



Instituto Tecnológico  
de Buenos Aires

18/OCTUBRE

**SELECCIÓN Y**

**EVALUACIÓN  
DE MODELOS**

—

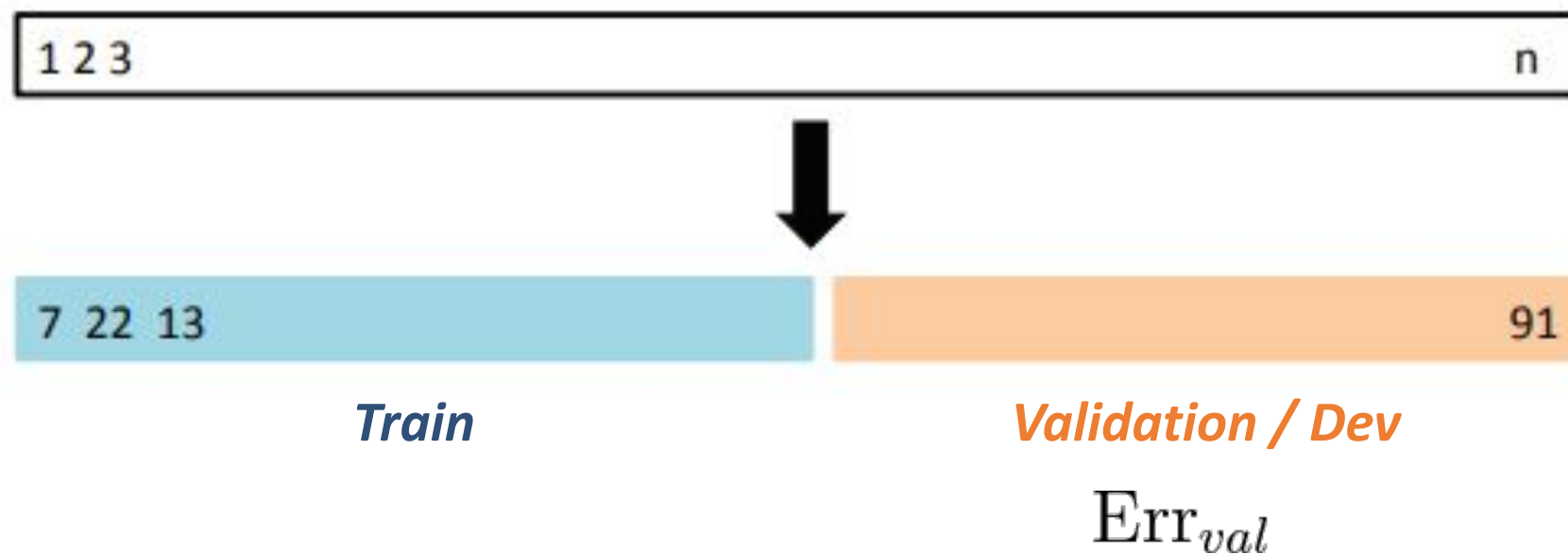
# Selección de modelos

Supongamos que ya partimos los datos en train - test  
(¿Por qué hacíamos esto?)

⇒ ¿Qué hacemos durante el desarrollo de modelos para  
**seleccionar modelos?**

# Train-Validation

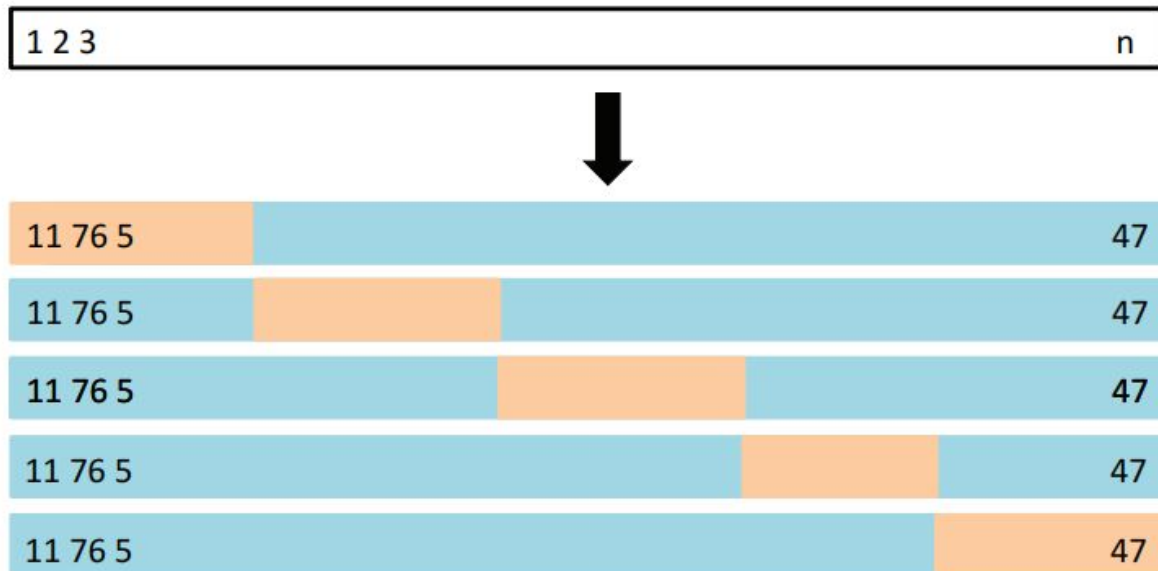
Entrenamos en unos datos (**train**) y evaluamos en otros (**validation/dev**)  
(de acá en adelante, vamos a asumir que el conjunto de test ya fue separado)



# k-Fold Cross-Validation

Cuando hay pocos datos, la **estimación en único validation set puede ser muy variable**  
 ⇒ ¡Intentemos usar todos los datos!

Ejemplo: k-fold CV con k=5



→ Para  $i=1, \dots, 5$ :

- Entrenamos en los **datos restantes** y predecimos en **fold  $i$**

→ Calculamos el **error**  $Err_{cvk}$  con la predicción de cada obs.

