

APRENDIZAJE AUTOMÁTICO

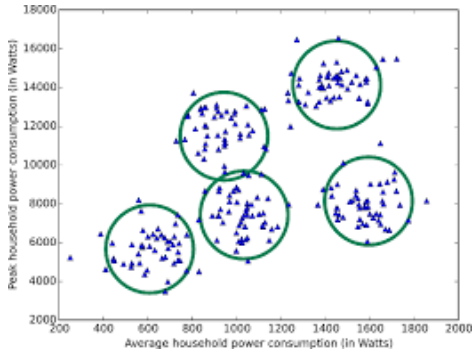
Javier Diaz Cely, PhD



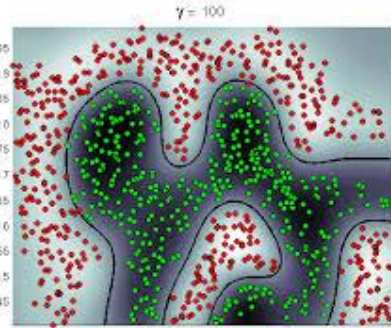
AGENDA



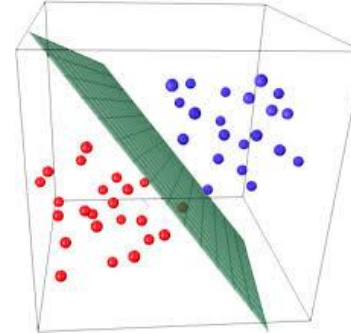
**Aprendizaje
automático**



**Aprendizaje
no supervisado**



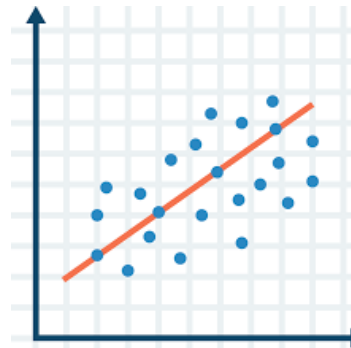
**Aprendizaje
supervisado**



Clasificación



**Métricas de
Evaluación de la
clasificación**

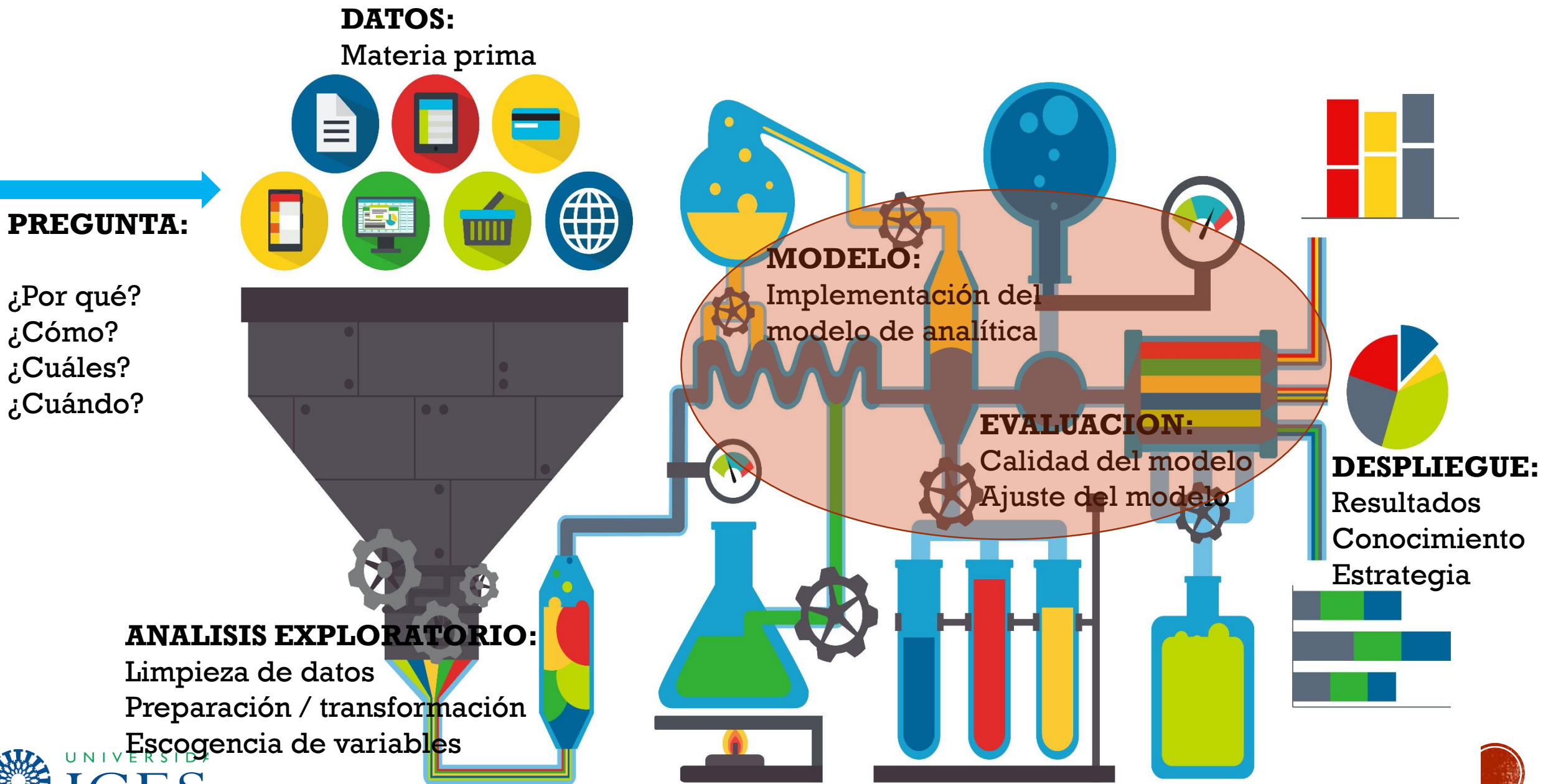


Regresión



**Métricas de
Evaluación de la
regresión**





Machine Learning



what society thinks I
do

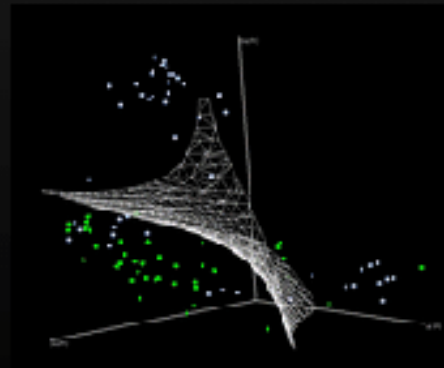


what my friends think
I do



what my parents think
I do

$$\begin{aligned}
 L_T &= \frac{1}{2} \|w\|^2 - \sum_{i=1}^n \alpha_i y_i (x_i \cdot w + b) + \sum_{i=1}^n \alpha_i \\
 \alpha_i &\geq 0, \forall i \\
 w &= \sum_{i=1}^n c_i y_i x_i, \sum_{i=1}^n c_i y_i = 0 \\
 \nabla \hat{J}(\theta_t) &= \frac{1}{n} \sum_{i=1}^n \nabla \ell(x_i, y_i; \theta_t) + \nabla r(\theta_t) \\
 \theta_{t+1} &= \theta_t - \eta_t \nabla \ell(x_{t+1}, y_{t+1}; \theta_t) - \eta_t \cdot \nabla r(\theta_t) \\
 \mathbb{E}_{x_i}[\ell(x_i, y_i; \theta_t)] &= \frac{1}{n} \sum_{i=1}^n \ell(x_i, y_i; \theta_t)
 \end{aligned}$$



```

1 library(ggplot2)
2 library(caret)
3
4 canciones <- read.table('
5 str(canciones)
6 summary(canciones)
7 head(canciones)
8

