

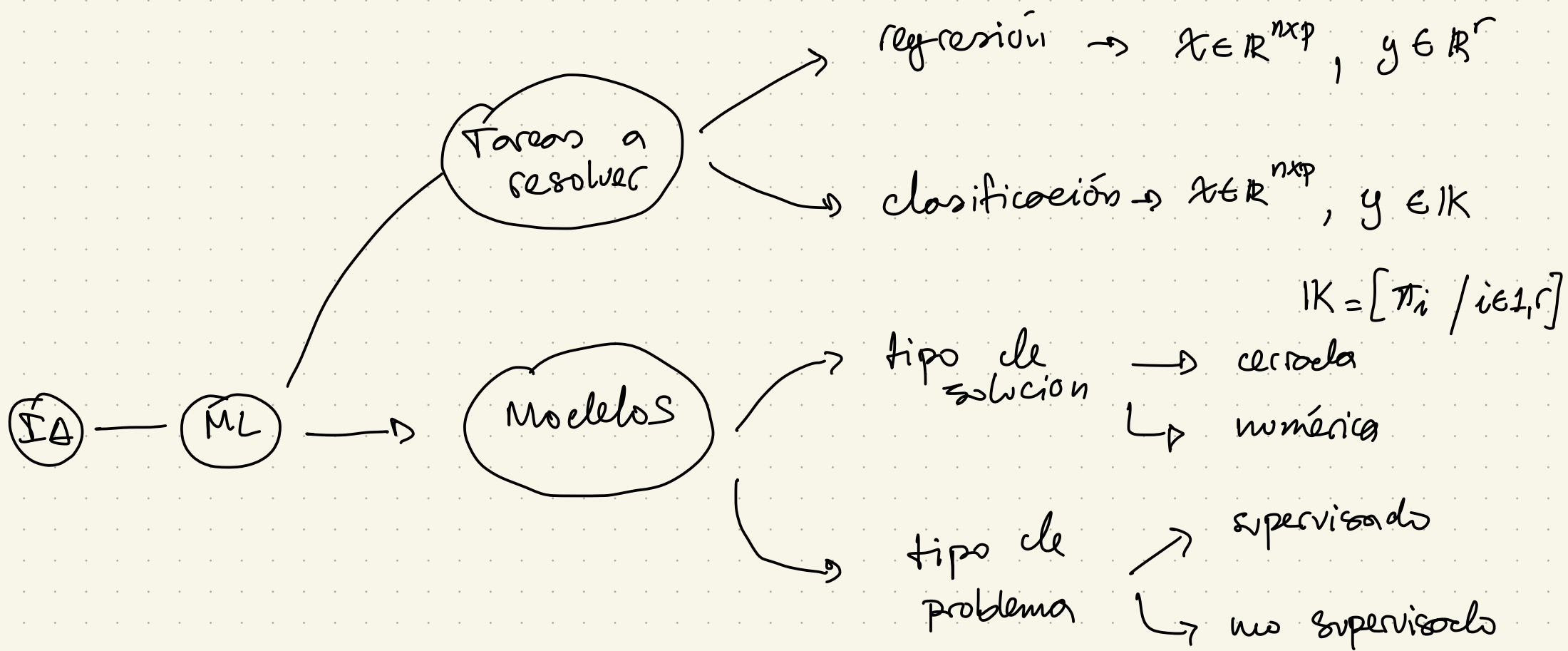
Introducción a la Inteligencia Artificial  
Facultad de Ingeniería  
Universidad de Buenos Aires



## Clase 5

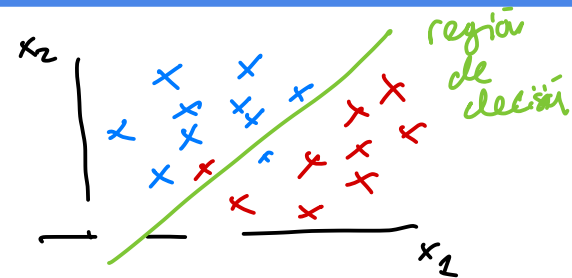
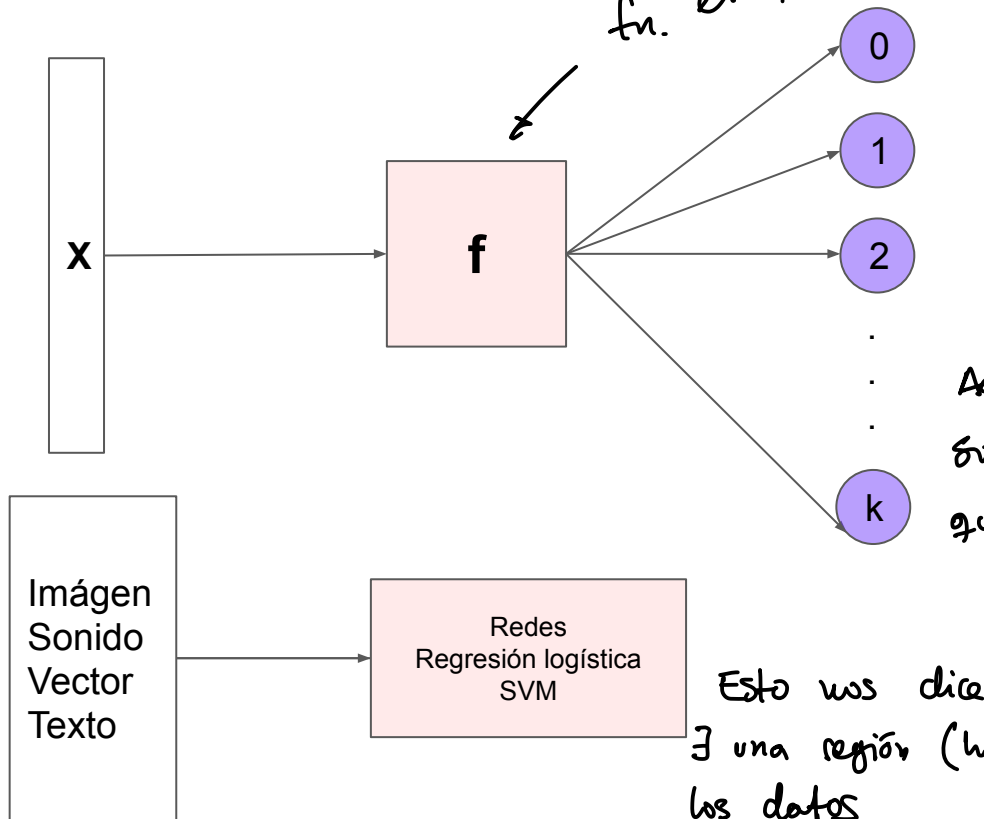
1. Clasificación Binaria
  - a. Motivación
  - b. Regresión Logística - Ejercicio de Aplicación
  - c. Regresión Logística - Teoría
2. Clasificación Multiclase
  - a. Motivación
  - b. Softmax
  - c. Ejercicio de Aplicación
3. Ejercicio integrador





