

ESCUELA SUPERIOR POLITÉCNICA DEL LITORAL



## Introducción a los modelos de clasificación y a su evaluación

Andrés G. Abad, Ph.D.

# Agenda

Introducción al problema de clasificación

Evaluando la precisión de la predicción

Validación cruzada

Matriz de confusión

Curva ROC (receiver operating characteristic)

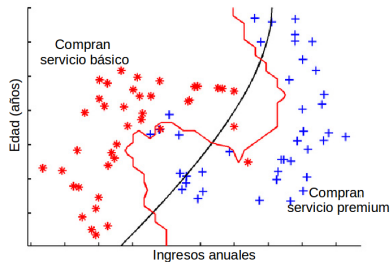
# Definición del problema de clasificación I

- ▶ Un objeto  $\mathbf{x} = [x_1, \dots, x_p]$ , con características  $x_i$ , pertenece exactamente a una clase  $c \in \{1, 2, \dots, C\}$ .
- ▶ Asumimos que tenemos un conjunto de datos

$$\mathcal{D} = \{(\mathbf{x}^{(1)}, c^{(1)}), \dots, (\mathbf{x}^{(n)}, c^{(n)})\}$$

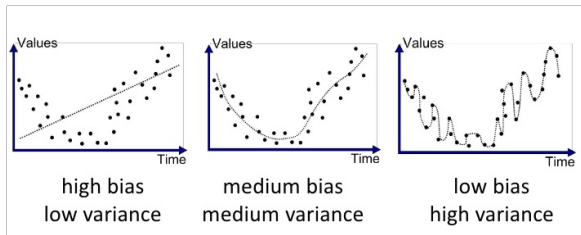
- ▶ Buscamos una función  $\hat{f}$  que asigne  $\mathbf{x}^{(i)}$  a  $c^{(i)}$  lo mejor posible:

$$\hat{f} = \arg \min_f \mathbb{P}_{(\mathbf{x}, c)}[\mathbb{1}(f(\mathbf{x}) \neq c)]$$



- ▶ Objeto  $\mathbf{x}$  pertenece a una de dos clases:  $\{\text{Basico}, \text{Premium}\}$
- ▶ Objeto  $\mathbf{x}$  medidos en dos características:  $x_1$  ingresos anuales, y  $x_2$  edad en años
- ▶ Dos clasificadores  $\hat{f}$ 's: convexo-cuadrático (línea negra) y no-convexo (línea roja)

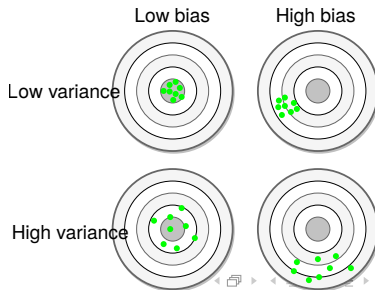
# El tradeoff entre Sesgo y Variance I



$$\mathbb{E} \left[ (y - \hat{f}(x))^2 \right] = (\text{Bias} [\hat{f}(x)])^2 + \text{Var} [\hat{f}(x)] + \sigma^2$$

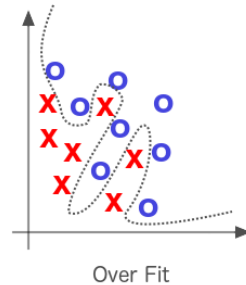
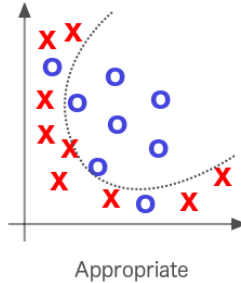
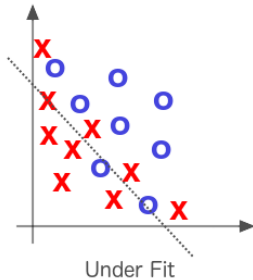
$$\text{Bias} [\hat{f}(x)] = \mathbb{E} [\hat{f}(x) - f(x)]$$

$$\text{Var} [\hat{f}(x)] = \mathbb{E} [\hat{f}(x)^2] - (\mathbb{E} [\hat{f}(x)])^2$$



# El tradeoff entre Sesgo y Variance II

## Clasificación

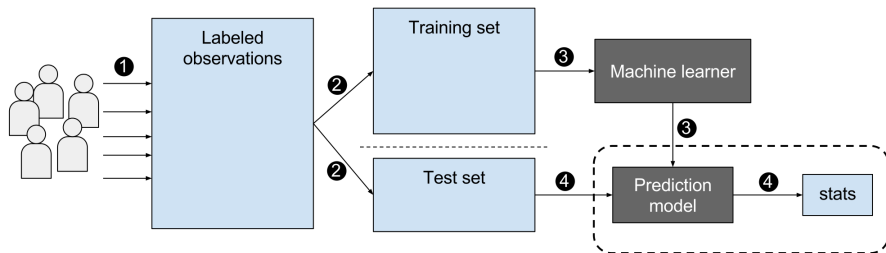


# Evaluando desempeño del modelo I

Para evaluar el desempeño de un modelo generalmente se utilizan algunas de las dos siguientes medidas

