

28 oct 2025

1-) Construction d'un modele: Selection de A56
predicteurs.

2-) Dev Models dans un contexte A57
Big Data: Spark ML.

3-) Deployment de Models A61
Concepts & Outils

TP: * Finalisation

* PPT

axis 1

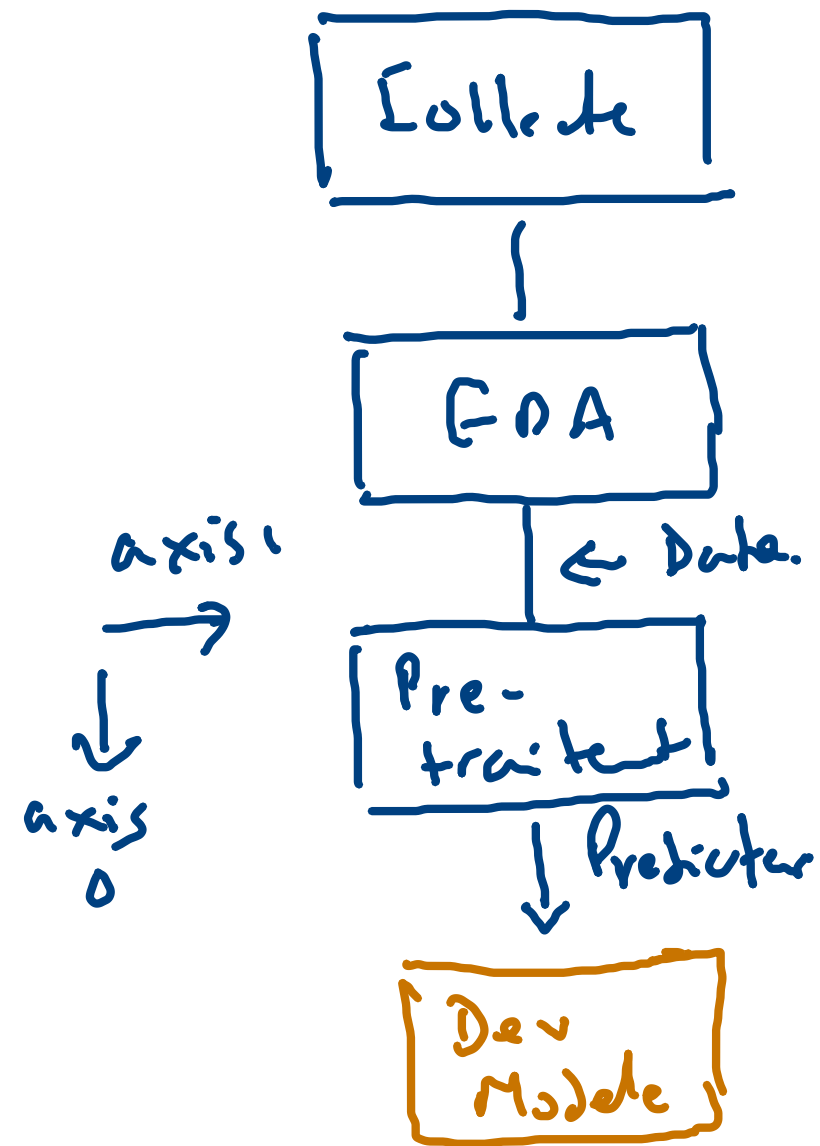
Descripteurs $\xrightarrow[\text{Selection}]{\text{Extraction}}$ Predictors

P.C.A : Reduction de Dimensions

$\Sigma P \rightarrow P$
 \downarrow
 Dev Model.