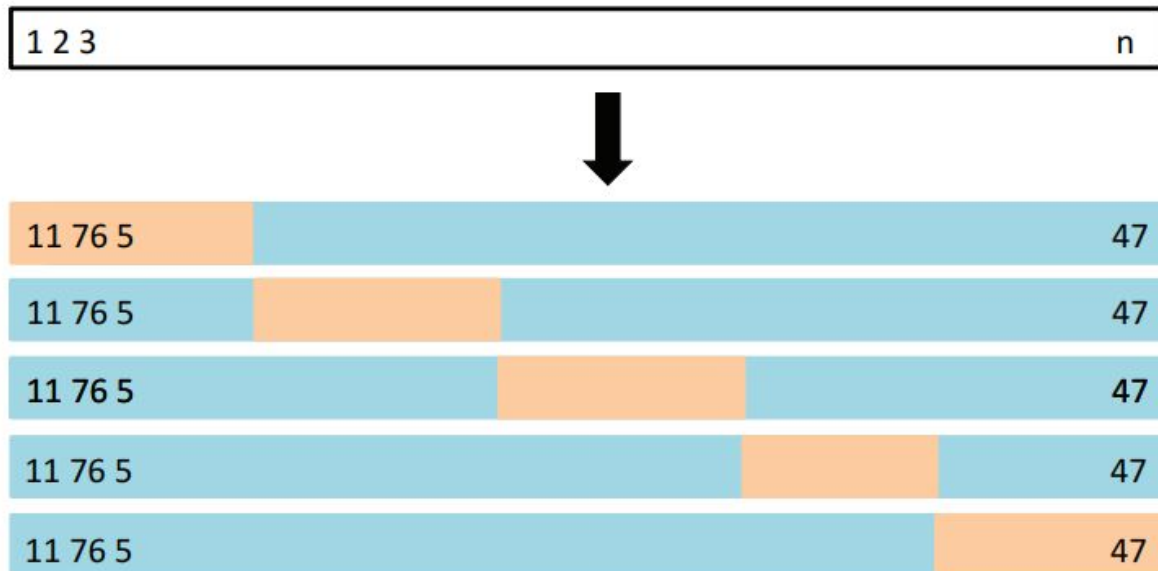
$Err_{cvk}$

¿Es una estimación *confiable* del error de mi mejor modelo?

# k-Fold Cross-Validation

Cuando hay pocos datos, la **estimación en único validation set puede ser muy variable**  
 ⇒ ¡Intentemos usar todos los datos!

Ejemplo: k-fold CV con k=5



→ Para  $i=1, \dots, 5$ :

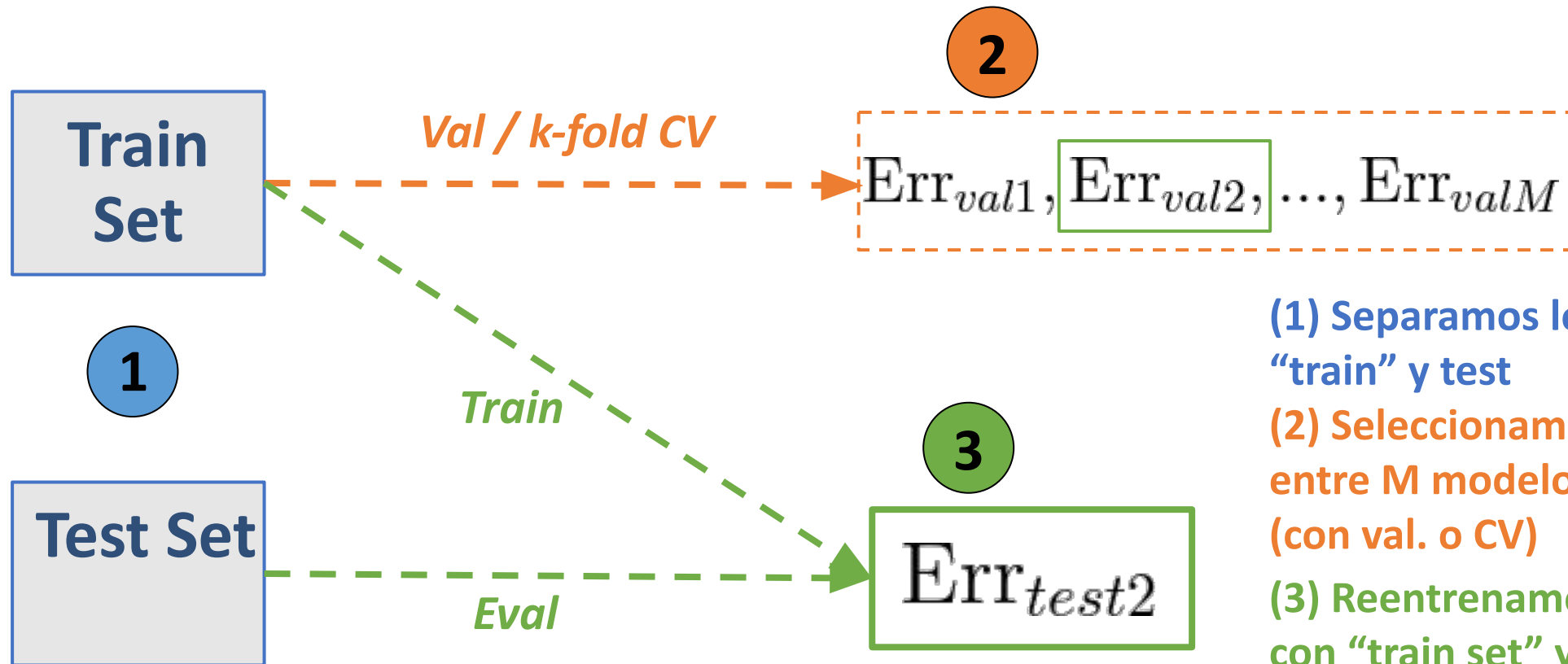
- Entrenamos en los **datos restantes** y predecimos en **fold  $i$**

→ Calculamos el **error**  $Err_{cvk}$  con la predicción de cada obs.

$Err_{cvk}$

¿Es una estimación *confiable* del error de mi mejor modelo?  
**¡NO!** solo nos interesa para encontrar el mínimo

