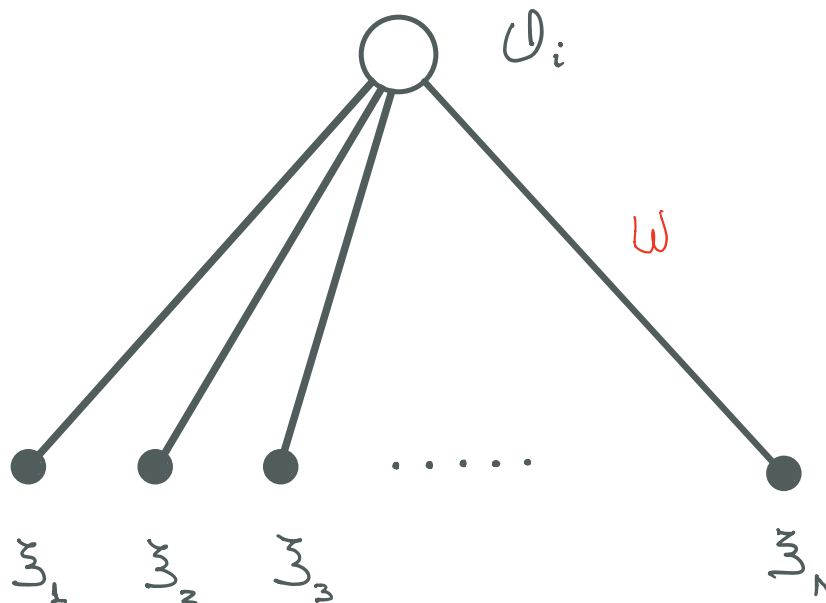
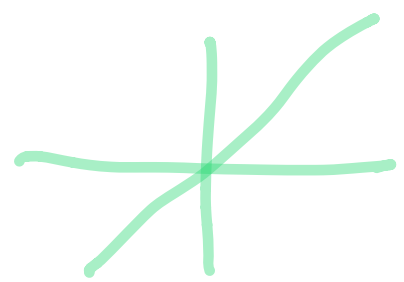


## PERCEPTRON SIMPLE CON NEURONA LINEAL

Cuando describí las posibles funciones de activación  $g(z)$ , las cuales debemos recordar, definí una neurona de salida, no contemplé una función lineal como posibilidad, pero ahora analizaremos el caso del perceptrón simple con una neurona de salida lineal. Recuerden que si la capa de salida tiene más de una neurona, estas siempre se pueden desacoplar.



$$g(z) = z$$



$$\underline{U_i} = \underline{g(h_i)} = h_i = \sum_k w_{ik} \sum_k$$

Como antes, suponemos que tenemos  $p$  entradas previamente etiquetadas

$$\sum_k^H \longmapsto \sum_i^H \quad (\text{valor de salida})$$

①

$$\sum_i^H = \sum_k w_{ik} \sum_k^H \quad (\text{valor deseado})$$

Nota: la salida puede tomar cualquier valor real, por lo cual no estamos haciendo una clasificación sino una REGRESIÓN

Existe una solución explícita llamada **PSEUDO INVERSA** que nos da los valores de los componentes del vector  $\vec{w}_i = (w_{i1}, w_{i2}, \dots, w_{in})$  y tiene la forma

$$w_{ik} = \frac{1}{N} \sum_{\mu\nu} \xi_i^\mu (Q^{-1})_{\mu\nu} \xi_k^\nu$$

$$Q = \begin{pmatrix} q_{11} & q_{1p} \\ q_{p1} & q_{pp} \end{pmatrix} = \frac{1}{N} (\xi_i^1, \dots, \xi_i^p) Q^{-1} \begin{pmatrix} \xi_i^1 \\ \vdots \\ \xi_i^p \end{pmatrix}$$

donde

$$Q_{\mu\nu} = \frac{1}{N} \sum_{k=1}^N \xi_k^\mu \xi_k^\nu = \frac{1}{N} \sum_k \xi_k^\mu \cdot \xi_k^\nu$$

es la matriz de superposiciones entre los elementos del conjunto de entrenamientos.

Se puede probar que si usamos esta expresión y la reemplazamos en (1) se cumple la condición.

