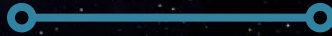




# **ASTRONOMÍA Y CIENCIA DE DATOS: DE LAS ESTRELLAS A LOS NÚMEROS**

Clase 8: Introducción al Aprendizaje Supervisado II:  
Random Forest



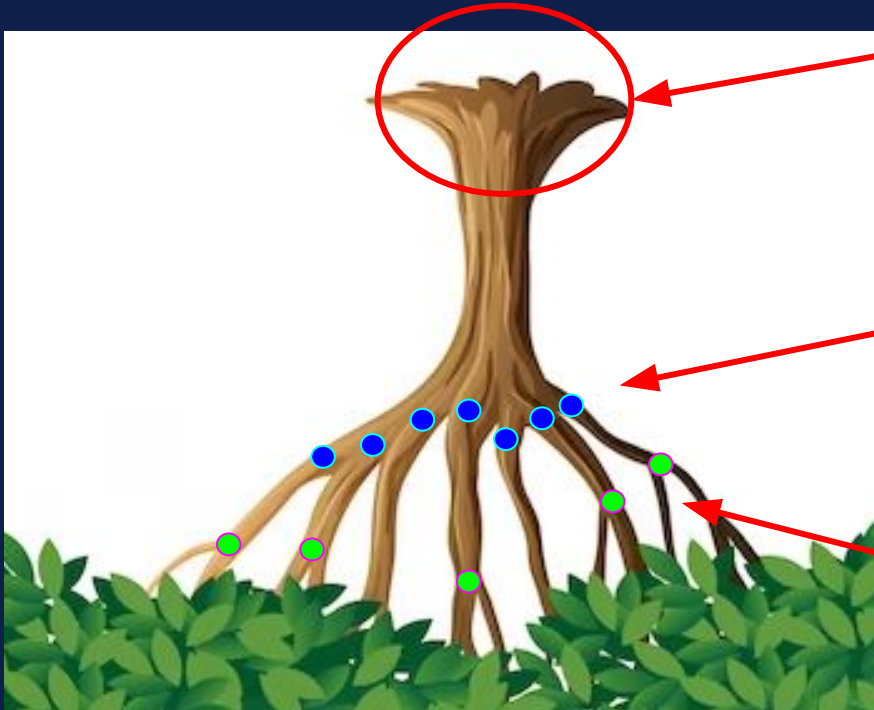
# Árboles de decisión



# Árboles de decisión



# Árboles de