# Selección de modelos

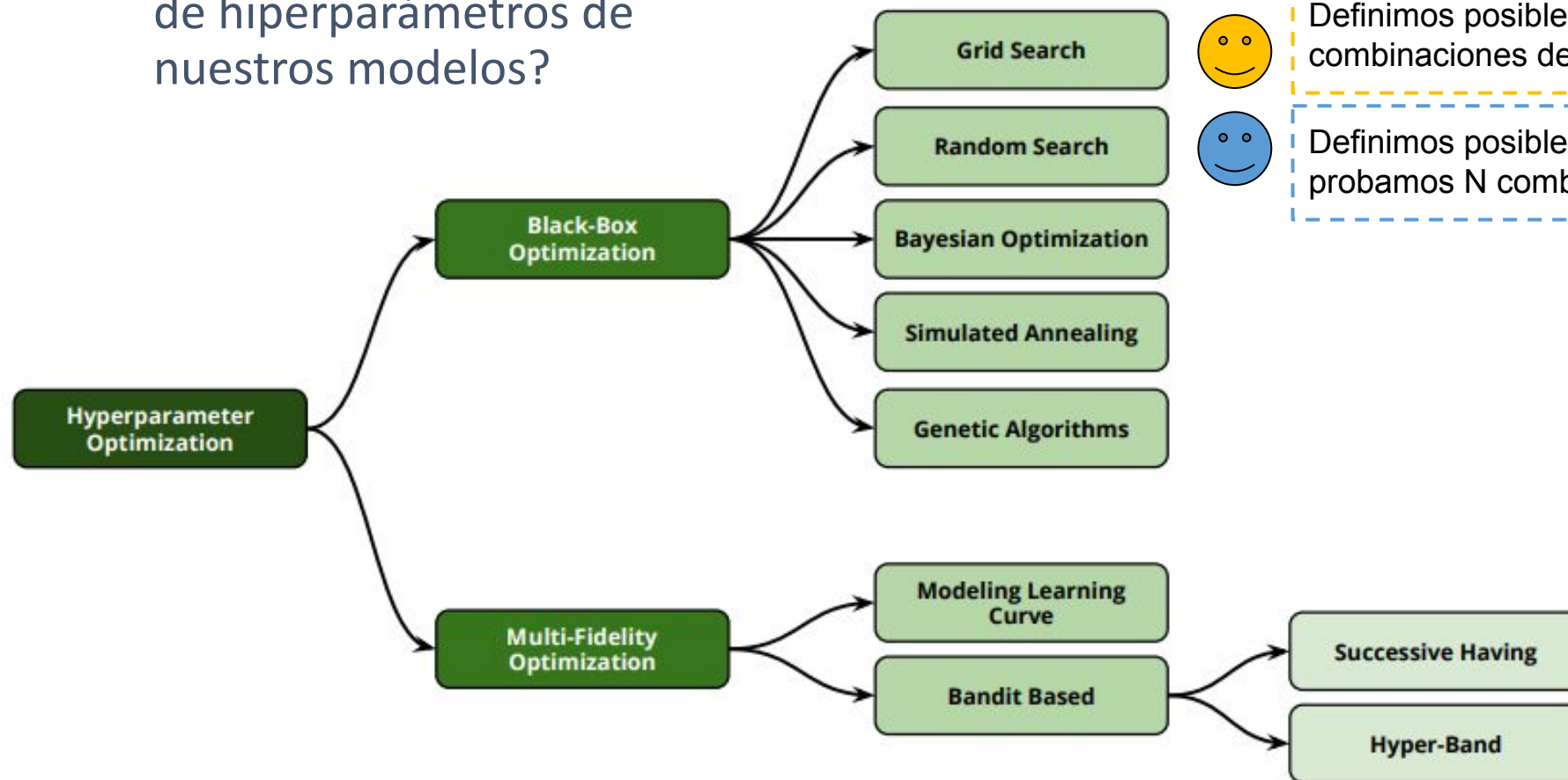


- (1) Separamos los datos en “train” y test
- (2) Seleccionamos uno de entre M modelos según  $Err_{vali}$  (con val. o CV)
- (3) Reentrenamos el modelo con “train set” y lo evaluamos en test

**IMPORTANTE:** el criterio de partición es el mismo en (1) y en (2)

# Optimización de Hiperparámetros (HPO)

¿Cómo recorremos el espacio de hiperparámetros de nuestros modelos?



# Evaluación de modelos

OK pero... ¿qué es el *error* de un modelo...?

¿Cómo medimos el rendimiento / performance de un modelo?  
(ya sea en *validation*, *CV*, o en *test*)



# Métricas de clasificación

$$y \in \{0, 1\} \quad \hat{y} \in \{0, 1\}$$

Supuesto: nuestro sistema **clasifica** entre 0 y 1

		<i>True class</i>		
		– or Null	+ or Non-null	Total
<i>Predicted class</i>	– or Null	True Neg. (TN)	False Neg. (FN)	N*
	+ or Non-null	False Pos. (FP)	True Pos. (TP)	P*
	Total	N	P	

**Accuracy** (*tasa de aciertos*)

$$\text{Acc} = \frac{\text{TP} + \text{TN}}{\text{P} + \text{N}} = 1 - \text{Err} = \frac{1}{n} \sum I(y_i = \hat{y}_i)$$

## Métricas de clasificación

Ejemplo: clasificación de default (paga / no paga)

		Default obs.		
		0	1	Total
Default pred.	0	4 988	3 784	8 772
	1	12	1 216	1 228
	Total	5 000	5 000	10 000

$FP = \dots$

$FN = \dots$

$Acc = \dots$

## Métricas de clasificación

Ejemplo: clasificación de default (paga / no paga)

		Default obs.		
		0	1	Total
Default pred.	0	4 988	3 784	8 772
	1	12	1 216	1 228
	Total	5 000	5 000	10 000

$$FP = 12$$

$$FN = 3784$$

$$Acc = 62\%$$

$$FPR = FP / N = 0.24\%$$

(tasa de falsos pos.)

$$FNR = FN / P = 75.7\%$$

(tasa de falsos neg.)

## Métricas de clasificación

Ejemplo: clasificación de default (paga / no paga) pero un poco más realista...

		Default obs.		
		0	1	
Default pred.	0	9 644	252	9 896
	1	23	81	104
		9 667	333	10 000

$$FP = \dots$$

$$FN = \dots$$

$$Acc = \dots$$

$$FPR = \dots$$

$$FNR = \dots$$

## Métricas de clasificación

Ejemplo: clasificación de default (paga / no paga) pero un poco más realista...

		Default obs.		
		0	1	
Default pred.	0	9 644	252	9 896
	1	23	81	104
		9 667	333	10 000

$$FP = 23$$

$$FN = 252$$

$$Acc = 97.3\%$$

$$FPR = 0.24\%$$

$$FNR = 75.7\%$$

En problemas “desbalanceados” (con una clase mayoritaria) el accuracy es poco informativo

¿Cuánto es Acc si nuestro modelo siempre predice 0?

## Métricas de clasificación

Ejemplo: clasificación de default (paga / no paga) pero un poco más realista...

		Default obs.		
		0	1	
Default pred.	0	9 644	252	9 896
	1	23	81	104
		9 667	333	10 000

$$FP = 23$$

$$FN = 252$$

$$Acc = 97.3\%$$

$$FPR = 0.24\%$$

$$FNR = 75.7\%$$

En problemas “desbalanceados” (con una clase mayoritaria) el accuracy es poco informativo

¿Cuánto es Acc si nuestro modelo siempre predice 0?

$$\rightarrow Acc = (9644 + 23 + 0) / 1000 = 96.67\%$$

## Métricas de clasificación

La mayor parte de las veces nos va a interesar evaluar **otras métricas** en lugar de la tasa de aciertos

