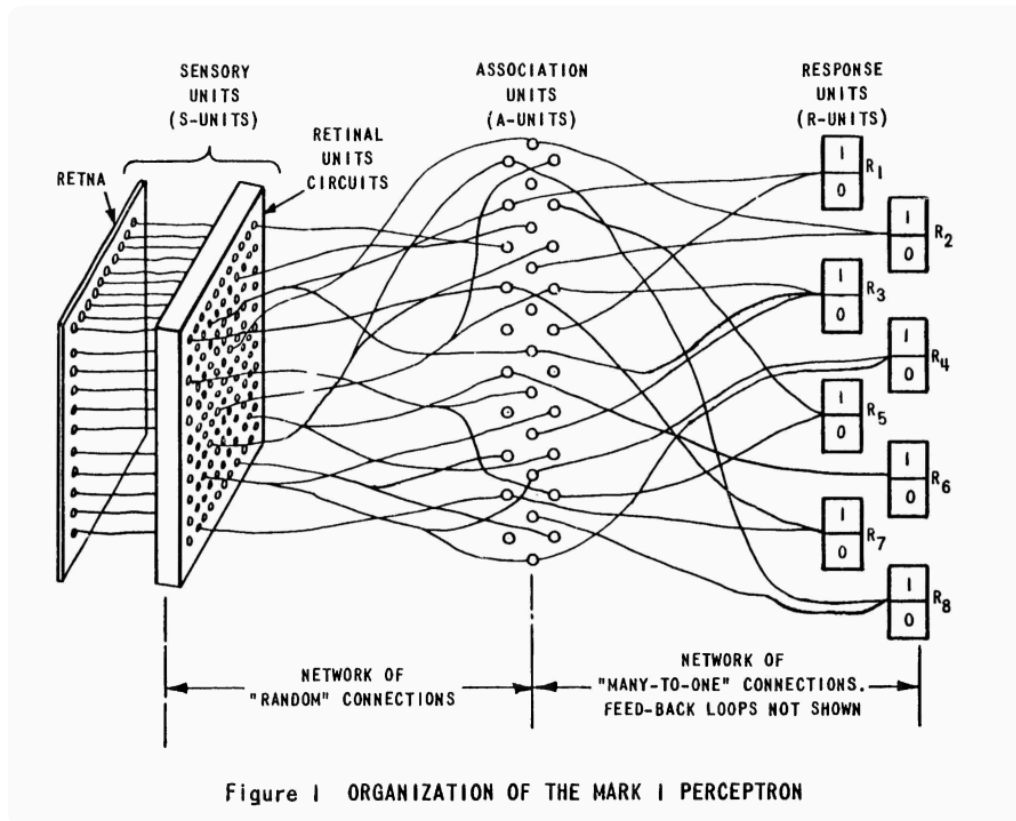
Veremos que Rosenblatt generó el primer programa convencional que podía aprender iterativamente los pesos sinápticos de un perceptron simple por medio de la realización de pruebas y las evaluaciones de los correspondientes errores, como hacemos los humanos y como hacemos los científicos naturales.

Primero implementó el algoritmo por software en una IBM 704, la primera computadora comercial de la historia de la cual se vendieron apenas 120 ejemplares. Esta computadora hacia solo 40000 operaciones de punto flotante (de números reales) por segundo. Hoy un procesador I7 realiza 100.000,000.000 operaciones por segundo, o sea, poco más de dos millones de veces más rápido y mide unos pocos milímetros cuadrados.

Rosenblatt no se conformó con desarrollar el algoritmo que veremos hoy. Él implementó el algoritmo físicamente en una máquina llamada **MARK I PERCEPTRON**, dedicada exclusivamente a clasificar imágenes de 400 píxeles (20 x 20). En particular, se eligió clasificar imágenes de hombres y mujeres.

En 1958 Rosenblatt escribió en un artículo de opinión en el New York Times en el cual afirmaba: *“Este es el embrión de una nueva computadora electrónica que espera poder caminar, hablar, ver, escribir, reproducirse y tener consciencia de su existencia”*. Como es de imaginar, le llovieron las burlas. Hoy sabemos que él estaba creando las bases de la algorítmica del actual Deep Learning.

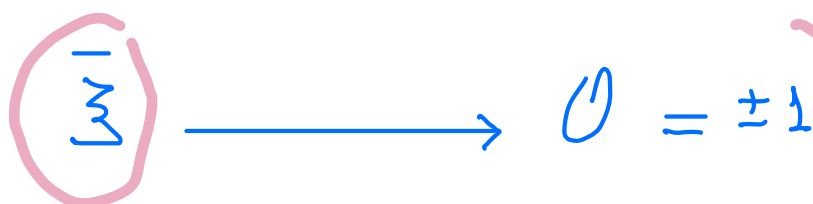




LA REGLA DE APRENDIZAJE DEL PERCEPTRON

Veamos entonces el primer algoritmo de aprendizaje de máquina.

Supongamos que tenemos un **PERCEPTRÓN** simple con N entradas y una neurona de salida binaria.