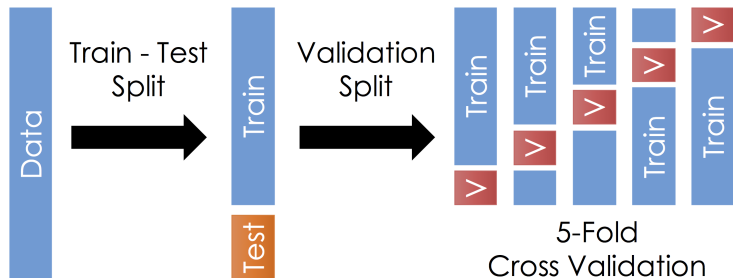
- ▶ **Modelo explicativo:**  $R^2$  (o considerando complejidad del modelo  $R^2_{adj}$ )
- ▶ **Modelo predictivo:**  $RSS$

Para estimar el  $RSS$  necesitamos particionar los datos en: (1) datos de entrenamiento y (2) datos de prueba



# Validación cruzada (cross validation) I

Para obtener estimaciones de la distribución de los estimadores (e.g. precisión de la predicción) podemos utilizar la validación cruzada (cross validation)



Esto nos permite tener estimaciones de intervalos (como intervalos de confianza)

# Evaluando desempeño en clasificación binaria I

		Condition (as determined by "Gold standard")			
Total population		Condition positive	Condition negative	Prevalence = $\frac{\Sigma \text{ Condition positive}}{\Sigma \text{ Total population}}$	
Test outcome	Test outcome positive	True positive	False positive (Type I error)	Positive predictive value (PPV, Precision) = $\frac{\Sigma \text{ True positive}}{\Sigma \text{ Test outcome positive}}$	False discovery rate (FDR) = $\frac{\Sigma \text{ False positive}}{\Sigma \text{ Test outcome positive}}$
	Test outcome negative	False negative (Type II error)	True negative	False omission rate (FOR) = $\frac{\Sigma \text{ False negative}}{\Sigma \text{ Test outcome negative}}$	Negative predictive value (NPV) = $\frac{\Sigma \text{ True negative}}{\Sigma \text{ Test outcome negative}}$
Positive likelihood ratio (LR+) = $\frac{\text{TPR}}{\text{FPR}}$		True positive rate (TPR, Sensitivity, Recall) = $\frac{\Sigma \text{ True positive}}{\Sigma \text{ Condition positive}}$	False positive rate (FPR, Fall-out) = $\frac{\Sigma \text{ False positive}}{\Sigma \text{ Condition negative}}$	Accuracy (ACC) = $\frac{\Sigma \text{ True positive} + \Sigma \text{ True negative}}{\Sigma \text{ Total population}}$	
Negative likelihood ratio (LR-) = $\frac{\text{FNR}}{\text{TNR}}$		False negative rate (FNR) = $\frac{\Sigma \text{ False negative}}{\Sigma \text{ Condition positive}}$	True negative rate (TNR, Specificity, SPC) = $\frac{\Sigma \text{ True negative}}{\Sigma \text{ Condition negative}}$		
Diagnostic odds ratio (DOR) = $\frac{\text{LR+}}{\text{LR-}}$					



# Matriz de confusión, accuracy y coeficiente (cohen's) kappa I

		Real		Total
		Positivo	Negativo	
Predicción	Positivo	$a$	$b$	$a + b$
	Negativo	$c$	$d$	$c + d$
Total		$a + c$	$b + d$	$N$

Tenemos que el Accuracy ( $ACC$ )

$$ACC = \frac{a + d}{N}$$

$$\text{Error} = 1 - ACC$$

# Matriz de confusión, accuracy y coeficiente (cohen's) kappa II

El coeficiente (cohen's) kappa  $\kappa$  mide el acuerdo entre dos fuentes donde cada una clasifica  $N$  items en  $C$  clases mutuamente excluyentes.

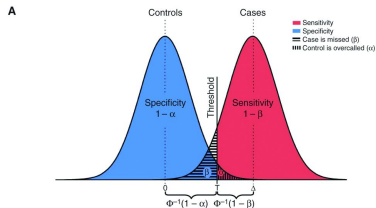
$$\kappa \equiv \frac{ACC - p_e}{1 - p_e},$$

donde  $p_e$  es la probabilidad hipotética de acuerdos al azar.

		A		Total
		Si	No	
B	Si	$a$	$b$	$a + b$
	No	$c$	$d$	$c + d$
Total		$a + c$	$b + d$	$N$

- ▶  $ACC = \frac{a+d}{N}$
- ▶  $p_{Si} = \frac{a+b}{N} \cdot \frac{a+c}{N}$
- ▶  $p_{No} = \frac{c+d}{N} \cdot \frac{b+d}{N}$
- ▶  $p_e = p_{Si} + p_{No}$

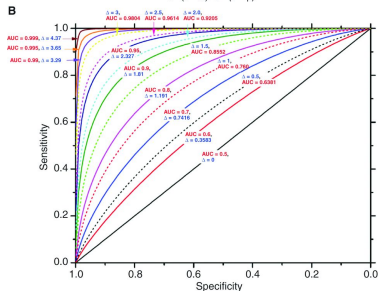
# Curva ROC (receiver operating characteristic) I



Predicción

Real

	Real		Total
	Positivo	Negativo	
Positivo	$a$	$b$	$a + b$
Negativo	$c$	$d$	$c + d$
Total	$a + c$	$b + d$	$N$



- Specificity =  $\frac{d}{b+d}$
- Sensitivity =  $\frac{a}{a+c}$