decisión



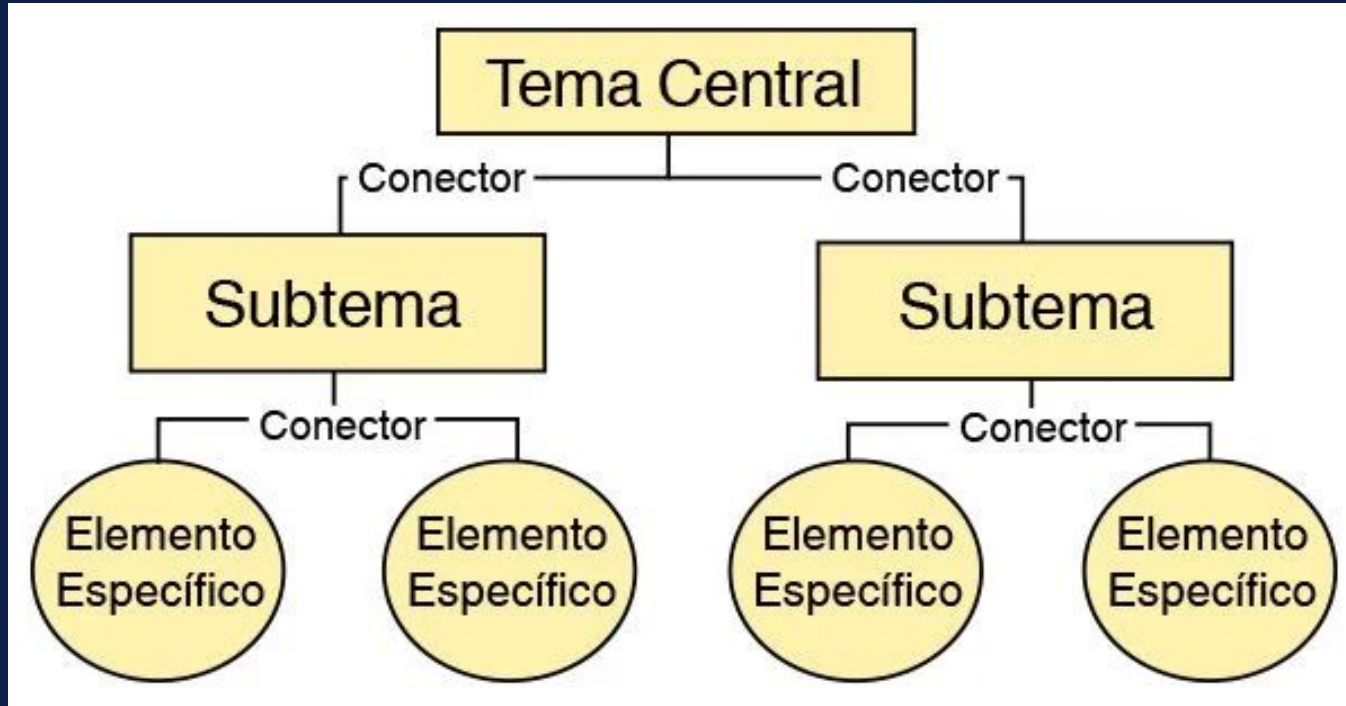
Base del tronco

Bifurcación en ramas principales

Bifurcación en ramas secundarias

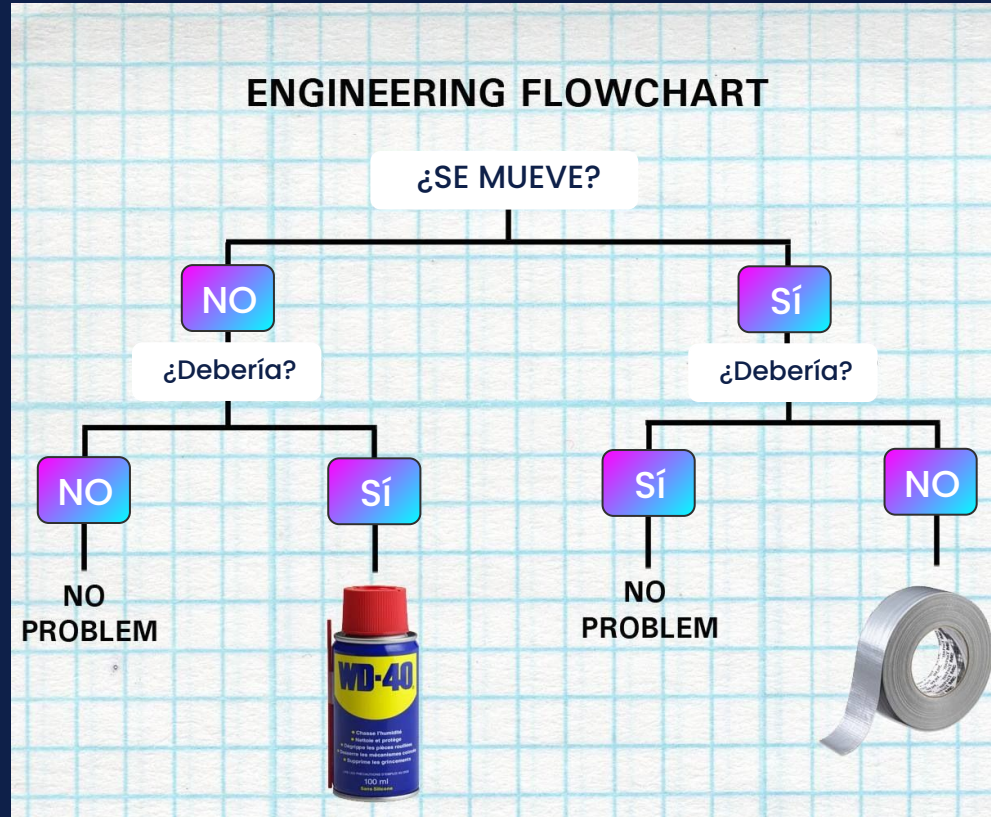
# Árboles de decisión

Son de hecho bastante similares a un mapa conceptual

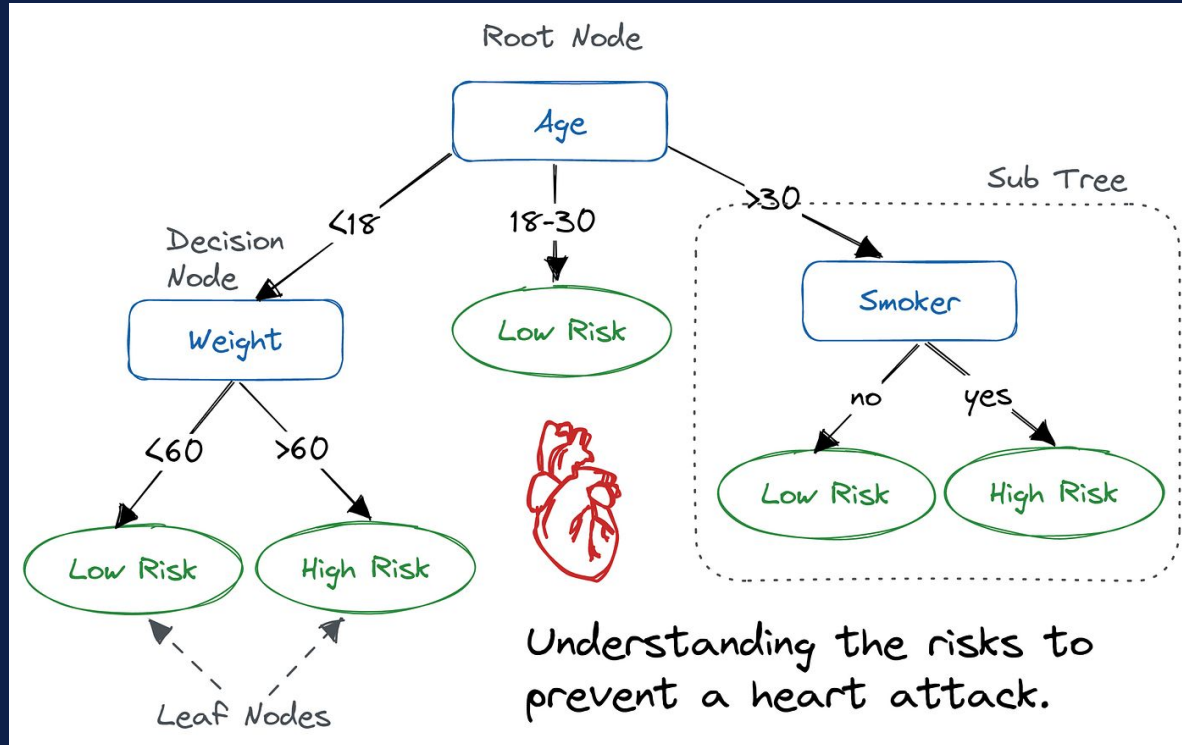




