



# APRENDIZAJE AUTOMÁTICO: TERCER CURSO DEL GRADO DE ING. INFORMÁTICA EN COMPUTACION

## Aprendizaje supervisado: Introducción y evaluación del rendimiento de un clasificador

**César Hervás-Martínez**  
**Grupo de Investigación AYRNA**

**Departamento de Informática y Análisis  
Numérico**  
**Universidad de Córdoba**  
**Campus de Rabanales. Edificio Einstein.**  
**Email: [chervas@uco.es](mailto:chervas@uco.es)**

**2019-2020**



# Clasificación automática y Reconocimiento de patrones



**Sistema de decisión automática** donde se reconocen los patrones de unos datos muestrales mediante modelos aprendidos (estimados) utilizando técnicas de computación inteligente.

Mediante esta metodología un **algoritmo** extrae información de un conjunto de datos etiquetados de forma tal que el **modelo aprendido** sea capaz de predecir las etiquetas asociadas a un conjunto de nuevos datos no utilizados en el aprendizaje

Estas técnicas han atraído la atención de la **sociedad de la información y de la industria.**

**Aplicaciones** en motores de búsqueda, diagnósticos médicos, detección de fraude, clasificación de secuencias de ADN, reconocimiento de caracteres, etc.



# UN PROBLEMA DE APRENDIZAJE



**Reconocimiento de caracteres (handwritten digit recognition)**

☐ ¿En qué consiste un problema de aprendizaje?

☐ **Conceptos básicos:**

- 1. Fase de entrenamiento (conjunto de entrenamiento)**
- 2. Fase de test o generalización (conjunto de test)**
- 3. Función de error**

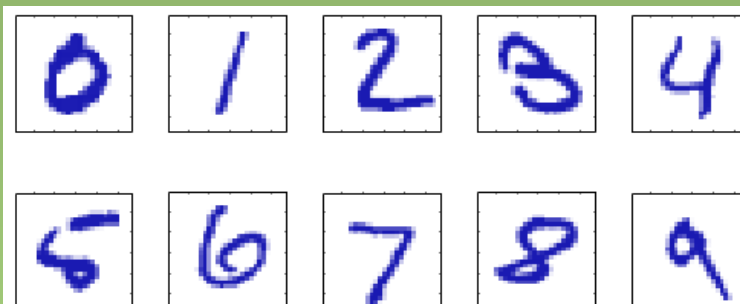


## Reconocimiento de caracteres escritos a mano

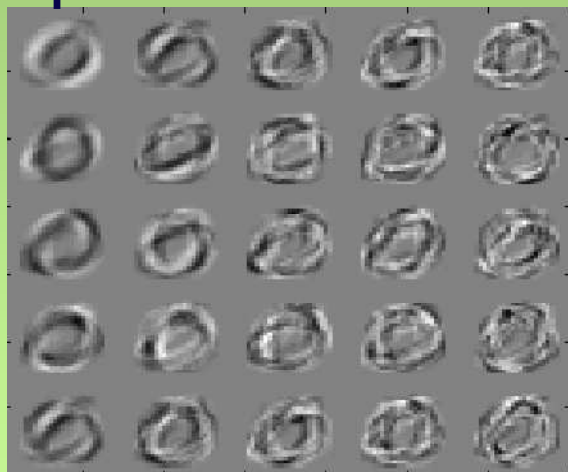
9 7 0 2 5 2 / 6 8 8  
6 7 8 8 8 7 0 \ /  
3 0 / 2 3 4 5 6 7 8  
9 0 \ 2 9 4 6 6 7 8  
0 / 2 5 6 0 9 9 6 /  
3 4 5 6 7 8 9 0 / 2  
6 7 8 9 0 / 2 3 4 5  
/ / / / 2 3 4 5 5 6  
4 6 4 5 6 7 6 8 9 9  
9 9 8 9 9 9 1 2 4 5  
6 7 8 9 4 5 6 3 3 2  
3 3 3 3 8 4 4 9 6 0  
9 0 0 6 5 8 9 5 6 8  
8 / 8 9 9 4 6 7 8 9  
7 2 2 5 7 8 3 6 4 1



# UN PROBLEMA DE APRENDIZAJE: Handwritten digits recognition



Muestra de los autovalores del número 0 generado por el clasificador PCA



Dígitos muestrales utilizados para entrenar al clasificador

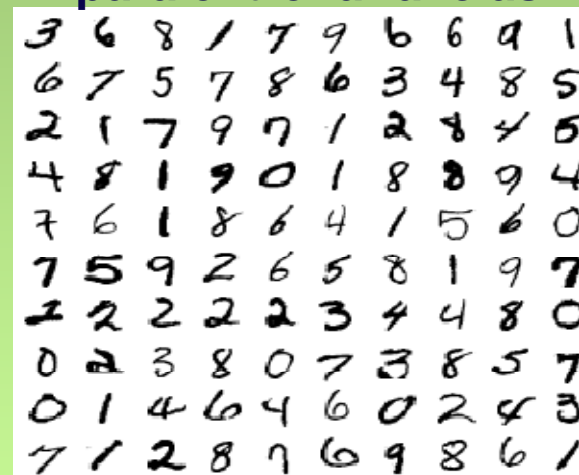


Imagen digital

Pre procesamiento

Extracción de características

Clasificación 1NN

Reconocimiento



## UN PROBLEMA DE APRENDIZAJE



### □ Formulación matemática del problema:

- Imágenes de 28x28 pixeles
- Representar cada imagen como un vector  $\mathbf{x} \in \mathbb{R}^{784}$
- Clasificador  $f(\mathbf{x})$

$$f : \mathbf{x} \rightarrow \{0, 1, 2, 3, 4, 5, 6, 7, 8, 9\}$$



## UN PROBLEMA DE APRENDIZAJE



- Extraemos una m.a.s o Conjunto de Entrenamiento:

$$D = \{(\mathbf{x}_i, y_i)\} \quad \text{para } i=1, \dots, n$$

- Siendo  $y_i$  el valor de la clase a la que pertenece el vector de pixeles  $\mathbf{x}_i$  correspondiente al patrón i-ésimo

