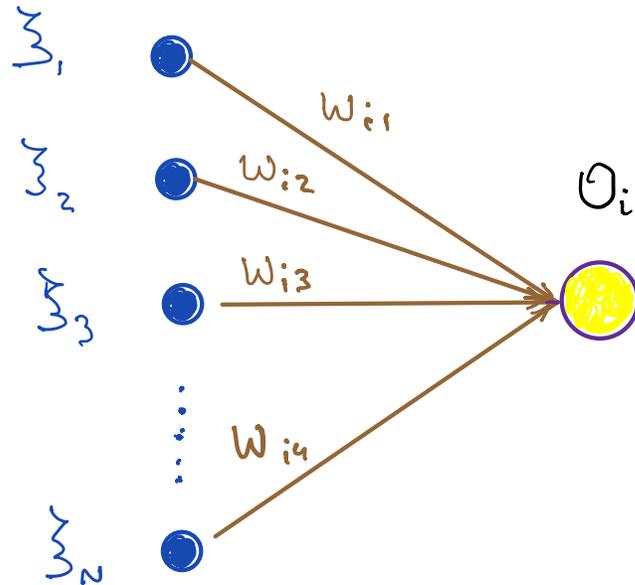


PERCEPTRON SIMPLE CON UNIDADES LINEALES

Volvamos a analizar un perceptron de una única capa representado en M perceptrones simples. La diferencia con el caso anterior, donde tenemos neuronas binarias de salida, es que ahora la neurona de salida es lineal



N neuronas
de entrada

1 neurona i
de salida

Los elementos del vector $\vec{\xi} = (\xi_1, \xi_2, \dots, \xi_N)$ son números reales. La salida es

$$O_i = h_i = \sum_k w_{ik} \xi_k$$

Notem que la función de activación es

$$g(h) = h$$

Supongamos nuevamente que tenemos P entredos $j=1, 2, \dots, P$ etiquetados, o sea, para los cuales conocemos la salida

Si entra \vec{X}^M debe salir \vec{Y}_i^M .

Separamos lo que sale, que es O_i^M

$$O_i^M = h_i = \sum_k W_{ik} \sum_k^M$$

¿ lo que debe salir, \vec{Y}_i^M .

Nota: los valores deseados pueden ser discretos
DISCUTIR

Este problema tiene una solución explícita, llamada PSEUDO-INVERSA. En este caso calculamos la matriz

$$Q_{\mu\nu} = \frac{1}{N} \sum_k \sum_k^M \sum_k^V$$

que define la similitud entre los elementos μ y ν del conjunto de entrenamientos.

Si Q es invertible, podemos definir los acoplamiento como

$$W_{ik} = \frac{1}{N} \sum_{\mu} \sum_{\nu} \xi_{i}^{\mu} (Q^{-1})_{\mu\nu} \xi_{\nu}^k$$

independientes

Para esto los p vectores ξ^{μ} deben ser linealmente
Si no lo son, no existe solución. Este es una
condición que aparece en todas las redes neuronales
lineales. Esto los hace poco interesantes.

Para que sean linealmente independientes debe
ser p menor o igual a N .

$$p \leq N$$

El método de descenso por el gradiente

Supongamos que mostramos a la red los
 p inputs y para cada uno de ellos calculamos

$$E(\vec{w}) = \frac{1}{2} \sum_{\mu} \sum_i (\xi_i^{\mu} - O_i^{\mu})^2$$

Noten que si $O_i^{\mu} = \xi_i^{\mu}$ (si la red resuelve bien

un ejemplo del conjunto de entrenamiento, ese ejemplo no contribuye a la suma, pues

$$\sum_i x_i = 0_i^M$$

↙
resultado
deseado

↘
resultado
real.

Esta función $E(\bar{w})$ se llama **función costo**. Es función de los N elementos de \bar{w} y no es función de los y_i^M , que serán considerados parámetros dados (fijos). Es un concepto que nos remite a la idea de función energía de los modelos de Hopfield, pero no es exactamente lo mismo. **DISCUTIR**

Por construcción $E(\bar{w}) \geq 0$ (es suma de NP términos mayores o iguales a cero).

Si $E(\bar{w}) = 0$, entonces la red resolvió exactamente el problema, pues significa que $y_i^M = 0_i^M$ para todo $M=1, 2, \dots, P$.