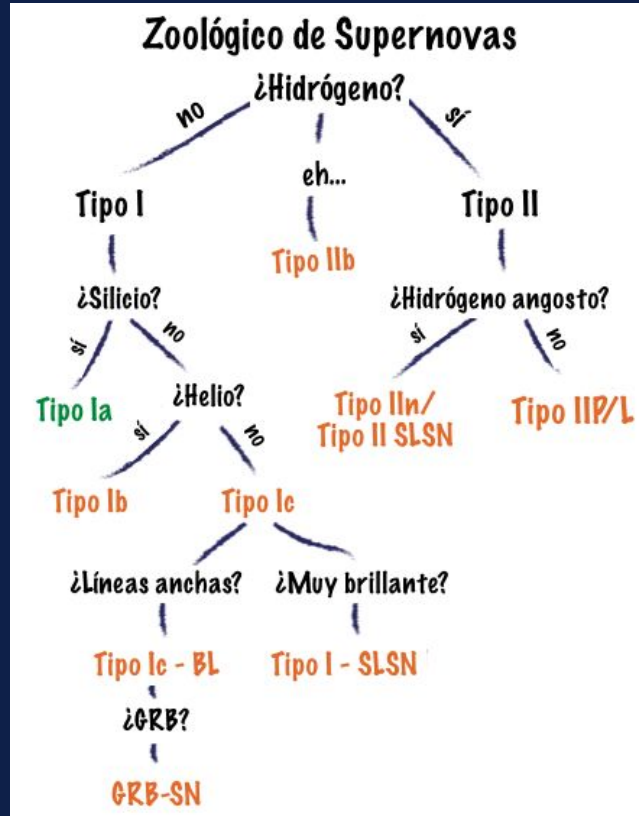
# Árboles de decisión



# Árboles de decisión



# Árboles de decisión





# Características de los DT

- Fáciles de interpretar
- Rápidos
- Deterministas → **(SÍ o NO)**
- Sensible a outliers (casos raros)
- Comúnmente **se sobreajusta a los datos.**

