

# Módulo 3. IA y grandes volúmenes de datos

#2. El problema de clasificación

# Clasificación binaria

- Disponemos de  $N$  pares de entrenamiento (observaciones)

$$\{(x_i, y_i)\}_{i=1}^N = \{(x_1, y_1), \dots, (x_N, y_N)\}$$

con  $x_i \in \mathbb{R}^n$ ,  $y_i \in \{-1, +1\}$ .

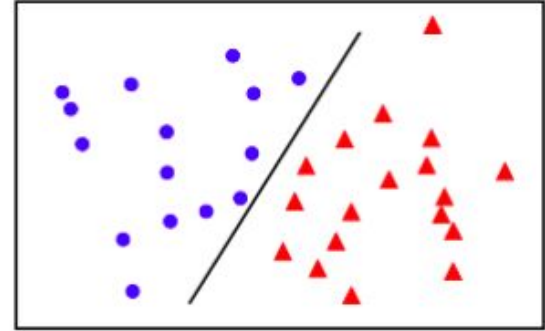
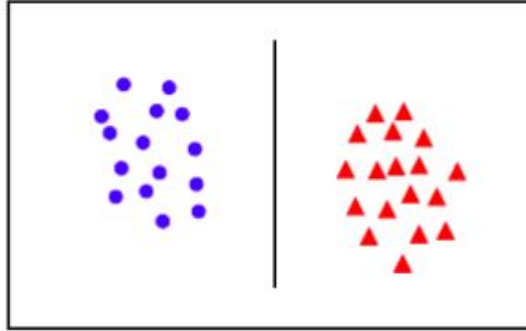
- Aprender una  $f(x)$  tal que

$$f(\mathbf{x}_i) \begin{cases} \geq 0 & y_i = +1 \\ < 0 & y_i = -1 \end{cases}$$

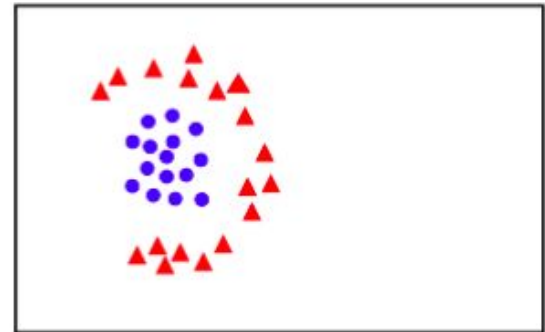
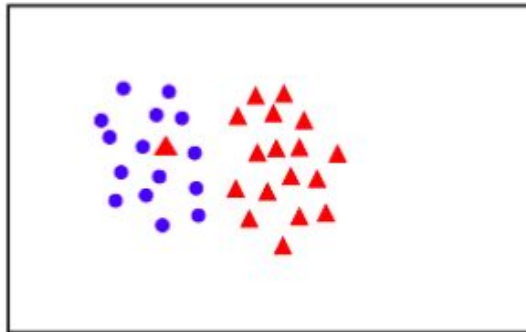
es decir:  $y_i f(x_i) > 0$  para una clasificación correcta.

# Separabilidad lineal

linealmente  
separable



**no**  
linealmente  
separable



# Clasificadores lineales

- La entrada es un vector  $\mathbf{x}_i$  de dimensionalidad  $n$
- La salida es una etiqueta  $y_i \in \{-1, +1\}$
- Clasificador = función de predicción + función de decisión

$$g(f(\mathbf{x})) \rightarrow \{-1, +1\}$$

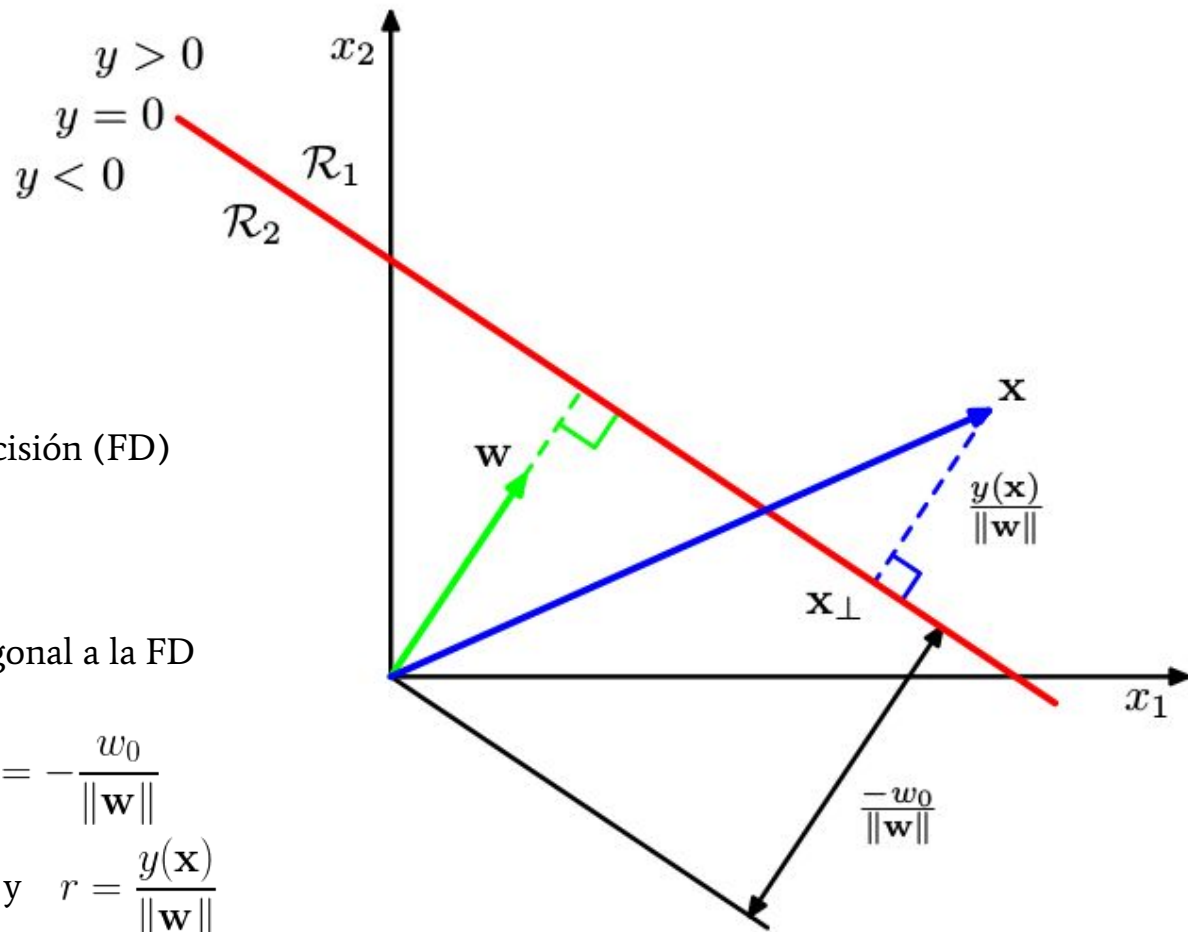
- Función de predicción **lineal**

$$f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$$

- Función de decisión

$$g(z) = \text{sign}(z)$$

$$g(f(\mathbf{x})) = \text{sign}(\mathbf{w}^T \mathbf{x} + w_0)$$



$$y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$$

$\mathbf{x}_A, \mathbf{x}_B$  : puntos en la frontera de decisión (FD)

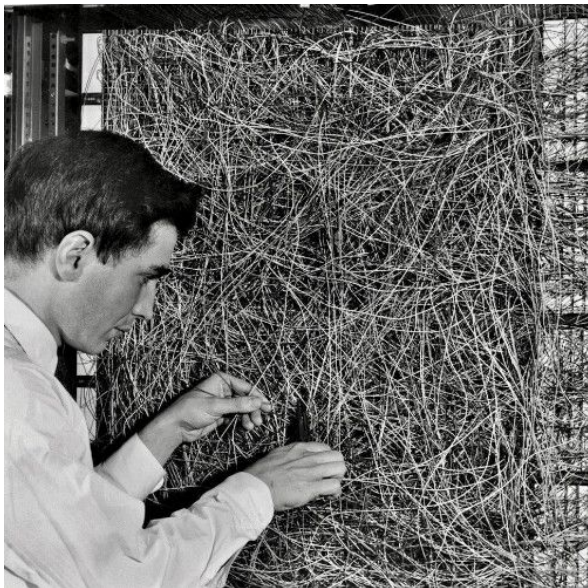
$$y(\mathbf{x}_A) = y(\mathbf{x}_B) = 0$$

$\mathbf{w}^T(\mathbf{x}_A - \mathbf{x}_B) = 0 \implies \mathbf{w}$  es ortogonal a la FD

$$\text{si } \mathbf{x} \text{ en la FD, } y(\mathbf{x}) = 0 \rightarrow \frac{\mathbf{w}^T \mathbf{x}}{\|\mathbf{w}\|} = -\frac{w_0}{\|\mathbf{w}\|}$$

$$\text{si } \mathbf{x} = \mathbf{x}_\perp + r \frac{\mathbf{w}}{\|\mathbf{w}\|} \rightarrow y(\mathbf{x}_\perp) = 0 \text{ y } r = \frac{y(\mathbf{x})}{\|\mathbf{w}\|}$$

# El algoritmo del "perceptrón"



- Propuesto por Roseblatt en 1958
- El objetivo es encontrar un hiperplano de separación. **Si los datos son linealmente separables, lo encuentra.**
- Es un algoritmo *online* (procesa un ejemplo a la vez)
- Muchas variantes ...

