

[Tu x a e]



[Ateliers Data] / [3. Régression. Prédire le prix des maisons ?]

Introduction



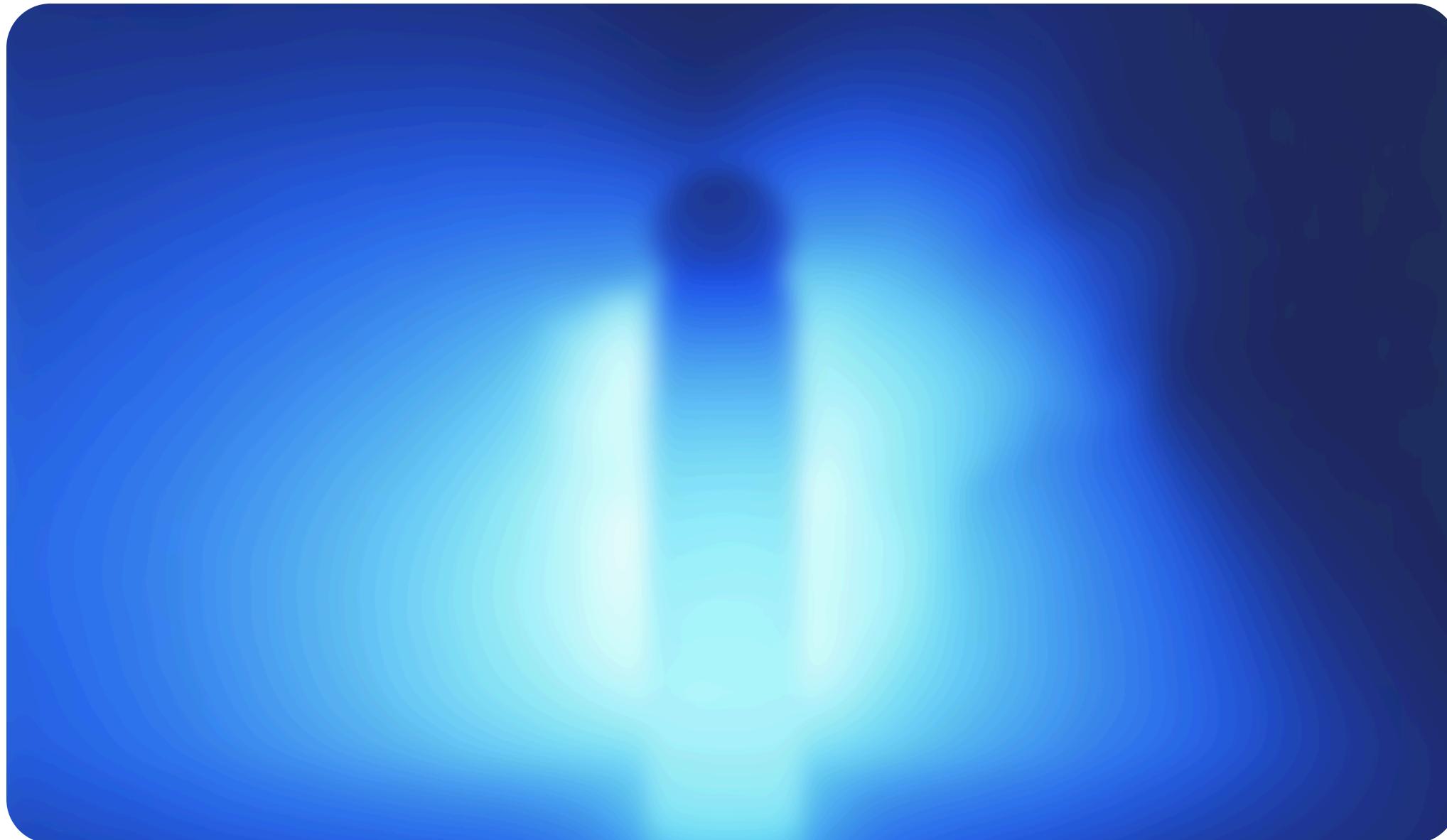
Jupyter-python

The JupyterLab IDE with Python, Julia, and a collection of standard data science packages.