## Clasificación

considero  $X \in \mathcal{D} \subset \mathbb{R}^{n \times p}$



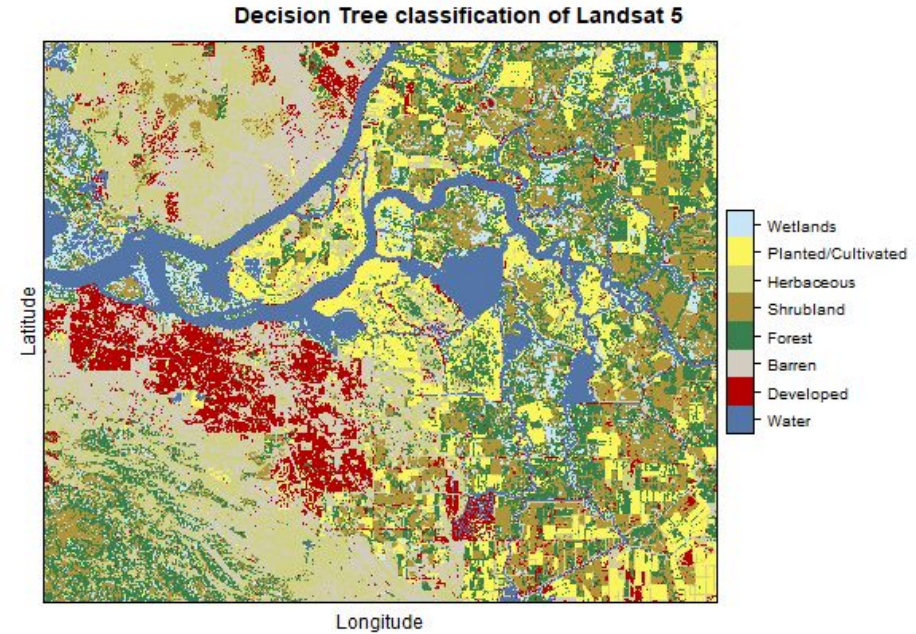
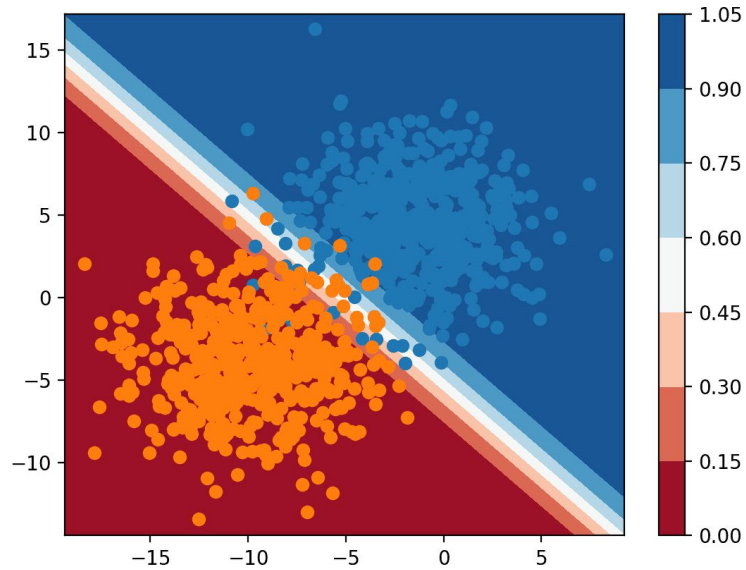
función etiquetadora  
 $f: \mathbb{R}^D \rightarrow \{0, 1, 2, \dots, k\}$

Ara estamos haciendo un supuesto  
SUPER importante, Estamos diciendo  
que  $Y$  es separable:

$$\forall x_i \exists y_i / y_i \in K \wedge f(y_i)!$$

Esto nos dice que en el espacio de features  
 $\exists$  una región (hiperplano) que separa  
los datos

## Clasificación



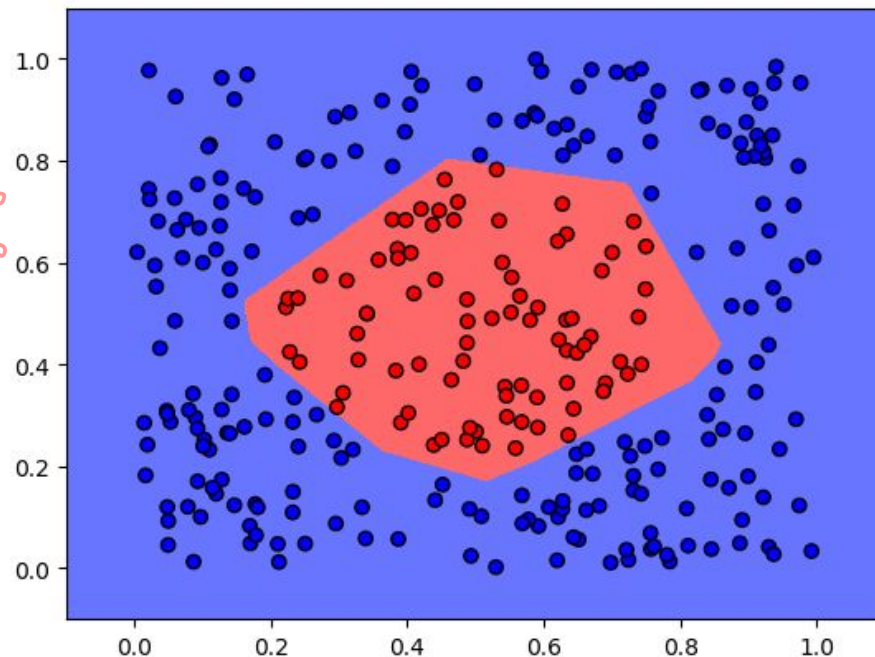
## Clasificación binaria

$$f(x_i) = \begin{cases} 1 & \text{si } y_i \in C_1 \\ 0 & \text{o.w.} \end{cases}$$

## Clasificación Binaria - Ejemplos

- Detección de fraudes
- Diagnóstico médico
- Detección de spam
- Sentiment Analysis
- Detección de objetos
- Outliers

*es fraude*  
*es legal*  
*es maligno*  
*es benigno*



## Clasificación Binaria - Diagnóstico médico

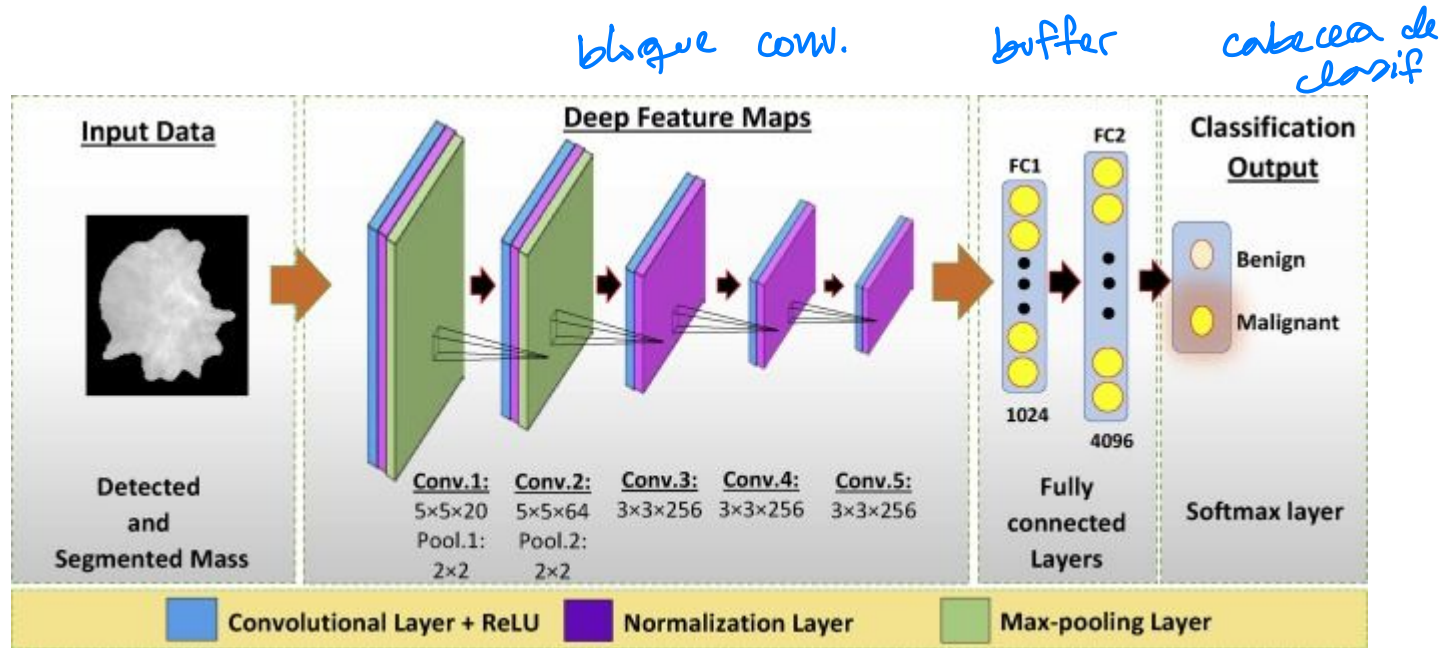


Imagen de: "A fully integrated computer-aided diagnosis system for digital X-ray mammograms via deep learning detection, segmentation, and classification"