```

what other programmers
think I do

what I think I do

what I really do

APRENDIZAJE AUTOMÁTICO

- **¿Por qué es necesario?**
 - Tareas complejas extremadamente difíciles de programar
 - Poder computacional disponible para tratar grandes volúmenes de datos

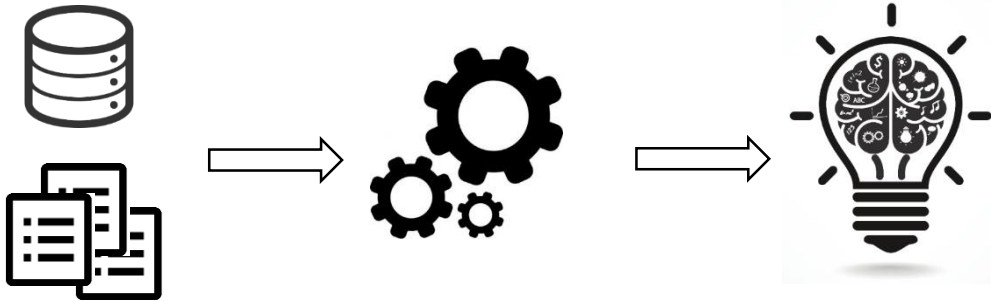
Las máquinas tienen que aprender por sí solas



APRENDIZAJE AUTOMÁTICO

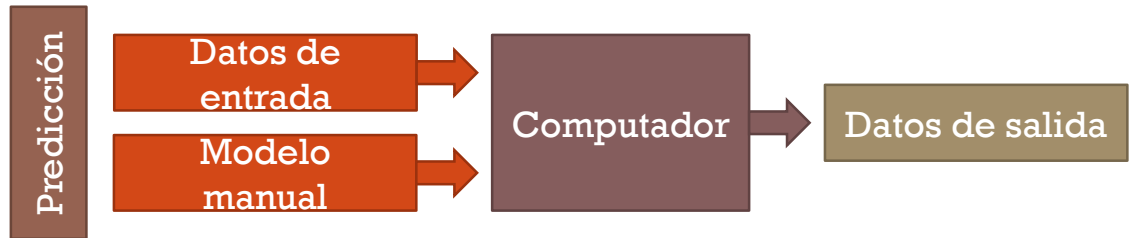
- Definición:

El aprendizaje automático es la ciencia que permite a los computadores aprender, sin ser explícitamente programados¹

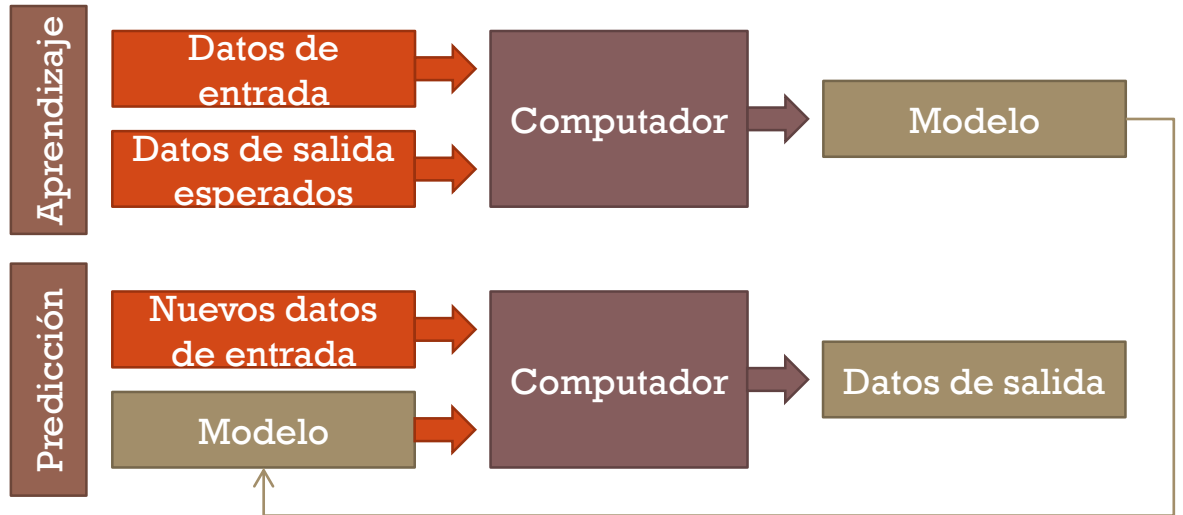


1. Andrew Ng, Stanford University, 2014

Modelo tradicional



Ciencia de datos



APRENDIZAJE AUTOMÁTICO

Aprendizaje supervisado

- Aprender a partir de un “experto”
- Datos de entrenamiento **etiquetados** con una clase o valor:

$(x_1, x_2, \dots, x_n, y)$

Predictores, explicativos,
independientes

Dependiente, objetivo,
salida

- **Meta:** predecir una clase o valor

Aprendizaje no supervisado

- Sin conocimiento de una clase o valor objetivo
- Datos **no** están **etiquetados**

(x_1, x_2, \dots, x_n)

- **Meta:** descubrir factores no observados, estructura, o una representación mas simple de los datos



APRENDIZAJE AUTOMÁTICO

Aprendizaje supervisado

Edad	Ingresos	Tiene carro?
24	1'200.000	NO
23	4'500.000	SI
45	1'250.000	SI
32	1'100.000	NO

Datos etiquetados:
"Respuestas correctas" disponibles

Factores/atributos/variables independientes, predictores, explicativos

Dependiente, objetivo, respuesta, salida

34	3'500.000
----	-----------

?

¿Cuál es el valor predicho para una instancia dada?

Aprendizaje no supervisado

Edad	Ingresos
24	1'200.000
23	4'500.000
45	1'250.000
32	1'100.000

Factores/atributos/variables

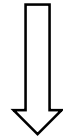
¿Se puede encontrar alguna estructura en los datos?



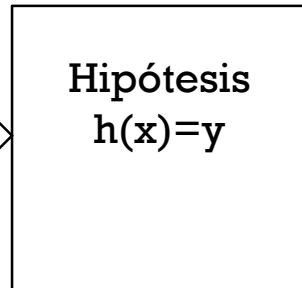
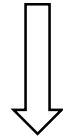
APRENDIZAJE AUTOMÁTICO

Aprendizaje supervisado

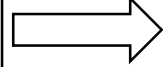
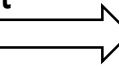
Set de entrenamiento(x_1, x_2, \dots, x_n, y)



Algoritmo de aprendizaje,
estimación de parámetros



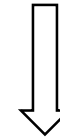
Set de text de test
(x_1', x_2', \dots, x_n')



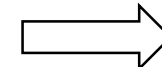
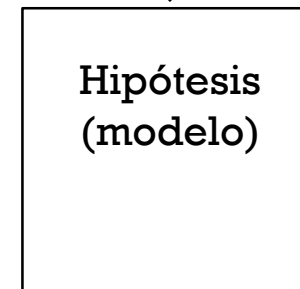
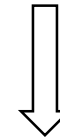
Resultado
(y')

Aprendizaje no supervisado

Set de entrenamiento(x_1, x_2, \dots, x_n)



Algoritmo de aprendizaje,
estimación de parámetros



Resultado
(estructura)



MÉTRICAS DE EVALUACIÓN

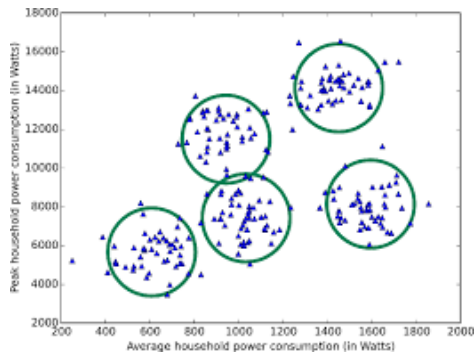
- Necesidad de evaluar la calidad de los modelos de aprendizaje automático
- Diferentes criterios a tener en cuenta:
 - Correctitud de la predicción
 - Simplicidad (parsimonia)
 - Interpretabilidad
 - Tiempo de aprendizaje o de predicción
 - Escalabilidad (importante para Big Data)



