

PROBLEMAS DE CLASIFICACIÓN

CLASIFICACIÓN: Cuando aparece una nueva observación, identificar a qué categoría pertenece.

Ejemplo: Nos llega un email. Identificar si es spam o no es spam.

Dada una nueva observación, nos preguntamos si pertenece (1) o no pertenece (0) a una categoría dada.

La predicción de si pertenece o no se da con una probabilidad

Posibles errores:

Matriz de confusión:

		CONDICIÓN REAL	
		1	0
CONDICIÓN PREDICHA	1	Verdadero positivo (TP)	Falsa positivo (FP) [ERROR TIPO I]
	0	Falsa negativo (FN) [ERROR TIPO II]	Verdadero negativo (TN)

Si hacemos predicciones sobre una muestra de N elementos podemos obtener varios parámetros:

Precisión (Accuracy, ACC)

$$ACC = \frac{TP + TN}{N}$$

Tasa de detección (TPR)

$$TPR = \frac{TP}{TP + FN}$$

Tasa de verdaderos negativos (TNR)

$$TNR = \frac{TN}{TN + FP}$$

Tasa de falsa alarma (FPR)

$$FPR = \frac{FP}{FP + TN}$$

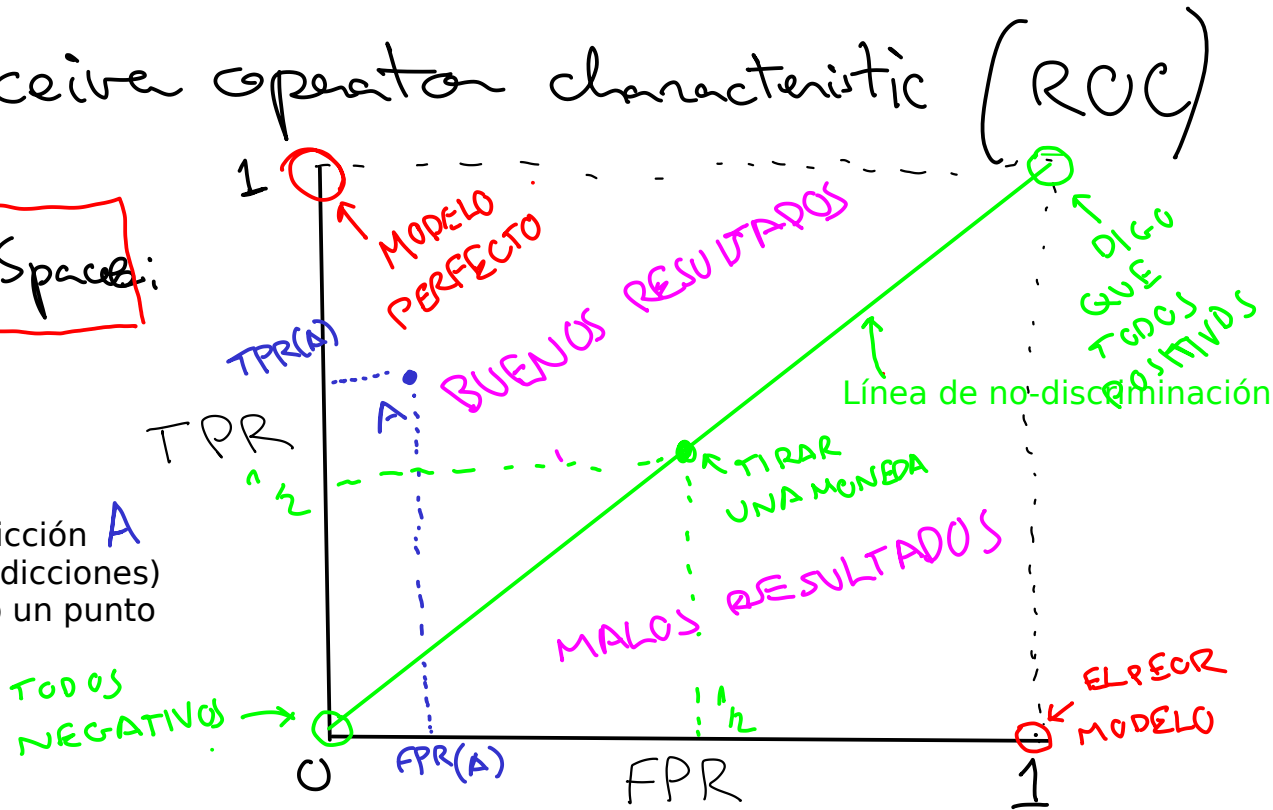
		True condition			
		Condition positive	Condition negative		
Predicted condition	Total population			Prevalence = $\frac{\sum \text{Condition positive}}{\sum \text{Total population}}$	Accuracy (ACC) = $\frac{\sum \text{True positive} + \sum \text{True negative}}{\sum \text{Total population}}$
	Predicted condition positive	True positive	False positive, Type I error	Positive predictive value (PPV), Precision = $\frac{\sum \text{True positive}}{\sum \text{Predicted condition positive}}$	False discovery rate (FDR) = $\frac{\sum \text{False positive}}{\sum \text{Predicted condition positive}}$
Predicted condition negative		False negative, Type II error	True negative	False omission rate (FOR) = $\frac{\sum \text{False negative}}{\sum \text{Predicted condition negative}}$	Negative predictive value (NPV) = $\frac{\sum \text{True negative}}{\sum \text{Predicted condition negative}}$
		True positive rate (TPR), Recall, Sensitivity, probability of detection, Power = $\frac{\sum \text{True positive}}{\sum \text{Condition positive}}$	False positive rate (FPR), Fall-out, probability of false alarm = $\frac{\sum \text{False positive}}{\sum \text{Condition negative}}$	Positive likelihood ratio (LR+) = $\frac{TPR}{FPR}$	Diagnostic odds ratio (DOR) = $\frac{LR+}{LR-}$
		False negative rate (FNR), Miss rate = $\frac{\sum \text{False negative}}{\sum \text{Condition positive}}$	Specificity (SPC), Selectivity, True negative rate (TNR) = $\frac{\sum \text{True negative}}{\sum \text{Condition negative}}$	Negative likelihood ratio (LR-) = $\frac{FNR}{TNR}$	

