



IIC 2433 Minería de Datos

<https://github.com/marcelomendoza/IIC2433>

- CLASIFICACIÓN -

Clasificación

¿Cuáles de estos ejemplos son árboles?



Clasificación

input $\mathbf{x} \in \mathbb{R}^d = \mathcal{X}$.

output $y \in \{-1, +1\} = \mathcal{Y}$. (Clasificación binaria)

target function $f : \mathcal{X} \mapsto \mathcal{Y}$. $\leftarrow f$ desconocida

data set $\mathcal{D} = (\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)$. \leftarrow Escenario supervisado

Objetivo de aprendizaje:

Queremos **aprender** f usando \mathcal{D} .

Clasificación

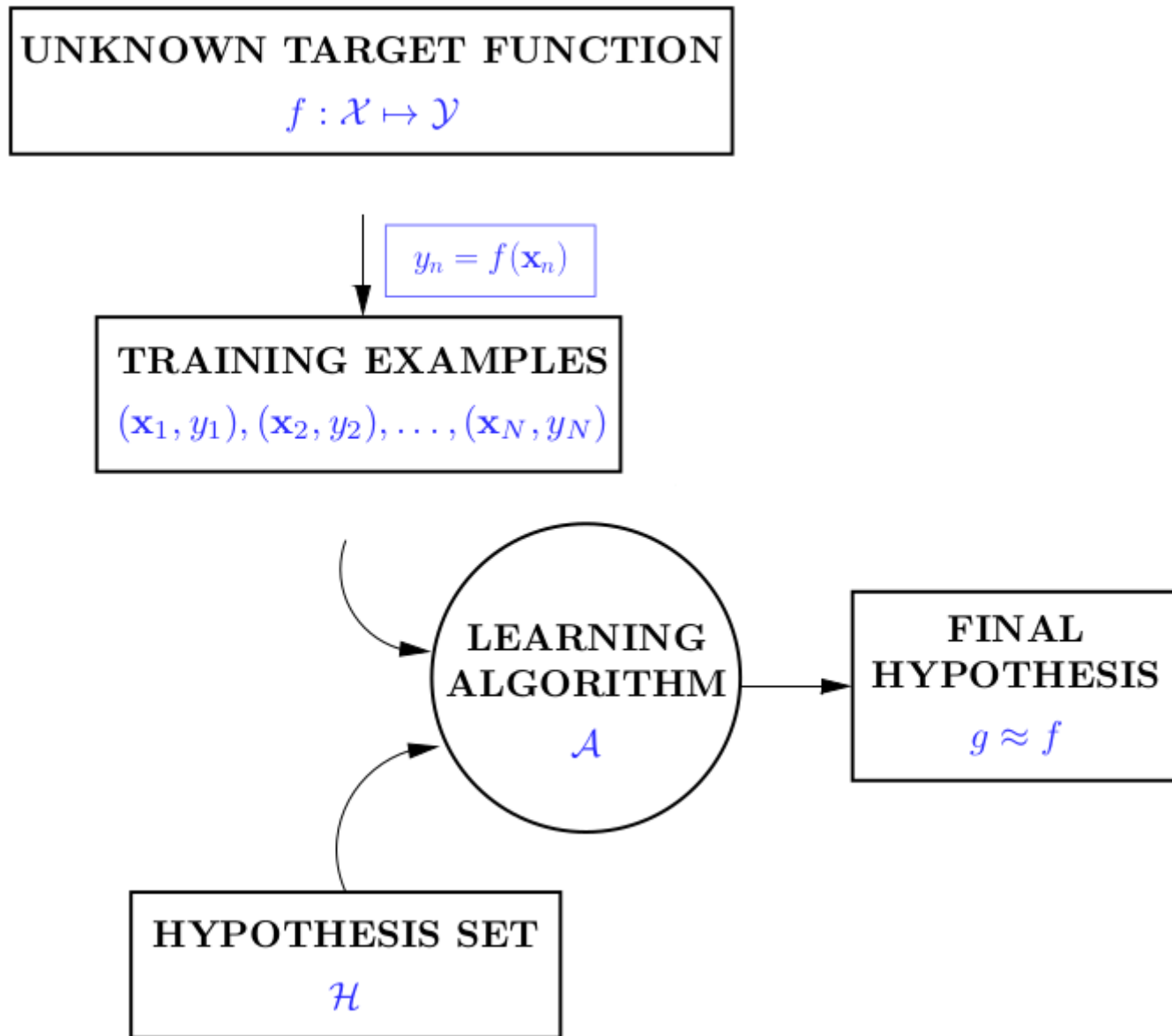
El algoritmo de aprendizaje

- Comenzamos con un conjunto de hipótesis candidatas \mathcal{H} que son factibles de representar f .

$$\mathcal{H} = \{h_1, h_2, \dots, \} \leftarrow \text{Conjunto de hipótesis}$$

- Un algoritmo A selecciona una hipótesis g desde \mathcal{H} . El algoritmo opera sobre datos etiquetados disponibles para seleccionar g (datos de entrenamiento).
- Para nuevos datos no etiquetados (i.e., y desconocido), usamos g para inferir y .

Clasificación



- PERCEPTRÓN LINEAL -

Un modelo simple (lineal)

age	32 years
gender	male
salary	40,000
debt	26,000
years in job	1 year
years at home	3 years
...	...

- Vector de características: $\mathbf{x} = [x_1, \dots, x_d]^T$.

- Aprenderemos distintos ‘pesos’ para las variables de entrada:

$$\text{“Credit Score”} = \sum_{i=1}^d w_i x_i.$$

¿Sujeto de crédito?

- Usaremos un **umbral** para decidir si aprobamos el crédito:

Aprobamos si: $\sum_{i=1}^d w_i x_i > \text{threshold}$, (score de crédito suficiente)

Rechazamos si: $\sum_{i=1}^d w_i x_i < \text{threshold}$. (score de crédito insuficiente)

Variable útil para mejorar el score: $w_i > 0$

Variable útil para disminuir el score: $w_i < 0$

Un modelo simple (lineal)

La hipótesis puede escribirse formalmente:

$$h(\mathbf{x}) = \text{sign} \left(\left(\sum_{i=1}^d w_i x_i \right) + w_0 \right)$$

Usamos el signo para tomar la
decisión

Bias (el umbral)

- El conjunto de hipótesis es:

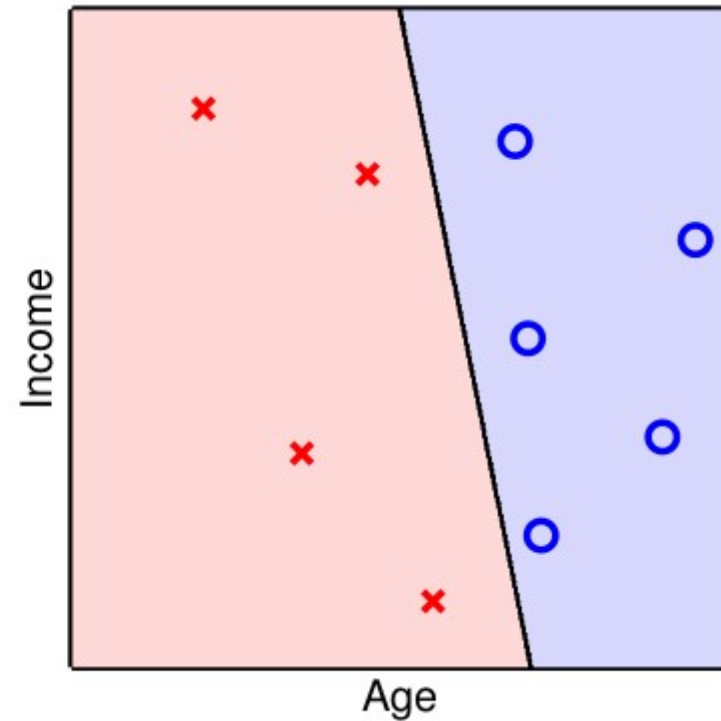
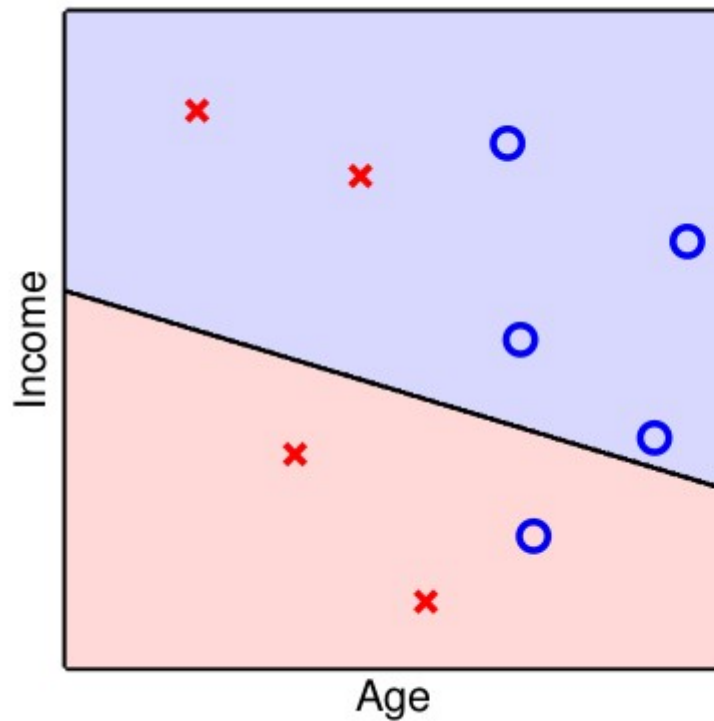
$$\mathcal{H} = \{h(\mathbf{x}) = \text{sign}(\mathbf{w}^T \mathbf{x})\} \quad \leftarrow \text{Infinito si los pesos están en } \mathbb{R}$$

- Este modelo se llama **Perceptrón**:

$$\mathbf{w} = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_d \end{bmatrix} \in \mathbb{R}^{d+1}, \quad \mathbf{x} = \begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_d \end{bmatrix} \in \{1\} \times \mathbb{R}^d.$$

El Perceptrón

- El Perceptrón usa los datos para encontrar una separación lineal según el atributo de clase y



- Como los datos están en un espacio de representación d -dimensional, el separador es un hiperplano.

El Perceptrón

Clasificación binaria: $\mathcal{Y} = \{+1, -1\}$

- 1: $\mathbf{w}(1) = \mathbf{0}$
- 2: **for** iteration $t = 1, 2, 3, \dots$
- 3: the weight vector is $\mathbf{w}(t)$. \leftarrow Random
- 4: From $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)$ pick any misclassified example.
- 5: Call the misclassified example (\mathbf{x}_*, y_*) ,

$$\text{sign}(\mathbf{w}(t) \cdot \mathbf{x}_*) \neq y_*.$$

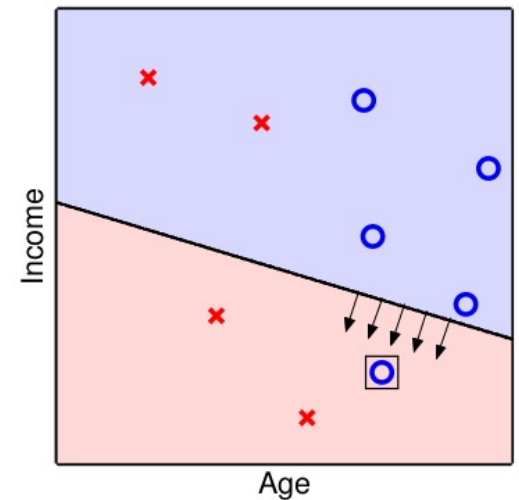
- 6: Update the weight:

$$\mathbf{w}(t+1) = \mathbf{w}(t) + y_* \mathbf{x}_*.$$

- 7: $t \leftarrow t + 1$



Update rule



Enfoque iterativo

El Perceptrón

Clasificación binaria: $\mathcal{Y} = \{+1, -1\}$

- Update rule (en ejemplos mal clasificados):

$$\mathbf{w}(t + 1) = \mathbf{w}(t) + y(t)\mathbf{x}(t).$$

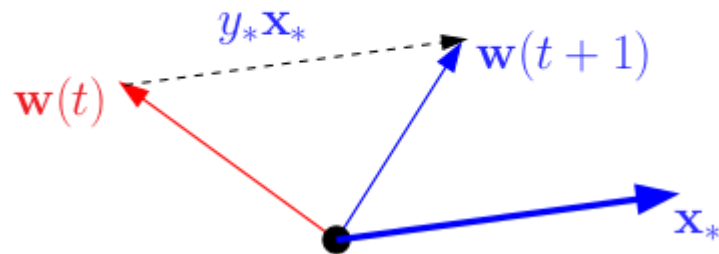
El Perceptrón

Clasificación binaria: $\mathcal{Y} = \{+1, -1\}$

- Update rule (en ejemplos mal clasificados):

$$\mathbf{w}(t+1) = \mathbf{w}(t) + y(t)\mathbf{x}(t).$$

Ground truth \longrightarrow $y_* = +1$



Orienta \mathbf{w} hacia la dirección de \mathbf{x}

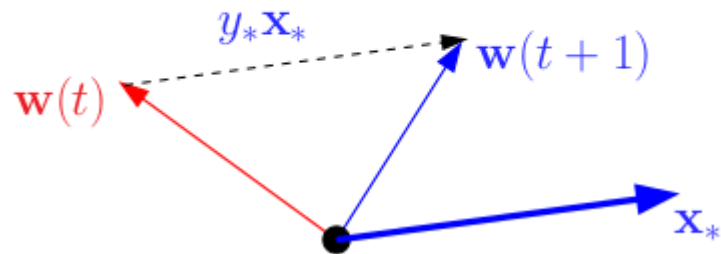
El Perceptrón

Clasificación binaria: $\mathcal{Y} = \{+1, -1\}$

- Update rule (en ejemplos mal clasificados):

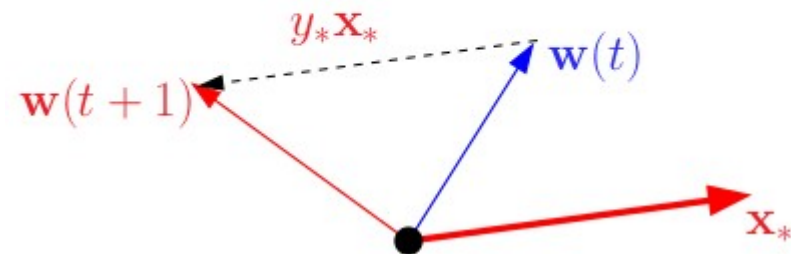
$$\mathbf{w}(t+1) = \mathbf{w}(t) + y(t)\mathbf{x}(t).$$