Para que exista la solución, debe existir  $Q^{-1}$ , y para esto los vectores  $\xi_k^\mu$  ( $\mu=1, 2, \dots, p$ ) deben ser

linealmente independientes. O sea, si existe una relación lineal de la forma

$$a_1 \sum_k^1 + a_2 \sum_k^2 + \dots + a_p \sum_k^p = 0 \quad \text{para todo } k \quad (k=1, 2, \dots, 10)$$

el problema no tiene solución.

## Descenso por el gradiente

Hoy presentemos este método por primera vez pero nos acompañara **SIEMPRE** en los modelos neuronales de aprendizaje automático.

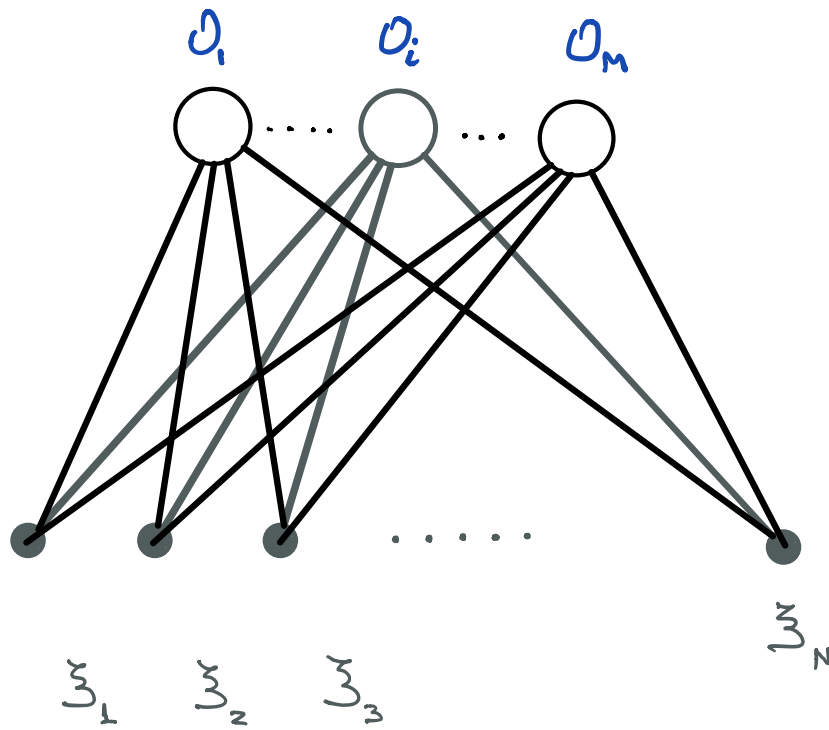
Definimos una función error o función costo de la forma

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

suma sobre  
todas las neuronas de  
la capa de salida

suma sobre los  
elementos del conjunto  
de entrenamiento





Cuando le presentamos la entrada  $\mu$  obtendremos el resultado dado por la regla del perceptron lineal

$$O_i^k = \sum_k w_{ik} z_k^k \quad g(z) = z$$

La idea es comenzar con un vector inicial arbitrario  $\vec{w}$  y presentarle el ejemplo  $\mu$

$$E : \mathbb{R}^N \rightarrow \mathbb{R}$$

$$\nabla E : \mathbb{R}^N \rightarrow \mathbb{R}^N \quad \nabla E_{ik} = \frac{\partial E}{\partial w_{ik}}$$

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

$$\nabla E = \left( \frac{\partial E}{\partial w_{i1}}, \frac{\partial E}{\partial w_{i2}}, \dots, \frac{\partial E}{\partial w_{in}} \right)$$