$$x_i = (1, 1, 0, 0, \dots, 1, 0, \dots, 0, 0, 1)$$

$$y_i \rightarrow \{0, 1, 2, 3, 4, 5, 6, 7, 8, 9\}$$



# UN PROBLEMA DE APRENDIZAJE



- ❑ **Funciones de error en el entrenamiento de un clasificador:**

**Métodos de estimación de la probabilidad de clasificación correcta: “Porcentaje de patrones mal clasificados”**

**“Minima sensibilidad”**

**“RMSE”**

**La curva ROC y el área bajo la curva ROC “AUC”**

**Puntuación de Brier**

- ❑ **Evaluación sensible al coste**

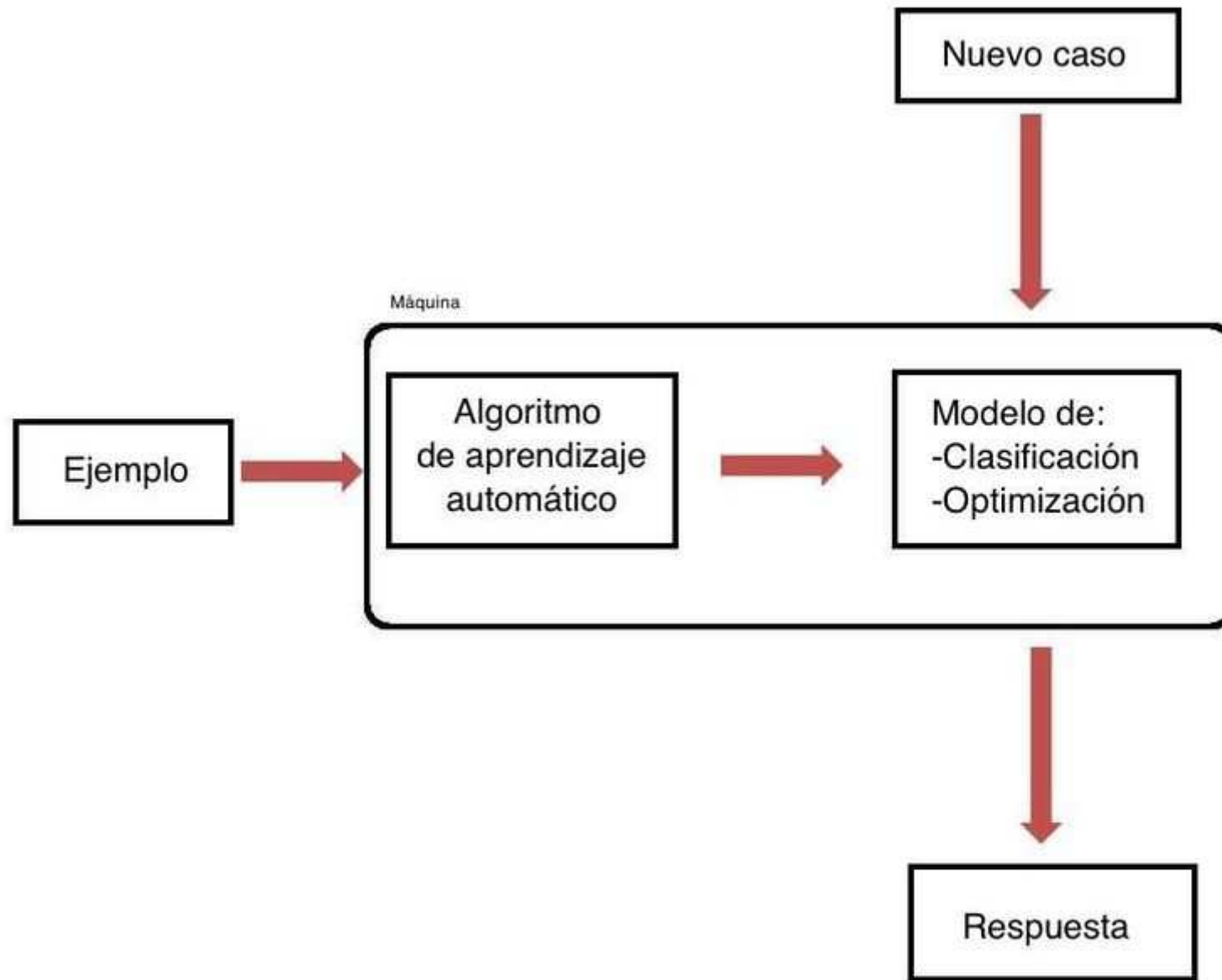
**Objetivo:**

**“Minimizar la función de error en el conjunto de test”**





# ¿QUE ES EL APRENDIZAJE DE MÁQUINAS?



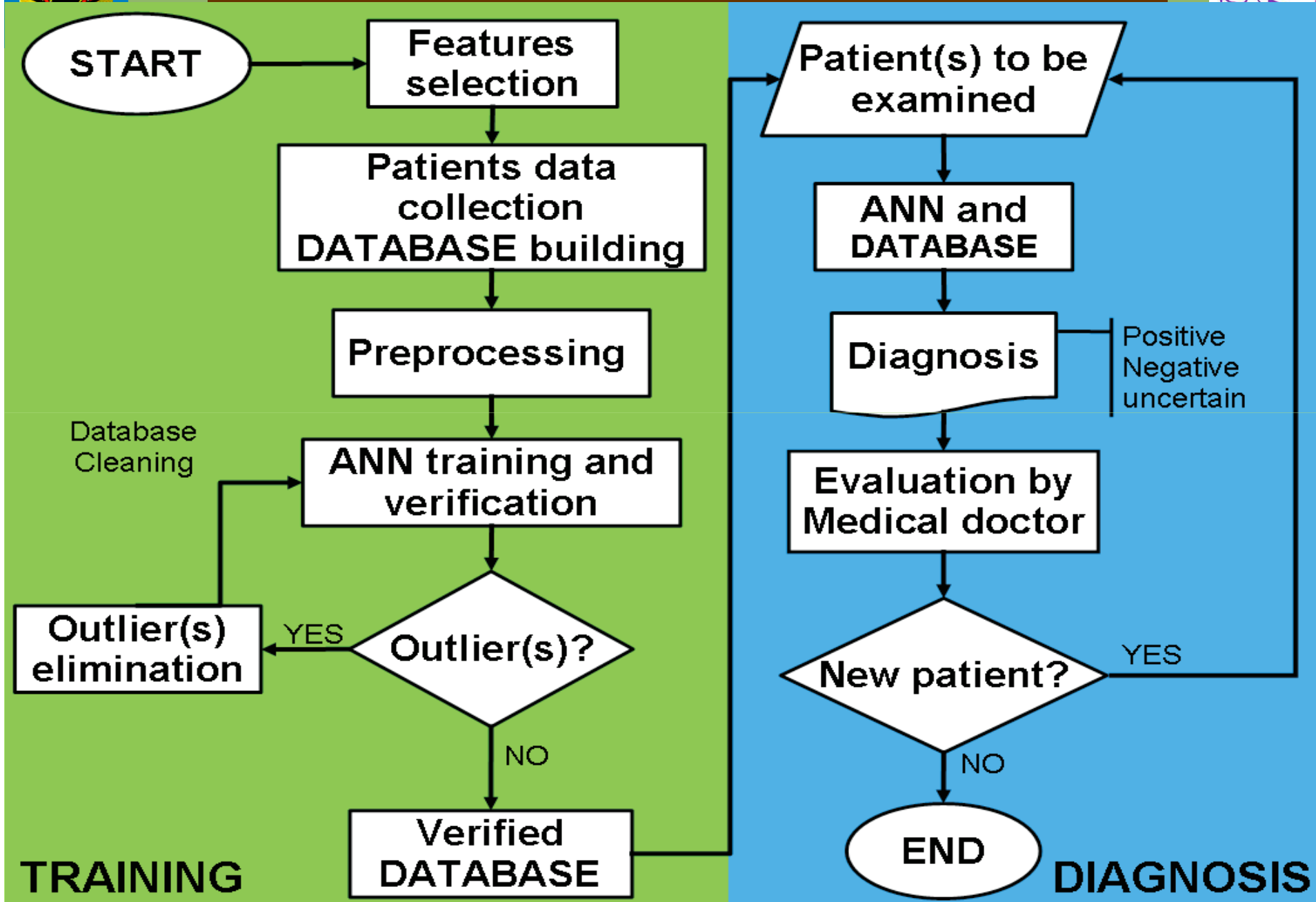


# ¿QUE ES EL APRENDIZAJE DE MÁQUINAS?





# ¿QUE ES EL APRENDIZAJE DE MÁQUINAS?





# APRENDIZAJE SUPERVISADO



□ Entrenamiento

Funciones (modelo)