Supongamos que tenemos p entradas etiquetadas correctamente, a las cuales le asignamos la salida (el valor de la neurona) correcta.

$$\bar{x}^\alpha \longrightarrow y^\alpha = \pm 1 \quad (\alpha = 1, 2, \dots, p)$$

Finalmente, supongamos que el problema es linealmente separable. Nos preguntamos entonces

¿Cómo encontrar un \bar{W} adecuado?

Elegimos, ante la falta de mejor criterio, los parámetros iniciales de nuestro perceptrón simple (las $N+1$ componentes de W) de forma aleatoria. Por ejemplo, podemos elegir cada componente W en forma independiente entre sí y con distribución gaussiana centrada en cero y con desviación 1 .

La idea central del algoritmo consiste en presentar sucesivamente a la red cada elemento del conjunto de entrenamiento, ya sea en un orden pre-establecido o aleatorio, y, dependiendo de que si el resultado es correcto o no, modificar levemente los parámetros.

$$y^\alpha = g(h^\alpha)$$

Si $0^\alpha = \xi^\alpha$ no cambiamos nada, pero

si $0^\alpha = -\xi^\alpha$ castigamos al perceptrón cambiando
levemente todos los acoplamientos.

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

$$\Delta W_k^\alpha = \begin{cases} 2 \cdot \xi^\alpha \xi_k^\alpha & 0^\alpha \neq \xi^\alpha \\ 0 & 0^\alpha = \xi^\alpha \end{cases}$$

$$\Delta W_k^\alpha = \eta (1 - \xi^\alpha 0^\alpha) \xi_k^\alpha \xi^\alpha$$

$$= \eta (\xi^\alpha - \xi^\alpha 0^\alpha) \xi_k^\alpha = \eta (\xi^\alpha - 0^\alpha) \xi_k^\alpha$$

$$\bar{W}^{\text{nuevo}} = \bar{W}^{\text{viejo}} + \Delta \bar{W}^\alpha$$

Y

En forma vectorial: $\Delta \bar{W}^\alpha = 2 \eta \xi^\alpha \bar{\xi}^\alpha$ si $\xi^\alpha \neq 0^\alpha$

Si la respuesta es la equivocada, llevamos la componente k del vector \bar{W} hacia el vector $\bar{\xi}^\alpha$ si $\xi^\alpha = 1$, y hacia $-\bar{\xi}^\alpha$ si $\xi^\alpha = -1$.

El parámetro η se llama *razón de aprendizaje* y hoy lo presentamos para que siempre nos acompañe. Él regulará, como veremos, la velocidad de aprendizaje y deberemos invertir mucho esfuerzo, en cada caso, para asignarle un valor adecuado.

Es un *hiperparámetro*, pues una vez determinado, en general (no siempre) estará fijo durante el aprendizaje. Para llegar al aprendizaje profundo actual, fue preciso invertir enormes esfuerzos en entender el sentido y función de este hiperparámetro. Ya lo veremos.

Recordemos la condición de estabilidad que debe cumplir el conjunto de pesos sinápticos:

$\mathcal{U}^\alpha = \text{sign}(\mathbf{h}^\alpha) = \mathcal{J}^\alpha$

↑ Lo que sale

↑ Lo que quiero que salga

$$\int^a \text{sign} \theta(h^2) = \int^a \int^d = 1$$

$$\text{sign}_{\Delta}(\sum^{\alpha} h^{\alpha}) = 1$$

$$\sum^\alpha h^\alpha > 0 \quad \left\{ \begin{array}{l} \text{con esto} \\ \text{elección} \end{array} \right.$$

Ya tenemos una regla de aprendizaje automático que encontrará una buena solución (no necesariamente la mejor) en un tiempo finito siempre y cuando el conjunto de entrenamiento sea linealmente separable.

Un pseudo código

Entrada:

η y una secuencia de ejemplos de entrenamiento, o sea, de entradas \bar{x}^α y sus correspondientes etiquetas y^α correctas:

$$\mathcal{C}(\{\bar{x}^\alpha, y^\alpha\}) \quad \alpha = 1, 2, \dots, p$$

Código:

Inicializamos \bar{w}_0 aleatorio (en \mathbf{R}^n)

Para cada ejemplo α (\bar{x}^α, y^α) en el conjunto de entrenamiento:

$$o^\alpha = \text{signo}(\bar{x}^\alpha \cdot \bar{w})$$

Si

$$o^\alpha \neq y^\alpha$$

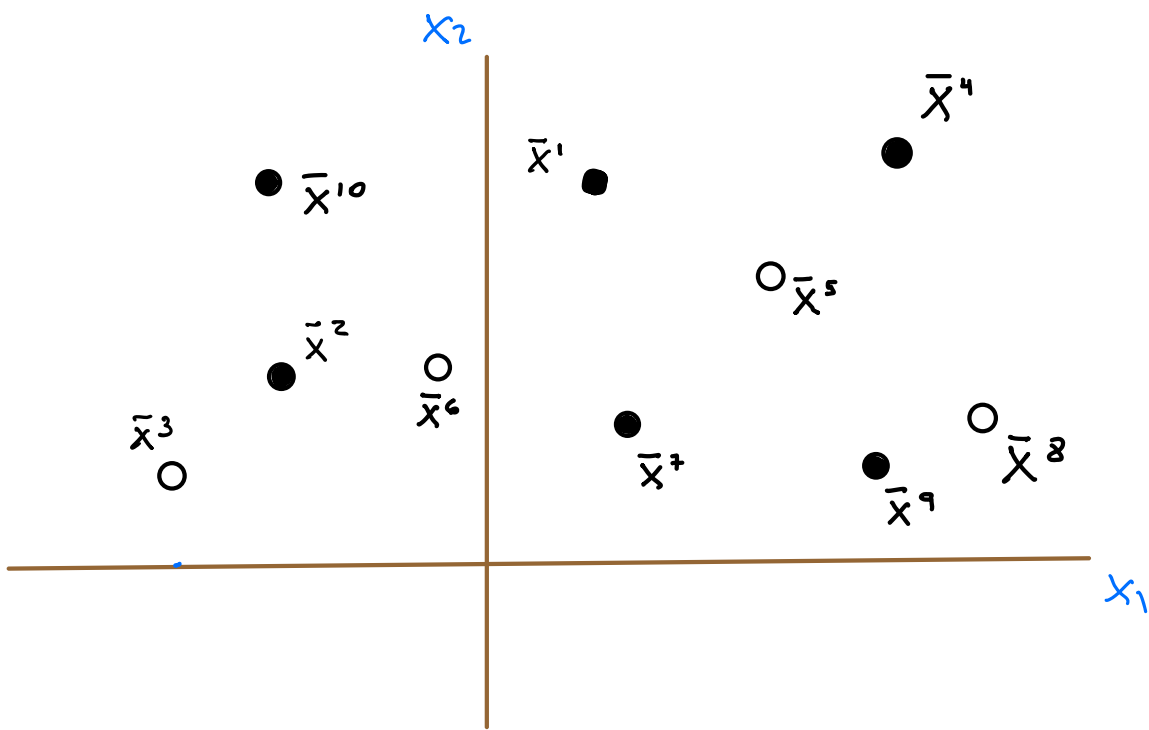
$$\Delta \bar{w} = 2\eta y^\alpha \bar{x}^\alpha$$

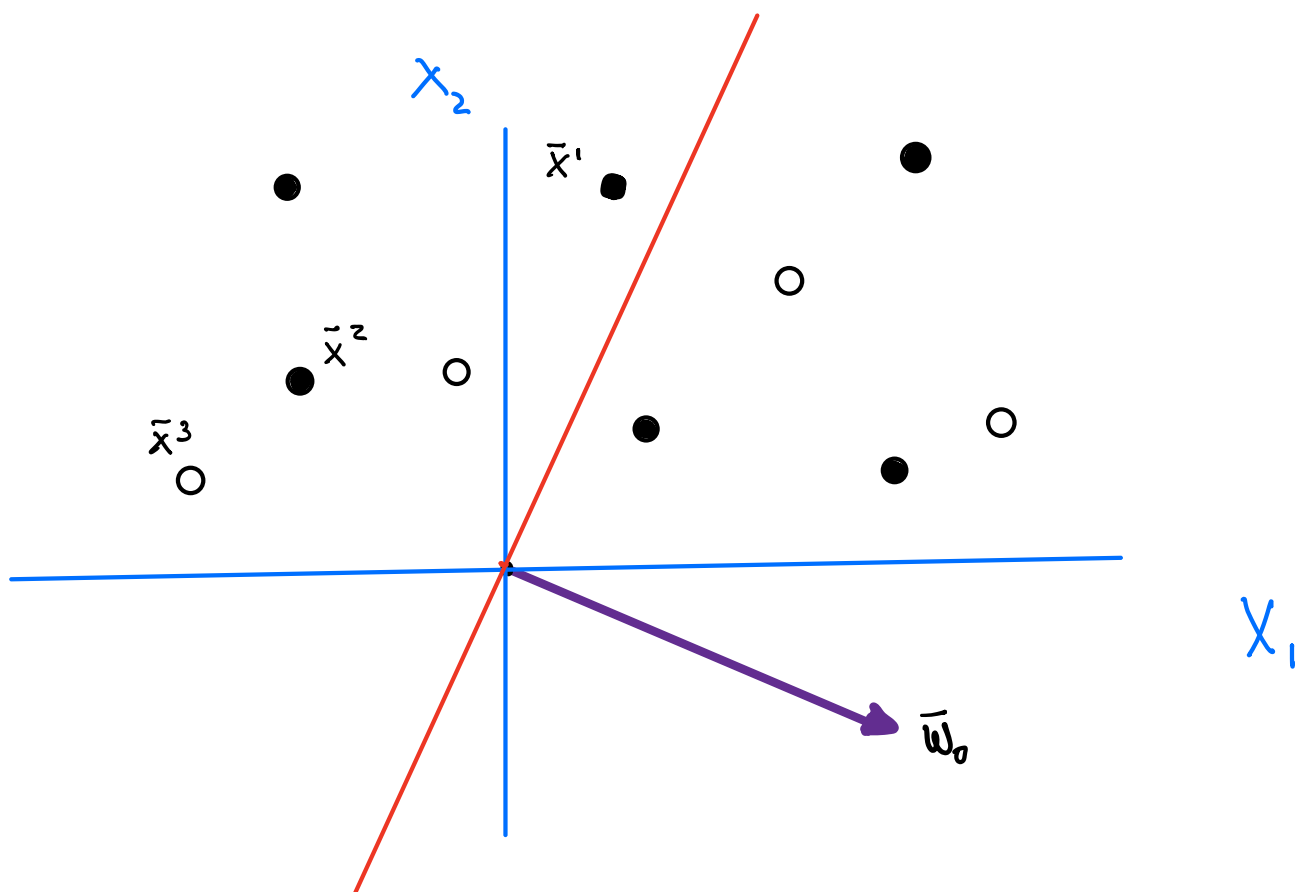
$$\bar{w}_{t+1} \leftarrow \bar{w}_t + \Delta \bar{w}$$