Tradeoff $\overline{\text{Biais-variance}}$
 du model

objectif : Prendre des predicteurs physiques
 \Rightarrow offrir un nombre plus réduit de
 predicteurs



All-in

~~10~~ ~~non~~

points

Taille

N
Inter

QI

y
(note)

All-in

? Question:

Poids
Taille

n'ont pas d'effet sur le y
(note)

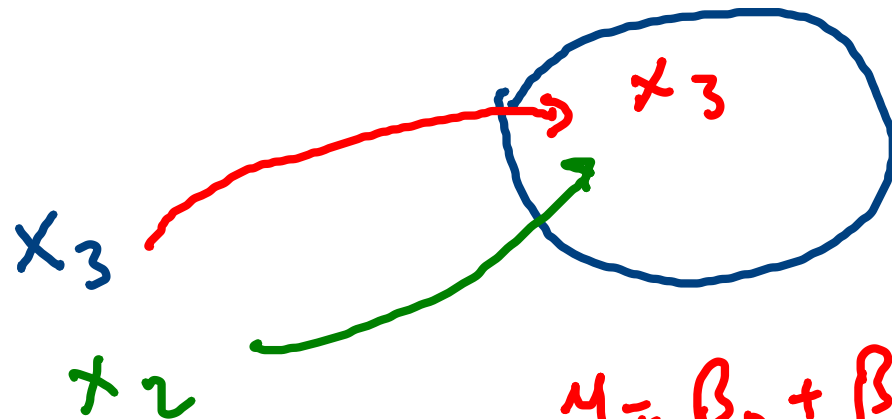
- Correlation.
- Expertise du domaine

Performance du modèle: À voir.

Stepwise Forward

x_1 x_2 x_3 x_4

Model. R.L



$$y = \beta_0 + \beta_1 \underline{x_3}$$

(?) Qualité < BL

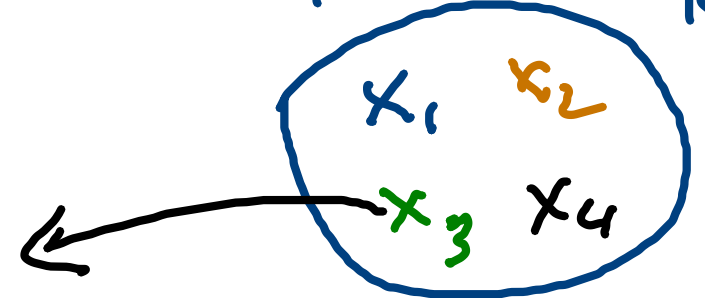
$$y = \beta_0 + \beta_1 x_3 + \beta_2 x_2$$

Qualité < BL
>

Stepwise Backward

x_1 x_2 x_3 x_4

Model. R.L



$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \beta_4 x_4$$

Qualité < BL

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_4 x_4$$

Qualité > BL

Backward x_1, x_2, x_3, x_4
p-value. 0.01 0.015 0.07 0.03

$$\alpha = 0.05$$

it2 Model $\rightarrow x_1, x_2, x_4$
 p-value 0.013 0.053 0.02

Model $x_1, x_4 \Rightarrow$ Model
 0.027 0.032 $\hookrightarrow x_1, x_4.$

E-D.A

Dist $x_1 | x_2 | \dots | x_n | y$

Standard Scale
 Robust Scale
 Min-Max Scale

Test d'Hypothèse.

H_0 : X_i n'a pas d'effet sur y Absence.

H_1 : X_i a un effet sur y Non Absence.

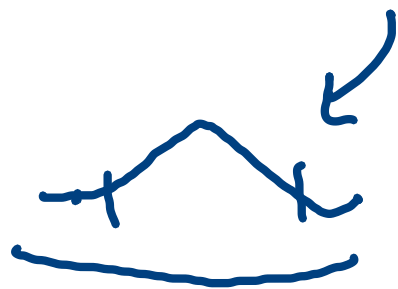
Calculer une statistique \rightarrow Dépend de la dist X_i
 $\alpha = 0.05$

Normal
centre



Z-test

t-test



p-value: Mesure de prob qui indique si on est dans les conditions de H_1 .

$p\text{-value (calculée)} < \underline{0.05 (\alpha)}$
 $\Rightarrow H_1$

Forward.

iteration 1.

x_1 — MOD 1

p-value 0.03

x_2 — MOD 2

p-value 0.04

x_3 — MOD 3

p-value 0.07

x_4 — MOD 4

p-value 0.01

min

MOD 1 T 1

x_4

$x_4 x_1$

0.02

$x_4 x_2$

0.06 \otimes

✓ $x_4 x_3$??