Por otro lado, cuanto menor es $E(\bar{w})$, menor será el error o el costo, pues más se parecerán los J_i a los D_i .

Nota: Dada una función

$$f: \mathbb{R}^n \rightarrow \mathbb{R} \quad (\text{varias variables})$$

su gradiente en un punto $\bar{x} \in D(f)$ se define como un campo vectorial, o sea, un vector en \mathbb{R}^n asignado a cada punto del dominio el cual pertenece a \mathbb{R}^n . Se lo denota

∇f y se define como:

$$\nabla f(\bar{x}) = \left(\frac{\partial f(\bar{x})}{\partial x_1}, \frac{\partial f(\bar{x})}{\partial x_2}, \dots, \frac{\partial f(\bar{x})}{\partial x_n} \right)$$

O sea, a cada vector \bar{x} asignamos un vector ∇f (esto se llama campo vectorial)

Generaliza el concepto de derivada en el caso de funciones de una variable $f: \mathbb{R} \rightarrow \mathbb{R}$ y geométricamente se puede ver que ∇f apunta en la dirección en la cual f

más rápidamente en el punto \bar{x} .

Volvamos a nuestro problema: el perceptron simple lineal.

$E(\bar{w}) = E(w_1, w_2, \dots, w_N)$ es una función

$$E: \mathbb{R}^n \longrightarrow \mathbb{R} \quad (\text{campo escalar})$$

Como es bien comportada podemos definir su gradiente en el vector \bar{w}

$$\nabla E = \left(\frac{\partial E}{\partial w_1}, \frac{\partial E}{\partial w_2}, \dots, \frac{\partial E}{\partial w_N} \right)$$

Este vector ∇E nos indica hacia donde crece más rápido E al variar \bar{w} . Entonces

$$-\nabla E$$

indica la dirección hacia donde decrece más rápido $E(w)$. Como buscamos un \bar{w}^{optimo} para el cual $E(\bar{w}^{\text{optimo}}) = 0$ construiremos la siguiente regla de aprendizaje

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

$$\text{con } \Delta W_{ik} = -\eta \frac{\partial E}{\partial W_{ik}} \quad \vec{w}^{\text{nuevo}} = \vec{w}^{\text{viejo}} + \vec{\Delta w}$$

O sea, ΔW_{ik} apunta hacia la dirección de más rápido descenso de $E(\vec{w})$.

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

Hagamos algunos cálculos

$$\begin{aligned} \frac{\partial E(\vec{w})}{\partial W_{ik}} &= \frac{\partial}{\partial W_{ik}} \left(\frac{1}{2} \sum_v \sum_j (y_j^v - o_j^v)^2 \right) \\ &= \frac{\partial}{\partial W_{ik}} \left(\frac{1}{2} \sum_v \sum_j (y_j^v - \sum_e W_{je} x_e^v)^2 \right) \\ &= -\frac{1}{2} \sum_v (y_i^v - o_i^v) \sum_k x_k^v \quad \text{batch} \end{aligned}$$

Entonces

$$\Delta W_{ik} = \eta (y_i^M - o_i^M) \sum_k x_k^M \xrightarrow{\text{online}}$$

$$\text{con } \delta_i^M = (y_i^M - o_i^M)$$

Esta regla se conoce como

- Regla Delta (Delta rule)
- Regla ADALINE
- Regla Widrow-Hoff
- Regla LMS (Least mean square)

ADALINE viene de ADAPtative Linear Neuron y fue definida por primera vez en 1960 por Bernard Widrow y su estudiante de doctorado Ted Hoff en Stanford University.

El parámetro η se llama "razón de aprendizaje" y ya había aparecido en el perceptrón binario. Veremos que es "FUNDAMENTAL" y que nos acompañará siempre. Digamos por ahora que debe ser "pequeño" sin abundar mucho más.

Noten que es una regla hebbiana, derivada de un gradiente.

En el caso bueno el equivalente a λ_i era ± 2 o 0 . Ahora toma un continuo de valores.

E es una función de \bar{w} y tiene un mínimo $E=0$.