# El algoritmo del "perceptrón"

Entrada:

- una secuencia de pares de entrenamiento  $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2) \dots$
- Una tasa de aprendizaje  $r$  (número pequeño y menor a 1)

Algoritmo:

- Inicializar  $\mathbf{w}^{(0)} \in \mathbb{R}^n$
- Para cada ejemplo  $(\mathbf{x}_i, y_i)$ 
  - Predecir  $y_i' = \text{sign}(\mathbf{w}^T \mathbf{x}_i)$
  - Si  $y_i' \neq y_i$ :  
$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} + r (y_i \mathbf{x}_i)$$

# El algoritmo del "perceptrón"

Entrada:

- una secuencia de pares de entrenamiento  $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2) \dots$
- Una tasa de aprendizaje  $r$  (número pequeño y menor a 1)

Algoritmo:

- Inicializar  $\mathbf{w}^{(0)} \in \mathbb{R}^n$
- Para cada ejemplo  $(\mathbf{x}_i, y_i)$ 
  - Predecir  $y_i' = \text{sign}(\mathbf{w}^T \mathbf{x}_i)$
  - Si  $y_i' \neq y_i$ :  
$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} + r (y_i \mathbf{x}_i)$$

Actualiza solo cuando comete un error

Error en positivos:

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} + r \mathbf{x}_i$$

Error en negativos:

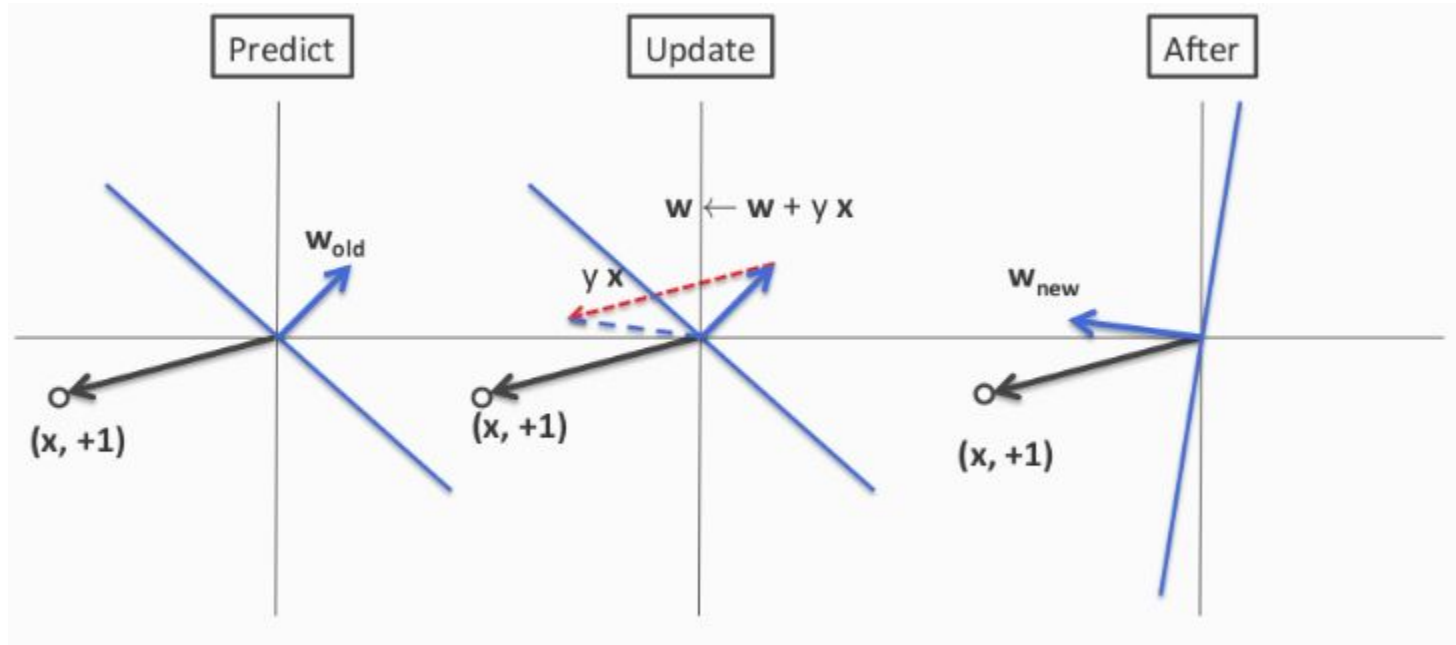
$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - r \mathbf{x}_i$$

Si  $y_i \mathbf{w}^T \mathbf{x}_i \leq 0 \rightarrow \text{error}$



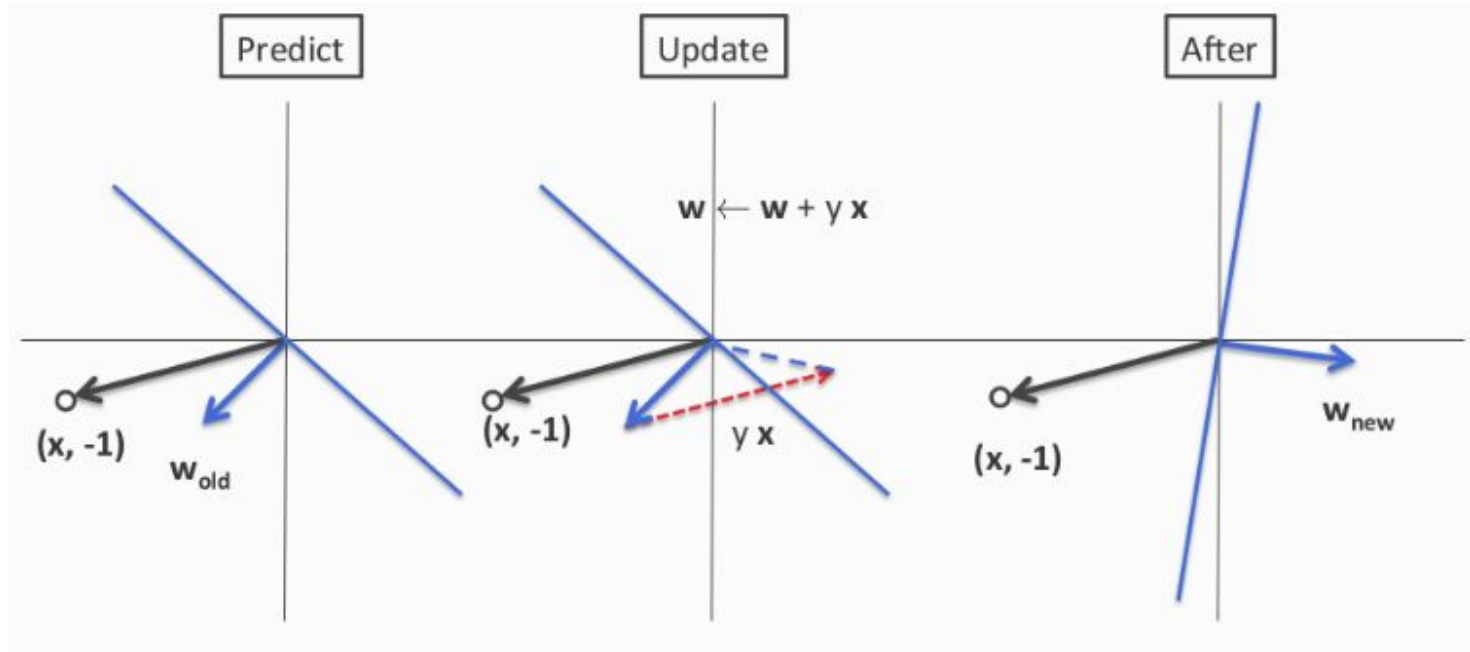
# Dinámica de actualización

Error en ejemplo **positivo**:



# Dinámica de actualización

Error en ejemplo **negativo**:



# El algoritmo estándar

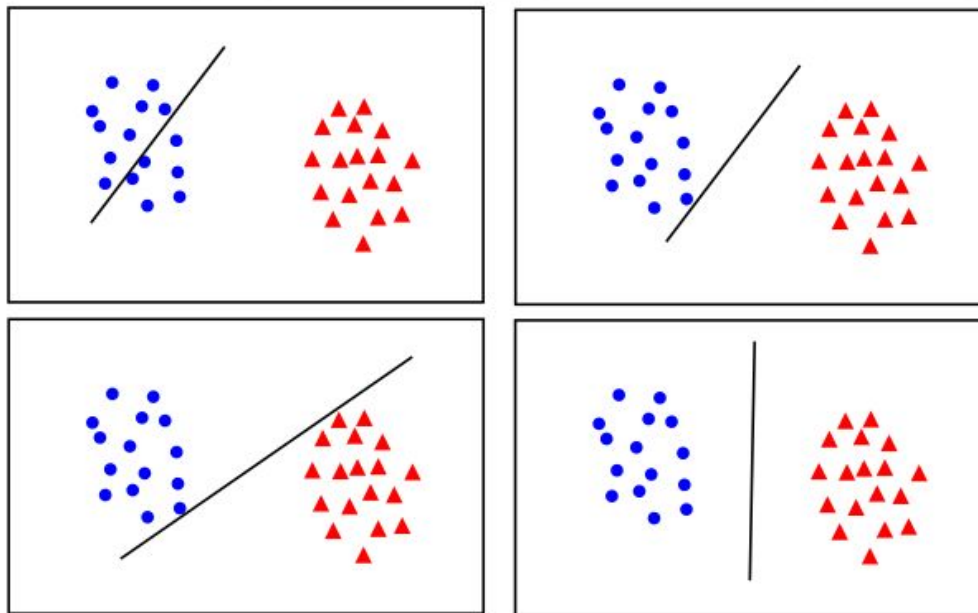
Dado un conjunto  $D=\{(\mathbf{x}_i, y_i), i=1, \dots, N\}$ ,  $y_i \in \{-1, +1\}$ , tasa de entrenamiento  $r$  y número de épocas  $T$

1. Inicializar  $\mathbf{w}^{(0)}$
2. Para época  $t=1, \dots, T$ 
  - a. *barajar* el conjunto de entrenamiento  $D$
  - b. Para cada muestra de entrenamiento  $(\mathbf{x}_i, y_i) \in D$ 
    - si  $y_i \mathbf{w}^{(t)T} \mathbf{x}_i \leq 0$ , actualizar  $\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} + r (y_i \mathbf{x}_i)$
3. Retornar  $\mathbf{w}^{(T)}$

$r, T$ : hiperparámetros

**Predicción:**  $\text{sgn}(\mathbf{w}^T \mathbf{x})$

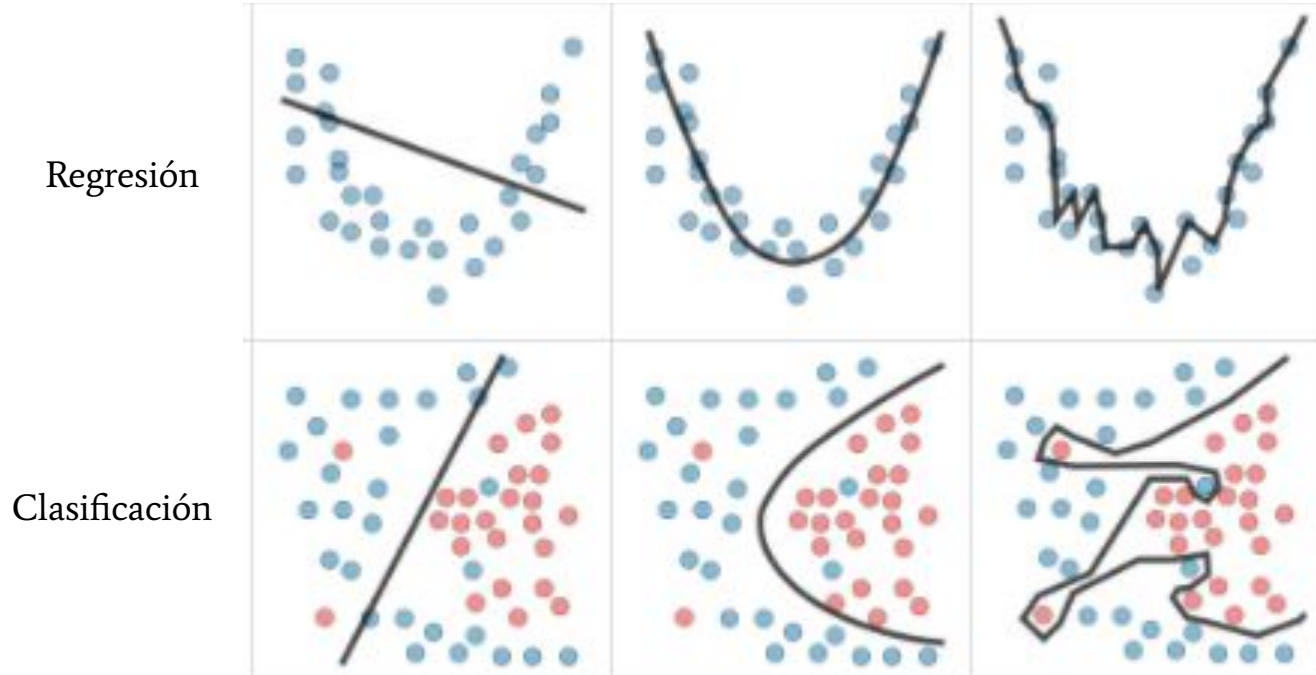
¿Cuál es el mejor  $w$ ?



Solución de **margen máximo**: el hiperplano más estable ante perturbaciones de la entrada

# Generalización en clasificación

- Complejidad del modelo  $\Leftrightarrow$  complejidad de la frontera de decisión



# Regresión logística

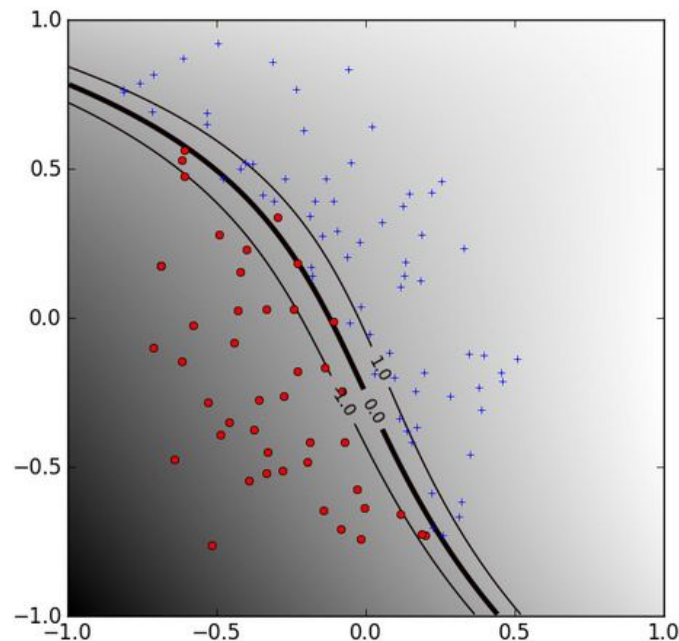
# Clasificación basada en probabilidades

- Objetivo: dar una estimación de probabilidad de que una instancia  $x$  sea de una clase  $y$ , es decir,  $p(y|x)$

- Recordar:

$$0 \leq p(\text{evento}) \leq 1$$

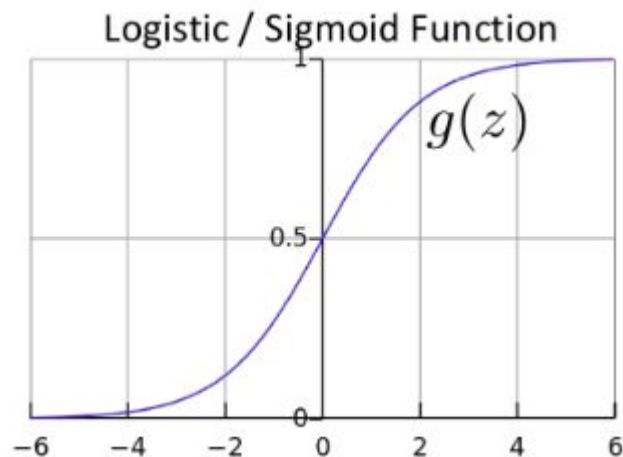
$$p(\text{evento}) + p(\neg \text{evento}) = 1$$