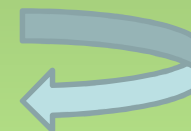
$$D = \{(\mathbf{x}_i, y_i) \in X \times Y\}$$

$$F = \{f : X \rightarrow Y\}$$

□ Aprendizaje



$$\hat{f} \in F$$

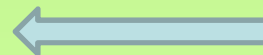


□ Predicción

$$\text{tal que } y_i \approx \hat{f}(\mathbf{x}_i)$$

Nuevos datos

$$\hat{y} = \hat{f}(\mathbf{x})$$



$\mathbf{x}$



# APRENDIZAJE SUPERVISADO Y NO-SUPERVISADO



Ejemplos canónicos de aprendizaje:

1. **Regresión** (aprendizaje supervisado), la variable dependiente es continua
2. **Clasificación** (aprendizaje supervisado), la variable dependiente es nominal
3. **Clasificación ordinal** (aprendizaje supervisado) la variable dependiente es ordinal
4. **Agrupamiento “Clustering”** (aprendizaje no supervisado) no se conoce el valor de la variable dependiente



# APRENDIZAJE SUPERVISADO Rendimiento de un clasificador



Consideremos un problema de clasificación con  $J$  clases y  $n$  patrones de entrenamiento o test

El rendimiento de un clasificador  $g$  se puede obtener a partir de su matriz de contingencia o confusión definida en la forma

$C/C^*$	$C^*_1$	...	$C^*_j$	...	$C^*_J$	
$C_1$	$n_{11}$	...	...	...	...	$n_{1.}$
...						
$C_i$			$n_{ij}$			$n_{i.} = \sum_{j=1}^J n_{ij}$
...						
$C_J$	$n_{k1}$				$n_{k1}$	$n_{k.}$
	$n_{.1}$	.....	$n_{.j} = \sum_{j=1}^J n_{ij}$	...	$n_1$	$n = \sum_{i=1}^J \sum_{j=1}^J n_{ij}$

$$M(g) = \left( n_{ij} / \sum_{i=1}^J \sum_{j=1}^J n_{ij} = n \right)$$

Donde  $n_{ij}$  representa el número de veces que los patrones son predichos por el clasificador  $g$  en la clase  $C^*_j$  cuando realmente pertenecen a la clase  $C_i$



## Coste de una mala clasificación



No podemos comparar numéricamente las clases  $C_i$  y  $C_j$  pero podemos asignar un coste artificial cuando confundimos, con el clasificador, la clase  $C_i$  por la clase  $C_j$ . Los costes se pueden definir mediante una matriz de costes  $C_{ij}$ , donde este valor indica el coste de clasificar una instancia o patrón de la clase  $C_i$  como perteneciente a la clase  $C_j$ .

**Algunas elecciones de matrices de costes y de transformación de variables**

$$C_{01} =$$

0	1	1	1	1
1	0	1	1	1
1	1	0	1	1
1	1	1	0	1
1	1	1	1	0

Coste 0,1

Medidas asociadas,  
CCR, MS

0	1	1	1	1
-1	0	1	1	1
-1	-1	0	1	1
-1	-1	1	0	1
-1	-1	-1	-1	0

Matriz de transformación

Medida asociada,

$\tau_b$

$$C_{abs} =$$

0	1	2	3	4
1	0	1	2	3
2	1	0	1	2
3	2	1	0	1
4	3	2	1	0

Coste absoluto

Medidas asociadas,  
RMSE, MAE, AMAE, MMAE, mMAE  
 $r_s$



## Clasificación binaria:

- Supongamos que tenemos un conjunto de  $N$  observaciones (conjunto de entrenamiento)

$$(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N)$$

$$(y_1, y_2, \dots, y_N)$$

$$\mathbf{x}_i \in \mathbb{R}^d, y_i \in \{-1, 1\}$$

- Problema de clasificación: estimar  $f(\mathbf{x})$  tal que

$$f(\mathbf{x}_i) = y_i$$





# CLASIFICACIÓN



- ❑ **Función de error: Porcentaje de patrones mal clasificados**

$$\frac{1}{N} \sum_{i=1}^N [y_i \neq f(\mathbf{x}_i)]$$

- ❑ **Otras denominaciones: Función de pérdida o Función de riesgo.**
- ❑ **Función de evaluación: Porcentaje de patrones correctamente clasificados**

$$C = \frac{1}{N} \sum_{i=1}^N [y_i = f(\mathbf{x}_i)]$$



# METRICAS EN CLASIFICACIÓN BINARIA



❑ Matriz de confusión (caso binario):

Clase asignada o predicha	Clase real de pertenencia	
	TP	FP
	FN	TN

Eficacia, o tasa de aciertos,

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

**TP= Verdadero Positivo,    FP= Falso Positivo,**

**FN= Falso Negativo , TN= Verdadero Negativo**



# MEDIDAS DE MERITO EN CLASIFICACIÓN BINARIA



	Clase asignada o predicha	
Clase real de pertenencia	TP	FN
	FP	TN

$$\text{Precision} = \frac{TP}{TP+FP}$$

$$\text{Especificidad} = \frac{TN}{TN+FP}$$

$$\text{Recall} = \text{Sensibilidad} = \frac{TP}{TP+FN}$$

$$F = \frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

$$\text{Kappa} = \frac{TP+TN-E(TP+TN)}{TP+TN+FP+FN-E(TP+TN)}$$



## METRICAS EN CLASIFICACIÓN BINARIA



□ Ejemplo de matriz de confusión (caso binario):

Clase asignada o predicha

Clase real de pertenencia

$$\begin{pmatrix} TP & FN \\ FP & TN \end{pmatrix} = \begin{pmatrix} 20 & 10 \\ 15 & 15 \end{pmatrix}$$

Sensibilidad      Especificidad

$$\text{Porcentaje de aciertos, } C = \frac{35}{60}$$

$$\text{Especificidad} = \frac{TN}{TN+FP}$$

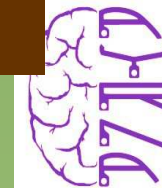
$$\text{Recall} = \text{Sensibilidad} = \frac{TP}{TP+FN}$$

$$\text{Especificidad} \equiv \text{TNR} = \frac{15}{15+15} = 0.5$$

$$\text{Sensibilidad} \equiv \text{TPR} = \frac{20}{20+10} = \frac{2}{3}$$



## METRICAS EN CLASIFICACIÓN BINARIA



- ❑ Problemas no-balanceados
- ❑ Ejemplo: enfermos de cáncer (10) e individuos sanos (991):