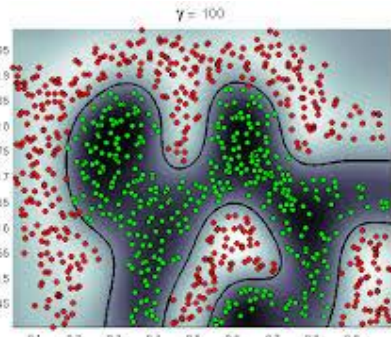
AGENDA



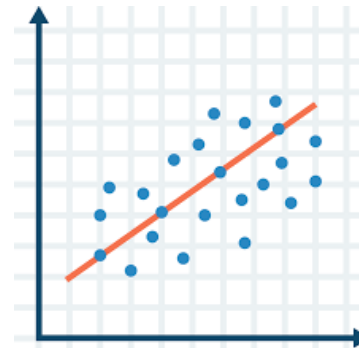
**Aprendizaje
automático**



**Aprendizaje
no supervisado**



**Aprendizaje
supervisado**



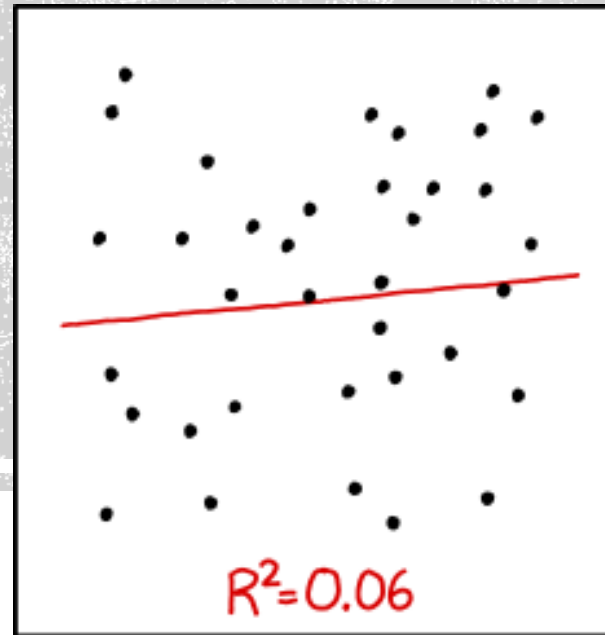
Regresión



**Métricas de
Evaluación de la
regresión**



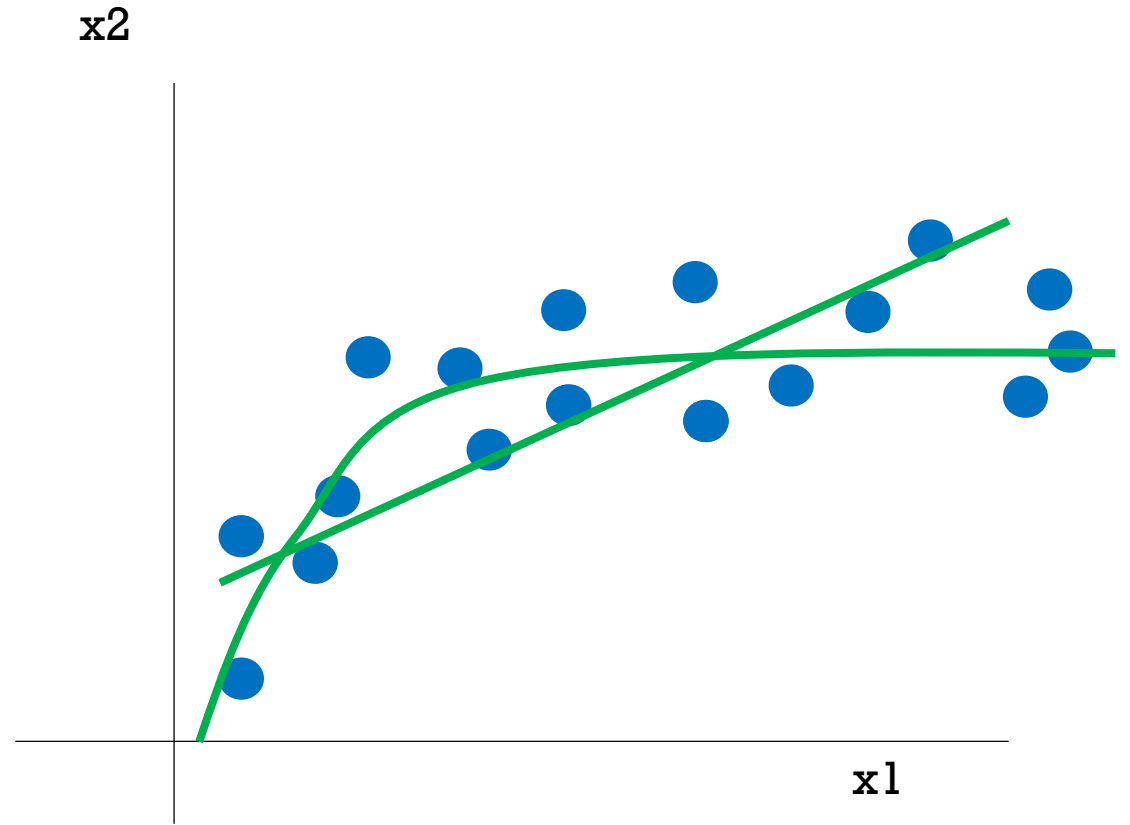
REGRESIÓN



I DON'T TRUST LINEAR REGRESSIONS WHEN IT'S HARDER TO GUESS THE DIRECTION OF THE CORRELATION FROM THE SCATTER PLOT THAN TO FIND NEW CONSTELLATIONS ON IT.

REGRESIÓN

- Encontrar modelos que permitan predecir valores continuos:
 - KNN
 - Regresión lineal
 - Regresión polinómica
 - Árboles de regresión
 - ...
- Valores **continuos** de la variable objetivo
- **Baseline**: medida de evaluación dada por un modelo que predice una medida de tendencia central (e.g. el promedio)



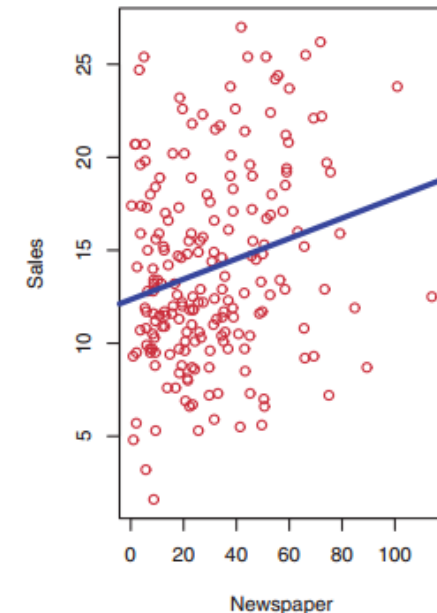
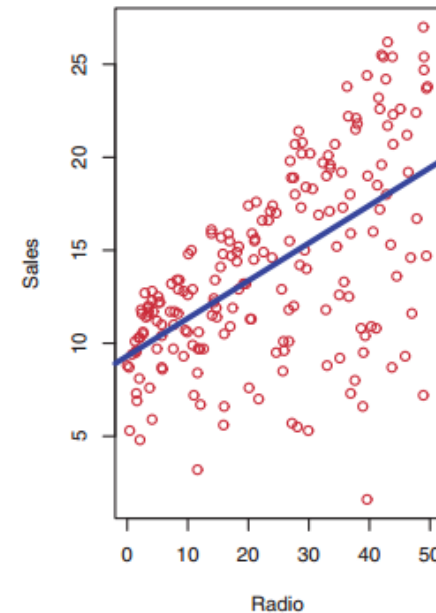
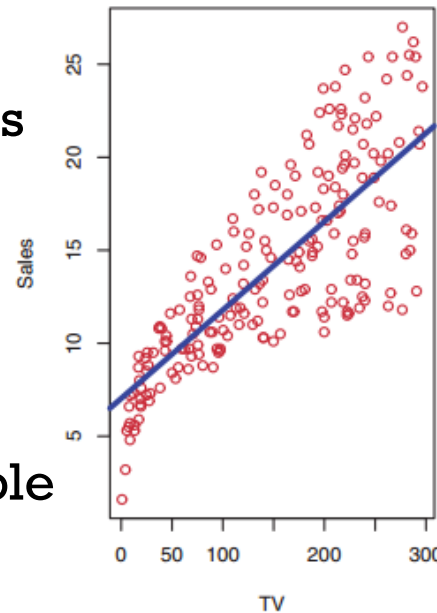
REGRESIÓN

■ Predicción:

- Procesos de caja negra
- Estimar el valor objetivo Y dado los valores de los predictores X

■ Inferencia:

- ¿Cuáles son los predictores asociados con la respuesta?
- ¿Cuál es la relación entre la variable respuesta y cada uno de los predictores?
- ¿Se puede resumir esa relación linealmente o se trata de una relación mas compleja?