# Regresión logística

- Aproximación probabilística al problema de clasificación
- La función de predicción  $h_w(x)$  debe dar una aproximación de  $p(y=1|x,w)$
- $0 \leq h_w(x) \leq 1$

$$h_w(x) = g(w^T x) = \frac{1}{1 + \exp(-w^T x)}$$





# Regresión logística

- Datos  $\left\{ \left( \mathbf{x}^{(1)}, y^{(1)} \right), \left( \mathbf{x}^{(2)}, y^{(2)} \right), \dots, \left( \mathbf{x}^{(n)}, y^{(n)} \right) \right\}$   
donde  $\mathbf{x}^{(i)} \in \mathbb{R}^d$ ,  $y^{(i)} \in \{0, 1\}$

- Modelo:  $h_{\boldsymbol{\theta}}(\mathbf{x}) = g(\boldsymbol{\theta}^{\top} \mathbf{x})$

$$g(z) = \frac{1}{1 + e^{-z}}$$

$$\boldsymbol{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_d \end{bmatrix}$$

$$\mathbf{x}^{\top} = \begin{bmatrix} 1 & x_1 & \dots & x_d \end{bmatrix}$$

# Regresión logística. Función de costo

- Conjunto de entrenamiento  $\{(\mathbf{x}^1, y^1), \dots, (\mathbf{x}^N, y^N)\}$ ,  $\mathbf{x} \in R^M$ ,  $y \in \{0, 1\}$
- $y$ : observaciones discretas  $\rightarrow$  muestras de una distribución Bernoulli

$$P(y = 1|\mathbf{x}, \mathbf{w}) = f(\mathbf{x}, \mathbf{w})$$

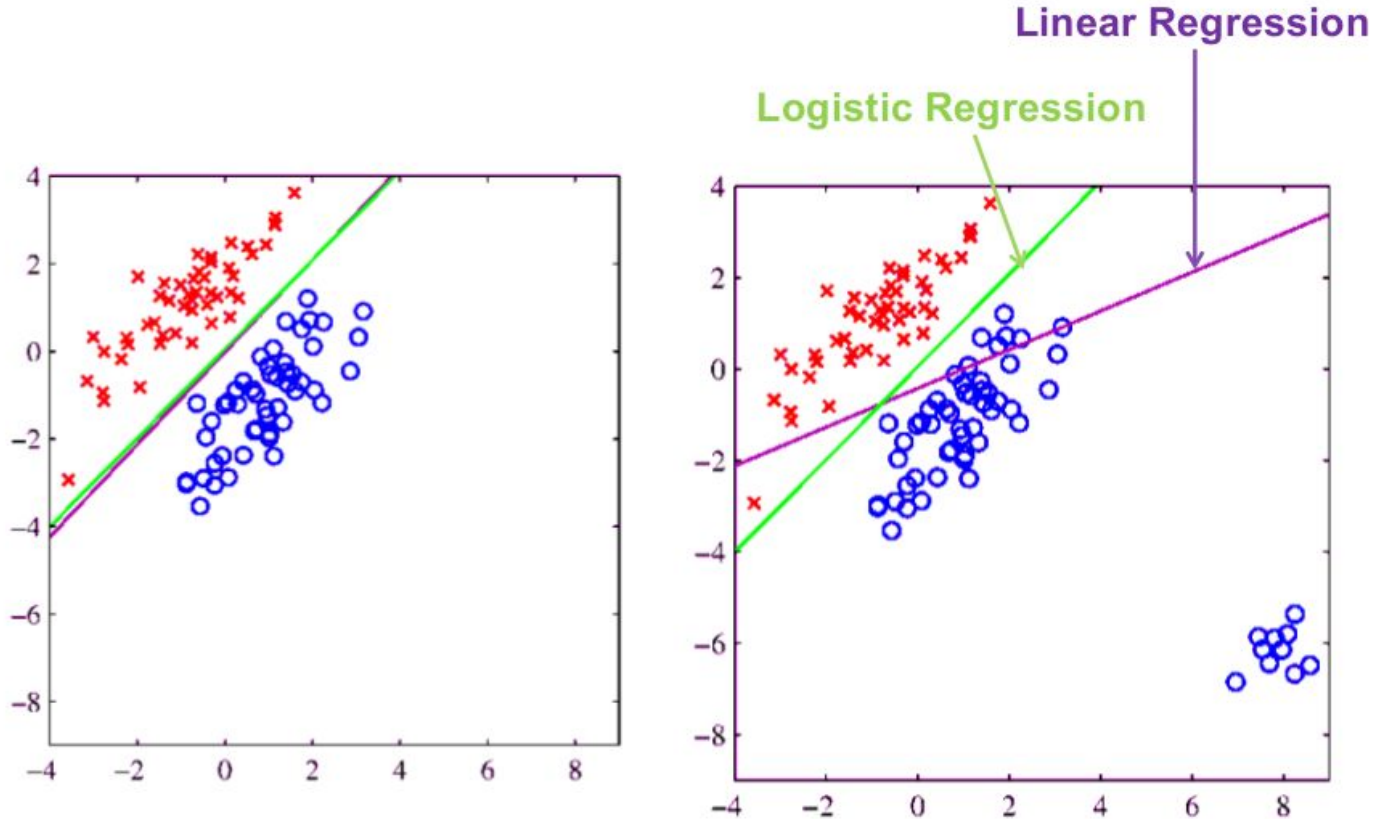
$$P(y = 0|\mathbf{x}, \mathbf{w}) = 1 - f(\mathbf{x}, \mathbf{w})$$

$$P(y|\mathbf{x}) = (f(\mathbf{x}, \mathbf{w}))^y (1 - f(\mathbf{x}, \mathbf{w}))^{1-y}$$

- Encontrar el  $\mathbf{w}$  que maximice la verosimilitud de las etiquetas en el conjunto de entrenamiento

$$\begin{aligned} -L(\mathbf{w}) = C(\mathbf{w}) &= \log P(\mathbf{y}|\mathbf{X}, \mathbf{w}) = \sum_{i=1}^N \log P(y^i|\mathbf{x}^i, \mathbf{w}) \\ &= \sum_i y^i \log f(\mathbf{x}^i, \mathbf{w}) + (1 - y^i) \log(1 - f(\mathbf{x}^i, \mathbf{w})) \end{aligned}$$

# Regresión lineal vs. regresión logística



Problemas multiclase

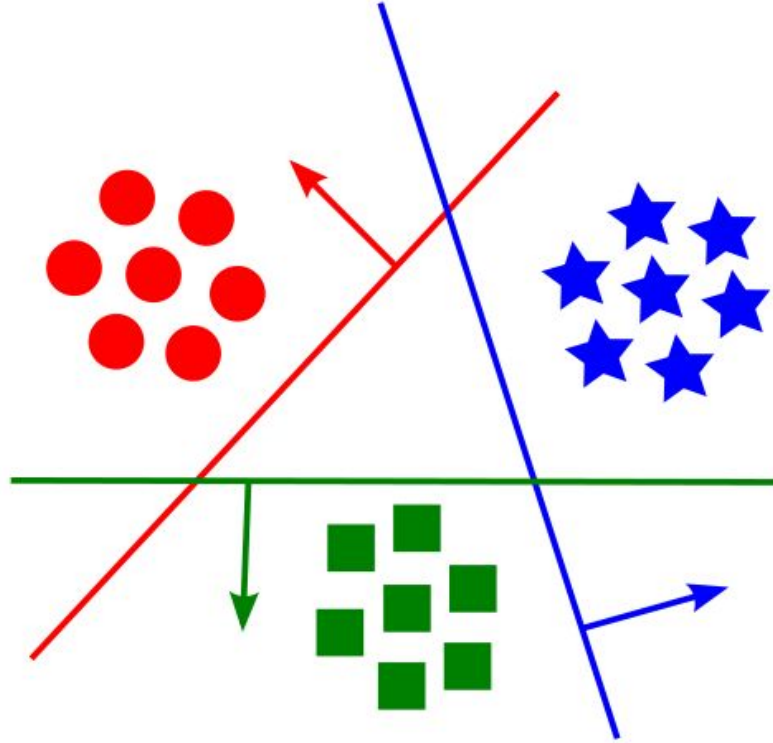
# Clasificación multiclase

- Una muestra puede pertenecer a 1 (o más) de  $K$  clases
  - Datos de entrenamiento  $\{(\mathbf{x}_i, y_i)\}, y_i=1, \dots, K$
- Distintos tipos de problemas:
  - multiclase:  $\mathbf{x}$  pertenece solo a una categoría
  - multietiqueta:  $\mathbf{x}$  puede pertenecer a más de una categoría
- A veces es más fácil descomponer el problema multiclase en una serie de problemas binarios. **Distintas estrategias: OVA, AVA, ...**

