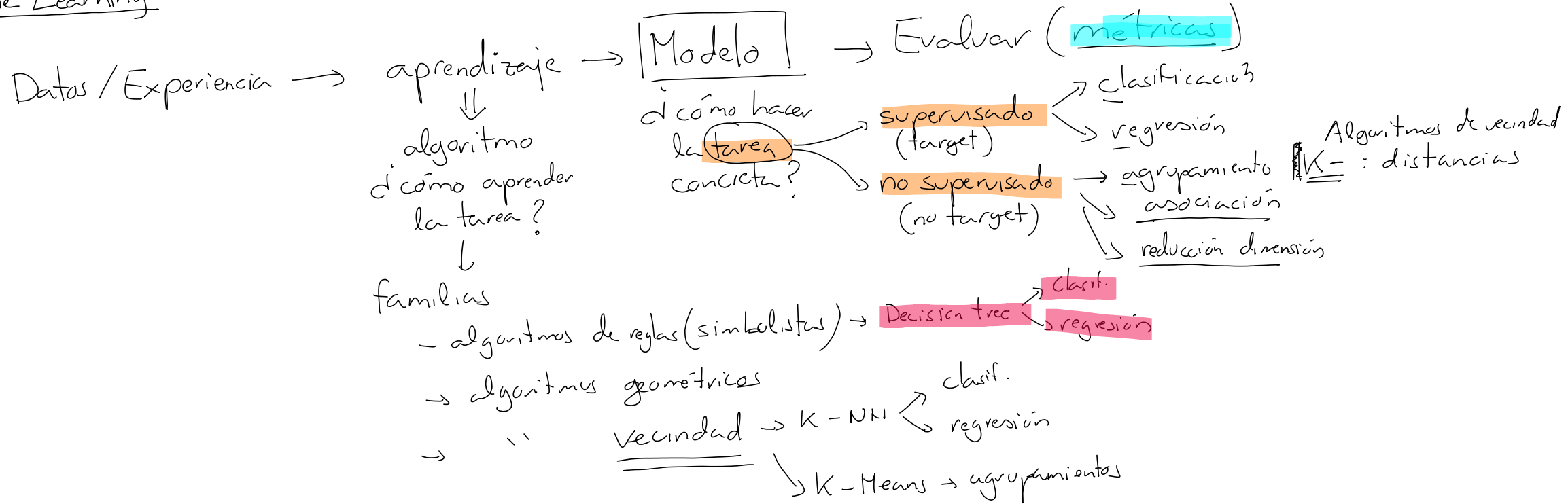
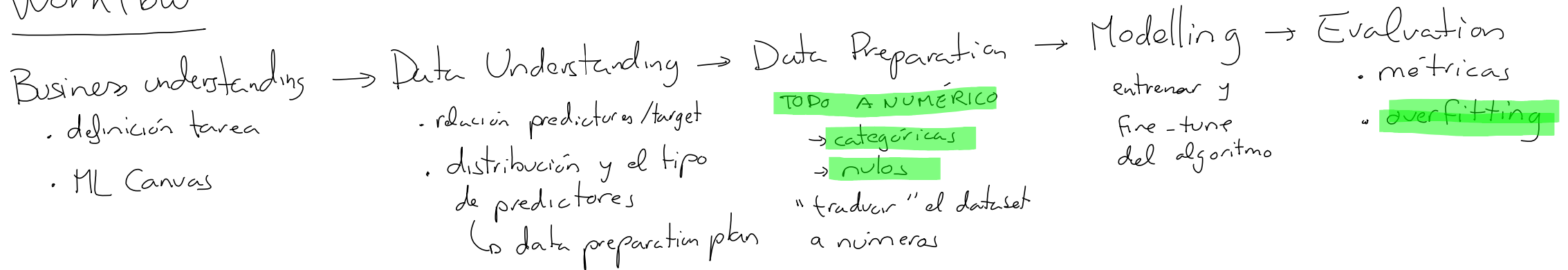