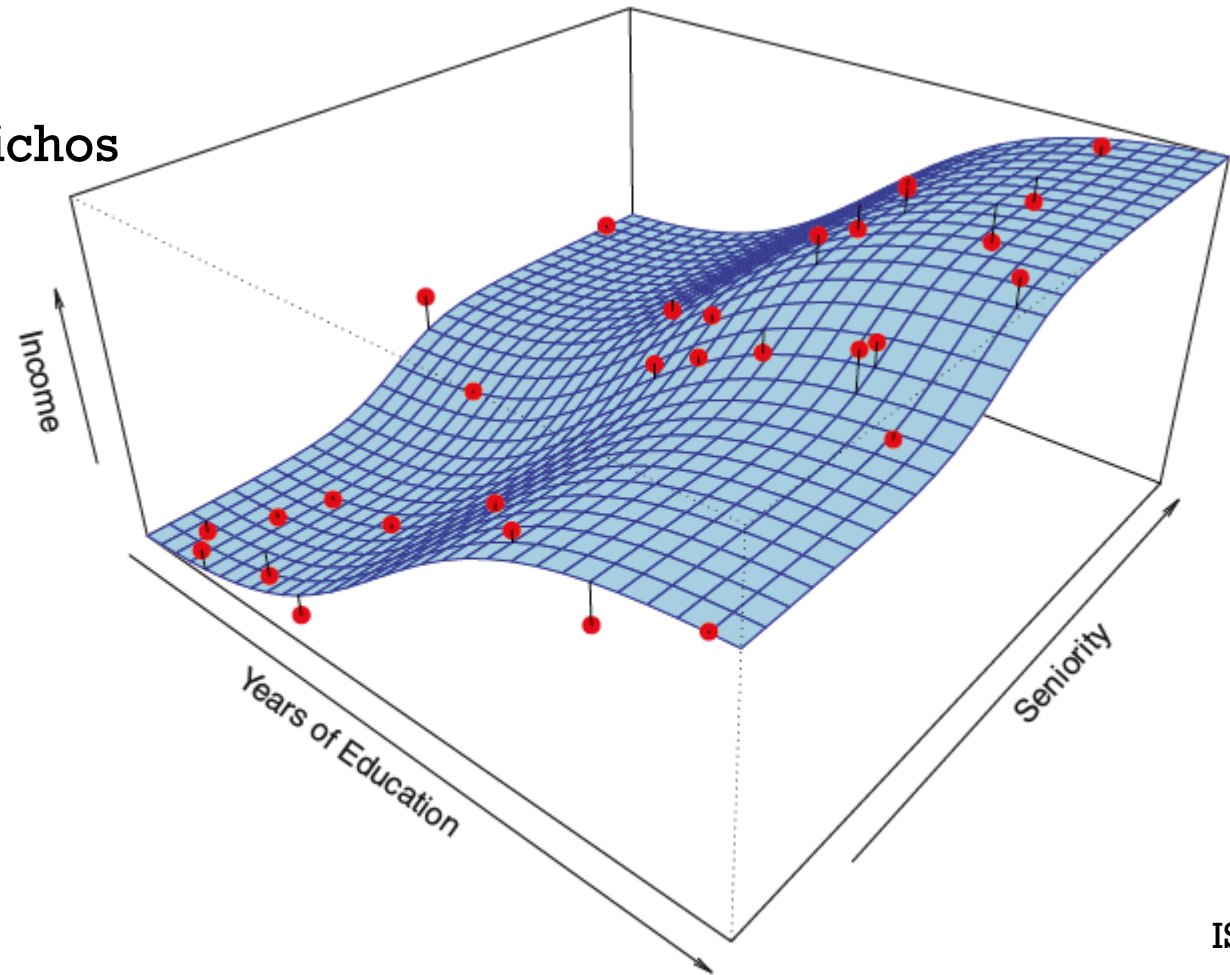
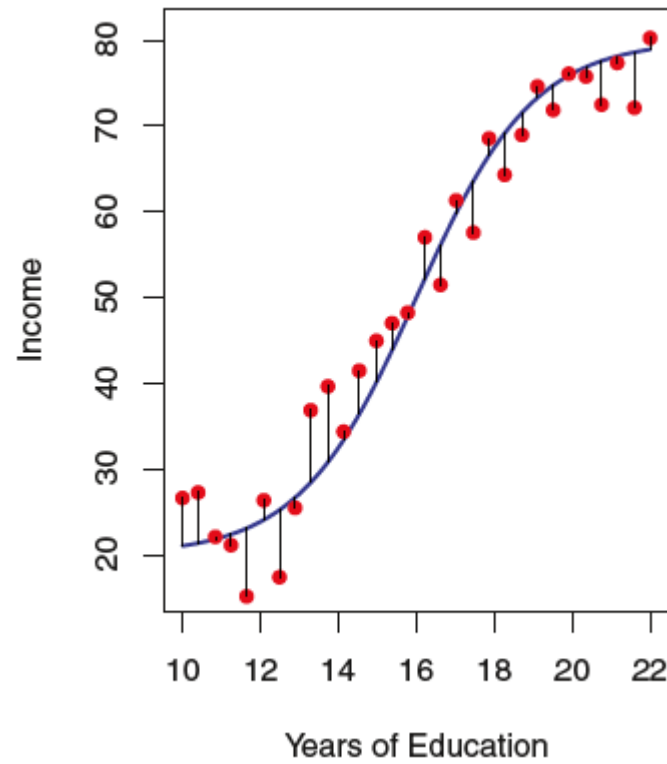
$$\text{Ventas} = f(\text{TV}, \text{Radio}, \text{Periódicos})$$

ISLR, 2013



RESIDUOS

Residuos: diferencia entre los valores reales y los valores predichos



ISLR, 2013



MÉTRICAS DE REGRESIÓN

Coeficiente de correlación (Pearson $\rho \in [-1;1]$): indica la fuerza de la relación lineal entre los predictores y la variables objetivo, que puede ser positive o negativa

- $|\rho| = 0$ no hay correlación
- $|\rho| = 0.10$ correlación muy débil
- $|\rho| = 0.25$ correlación débil
- $|\rho| = 0.50$ correlación media
- $|\rho| = 0.75$ correlación fuerte
- $|\rho| = 0.90$ correlación muy fuerte
- $|\rho| = 1$ correlación perfecta

$$\rho_{x,y} = \frac{Cov(x, y)}{\sigma_x \sigma_y}$$

Coeficiente de determinación ($R^2 = \rho^2$): indica el porcentaje de la varianza de la variable objetivo que puede ser explicada por los predictores a partir de la relación lineal



MÉTRICAS DE REGRESIÓN

- MAE (mean absolute error):

$$\frac{1}{m} \sum_{i=1}^m |h_{\theta}(x_i) - y_i|$$

- MSE (mean square error):