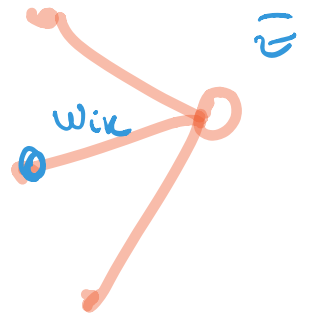
$$\Delta w_{ik} = -\eta \frac{\partial E}{\partial w_{ik}}$$

$$= -\eta \frac{\partial}{\partial w_{ik}} \left( \frac{1}{2} \sum_{\mu=1}^P (\sum_i^M - \underline{O_i^M})^2 \right)$$

$$= -\eta \sum_{\mu=1}^P \frac{1}{2} \frac{\partial}{\partial w_{ik}} \left( \sum_i^M - \sum_j w_{ij} \sum_j^M \right)^2$$

$$= -\eta \sum_{\mu=1}^P \frac{2}{2} \left( \sum_i^M - \sum_j w_{ij} \sum_j^M \right) \frac{\partial}{\partial w_{ik}} \left( - \sum_j w_{ij} \sum_j^M \right)$$

$$= \eta \sum_{\mu=1}^P \left( \sum_i^M - O_i^M \right) \sum_k^M$$



$$\Delta w_{ik} = \eta \sum_{\mu=1}^P (\sum_i^M - O_i^M) \sum_k^M$$

$$\Delta \bar{w}_i = \eta \sum_{\mu=1}^P (\sum_i^M - O_i^M) \sum_k^M$$

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

$$\bar{w}_i^{\text{nuevo}} = \bar{w}_i^{\text{anterior}} + \Delta \bar{w}_i$$

Si definimos

$$\delta_i^M = \sum_i^M - O_i^M \quad (\text{delta } i^M)$$

entonces

$$W_i = W^{\text{anterior}} + \eta \sum \delta_i^M \sum^M$$

$$W_{ik}^{\text{nuevo}} = W_{ik}^{\text{anterior}} + \eta \sum_M \delta_i^M \sum_k^M$$

Esta regla de descenso por el gradiente se conoce como

- Regla delta o ADELPH (Widrow y Hoff) (1960)
- Regla best mean square (LMS) (1986)
- Regla Rescorla-Wagner (1972)

Tenemos  $M$  vectores sinápticos  $\vec{w}_i$  con  
 $i = 1, 2, \dots, M$

$$\vec{w}_1 = (w_{11}, w_{12}, \dots, w_{1N})$$

$$\vec{w}_2 = (w_{21}, w_{22}, \dots, w_{2N})$$

$\vdots$

$$\vec{w}_M = (w_{M1}, w_{M2}, \dots, w_{MN})$$

Entonces la función costo depende de  $M \cdot N$  coeficientes para un dado problema.

Recordemos que el problema es encontrar los  $M \cdot N$  coeficientes con los cuales el perceptron lineal resuelve bien todos los ejemplos del conjunto de entrenamiento

$$\left\{ \left( \vec{x}^\mu, y^\mu \right) / \mu = 1, 2, \dots, P \right\}$$

El gradiente de  $E(\{\bar{w}\})$  es un vector en  $\mathbb{R}^{M \times N}$

$$\nabla E = \overbrace{\left( \frac{\partial E}{\partial w_{11}} \dots \dots \dots \frac{\partial E}{\partial w_{MN}} \right)}^{M \times N}$$

← M · N elementos →

El gradiente es un vector que apunta en la dirección de mayor crecimiento de  $E(\{\bar{w}\})$  en  $\mathbb{R}^{M \times N}$ .

Comenzamos eligiendo los  $M \times N$  vectores  $\bar{w}_i$  ( $i=1, \dots, M$ ) al azar.

## APRENDIZAJE EN BATCH

- Sea  $\mu = 1, 2, 3, \dots, P$
- le presentamos  $\sum^M$  a la red y obtenemos los  $M$  valores  $O_i^\mu$ .
- con estos  $M$  valores calculo cada uno de los componentes del gradiente de  $E$  ( $M \times N$ )

- con los elementos del gradiente calculo las variaciones que aplicaremos a cada  $W_{ik}$

$$\underline{\Delta W_{ik}} = -\eta \frac{\partial E}{\partial W_{ik}} = \eta \sum_M \delta_i^M \sum_k^M$$