Repaso clase 1

Machine Learning



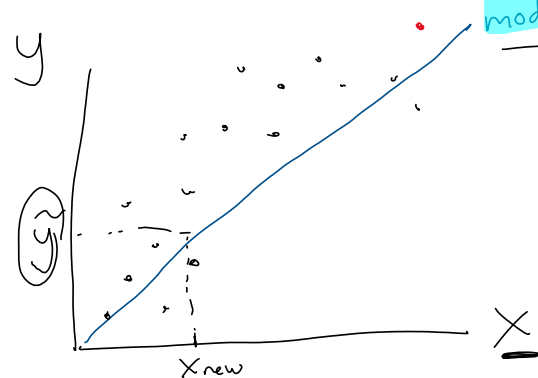
Workflow



Overfitting = memorización \neq Generalización

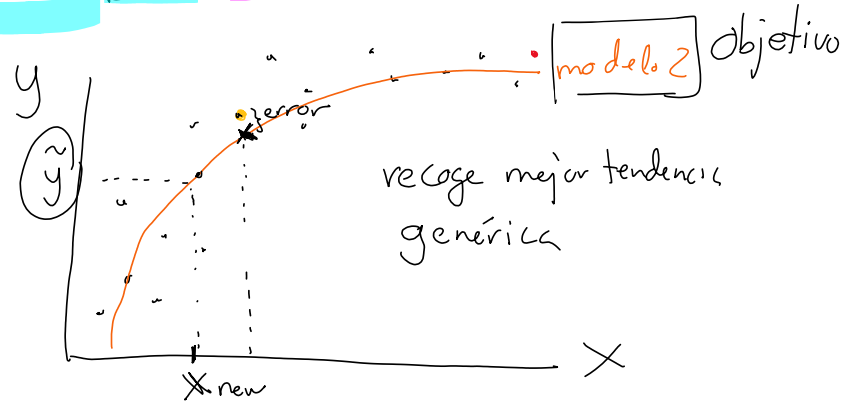
X = predictor • : datos
 y = target

Mismo dataset:



Modelo 1

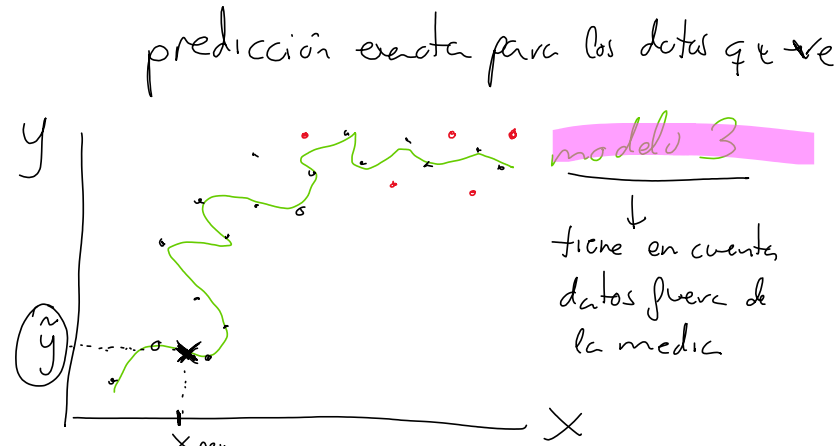
Bias - variance trade-off



Modelo 2 Objetivo

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X^2$$

Se equivoca un poco
 $=$
 se equivoca un poco



Modelo 3

tiene en cuenta
 datos fuera de
 la media

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X^2 + \beta_3 X^3$$

Overfitting
 no se equivoca
 \neq
 sí que se equivoca bastante

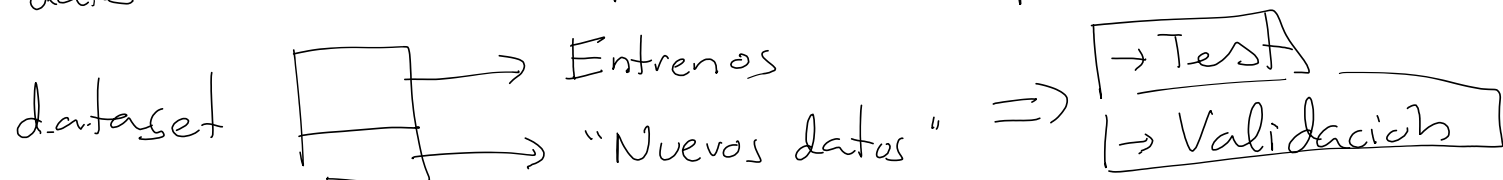
Algoritmo: $Y = \beta_0 + \beta_1 X$
 Underfitting