ROC space, edit 1

Receiver operator characteristic (ROC)

ROC Space:

Cada posible predicción A (o conjunto de predicciones) puede verse como un punto en el ROC space.



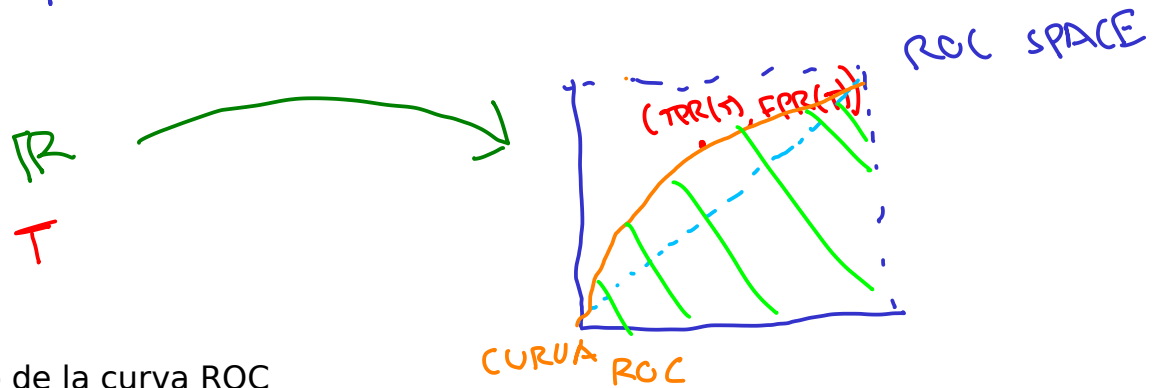
Curva ROC

En general las predicciones se dan como una variable aleatoria continua X de forma que

$$TPR(T) = \int_T^{\infty} f_1(x) dx$$

$$FPR(T) = \int_T^{\infty} f_0(x) dx$$

$$\begin{cases} X > T \rightarrow 1 & \leftarrow f_1 \\ X \leq T \rightarrow 0 & \leftarrow f_0 \end{cases}$$



El área debajo de la curva ROC informa de lo bueno que es el modelo predictivo:

- A > 0.5 BIEN
- A = 0.5 RANDOM
- A < 0.5 MAL

REGRESIÓN LOGÍSTICA

Odds ratio

Ejemplo:

Porcentaje de participación:

$$\frac{467}{2484} = 18,8\%$$

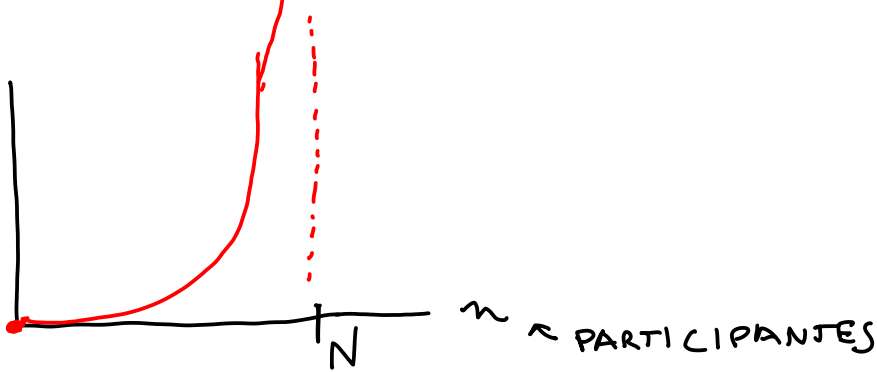
Odd (los que sí, frente a los que no):