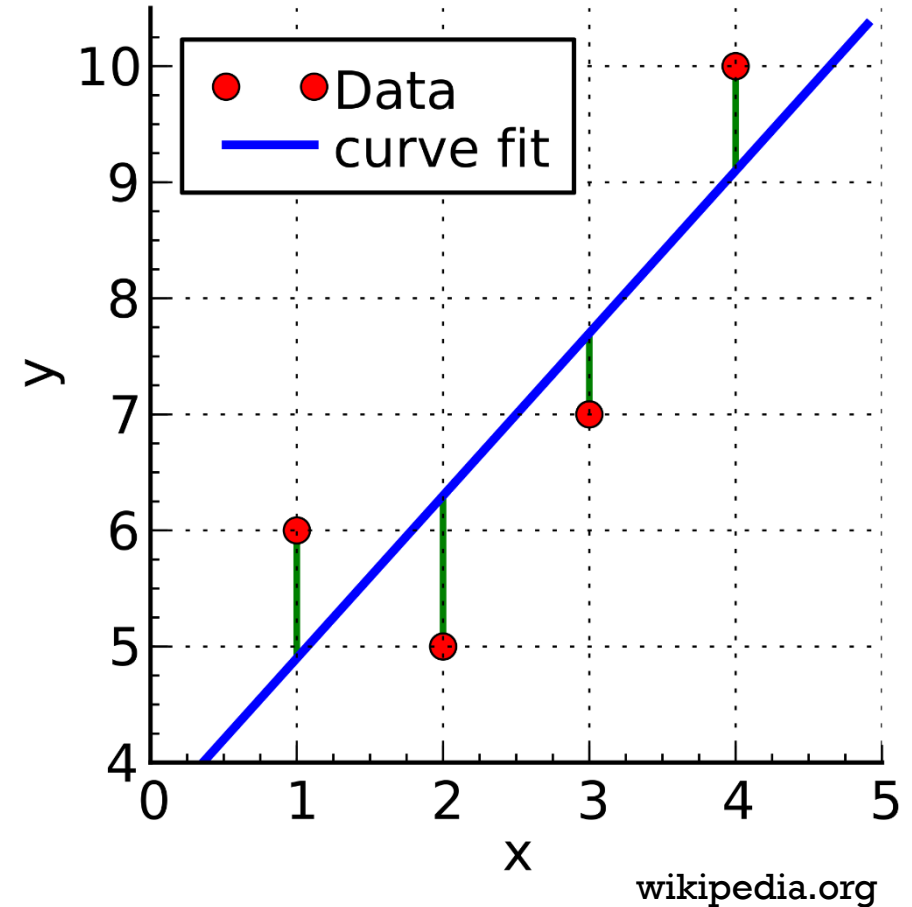
$$\frac{1}{m} \sum_{i=1}^m (h_{\theta}(x_i) - y_i)^2$$

- RMSE (root mean square error):

$$\sqrt{\frac{1}{m} \sum_{i=1}^m (h_{\theta}(x_i) - y_i)^2}$$

- R^2 (coeficiente de determinación):

$$1 - \frac{\sum_{i=1}^m (h_{\theta}(x_i) - y_i)^2}{\sum_{i=1}^m (y_i - \bar{y})^2}$$



wikipedia.org



VARIABLES CATEGÓRICAS

- Las variables predictoras deben ser numéricas.
- Las variables categóricas debes ser convertidas en numéricas:
 - One hot encoding: se crea una variable para cada valor posible de cada variable categórica
 - Contraste o “dummy”: se crea una variable para cada valor posible menos 1 de cada variable categórica.

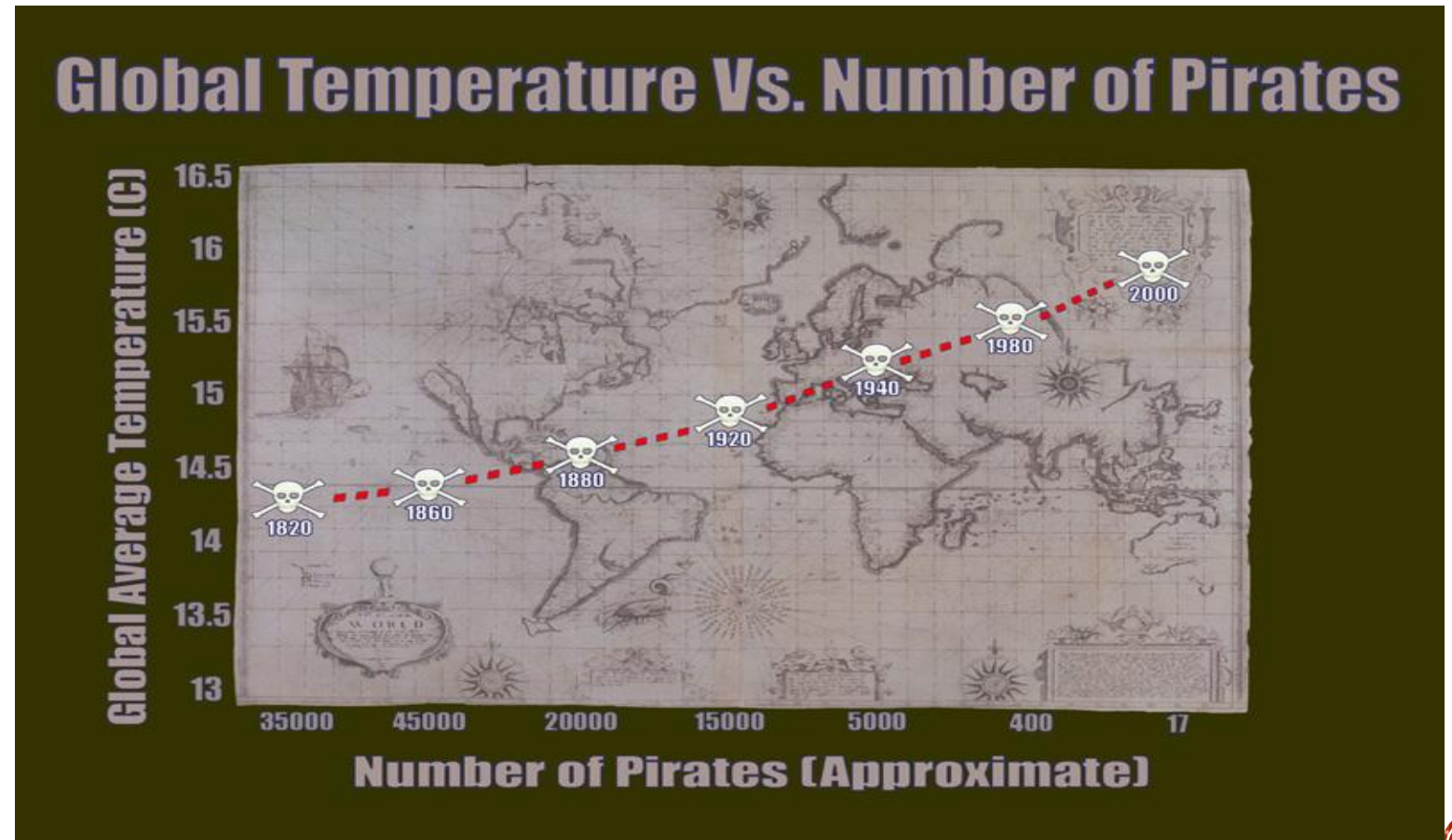
Ejemplo: variable estrato con 3 valores posibles (bajo, medio y alto)

	Estrato_bajo	Estrato_medio
Valor = bajo	1	0
Valor = medio	0	1
Valor = alto	0	0



REGRESIÓN — CUIDADO!

