

- analysis data exploration
- MANIPULACIÓN DATOS - PANDAS
visualización
 - [dataframes
drop
loc
seaborn
 - AGROUPAMIENTO - CLUSTERING
 - [k-means
jerárquico
GMM
GMM + BIC → How many clusters?
 - Reducción Dimensionalidad -
 - [PCA
Selección atributos
Proyección datos multi-dimensionales]

TÉCNICAS DE REGRESIÓN - supervisadas

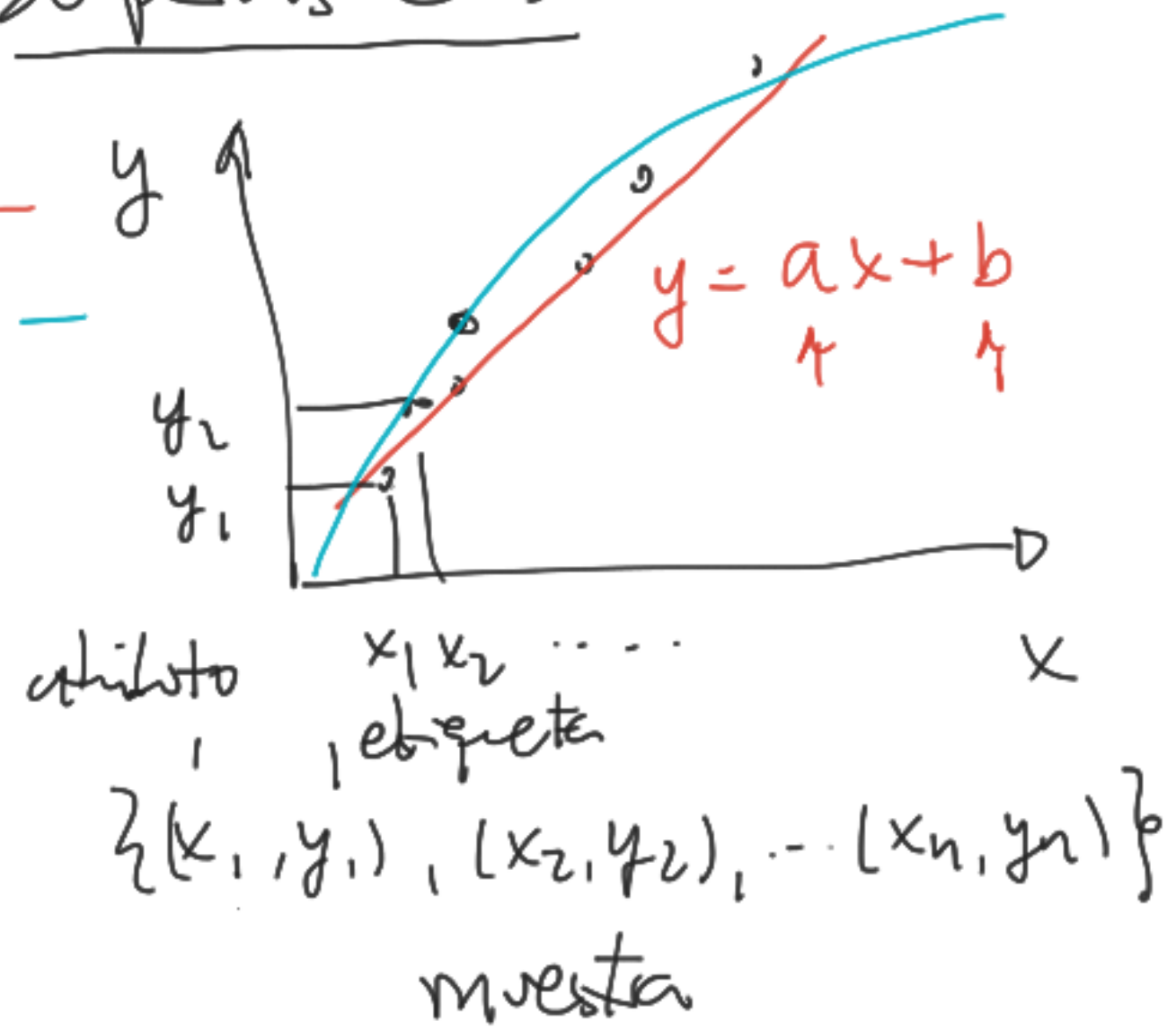
- datos univariados

linear —
no linear —

- selección modelos

(BIC - jerárquico)