Ground truth \longrightarrow $y_* = +1$



Orienta \mathbf{w} hacia la dirección de \mathbf{x}

$y_* = -1$ \longleftarrow Ground truth



Orienta \mathbf{w} en la dirección contraria de \mathbf{x}

- REGRESIÓN LOGÍSTICA -

Modelos lineales

señal \rightarrow $\mathcal{S} = \mathbf{w}^T \mathbf{x}$

entrada \downarrow

\uparrow parámetros

The diagram shows the equation $\mathcal{S} = \mathbf{w}^T \mathbf{x}$. A blue arrow points from the word 'señal' to the variable \mathcal{S} . A blue arrow points from the word 'entrada' down to the variable \mathbf{x} . A red arrow points from the word 'parámetros' up to the variable \mathbf{w} . The variable \mathbf{w} is colored red, and the variable \mathbf{x} is colored blue.

Modelos lineales

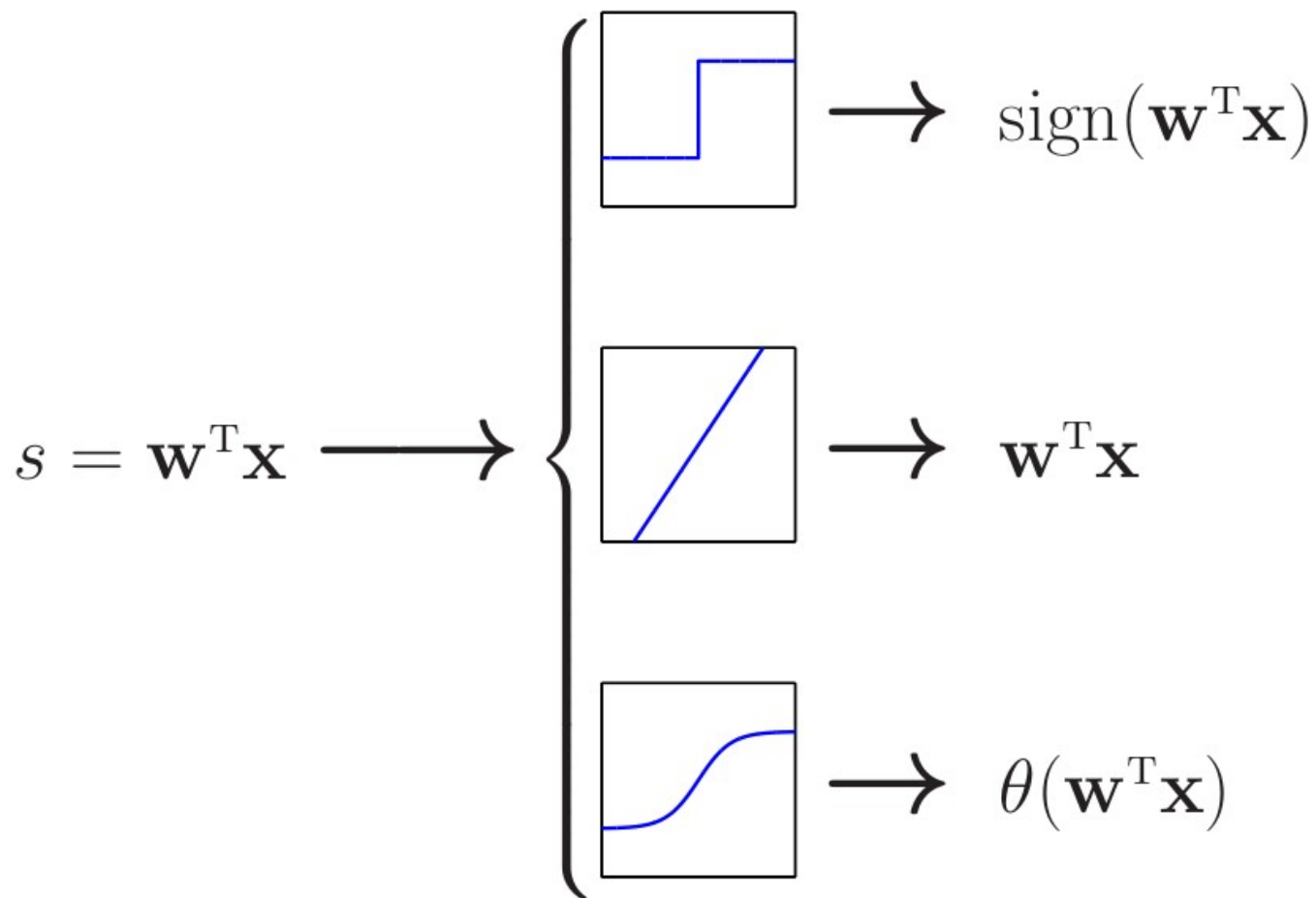
activación

$$y = \theta(s)$$

entrada

$$s = \mathbf{w}^T \mathbf{x}$$

parámetros



$$\{-1, +1\}$$

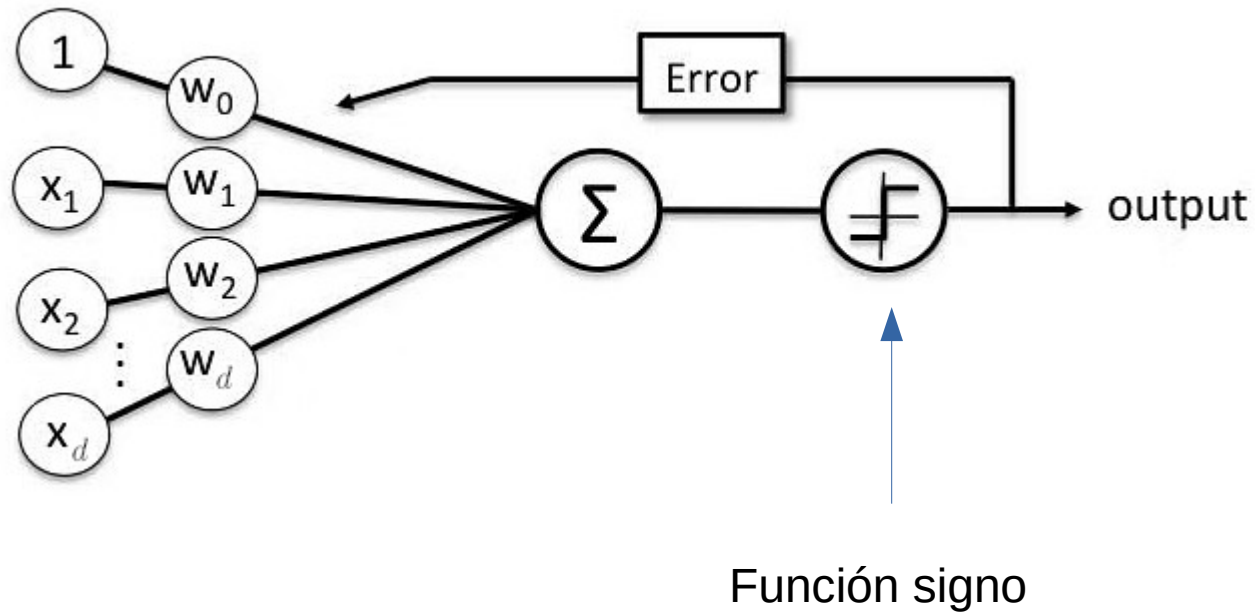
$$\mathbb{R}$$

$$[0, 1]$$

Perceptrón lineal

$$\mathcal{H}_{\text{lin}} = \{h(\mathbf{x}) = \text{sign}(\mathbf{w}^T \mathbf{x})\}$$