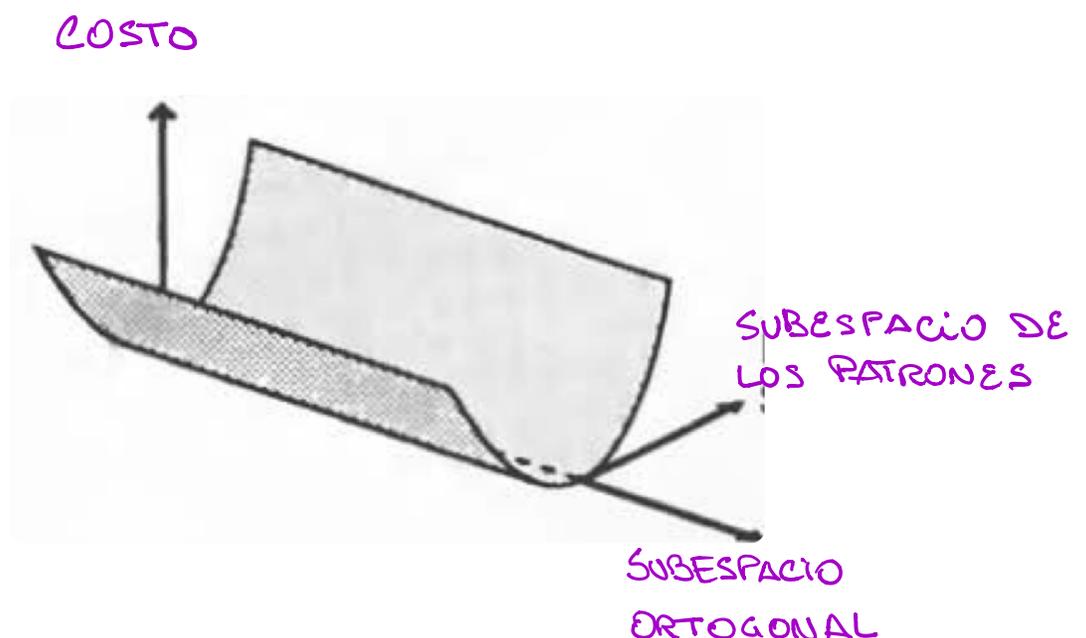
Supongamos que existe solución, entonces $p \leq N$.

Luego los p vectores definen un subespacio.

En el subespacio ortogonal al subespacio expandido

la función debe valer también cero, pues no depende de esas dimensiones. El siguiente es

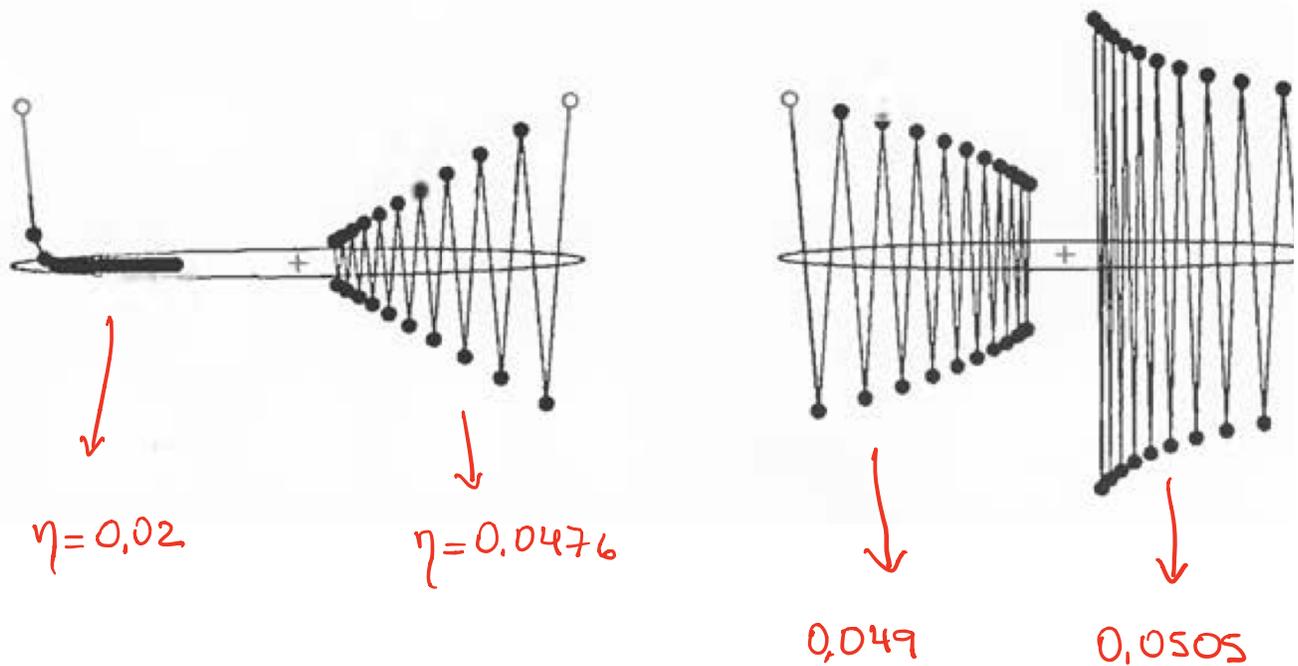
un esquema de la forma de $E(\bar{w})$ en \mathbb{R}^N .