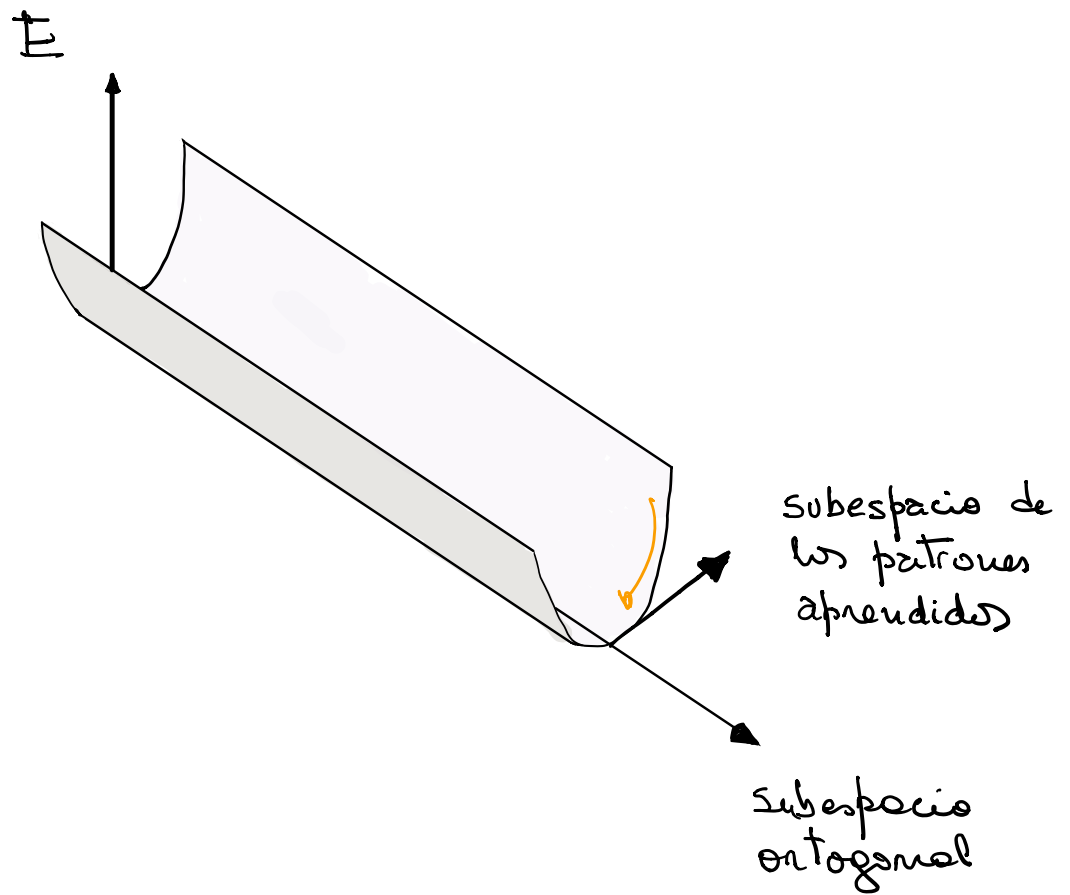
$$\delta_i^M = (z_i^M - O_i^M)$$

- con los  $\Delta W_{ik}$  actualizemos **TODOS** los acoplamiento sinápticos

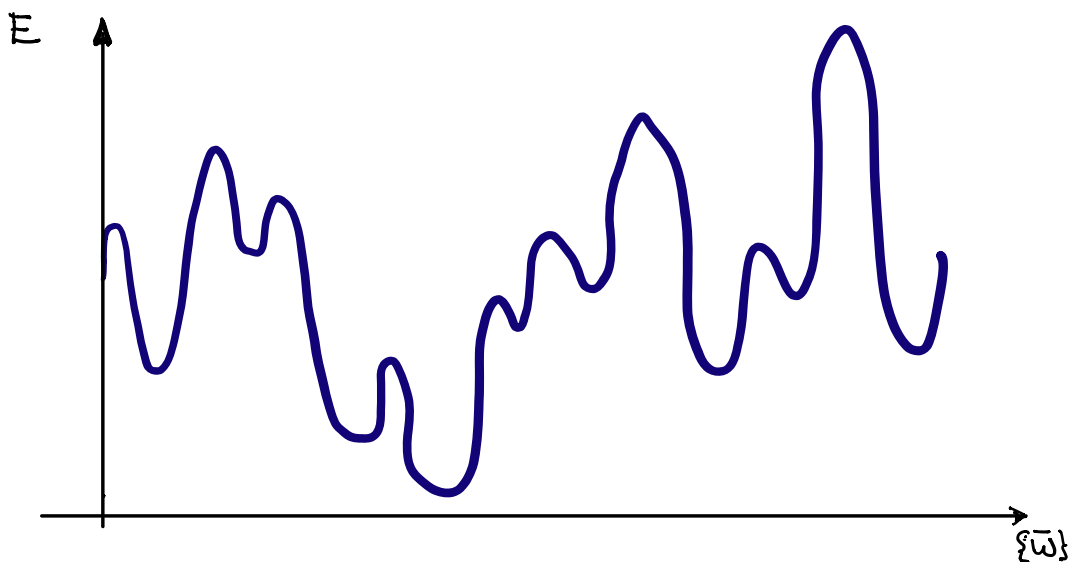
$$\begin{aligned} W_{ik}^{\text{nuevo}} &= W_{ik}^{\text{anterior}} + \Delta W_{ik} \\ &= W_{ik}^{\text{anterior}} + \eta \sum_M \delta_i^M \sum_k^M \end{aligned}$$