- Árboles de decisión

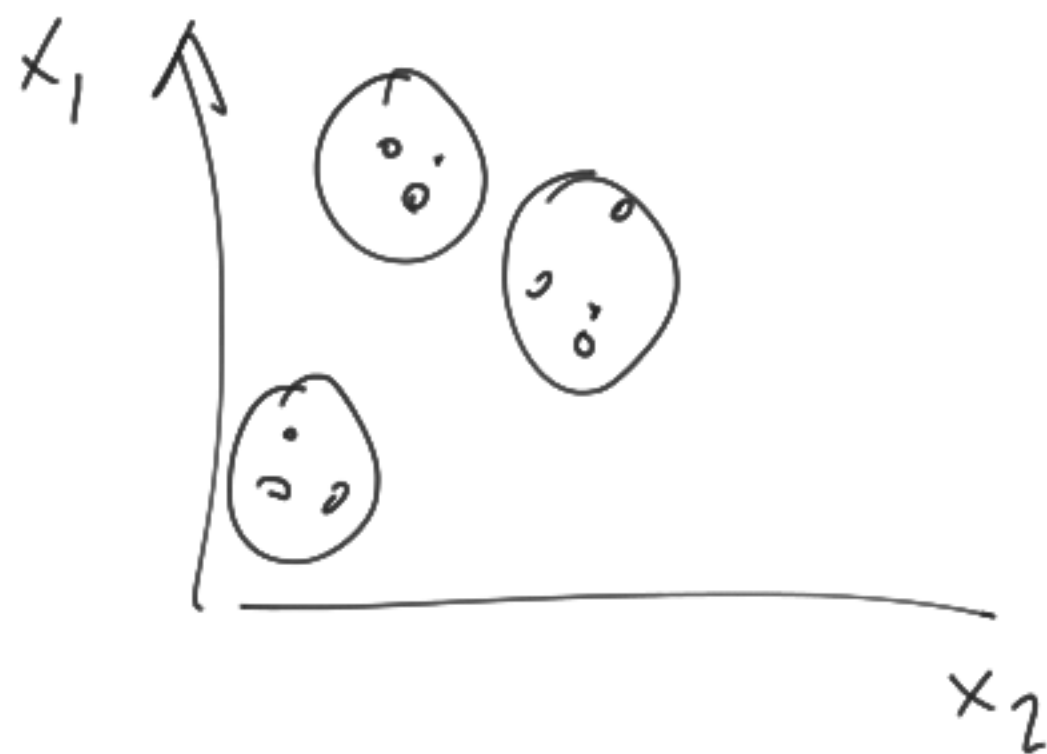


↓
estimación estadística
parámetros modelo
 a, b

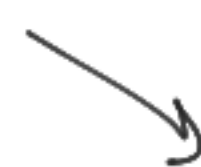
no-supervised



CLUSTERING



supervised

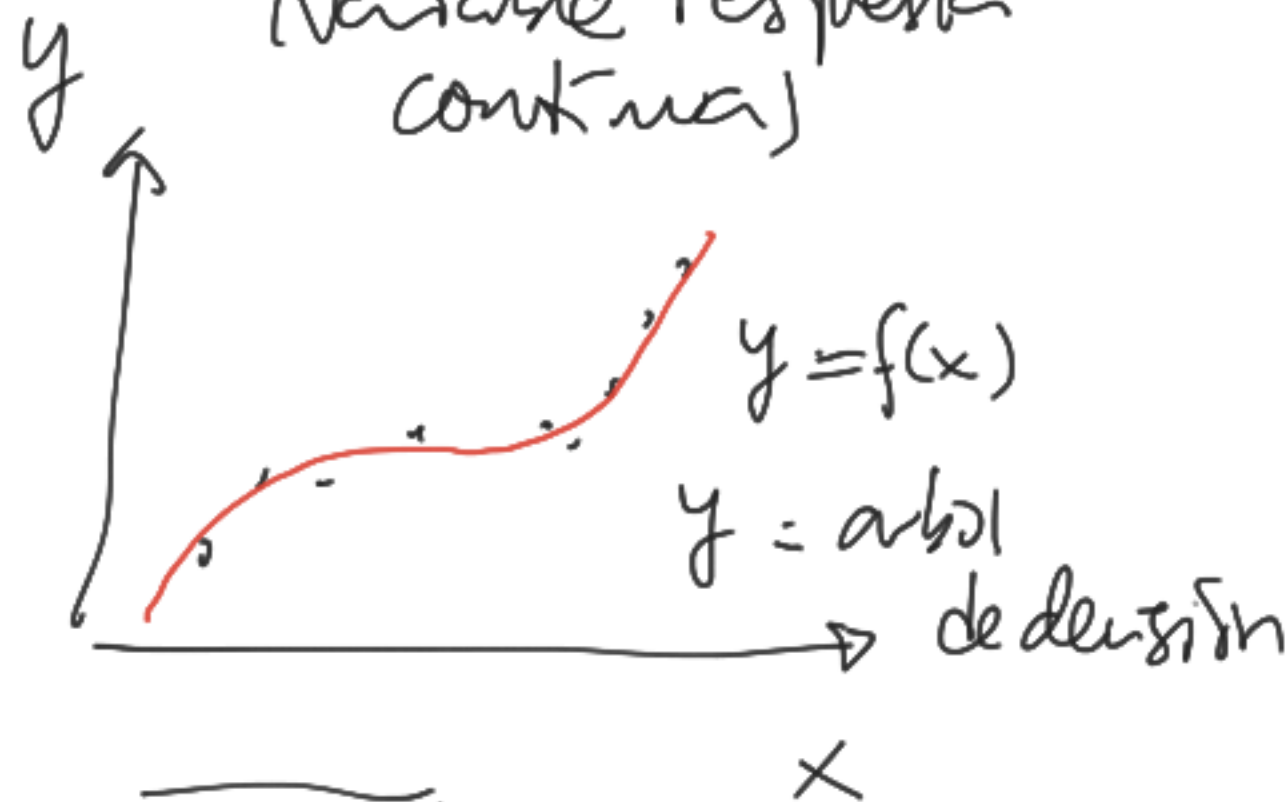


Regression
(variable response
continuous)

Classification

$A_{N \times M}$

w - etiquetas
clase



$x \rightarrow \boxed{\text{algoritmo de densidad}} \rightarrow \hat{y}$

$\vec{x} \rightarrow \boxed{C} \rightarrow \hat{w}$

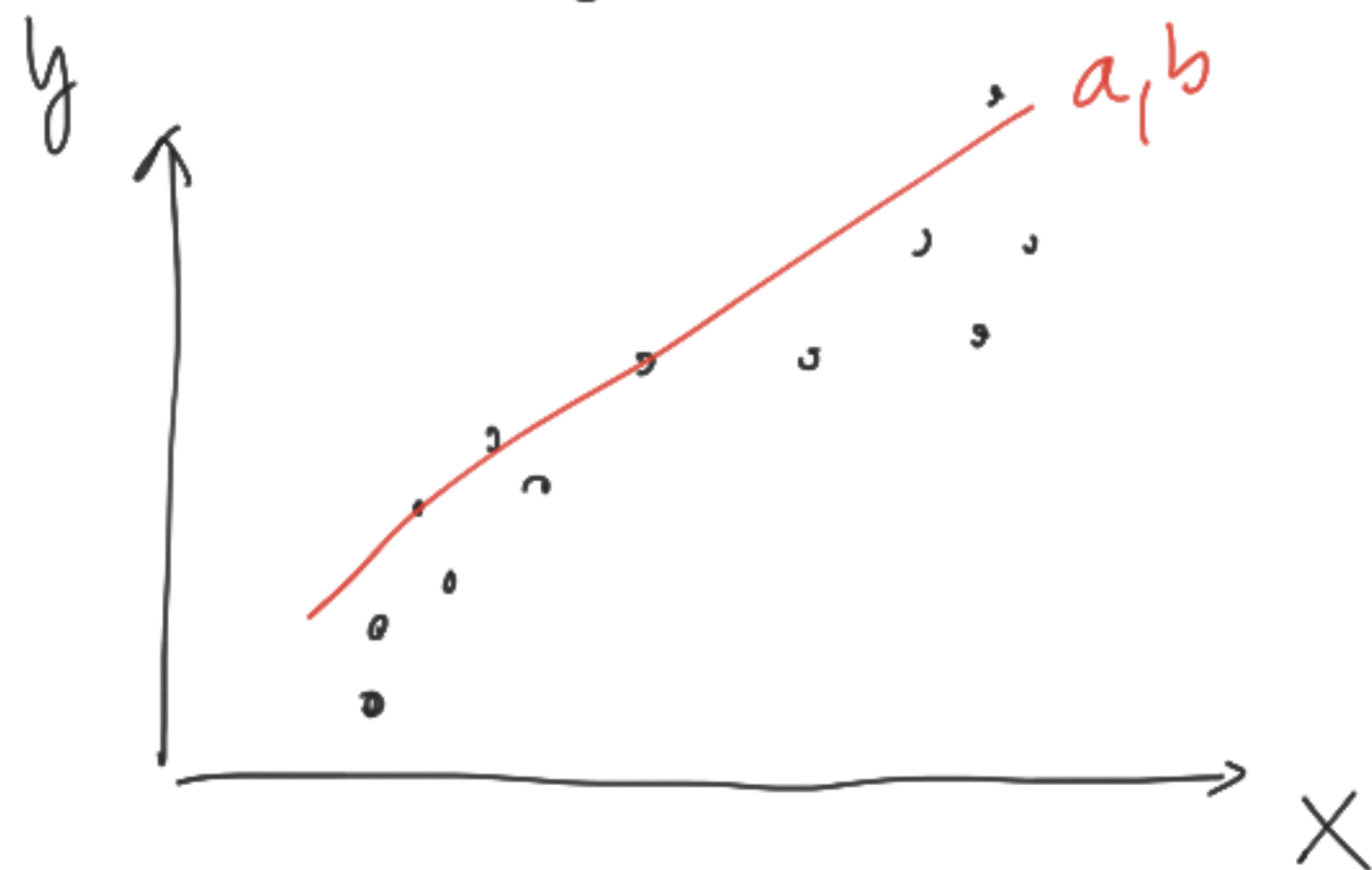
DATOS NO LINEALES

y | x
| |
| |
| |
| |
| |

- estimar
- entrenar
- ajustar

fit

a, b



$x \rightarrow$ modelo $\rightarrow \hat{y}$
predict

DATOS NO LINEALES

atributos

y | x_1, \dots, x_M
~
 N | $A_{N \times M}$
respuestas

datos $x_1, x_2, x_3, \dots, x_N \rightarrow \hat{\mu} = \frac{1}{N} \sum_{i=1}^N x_i$

$$\tilde{x}_1 = x_1 - \hat{\mu}$$

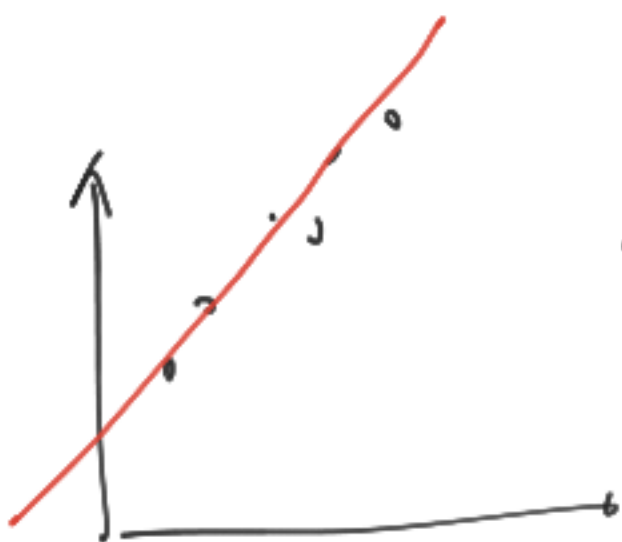
$$\tilde{x}_2 = x_2 - \hat{\mu}$$

\rightarrow

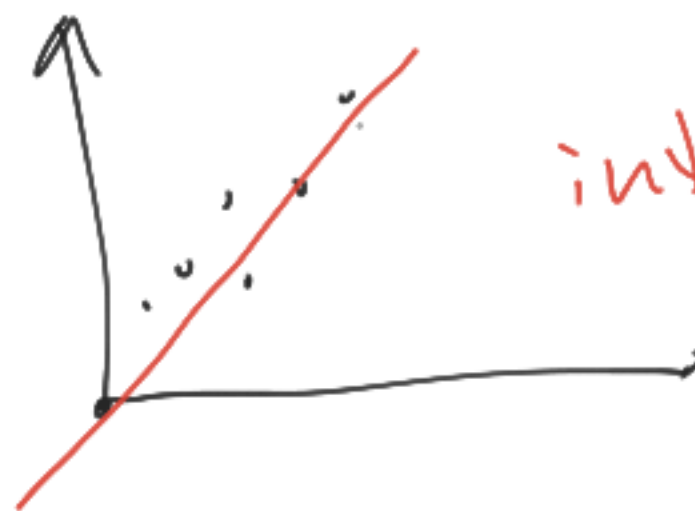
$$\{\tilde{x}_1, \tilde{x}_2, \tilde{x}_3, \dots, \tilde{x}_N\}$$

datos centrados (a)