0.01

MOD 1 T 2

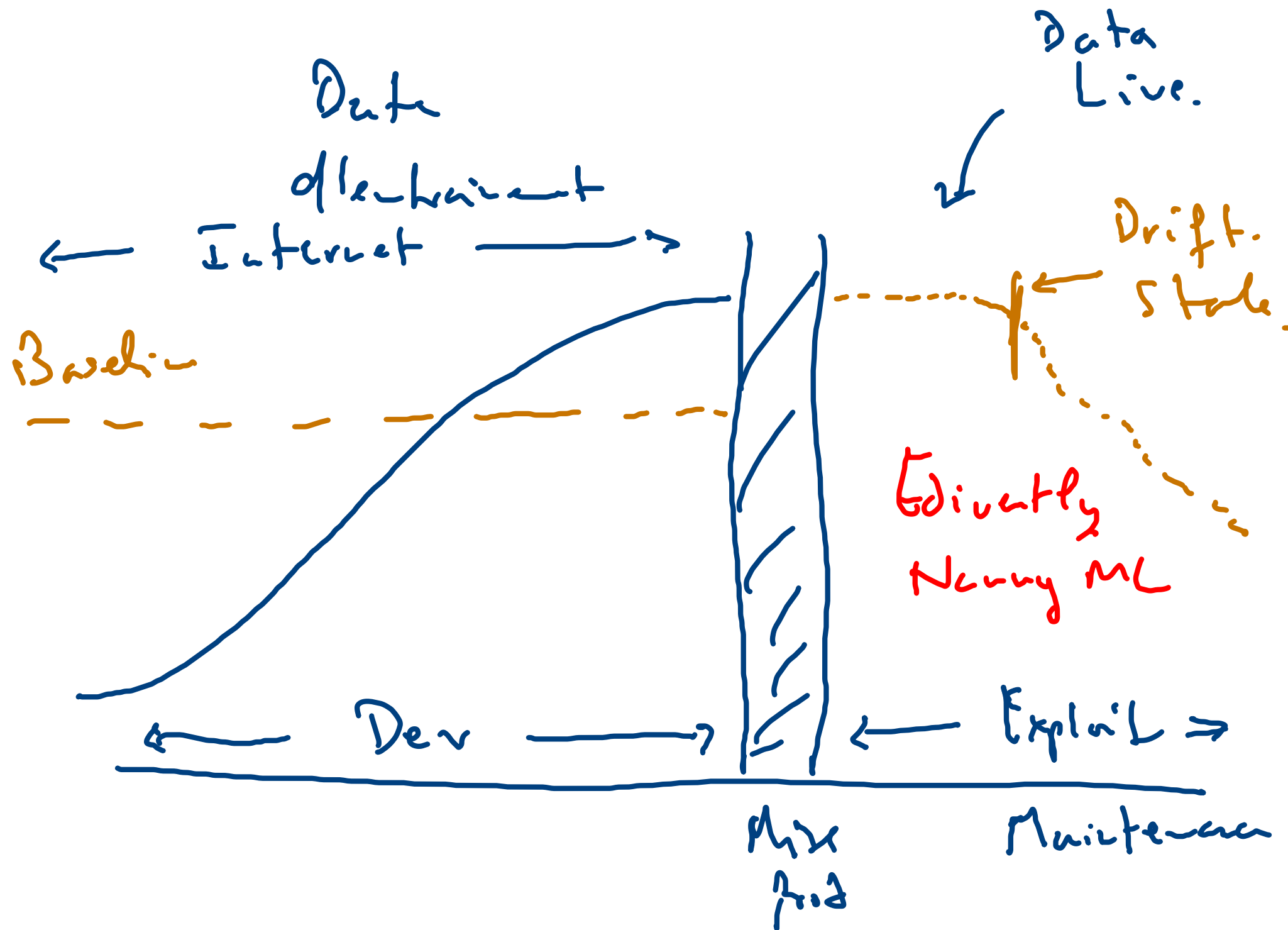
$x_4 x_3$

FIN

$x_4 x_3 x_1$
 $x_4 x_3 x_2$

0.055 $\geq \alpha$

0.07 $\geq \alpha$



	x_1	x_2	x_N	y	<u>train-test</u>
test	-	-	-	-	y_{test} connu.