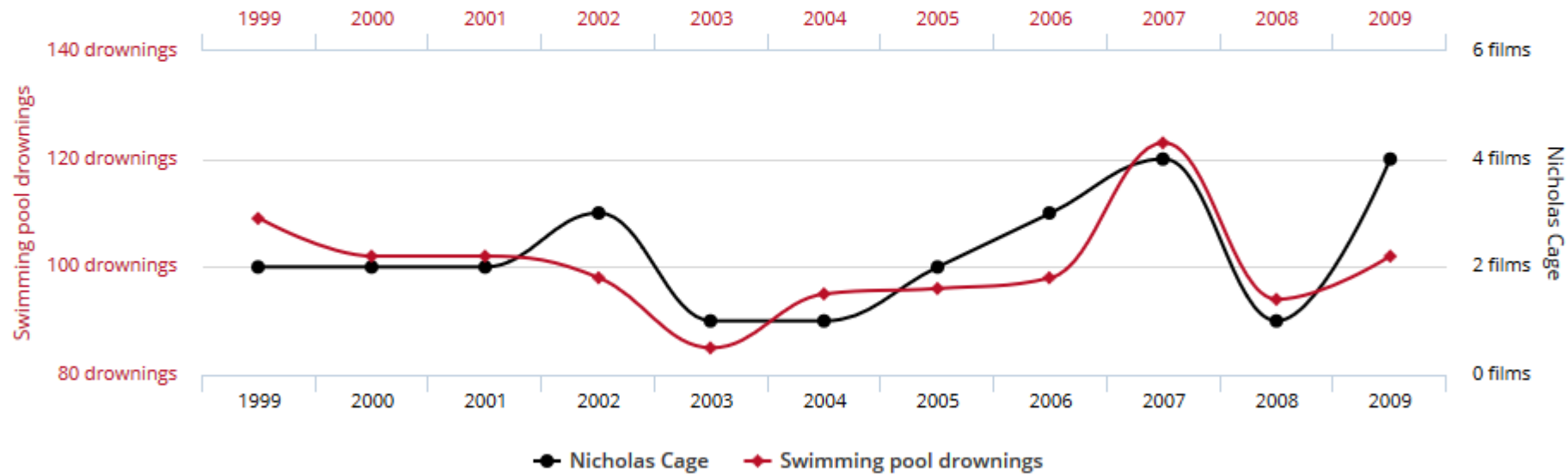
**Correlación y
causalidad son dos
cosas muy diferentes**



REGRESIÓN — CUIDADO!

Number of people who drowned by falling into a pool
correlates with
Films Nicolas Cage appeared in

Correlation: 66.6% ($r=0.666004$, $p>0.05$)



tylervigen.com

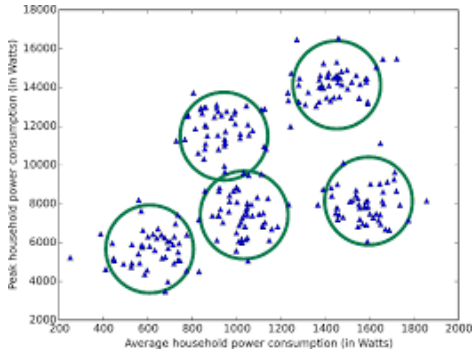
Data sources: Centers for Disease Control & Prevention and Internet Movie Database



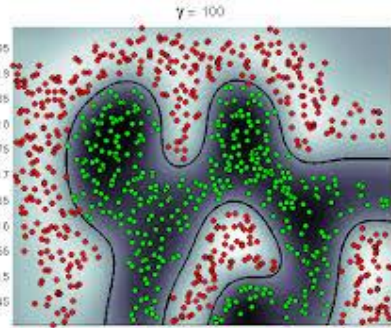
AGENDA



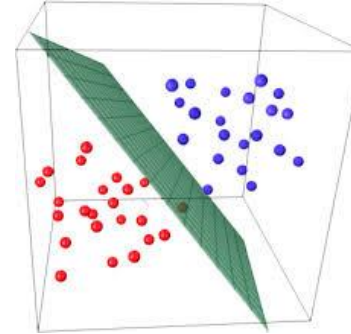
**Aprendizaje
automático**



**Aprendizaje
no supervisado**



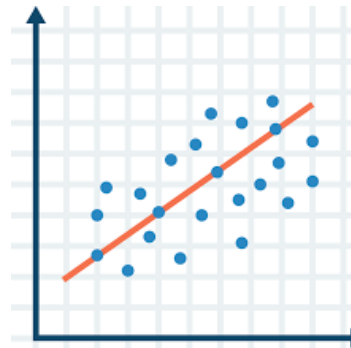
**Aprendizaje
supervisado**



Clasificación



**Métricas de
Evaluación de la
clasificación**



Regresión



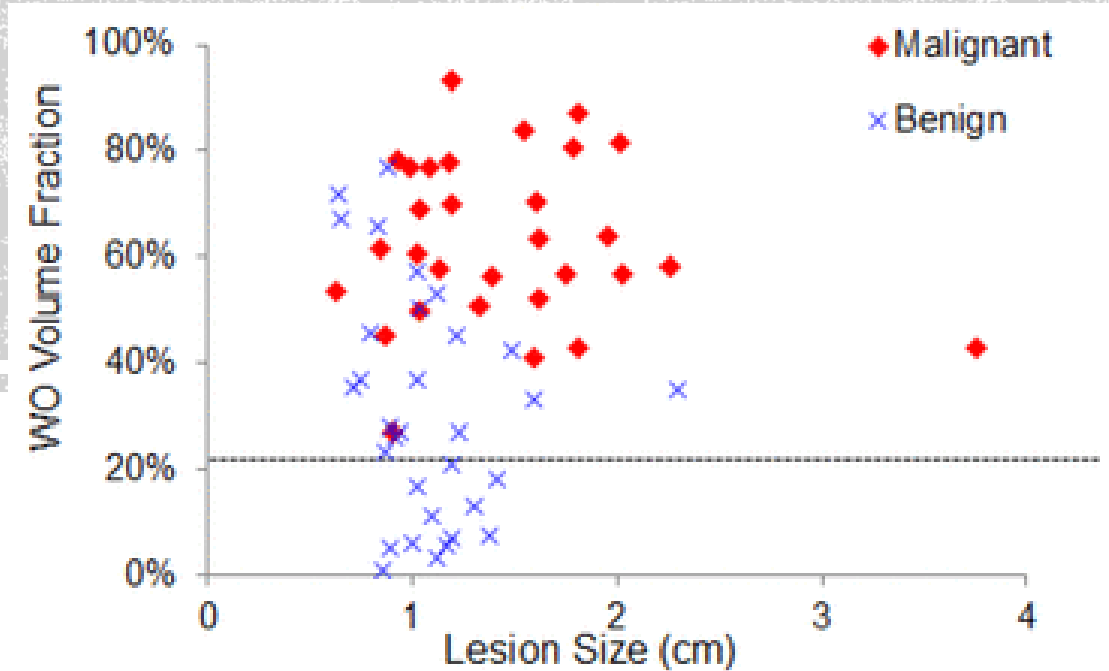
**Métricas de
Evaluación de la
regresión**



CLASIFICACIÓN

A scatter plot titled 'CLASIFICACIÓN' showing the distribution of 'Volume Fraction' (y-axis, 40% to 100%) for two classes: 'Malignant' (red diamonds) and 'Benign' (blue crosses). The x-axis is unlabeled. Malignant cases are generally clustered at higher volume fractions (50-95%), while Benign cases are more spread out (35-75%).