(media nula)

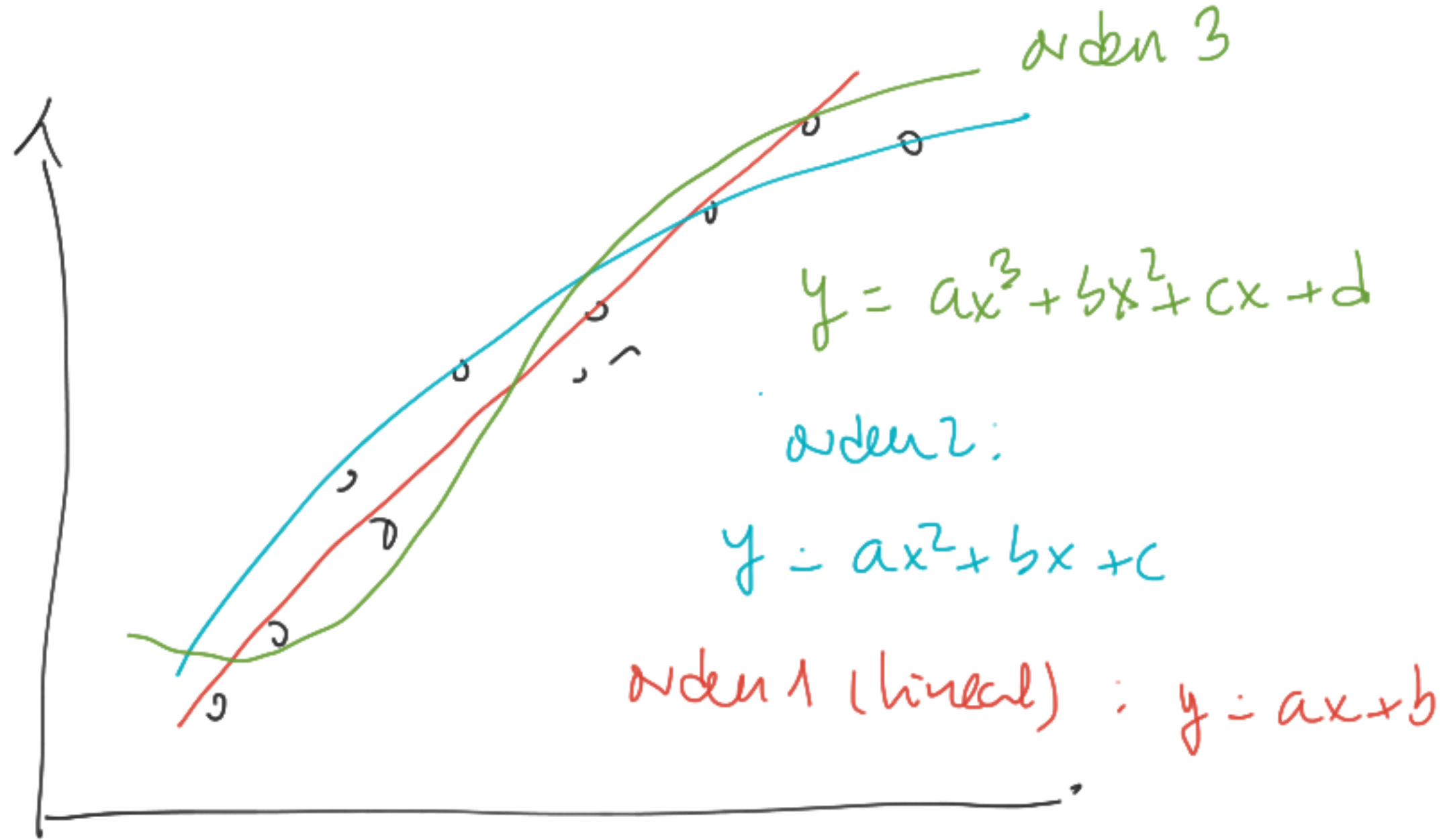


(a, b)

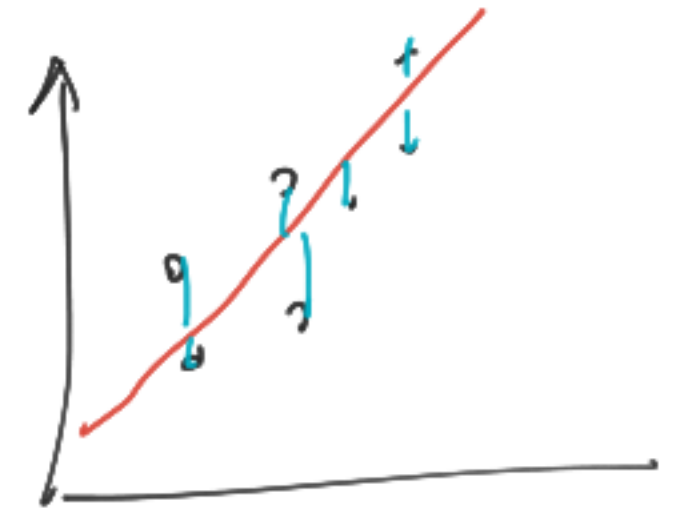


intersección = 0 (b=0)

SELECCIÓN MODELOS REGRESIÓN → INDICES PARSIMONIA



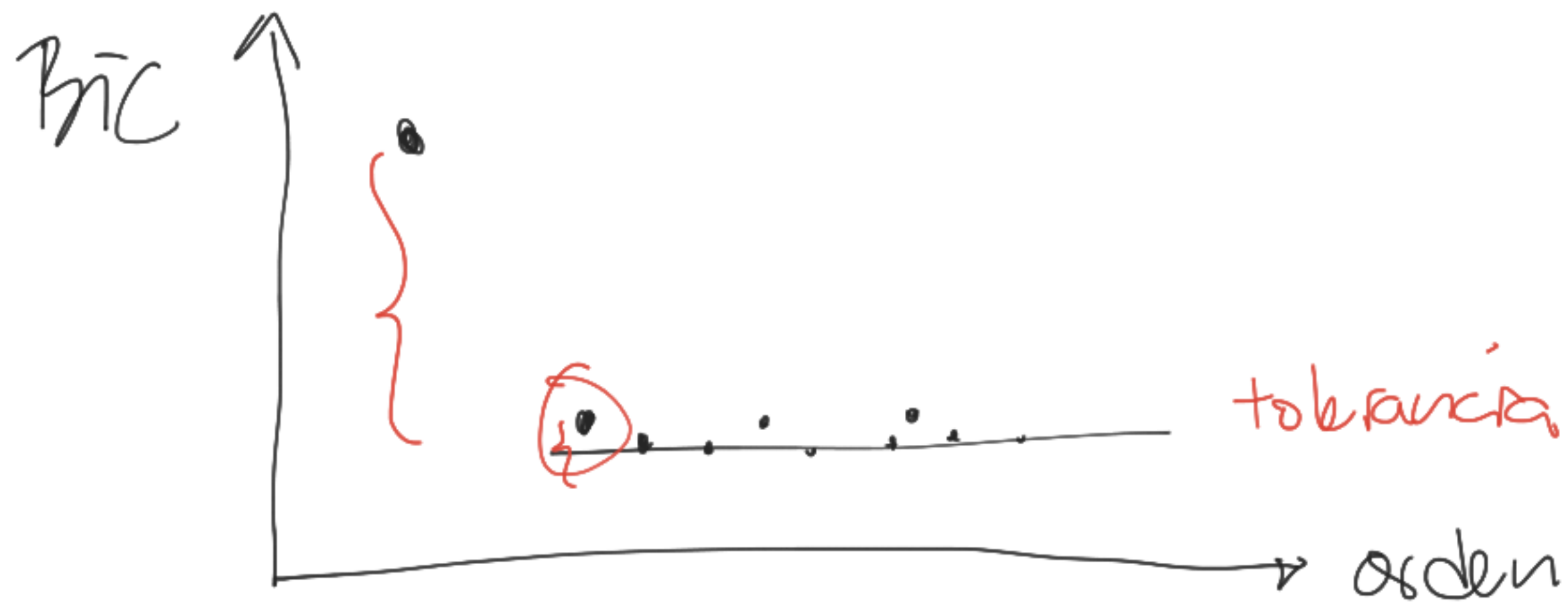
ERROR + COMPLEJIDAD
(# parámetros)



