

De la estadística tradicional al *Machine Learning*

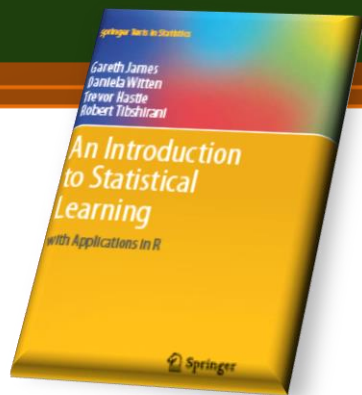


Daniel Vélez Serrano
29 mayo 2019

Un *trade-off* habitual en la modelización analítica

$$Y = f(x)$$

if?



Interpretar vs Predecir

Interpretar vs Predecir

INTERPRETAR



DEDUCIR



ENTENDER



APRENDIZAJE MANUAL

PREDECIR



ACERTAR



!!!GANAR!!!



APRENDIZAJE AUTOMÁTICO

if?

Interpretar vs Predecir

INTERPRETAR

interpretar

Conjugar

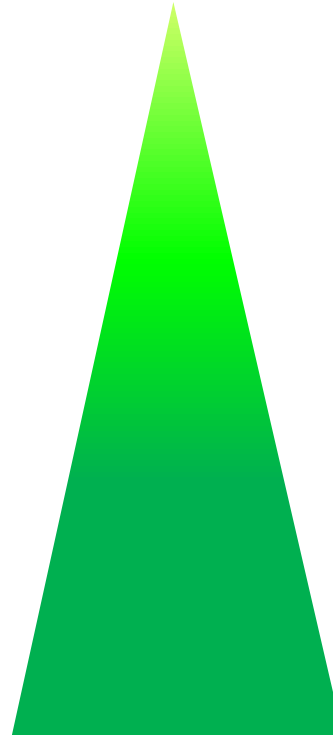
Del lat. *interpretāri*.

1. Explicar y declarar el sentido de algo



PREDECIR

if?

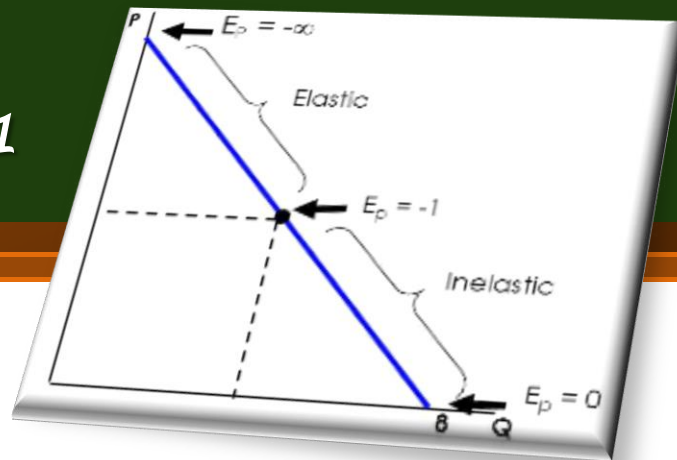


Modelización tradicional: Regresión lineal simple

Ejemplo 1: PRICING

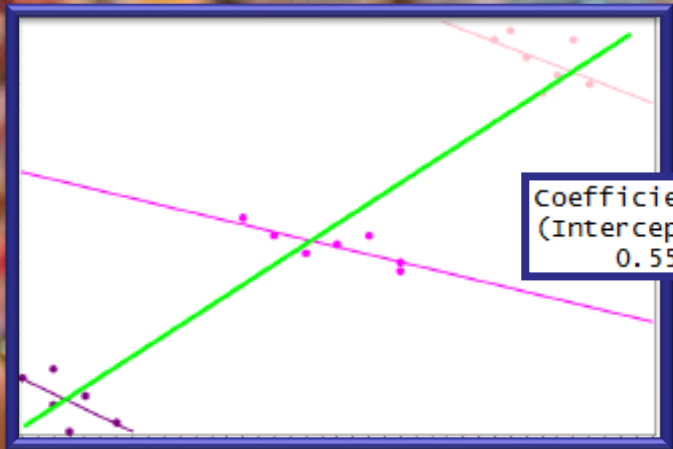
$$\text{Log}(Q) = \beta_0 + \beta_1 \text{log}(P)$$

$$\text{elasticidad} = \frac{\% \Delta Q}{\% \Delta P} \sim \beta_1$$





Sentido lógico. Buscando interpretabilidad: Datos de panel



Coefficients:
 (Intercept) 0.5523
 logPrecio 1.5800

```
(Intercept) 21.8771364
logPrecio -2.5518393
as.factor(referencia)24 0.3823752
as.factor(referencia)32 -0.2284309
as.factor(referencia)33 -2.9062072
as.factor(referencia)43 -1.9563526
as.factor(referencia)46 -3.0952390
as.factor(referencia)48 -3.9440195
as.factor(referencia)54 -0.9276794
as.factor(referencia)62 -0.8460243
```

Familia β_1

1	-4.2661328
2	-2.1283772
3	-1.4086289
4	-1.3024020
5	0.7617251
6	1.5799841
7	1.0752644
8	-2.8272297

Familia β_1

1	-10.6285540
2	-2.2364866
3	-2.2379316
4	-1.7093559
5	-1.9607779
6	-2.5518393
7	-2.2429309
8	-2.8272297