Clase real

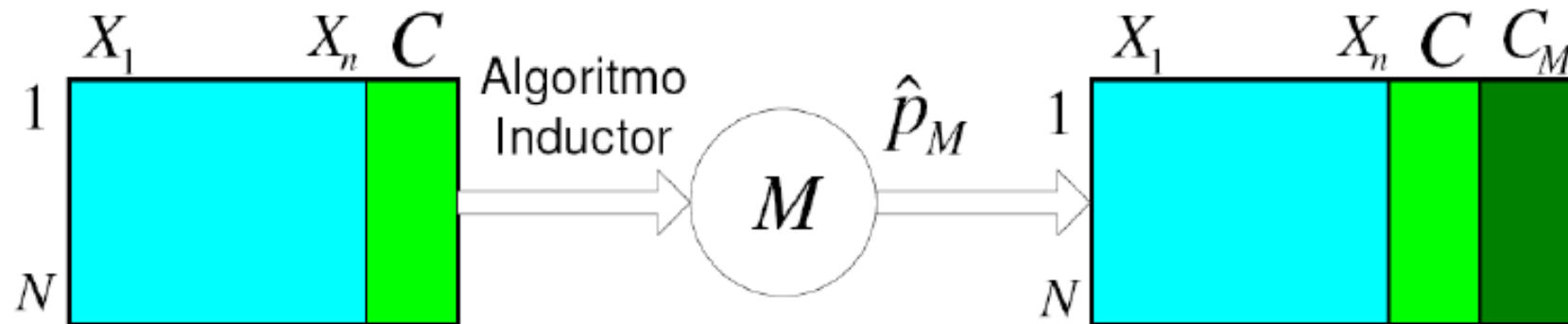
$$\begin{pmatrix} TP & FN \\ FP & TN \end{pmatrix} = \begin{pmatrix} 0 & 10 \\ 1 & 990 \end{pmatrix}$$

$$C = \frac{990}{1001} \quad Sensi = 0, Espec = \frac{990}{991}$$

$$\text{Recall} = \text{Sensibilidad} = \frac{TP}{TP + FN} \quad \text{Especificidad} = \frac{TN}{TN + FP}$$



# CLASIFICACIÓN DESHONESTA

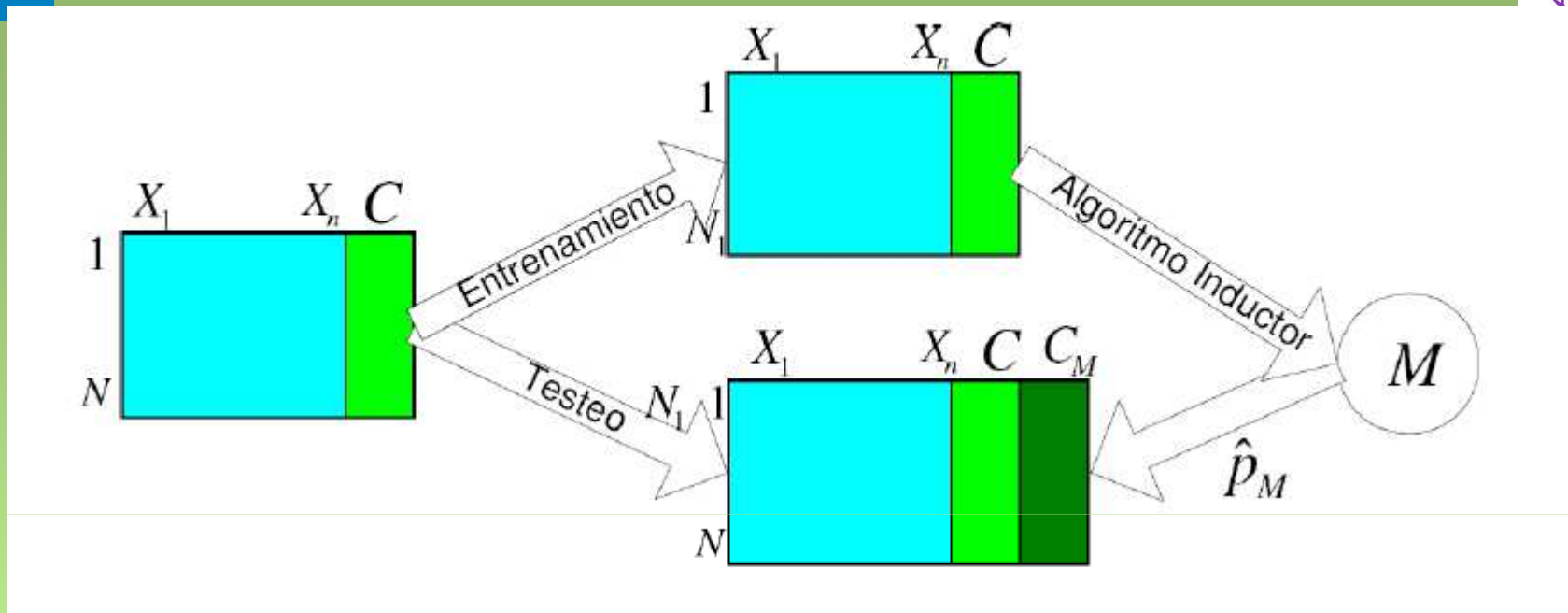
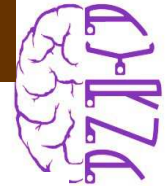


Estimación de la Precisión  $\hat{p}_M = \frac{1}{N} \sum_{i=1}^N \delta(C_i = C_{i,M})$

**C= Clase real,  $C_M$ = Clase estimada por el modelo**



# CLASIFICACIÓN MEDIANTE HOLDOUT 75-25



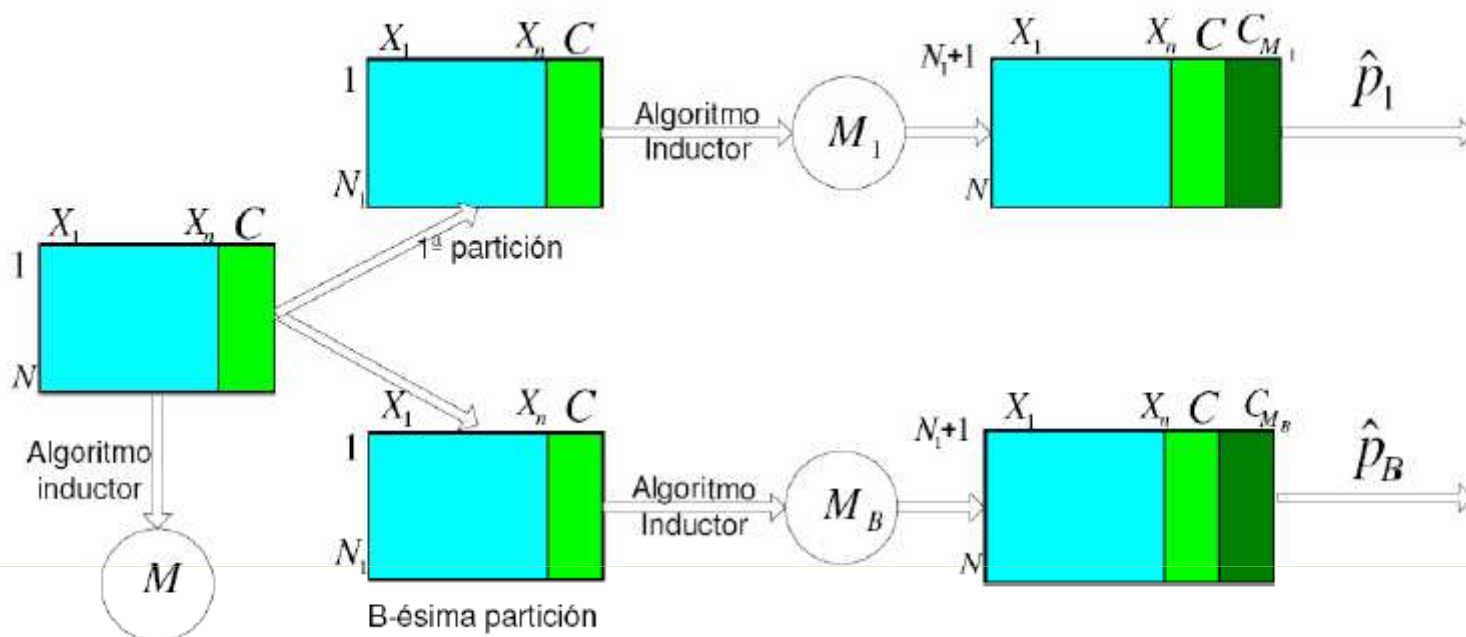
## Estimación de la Precisión

$$\hat{p}_M = \frac{1}{N - N_1} \sum_{i=1}^{N - N_1} \delta(C_{N_1+i} = C_{N_1+i, M})$$