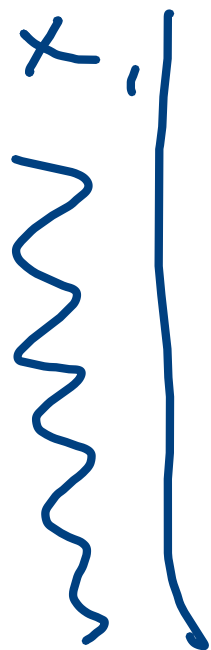
$y_{\text{pres}} \rightarrow \hat{a}$ partir du modèle

Calcul $\leftarrow \text{Met}(y_{\text{pres}}, y_{\text{test}})$

<u>Prod</u>	x_1	x_2	x_N	y_{pres}	y_{pres} connu
	-	-	-		$y_{\text{real}} ??$

Approach Calcul Métrique

train





D.P_{train}

Prod



← Outfit → D.P_{prod}

Problems

x_1	x_2	...	x_n	y
				

axis 0

E.D.A 1) Analyse
statistique
Descriptive

moy, median.
ecart type, etendue
min, max

3) Apports de
correction.

2) Dist Prob $\{x_i\}$



4) Tests d'hypothese

H_0 : x_i pas d'effet sur y
 H_1 : x_i a effet sur y

\$\$

Statistique:

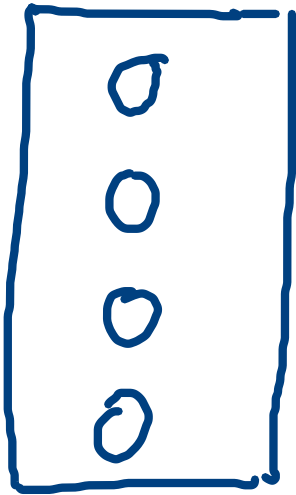
Test paramétrique — Z-test
— t-test

Test non paramétrique



Spark

ma_liste

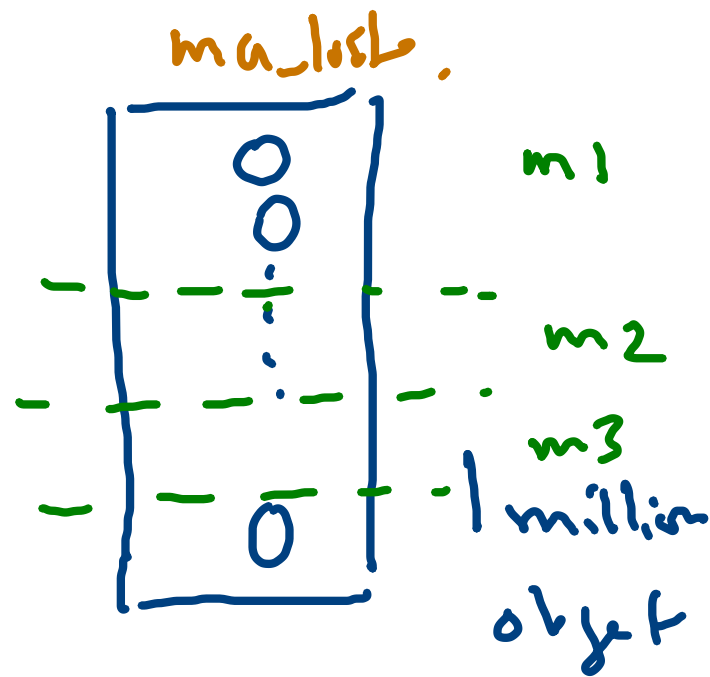


for tmp in ma_liste :