-2

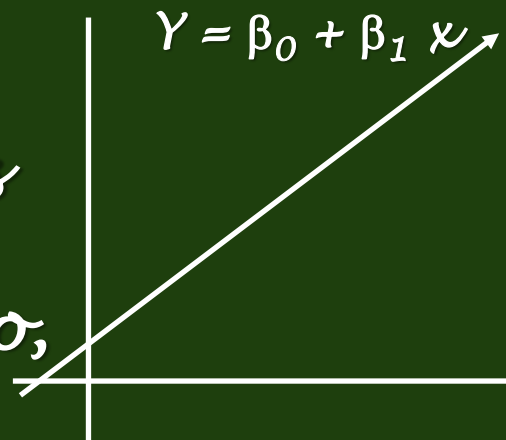
Ejemplo 2

$$Y = f(x) = \beta_0 + \beta_1 x + \varepsilon$$

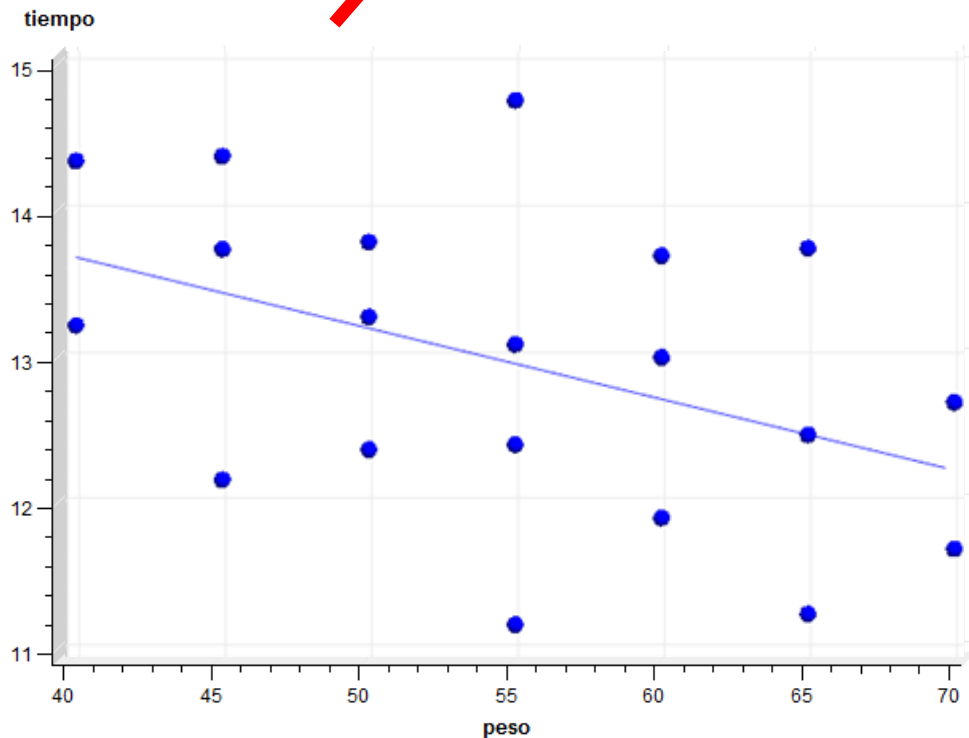
Y = tiempo en correr 100 metros

x = peso

β_1 = por cada kg de más de peso,
¿cuántos segundos más
se tarda en correrlos?