## Clasificación Binaria - Spam detection

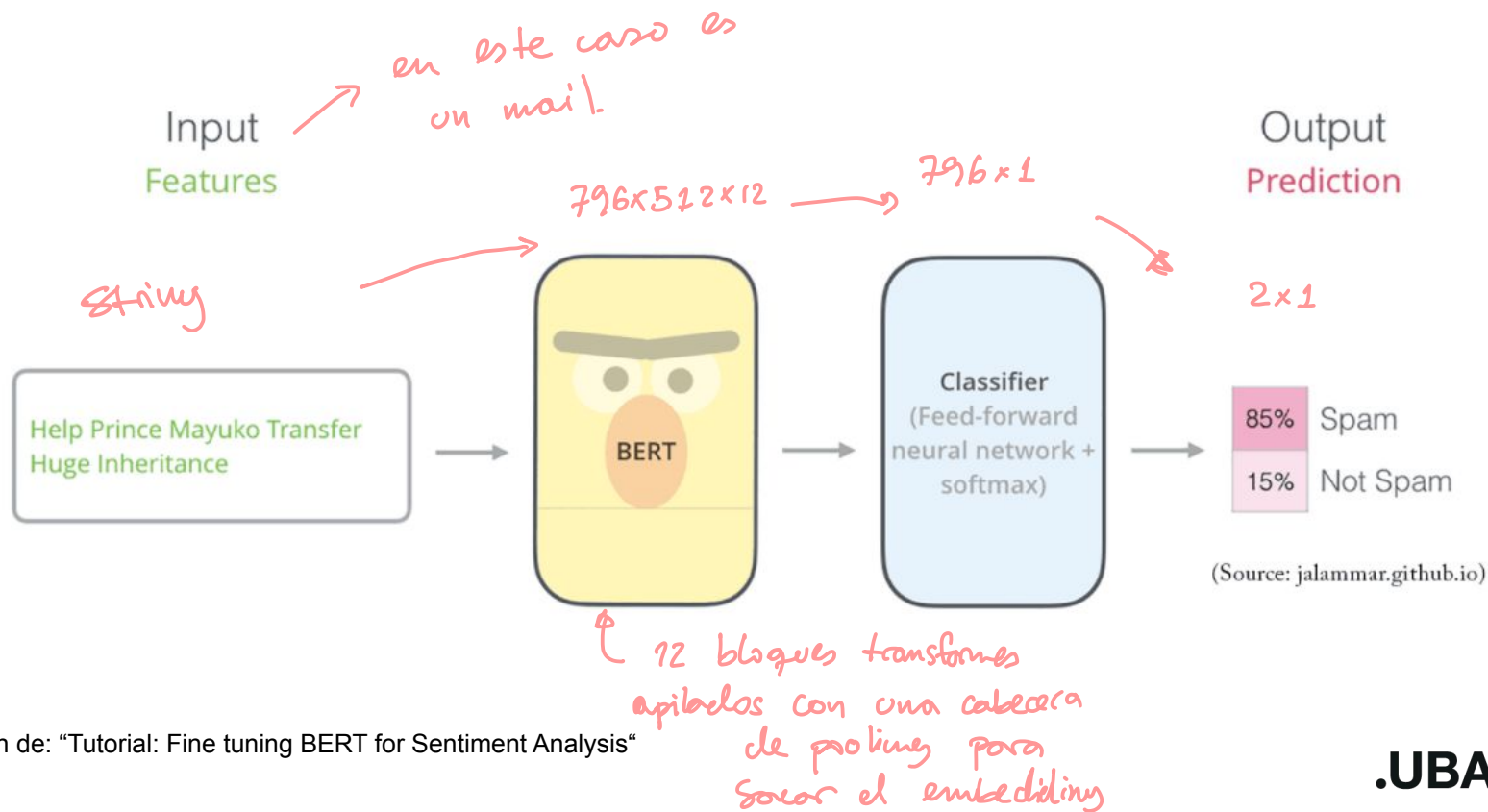
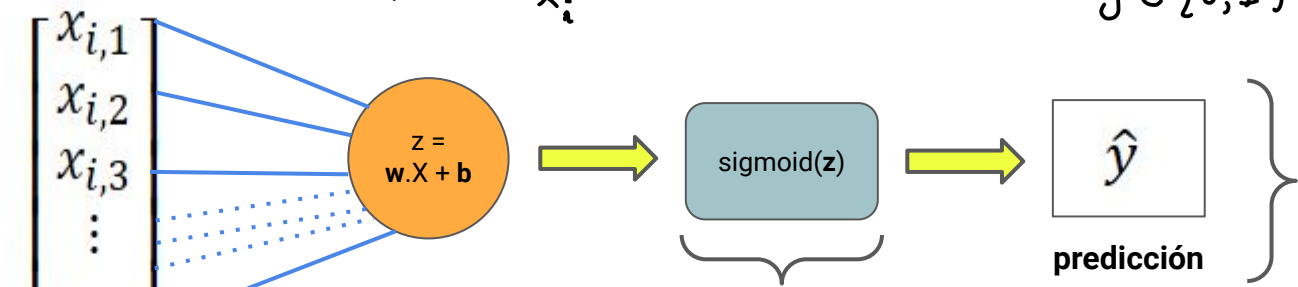
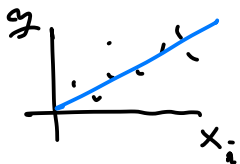


Imagen de: "Tutorial: Fine tuning BERT for Sentiment Analysis"

## Regresión Logística



¿Función de Costo?  
 $\mathcal{L}(y, \hat{y})$

$X$  muestra  $i$   
( $m$  features)

$$x_i \in \mathbb{R}^m$$

$$x_i \in \mathbb{R}$$

$$\bar{w}_i \in \mathbb{R}^{2 \times m}$$

$$b \in \mathbb{R}$$

$$\phi(z)$$

$$\phi: \mathbb{R} \mapsto [0, 1]$$

función de  
squashing

Métricas: en clasif. no podemos usar lo que veníamos usando

MSE (gato, perro)?  $\rightarrow$  en clasif. usamos la matriz de confusión.

real \ Pred		← prediction	
		T	F
GL* {	T	TP	FN
	F	FP	TN

$$y = \mathbb{I}_{C=K} = \begin{cases} 1 & \text{si } y_i = K \\ 0 & \text{si } y_i \neq K \end{cases}$$

\* GL: Golden Label  
Ground truth  
Real values

• Accuracy =  $\frac{TP + TN}{\# \text{muestras}} \in [0, 1]$

• Specificity =  $\frac{TN}{TN + FP}$   
(TN rate)

• precision =  $\frac{TP}{TP + FP}$

• recall =  $\frac{TP}{TP + FN}$

•  $F_1 = 2 \cdot \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}} = \frac{TP}{TP + \frac{1}{2}(FP + FN)}$

•  $F_\beta = (1 + \beta)^2 \frac{\text{prec.} \cdot \text{recall}}{\beta^2 \text{prec} + \text{recall}} \quad (F_{0.25})$

## Formas de clasificar:

+ Definimos un proc. de scoring  $\rightarrow$  como resolvemos una clasif.

$\uparrow$  complejidad  $\Rightarrow \uparrow$  computo  $\Rightarrow$  mejor modelo  
?

1. Modelos generativos: modelar las distrib. de I/O ms generar la salida (permite generar información sintética)  $\rightarrow$  GAN's

2. Modelos discriminantes: plantear  $P(C_k | \bar{X})$  ms utilizar métodos de inferencia (bayesiana) para estimar  $P$ .

3. modelos de función discriminadora: buscamos  $f: \mathbb{R}^{n \times m} \mapsto \{C_1, \dots, C_k\}$

$\hookrightarrow$  Regresión logística (la usamos porque queremos explotar lo conocido).

partimos de suponer  $y \sim \text{Bernoulli}(1, \pi_i) \rightarrow \begin{cases} E(y) = \pi_i \\ \text{Var}(y) = \pi_i \cdot (1 - \pi_i) \end{cases}$

me gustaría poder obtener algo así:

$$\pi_i \sim \bar{X}^t \cdot \bar{\beta} \quad ; \quad \bar{X} \in \mathbb{R}^{n \times p}, \bar{\beta} \in \mathbb{R}^{1 \times p}$$

$\pi_i \in [0,1]$ , pero  $\bar{X}^t \cdot \bar{\beta} \in \mathbb{R} \Rightarrow$  tengo que plantear una fn.  
de mapeo.

buscamos  $f: [0,1] \mapsto \mathbb{R}_{>0} \leadsto \text{odds}_i = \frac{\pi_i}{1 - \pi_i}$  razón de probabilidad

$$\text{odds}_i = \frac{P(A)}{P(\bar{A})}$$

Vamos a tomar el log de odds, esto me permite definir una función biyectiva  $\mathbb{R}_{>0} \mapsto \mathbb{R}$ . Esta transf. se llama logit

