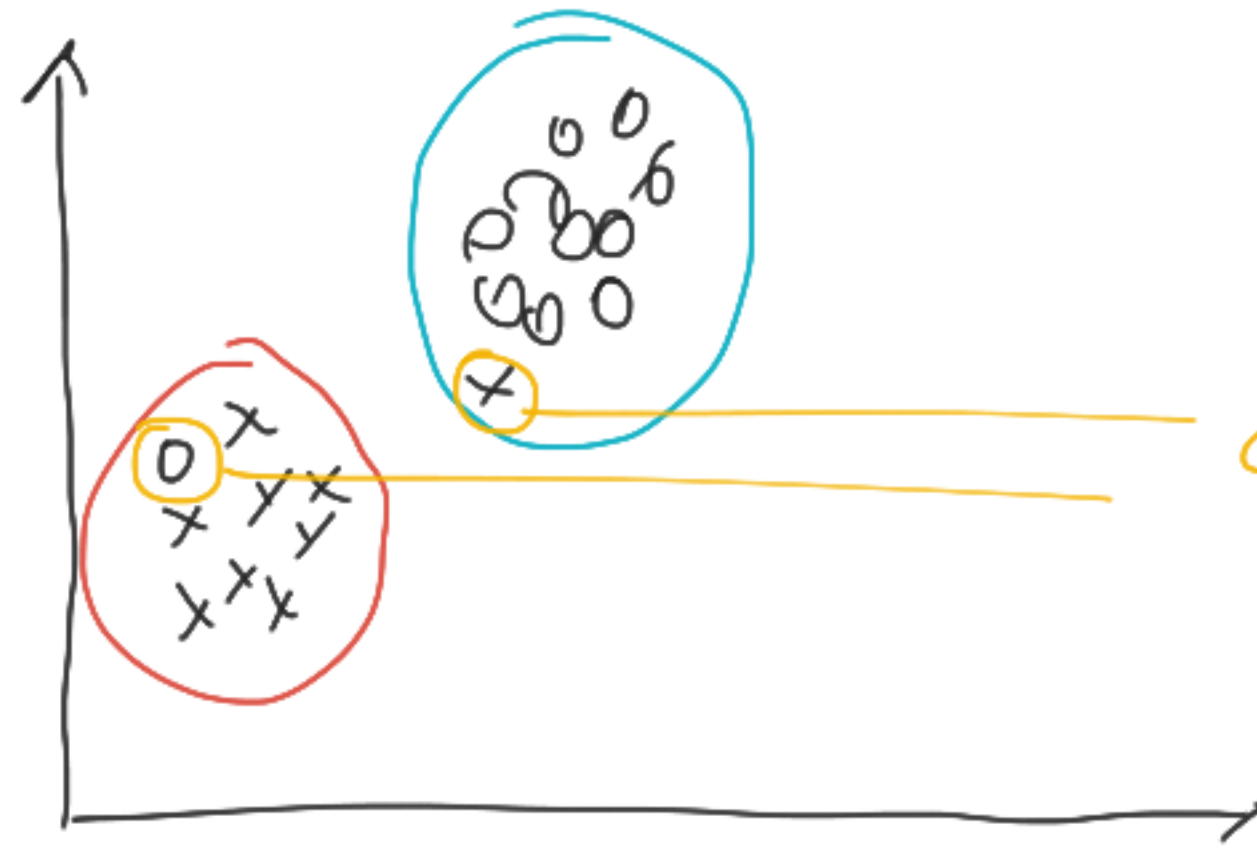


# EVALUACIÓN ALGORITMOS DE AGRUPAMIENTO



errores de agrupamiento

↳ Adjusted Rand Score

[https://scikit-learn.org/stable/modules/generated/sklearn.metrics.adjusted\\_rand\\_score.html](https://scikit-learn.org/stable/modules/generated/sklearn.metrics.adjusted_rand_score.html)

Ejercicio: Generar 1000 observaciones 2D  
en dos clusters utilizando make\_blobs  
↓ sklearn

Representar el Adjusted Rand Index  
a medida que los centroides se acercan



Ex: Titanic

age sex fare ...  
 $x_1$   $x_2$  ...