⇒ hoy vamos ver algunas y más adelante vamos a ver cómo elegir

# Métricas de clasificación

		Obs		
		0	1	
Pred	0	TN	FN	N*
	1	FP	TP	P*
		N	P	

		Default obs.		
		0	1	
Default pred.	0	9 644	252	9 896
	1	23	81	104
		9 667	333	10 000

$$TPR = TP/P = 1 - FNR$$

*Recall /  
Sensibilidad*

$$TNR = TN/N = 1 - FPR$$

*Especificidad*

$$PPV = TP/P^*$$

*Precision  
(Positive predictive value)*

$$Acc = 97.3\%$$

$$Recall = 24.3\% = 1 - 0.757$$

$$Specificity = 99.8\% = 1 - 0.0024$$

$$Precision = 77.9\%$$



# Métricas de clasificación

		Obs		
		0	1	
Pred	0	TN	FN	N*
	1	FP	TP	P*
		N	P	

$$TPR = TP/P = 1 - FNR$$

*Recall /  
Sensibilidad*

$$TNR = TN/N = 1 - FPR$$

*Especificidad*

$$PPV = TP/P^*$$

*Precision  
(Positive predictive value)*

*¿Qué modelo optimiza el recall?*

*¿Qué modelo optimiza la especificidad?*

		Default obs.		
		0	1	
Default pred.	0	9 644	252	9 896
	1	23	81	104
		9 667	333	10 000

## Métricas de clasificación

		Obs		
		0	1	
Pred	0	TN	FN	N*
	1	FP	TP	P*
		N	P	

$$\text{TPR} = \text{TP}/P = 1 - \text{FNR}$$

*Recall /  
Sensibilidad*

$$\text{TNR} = \text{TN}/N = 1 - \text{FPR}$$

*Especificidad*

$$\text{PPV} = \text{TP}/P^*$$

*Precision  
(Positive predictive value)*

		Default obs.		
		0	1	
Default pred.	0	9 644	252	9 896
	1	23	81	104
		9 667	333	10 000

*¿Qué modelo optimiza el recall?*

Predigo siempre 1  $\rightarrow (252+81)/333 = 1$

*¿Qué modelo optimiza la especificidad?*

Predigo siempre 0  $\rightarrow (9644+23)/9667 = 1$

$\Rightarrow$  hay un **trade-off** entre *capturar 1s (recall)* y *capturar 0s (especificidad)*

## Métricas de clasificación

Métricas de la matriz de confusión que pueden reemplazar a accuracy (más adelante vamos a ver cómo elegir)

**F-score** (combinación entre precision y recall)

$$F_{\beta} = \frac{(\beta^2 + 1)PR}{\beta^2 P + R} \quad F_1 = \frac{2PR}{P + R} \quad \begin{array}{ll} \beta = 2 & + \text{recall} \\ \beta = 0.5 & + \text{precision} \end{array}$$

## Función de costos

(a veces podemos definir los costo relativos/absolutos de acertar y/o error)

$$C = C_{TN}TN + C_{TP}TP + C_{FN}FN + C_{FP}FP$$

## Métricas de clasificación

En general los modelos nos devuelven un **puntaje** o **score**, no una clase

$$\hat{p} \in [0, 1]$$

$$\hat{y} = \underset{y}{\operatorname{argmax}} \widehat{p(y|x)}$$
$$(\hat{y} = 1 \iff \hat{p} > 0.5)$$

Esto no necesariamente es óptimo!

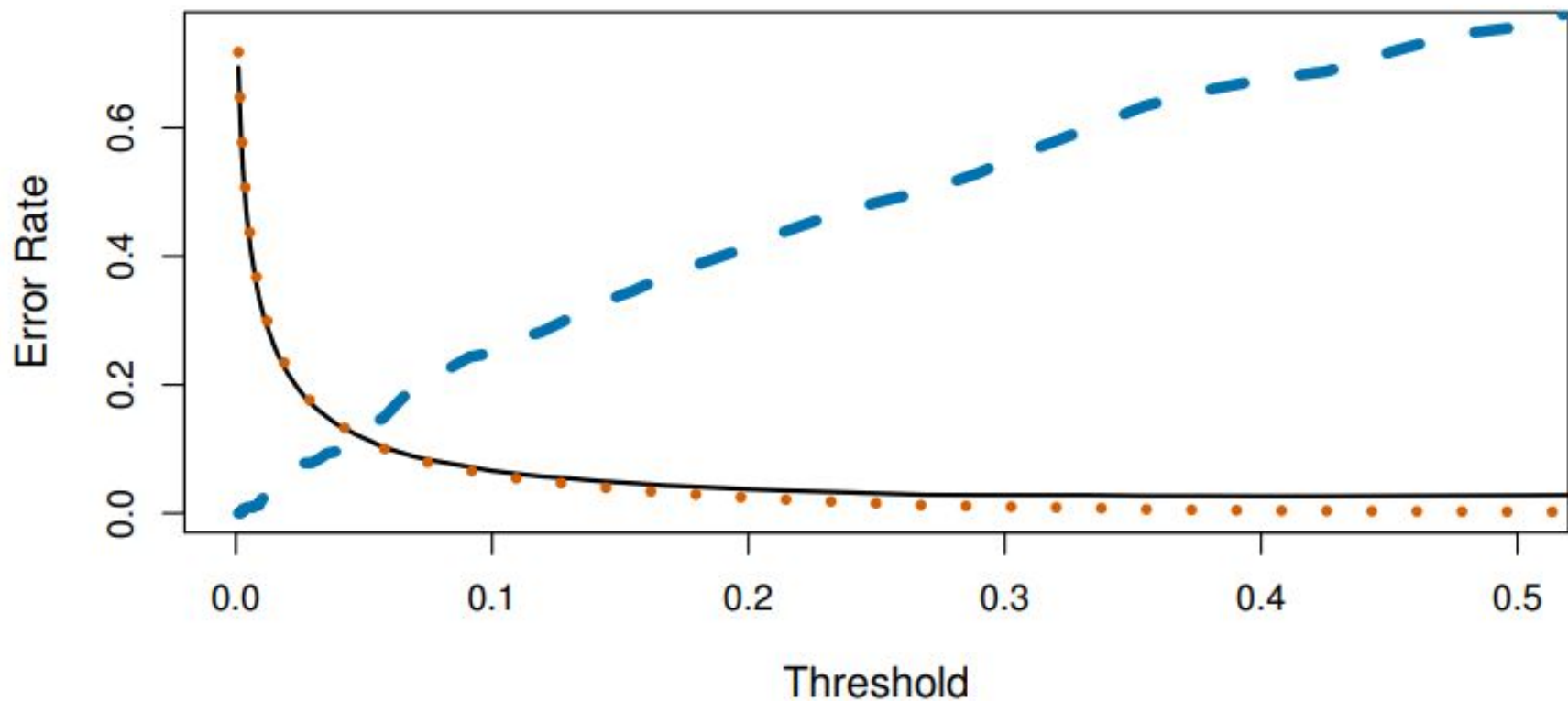
- porque el score puede estar *descalibrado*, y/o
- porque ese umbral puede no optimizar nuestra métrica