(log odds ratio):

$$\eta_i = \text{logit}(\pi_i) = \bar{X}^t \cdot \bar{\beta}$$

existe antilogit:

$$\pi_i = \text{antilogit}(\eta_i) =$$

$$\frac{e^{\eta_i}}{1 + e^{\eta_i}}$$

fn. sigmoide  
fn. logística

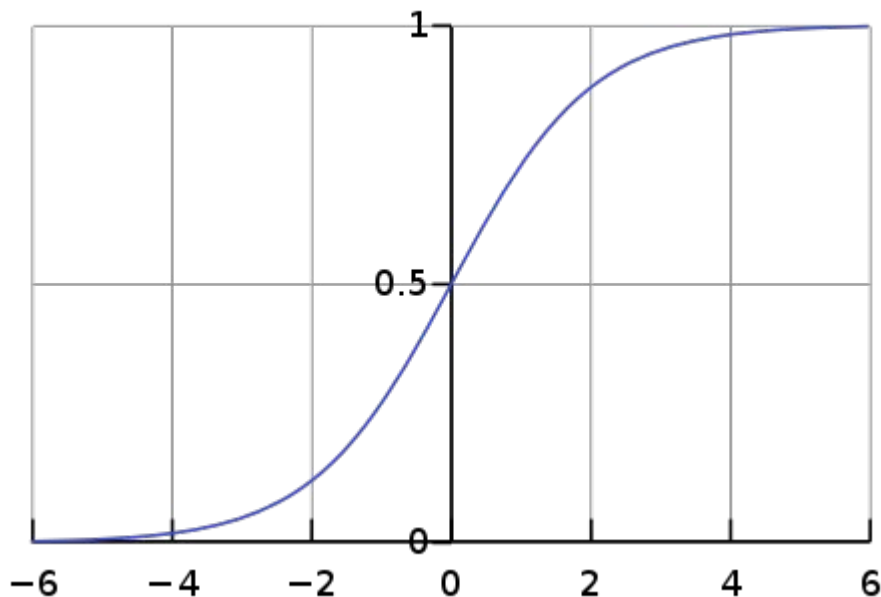
## Logistic function

$$\sigma(x)$$
$$S(x) = \frac{1}{1 + e^{-x}} = \frac{e^x}{e^x + 1} = 1 - S(-x).$$

$$S(x) = \sigma(x) \in [0, 1]$$

$$S, \sigma: \mathbb{R} \mapsto [0, 1]$$

$$\sigma(x, a)$$



## Modelo de regresión logística:

Sean  $y_1, \dots, y_n$  realizaciones de un proc. Bernoulli  $(1, \pi_i)$ , asumimos que existe una relación **lineal** entre  $x_i$  (datos) y el **logit** de la **razón de prob.** de  $\pi_i$ :

$$\text{logit}(\pi_i) = \beta_0 + \sum_{j=0}^p \beta_j \cdot x_{ij}$$

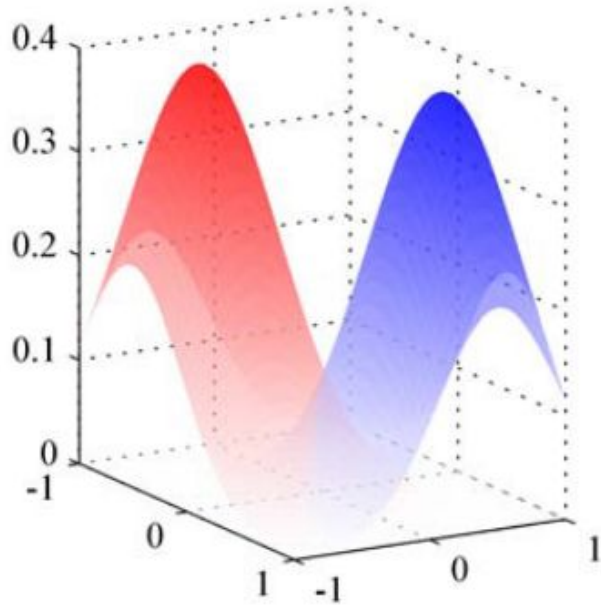
Este modelo es parte de la familia de **modelos lineales generalizados**

→ Es equivalente a decir **modelo lineal generalizado con respuesta binomial** y **fn. de enlace logit**.

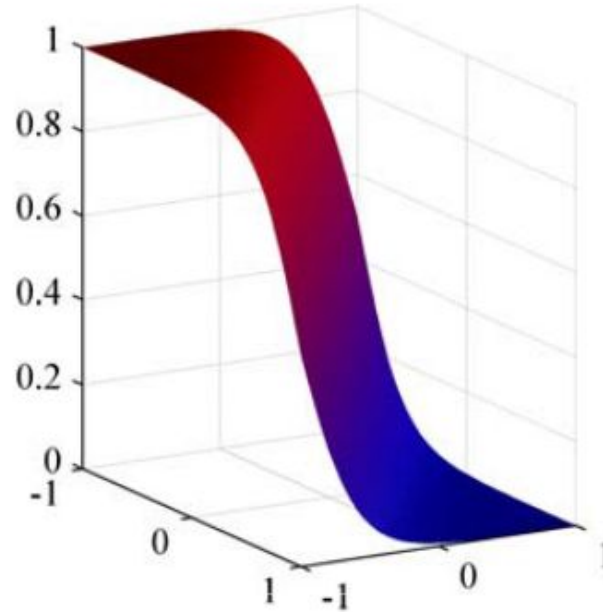
$$g(E(y/x)) = \beta_0 + \sum_i \beta_i x_i$$

$y \sim f$ ;  $f$  distrib. de la respuesta  $\wedge g$  es la función de enlace

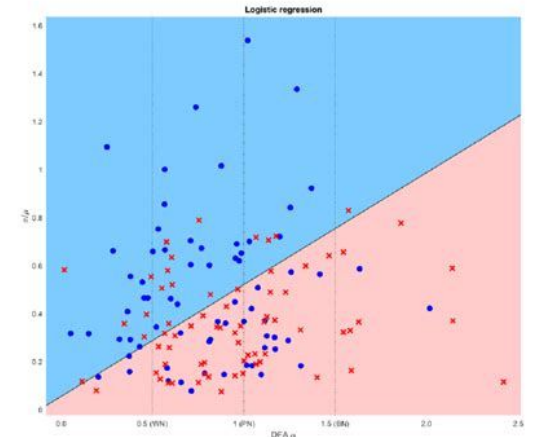
## Regresión Logística



Class-conditional -  $P(x|C_n)$



Posterior -  $P(C_n|x)$





¿Cómo obtenemos  $\beta$ ?

queremos mapear  $\sigma(z) = \sigma(w^t x)$  prop tal  $\partial_a \sigma(a) = \sigma(a)(1 - \sigma(a))$

$$\partial_w \sigma(z) = \sigma(w^t x) (1 - \sigma(w^t x)) \quad (1)$$

planteamos la fn. de verosimilitud:  $P(y/x) = \prod_{i=1}^N \hat{y}_i^{y_i} (1 - \hat{y}_i)^{1-y_i} = \mathcal{L}$

$\begin{cases} N: \text{cant. de muestras} \\ y: \text{golden label} \\ \hat{y}: \text{pred (prob)} \end{cases}$

$$\hat{y}_i = P(c/x)$$

$$\mathcal{L} = \sum_{i=1}^N \ln(P_w(y = y_i | X = x_i)) \rightarrow \text{buscamos maximizar } \mathcal{L}$$

$$\max_w \sum_i \ln(\sigma(w^t x)^{y_i} \cdot (1 - \sigma(w^t x))^{1-y_i})$$

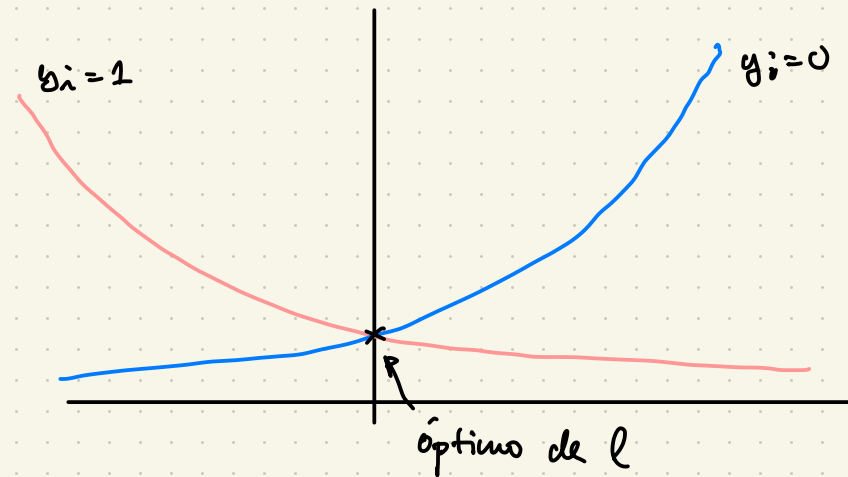
$$\max_w \sum_i y_i \underbrace{\ln \sigma(w^t x)}_{(A)} + (1-y_i) \underbrace{\ln (1 - \sigma(w^t x))}_{(B)}$$

multiplicar  $\times (-1)$ :

$$\min_w \sum_i -y_i \textcircled{A} - (1-y_i) \textcircled{B} \quad \text{no si minimizo } \mathcal{L} \text{ encuentro el óptimo}$$

podemos buscar el mínimo usando la entropía binaria cruzada (binary crossentropy):

$$\mathcal{J}(w) = \frac{1}{N} \sum_i -y_i \textcircled{A} - (1-y_i) \textcircled{B}$$