# Estrategia uno contra todos (OVA)

- **Asumimos que cada clase es separable del resto**
- Dado un conjunto de entrenamiento  $D=\{(\mathbf{x}_i, y_i)\}$ ,  $y_i=1,\dots,K$ 
  - Descomponer el problema en  $K$  problemas binarios. Para la clase  $k$ , crear un problema tal que:
    - Ejemplos cuya etiqueta es  $y_i=k$  son ejemplos positivos
    - Ejemplos cuya etiqueta es  $y_i \neq k$  son ejemplos negativos
  - Generar  $K$  clasificadores binarios con **función de predicción**  
 $f_k(\mathbf{x})$ ,  $k=1,\dots,K$ .
- Predicción (*winner takes all*):  $k^* = \operatorname{argmax}_k f_k(\mathbf{x})$

## Estrategia uno contra todos (OVA)

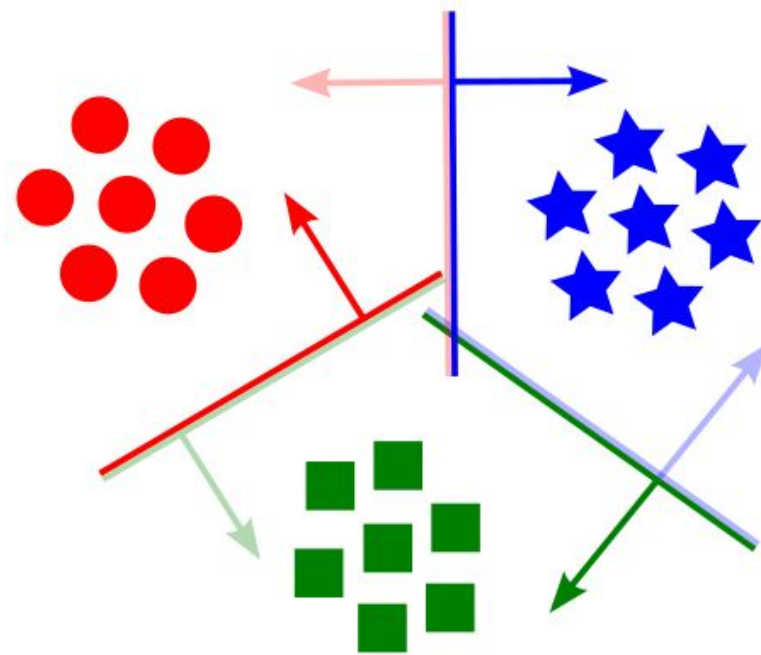


# Estrategia todos contra todos (AVA)

- **Asumimos que cada clase par de clases es separable**
- Dado un conjunto de entrenamiento  $D=\{(\mathbf{x}_i, y_i)\}$ ,  $y_i=1,\dots,K$ 
  - Descomponer el problema en  $K(K-1)/2$  problemas binarios.  
Para el par de clases  $(i, j)$ ,  $i \neq j$ , crear un problema tal que:
    - Ejemplos cuya etiqueta es  $y_i=i$  son ejemplos positivos
    - Ejemplos cuya etiqueta es  $y_i=j$  son ejemplos negativos
  - Generar  $K(K-1)/2$  clasificadores binarios con **función de decisión**  $g_{(i,j)}(\mathbf{x})$
- Predicción (*voting*): cada clase recibe  $K-1$  “votos”



## Estrategia todos contra todos (AVA)



# Regresión logística multiclase

- Para dos clases:

$$h_{\theta}(x) = \frac{1}{1 + \exp(-\theta^T x)} = \frac{\exp(\theta^T x)}{\boxed{1} + \boxed{\exp(\theta^T x)}}$$

Diagram illustrating the weights assigned to the classes in the two-class logistic regression model:

- The term  $1$  is labeled "peso asignado a  $y=0$ ".
- The term  $\exp(\theta^T x)$  is labeled "peso asignado a  $y=1$ ".

- Para  $C$  clases ( $c=1, \dots, C$ ):

$$p(y = c \mid x; \theta_1, \dots, \theta_C) = \frac{\exp(\theta_c^T x)}{\sum_{c=1}^C \exp(\theta_c^T x)}$$

(función **softmax**)

# Métricas en clasificación

# Importancia de las métricas

- La función de costo es solo un *proxy* al problema en el mundo real
- Las métricas ayudan a capturar objetivos reales en forma cuantitativa (no todos los errores son iguales)
- Ayudan a la organización el trabajo de los equipos en función de los requerimientos del problema
- Permiten cuantificar diferencias en:
  - performance deseada vs modelo base
  - performance deseada vs actual
  - evolución en el tiempo
- Deberían ser el objetivo del entrenamiento, pero a veces es difícil.

# Clasificación binaria

- Entrada:  $x$ , salida:  $y$  (valores 0/1 o -1/+1)
- Predicción del modelo:  $\hat{y}=h(x)$
- Dos tipos de modelos:
  - Modelos que predicen directamente una variable categórica (kNN, árboles de decisión)
  - Modelos que predicen un puntaje (*score*) (SVM, regresión logística)
    - Se necesita elegir un umbral (func. de decisión)
    - Nos enfocaremos en esta última clase de modelos. Los anteriores se pueden ver como un caso especial.

# Modelos basados en *scores*

Score = 1



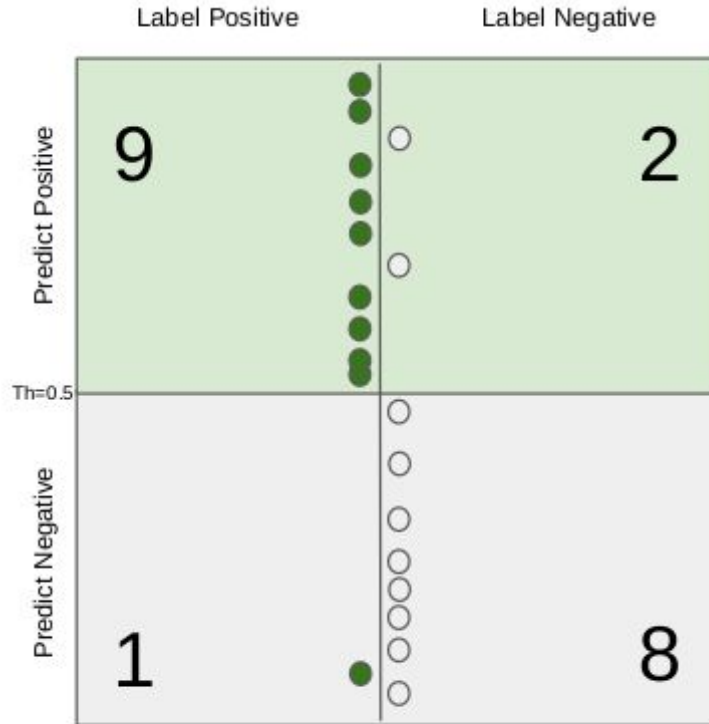
Score = 0

|   |                  |
|---|------------------|
| ● | Positive example |
| ○ | Negative example |

# Umbral → clasificador → métrica puntual



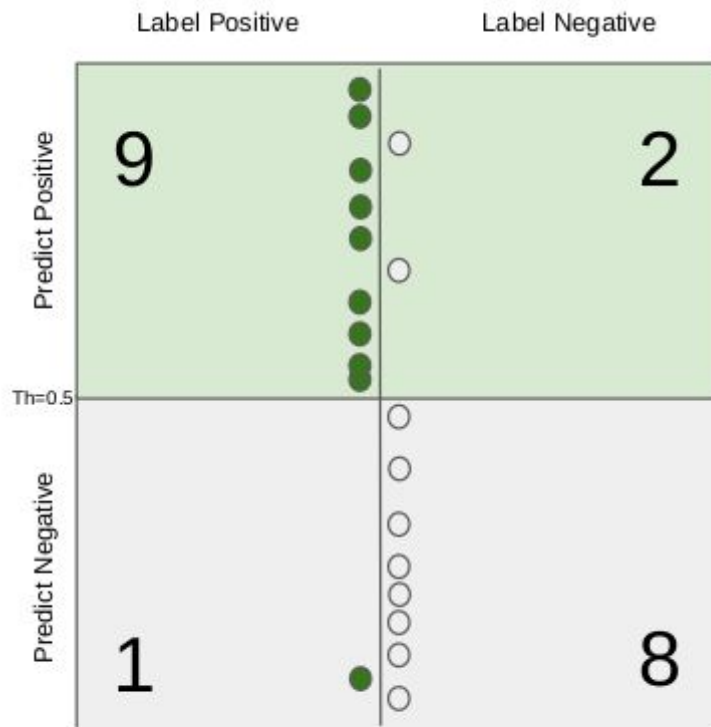
# Matriz de confusión



- la suma total es fija (muestra)
- la suma por columnas es fija (muestras por clase)
- la calidad del modelo y el valor de umbral deciden el agrupamiento de filas
- queremos que los elementos diagonales tengan valores grandes y los no diagonales valores chicos