$$N_1 \simeq 0,75 \times N$$



# CLASIFICACION MEDIANTE UN HOLDOUT REPETIDO, B veces



En realidad es un *ensemble* (modelo de modelos) promedio de los resultados de los modelos  $M_1, M_2, \dots, M_B$

$$\hat{p}_M = \frac{1}{B} \sum_{i=1}^B \hat{p}_i$$

$$E(\hat{p}_M) = \frac{1}{B} \sum_{i=1}^B E(\hat{p}_i); \quad V(\hat{p}_M) = \frac{1}{B^2} \left( \sum_{i=1}^B V(\hat{p}_i) + 2 \sum_{i=1 < j}^B Cov(\hat{p}_i, \hat{p}_j) \right)$$

Luego cuanto más dependientes sean los estimadores  $\hat{p}_i$  más varianza tendrá  $\hat{p}_M$

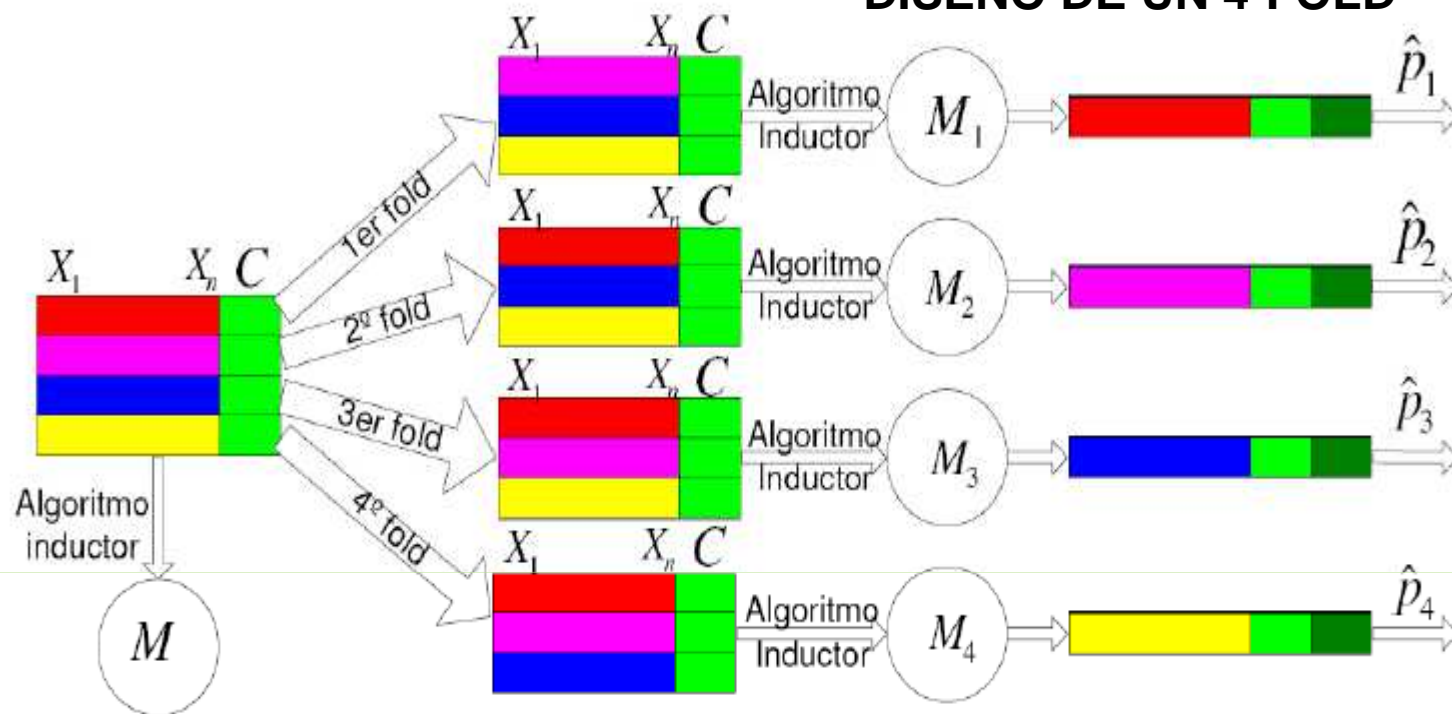




# CLASIFICACIÓN MEDIANTE UN k-FOLD



## DISEÑO DE UN 4-FOLD



$$\hat{p}_M = \frac{1}{k} \sum_{i=1}^k \hat{p}_i$$

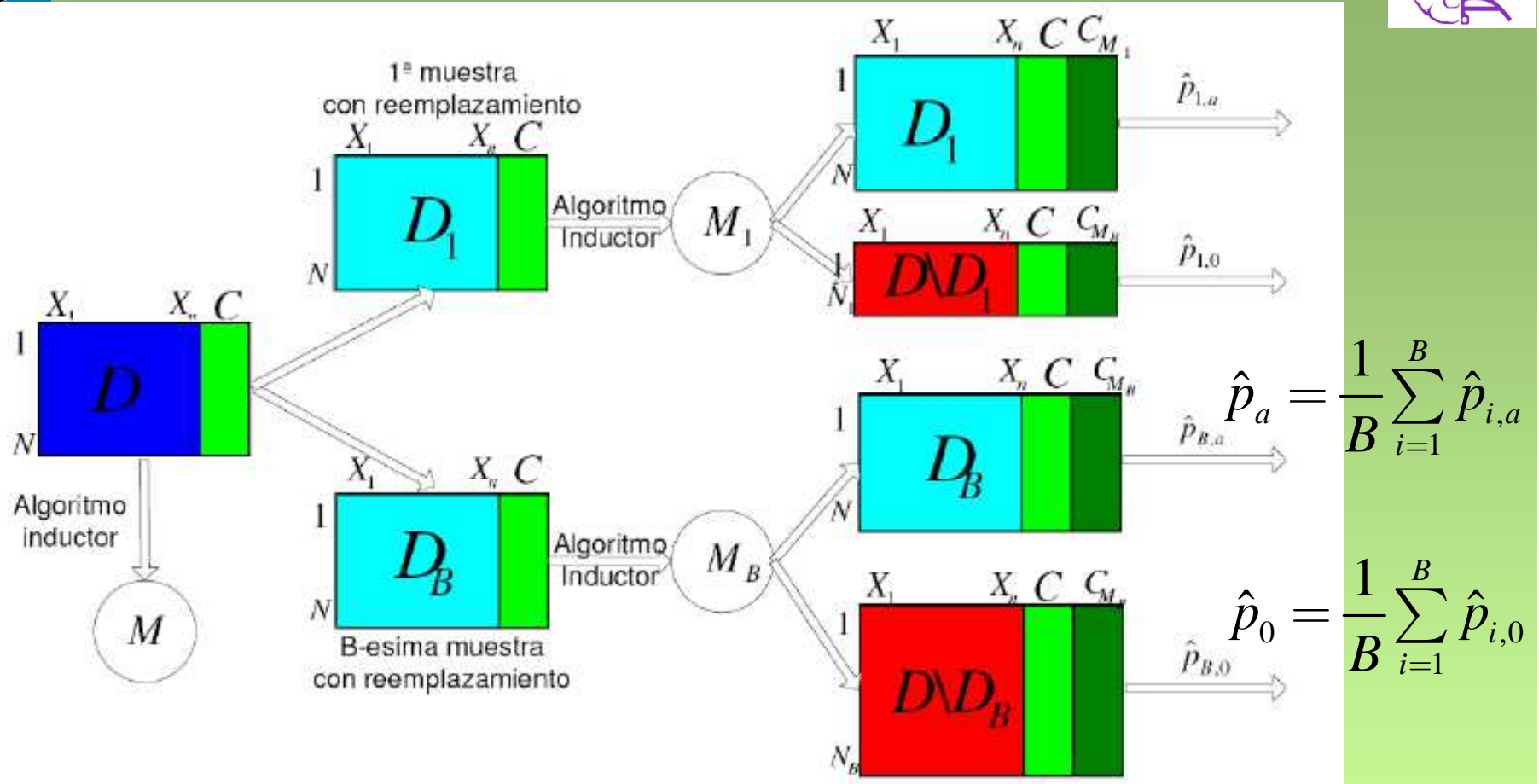
En este caso k=4

$$E(\hat{p}_M) = \frac{1}{k} \sum_{i=1}^k E(\hat{p}_i); \quad V(\hat{p}_M) = \frac{1}{k^2} \left( \sum_{i=1}^k V(\hat{p}_i) + 2 \sum_{i=1 < j}^k Cov(\hat{p}_i, \hat{p}_j) \right)$$

Luego cuanto más dependientes sean los estimadores  $\hat{p}_i$  más varianza tendrá  $\hat{p}_M$



# METODO DE ESTIMACION 0,632 BOOTSTRAPPING



$$\hat{p}_M = \hat{p}_{0,632B_0} = (0,368 \hat{p}_a + 0,632 \hat{p}_0)$$

Es por tanto un *ensemble de ensembles*



# CLASIFICACIÓN



## Consejos de uso de los distintos métodos

**Método Holdout:** utilizarlo con N grande.

**Método Holdout repetidas veces:** no hay control sobre los casos usados como entrenamiento/testeo.

**Método de estimación basado en (k-fold cross validation):** estimación insesgada de la probabilidad de acierto, pero con alta varianza.

**Método de estimación 0,632 bootstraping:** insesgada en el límite y con baja varianza.



# PUNTUACIÓN DE BRIER PARA CLASIFICACIÓN BINARIA



Salida probabilística del modelo

	$X_1$	...	$X_n$	$C$	$p(C_M = 0 \mathbf{x})$	$p(C_M = 1 \mathbf{x})$