Participación en una huelga por sexo:

	HOMBRE	MUJER	TOTAL
NO	960	1057	2017
SÍ	261	206	467
	1221	1263	2484

$$\frac{467}{2017} = 0,232$$

$$\frac{m}{N-m}$$



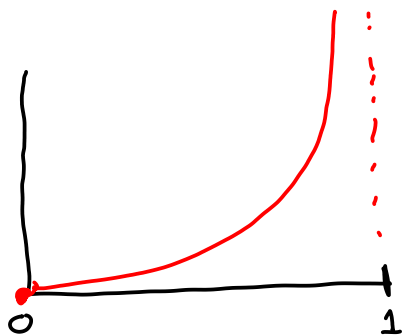
$$\frac{m}{N-m}$$

PROB. PARTICIPA

\downarrow
 $\sqrt{N}, P = \frac{m}{N}$

$$\frac{P}{1-P}$$

$$\frac{P}{1-P}$$



Podemos sacar los ODDs por separado:

ODD(participa/no participa) HOMBRE 0,272

MUJER 0,195

ODD RATIO:

$$OR = \frac{ODD_A}{ODD_B} = \frac{\frac{P_A}{1-P_A}}{\frac{P_B}{1-P_B}}$$

Para estudiar los odds, nos interesa considerar la función LOGIT

$$ODD(p) = \frac{P}{1-P}$$

$$Logit(p) = \log(ODD) = \log\left(\frac{P}{1-P}\right)$$

$$ODD = P/q \rightarrow \logit = \ln(P/q)$$

$$ODD = q/p \rightarrow \logit = \ln(q/p) = -\ln(P/q)$$



LOGIT



Volviendo al caso anterior:

	HOMBRE	MUJER	TOTAL
NO	0,7862	0,8368	0,8119
SÍ	0,2137	0,1631	0,1884
	1	1	1
ODD	0,272	0,195	
LOGIT	-1,30	-1,63	

Regresión logística

Training set

Datos, que pertenecen o no pertenecen a una variable categórica

$$y = \begin{cases} 1 & \text{si pertenece} \\ 0 & \text{si no pertenece} \end{cases}$$

X = Variables que vamos a predecir

Idea: Hacer un ajuste lineal a Logit (p)

$$\text{logit}(p) = \alpha + \beta x$$

$$p = \frac{1}{1 + e^{-(\alpha + \beta x)}}$$

En reg. lineal

$$y = \alpha + \beta x$$

En reg. logística

$$p = \frac{1}{1 + e^{-(\alpha + \beta x)}}$$

Test

↓
Overtax

p ← PROB. DE QUE
PERTENEZCA

función logística:

$$\sigma(t) = \frac{1}{1 + e^{-t}}$$

En el ejemplo anterior:

Mi target es estimar p , la probabilidad de hacer huelga.

Mi feature es

$$X = \begin{cases} 0 & \text{si } \underline{\text{hombre}} \\ 1 & \text{si } \underline{\text{mujer}} \end{cases}$$

$$\text{logit}(p) = \alpha + \beta x$$

$$X=0 \rightarrow \underset{\substack{|| \\ -1,30}}{\text{logit(HOMBRE)}} = \alpha \Rightarrow \alpha = -1,30$$

$$X=1 \rightarrow \underset{\substack{|| \\ -1,63}}{\text{logit(MUJER)}} = \alpha + \beta \Rightarrow \beta = -0,33$$

$$P = \frac{1}{1 + e^{1,30 + 0,33x}}$$

$$P(0) = 0,214$$

$$P(1) = 0,163$$

OBSERVACIÓN:

$$e^{\alpha} = 0,272 = \text{ODD}_{\text{HOMBRES}}$$

$$e^{\alpha + \beta} = \text{ODD}_{\text{MUJERES}}$$

$$e^{\beta} = \text{ODD}_M / \text{ODD}_H = \text{OR}$$

EJERCICIO:

Participación en la huelga por nivel de estudios

	Básicos	Medios	Universitarios
NO	1116	554	344
SI	138	182	145
	1254	736	489

Se pide: Hallar los ODDS, los logit y los parámetros de la regresión logística.

	x_1	x_2
BÁSICOS	0	0
MEDIOS	1	0
UNIVERSITARIOS	0	1

El ajuste es a:

$$p = \frac{1}{1 + e^{(-\alpha + \beta_1 x_1 + \beta_2 x_2)}}$$

Función de coste:

La regresión logística tiene la forma:

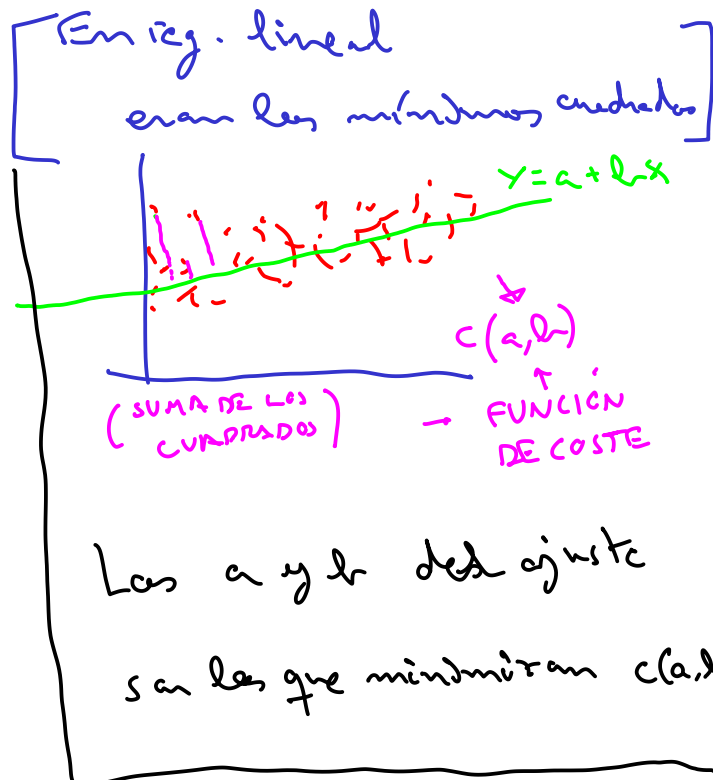
$$p = \frac{1}{1 + e^{-\vec{\theta} \cdot \vec{x} - \vec{\eta}}}$$

$$\vec{\theta} = (\theta_1, \dots, \theta_n)$$

$$\vec{x} = (x_1, \dots, x_n)$$

$$\vec{\theta} \cdot \vec{x} = \theta_1 x_1 + \dots + \theta_n x_n$$

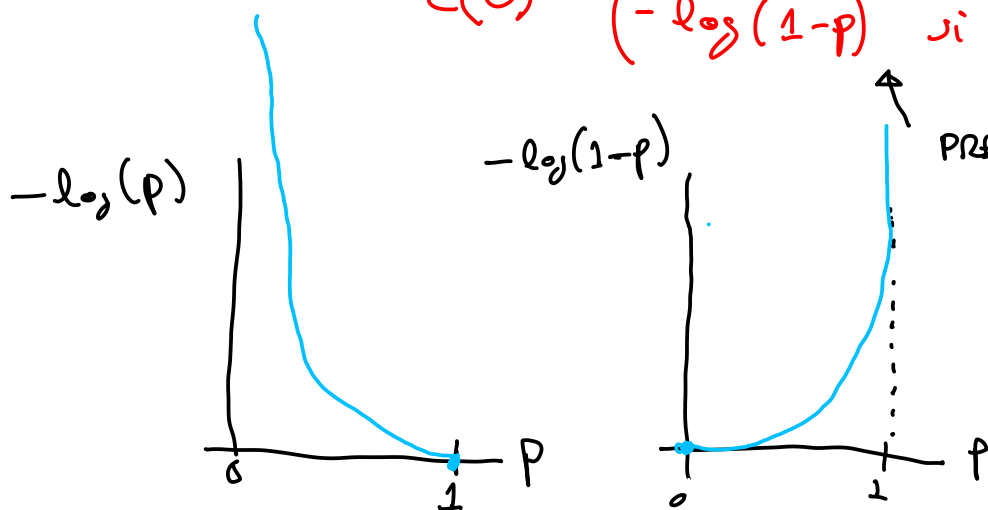
$$\vec{\eta} = (\eta_1, \dots, \eta_n)$$



• Para una: θ, x

PREMIO LAS
PROBABILIDADES ALTAS

$$c(\theta) = \begin{cases} -\log p & \text{si } y=1 \\ -\log(1-p) & \text{si } y=0 \end{cases}$$



PREMIO LAS PROBABILIDADES BAJAS

Para $\vec{x}, \vec{\theta}$, promedio:

$$\vec{y} = (y_1, \dots, y_m)$$

$$J(\vec{\theta}) = -\frac{1}{m} \sum_{i=1}^m \left[y_i \log(p_i) + (1-y_i) \log(1-p_i) \right]$$

↑
ES CONVEXA

→ Los algoritmos de optimización
(p.ej. el "gradient descent") hallan
el MÍNIMO GLOBAL ..