$(a, b) - 2$

$\sum \text{residuos}^2$

$\left. \begin{array}{l} Aic \\ Bic \end{array} \right\}$



$$\{x_1, x_2, \dots, x_n\} \rightarrow \hat{\mu} = \frac{1}{n} \sum_{i=1}^n x_i$$

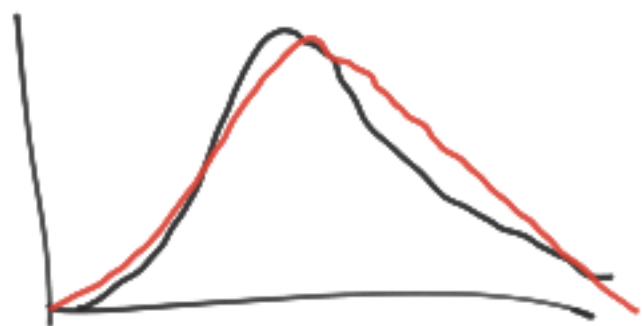
error estándar estimación media $\hat{\mu}$

$$S.e.m = \frac{s}{\sqrt{n}}$$

$$\alpha = 0.05 \quad (95\%)$$

Intervalo
Confianza

$$\left[\hat{\mu} - \underbrace{t_{\frac{\alpha}{2}, n-1}}_{\text{tablas distribución } t\text{-Student}} \cdot \frac{s}{\sqrt{n}} , \hat{\mu} + t_{\frac{\alpha}{2}, n-1} \cdot \frac{s}{\sqrt{n}} \right]$$



tablas distribución
t-Student

(~Gaussiana con
tamaño n finito)

$$y = ax + b$$

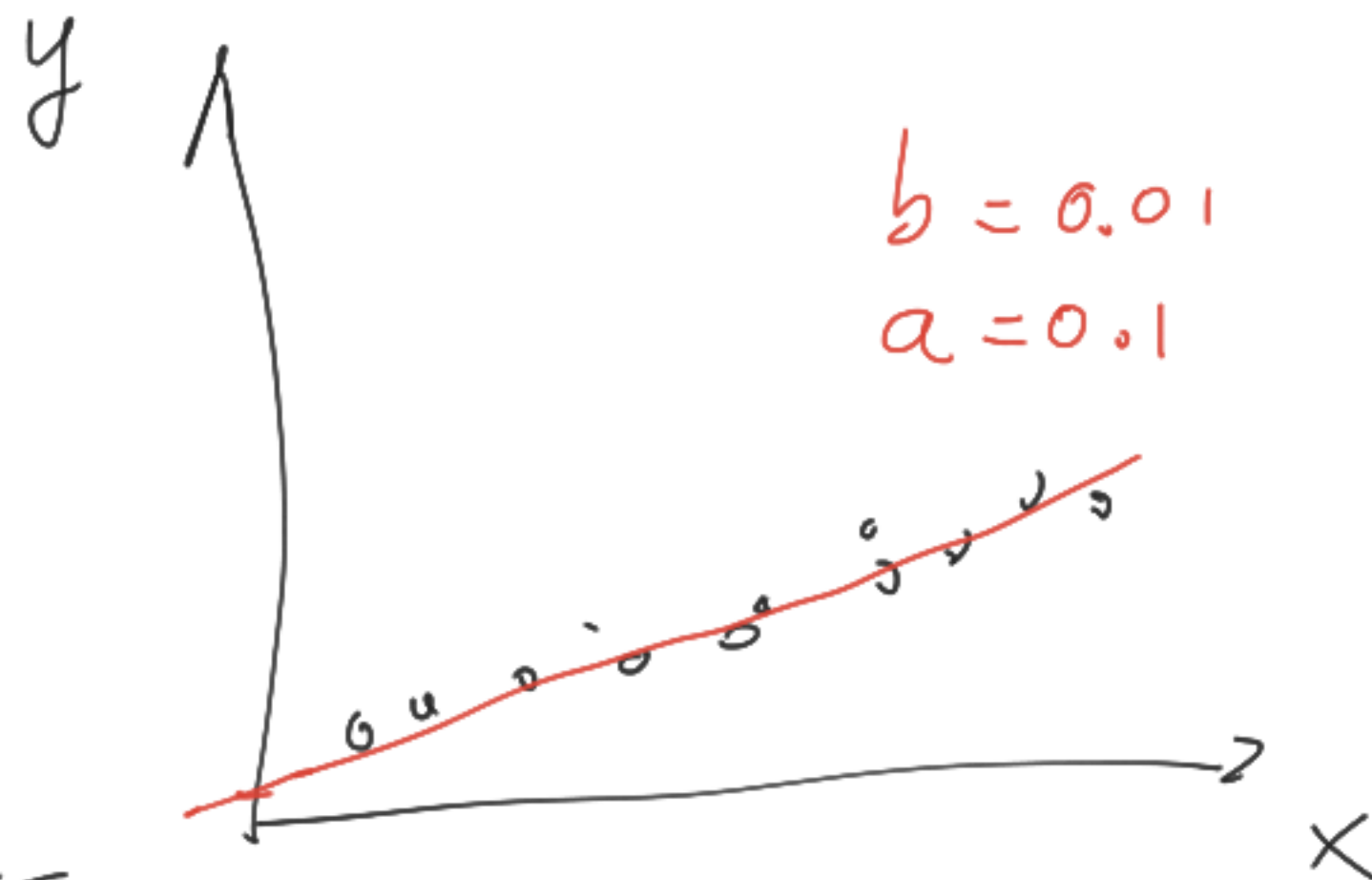
$$\{(x_1, y_1), (x_2, y_2) \dots, (x_n, y_n)\}$$

$\hat{a} \rightarrow$ intervalo confianza 95%

$\hat{b} \rightarrow [-0.02, 0.02]$

no se puede decir
que no hay
crecimiento de y con x

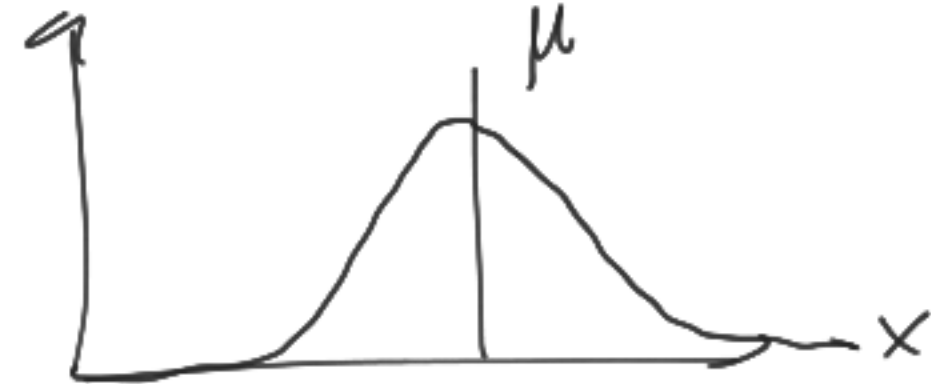
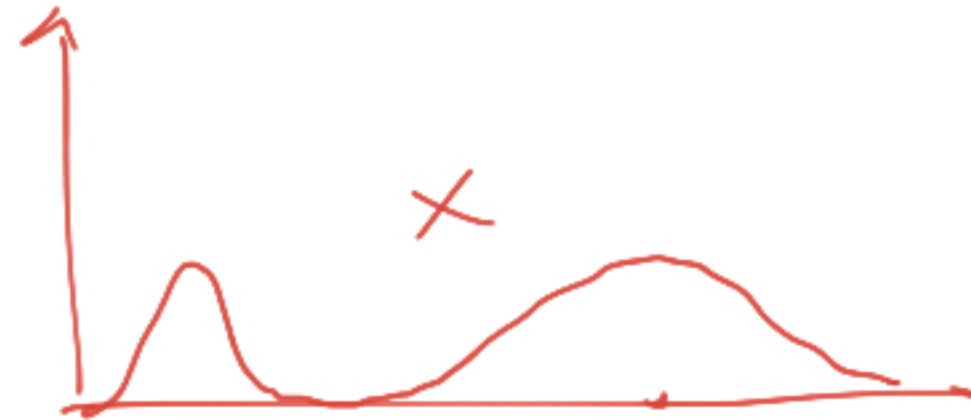
\times
 $[-0.08, 0.13]$ (??)