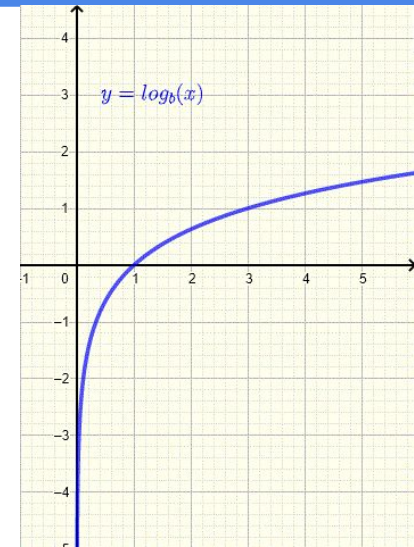
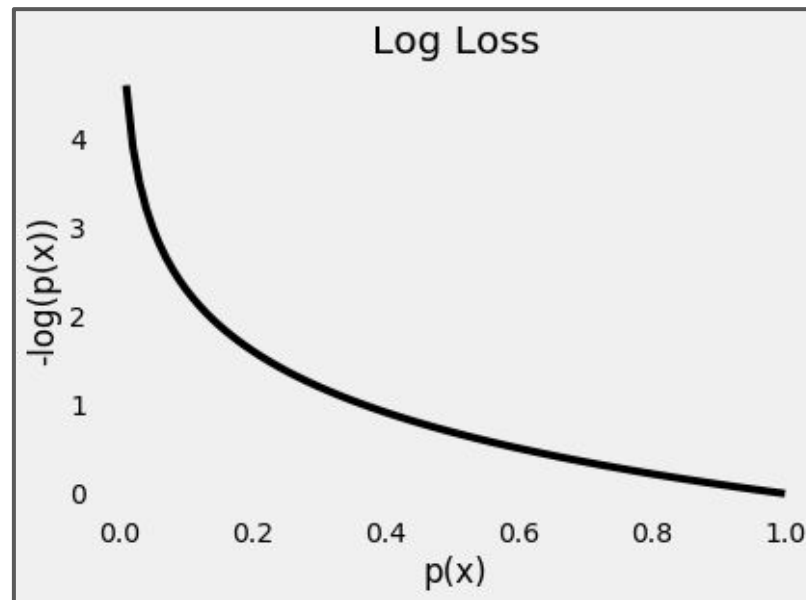


$$\nabla \mathcal{J}_w = 0$$

$\hookrightarrow$  GD  
SGD  
GD MB

## Función de costo - Binary cross entropy

$$H_p(q) = -\frac{1}{N} \sum_{i=1}^N y_i \cdot \log(p(y_i)) + (1 - y_i) \cdot \log(1 - p(y_i))$$

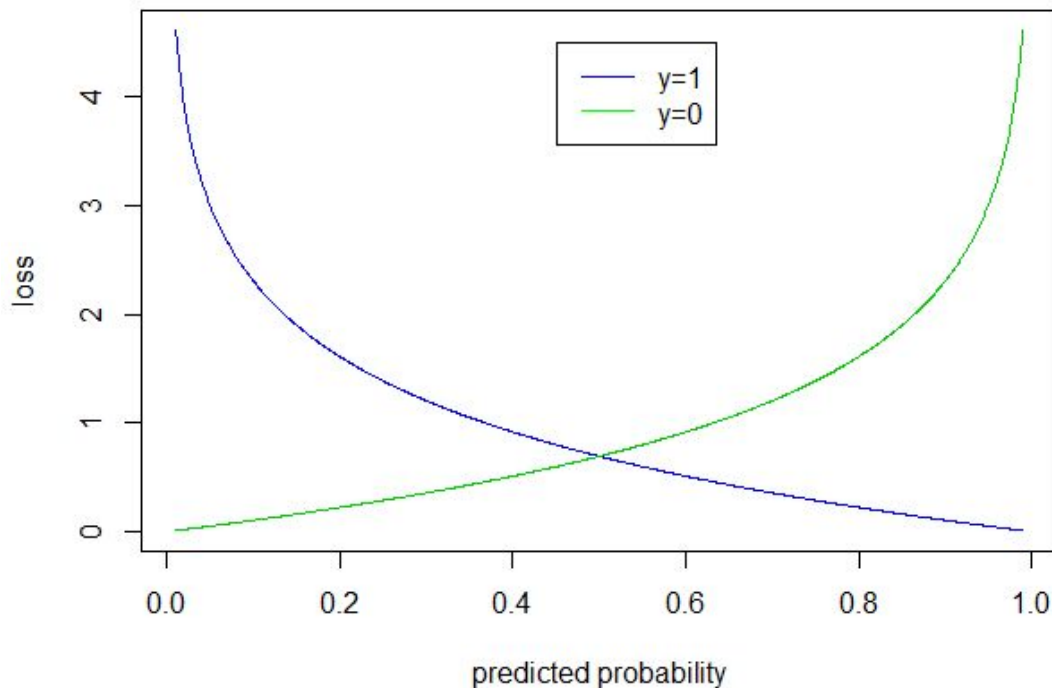


## Función de costo - Binary cross entropy

$$\log \phi(y_i) = \log(\sigma(w^T x))$$

↑  
↑  
data  
to predict output

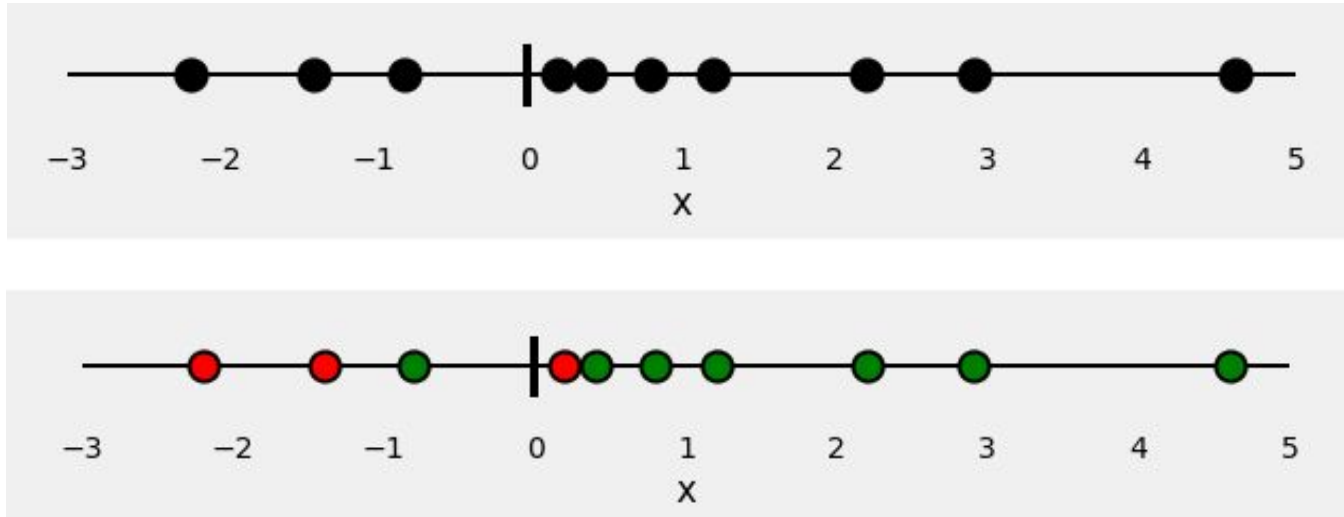
$$H_p(q) = -\frac{1}{N} \sum_{i=1}^N y_i \cdot \log(p(y_i)) + (1 - y_i) \cdot \log(1 - p(y_i))$$