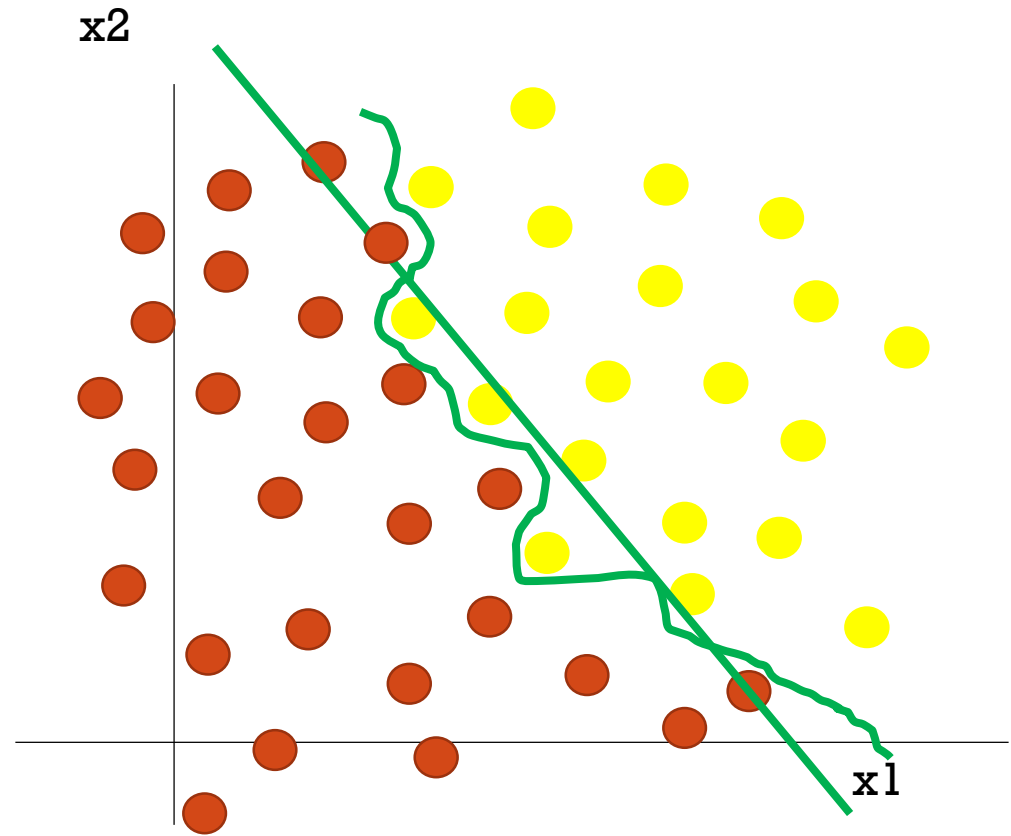
Class	Volume Fraction (%)
Malignant	42
Malignant	50
Malignant	52
Malignant	55
Malignant	56
Malignant	57
Malignant	58
Malignant	58
Malignant	59
Malignant	60
Malignant	62
Malignant	63
Malignant	64
Malignant	65
Malignant	66
Malignant	67
Malignant	68
Malignant	70
Malignant	72
Malignant	73
Malignant	74
Malignant	75
Malignant	76
Malignant	77
Malignant	78
Malignant	78
Malignant	79
Malignant	80
Malignant	81
Malignant	82
Malignant	83
Malignant	84
Malignant	85
Malignant	86
Malignant	87
Malignant	88
Malignant	89
Malignant	90
Malignant	92
Malignant	93
Malignant	94
Malignant	95
Benign	35
Benign	36
Benign	37
Benign	38
Benign	39
Benign	40
Benign	41
Benign	42
Benign	43
Benign	44
Benign	45
Benign	46
Benign	47
Benign	48
Benign	49
Benign	50
Benign	51
Benign	52
Benign	53
Benign	54
Benign	55
Benign	56
Benign	57
Benign	58
Benign	59
Benign	60
Benign	61
Benign	62
Benign	63
Benign	64
Benign	65
Benign	66
Benign	67
Benign	68
Benign	69
Benign	70
Benign	71
Benign	72
Benign	73
Benign	74
Benign	75
Benign	76
Benign	77
Benign	78
Benign	79
Benign	80
Benign	81
Benign	82
Benign	83
Benign	84
Benign	85
Benign	86
Benign	87
Benign	88
Benign	89
Benign	90
Benign	91
Benign	92
Benign	93
Benign	94
Benign	95



http://www.jacmp.org/index.php/jacmp/article/view/5187/html_374

CLASIFICACIÓN

- Encontrar modelos que describan clases para futuras predicciones:
 - KNN
 - Árboles de decisión
 - Regresión logística
 - Redes neuronales
 - ...
- Valores **discretos** de la variable objetivo
- Incluye la estimación de **probabilidades** de clase
- **Baseline**: medida de evaluación dada por un clasificador que escoge siempre la clase mayoritaria



MÉTRICAS DE CLASIFICACIÓN

- Se usa una **matriz de confusión** para evaluar diferentes métricas de correctitud/error
- Se utilizan dos calificadores para describir cada una de sus casillas:
 - Un calificador de la correctitud de la predicción con respecto a la realidad: Verdadero o Falso
 - Un calificador del tipo de la predicción: Positivo o Falso, con respecto a cada clase de interés (i.e churn)
- Dependiendo del contexto los tipos de error pueden ser mas graves que otros (los diferentes)

		Predicción	
		Churn ^P	No churn ^N
Realidad	Churn ⁺	VP	FN - Tipo II
	No churn ⁻	FP - Tipo I	VN

- La diagonal (en verde) muestra las instancias correctamente clasificadas. Las demás casillas resume diferentes tipos de error:
 - Tipo I: Falsos positivos
 - Tipo II: Falsos negativos

¿Qué pasa cuando hay mas de dos clases?



MÉTRICAS DE CLASIFICACIÓN

- Interpretarían el caso de la detección de un email spam

TP, TN:

FP: , consecuencia:

FN: , consecuencia:

- Interpretar el caso del diagnóstico de una enfermedad grave?

TP, TN:

FP: , consecuencia:

FN: , consecuencia:

		Predicción	
		Churn ^P	No churn ^N
Realidad	Churn ⁺	VP	FN - Tipo II
	No churn ⁻	FP - Tipo I	VN

- Interpretar el caso de la prospección de clientes de un crédito de consumo (baja aceptación)

TP, TN:

FP: , consecuencia:

FN: , consecuencia:



MÉTRICAS DE CLASIFICACIÓN

- Tasa de correctitud (*accuracy*) = $(VP+VN)/(VP+VN+FP+FN)$
- Error de mala clasificación (contrario de *accuracy*) = $(FP+FN)/(VP+VN+FP+FN)$: probabilidad de error
- Precisión = $VP / (VP+FP)$: valor de predicción positiva, $P(\text{Real+} | \text{Predicho+})$
- *Recall* (o TPR o sensibilidad) = $VP / (VP+FN)$: qué proporción de todos los positivos reales pude identificar como tal, $P(\text{Predicho+} | \text{Real+})$
- Especificidad (o TNR): $VN / (VN+FP)$: qué proporción de todos los negativos reales pude identificar como tal, $P(\text{Predicho-} | \text{Real-})$
- Valor de predicción negativa (FPR) = $VN / (VN+FN)$
- F1-Measure = $2 * \frac{\text{precision} * \text{recall}}{\text{precision} + \text{recall}}$ (promedio armónico)

		Predicción	
		Churn ^P	No churn ^N
Realidad	Churn ⁺	VP	FN - Tipo I
	No churn ⁻	FP - Tipo I	VN