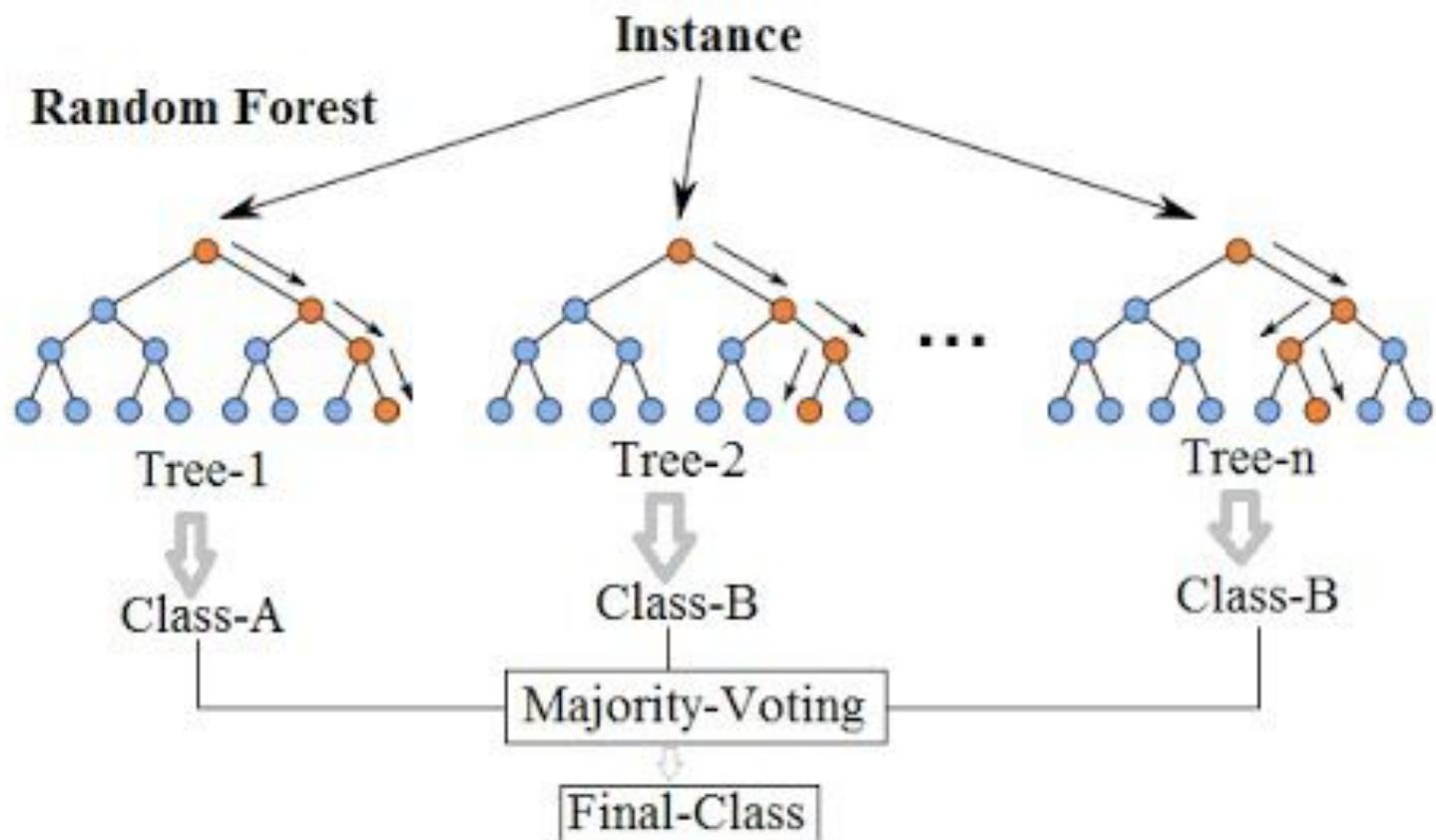
# Características de los DT

- Fáciles de interpretar
- Rápidos
- Deterministas → (SÍ o NO)
- Sensible a outliers (casos raros)
- Comúnmente **se sobreajusta a los datos.**

Cómo podemos  
compensar los  
e  
u  
CON MÁS  
ÁRBOLES!  
perder sus  
ventajas?