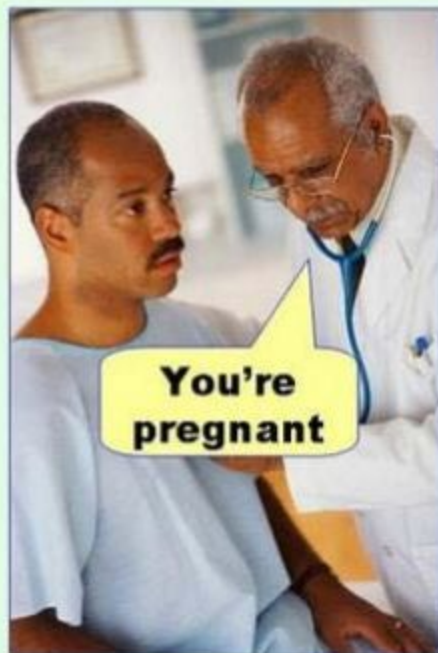


# Matriz de confusión

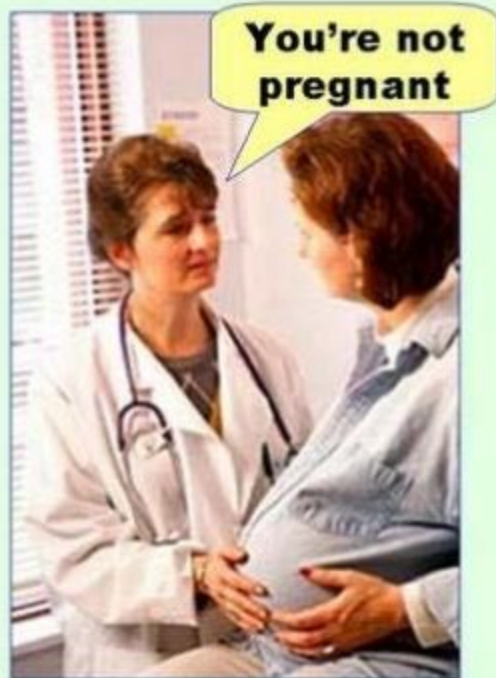


- *true positives* (TP) = 9
- *true negatives* (TN) = 8
- *false positives* (FP) = 2
- *false negatives* (FN) = 1

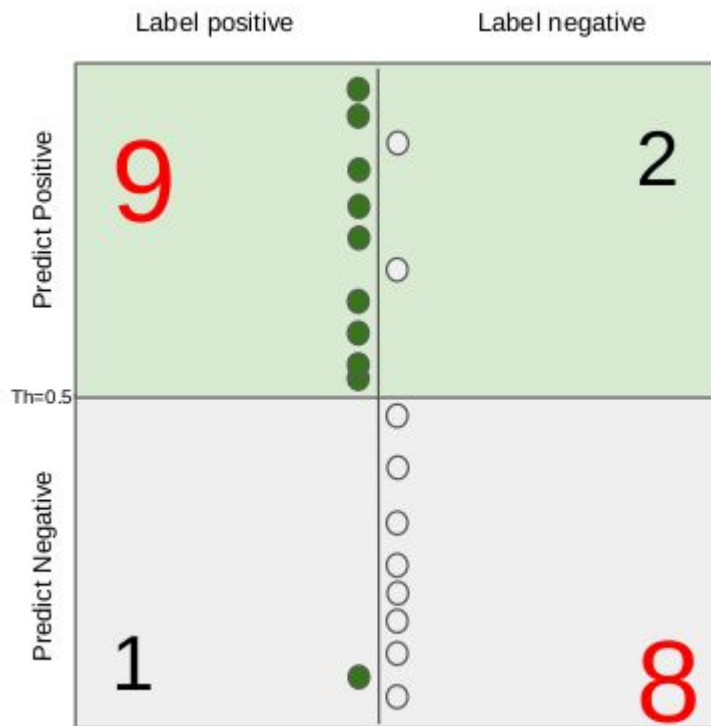
**Type I error**  
(false positive)



**Type II error**  
(false negative)



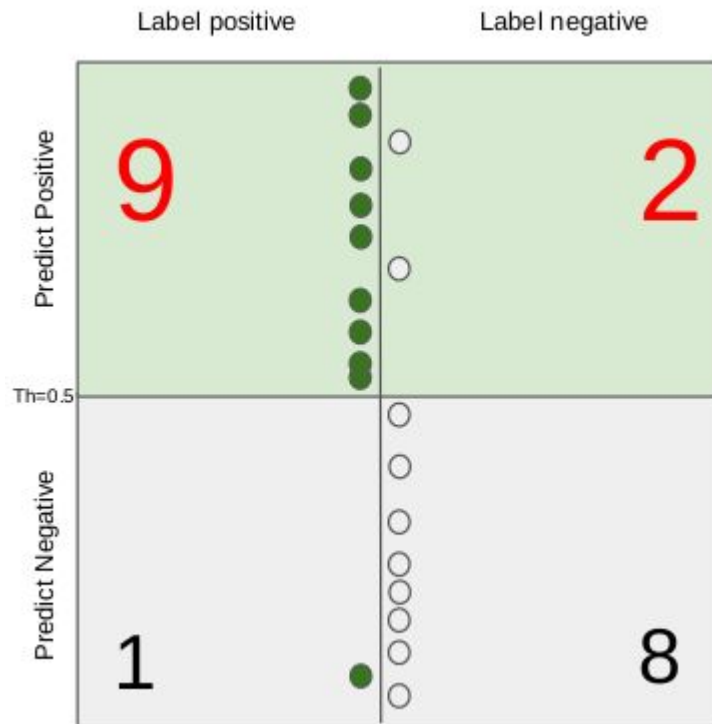
# Métricas puntuales: exactitud (*accuracy*)



| Th  | TP | TN | FP | FN | Acc |
|-----|----|----|----|----|-----|
| 0.5 | 9  | 8  | 2  | 1  | .85 |

- $acc = (TP + TN) / (TP + FP + TN + FN)$
- equivalente al costo 0/1

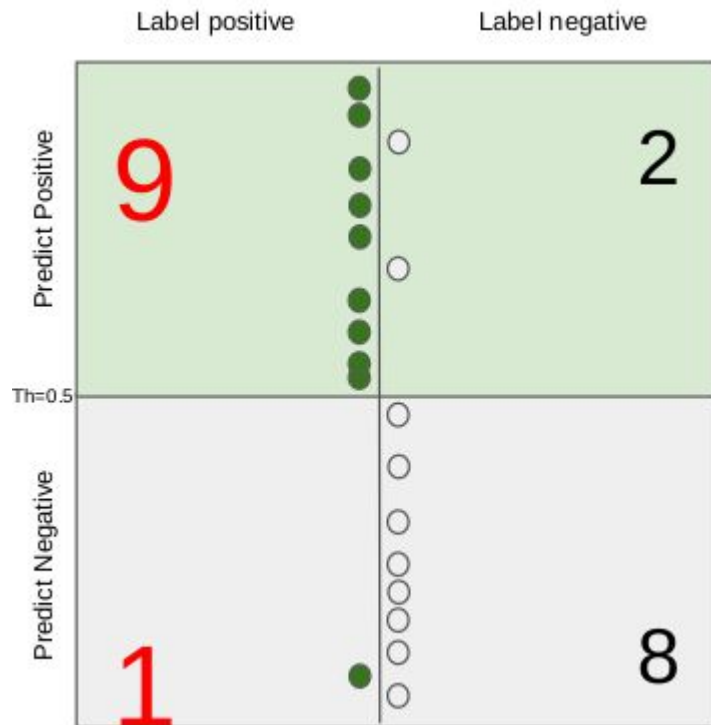
# Métricas puntuales: precisión



| Th  | TP | TN | FP | FN | Acc | Pr  |
|-----|----|----|----|----|-----|-----|
| 0.5 | 9  | 8  | 2  | 1  | .85 | .81 |

- $Prec = TP / (TP + FP)$
- Prec 100% = todos bajo el umbral salvo el de score más alto (siempre que sea correcto)