⇒ Tenemos que elegir un **umbral** que optimice la métrica que nos interesa

## Métricas de clasificación

$$\hat{y} = \begin{cases} 1 & \text{si } \hat{p} \geq \theta \\ 0 & \text{si } \hat{p} < \theta \end{cases}$$

*Por ejemplo, bajar el umbral puede reducir Accuracy pero mejorar Recall  $\rightarrow$  y entonces, quizás, la métrica que nos interese*

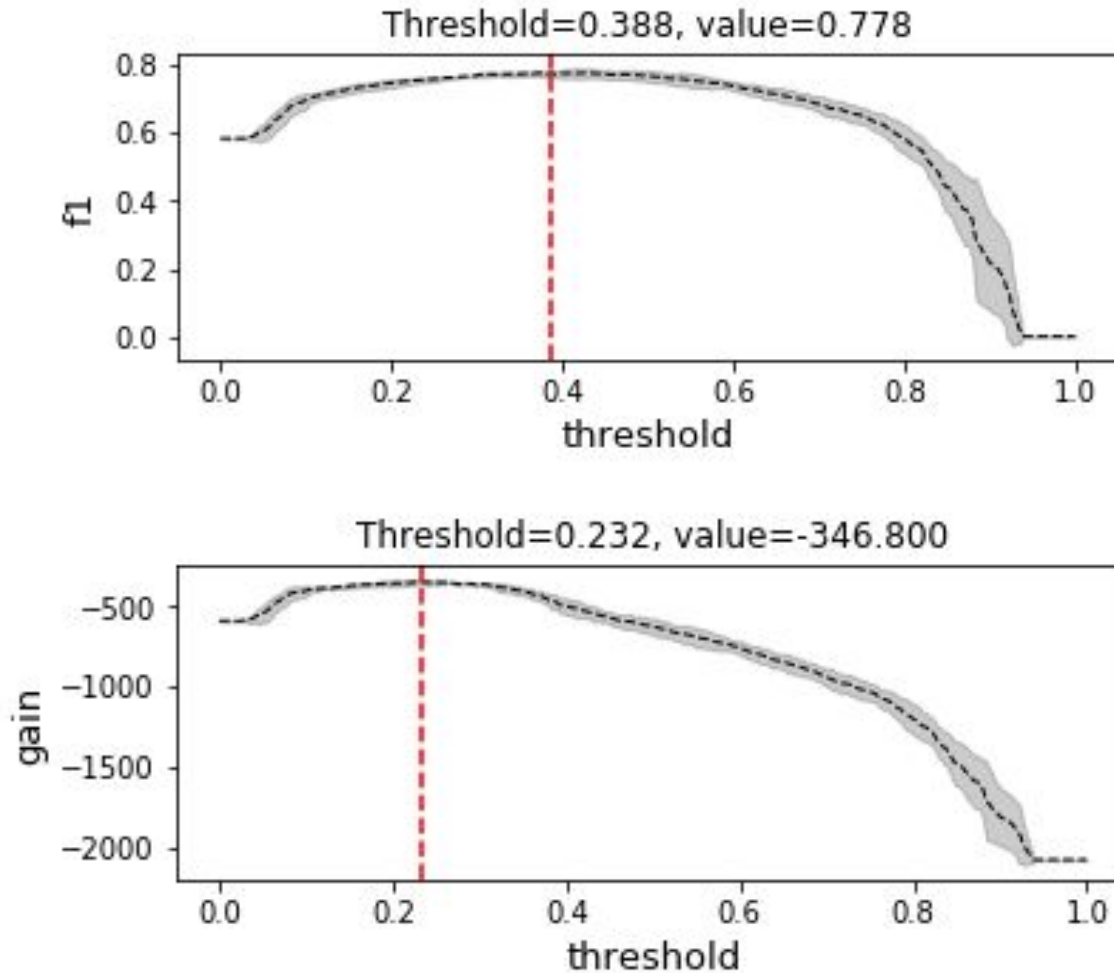


**FNR = 1 - Recall**

**Error rate = 1 - Accuracy**

**FPR = 1 - Specificity**

# Métricas de clasificación



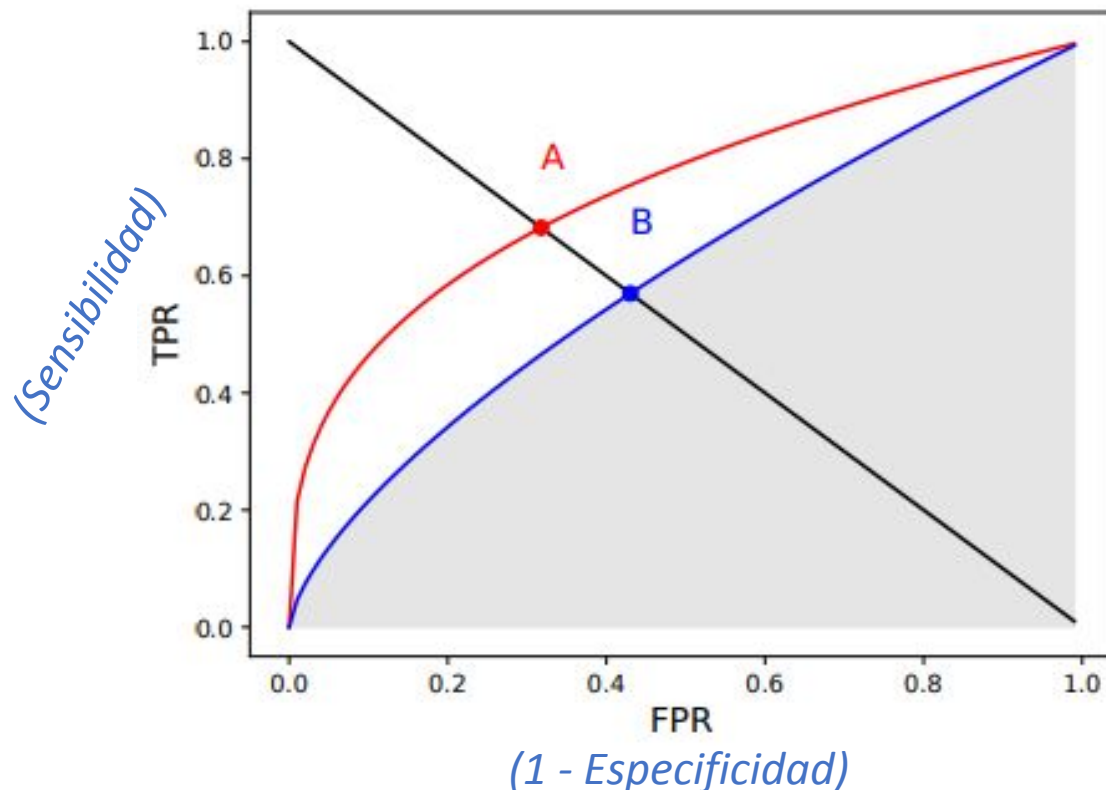
$$\hat{y} = \begin{cases} 1 & \text{si } \hat{p} \geq \theta \\ 0 & \text{si } \hat{p} < \theta \end{cases}$$