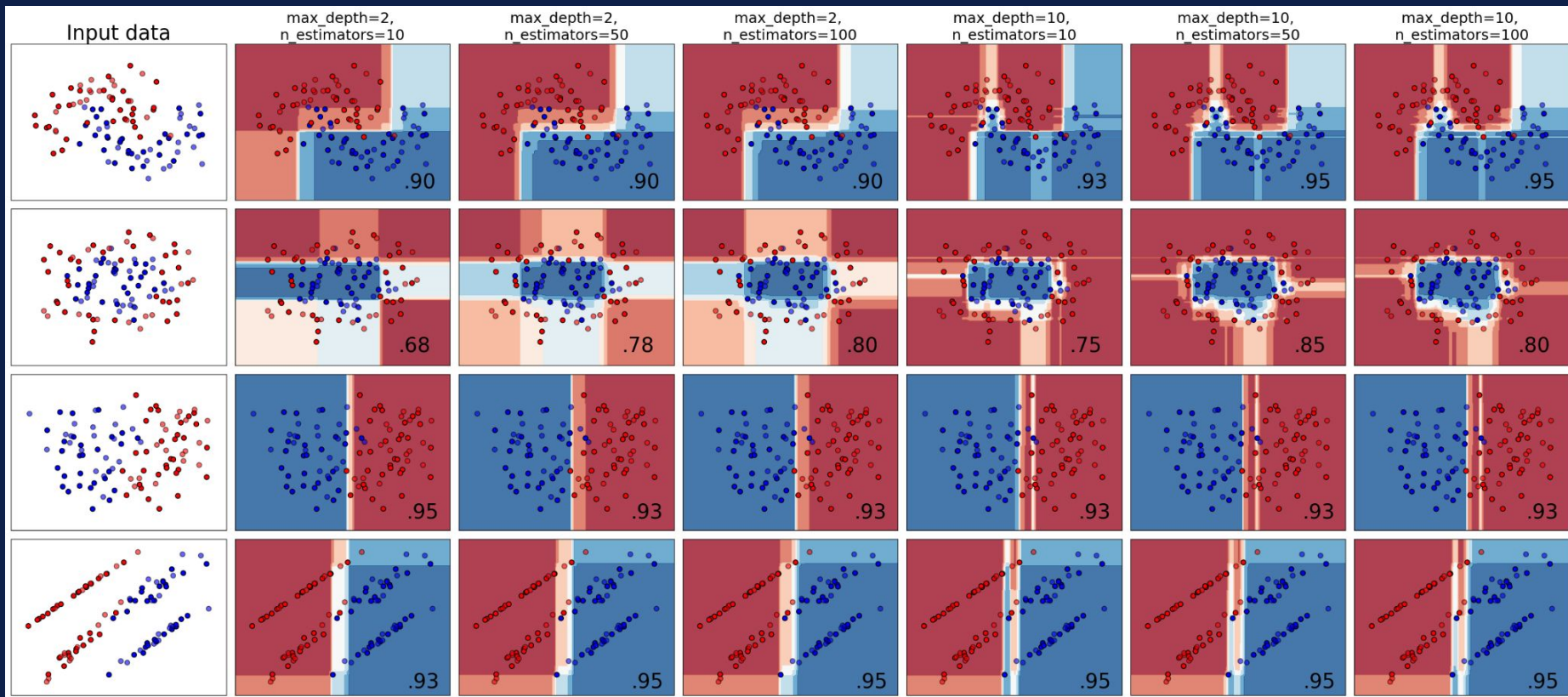
# Random Forest Simplified



# Random Forest

- número de árboles (*n\_estimators*)
- Características a considerar (*max\_features*)
- Número máximo de sub-ramas (*max\_depth*)

# Random Forest





# Características de los DT

- Fáciles de interpretar
- Rápidos (en entrenamiento y aplicaciones)
- Deterministas → (SÍ o NO)
- Sensible a outliers (casos raros)
- Comúnmente **se sobreajusta a los datos.**

# Características de RF

- Se complica la interpretación
- Rápidos (en entrenamiento y aplicaciones)
- No Deterministas → (**entrega un porcentaje de confianza**, ej 83%)
- Más **robusto** con outliers
- El modelo se vuelve **más general a los datos**.

# Árboles v/s Random Forest

	Árboles	Random Forest
<i>Interpretabilidad</i>	✓ ✓	✓ ✗
<i>Precisión</i>	✓ ✗	✓ ✓
<i>Overfitting</i>	✗	✓ ✓
<i>Robustez</i>	✗ ✗	✓ ✓
<i>Velocidad</i>	✓ ✓	✓
<i>Clasificación</i>		
<i>Regresión</i>		

y aún quedan un montón de  
modelos que no hemos  
mencionado!

# Random Forest

Un montón de  
modelos que no hemos  
mencionado!

Random  
Forest

un montón de  
modelos que no hemos  
mencionado!

Redes  
Neuronales



Random  
Forest

un montón de  
modelos que no hemos  
hecho!

Super Vector  
Machine

Redes  
Neuronales

Random  
Forest

modelos de

Super Vector  
Machine

Proceso  
Gaussiano

hado!

Redes  
Neuronales

Random  
Forest

Proceso  
Gaussiano

modelos de  
Super Vectorizado!

M

QDA y LA

des  
onales

Random  
Forest

Proceso

modelos

Super Vec

M

Nearest  
Neighbors

QDA y LA

des

onales

AdaBoost

Random

Proceso

Super Vec

Nearest  
Neighbors

M

QDA y LA

des

onales

Random

AdaBoost

Processo

Super Vec

Nearest

Transformers

QDA

naïves



Teniendo tantas  
opciones, ¿cómo  
**evaluamos** qué  
modelo es mejor?