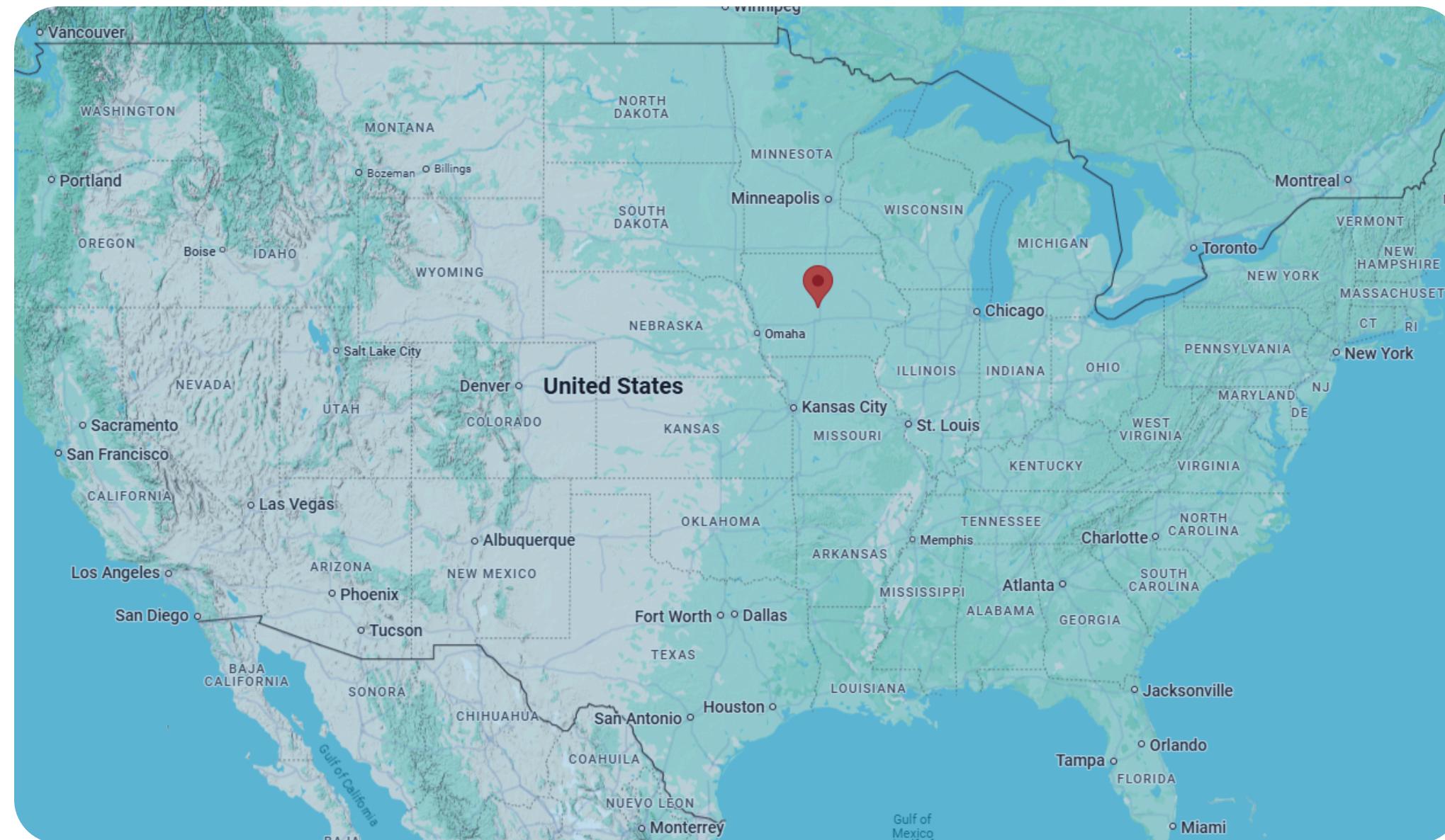
[Learn more](#) [Launch](#)