*Por ejemplo, bajar el umbral puede reducir Accuracy pero mejorar Recall  
→ y entonces, quizás, la métrica que nos interese*

# Métricas de clasificación

$$y \in \{0, 1\} \quad \hat{p} \in [0, 1]$$

A veces nos interesa evaluar *globalmente* los scores i.e. para cualquier umbral posible



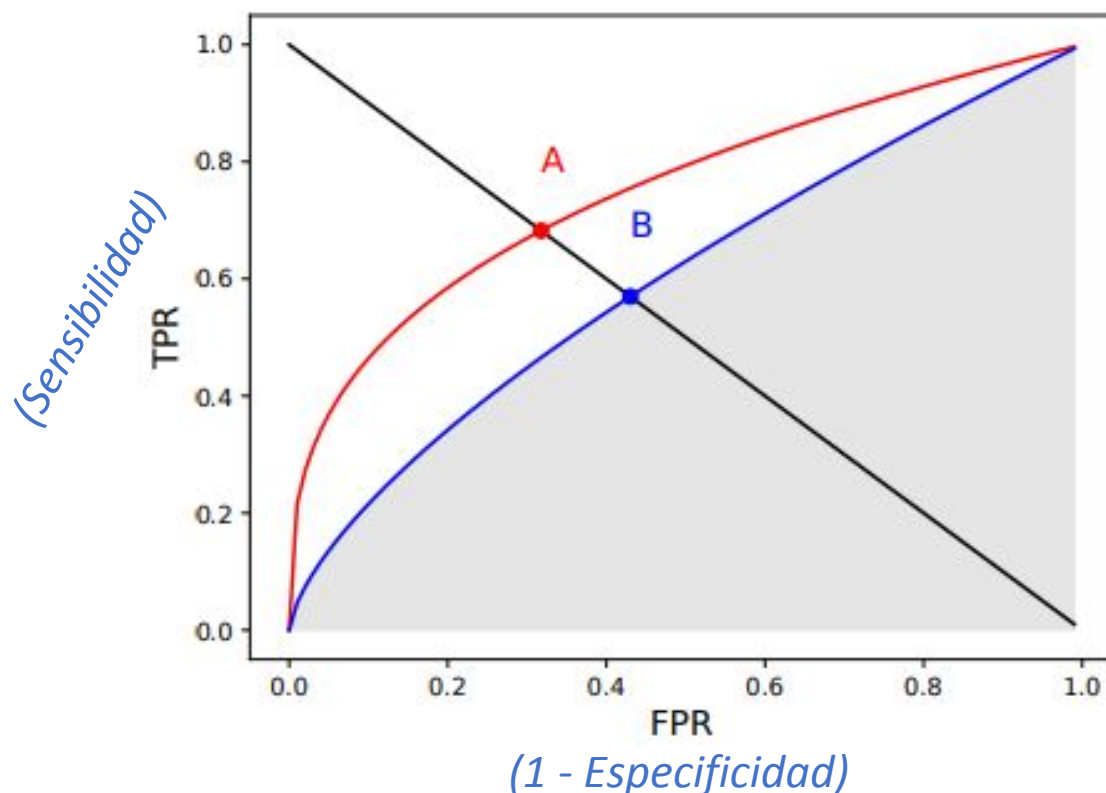
## AUROC (Area Under the ROC Curve)

- Mide la **capacidad global de discriminar** entre 0s y 1s
- Se puede interpretar como la probabilidad de que el score de una obs. positiva seleccionada al azar sea más alto que el de una negativa

# Métricas de clasificación

$$y \in \{0, 1\} \quad \hat{p} \in [0, 1]$$

A veces nos interesa evaluar *globalmente* los scores i.e. para cualquier umbral posible



**AUROC**  
(Area Under the ROC Curve)

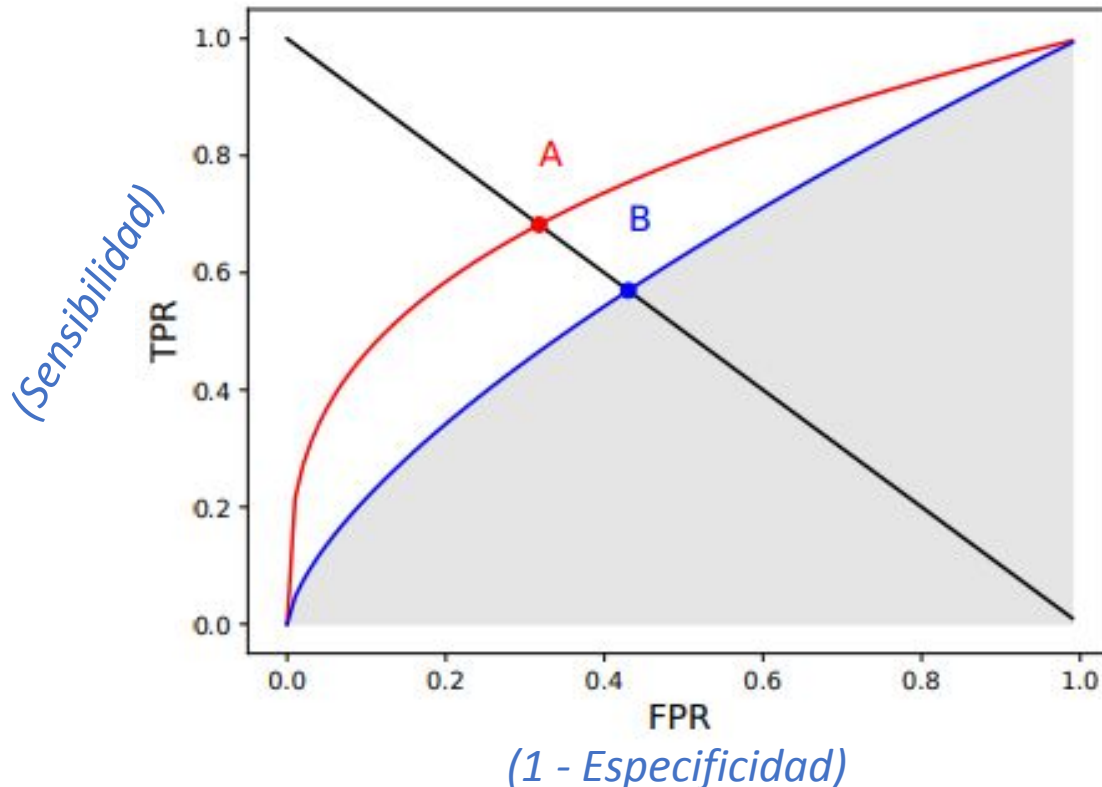
¿A qué valor tiende el AUROC de un score random?