entrada
↓
 $\mathcal{S} = \mathbf{w}^T \mathbf{x}$
↑
parámetros



Regresión logística

Clase objetivo

Objetivo:

$$f(\mathbf{x}) = \mathbb{P}[y = +1 \mid \mathbf{x}].$$

Regresión logística

Clase objetivo

Objetivo:

$$f(\mathbf{x}) = \mathbb{P}[y = +1 \mid \mathbf{x}].$$

Modelo:

$$h(\mathbf{x}) = \theta \left(\sum_{i=0}^d w_i x_i \right) = \theta(\mathbf{w}^T \mathbf{x})$$

Regresión logística

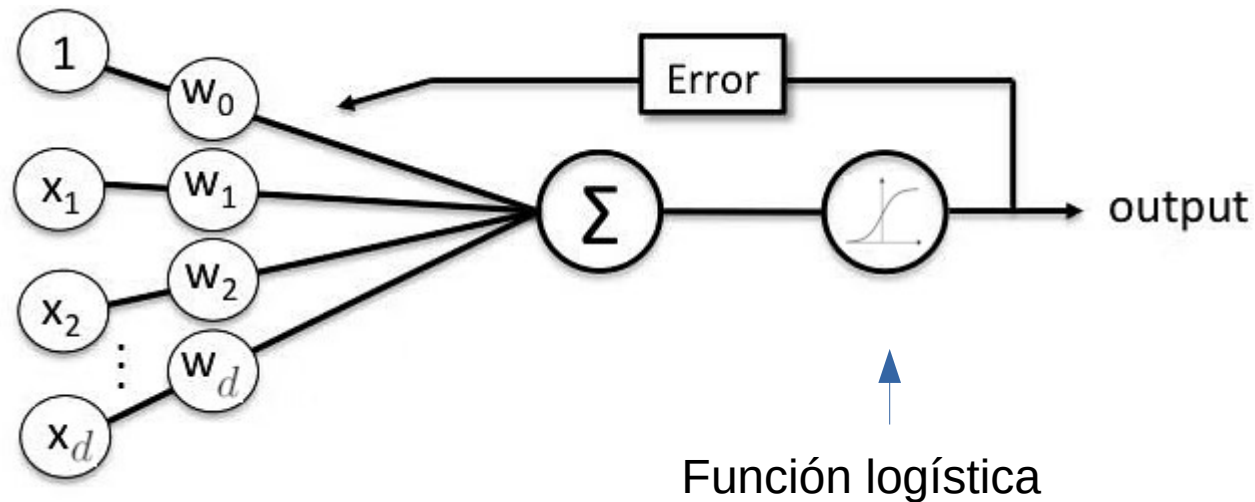
Clase objetivo

Objetivo:

$$f(\mathbf{x}) = \mathbb{P}[y = +1 \mid \mathbf{x}].$$

Modelo:

$$h(\mathbf{x}) = \theta \left(\sum_{i=0}^d w_i x_i \right) = \theta(\mathbf{w}^T \mathbf{x})$$



Regresión logística

Clase objetivo

Objetivo:

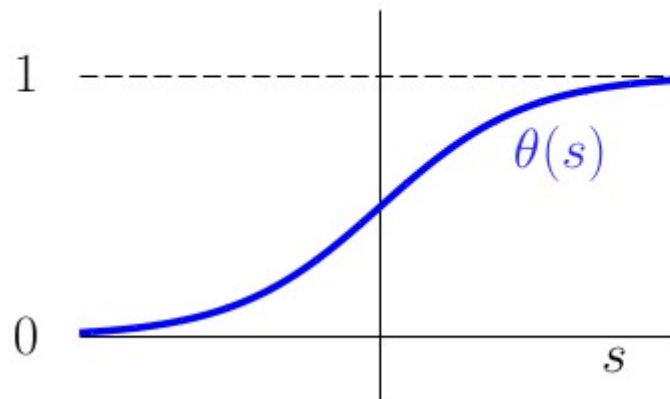
$$f(\mathbf{x}) = \mathbb{P}[y = +1 \mid \mathbf{x}].$$

Modelo:

$$h(\mathbf{x}) = \theta \left(\sum_{i=0}^d w_i x_i \right) = \theta(\mathbf{w}^T \mathbf{x})$$

Función logística:

$$y \in [0, 1]$$



$$\theta(s) = \frac{e^s}{1 + e^s} = \frac{1}{1 + e^{-s}}.$$

$$\theta(-s) = \frac{e^{-s}}{1 + e^{-s}} = \frac{1}{1 + e^s} = 1 - \theta(s).$$

Regresión logística

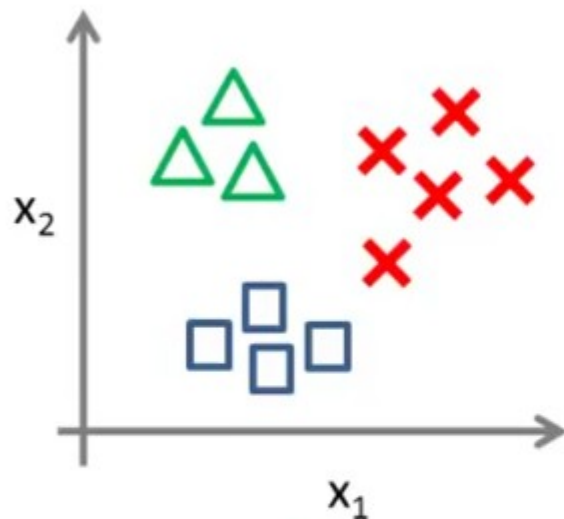
Notar que: $\mathcal{D} = (\mathbf{x}_1, y_1 = \pm 1), \dots, (\mathbf{x}_N, y_N = \pm 1)$


Un buen modelo logra lo siguiente:

$$\begin{cases} h(\mathbf{x}_n) \approx 1 & \text{si } y_n = +1; \\ h(\mathbf{x}_n) \approx 0 & \text{si } y_n = -1. \end{cases}$$


Clasificación multiclase

One-vs-all (one-vs-rest):