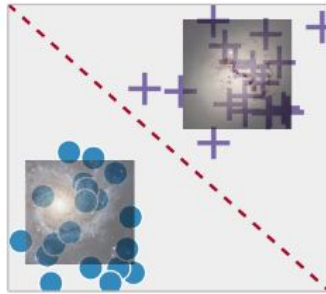
# Métricas de evaluación

Dependen del tipo de problema que  
queremos resolver

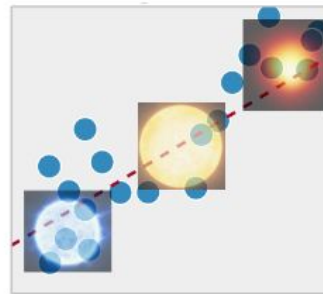
Supervisado



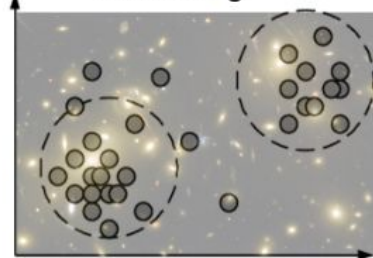
Classification



Regression



Clustering



No supervisado



# Métricas de regresión

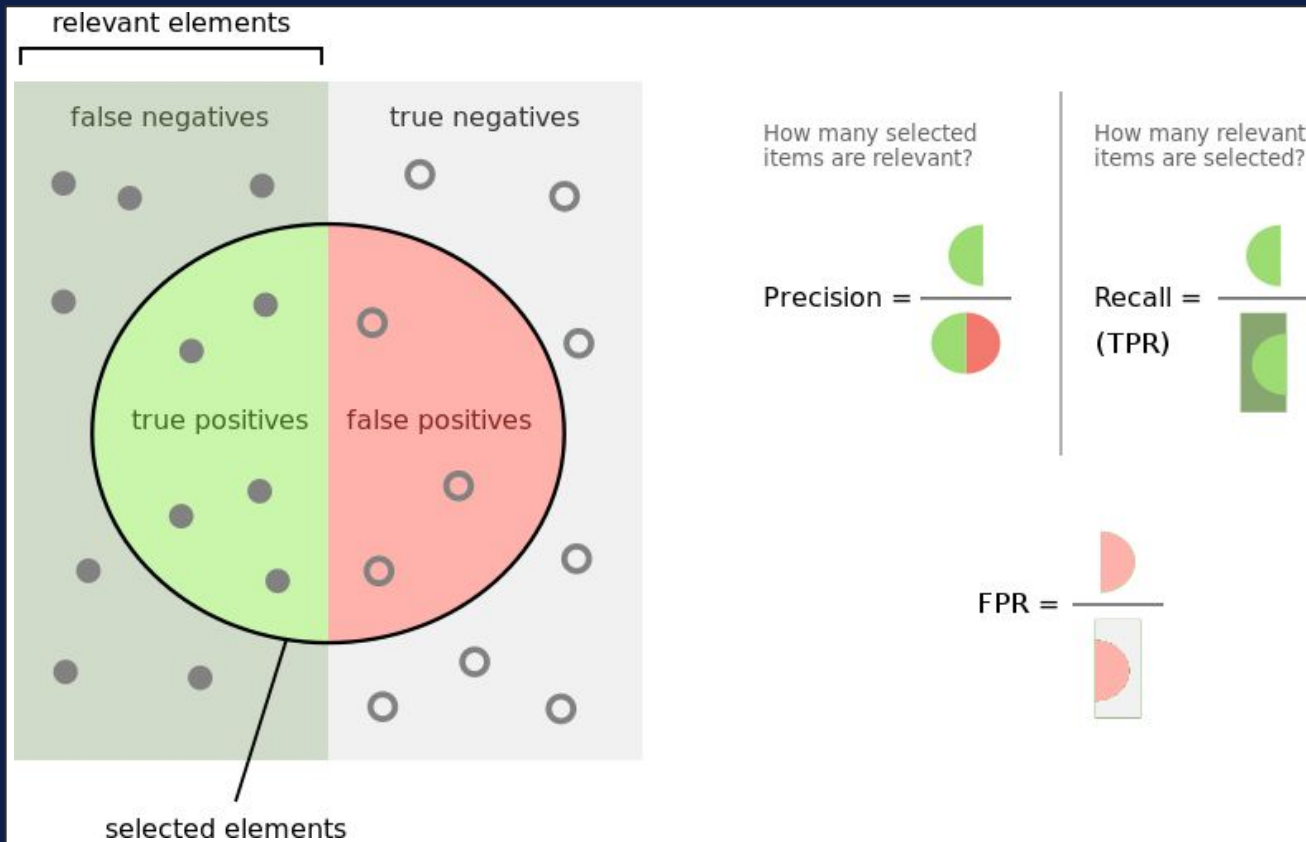
Buscamos una buena representación  
de comportamiento

$$\text{MAE: } \frac{\sum_{i=1}^n |y_i - f(x_i)|}{n} \quad \text{RMSLE: } \sqrt{\frac{\sum_{i=1}^n (\log(1 + y_i) - \log(1 + f(x_i)))^2}{n}}$$

$$\text{MSE: } \frac{\sum_{i=1}^n (y_i - f(x_i))^2}{n} \quad \text{RMSE: } \sqrt{\frac{\sum_{i=1}^n (y_i - f(x_i))^2}{n}}$$

¿Cuándo se usa cada una?