<https://shorturl.at/LbZDW>

Le problème



Le problème



Ames, Iowa

Le problème



Ames, Iowa

Le problème



(la dernière fois)
Classification

Le problème



Régression

Le dataset



Le dataset

Boston Housing Dataset
(Harrison and Rubinfeld, 1978)

506 observations, 14 variables

Le dataset



Dean De Cock, 2011
(<https://jse.amstat.org/v19n3/decock.pdf>)

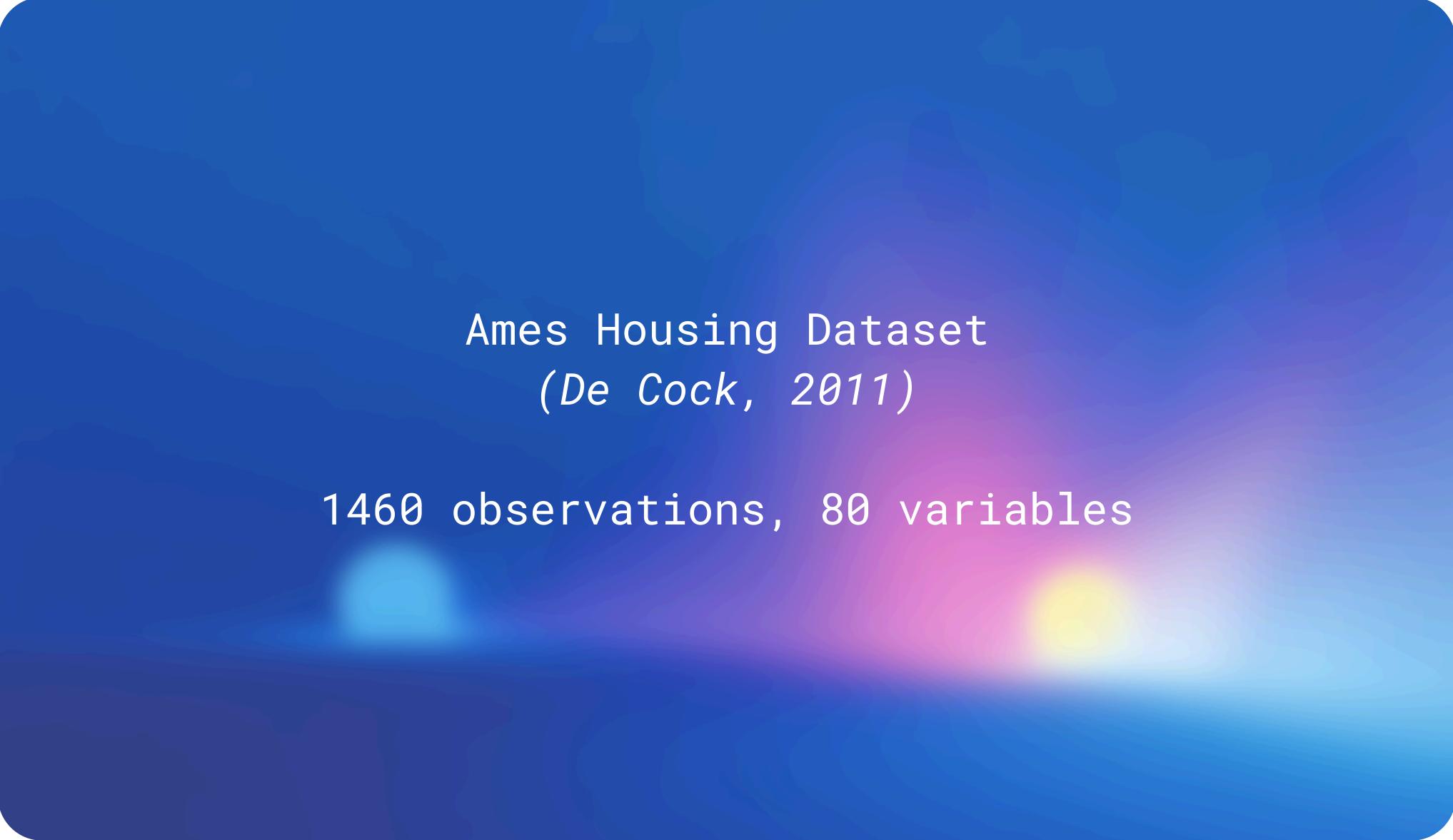
Le dataset

Dean De Cock, 2011

(<https://jse.amstat.org/v19n3/decock.pdf>)