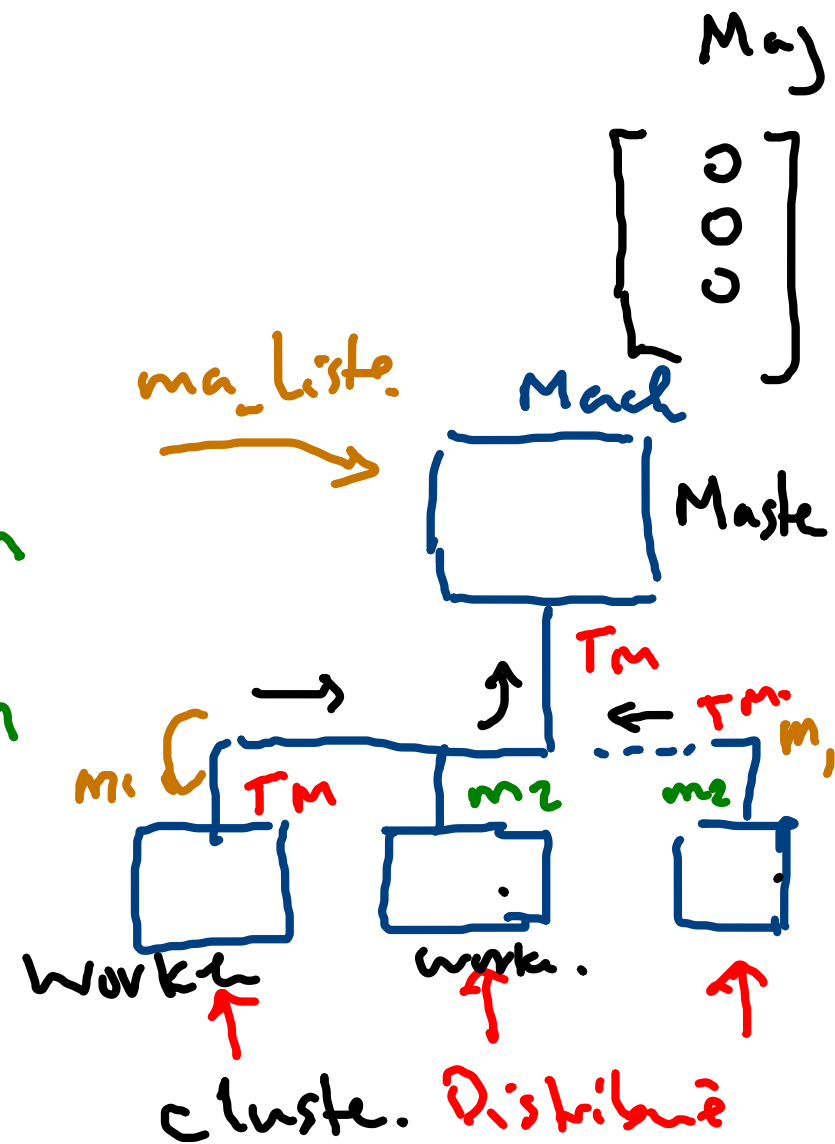
tmp.nom = tmp.nom.upper() } T_M

non
salvée.