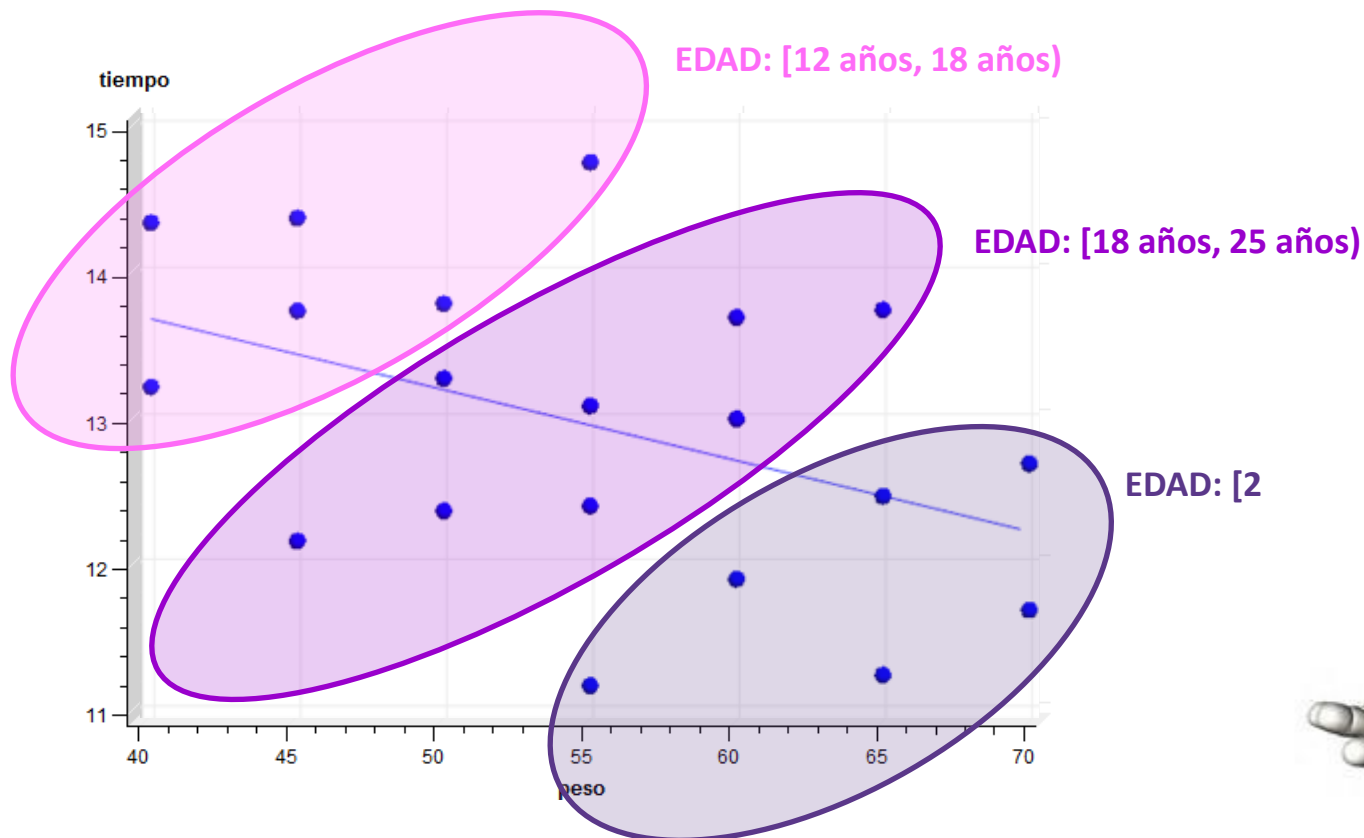


Explicar Sentido lógico. Fenómeno de confusión



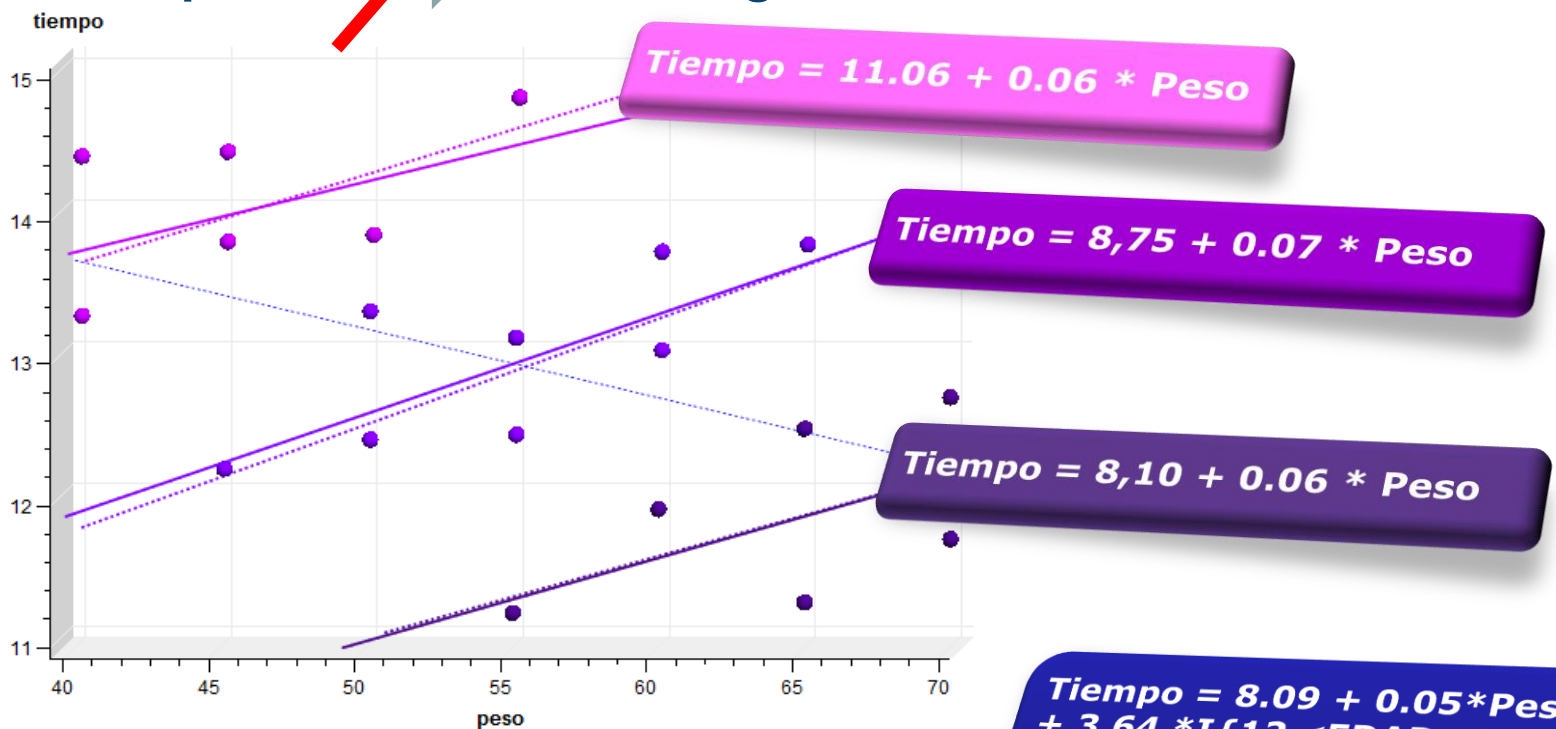
$$\text{Tiempo} = 15.61 - 0.05 \times \text{Peso}$$

Explicar Sentido lógico. Fenómeno de confusión



Incluir una variable *dummy* asociada a cada persona no parece buena opción ¹⁰

Explicar Sentido lógico. Fenómeno de confusión



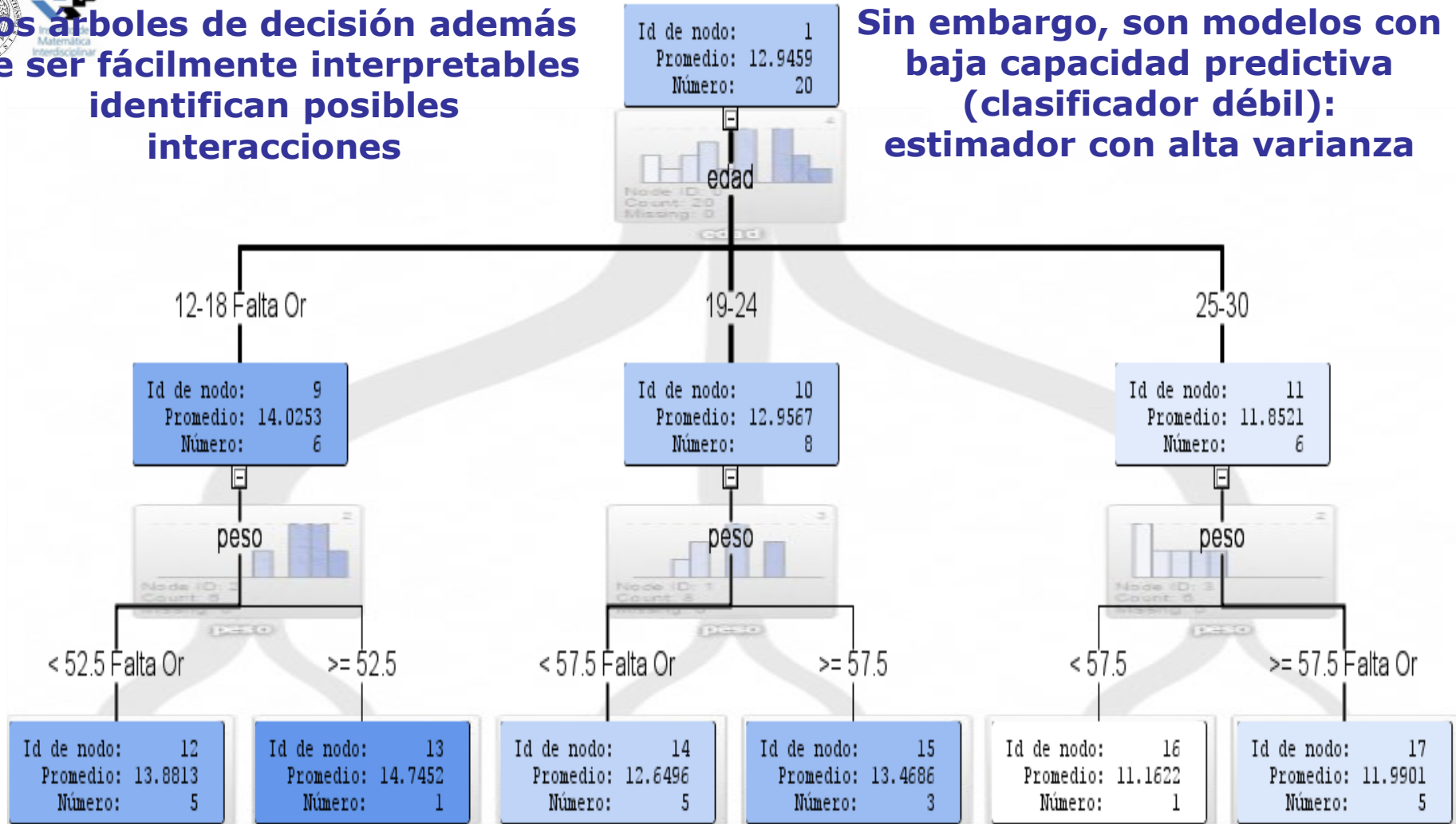
**Plantear la interacción adecuada
permite explicar bien
¿Y si tengo cientos de variables?
¿Qué interacciones contemplo?**