$$w_{n+1} = w_n + \alpha \nabla w$$

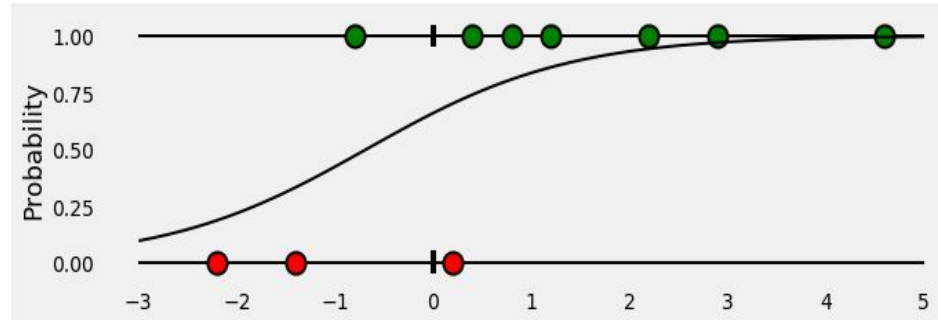
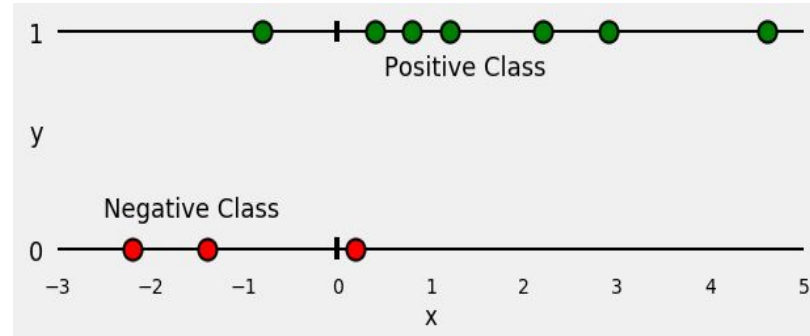
$\nabla w$ ?

## Regresión Logística

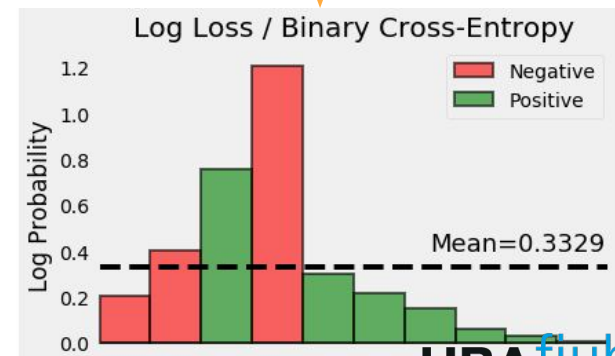
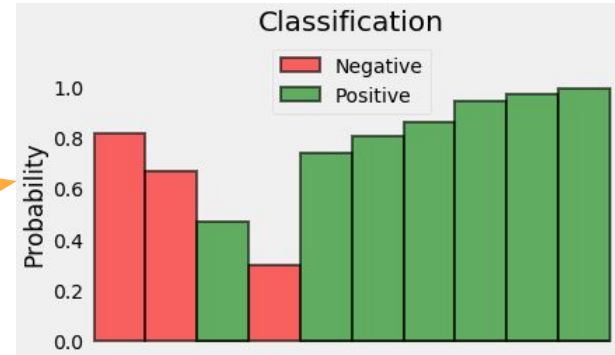
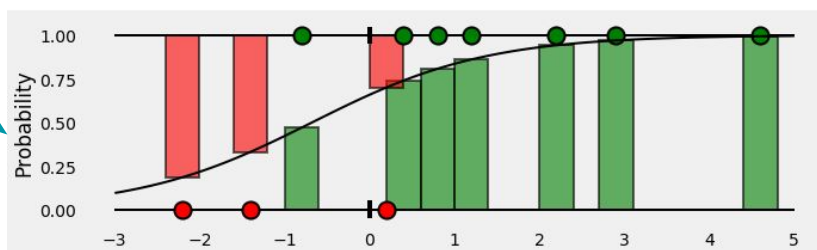
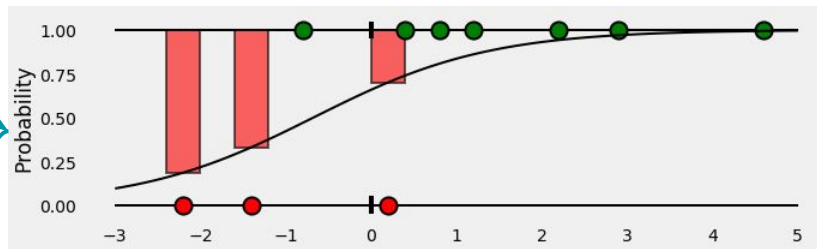
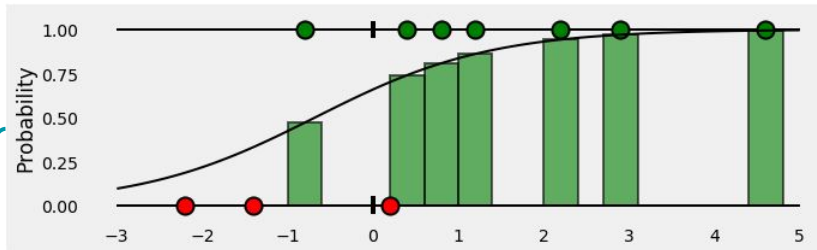