# Métricas de clasificación

$$y \in \{0, 1\} \quad \hat{p} \in [0, 1]$$

A veces nos interesa evaluar *globalmente* los scores i.e. para cualquier umbral posible



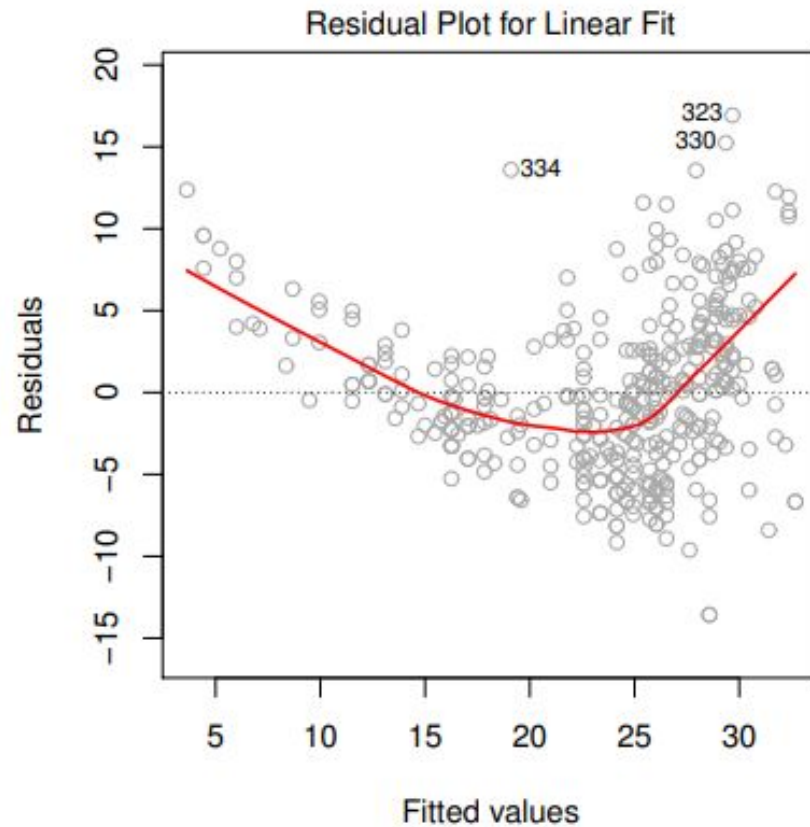
**AUROC**  
(Area Under the ROC Curve)

¿A qué valor tiende el AUROC de un score random?

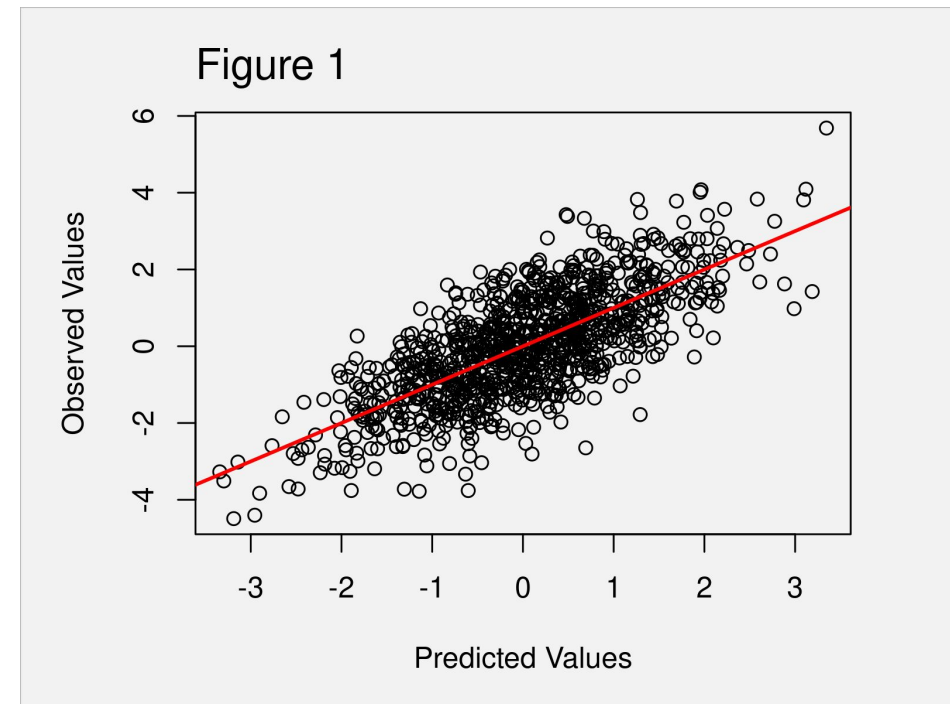
0.5 (para cualquier punto de corte, me equivoco igual  $\rightarrow$  FPR  $\approx$  FNR  $\rightarrow$  diagonal imaginaria en el gráfico)

# Métricas de regresión

## Gráfico de residuos



## Gráfico de correlación



## Métricas de regresión

$$\text{RMSE}(y, \hat{y}) = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}$$

**Root mean squared error**  
(error cuadrático)

$$\text{MAE}(y, \hat{y}) = \frac{1}{n} \sum_{i=0}^n |y_i - \hat{y}_i|$$

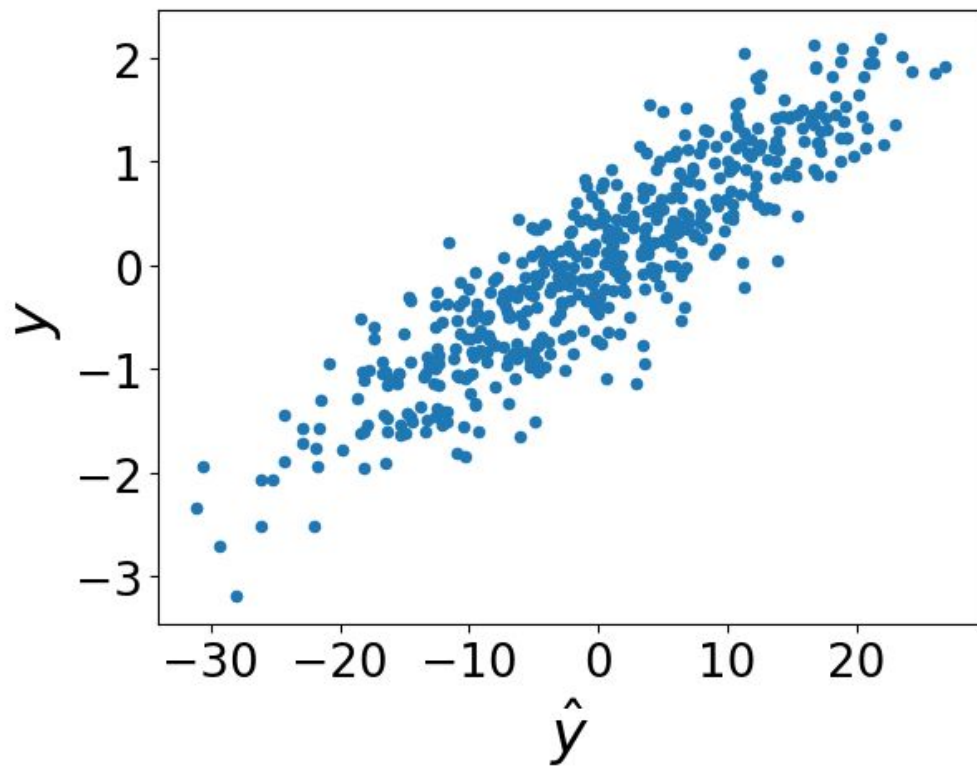
**Mean absolute error**  
(error absoluto)