Métrica: error medio

↳ error datos \rightarrow se equivoca mucho
 ↳ error nuevo dato \rightarrow se equivoca mucho

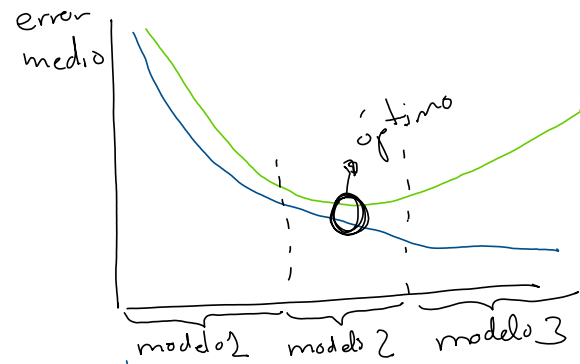
Objetivo: métrica en datos nuevos y antiguos debe ser igual y lo más alta posible

"Nuevos datos": no vanos a esperar al futuro, partimos los datos disponibles



Métrica

error medio \rightarrow objetivo : lo más pequeño posible



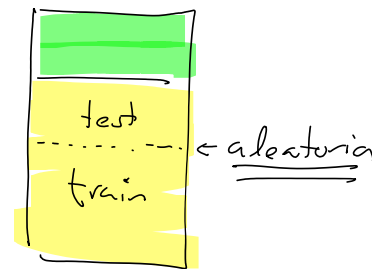
datos "antiguos"

datos "nuevos"

Table resumen

Validation: siempre a parte (2k-3k filas) y se parece a la realidad

development: $\left. \begin{array}{l} \text{test (>2k/3k filas)} \\ \text{train} \end{array} \right\} \text{aleatoria}$



| | | |
|-------------------------|--------|----------------|
| $\parallel \rightarrow$ | filas | |
| | > 60k | Random Holdout |
| | 60-20k | K-fold |
| | < 20k | Bootstrap |

Note: ¿cómo escoger validation?

target

\rightarrow estacional: ciclo (Black Friday, venta helados, ...) \Rightarrow datos ciclo anterior

\rightarrow trend/tendencia: (precio vivienda) \Rightarrow datos más recientes

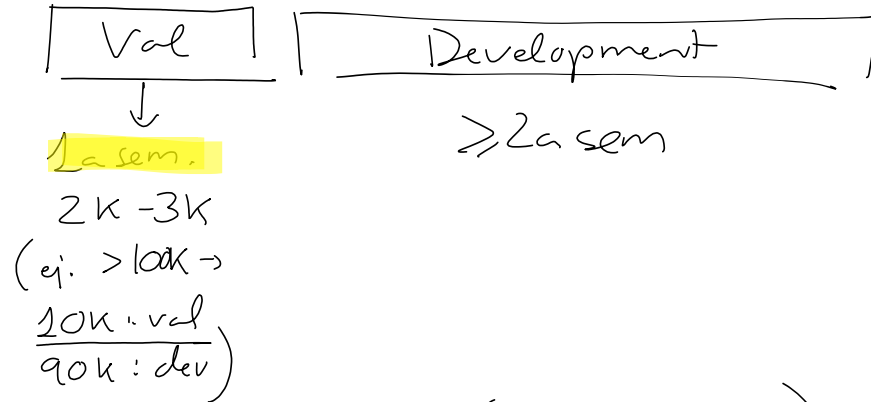
Ejemplo: nuevo lanzamiento producto (iphone 16) ^{new}

Business understanding: modelo que prediga quién es más probable que compre el iphone 16 en la primera semana de lanzamiento, de entre mi base de clientes.