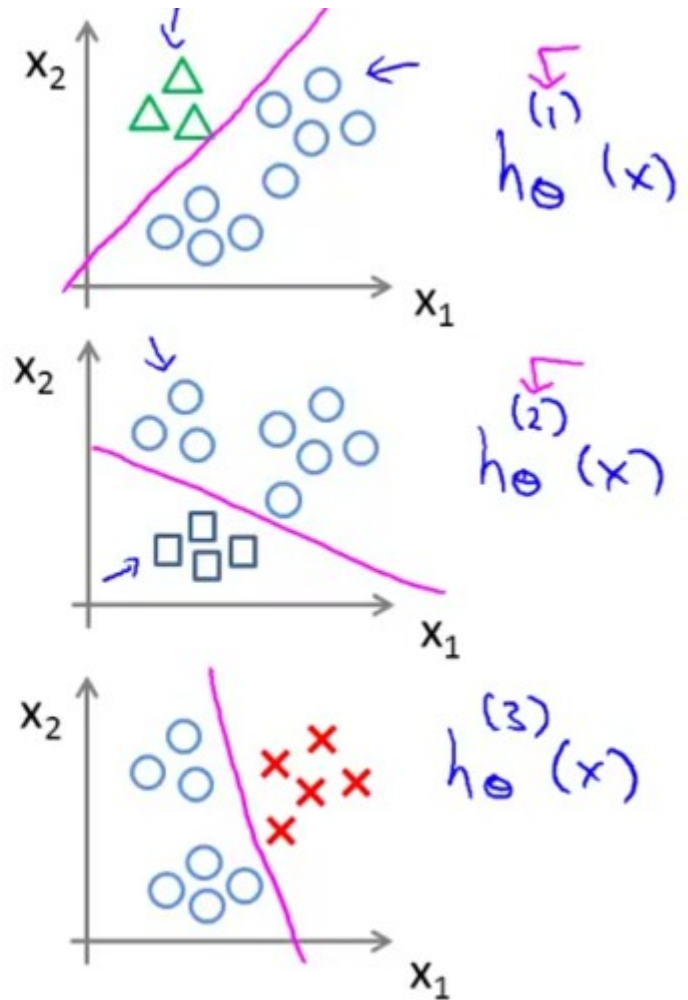


Class 1:  \leftarrow

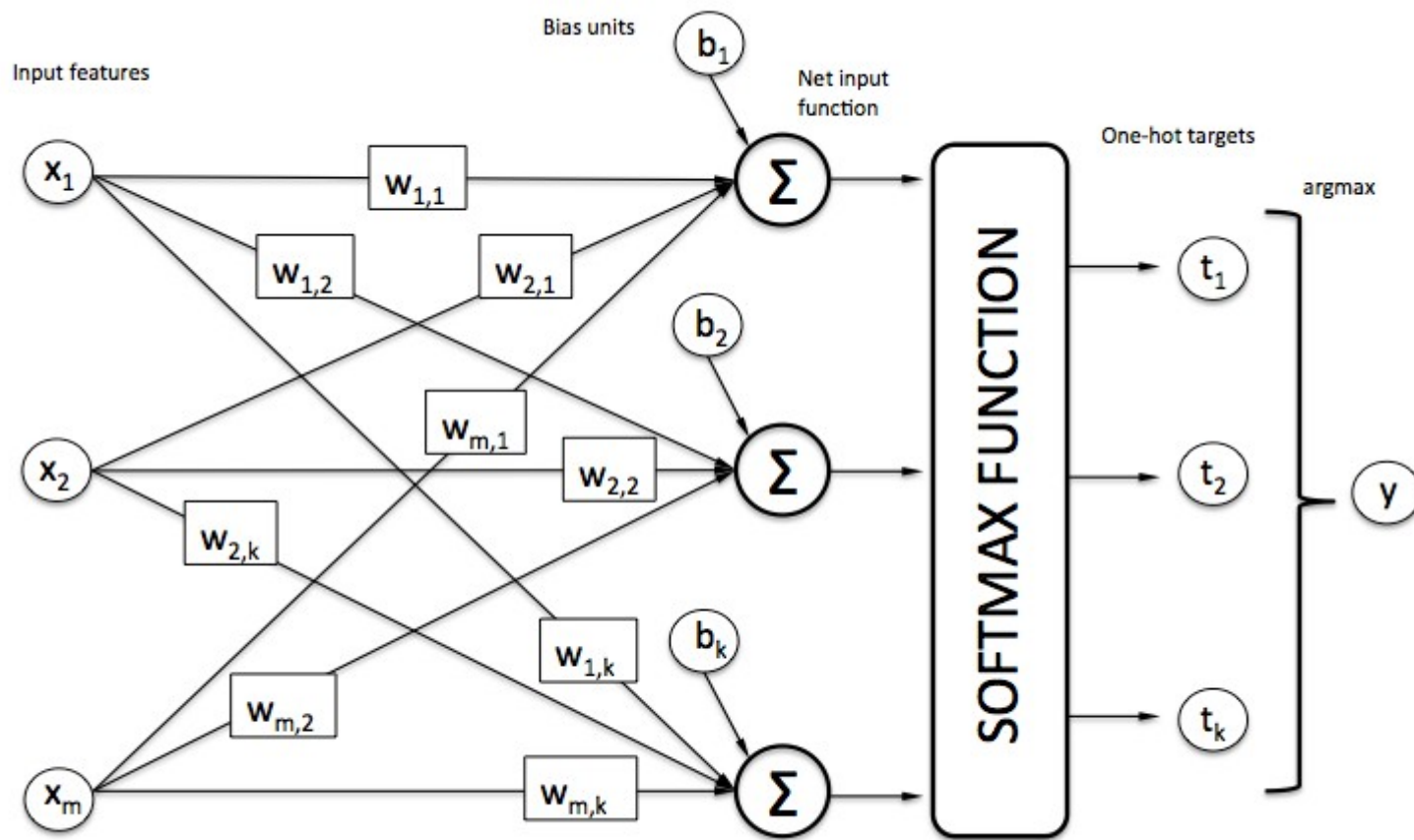
Class 2:  \leftarrow

Class 3:  \leftarrow

$$h_{\theta}^{(i)}(x) = P(y = i|x; \theta) \quad (i = 1, 2, 3)$$



Clasificación multiclase



$$P(y = j \mid z^{(i)}) = \phi_{softmax}(z^{(i)})$$

6
11
7

$$\frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}}$$

.01
.97
.02

- MÉTRICAS DE RENDIMIENTO -

Métricas

		+1	-1		
VALORES PREDICCIÓN	+1	Verdaderos positivos	Falsos Positivos	+1	
	-1	Falsos Negativos	Verdaderos Negativos	-1	
		VALORES REALES			

Métricas

verdaderos positivos (TP)	falsos positivos (FP)
falsos negativos (FN)	verdaderos negativos (TN)

$$\text{accuracy} = (TP + TN) / (TP + FP + FN + TN)$$

$$P = TP / (TP + FP)$$

$$R = TP / (TP + FN)$$

$$F = \frac{(\beta^2 + 1)PR}{\beta^2 P + R}$$

F-measure balanceada con $\beta = 1$

- GRADIENTE DESCENDENTE -

Regresión logística

Si el ajuste es adecuado:

$$P(y \mid \mathbf{x}) = \theta(y \cdot \mathbf{w}^T \mathbf{x}) \quad \leftarrow$$



Luego, podemos expresar la función de verosimilitud:

$$P(y_1, \dots, y_N \mid \mathbf{x}_1, \dots, \mathbf{x}_n) = \prod_{n=1}^N P(y_n \mid \mathbf{x}_n).$$

Regresión logística

Maximizamos la función de verosimilitud:

$$\begin{aligned} & \max \quad \prod_{n=1}^N P(y_n \mid \mathbf{x}_n) \\ \Leftrightarrow & \max \quad \ln \left(\prod_{n=1}^N P(y_n \mid \mathbf{x}_n) \right) \\ \equiv & \max \quad \sum_{n=1}^N \ln P(y_n \mid \mathbf{x}_n) \end{aligned}$$