$$\begin{aligned} &Tiempo = 8.09 + 0.05 * Peso \\ &+ 3.64 * I\{12 \leq EDAD < 18\} \\ &+ 0.99 * I\{18 \leq EDAD < 25\} \\ &- 0.008 * Peso * I\{12 \leq EDAD < 18\} \\ &+ 0.011 * Peso * I\{18 \leq EDAD < 25\} \end{aligned}$$



Los árboles de decisión además de ser fácilmente interpretables identifican posibles interacciones

Sin embargo, son modelos con baja capacidad predictiva (clasificador débil): estimador con alta varianza



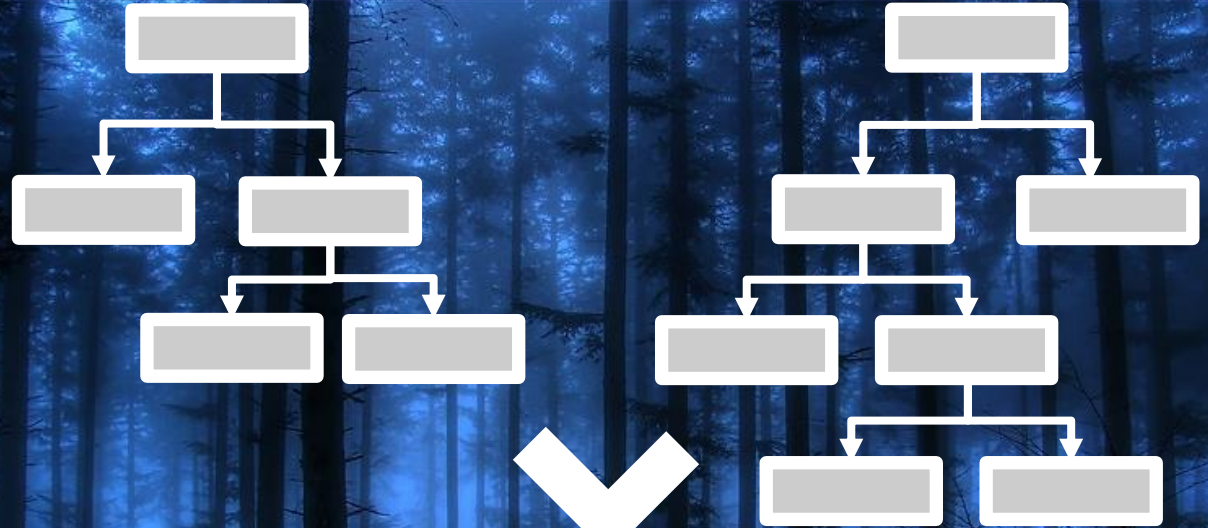
High Variance



Leo Breiman
Random Forest

“Si se dispone de X_1, \dots, X_n independientes, cada una de ellas con varianza σ^2 , su media \bar{X} tiene varianza $\frac{\sigma^2}{n}$ ”

➔ Promediar predicciones asociadas a diferentes muestras y basadas en variables independientes



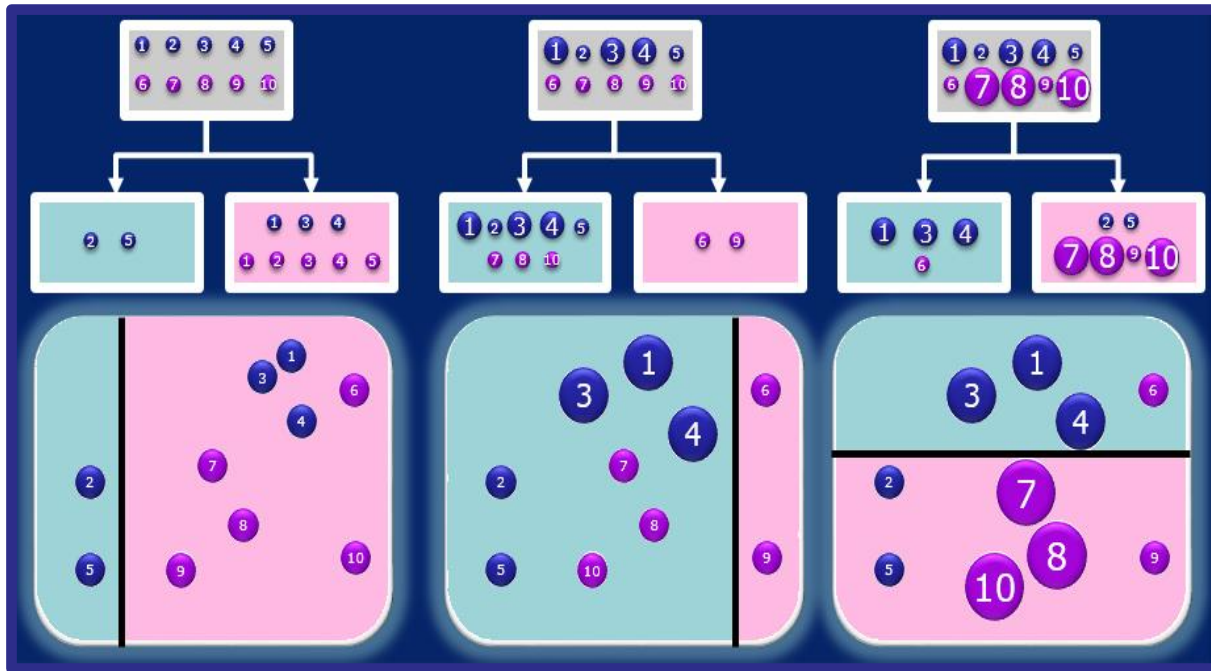
$$f(x) = \sum \alpha_i f_i(x)$$

¿Puede un conjunto de clasificadores débiles crear un clasificador robusto?