$x_n$  label

's'

'd'

's' - survived

'd' - dead

Adjusted  
Rand Score

post-hoc

PCA + clustering



# SISTEMAS de RECOMENDACIÓN (RECOMMENDERS)

usuarios      items      valoraciones

$r_{ij}$  : valoración películas      valoración del usuario 1 de la película 2

Movielens  
data

100.000 valoraciones

1000 usuarios

1700 películas

Id usuario \ Id película	1	2	3	4	5	6	7	8
1		3	1		4	3	5	
2	4	1	3		5			2
3	2	1		5				1
4	3		2			5		4

(rango  
1-5)

← valores  
comunes

<https://grouplens.org/datasets/movielens/100k/>

Recomendadores basados en distancia entre usuarios: (KNN)

K - primeros vecinos de un usuario para recomendar

↳ distancia entre usuarios → <sup>valoraciones</sup> comunes  
→ similitud Pearson (correlación) <sup>val</sup><sub>com</sub>  
→ similitud euclídea (a, b):



ex: similitud entre usuarios 2 y 3

usuario 2 →  $a = [4, 1, 2]$

usuario 3 →  $b = [2, 1, 1]$

$d(a, b)$  = distancia euclídea entre las valoraciones comunes

$$S(a, b) = \frac{1}{1 + d(a, b)}$$

$$d(2, 3) = \sqrt{(4-2)^2 + (1-1)^2 + (2-1)^2} = \sqrt{4+1} = \sqrt{5} \rightarrow S(2, 3) = \frac{1}{1+\sqrt{5}} = \underline{\underline{0.31}}$$

Medir distancias a pares entre todos los usuarios

$d_{ij}$

Recomendar una película a usuario  $r \rightarrow$

buscar los  $k$  usuarios más cercanos a  $r$

$\rightarrow$  RECOMENDACIÓN!

LIBERÍA  $\rightarrow$  SURPRISE

<http://surpriselib.com/>

fill with  
mean  
of  
nearest  
sets

1 2 3 . . .

1700

opano 6



$k=2$

$\left\{ \begin{array}{l} 138 \\ 421 \end{array} \right.$

$\{m_1$

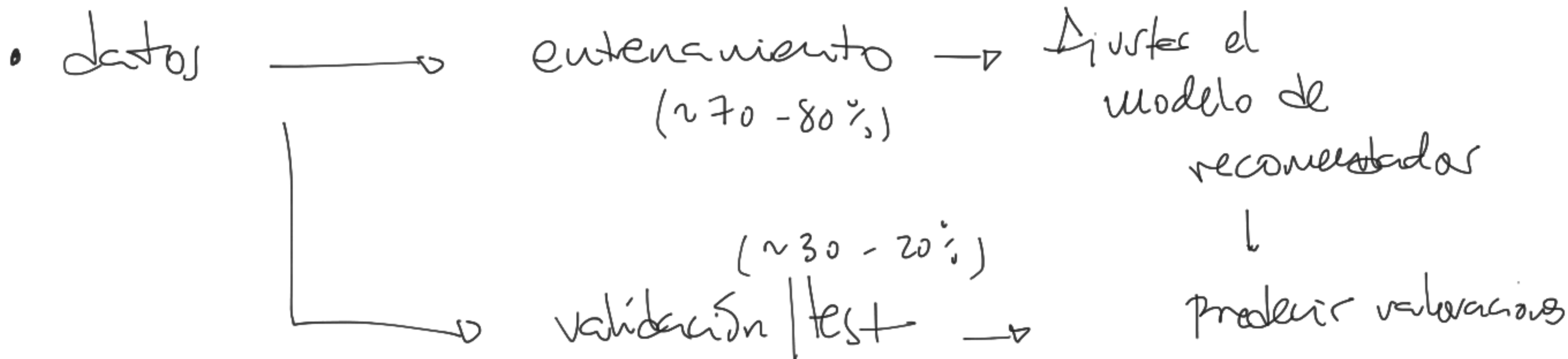
$m_2$

. . .