- empezamos de nuevo

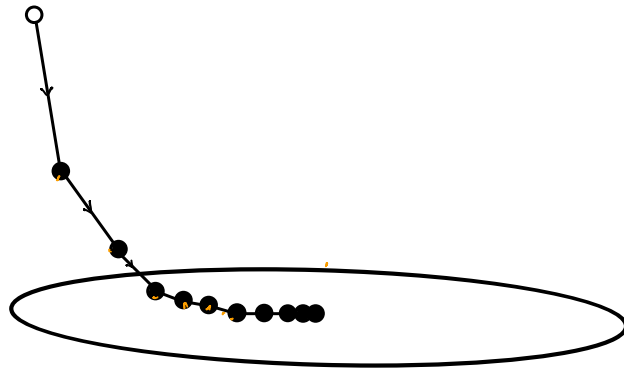
y volvemos al inicio a menos que  $\eta = P$ .



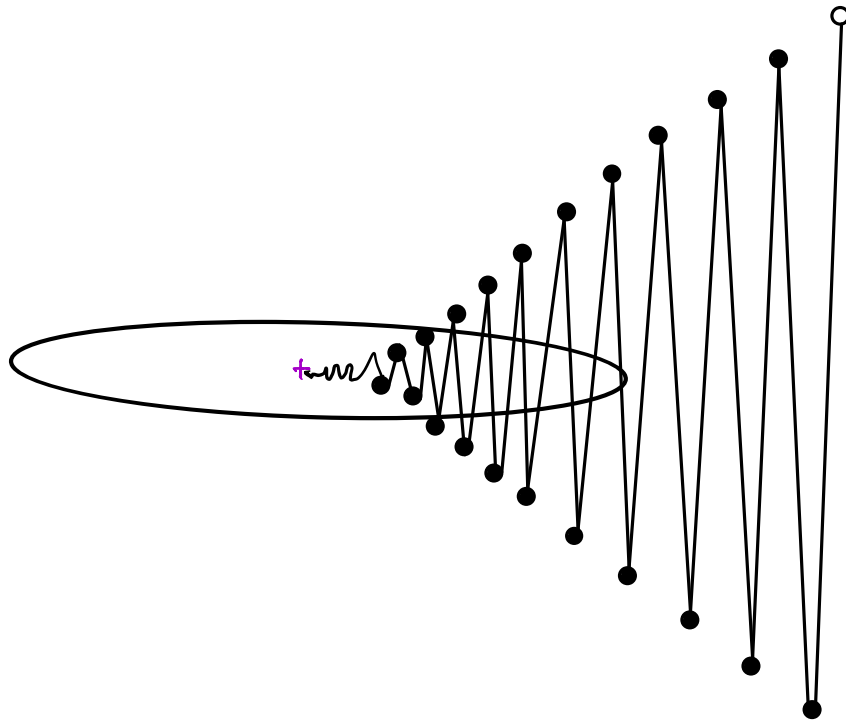
Ahora no vale la reproducibilidad lineal. Tenemos que intentar aprender descendiendo por la función error, pero  $E$  es una función muy rugosa