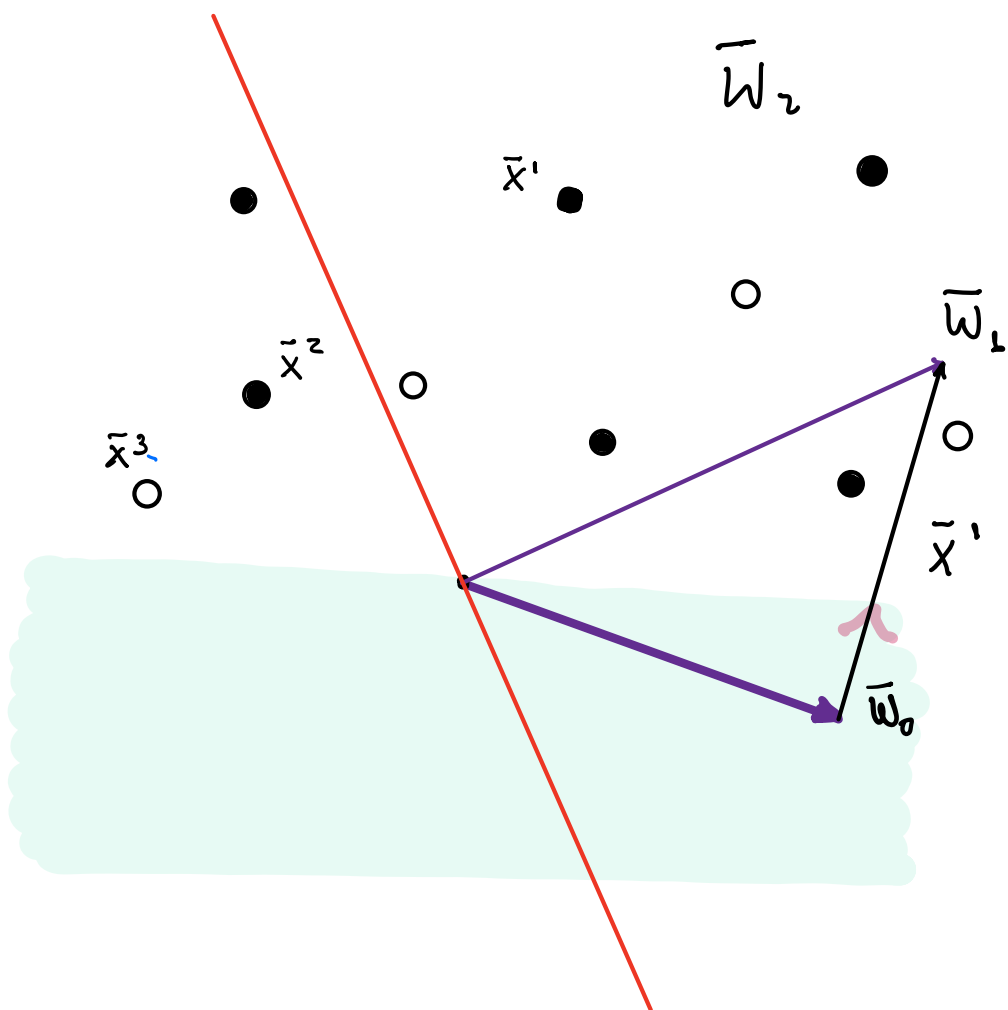
Salida: el vector final de acoplamiento \bar{w}_f .