$m$

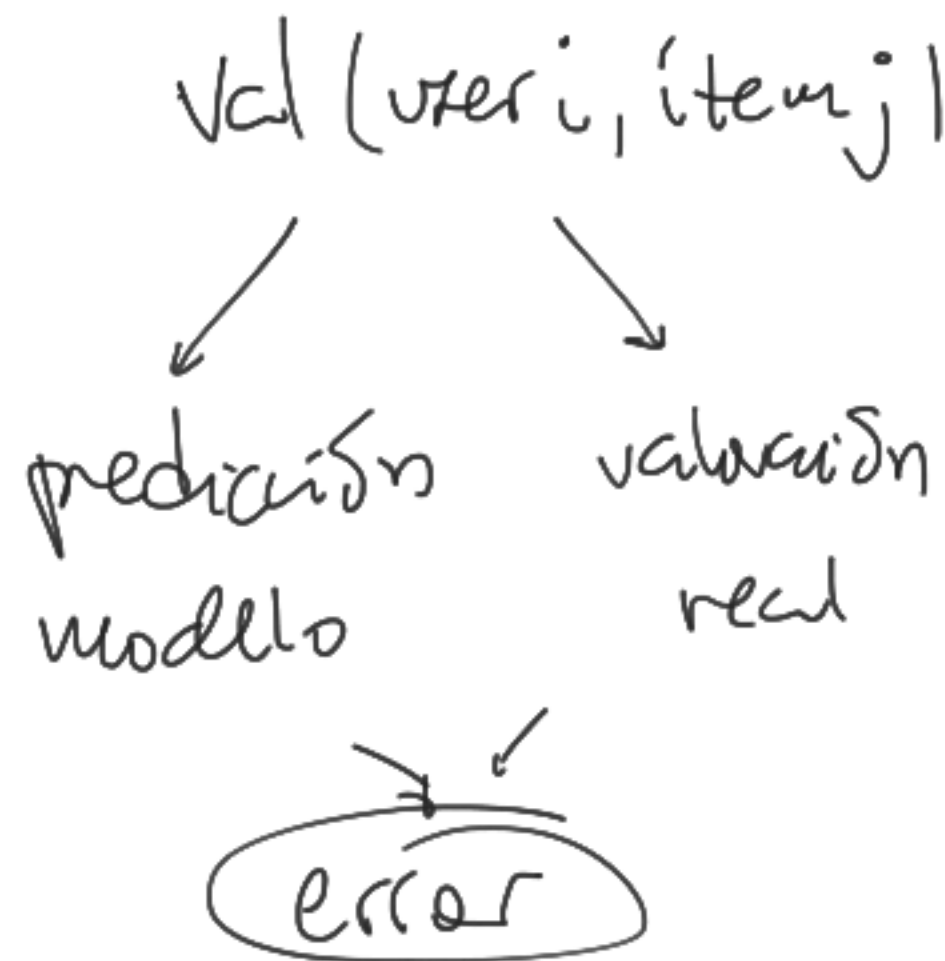
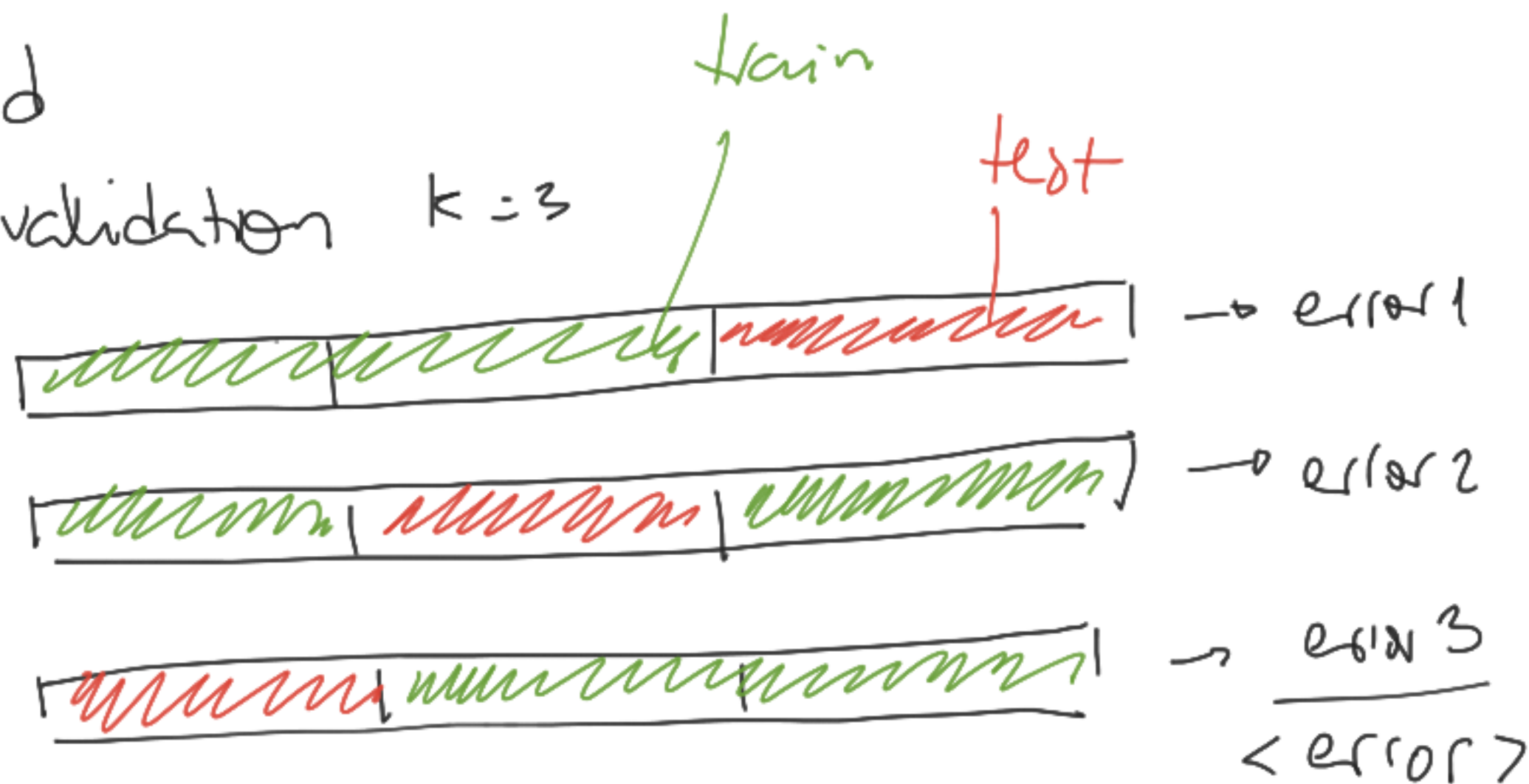
$\{m_{1700}$