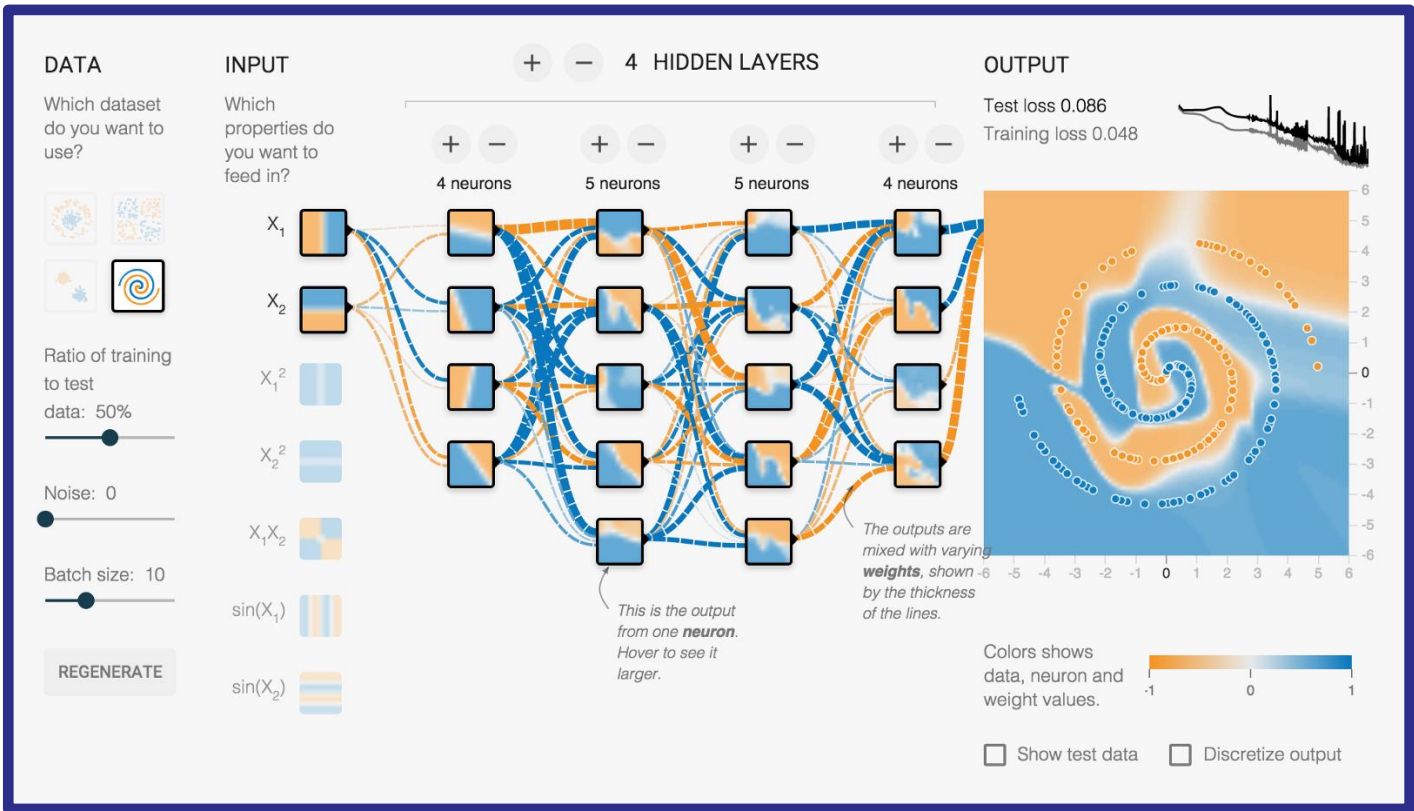
Robert Schapire Yoav Freund

Adaboost
Premio Gödel 2003



Los modelos tipo “*boosting*” se han convertido en un referente dentro de la modelización automática (*Machine Learning*) por su alto nivel de predictibilidad.

Las redes neuronales es otra de las estrategias de modelización más competitivas dentro del *Machine Learning*



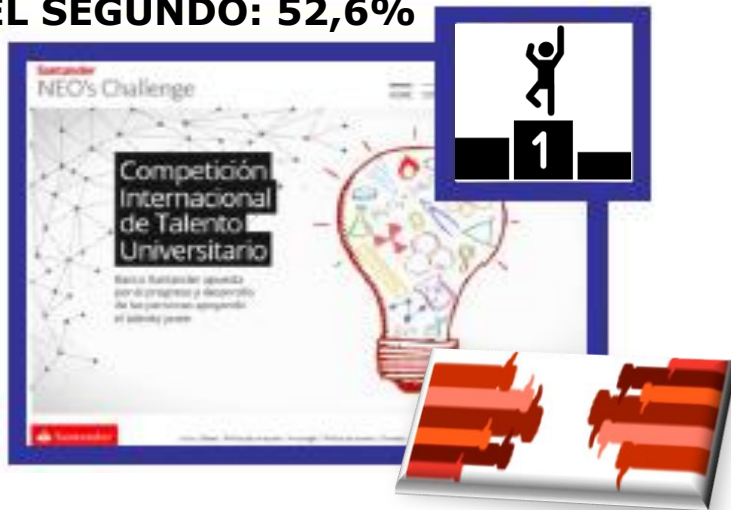
Hasta qué punto sacrificar la capacidad de interpretación de un modelo por su capacidad predictiva?



Casos de éxito de estrategias tipo “boosting”

Net Promoter Score

ESPAÑA (UCM): 56,5%
EL SEGUNDO: 52,6%



Predecir el grado de satisfacción de un cliente

Regresión: 52% → Boosting: 54%

Fraude energético

INNOVA-TSN: 80%
EL SEGUNDO: 60%



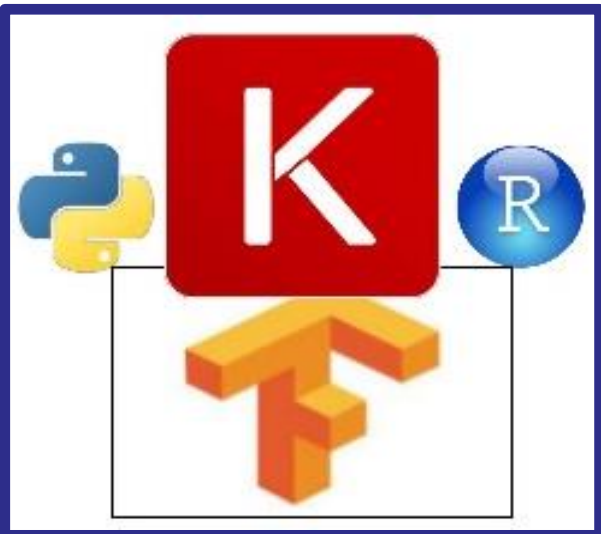
Seleccionar clientes sospechosos de manipular su contador eléctrico