$(\mathbf{x}^{(1)}, c^{(1)})$	$x_1^{(1)}$	...	$x_n^{(1)}$	1	0,18	0,82
$(\mathbf{x}^{(2)}, c^{(2)})$	$x_1^{(2)}$	...	$x_n^{(2)}$	0	0,51	0,49
...	...	...	...	...	...	...
$(\mathbf{x}^{(N)}, c^{(N)})$	$x_1^{(N)}$	...	$x_n^{(N)}$	1	0,55	0,45

Salida real

La puntuación de Brier es

$$B = \frac{1}{N} \sum_{i=1}^N \sum_{C=0}^1 \left[ p(C_M = c | x_i) - \delta(C_i, C_{m,i}) \right]^2$$

Y en este ejemplo concreto

$$B = \frac{1}{N} \left[ (0,18 - 0)^2 + (0,82 - 1)^2 + (0,51 - 0)^2 + (0,49 - 1)^2 + \dots \right. \\ \left. + (0,55 - 0)^2 + (0,45 - 1)^2 \right]$$



# CLASIFICACIÓN



Medida de la calibración para un clasificador que asigne, para cada patrón, probabilidades a posteriori a cada valor de la clase.

Suponiendo que la clase real del patrón  $\mathbf{x}$  es 0, se trata de distinguir entre:

$$p(C_M = 0 | \mathbf{x}) = 0,51 \quad y \quad p(C_M = 0 | \mathbf{x}) = 0,97$$

Interesa clasificadores con bajo valor de Brier (bastante seguros en sus predicciones, **dado que las probabilidades a posteriori son altas**)

Para problemas con 2 clases:

$$0 \leq B \leq 2$$



# CLASIFICADOR CON MENOR COSTE



- Ejemplos:

Matrices de confusión

Pred Real

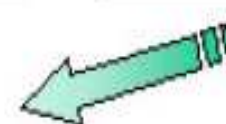
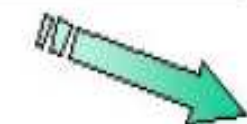
$G_1$	abrir	cerrar
ABRIR	300	500
CERRAR	200	99000

Real

$G_2$	abrir	cerrar
ABRIR	0	0
CERRAR	500	99500

Real

$G_3$	abrir	cerrar
ABRIR	400	5400
CERRAR	100	94100



Real

Predicho

	abrir	cerrar
ABRIR	0	100€
CERRAR	2000€	0

Matriz de  
coste

Matrices resultado

$G_1$	abrir	cerrar
ABRIR	0€	50.000€
CERRAR	400.000€	0€

$G_2$	abrir	cerrar
ABRIR	0€	0€
CERRAR	1.000.000€	0€

$G_3$	abrir	cerrar
ABRIR	0€	540.000€
CERRAR	200.000€	0€

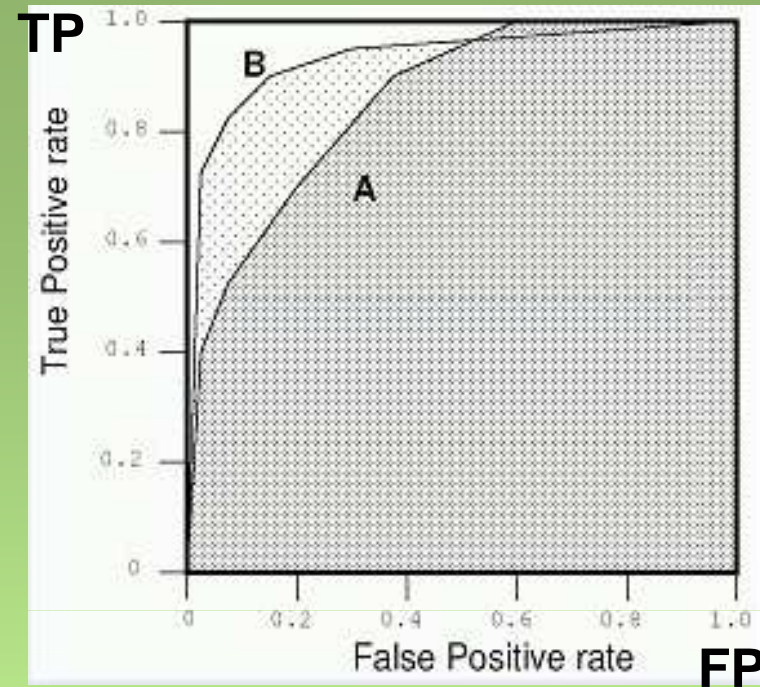
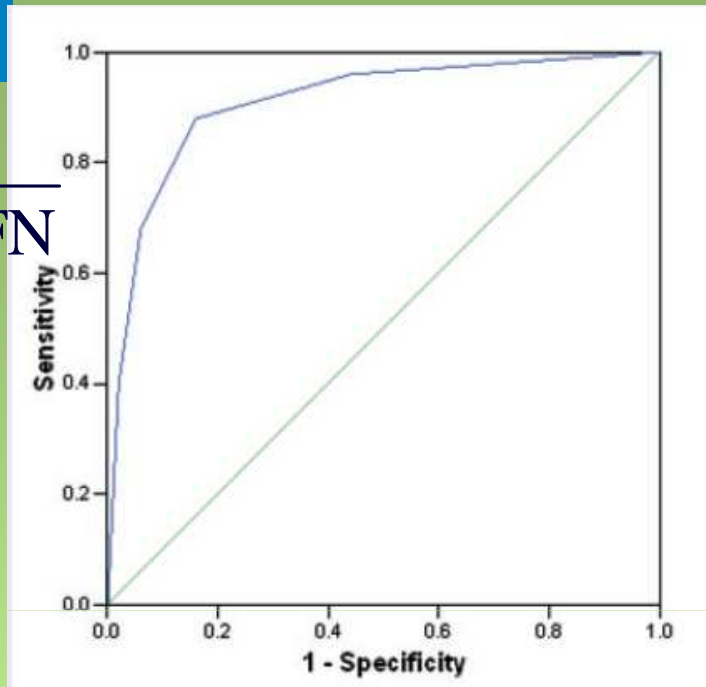