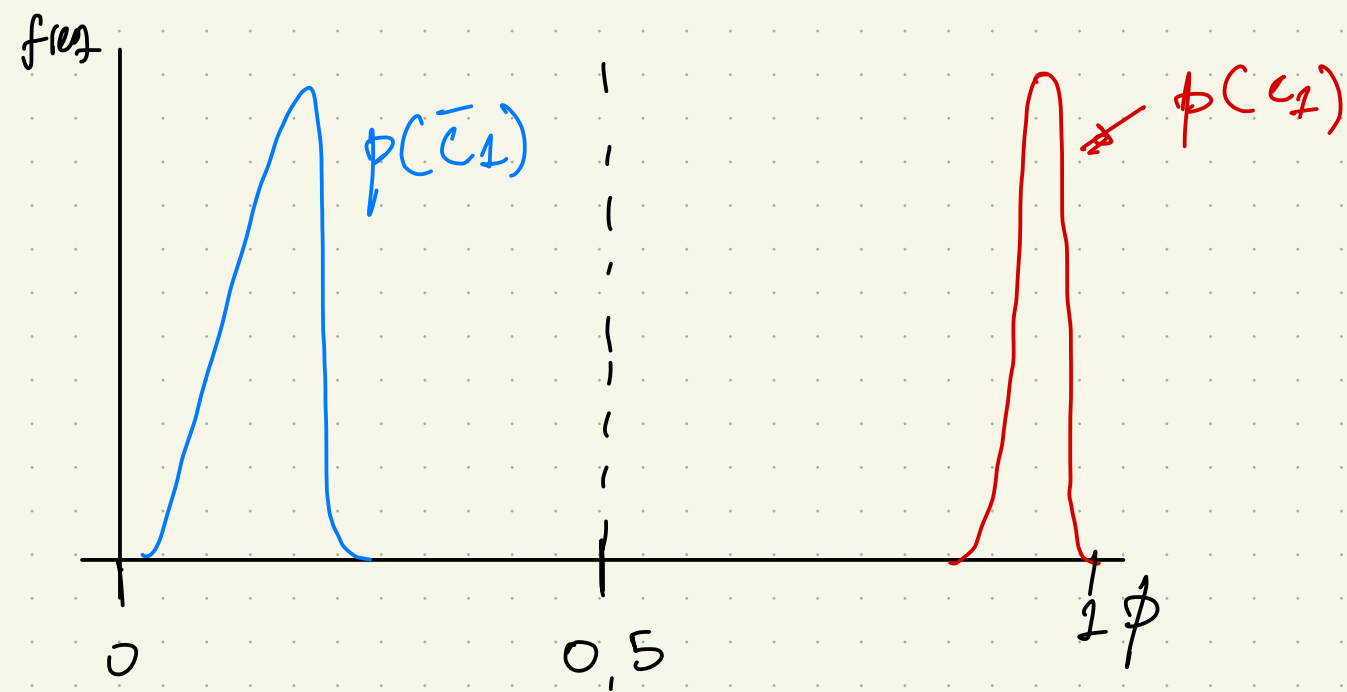
1: Verde, 0: Rojo

## Regresión Logística

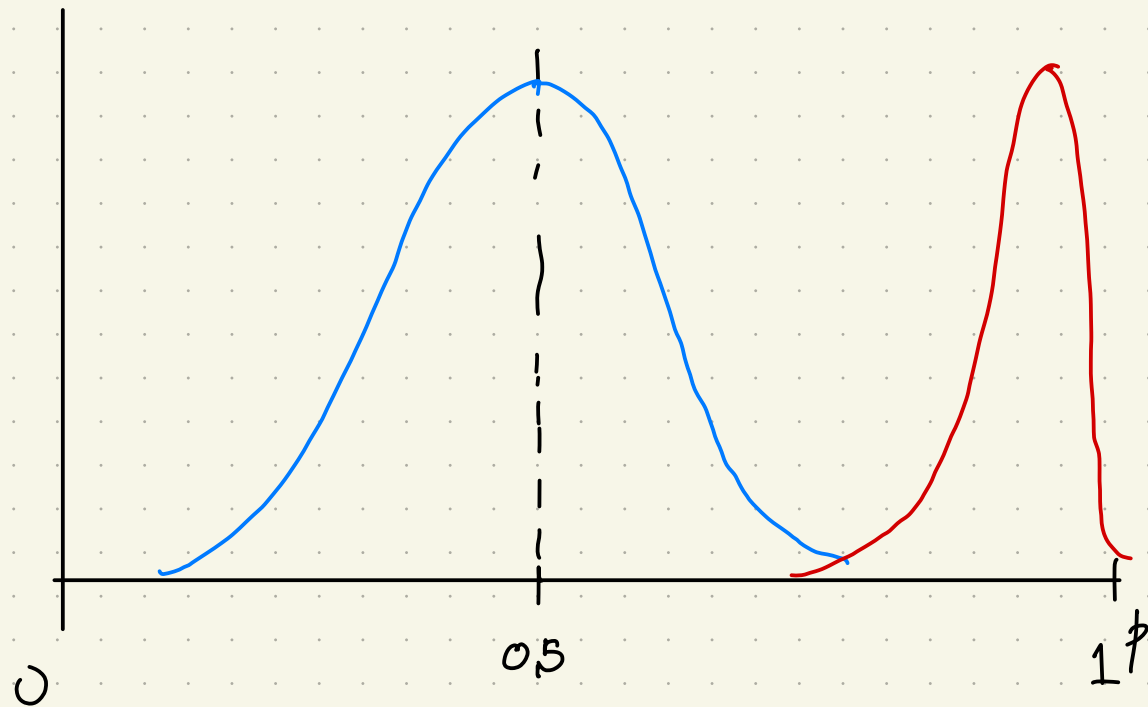


## Regresión Logística





$$\hat{y} = \begin{cases} 1 & \text{si } p \geq 0,5 \\ 0 & \text{o.w} \end{cases}$$



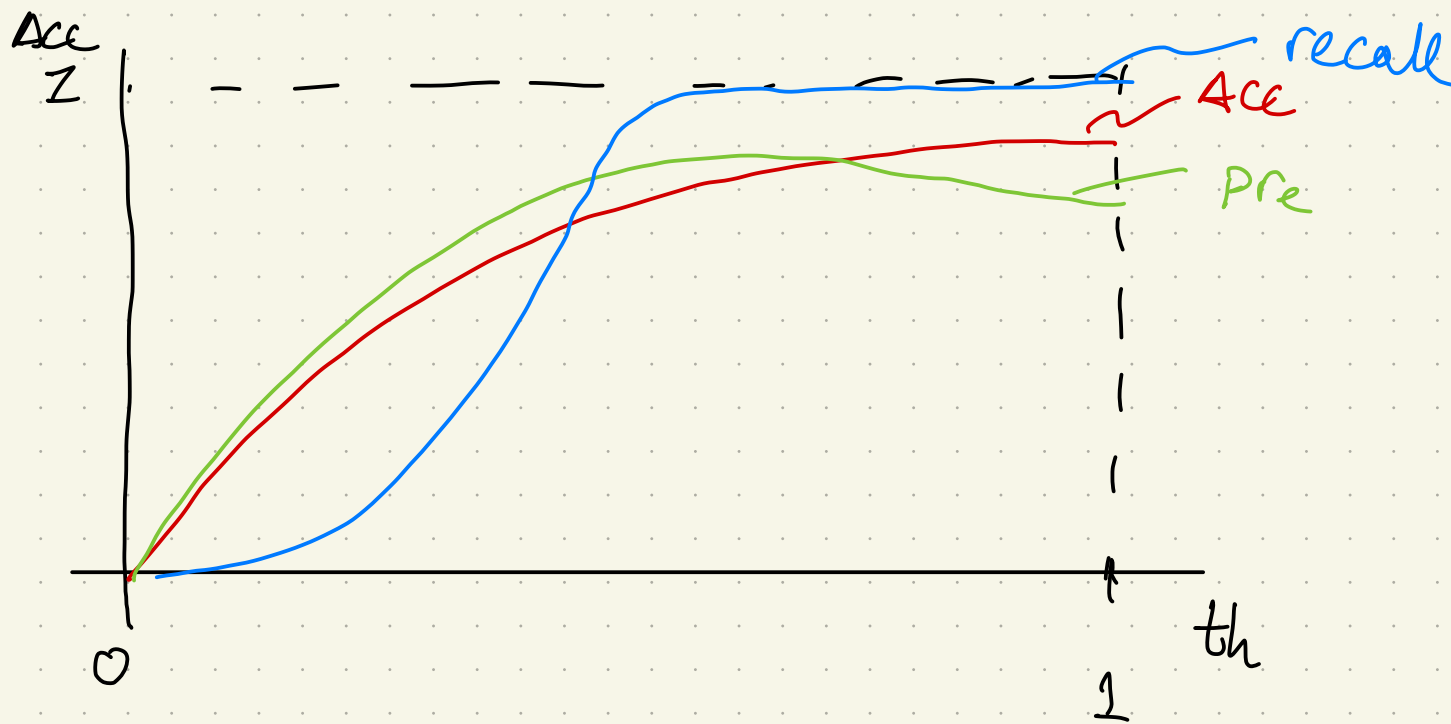
$$y = \begin{cases} 1 & \text{si } p > th \\ 0 & \text{o.w} \end{cases}$$

donde  $th$  lo voy a elegir  
de acuerdo a mis métricas

precisión ( $th$ ) acc( $th$ )

recall ( $th$ )





buscamos un threshold que maximice por ej.  $F_1$ .  
 y teniendo fijo el  $th$ . calculo la matriz de confusión final.

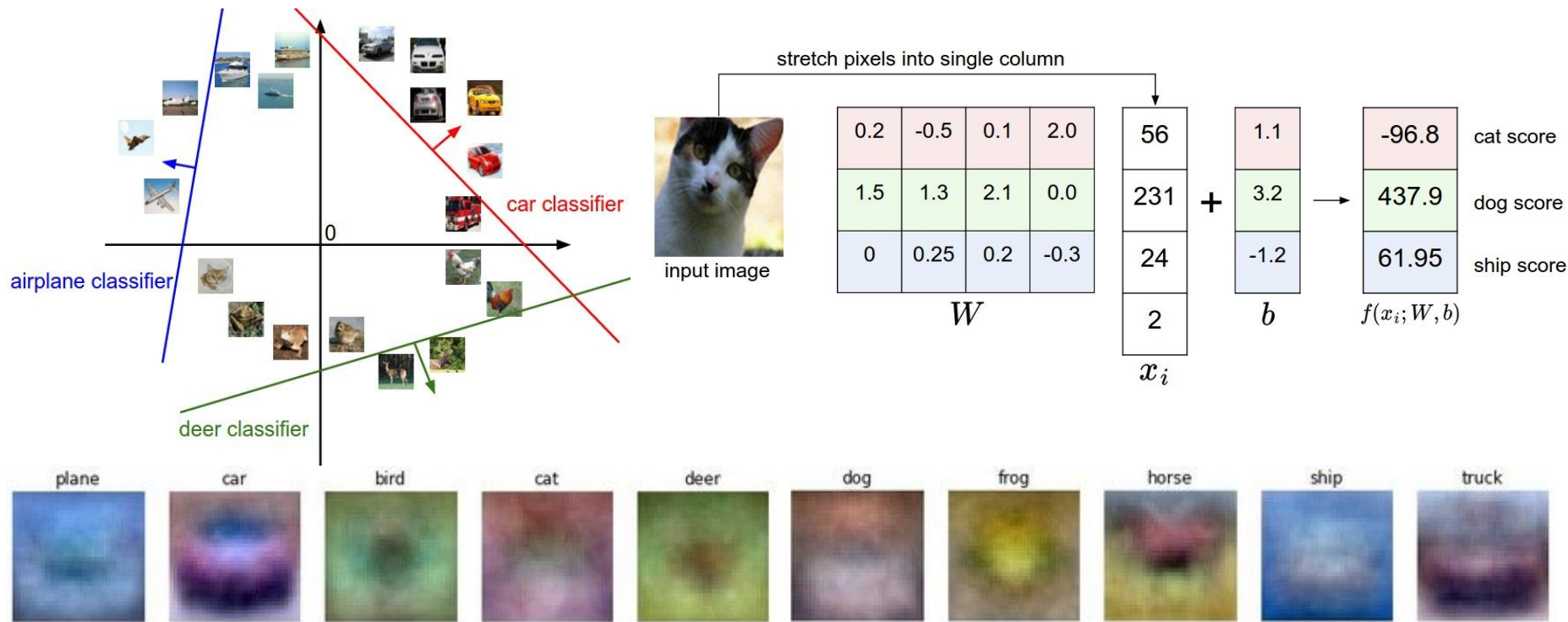
nota: cuando los prob  $p(y=0)$  y  $p(y=1)$  tienen mucho solapamiento. A veces conviene tener un tercer estado,