Regresión: 40% → Boosting: 80%



La alta capacidad predictiva de estos modelos ha justificado el desarrollo de softwares que implementan eficientemente este tipo de técnicas:

H₂O es un referente claro dentro del ámbito del Aprendizaje Automático
KERAS es un referente claro dentro de un ámbito más específico: *Deep Learning*

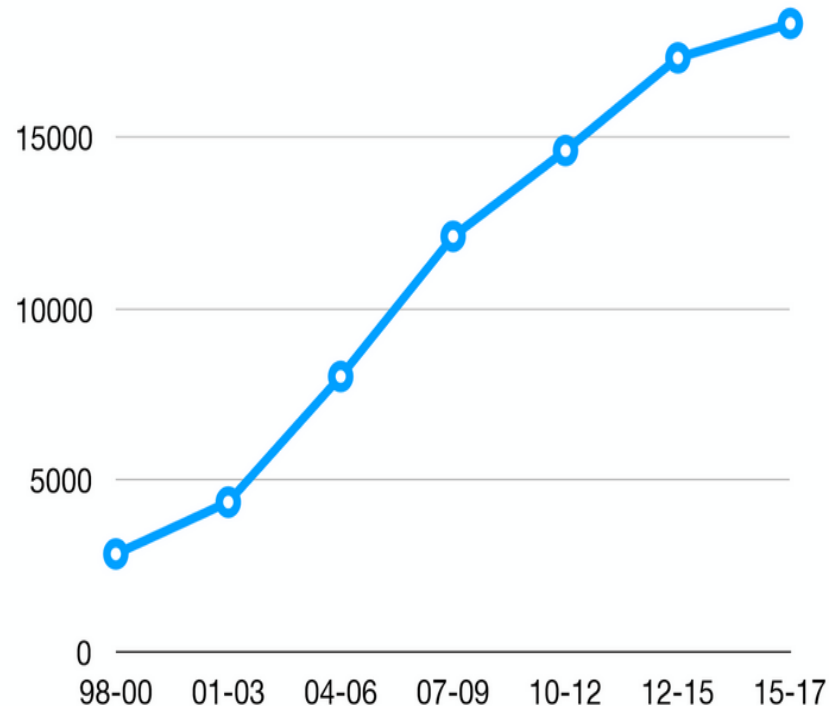


<https://www.youtube.com/watch?v=wcyMBRRLmq5>

En algunos ámbitos como finanzas o la salud, es preciso auditar el proceso de decisión y asegurar de que no sea discriminatorio ni viole la ley

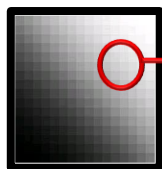


Number of papers on topic of interpretable machine learning



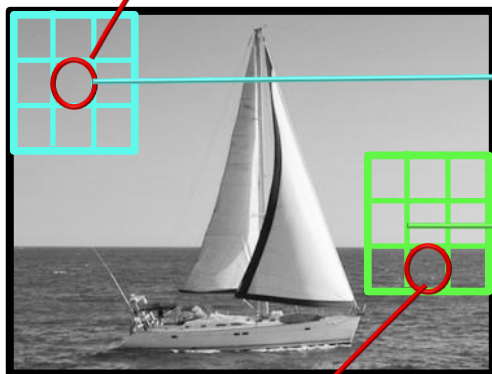
Por ello, en los últimos años han proliferado sensiblemente los algoritmos de *Machine Learning* que buscan la interpretabilidad

Algunos modelos, como las redes de convolución, ayudan a entender mediante la extracción de características, qué está teniendo en cuenta un modelo cuando identifica una imagen



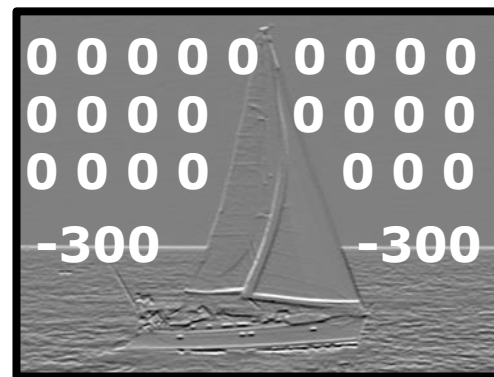
100

Máscara de Convolución

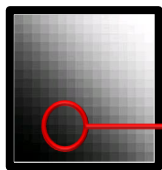


$$\begin{pmatrix} 1 & 1 & 1 \\ 0 & 0 & 0 \\ -1 & -1 & -1 \end{pmatrix} \longrightarrow 0$$