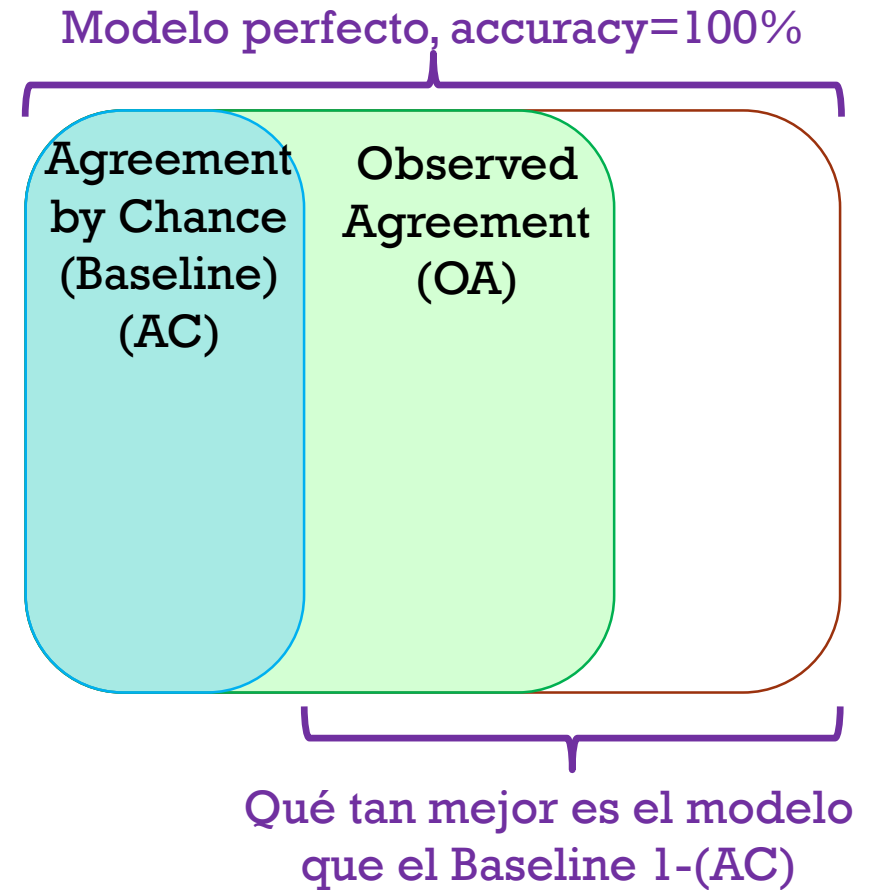
Imaginemos el problema de detección de spam mail e interpretemos cada métrica

Imaginemos el problema de diagnóstico de cáncer e interpretemos cada métrica



MÉTRICAS DE CLASIFICACIÓN

- Coeficiente de concordancia **Kappa**
 - Para datos nominales u ordinales
 - Concordancia entre las predicciones y las clases reales
 - Sustraer el efecto de concordancia por suerte (AC) del valor del **accuracy** (concordancia observada - OA)
 - Valores van de 0 a 1
 - Muy útil sobretodo cuando las clases no están balanceadas
 - Diagnóstico de enfermedades raras
 - Clientes que acepten productos de crédito)
 - $$\text{Kappa} = \frac{OA - AC}{1 - AC}$$



MÉTRICAS DE CLASIFICACIÓN

■ Coeficiente de concordancia **Kappa**

- Para datos nominales u ordinales
- Concordancia entre las predicciones y las clases reales
- Sustrae el efecto de concordancia por suerte (AC) del valor del **accuracy** (concordancia observada - OA)
- Valores van de 0 a 1
- Muy útil sobretodo cuando las clases no están balanceadas
 - Diagnóstico de enfermedades raras
 - Clientes que acepten productos de crédito)

		Predicciones		TOTAL
		+	-	
reales	+	10	4	14
	-	3	2	5
TOTAL		13	6	19

OA = 0,63

AC = 0,59

Kappa = 0,11

Accuracy (OA) = $(10+2)/19=0,63$

(AC) = $(13/19 * 14/19) + (6/19 * 5/19) = 0,59$

Kappa = $(OA-AC)/(1-AC) = 0,11$

		Predicciones		TOTAL
		+	-	
reales	+	0	3	3
	-	0	97	97
TOTAL		0	100	100

OA = 0,97

AC = 0,97

Kappa = 0,00

Accuracy (OA) = $(0+97)/100=0,97$

(AC) = $(0/100 * 3/100) + (100/100 * 97/100) = 0,97$

Kappa = $(OA-AC)/(1-AC) = 0$

		Predicciones		TOTAL
		+	-	
reales	+	1475	988	2463
	-	556	1981	2537
TOTAL		2031	2969	5000

OA = 0,69

AC = 0,50

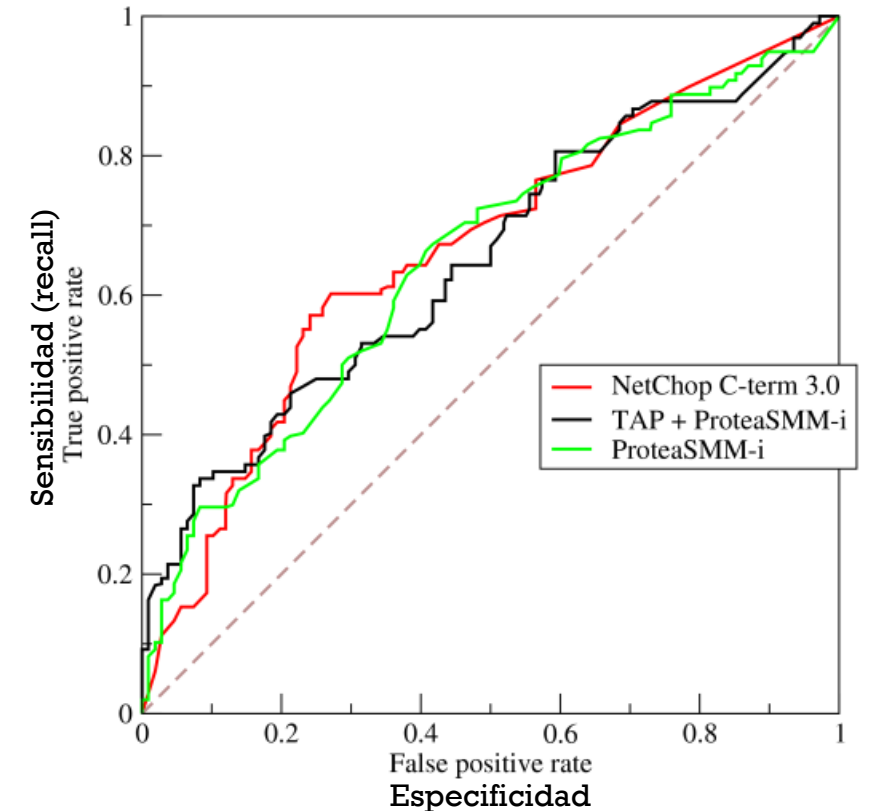
Kappa = 0,38



MÉTRICAS DE CLASIFICACIÓN

Comparación de varios modelos:

- Validación cruzada con accuracy o error de mala clasificación
- Medida-F=
$$2 * (\text{Precisión} * \text{Recall}) / (\text{Precisión} + \text{Recall})$$
- Área debajo de la curve ROC
- Todas las métricas son insensibles a los diferentes costos del error de las diferentes clases, que depende del contexto. Pueden tener impactos diferentes.



Wikipedia.org



MÉTRICAS DE CLASIFICACIÓN

TALLER: CÁLCULO DE MÉTRICAS

Calcule las métricas de evaluación de un modelo de clasificación cuyos resultados están reflejados en la tabla siguiente

- Error, Accuracy y Kappa global
- Precisión, Recall, especificidad, F-Measure de cada clase

	PREDICCIÓN				
REAL	Esporádico	Fiel	Parcial	Promocional	Total
Esporádico	61	8	1	0	70
Fiel	0	56	17	0	73
Parcial	0	0	15	0	15
Promocional	0	0	0	24	24
Total	61	64	33	24	182

REFERENCIAS

- *Introduction to Statistical Learning with Applications in R (ISLR)*, G. James, D. Witten, T. Hastie & R. Tibshirani, 2014
- *Data Mining (4th Edition)*, Ian Witten, Eibe Frank, Mark A. Hall & Christopher J. Pal, Elsevier, 2016
- *Machine Learning*, Tom M. Mitchell, McGraw-Hill, 1997
- *Data Science for Business*, Foster Provost & Tom Fawcett, O'Reilly, 2013
- *False positives, false negatives and confusion matrices*, Carlos Guestrin, 2017
- <http://www.cs.waikato.ac.nz/ml/weka/mooc/dataminingwithweka/>
- <https://www.ibm.com/developerworks/library/os-weka1/>