Validation: parecer a realidad

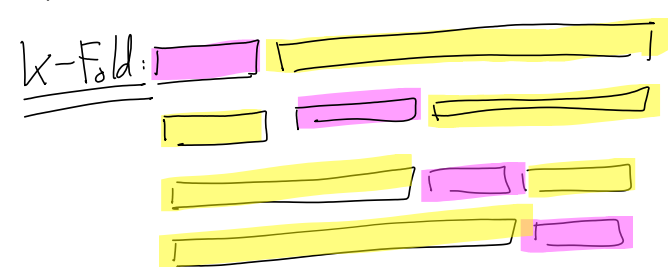
total ventas de iphone 15

Test: misma distribución que train \Rightarrow



¿Partición development?
8K filas
Dev

RH: train test >2k



mismo algoritmo

modelo 1
modelo 2
modelo 3
modelo 4

train₁ = 80% test₁ = 78%
train₂ = 79% test₂ = 78%
train₃ = 82% test₃ = 80%
train₄ = 78% test₄ = 78%

overfitting (memorización)

overfitting: 2pp
1pp
2pp
0pp } avg = 1,25pp

si fila 1 esta en versión 1 test,
no puede estar en otras
versiones de test

¿Cómo escoger modelo?

\rightarrow overfitting y performance medio

Data preparation → todo sea numérico

Nulos

1. No eliminar filas, porque cambia la distribución

si elimino si:

I. Nulo en el target

II. Afecta a una cantidad pequeña $< 0,5\%$
(1%)

III. Fila en concreto tiene mayoría de predictores con nulos ($> 70\%$)

2. Más kxo en eliminar predictores

↳ eliminar atributo/predictor si $> 95\%$ filas con nulos

Tratamientos

I. Según categórica / numérica

II. Según algoritmo

Técnicas:

(categ.) (numérica) } prio 3
- moda / media / mediana

- crear categórica nueva (valor fuera de rango) } prio 2

- valores agrupados / cercanos

- moda / media / mediana subset

- time series: día anterior

- ML: vecindad K-NN

- otras variables } prio 1

- data engineering

funciona bien en algoritmos simbólicos

- categórica: nueva etiqueta

- numérica: valor fuera de rango
(boolean) → -1
1/0

No funciona bien geométricos o vecindad

Transformer categoriales:

1. Ajuste de tipos: Python lee mal el atributo (ej. IDs)

1. Intrínseco → Label Encoding

→ respuesta encuesta

muy mal muy bien
1 5

→ versiones app:

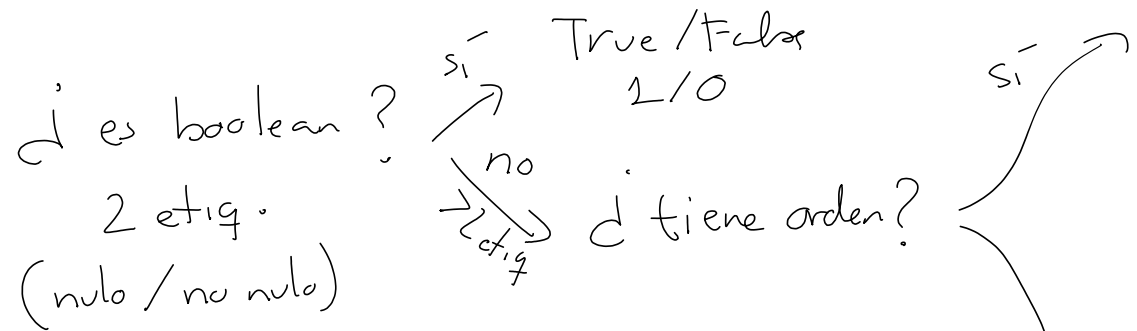
Android: 'Baklava' → 16 1
'Vanilla Ice Cream' → 15 2
'Upside down cake' → 14 3

2. Fuentes externas

weather 01/01/18 → temp.
mm lluvia

ciudad: población
PIB

matrícula: antigüedad



reducir # etiquetas
↳ agrupar etiq. minoritarias
en "Otros"

Frequency encoding

Conteo del número de
veces que aparece la
etiqueta en el dataset

| ID | Ciudad | ID | FE |
|----|----------|----|----|
| 1 | Madrid | 1 | 3 |
| 2 | Valencia | 2 | 2 |
| 3 | Valencia | 3 | 2 |
| 4 | Madrid | 4 | 3 |
| 5 | Madrid | 5 | 3 |