# AREA BAJO LA CURVA ROC



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



$$1 - \frac{TN}{TN+FP} = \frac{FP}{TN+FP}$$

Si cada punto de la curva ROC representa un clasificador: los puntos de la curva cuanto mas a la izquierda y mas arriba mejor será el clasificador

Si cada punto de la curva ROC corresponde a un umbral con el que se toma la decisión: seleccionar el clasificador con mayor área bajo la curva ROC (AUC)



# AREA BAJO LA CURVA ROC



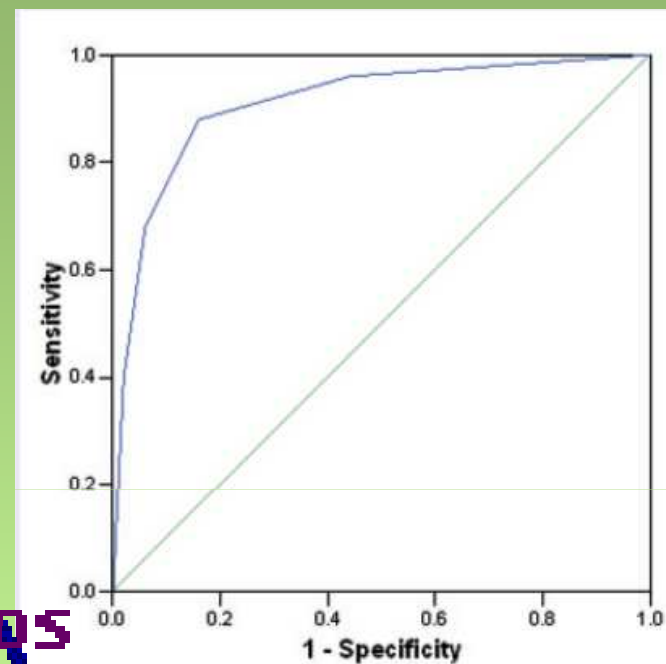
Puntos de corte  
Punto de corte

Normales

Diabéticos

FN

FP



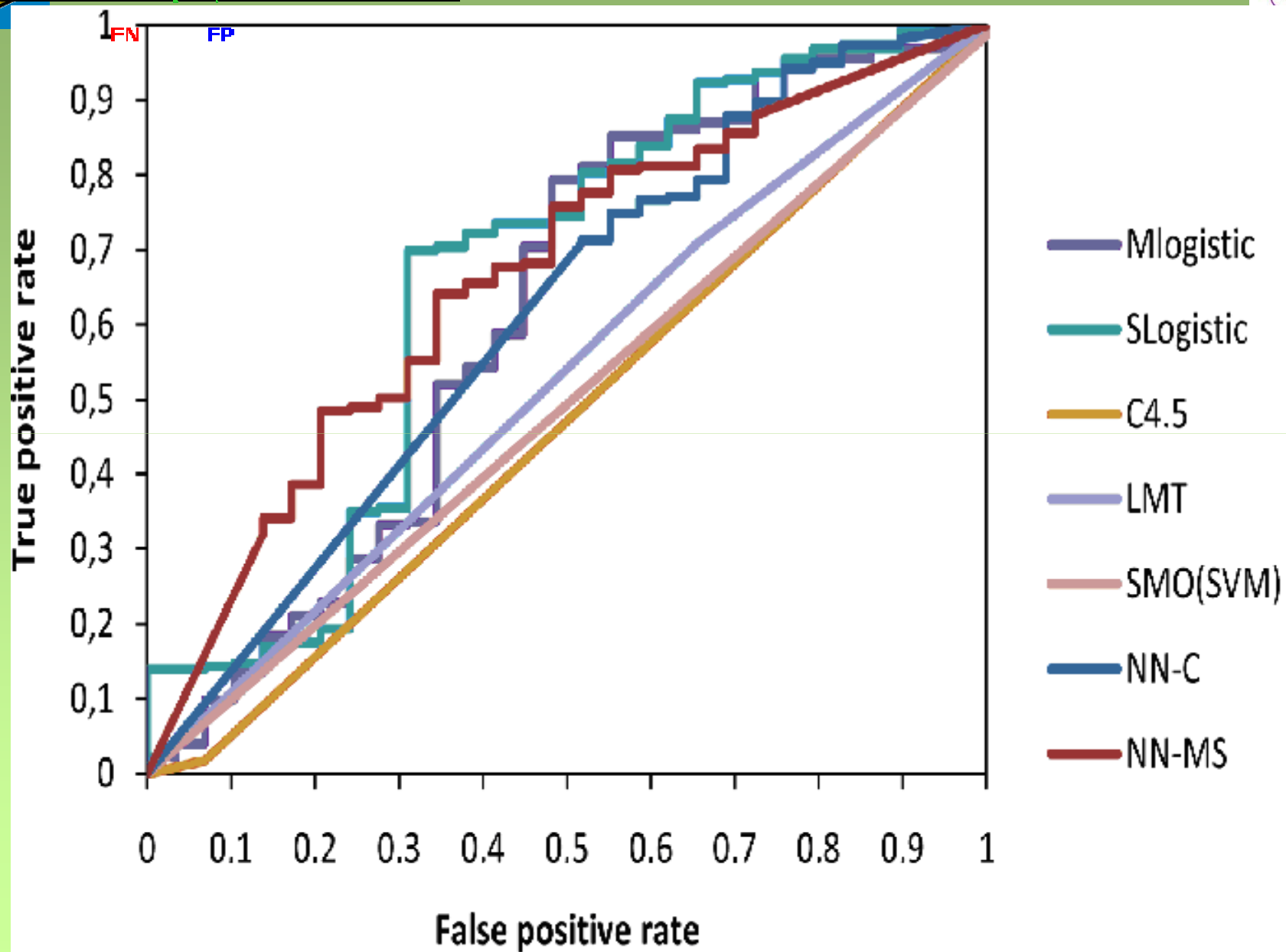




RESUMEN

Diabéticos

# CLASIFICACIÓN





## Clasificación multiclase: Ejemplo de base de datos de clasificación nominal



La base de datos iris es bien conocida en la investigación en estadística y recientemente en aprendizaje puesto que se utiliza para ilustrar la eficacia de algoritmos de clasificación nominal

**iris setosa**



**iris versicolor**



**iris virginica**



**La base de datos contiene medidas de las dimensiones de 50 muestras de flores de cada especie.**

Fisher, R.A. (1936). "The Use of Multiple Measurements in Taxonomic Problems". *Annals of Eugenics* 7: 179–188, disponible en:

<http://digital.library.adelaide.edu.au/coll/special//fisher/138.pdf>.