Regresión logística

Maximizamos la función de verosimilitud:

$$\max \quad \prod_{n=1}^N P(y_n \mid \mathbf{x}_n)$$

$$\Leftrightarrow \max \quad \ln \left(\prod_{n=1}^N P(y_n \mid \mathbf{x}_n) \right)$$

$$\equiv \max \quad \sum_{n=1}^N \ln P(y_n \mid \mathbf{x}_n)$$

$$\Leftrightarrow \text{min} \quad - \frac{1}{N} \sum_{n=1}^N \ln P(y_n \mid \mathbf{x}_n)$$

$$\equiv \min \quad \frac{1}{N} \sum_{n=1}^N \ln \frac{1}{P(y_n \mid \mathbf{x}_n)}$$

$$\equiv \min \quad \frac{1}{N} \sum_{n=1}^N \ln \frac{1}{\theta(y_n \cdot \mathbf{w}^T \mathbf{x}_n)}$$

Regresión logística

Maximizamos la función de verosimilitud:

$$\max \prod_{n=1}^N P(y_n \mid \mathbf{x}_n)$$

$$\Leftrightarrow \max \ln \left(\prod_{n=1}^N P(y_n \mid \mathbf{x}_n) \right)$$

$$\equiv \max \sum_{n=1}^N \ln P(y_n \mid \mathbf{x}_n)$$

$$\Leftrightarrow \text{min} - \frac{1}{N} \sum_{n=1}^N \ln P(y_n \mid \mathbf{x}_n)$$

$$\equiv \min \frac{1}{N} \sum_{n=1}^N \ln \frac{1}{P(y_n \mid \mathbf{x}_n)}$$

$$\equiv \min \frac{1}{N} \sum_{n=1}^N \ln \frac{1}{\theta(y_n \cdot \mathbf{w}^T \mathbf{x}_n)}$$

$$\equiv \min \frac{1}{N} \sum_{n=1}^N \ln(1 + e^{-y_n \cdot \mathbf{w}^T \mathbf{x}_n})$$

$$\theta(s) = \frac{1}{1 + e^{-s}}.$$



Regresión logística

Tenemos una expresión para:

Parámetros del modelo

$$E_{\text{in}}(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N \ln(1 + e^{-y_n \cdot \mathbf{w}^T \mathbf{x}_n})$$

Cross-entropy

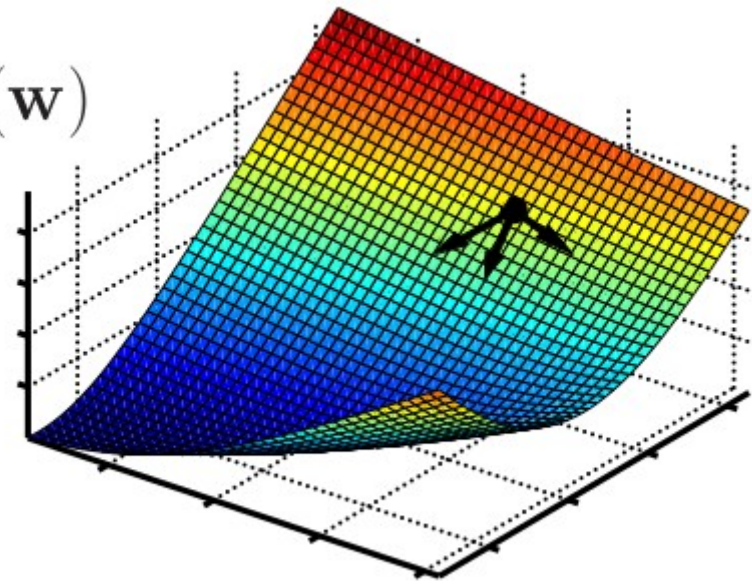
La función es convexa, por lo que podemos optimizarla de forma iterativa:

Idea del gradiente descendente:

$E_{\text{in}}(\mathbf{w})$

$$\mathbf{w}(t+1) = \mathbf{w}(t) + \eta \hat{\mathbf{v}}$$

¿En cuál dirección conviene moverse?



Regresión logística

Tenemos una expresión para:

Parámetros del modelo

$$E_{\text{in}}(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N \ln(1 + e^{-y_n \cdot \mathbf{w}^T \mathbf{x}_n})$$

Cross-entropy

La función es convexa, por lo que podemos optimizarla de forma iterativa:

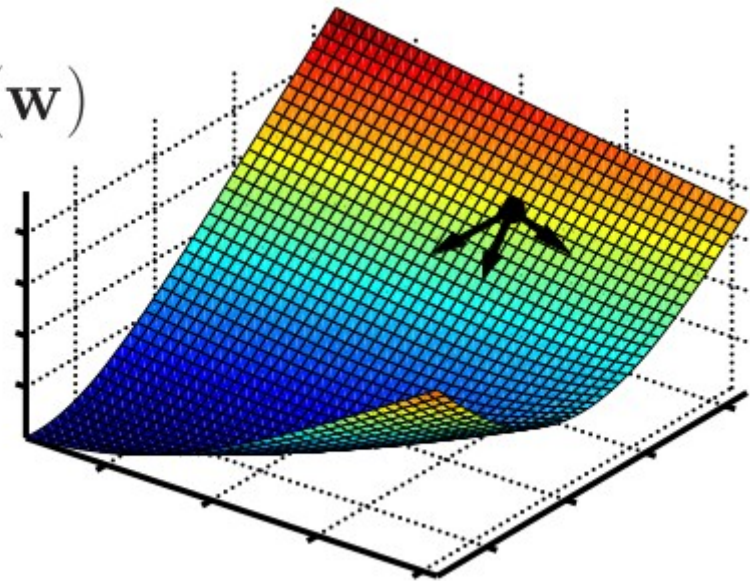
Idea del gradiente descendente:

$$\mathbf{w}(t+1) = \mathbf{w}(t) + \eta \hat{\mathbf{v}}$$

¿En cuál dirección conviene moverse?

Move in the direction $\mathbf{v}_t = -\mathbf{g}_t$.

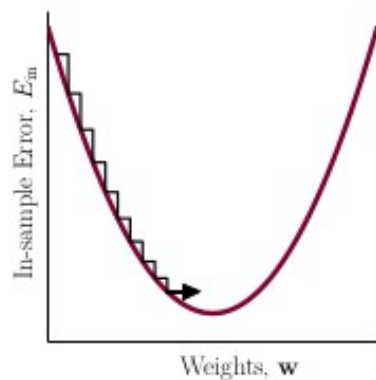
$E_{\text{in}}(\mathbf{w})$



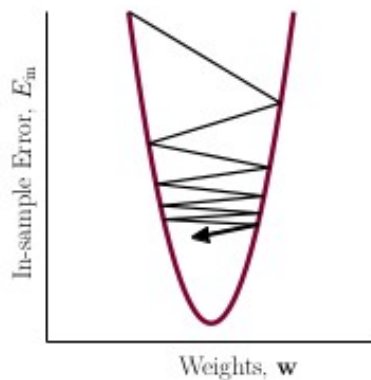
Gradiente descendente

El efecto del **learning rate**:

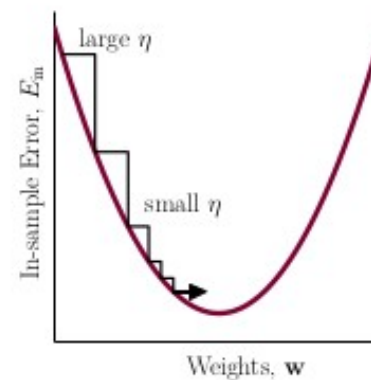
η muy pequeño



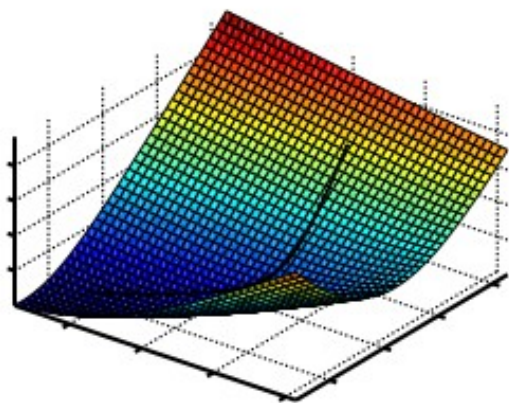
η muy grande



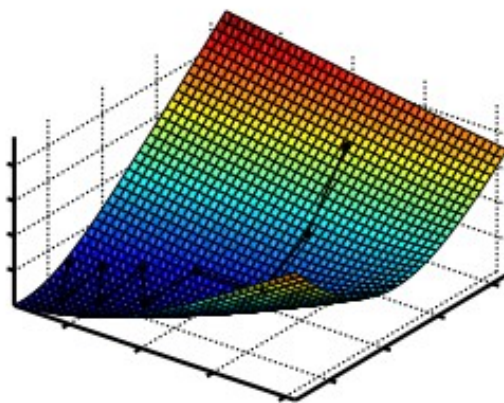
η_t



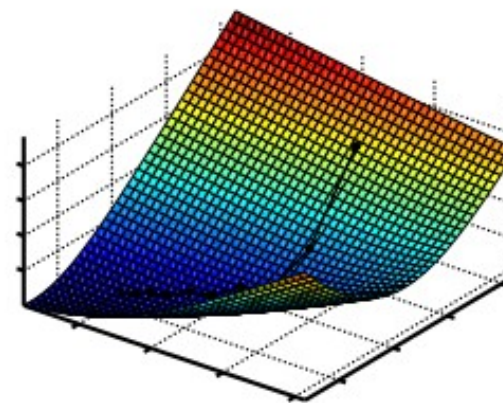
Adam



$\eta = 0.1$; 75 steps



$\eta = 2$; 10 steps



variable η_t ; 10 steps

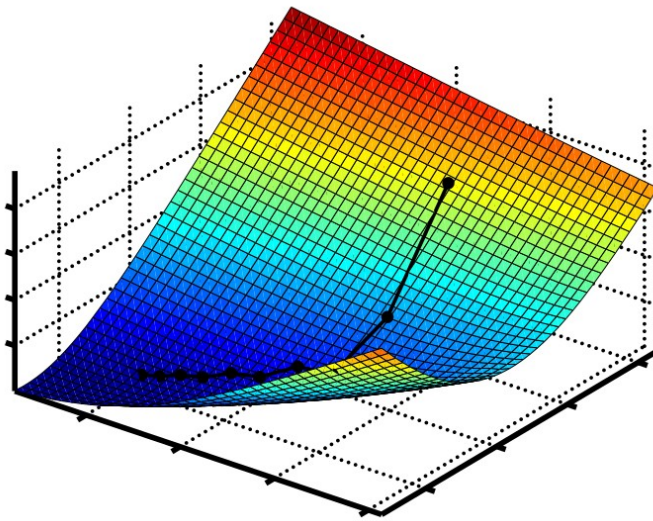
Gradiente descendente

Una variación de gradiente descendente considera la evaluación del gradiente en una muestra del training set.

$$\mathbf{w}(t+1) \leftarrow \mathbf{w}(t) - \eta \nabla_{\mathbf{w}} e(\mathbf{w}, \mathbf{x}_*, y_*)$$

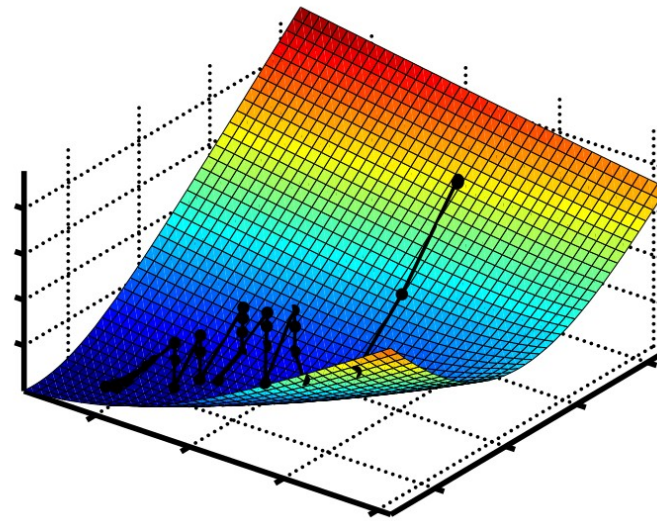
Dado que la muestra se toma al azar, se le denomina gradiente descendente estocástico.

GD



$\eta = 6$
10 steps
 $N = 10$

SGD



$\eta = 2$
30 steps

Puede ayudar a escapar
de óptimos locales

Explainers

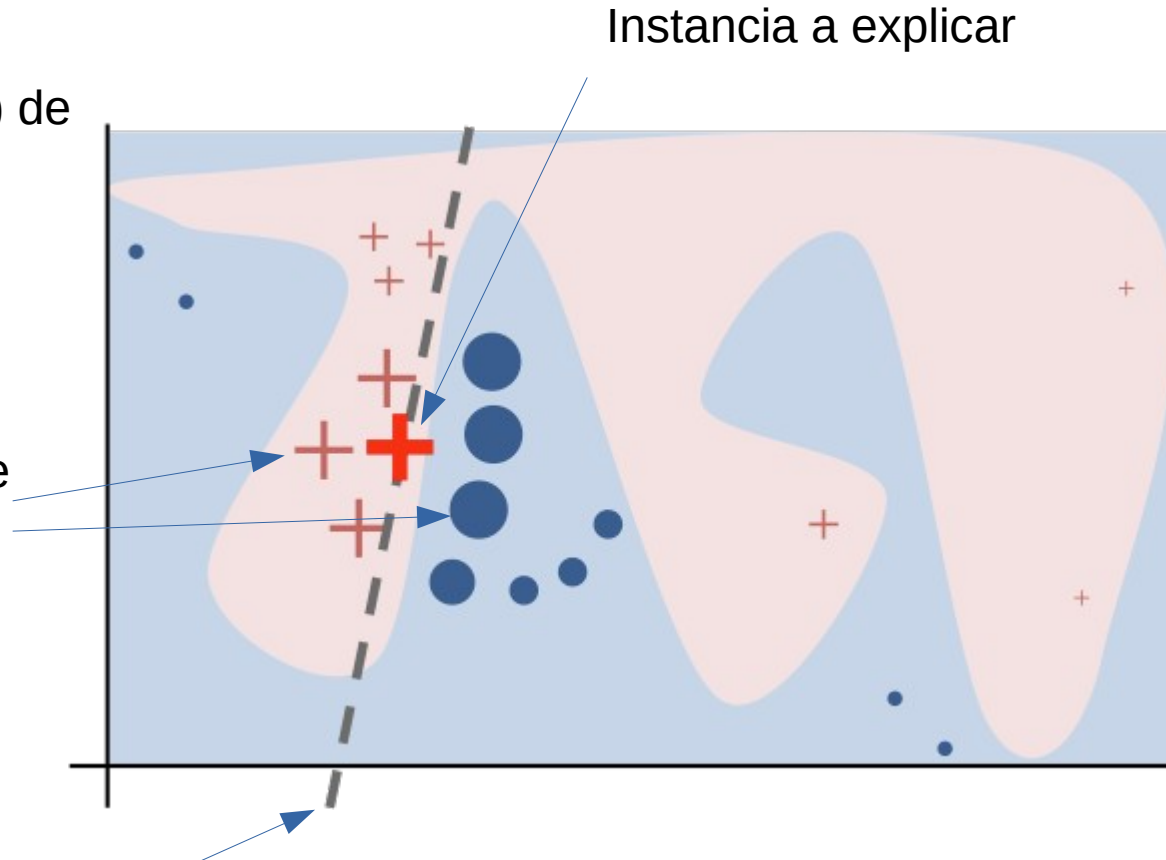
LIME (Local Interpretable Model-Agnostic Explanations)