Tenemos entonces una estructura tipo **concha**.

Esto nos dice que hay muchas soluciones, que viven en el fondo de la concha. En el subespacio expandido por los patrones de entrenamiento la regla ADALINE deuce la función costo. Por eso el valor de η se vuelve crucial.

No vamos a demostrar la convergencia



Noten cuán sensible es el método al valor de η .

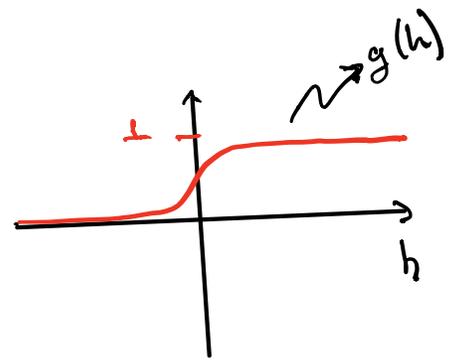
PERCEPTRON SIMPLE CON UNIDADES NO LINEALES

Ahora hacemos lo mismo cuenta en el caso en el cual la neurona de salida es no lineal

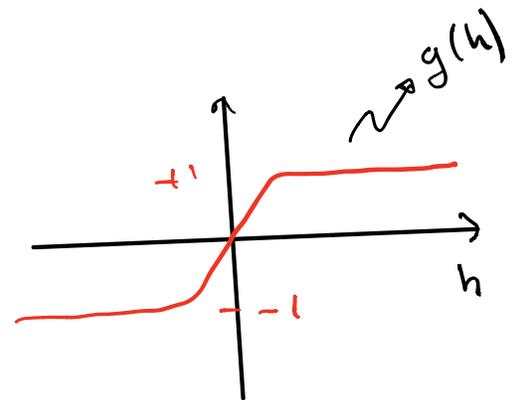
$$O_i = g(h_i)$$

Dos casos típicos son

$$g(h) = \frac{1}{1 + e^{-\beta h^2}}$$



$$g(h) = \tanh(\beta h)$$



Ahora la cuenta es muy sencilla.

$$E(\bar{w}) = \frac{1}{2} \sum_{\mu_i} (Y_i^\mu - O_i^\mu)^2$$

$$E(\bar{w}) = \frac{1}{2} \sum_{\mu_i} \left[y_i^\mu - g \left(\sum_k w_{ik} z_k^\mu \right) \right]^2$$

$$\frac{\partial E}{\partial w_{ik}} = - \sum \left[y_i^\mu - g(h_i^\mu) \right] g'(h_i^\mu) z_k^\mu$$

Definimos $\delta_i^\mu = (y_i^\mu - o_i^\mu) g'(h_i^\mu)$

de forma que

$$\begin{aligned} w_{ik}^{\text{nuevo}} &= w_{ik}^{\text{viejo}} + \Delta w_{ik} \\ &= w_{ik}^{\text{viejo}} - \eta \frac{\partial E}{\partial w_{ik}} \\ &= w_{ik}^{\text{viejo}} + \eta (y_i^\mu - o_i^\mu) g'(h_i^\mu) z_k^\mu \end{aligned}$$

$$w_{ik}^{\text{nuevo}} = w_{ik}^{\text{viejo}} + \eta \delta_i^\mu z_k^\mu$$

Un problema es la necesidad de evaluar $g(h_i^*)$ que a diferencia del caso lineal tiene un alto costo en tiempo de cálculo. Además hay que evaluar $g'(h_i^*)$.

$$\text{Si } g(h) = \tanh(\beta h), \text{ entonces } g'(h) = \beta(1-g^2)$$

$$\text{Si } g(h) = \frac{1}{1 + e^{-2\beta h}}, \text{ entonces } g'(h) = 2\beta g(1-g)$$

El hecho de que g' sea función de g facilita las cosas.