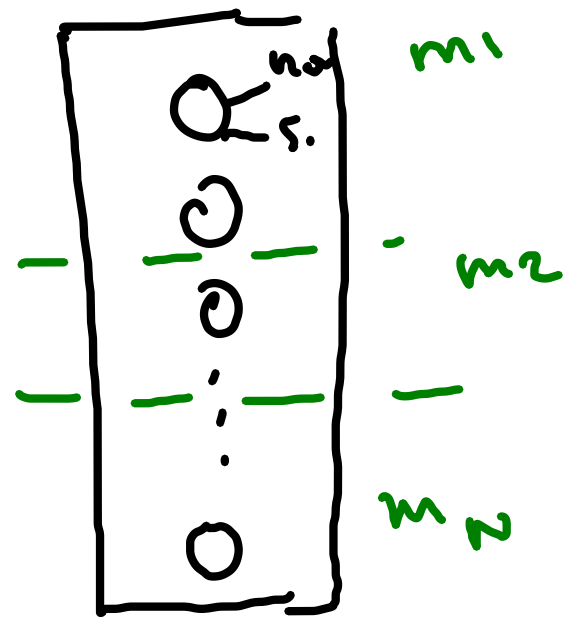
Driver heartbeat



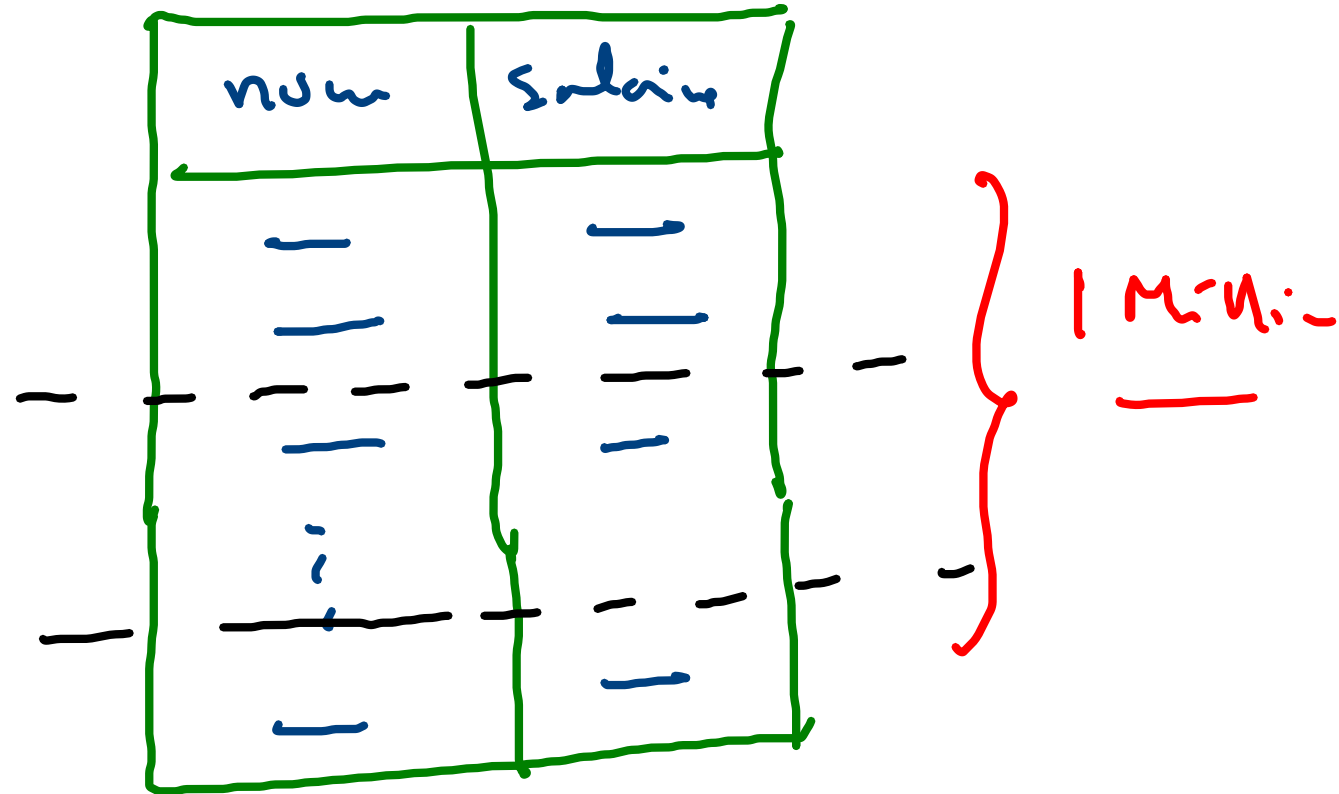
$[0 \ 0 \ \dots \ 0]$ $\leftarrow T_M$
 $[0 \ \dots \ 0]$ $\leftarrow T_M$
 $[0 \ \dots \ 0]$ $\leftarrow T_M$



Framework Spark



liste
RDD



DataFrame - Spark

↳ Representation
ML

SparkML

Scikit-learn
↓
Dataframe.
RL
KNN
⋮

RL

for loss

$$MSE = \frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i)^2$$

$$\Rightarrow \hat{\beta}_i$$

$$= \frac{1}{n} \sum_{i=1}^n () + \sum_{i=1001}^{\infty} + \dots$$