EJEMPLO :  $E = x^2 + 20y^2$

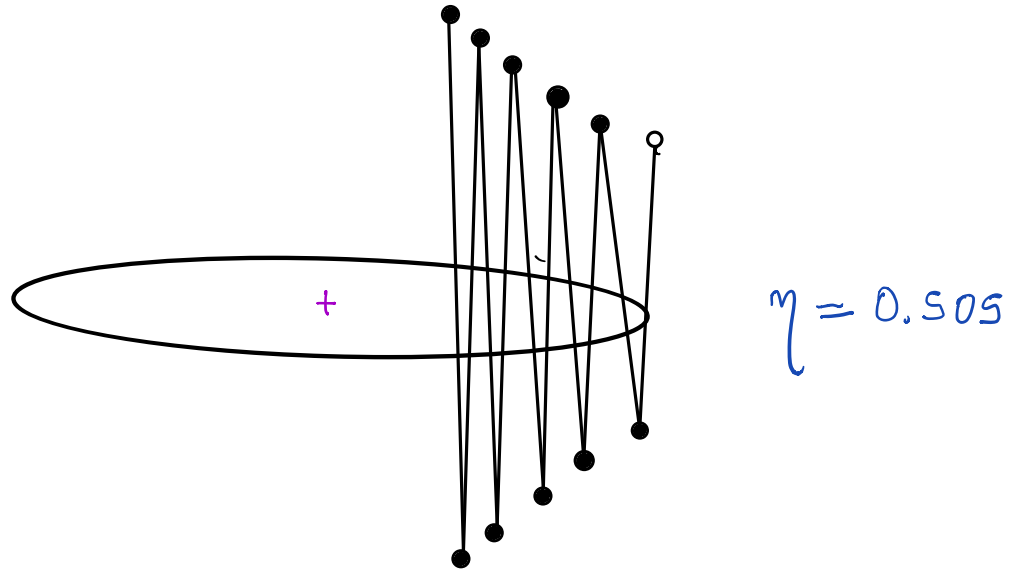


$\eta = 0.02$



$\eta = 0.476$





En el caso de neuronas de salida lineales, la forma cuadrática de la función  $E$  se puede diagonalizar si los ejemplos del conjunto de entrenamiento son linealmente independientes.

$$E = \sum_{\lambda=1}^{M \times N} \alpha_{\lambda} (w_{\lambda} - w_{\lambda}^0)^2$$

$w_{\lambda}$ : combinación lineal de los  $w_{ik}$ .

Por la forma cuadrática de  $E$  los autovalores  $\alpha_x$  son positivos o nulos. Los autovectores asociados a los autovalores de valor cero definen la parte de  $E$  independiente de  $\omega_x$ 's.

Si hacemos un descenso por el gradiente en la base diagonalizada

$$\Delta \omega_x = -\eta \frac{\partial E}{\partial \omega_x} = -2\eta \alpha_x (\omega_x - \omega_x^0)$$

Tenemos que llegar en el menor número de pasos a  $\omega_x^0$  (un número que surge de la diagonalización)

$$\delta \omega_x^{\text{opt}} = \omega_x - \omega_x^0 \quad \text{en dirección de } \lambda$$

$$\begin{aligned} \delta \omega_x^{\text{nuevo}} &= \delta \omega_x^{\text{old}} + \Delta \omega_x \\ &= \omega_x - \omega_x^0 - 2\eta \alpha_x (\omega_x - \omega_x^0) \\ &= (\omega_x - \omega_x^0) \cdot (1 - 2\eta \alpha_x) \\ &= \delta \omega_x^{\text{opt}} \cdot (1 - 2\eta \alpha_x) \end{aligned}$$

$\eta$  es tal que  $|1 - 2\eta \alpha_x| < 1$

El mayor  $\lambda$  define, a través de su autovector, la dirección de mayor curvatura de  $E$

$$\eta < \frac{1}{\lambda_{\max}}$$

Por otro lado, el menor autovector positivo define la dirección de más lenta aproximación al mínimo.

Así, si  $\lambda_{\max} / \lambda_{\min}$  es un número muy grande, estaremos en problemas y el método convergerá muy lentamente.

Note: \* no pedimos que el problema sea linealmente separable

\* no encontramos una solución perfecta

\* igual que con el problema de clasificación, estudiamos el error en  $n$  épocas. Solo que ahora la medida de error es la función

$E$

\* los computadores llaman a la

función eno como función de perdida  
o costo