# Métricas de Clasificación



# Métricas de Clasificación

## Tipos de errores (FP, FN)

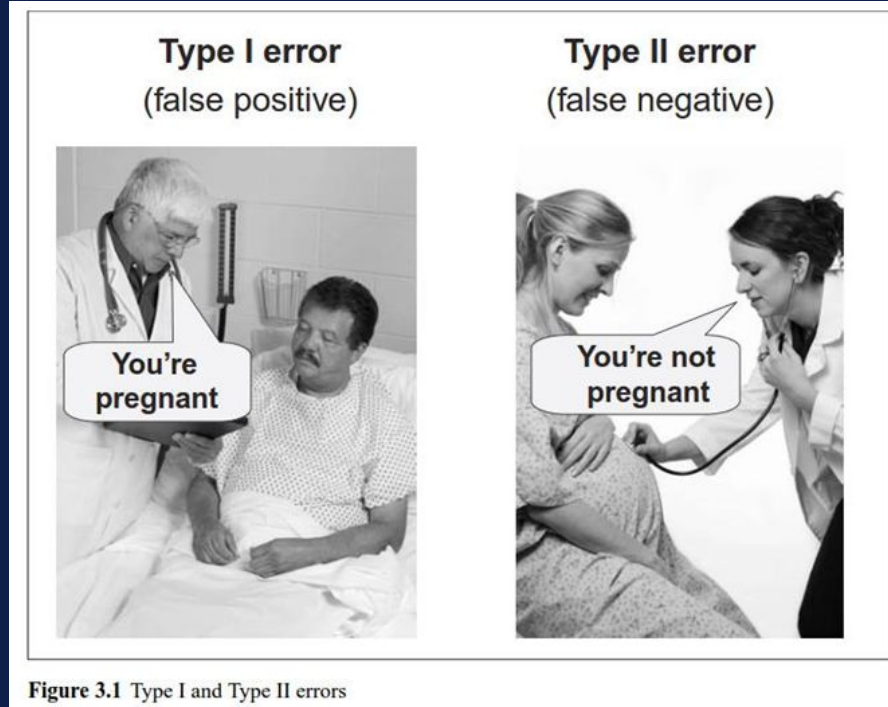
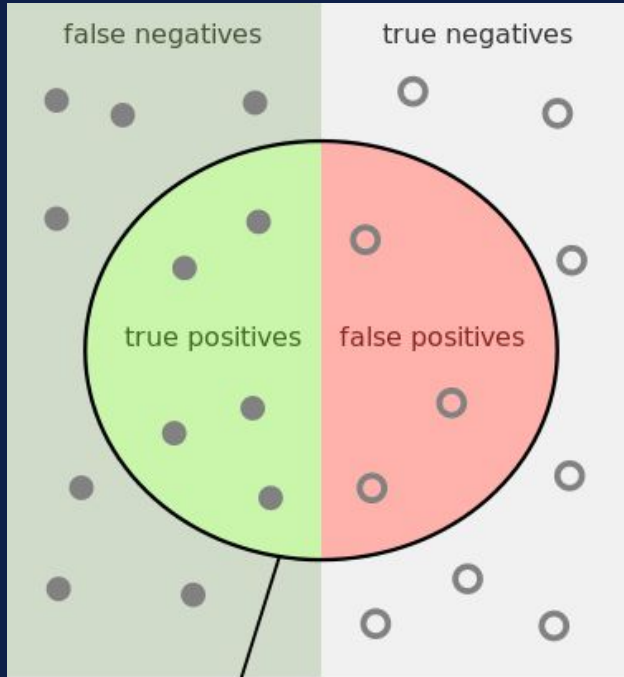


Figure 3.1 Type I and Type II errors

# Métricas de evaluación



$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN}$$

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

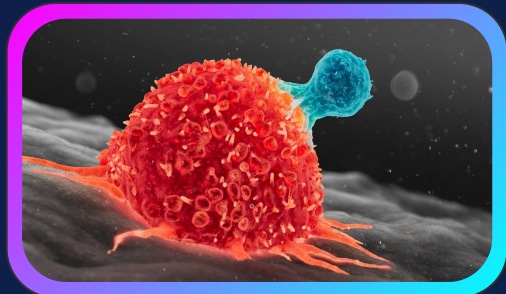
$$Recall = \frac{TP}{TP + FN}$$

$$F1 - score = \frac{2 * Precision * Recall}{Precision + Recall}$$



# Accuracy Paradox

Modelo de  
clasificación de cáncer



Accuracy: 0.96

	Cáncer	No cáncer
Clasific. cáncer	5	0
Clasific. no cáncer	45	950



# Usemos otras métricas...

	Cáncer	No cáncer
Clasific. cáncer	5	0
Clasific. no cáncer	45	950

TP = 5

TN = 950

FP = 0

FN = 45

Accuracy = 0.96

Precision = 1

Recall = 0.1

F1-score = 0.182

¿Cómo interpretamos este resultado?

# Matriz de confusión

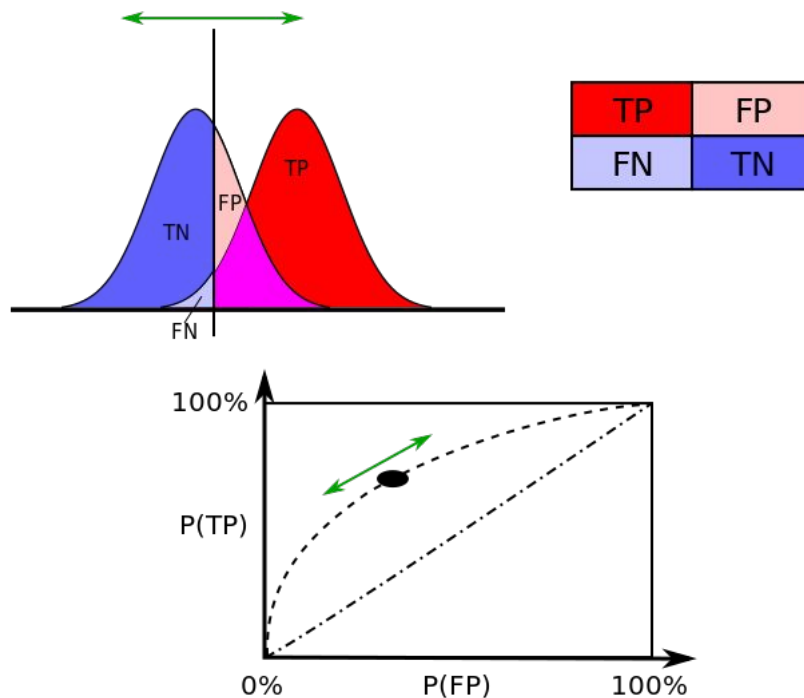
True Class	airplane	923	4	21	8	4	1	5	5	23	6
	automobile	5	972	2					1	5	15
	bird	26	2	892	30	13	8	17	5	4	3
	cat	12	4	32	826	24	48	30	12	5	7
	deer	5	1	28	24	898	13	14	14	2	1
	dog	7	2	28	111	18	801	13	17		3
	frog	5		16	27	3	4	943	1	1	
	horse	9	1	14	13	22	17	3	915	2	4
	ship	37	10	4	4		1	2	1	931	10
	truck	20	39	3	3			2	1	9	923
		airplane	automobile	bird	cat	deer	dog	frog	horse	ship	truck
		Predicted Class									