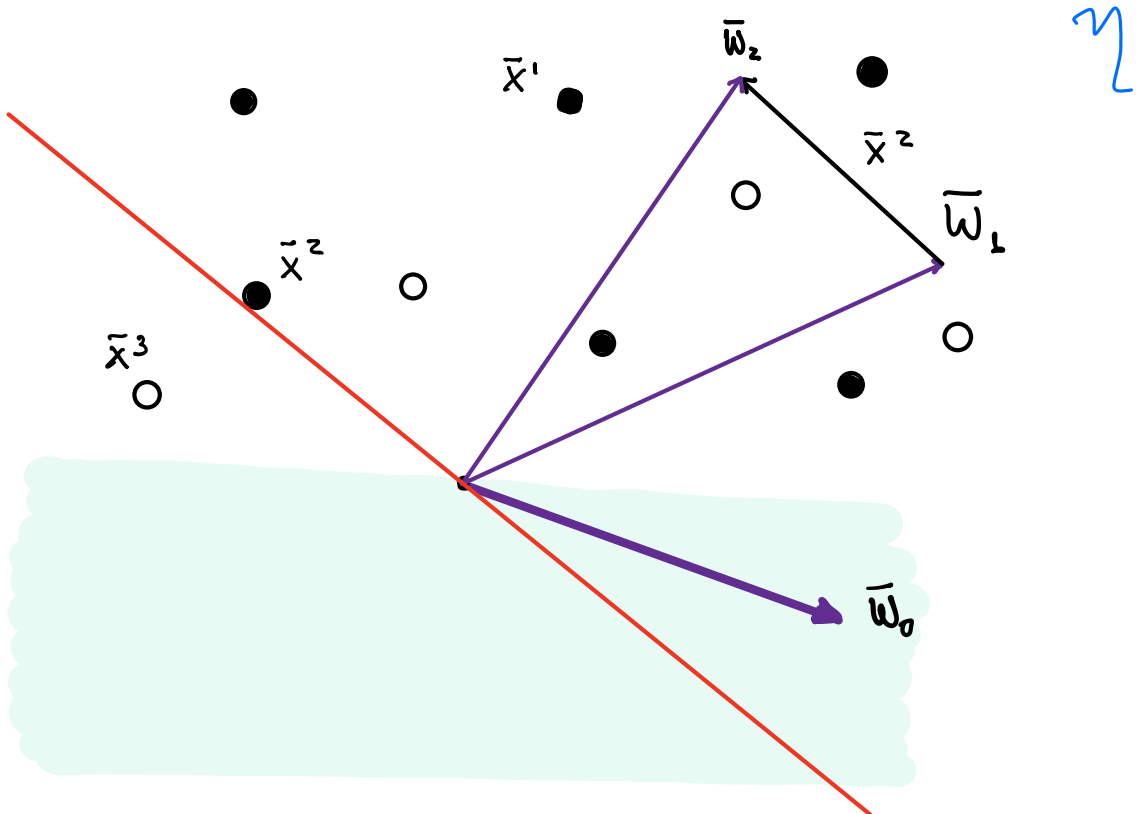
Miremos cómo funciona el algoritmo:

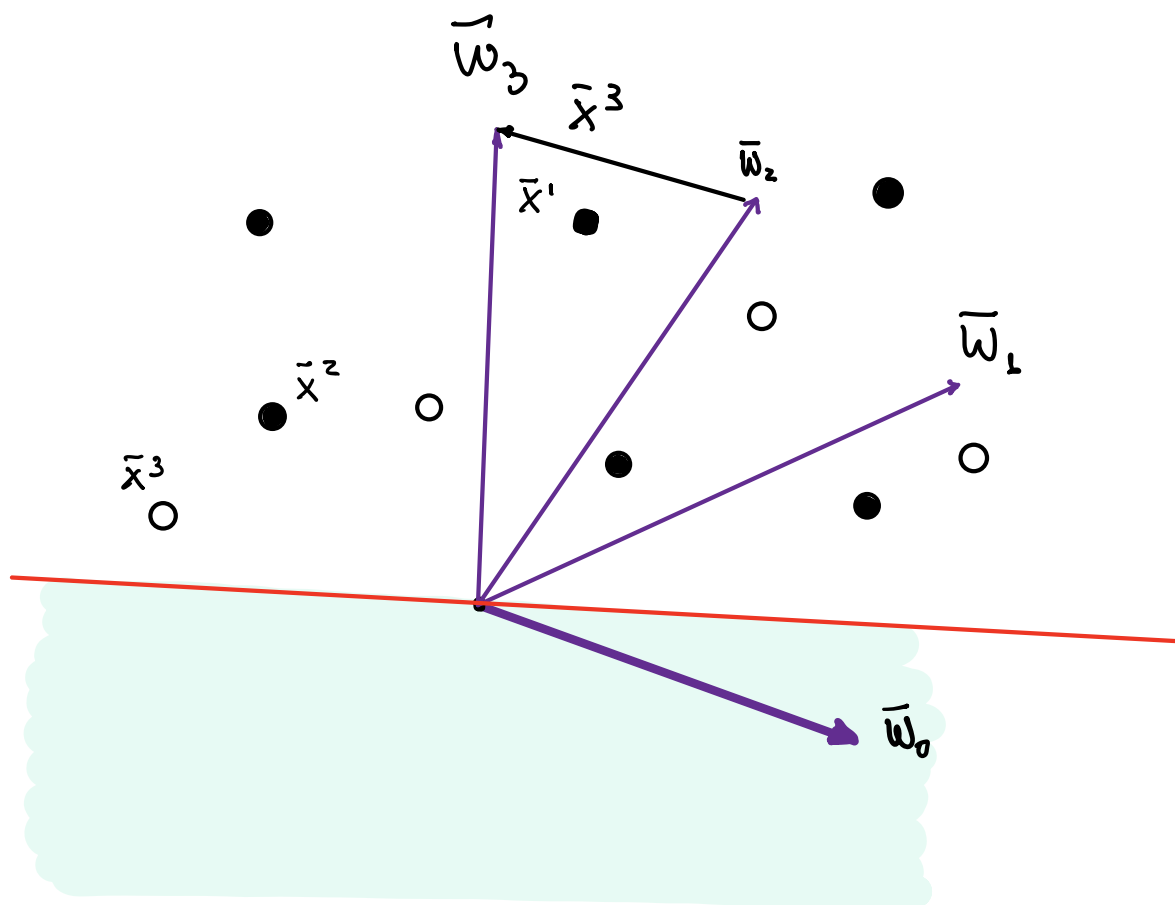
$$\eta = \frac{1}{2}$$

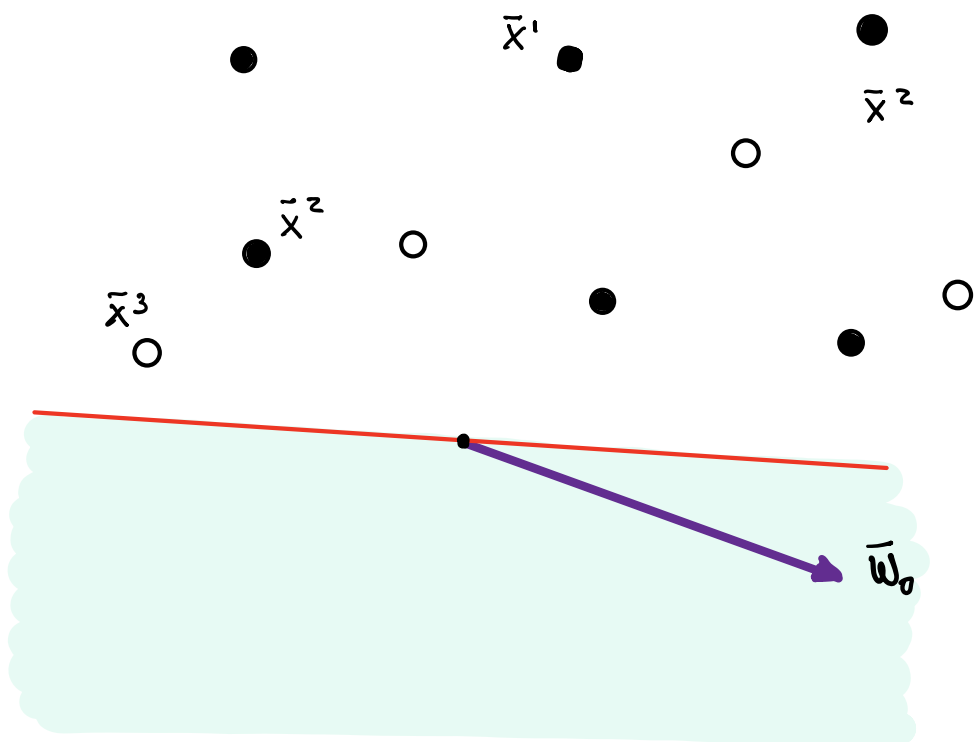




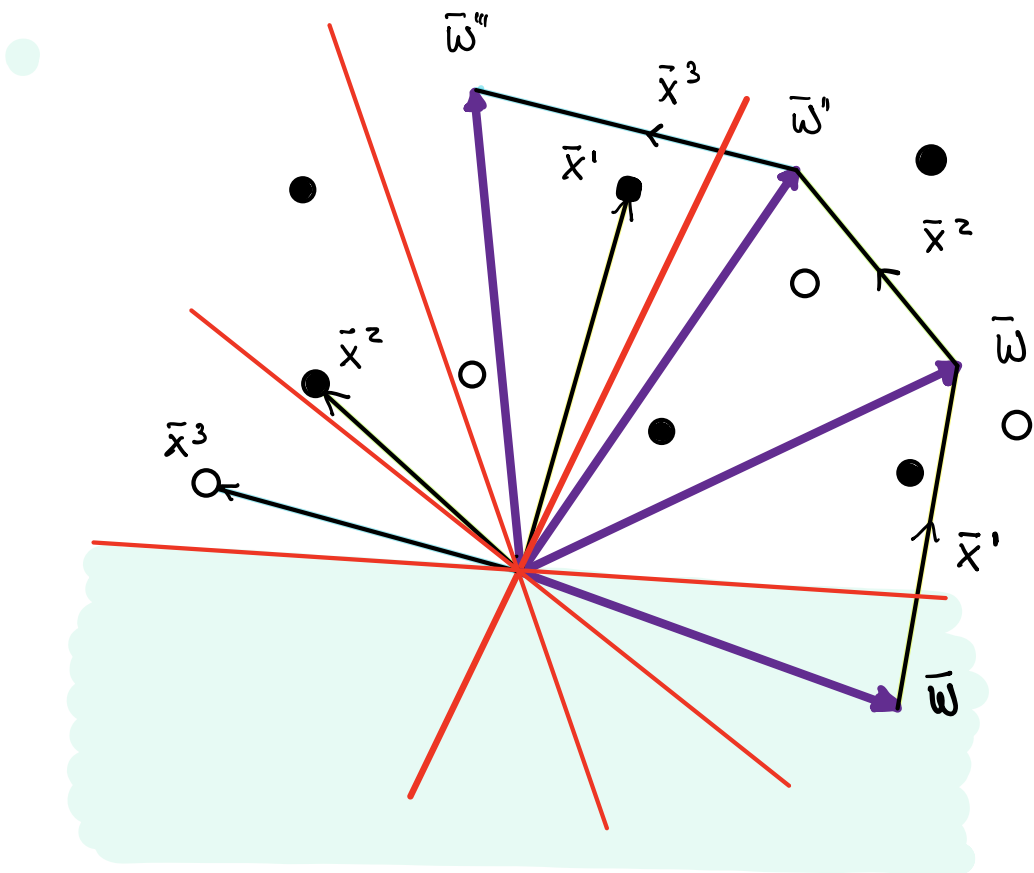