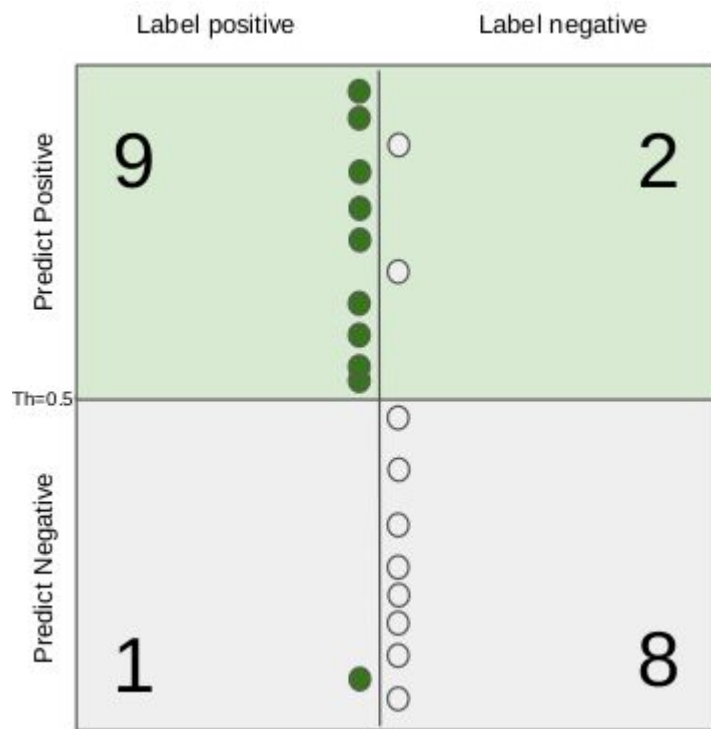
# Métricas puntuales: sensibilidad (*recall*)



| Th  | TP | TN | FP | FN | Acc | Pr  | Recall |
|-----|----|----|----|----|-----|-----|--------|
| 0.5 | 9  | 8  | 2  | 1  | .85 | .81 | .9     |

- $\text{Recall} = \text{TP} / (\text{TP} + \text{FN})$
- Recall 100% = todos los puntos por encima del umbral

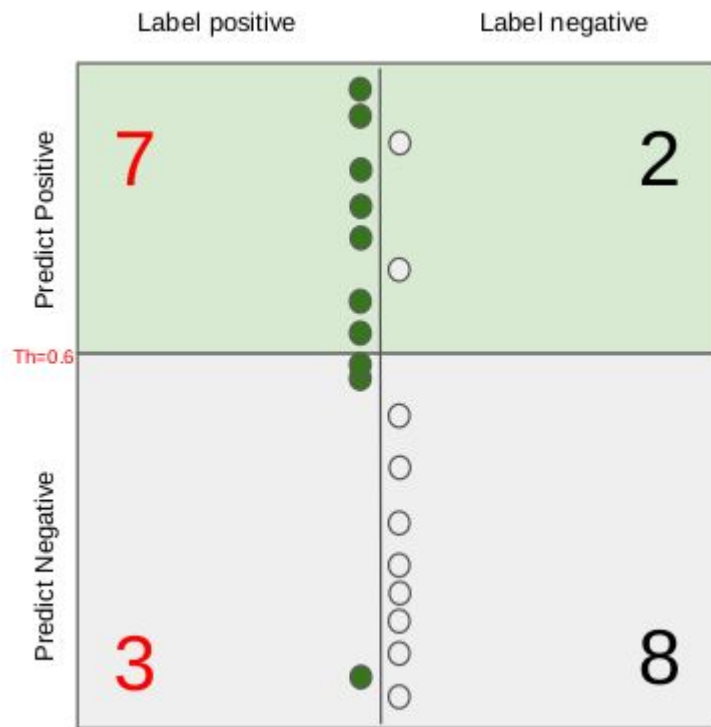
# Métricas puntuales: F1-score



| Th  | TP | TN | FP | FN | Acc | Pr  | Recall | F1   |
|-----|----|----|----|----|-----|-----|--------|------|
| 0.5 | 9  | 8  | 2  | 1  | .85 | .81 | .9     | .857 |

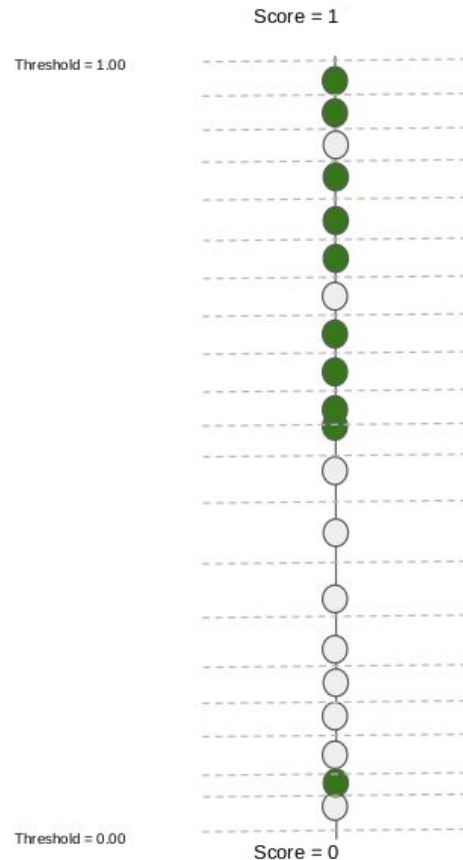
$$F_1 = \left( \frac{2}{\text{recall}^{-1} + \text{precision}^{-1}} \right) = 2 \cdot \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}}$$

# Métricas puntuales: cambio del umbral



| Th  | TP | TN | FP | FN | Acc | Pr  | Recall | F1   |
|-----|----|----|----|----|-----|-----|--------|------|
| 0.6 | 7  | 8  | 2  | 3  | .75 | .77 | .7     | .733 |

# umbrales efectivos = # ejemplos + 1

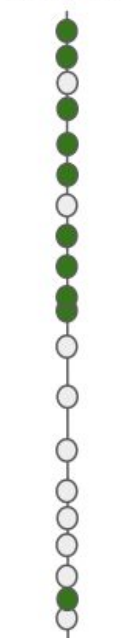


| Threshold | TP | TN | FP | FN | Accuracy | Precision | Recall | Specificity | F1    |
|-----------|----|----|----|----|----------|-----------|--------|-------------|-------|
| 1.00      | 0  | 10 | 0  | 10 | 0.50     | 1         | 0      | 1           | 0     |
| 0.95      | 1  | 10 | 0  | 9  | 0.55     | 1         | 0.1    | 1           | 0.182 |
| 0.90      | 2  | 10 | 0  | 8  | 0.60     | 1         | 0.2    | 1           | 0.333 |
| 0.85      | 2  | 9  | 1  | 8  | 0.55     | 0.667     | 0.2    | 0.9         | 0.308 |
| 0.80      | 3  | 9  | 1  | 7  | 0.60     | 0.750     | 0.3    | 0.9         | 0.429 |
| 0.75      | 4  | 9  | 1  | 6  | 0.65     | 0.800     | 0.4    | 0.9         | 0.533 |
| 0.70      | 5  | 9  | 1  | 5  | 0.70     | 0.833     | 0.5    | 0.9         | 0.625 |
| 0.65      | 5  | 8  | 2  | 5  | 0.65     | 0.714     | 0.5    | 0.8         | 0.588 |
| 0.60      | 6  | 8  | 2  | 4  | 0.70     | 0.750     | 0.6    | 0.8         | 0.667 |
| 0.55      | 7  | 8  | 2  | 3  | 0.75     | 0.778     | 0.7    | 0.8         | 0.737 |
| 0.50      | 8  | 8  | 2  | 2  | 0.80     | 0.800     | 0.8    | 0.8         | 0.800 |
| 0.45      | 9  | 8  | 2  | 1  | 0.85     | 0.818     | 0.9    | 0.8         | 0.857 |
| 0.40      | 9  | 7  | 3  | 1  | 0.80     | 0.750     | 0.9    | 0.7         | 0.818 |
| 0.35      | 9  | 6  | 4  | 1  | 0.75     | 0.692     | 0.9    | 0.6         | 0.783 |
| 0.30      | 9  | 5  | 5  | 1  | 0.70     | 0.643     | 0.9    | 0.5         | 0.750 |
| 0.25      | 9  | 4  | 6  | 1  | 0.65     | 0.600     | 0.9    | 0.4         | 0.720 |
| 0.20      | 9  | 3  | 7  | 1  | 0.60     | 0.562     | 0.9    | 0.3         | 0.692 |
| 0.15      | 9  | 2  | 8  | 1  | 0.55     | 0.529     | 0.9    | 0.2         | 0.667 |
| 0.10      | 9  | 1  | 9  | 1  | 0.50     | 0.500     | 0.9    | 0.1         | 0.643 |
| 0.05      | 10 | 1  | 9  | 0  | 0.55     | 0.526     | 1      | 0.1         | 0.690 |
| 0.00      | 10 | 0  | 10 | 0  | 0.50     | 0.500     | 1      | 0           | 0.667 |