$$100 \times 3 - 200 \times 3 \longrightarrow -300$$

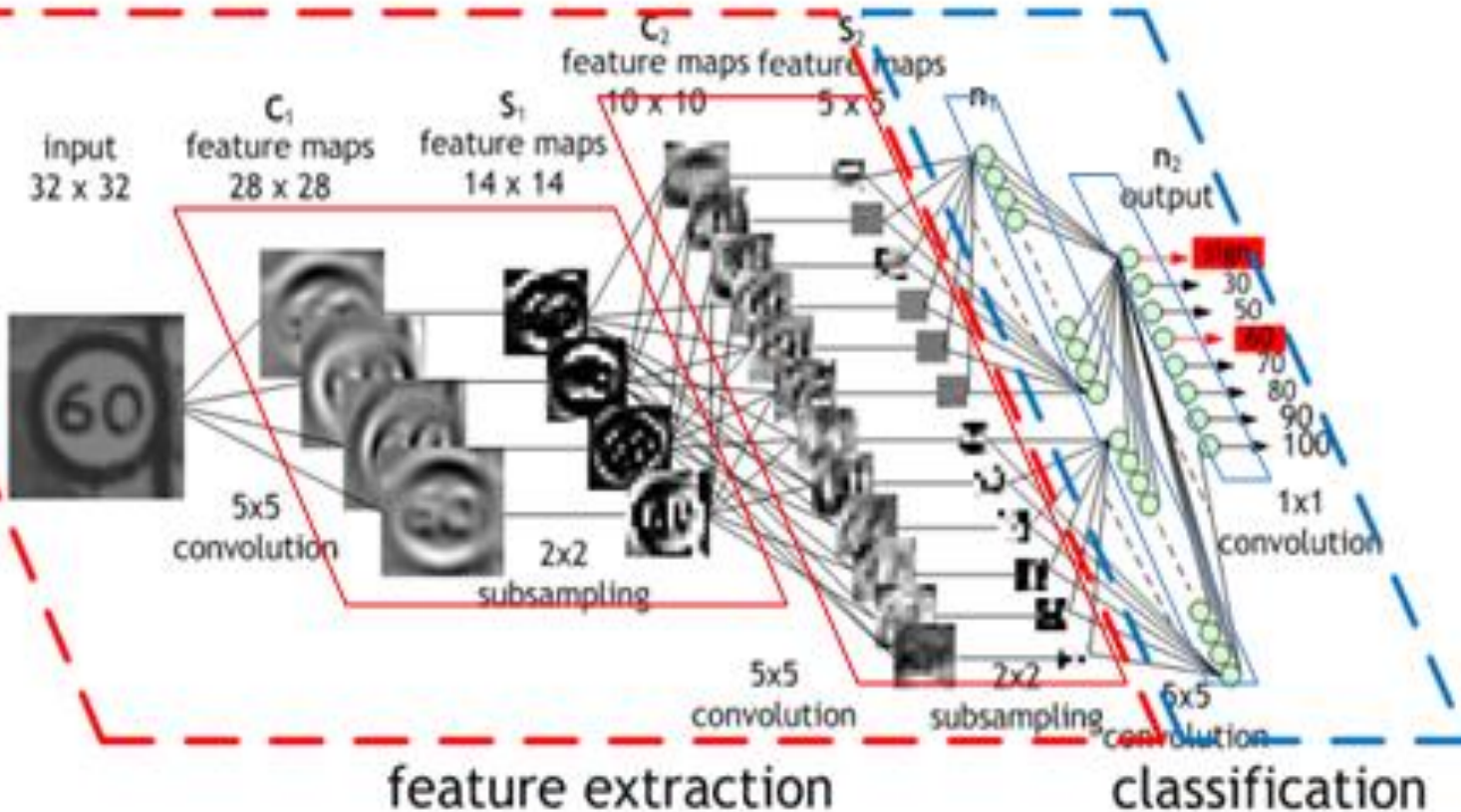


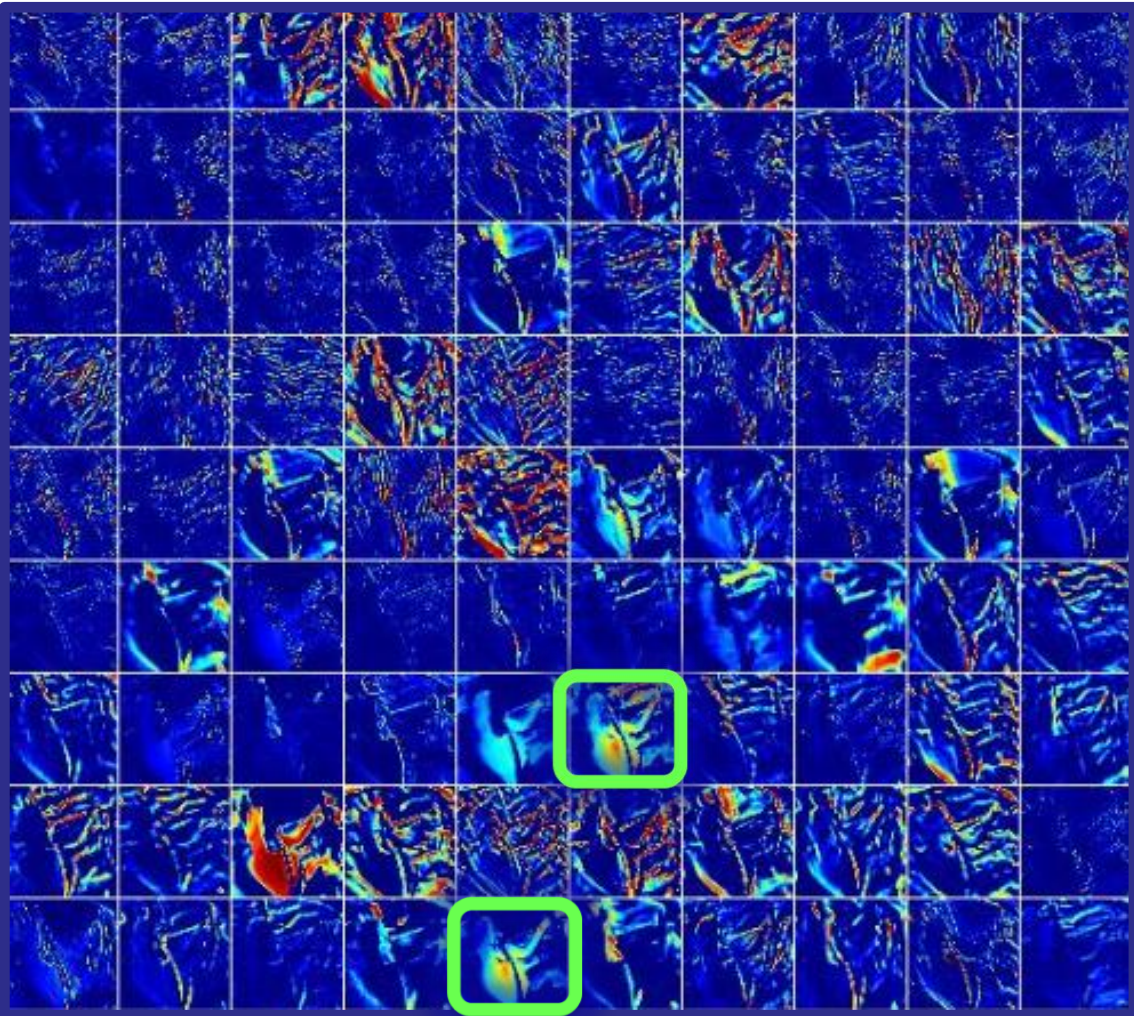
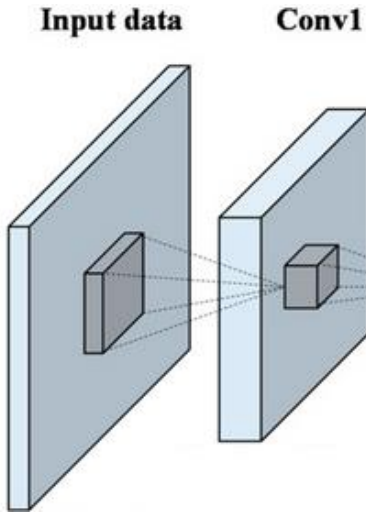
Característica extraída:
líneas horizontales



200

Las capas de convolución permiten la extracción de características como paso previo al ajuste del modelo que lleva a cabo la predicción/clasificación





**Ejemplo:
Identificación
mediante la imagen
captada por una
cámara de personas
hablando por el
móvil**





Sin embargo, lo ideal es la **búsqueda de metodologías** que permitan interpretar los resultados de un modelo con independencia de la naturaleza de éste