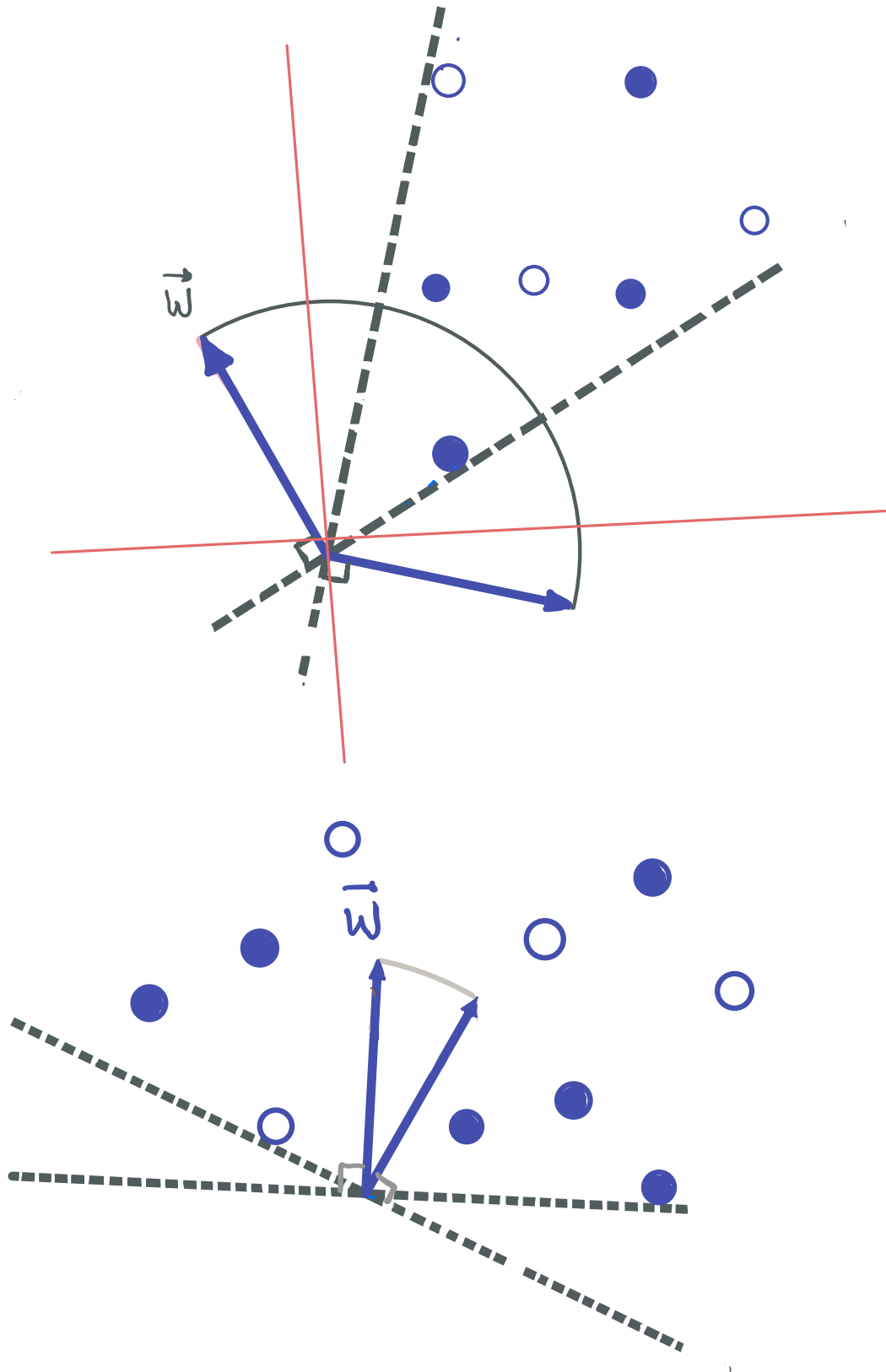
Miremos la evolución en el tiempo de la solución. Aquí llegamos en tres pasos pero para problemas reales necesitaremos millones de iteraciones.