$a > 0$ es estadística-
mente significativo
 \downarrow
 y crece con x

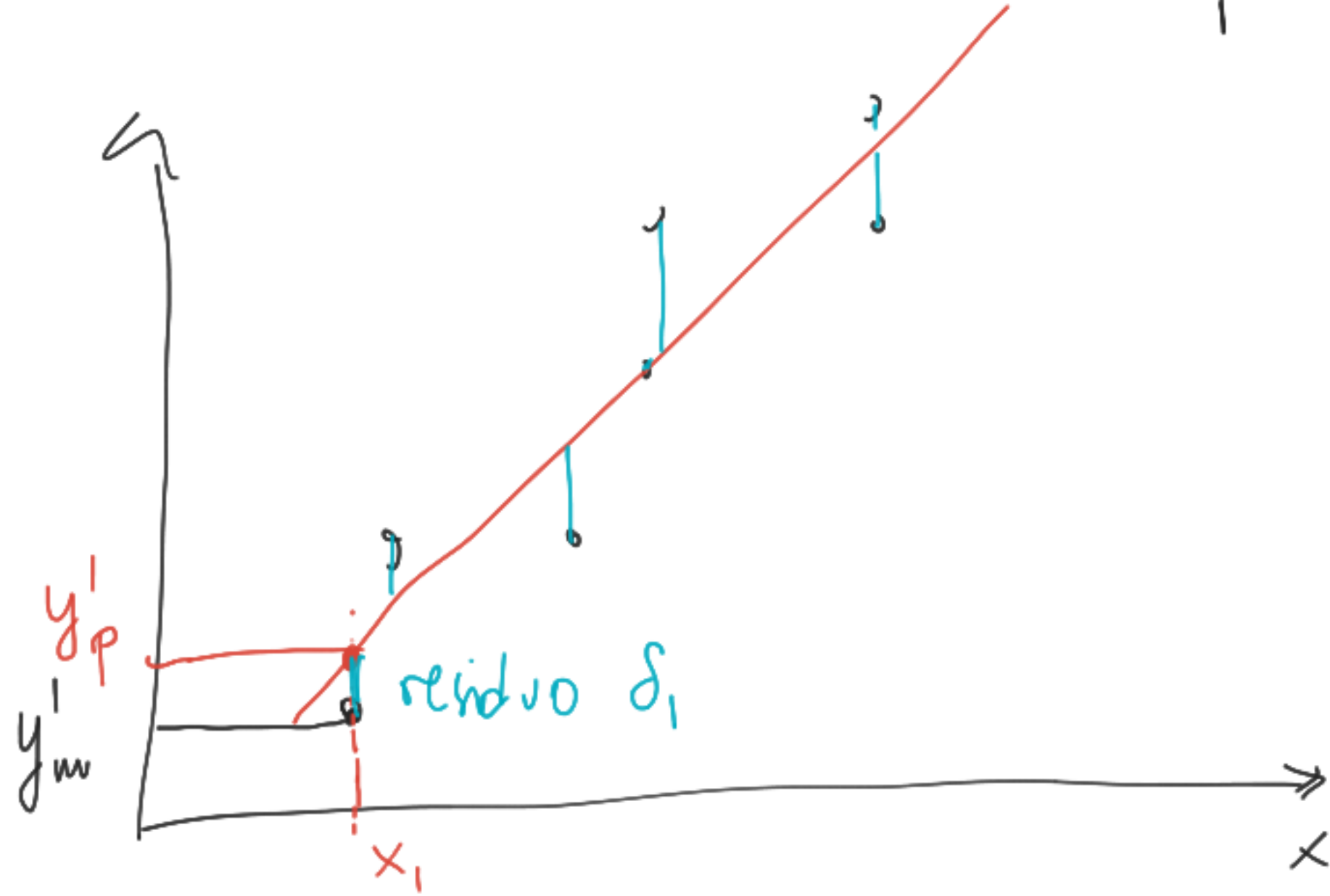


distribuição normal - gaussian

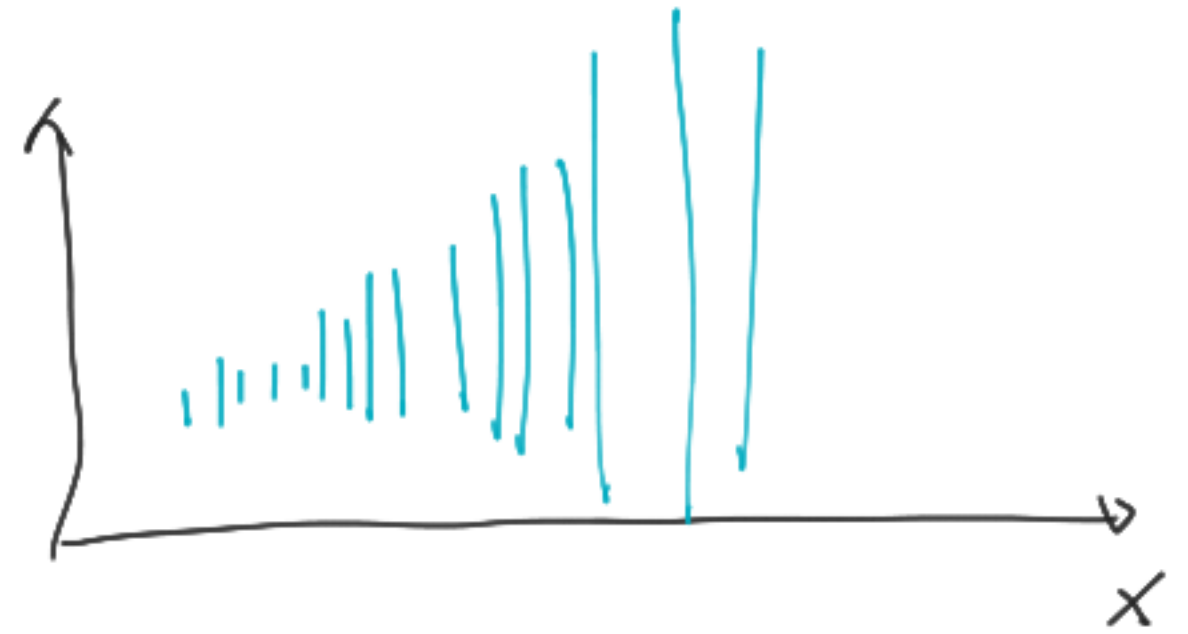


DISTRIBUCIÓN RESIDUOS