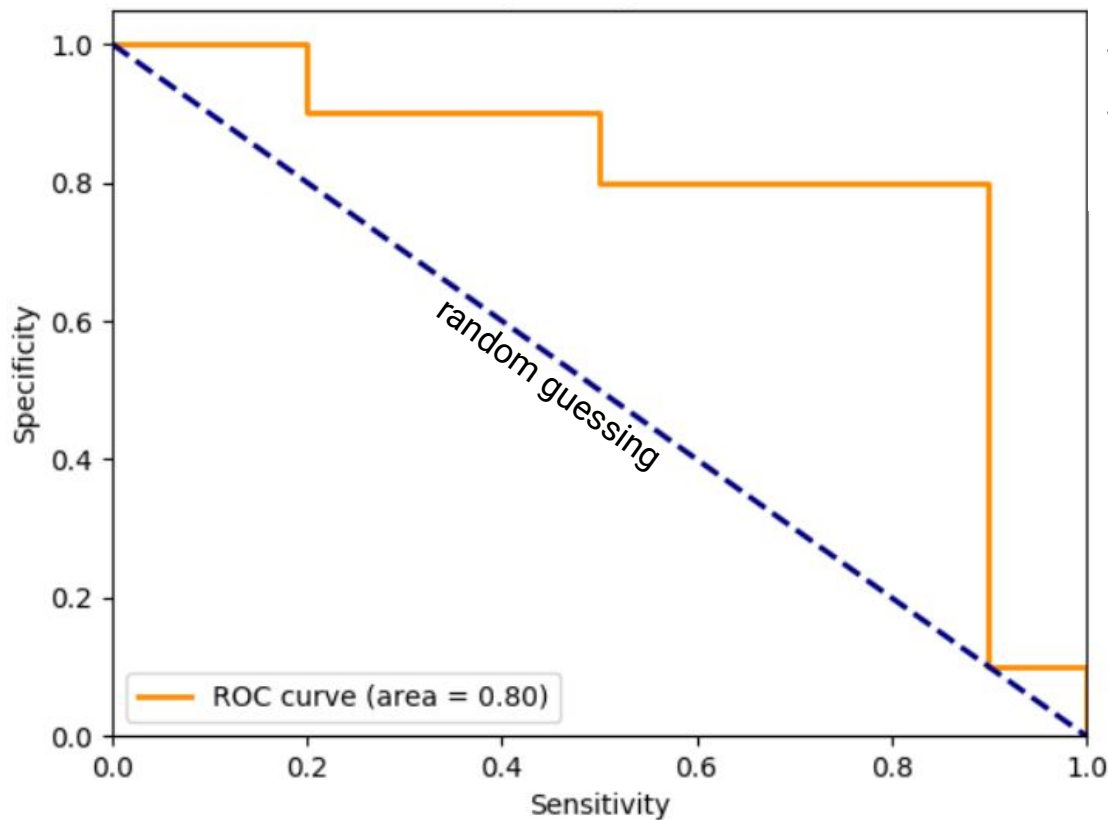


# Métricas resumen: curvas ROC (rotada)

Score = 1



Score = 0



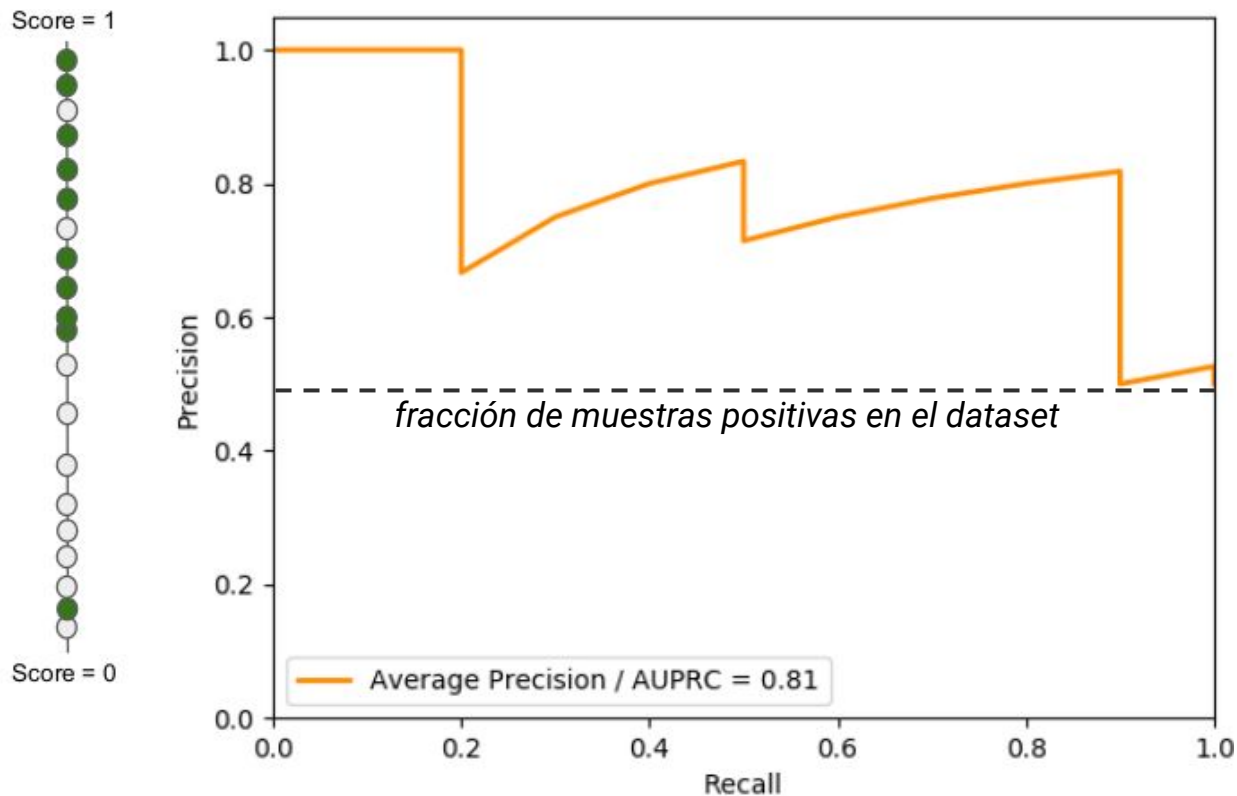
$specificity = tnr = TN / Neg = TN / (TN + FP)$

$sensitivity = tpr = TP / Pos = TP / (TP + FN)$

métrica AUC = área bajo la curva ROC

ROC = Receiver Operating Characteristic

# Métricas resumen: curvas PR



$precision = TP / (\text{pos predichos})$

$recall = TP / (\text{pos verdaderos})$

AUPR = área bajo la curva PR

# Curvas ROC y PR en validación cruzada

## Opción 1:

- Asumir que magnitudes de los scores son comparables entre corridas
- Acumular predicciones de todas las corridas
- Trazar la curva usando predicciones acumuladas

## Opción 2:

- Trazar las curvas individuales para cada partición
- Considerar la “curva promedio”

# Resumen curvas ROC y PR

- Permiten evaluación cuantitativa a distintos niveles de “confianza”
- Asumen problemas binarios
- Se pueden resumir en medidas del tipo “área bajo la curva”
- Las curvas ROC son insensibles a cambios en la distribución de clases en el conjunto de test
- Las curvas PR muestran la fracción de las predicciones que son FP
- Las curvas PR son útiles en problemas con una proporción de muestras negativas muy alta
- Permiten determinar umbrales óptimos para distintos puntos de operación

# Problemas multiclase

- Problema con N clases => matriz de confusión de NxN
  - La mayoría de las métricas se analizan como N problemas binarios (OVA)
    - El desbalance crece con el número de clases
  - Variantes multiclase de métricas AUC
    - micro vs. macro average
- activity recognition from video

|       |     |     |    |     |   |   |    |   |   |
|-------|-----|-----|----|-----|---|---|----|---|---|
| bend  | 100 | 0   | 0  | 0   | 0 | 0 | 0  | 0 | 0 |
| jack  | 0   | 100 | 0  | 0   | 0 | 0 | 0  | 0 | 0 |
| jump  | 0   | 0   | 89 | 0   | 0 | 0 | 11 | 0 | 0 |
| pjump | 0   | 0   | 0  | 100 | 0 | 0 | 0  | 0 | 0 |