$$\text{MAPE}(y, \hat{y}) = \frac{1}{n} \sum_{i=1}^n \frac{|y_i - \hat{y}_i|}{\max(\epsilon, |y_i|)}$$

**Mean absolute percentage error**  
(error porcentual)

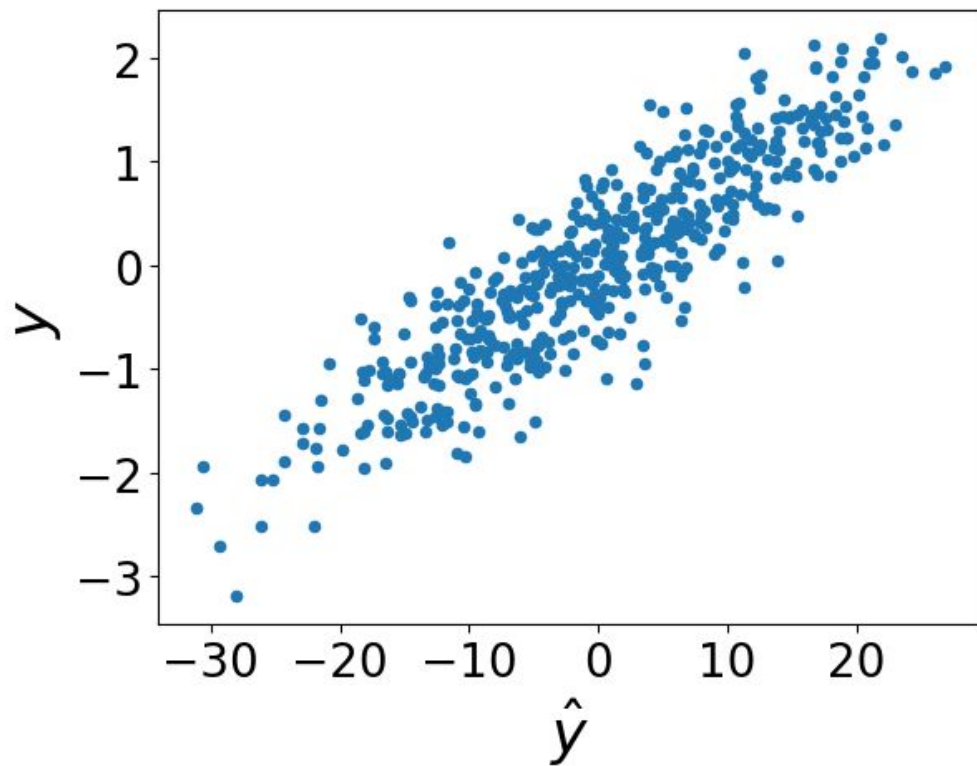
## Métricas de regresión

¿Cómo nos darían las métricas anteriores en este ejemplo?  
(bien / mal / regular / kcyo)



## Métricas de regresión

¿Cómo nos darían las métricas anteriores en este ejemplo?  
(bien / mal / regular / kcyo)



¿Y si usamos el viejo y querido R cuadrado?  
*Veamos... →*

# Métricas de regresión

## R cuadrado

$$(A) \quad R^2(y, \hat{y}) = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

$$(B) \quad R^2(y, \hat{y}) = \text{pearson}(y, \hat{y})^2$$

- Podemos interpretar  $R^2(A)$  como la comparación entre nuestro modelo ( $y_{\hat{}}$ ) y un modelo baseline ( $y_{\text{mean}}$ )
- $R^2$  está basado en MSE

*“(A) y (B) son iguales ...”*

*“ $0 < R^2(A) < 1$  ...”*

# Métricas de regresión

## R cuadrado

$$(A) \quad R^2(y, \hat{y}) = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

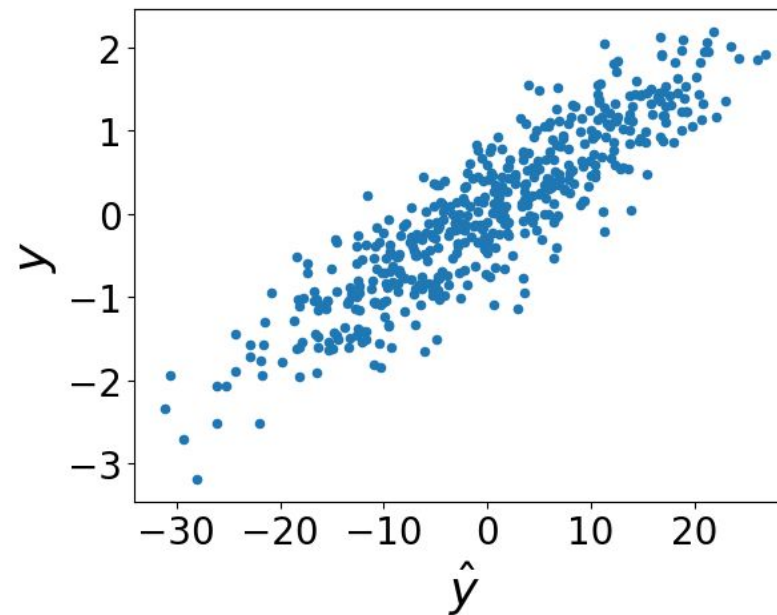
$$(B) \quad R^2(y, \hat{y}) = \text{pearson}(y, \hat{y})^2$$

- Podemos interpretar  $R^2(A)$  como la comparación entre nuestro modelo ( $y_{\hat{}}$ ) y un modelo baseline ( $y_{\text{mean}}$ )
- $R^2$  está basado en MSE

“(A) y (B) son iguales ...”

“ $0 < R^2(A) < 1$  ...”

en los datos de entrenamiento de MCO!



$$R^2(B) = 0.80$$
$$R^2(A) = -108$$