$$y_{\text{predicción}} - y_{\text{investida}} = \delta$$

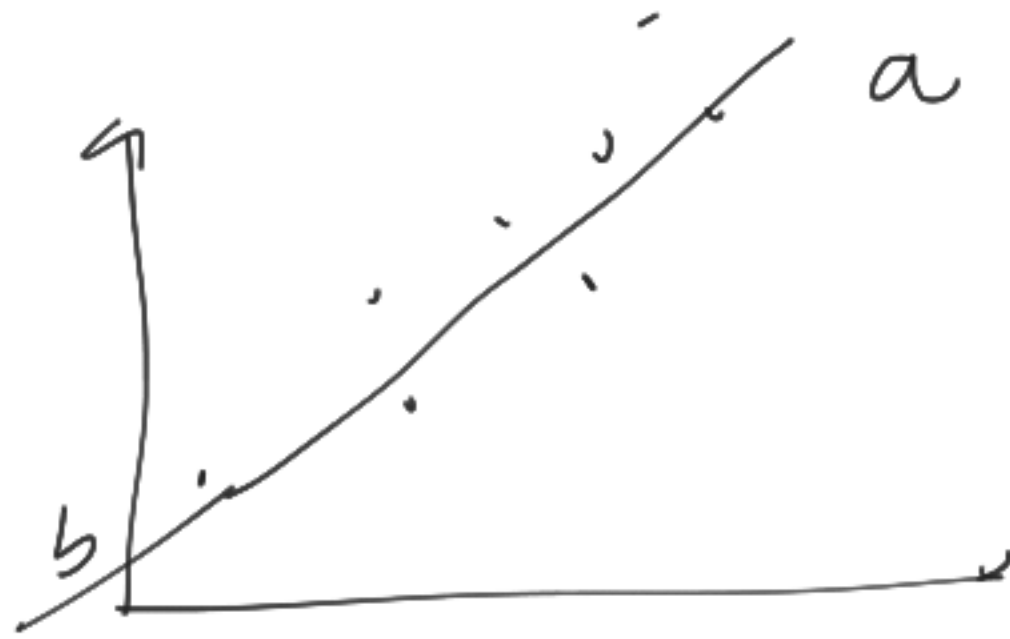


$$y'_p = \text{model.predict}(x_1)$$

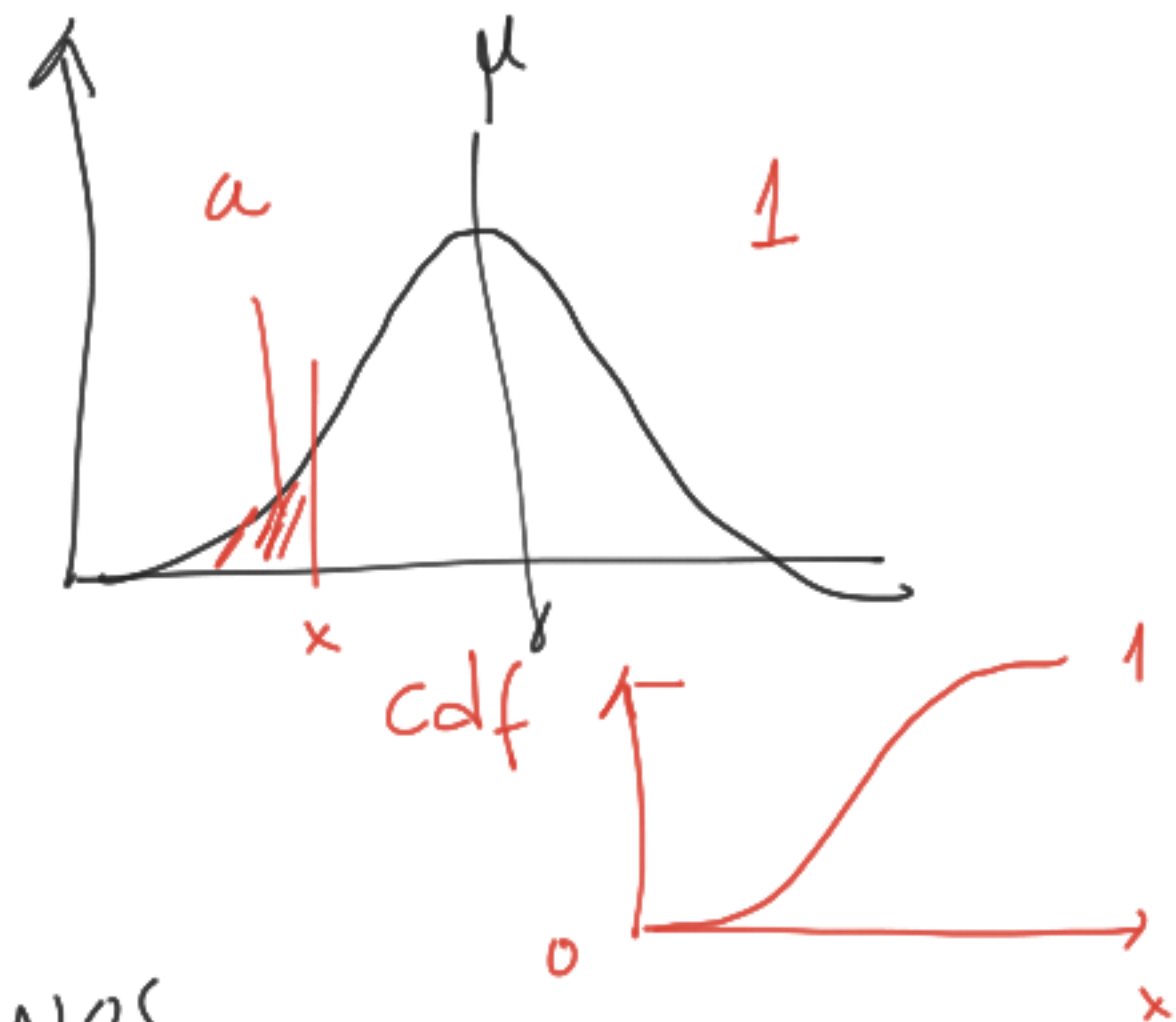
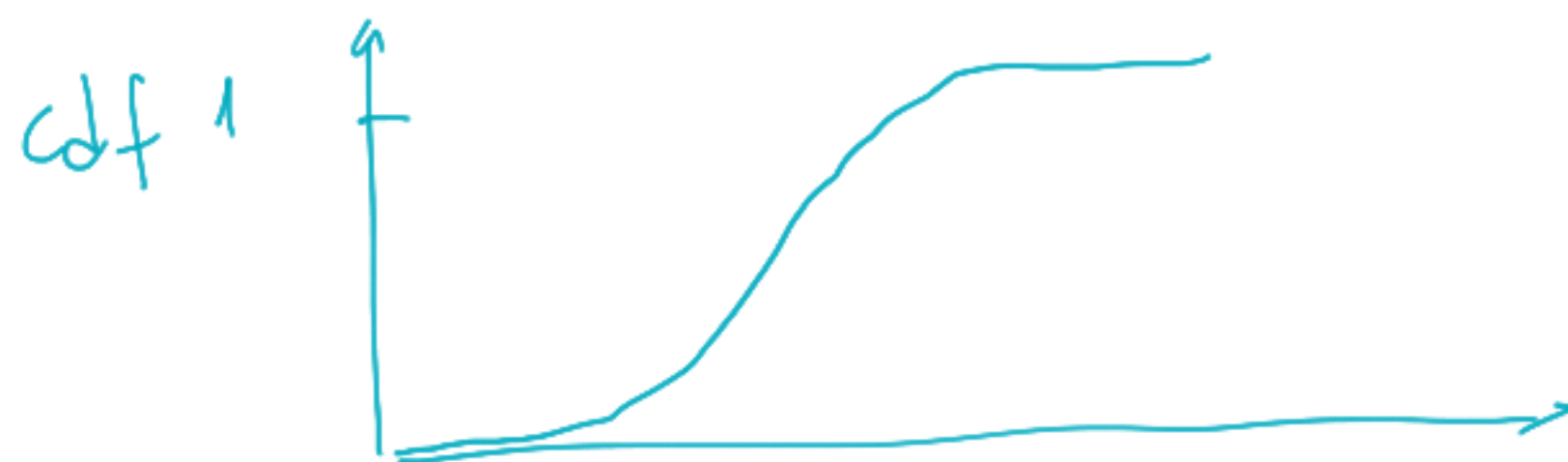
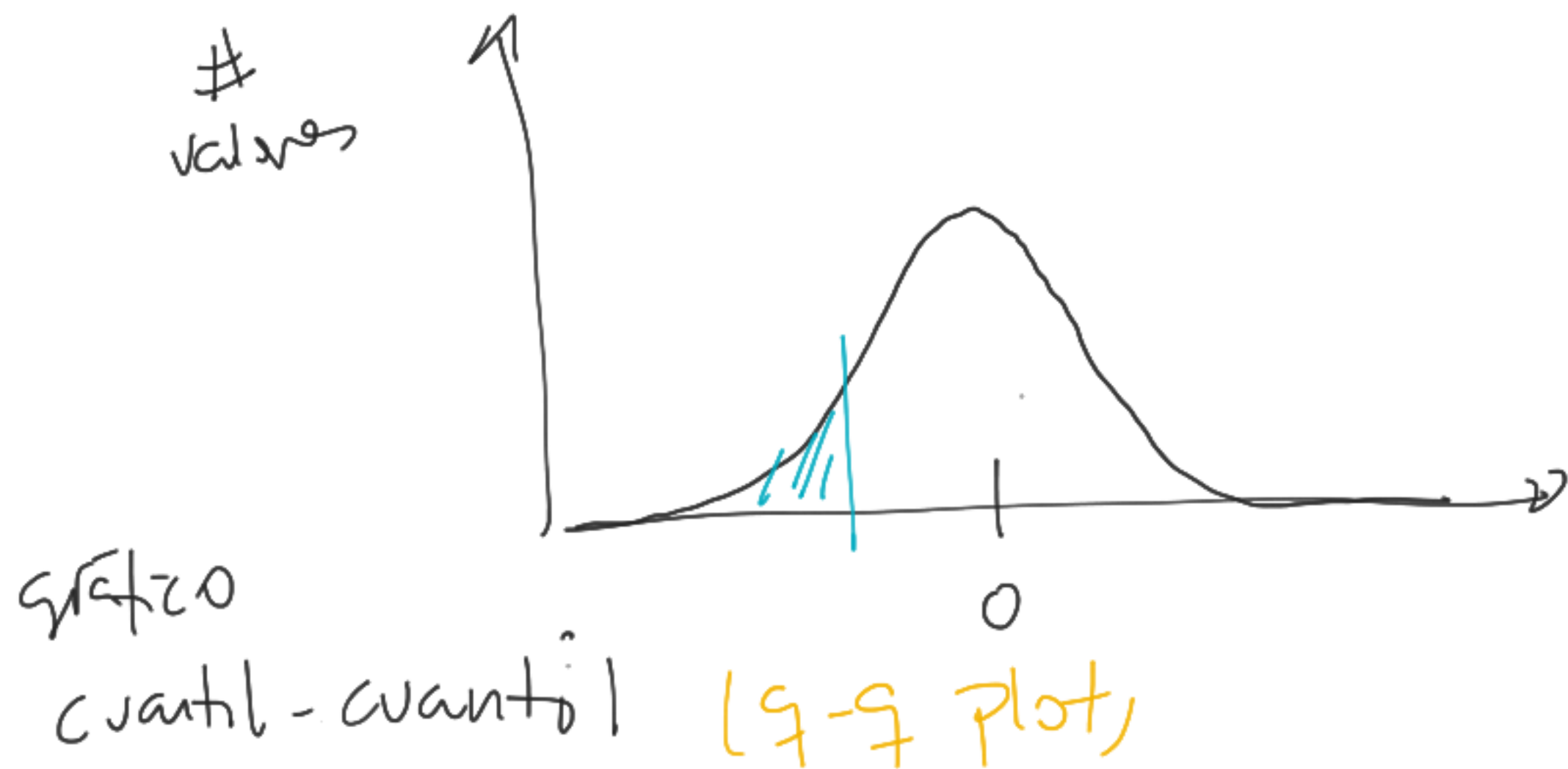