Vemos más adelante que estas expresiones para g y g' son responsables de la dificultad de construir redes de muchos capas (o sea, redes profundas).

Es importante saber que los neurones no lineales de salida se usan incluso para problemas binarios. **DISCUTIR**. Esto permite pensar en términos de descenso por el gradiente.

La cantidad de problemas que se pueden resolver con redes de una capa es pequeña.

Las condiciones para la existencia de solución son los mismos. Esto es fácil de verificar si se reemplaza la salida deseada ζ_i^M por $g^{-1}(\zeta_i^M)$. Lo que cambia es el "paseo" de la función costo $E(\bar{w})$. Pueden existir mínimos locales.

OTRAS FUNCIONES DE COSTO

$$E(\bar{w}) = \sum_i \sum_{\mu} \left[\frac{1}{2} (1 + \zeta_i^{\mu}) \log \left(\frac{1 + \zeta_i^{\mu}}{1 + O_i^{\mu}} \right) + \frac{1}{2} (1 - \zeta_i^{\mu}) \log \left(\frac{1 - \zeta_i^{\mu}}{1 - O_i^{\mu}} \right) \right]$$

Sara Solla, 1988 ENTROPÍA RELATIVA

Si $\zeta_i^{\mu} = O_i^{\mu}$ para todo μ y todo i , $E(\bar{w})$ es cero.

Por otro lado, con los $g(h)$ usados ($-1 < g(h) < +1$)

es una función definida positiva.

Ejercicio: mostrar que $\Delta W_{ik} = \eta \delta_i^{\mu} \zeta_k^{\mu}$ pero

con
$$\delta_i^{\mu} = \beta (\zeta_i^{\mu} - O_i^{\mu})$$

Esto nos ahorra tener que calcular g'

NEURONAS ESTOCÁSTICAS

Ahora que tenemos una función $E(\bar{w})$ que minimizar la cual tiene un mínimo global en $E(\bar{w})=0$ pero puede tener muchos mínimos locales, podemos introducir neuronas aleatorias en la red, como hicimos con la red de Hopfield.

Para ello calculamos $h_i^M = \sum_k W_{ik} \xi_i^k$, y con h_i^M definimos una probabilidad

$$\text{Prob}(S_i^M = +1) = \frac{1}{1 + e^{-2\beta h_i^M}}$$

$$\text{Prob}(S_i^M = -1) = 1 - \text{Prob}(S_i^M = +1) = \frac{1}{1 + e^{2\beta h_i^M}}$$

Nuevamente tenemos que

$$\langle S_i^M \rangle = (+1) \cdot \text{Prob}(S_i^M = +1) + (-1) \cdot \text{Prob}(S_i^M = -1)$$

$$= \frac{1}{1 + e^{-2\beta h_i^M}} - \frac{1}{1 + e^{2\beta h_i^M}}$$

$$= \frac{1}{e^{-\beta h_i^M} (e^{\beta h_i^M} + e^{-\beta h_i^M})} - \frac{1}{e^{\beta h_i^M} (e^{-\beta h_i^M} + e^{\beta h_i^M})}$$

$$= \frac{(e^{\beta h_i^M} - e^{-\beta h_i^M})}{(e^{\beta h_i^M} + e^{-\beta h_i^M})}$$

$$= \tanh(\beta h_i^M)$$

o

$$\langle S_i^M \rangle = \tanh\left(\beta \sum_k w_{ik} S_k^M\right)$$

Podemos calcular $\langle S_i^M \rangle$ en una simulación y el valor estimado lo podemos usar en la actualización de los pesos sinápticos

$$\Delta w_{ik} = \eta \delta_i^M S_k^M$$

$$\delta_i^M = (S_i^M - \langle S_i^M \rangle)$$

Se puede mostrar (está en el libro) que esta regla decrece el promedio de la función error

$$E = \frac{1}{2} \sum_i \sum_{\mu} (S_i^{\mu} - \langle S_i^{\mu} \rangle)^2$$