Nos permite visualizar un clasificador multiclase

¿Qué esperamos de un buen clasificador?

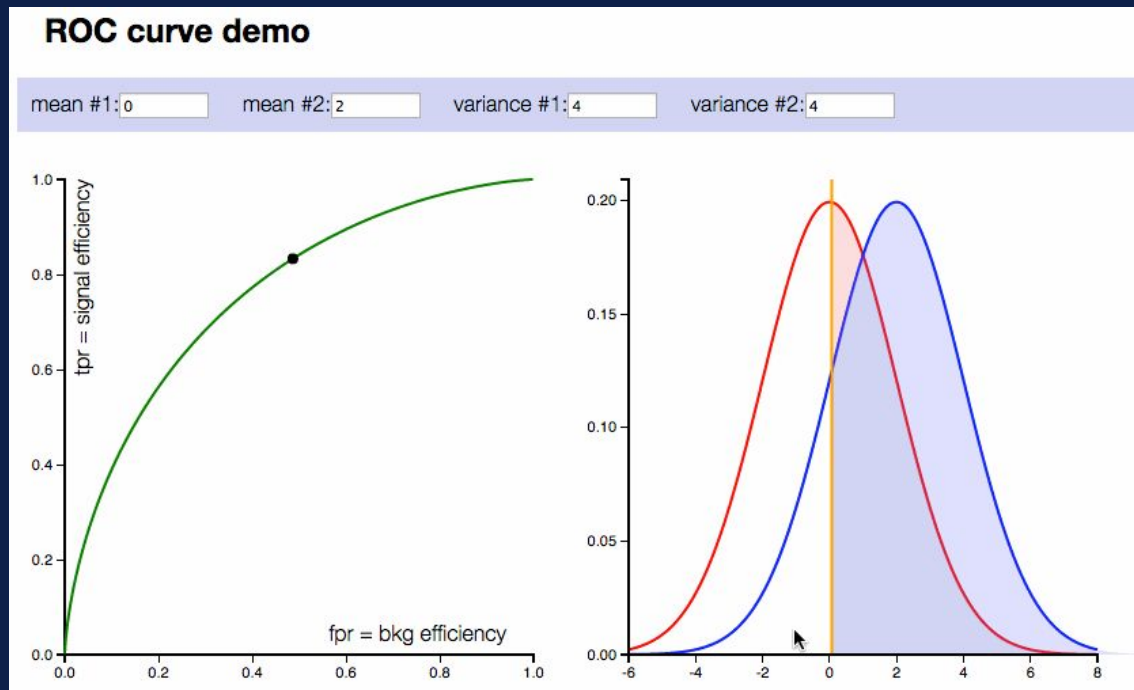
# Métricas de evaluación

## ROC Curve



# Métricas de evaluación

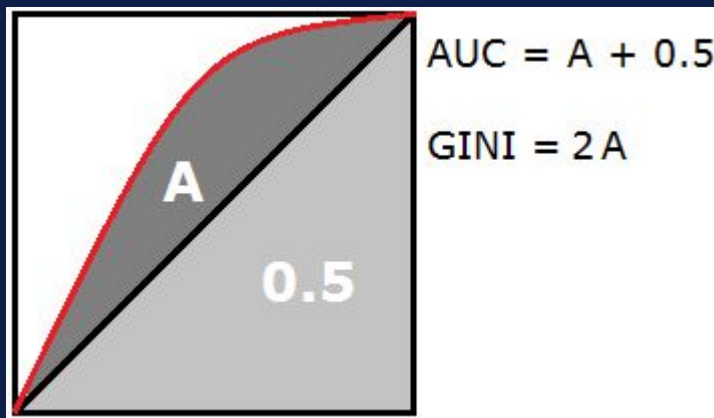
## ROC Curve



<https://arogozhnikov.github.io/2015/10/05/roc-curve.html>

# Métricas de evaluación

AUC (Área bajo la curva) y GINI



Comúnmente usada  
para comparar  
modelos!



Ay que lataaa!