$$y = \underset{1 \times 1}{(x \ 1)} \cdot \underset{2 \times 1}{\begin{pmatrix} a \\ b \end{pmatrix}} + \underset{\beta}{(\varepsilon)} = ax + b + \varepsilon$$

add-column
↓



Residuals distribution normal

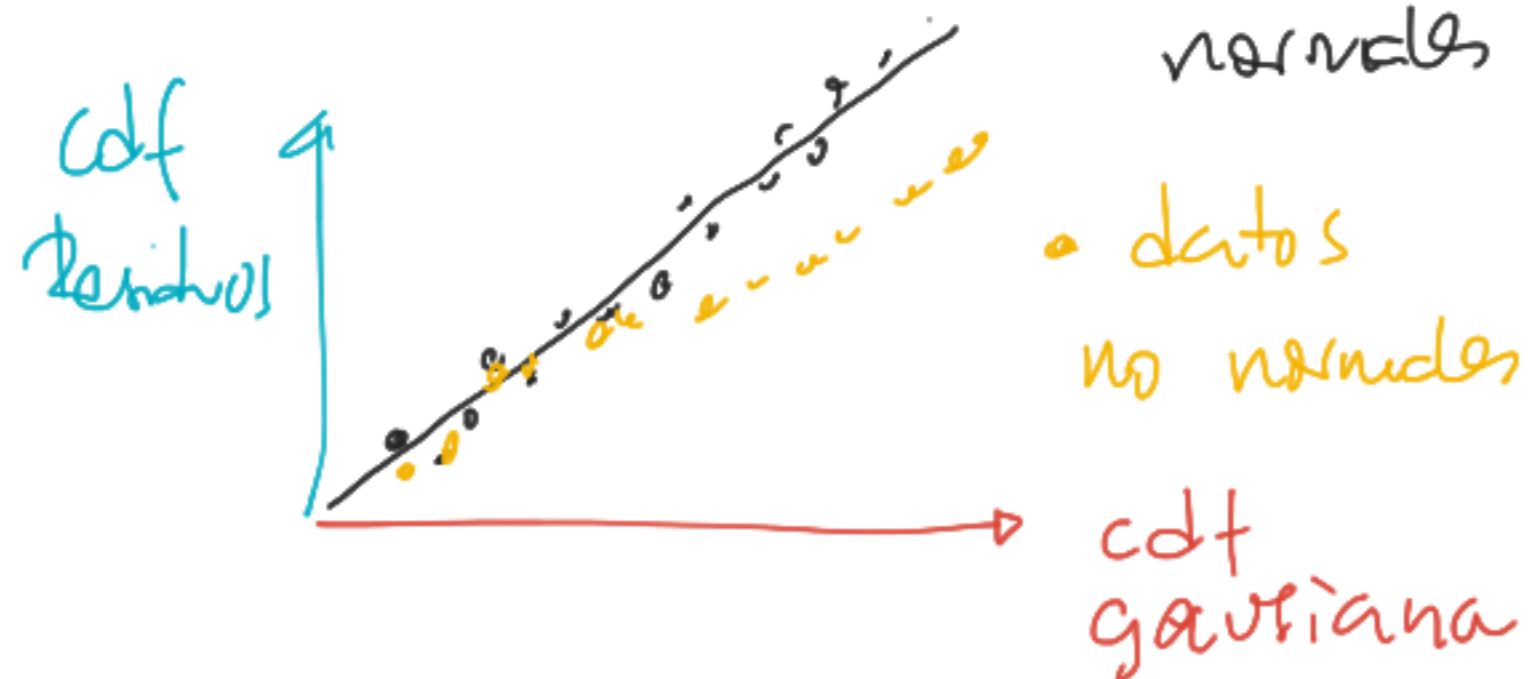


residuos

$$\Delta_i = y_i - \hat{y}_i$$

• datos normales

• datos no normales



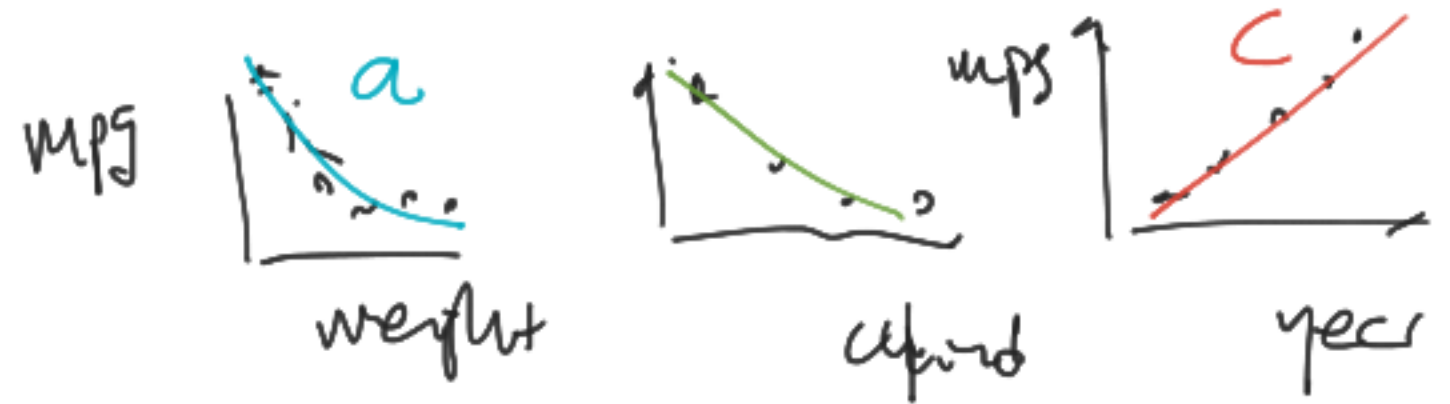
univariate $y = f(x)$

multivariate $y = f(x_1, x_2, x_3 \dots)$

Ex: mpg cars \rightarrow

mpg = $f(\text{weight}, \text{cylinders}, \overset{\text{model}}{\text{year}})$

pairplot



$$\text{mpg} = \underset{a < 0}{(a)} \cdot \text{weight} + \underset{b < 0}{(b)} \cdot \text{cylinders} + \underset{c > 0}{(c)} \cdot \text{year} + d$$