$$y \in \{True, False, NS/NC\}$$

## Clasificación multiclase

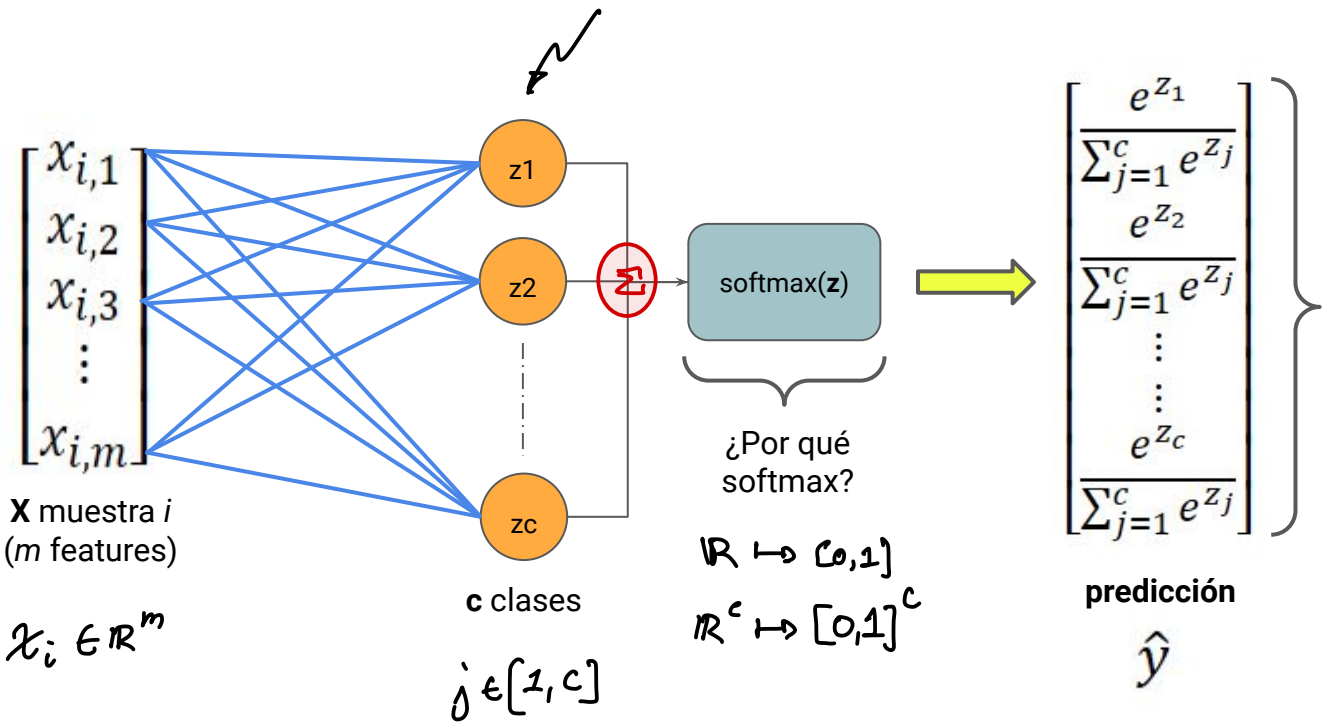
$$f = \Pi_{(y=\pi_i)}$$

## Clasificación Multiclase - Motivación



## Softmax

$z_j \in \mathbb{R} \quad \tilde{z}_j = w_j x_i + b$



¿Función de Costo?  
 $\mathcal{L}(y, \hat{y})$

matriz de conf. multilabel:

<div> <div>pred</div> <div>real</div> </div>	$C_1$	$C_2$	...	$C_c$
$C_1$	$TP_{C_1}$	$M_{C_2C_1}$	...	
$C_2$	$M_{C_1C_2}$	$TP_{C_2}$	...	
$\vdots$	$\vdots$	$\vdots$	$\ddots$	
$C_c$				$TP_{C_c}$

las métricas ahora son por clase:

# Acccy

$$\Delta c_1 c_2$$

、  
-  
、

Accccc

Ahora tenemos el overall del modelo:

ACC<sub>mean</sub> → principio

$\Delta C_{macro} \rightarrow n$  peso do por classe

Acc micro  $\rightarrow$  " " " soporte

## Softmax

$j$ : es la clase

$C$  clases  $\Rightarrow j \in \{1, \dots, C\}$

$$\text{odds} = \frac{P(A)}{P(\bar{A})}$$

$$P(y_i | x_i; W) = \frac{e^{f_{y_i}}}{\sum_j e^{f_j}}$$

$$\frac{e^{f_{y_i}}}{\sum_j e^{f_j}} = \frac{C e^{f_{y_i}}}{C \sum_j e^{f_j}} = \frac{e^{f_{y_i} + \log C}}{\sum_j e^{f_j + \log C}}$$

$q(x)$

$$H(p, q) = - \sum_x p(x) \log q(x)$$

1

[Softmax Forma Gráfica](#)

2

[Softmax Visualización 3D](#)

Los pesos  $w_{ij}$  se obtienen numéricamente con GD y amigos.

## Softmax

### Derivación Softmax

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

$$\frac{\partial p_i}{\partial z_k} = \frac{\partial \frac{e^{z_i}}{\sum_{j=1}^C e^{z_j}}}{\partial z_k}$$

$$\frac{\partial p_i}{\partial z_k} = p_i(\delta_{ik} - p_k) \quad \delta_{ik} = \begin{cases} 1, i = k \\ 0, i \neq k \end{cases}$$

### Derivación Cross-Entropy

$$\begin{aligned} L &= - \sum_i y_i \log(p_i) \\ \frac{\partial L}{\partial z_i} &= - \sum_j y_j \frac{\partial \log(p_j)}{\partial z_i} \\ &= - \sum_j y_j \frac{\partial \log(p_j)}{\partial p_j} \times \frac{\partial p_j}{\partial z_i} \\ &= - \sum_j y_j \frac{1}{p_j} \times \frac{\partial p_j}{\partial z_i} \end{aligned}$$

$$\begin{aligned} \frac{\partial L}{\partial z_i} &= -y_i(1 - p_i) - \sum_{j \neq i} y_j \frac{1}{p_j} (-p_j \cdot p_i) \\ &= -y_i(1 - p_i) + \sum_{j \neq i} y_j \cdot p_i \\ &= p_i \left( y_i + \sum_{j \neq i} y_j \right) - y_i \end{aligned}$$

$$\frac{\partial L}{\partial z_i} = p_i - y_i$$

$$\frac{\partial z_i}{\partial W} = x_i$$

$$\frac{\partial L}{\partial W} = \sum_{i=1}^N (p_i - y_i) x_i$$

Usar gradiente descendente para actualizar W !!!



## Bibliografía

- The Elements of Statistical Learning | Trevor Hastie | Springer
- An Introduction to Statistical Learning | Gareth James | Springer
- Deep Learning | Ian Goodfellow | <https://www.deeplearningbook.org/>
- Mathematics for Machine Learning | Deisenroth, Faisal, Ong
- Artificial Intelligence, A Modern Approach | Stuart J. Russell, Peter Norvig
- Understanding binary cross-entropy: a visual explanation | Daniel Godoy
- Visual Information Theory | [Link](#)
- <https://cs231n.github.io/>
- Classification and Loss Evaluation-Softmax and Cross Entropy Loss | Paras Dahal