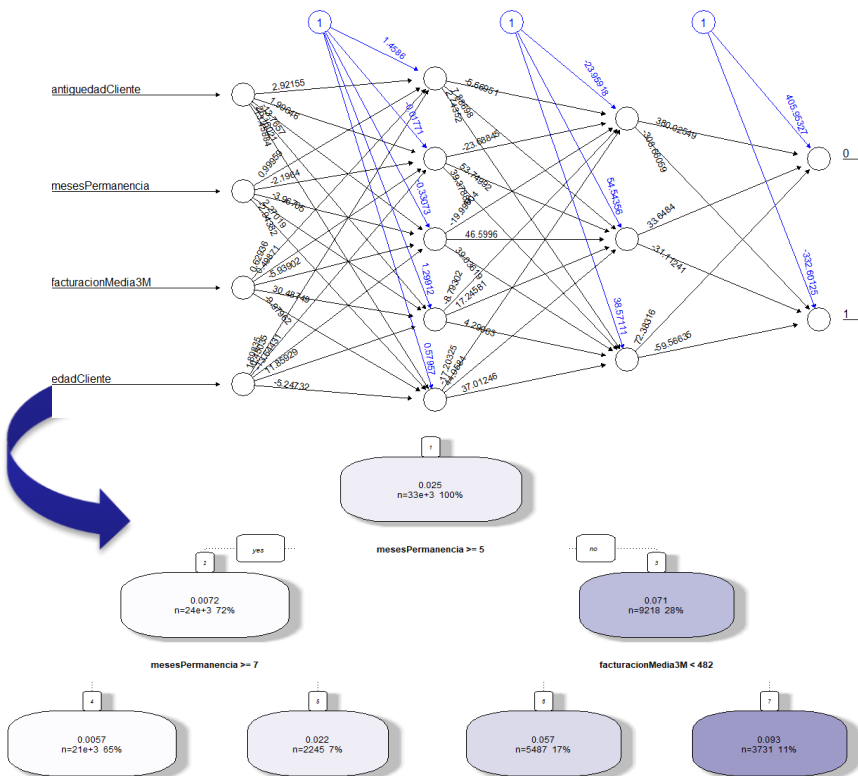
Metodologías globales: ¿Cuáles son las variables que más importancia tienen en la predicción? ¿Cómo influyen?

- Árbol de decisión
- Gráficos de importancia de variables

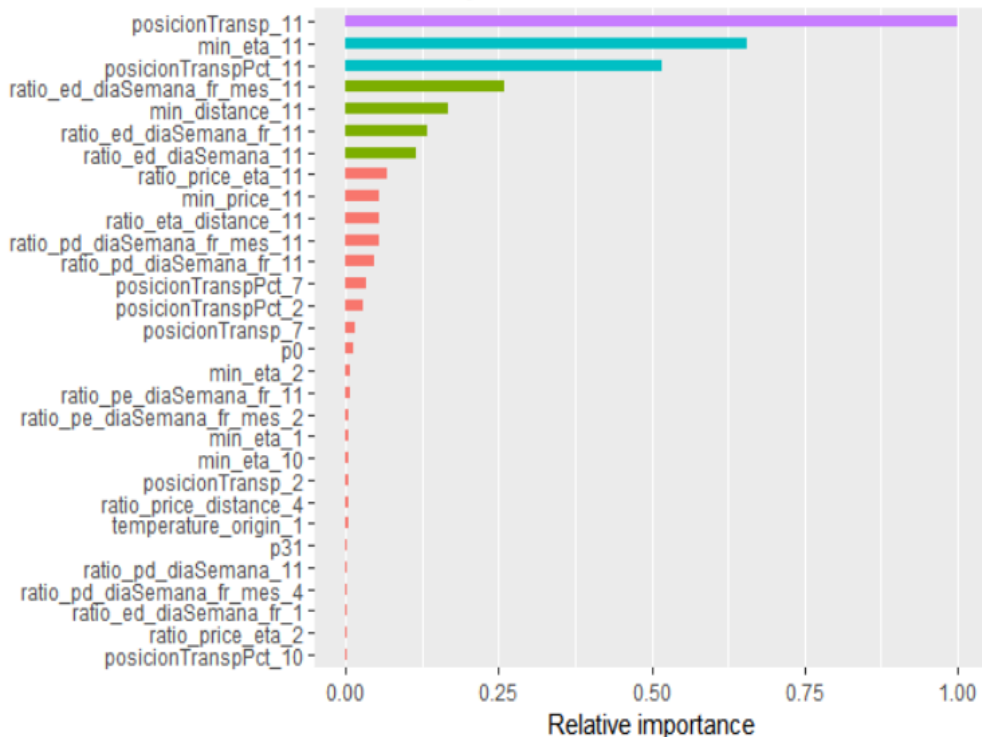
Metodologías locales: ¿Cuáles son las variables que más condicionan la predicción dada a un individuo? (ejemplo: ¿de qué sirve saber que un cliente va a dejar mi compañía si no conozco la razón (variable) que lo justifica?

- LIME: *Local Interpretable Model-agnostic Explanations*
- DALEX: *Descriptive mACHine Learning Explanations*
- IML: *Interpretable Machine Learning*

Metodologías globales



Feature importance: Gain



En un problema de clasificación binaria, un árbol de regresión permitiría explicar la probabilidad predicha por un modelo

Un gráfico de importancia de variables permite ver en cuántos de los modelos ensamblados entra cada una de las variables o su contribución

LIME: Una metodología local

Se trata de una metodología que busca justificar, registro a registro, cuál es la variable que tiene mayor contribución al valor predicho para él

Para ello, se procede de la siguiente manera:

1. Se selecciona un registro

2. Se generan registros nuevos en un entorno suyo mediante variaciones aleatorias del valor de las explicativas asociadas a dicho registro, dando más peso a las más similares

3. Se calcula la predicción asociada a los registros

4. Se ajusta un modelo sencillo e interpretable para separar los registros predichos en una y otra clase



Registros simulados:

- 1 "Me gusta película" ➔ predicción 1
- 2 "Me esta película" ➔ predicción 2
- 3 "Me película" ➔ predicción 3
- 4 "Me gusta" ➔ predicción 4

Ajustar un modelo interpretable (árbol) para explicar los registros en función de las palabras

Conclusiones

En algunas problemáticas, las predicciones proporcionadas por un modelo “*black box*” pueden ser sensiblemente mejores que las de uno explicable
Sin embargo, las decisiones tomadas por dicho modelo no siempre se entienden bien



©marketoonist.com

Es importante poder interpretar dichas predicciones para poder tomar decisiones justificadas

*¡Muchas
Gracias!*