Regresión con árboles de decisión (Regression trees)

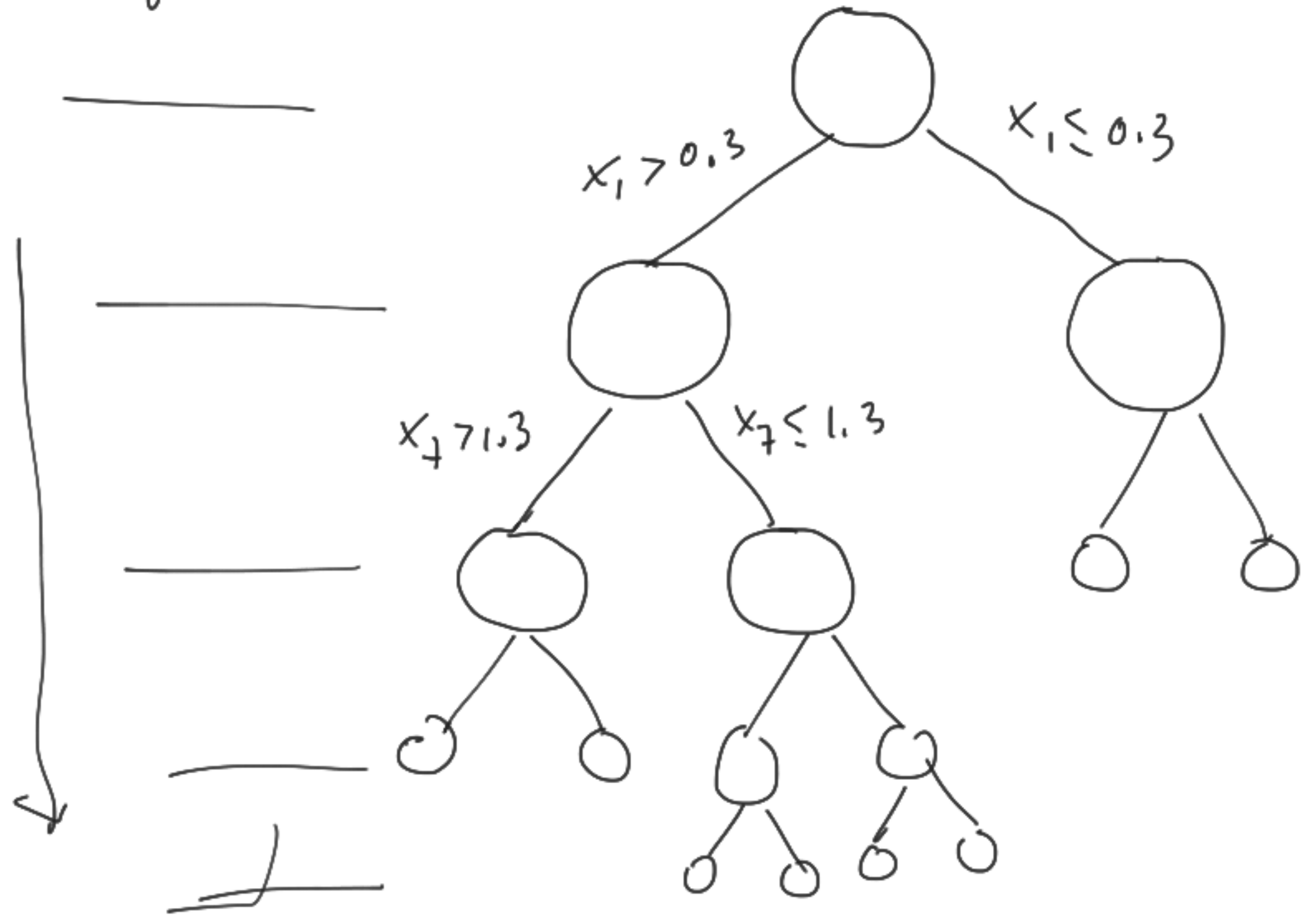
Datos x_1 $x_2 \dots x_n$ y

obs 1

⋮

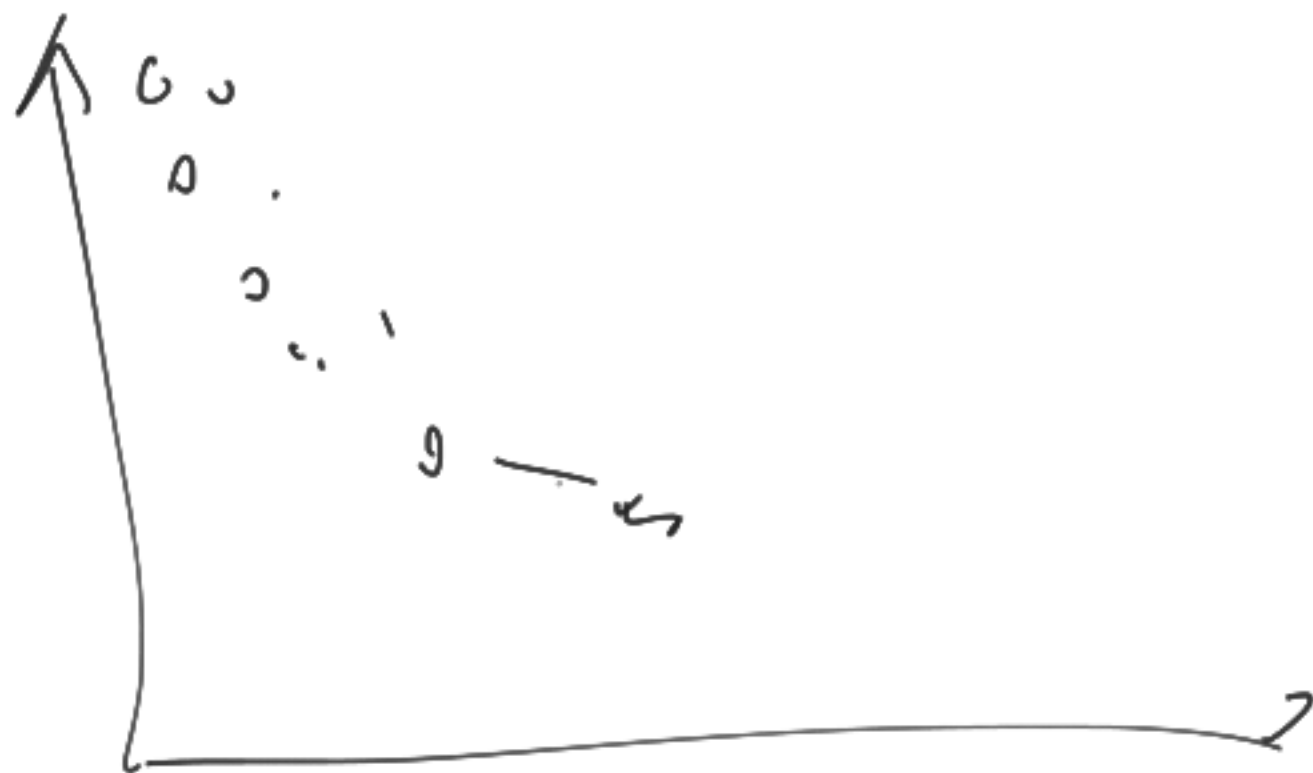
obs n

tree
depth



Decision tree regressor

mpg ~ weight



- Ejercicio:
- load multivariate data
 - Clean data → NANS, variable tipo string... (M)
→ (d)
 - Dimensionality Reduction (PCA)
(90-95% varianza)
 - Projectar datos espacio PCA reducido
 - Clustering (k-means, hierarchical, gmm)
 - Visualizar resultados