## Clasificación MULTICLASE: Ejemplo de base de datos de clasificación nominal



Anderson midió las siguientes dimensiones:

- Longitud de los sépalos
- Ancho de los sépalos
- Longitud de los pétalos
- Ancho de los pétalos
- Las cuatro dimensiones de la flor iris se denominan *atributos o variables de entrada o independientes*.
- Las tres especies ( $Q=3$ ) de la flor iris se denominan *clases o variable de salida o dependiente*. Cada ejemplo de flor iris se denomina una *muestra*, un *patrón* o una *instancia*.

Las clases se etiquetan en un formato “1 de Q”, de esta forma tendremos:

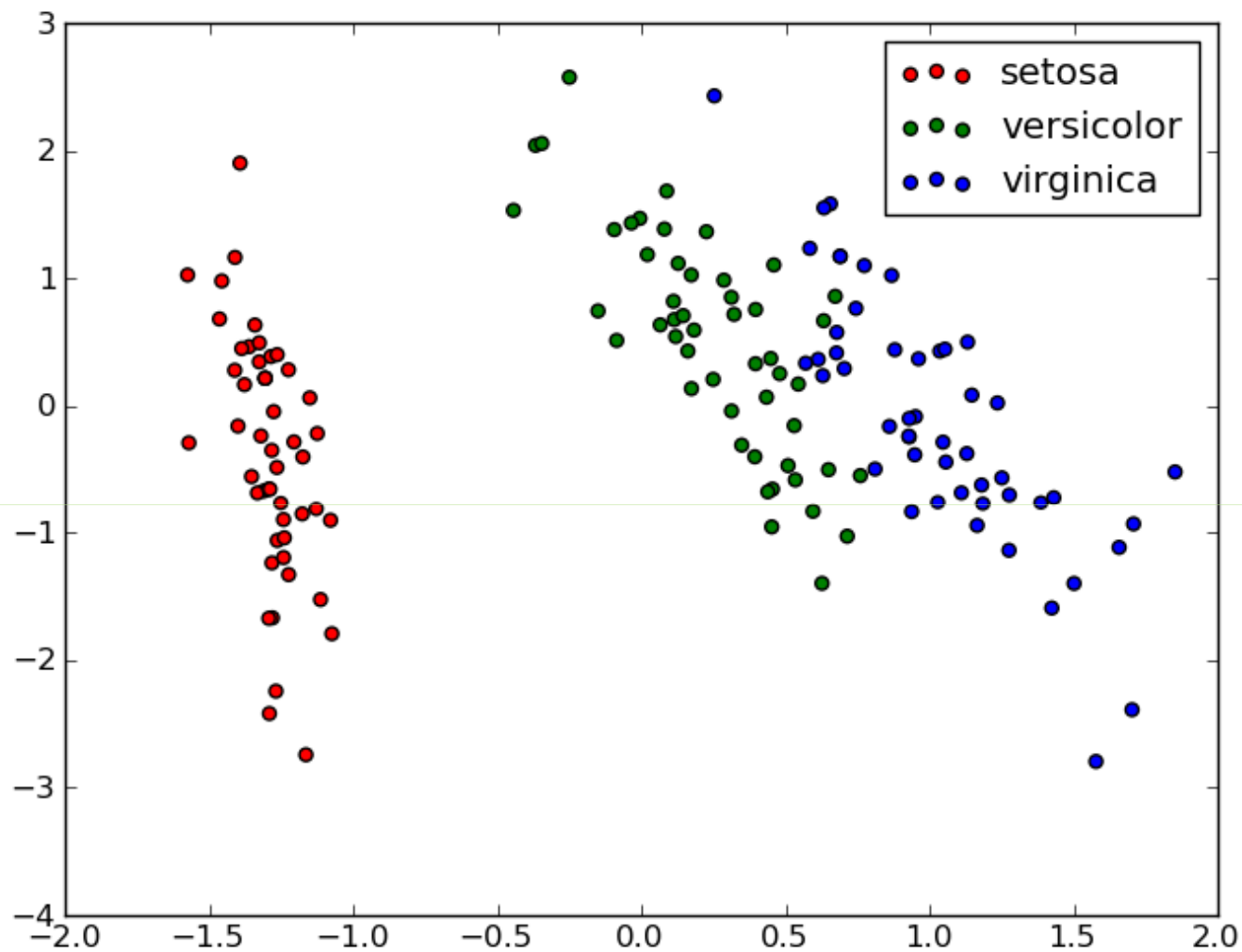
(1,0,0) para la clase setosa

(0,1,0) para la clase virgínica

(0,0,1) para la clase versicolor



## Clasificación MULTICLASE: Ejemplo de base de datos de clasificación nominal



**Ejemplo de proyección en dos dimensiones de las tres clases**



## Métricas de rendimiento



Existen multitud de métricas de rendimiento de un clasificador (multi-clase y binario): CCR, TPR, FPR, Especificidad, AUC, MSE, Entropía, RMSE, etc.

**Problema:** Un mismo valor de una medida puede representar clasificadores muy distintos, en especial en problemas no balanceados y/o con gran número de clases.

Un buen clasificador debería obtener un alto nivel de precisión global, así como un aceptable nivel de precisión para cada clase.

Medida bidimensional en MOEAS: Precisión-Mínima Sensibilidad.



# Métricas de rendimiento



Matriz de confusión de un clasificador y obtención de las medidas MS y C.

$$M = \begin{pmatrix} n_{11} & n_{12} & \cdots & n_{1n} \\ n_{21} & n_{22} & \cdots & n_{2Q} \\ \vdots & \vdots & \vdots & \vdots \\ n_{Q1} & n_{Q2} & \cdots & n_{QQ} \end{pmatrix} \begin{matrix} f_1 \\ f_2 \\ \vdots \\ f_Q \end{matrix}$$

**Sensibilidad de la**

$$S_i = \frac{\text{clase } i}{n_{ii}} / f_i$$

**Minima sensibilidad (MS)**

$$MS = \min \{ S_i; i = 1, \dots, Q \}$$

**Precisión (C)**

$$C = (1/N) \sum_{j=1}^Q n_{jj}$$



# APRENDIZAJE AUTOMÁTICO: TERCER CURSO DEL GRADO DE ING. INFORMÁTICA EN COMPUTACION



## Métodos de evaluación del rendimiento de un clasificador

**GRACIAS POR SU ATENCIÓN**

**César Hervás-Martínez**  
**Grupo de Investigación AYRNA**

**Departamento de Informática y Análisis  
Numérico**  
**Universidad de Córdoba**  
**Campus de Rabanales. Edificio Einstein.**  
**Email: [chervas@uco.es](mailto:chervas@uco.es)**

**2019-2020**