1. Se perturban
(eliminan features) de
la instancia

2. Se obtienen las
predicciones sobre
los ejemplos
perturbados

3. Se construye un modelo lineal que aproxima al
modelo original usando las predicciones (noisy label)

4. Se identifica cuales perturbaciones inclinan la decisión



LIME (Local Interpretable Model-Agnostic Explanations)



@Botcheckl: a service for anomalous activity detection in Twitter

You can enter the username or uid of the account you want to analyze

Input

recuerdamebot

Analyze activity

Prediction probabilities

Human	0.26
Likely bot	0.74

Human

Likely bot

EMOTICONS M...
0.06
0.11 < MENTIO...
0.06
NAME # HAPPY...
0.05
USER # ADP <=...
0.05
HAPPY EMOJIS...
0.04
NAME # EMOTIC...
0.03
DOMINANCE MI...
0.03
USER # RETWEE...
0.03
SURPRISE E...
0.01

VALENCE MIN ...
0.01

- ANEXO: GRADIENTE DESCENDENTE -

Gradiente descendente

$$\begin{aligned}\Delta E_{\text{in}} &= E_{\text{in}}(\mathbf{w}(t+1)) - E_{\text{in}}(\mathbf{w}(t)) \\ &= E_{\text{in}}(\mathbf{w}(t) + \eta \hat{\mathbf{v}}) - E_{\text{in}}(\mathbf{w}(t)) \\ &= \eta \nabla E_{\text{in}}(\mathbf{w}(t))^{\text{T}} \hat{\mathbf{v}} + O(\eta^2)\end{aligned}$$

$$f(x) = f(x_0) + f'(x_0)(x - x_0) + \frac{f''(x_0)}{2!}(x - x_0)^2 + \dots + \frac{f^{(k+1)}(c)}{(k+1)!}(x - x_0)^{k+1}.$$

$$\begin{aligned}\Delta E_{\text{in}} &= E_{\text{in}}(\mathbf{w}(t+1)) - E_{\text{in}}(\mathbf{w}(t)) \\ &\quad \underbrace{f(x_0 + \eta \cdot x) - f(x_0)}_{\text{(expansión de Taylor de 1er orden)}} \\ &= f(x_0) + \nabla f(x_0) \cdot (x_0 + \eta \cdot x - x_0) + \dots\end{aligned}$$

Gradiente descendente

$$\begin{aligned}
 \Delta E_{\text{in}} &= E_{\text{in}}(\mathbf{w}(t+1)) - E_{\text{in}}(\mathbf{w}(t)) \\
 &= E_{\text{in}}(\mathbf{w}(t) + \eta \hat{\mathbf{v}}) - E_{\text{in}}(\mathbf{w}(t)) \\
 &= \eta \nabla E_{\text{in}}(\mathbf{w}(t))^T \hat{\mathbf{v}} + O(\eta^2)
 \end{aligned}$$

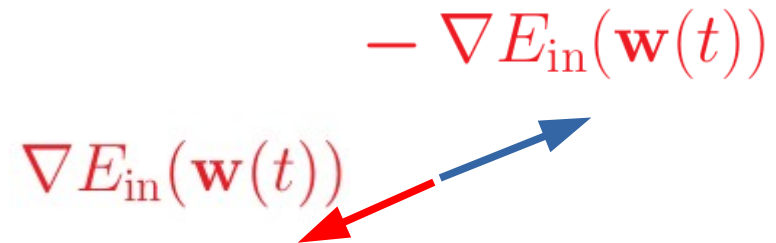
$$f(x) = f(x_0) + f'(x_0)(x - x_0) + \frac{f''(x_0)}{2!}(x - x_0)^2 + \dots + \frac{f^{(k+1)}(c)}{(k+1)!}(x - x_0)^{k+1}.$$

$$\begin{aligned}
 \Delta E_{\text{in}} &= E_{\text{in}}(\mathbf{w}(t+1)) - E_{\text{in}}(\mathbf{w}(t)) \\
 &\quad \underbrace{f(x_0 + \eta \cdot x) - f(x_0)}_{\text{(expansión de Taylor de 1er orden)}} \\
 &\quad \cancel{f(x_0)} + \nabla f(x_0) \cdot (\cancel{x_0} + \eta \cdot x - \cancel{x_0}) + \dots
 \end{aligned}$$

Gradiente descendente

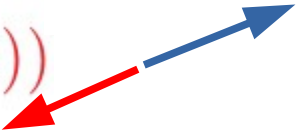
$$= \eta \underbrace{\nabla E_{\text{in}}(\mathbf{w}(t))^T}_{\text{Minimizado en}} \hat{\mathbf{v}} + O(\eta^2)$$

$$\text{Minimizado en } \hat{\mathbf{v}} = - \frac{\nabla E_{\text{in}}(\mathbf{w}(t))}{\|\nabla E_{\text{in}}(\mathbf{w}(t))\|}$$



Gradiente descendente

$$= \eta \underbrace{\nabla E_{\text{in}}(\mathbf{w}(t))^T \hat{\mathbf{v}}}_{\text{Minimizado en } \hat{\mathbf{v}} = -\frac{\nabla E_{\text{in}}(\mathbf{w}(t))}{\|\nabla E_{\text{in}}(\mathbf{w}(t))\|}} + O(\eta^2)$$

$$\nabla E_{\text{in}}(\mathbf{w}(t)) \quad -\nabla E_{\text{in}}(\mathbf{w}(t))$$


Gradiente descendente

1: Initialize at step $t = 0$ to $\mathbf{w}(0)$.

2: **for** $t = 0, 1, 2, \dots$ **do**

3: Compute the gradient

$$\mathbf{g}_t = \nabla E_{\text{in}}(\mathbf{w}(t)).$$

4: Move in the direction $\mathbf{v}_t = -\mathbf{g}_t$.

5: Update the weights:

$$\mathbf{w}(t+1) = \mathbf{w}(t) + \eta \mathbf{v}_t.$$

6: Iterate ‘until it is time to stop’.

7: **end for**

8: Return the final weights.