• K-fold

Cross-validation  $k=3$





① ELECCIÓN EL PROBLEMA - CONJUNTO DATOS

PROYECTO

② MANIPULACIÓN DATOS, FILTRADO, IMPUTACIÓN Y REPRESENTACIÓN GRÁFICA  
+ CÓDIGO + TEXTO COMENTARIOS

- tamaño
- filtrado - interpretación
- presencia missing data
- tipo de variables
- etiqueta de clase\*
- gráficas (scatter plots, histogramas, pair plots...)

python notebook  
- autoexplicativo  
- ejecutable

③ REDUCCIÓN DIMENSIONALIDAD (PCA)

④ CLUSTERING (k-means, aglomerativo, GMM, GMM+BIC)

⑤ SELECCIÓN ATRIBUTOS\* (RFE, univariada)

⑥ REGRESIÓN UNIVARIADA (LINEAL / ÁRBOLES DECISION)

DATOS



seaborn

sklearn

kaggle

UCI datasets

plotly

<https://plotly.github.io/datasets/>

<https://archive.ics.uci.edu/ml/datasets.php>

built-in  
datasets

## EJERCICIOS - MISCELÁNEA:

- ① Cuántos supervivientes mayores de 40 años hay en el titanic dataset?
- ② ¿Cuál es el adjusted Rand Score de un clustering k-means ( $k=3$ ) de los datos de mpg-cars con respecto a la etiqueta de origen de los coches (ground truth).
- ③ ¿Cuál es la variable que mejor predice la supervivencia en el titanic y la que mejor predice el origen en los datos mpg.
- ④ Cuántas observaciones del iris dataset tienen un petal length fuera del rango intercuartil  $\pm IQR = [Q_1, Q_3]$   
representar el histograma con