Notemos que una actualización del vector \vec{w} implica mostrarle todos los elementos del conjunto de entrenamiento a la red. A este número de iteraciones lo llamamos una **ÉPOCA**.

Salvo casos muy patológicos, no existe en general un único vector \vec{w} que resuelve el desafío, sino infinitos.

Miremos dos buenos ejemplos artificiales:



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

para cada posible \vec{w} . $D(\vec{w})$ es la mínima proyección de \vec{w} sobre los \vec{x}^{μ} del conjunto de entrenamiento. Como está dividido por $|\vec{w}|$ solo depende de la dirección de \vec{w} pero no de su magnitud.

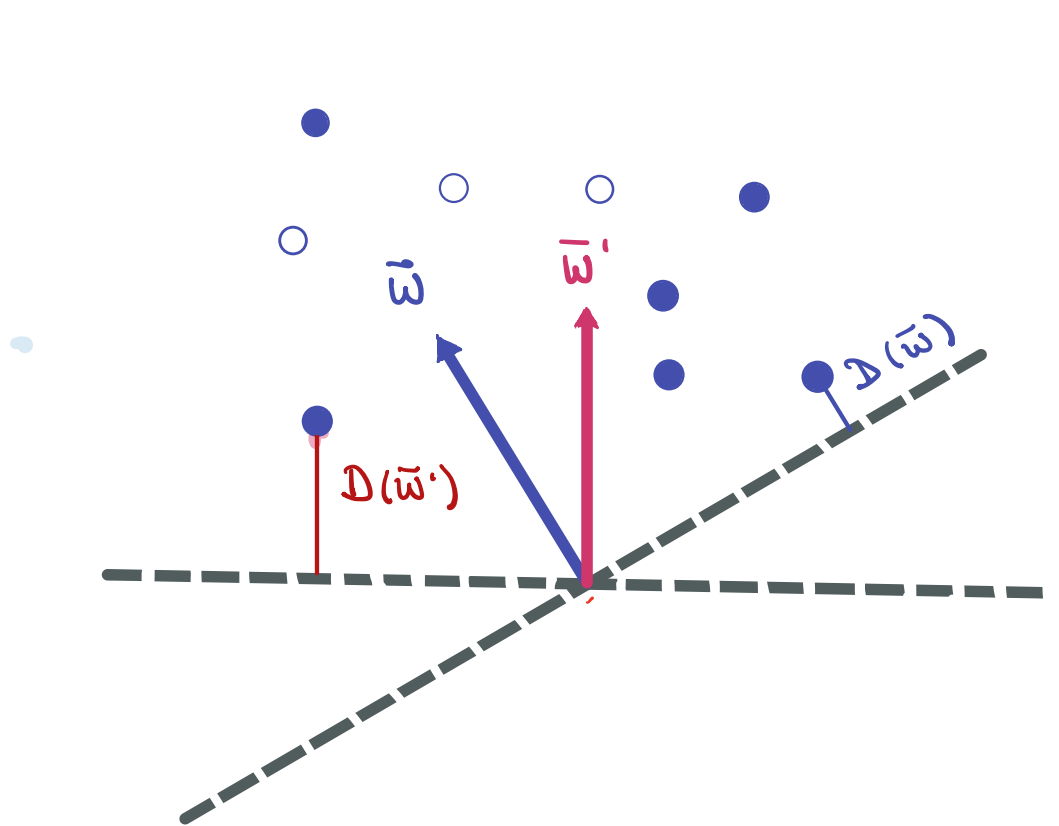
Para un dado valor de K , W puede variar en magnitud y en dirección. Podemos definir la cantidad:

$$\begin{aligned} D(\vec{w}) &= \min_{\mu} \frac{1}{|\vec{w}|} \cdot |\vec{w}| \cdot |\vec{x}^{\mu}| \cos(\theta_{\mu}) \\ &= \min_{\mu} |\vec{x}^{\mu}| \cos(\theta_{\mu}) \end{aligned}$$

Si $D(\vec{w})$, que es un número real (no un vector), es positivo, entonces todos los vectores \vec{x}^{μ} están en el lado correcto. Si es negativo, hay al menos un \vec{x}^{μ} para el cual no funciona el método.

Supongamos que todos los elementos del conjunto de entrenamiento están del lado correcto. Podemos definir el mejor D :

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



Pero no nos conformamos con esto. Para que la solución que buscamos sea más robusta (expliquemos esto) vamos a pedir que

$$\sum^N h^N > N \kappa \quad \text{kappa}$$

$$\Delta w_k = \eta \Theta (N\alpha - \sum^N h^N) \sum^N \xi_k^N$$

donde:

$$\Theta(z) = \begin{cases} 1 & z \geq 0 \\ 0 & z < 0 \end{cases}$$

Esta regla se llama **REGLA DEL APRENDIZAJE DEL PERCEPTRÓN** y fue publicada por Frank Rosenblatt en 1962.

Se puede demostrar rigurosamente que si nuestro problema por aprender es **LINEALMENTE SEPARABLE**, la regla converge a una solución \bar{w} en un número finito de pasos.