3970 observations, 113 variables

Le dataset



Ames Housing Dataset
(De Cock, 2011)

1460 observations, 80 variables

Le dataset

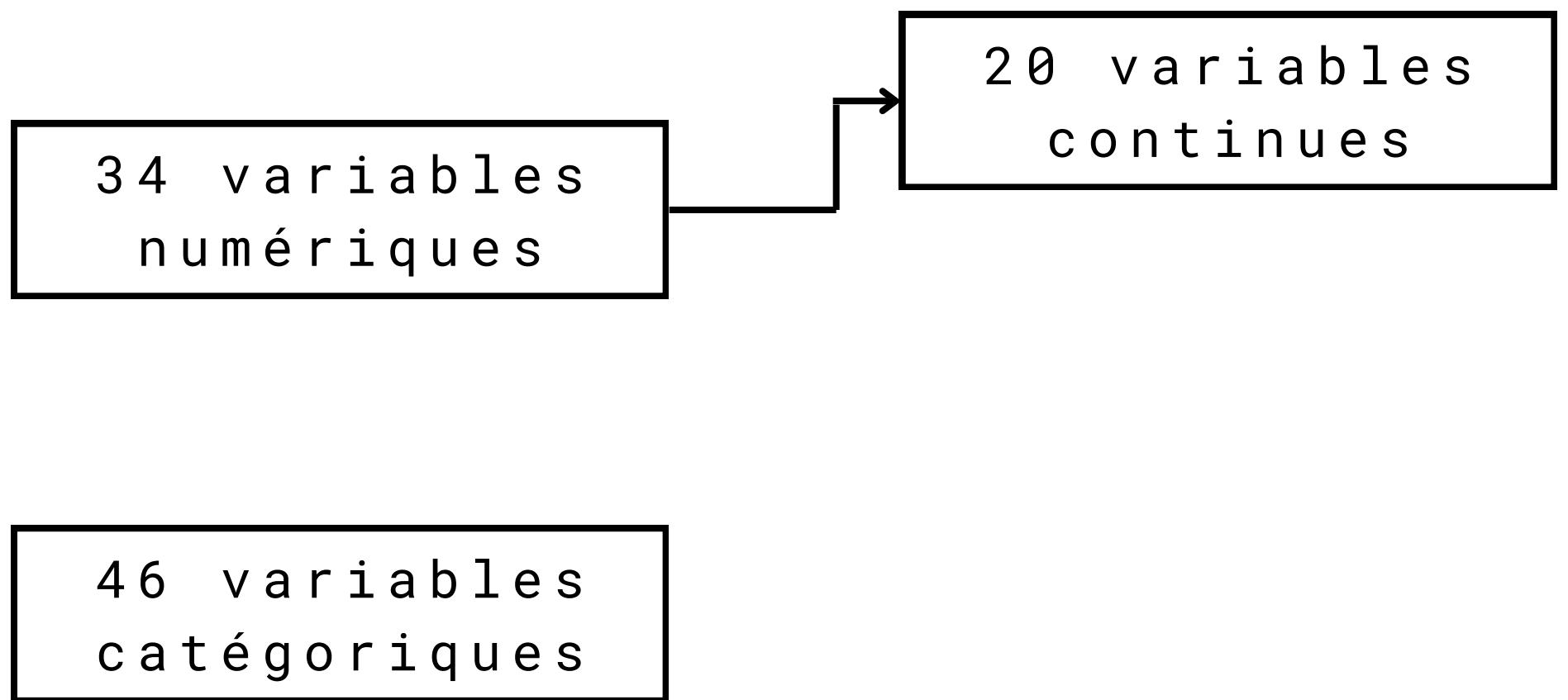
34 variables
numériques

Le dataset

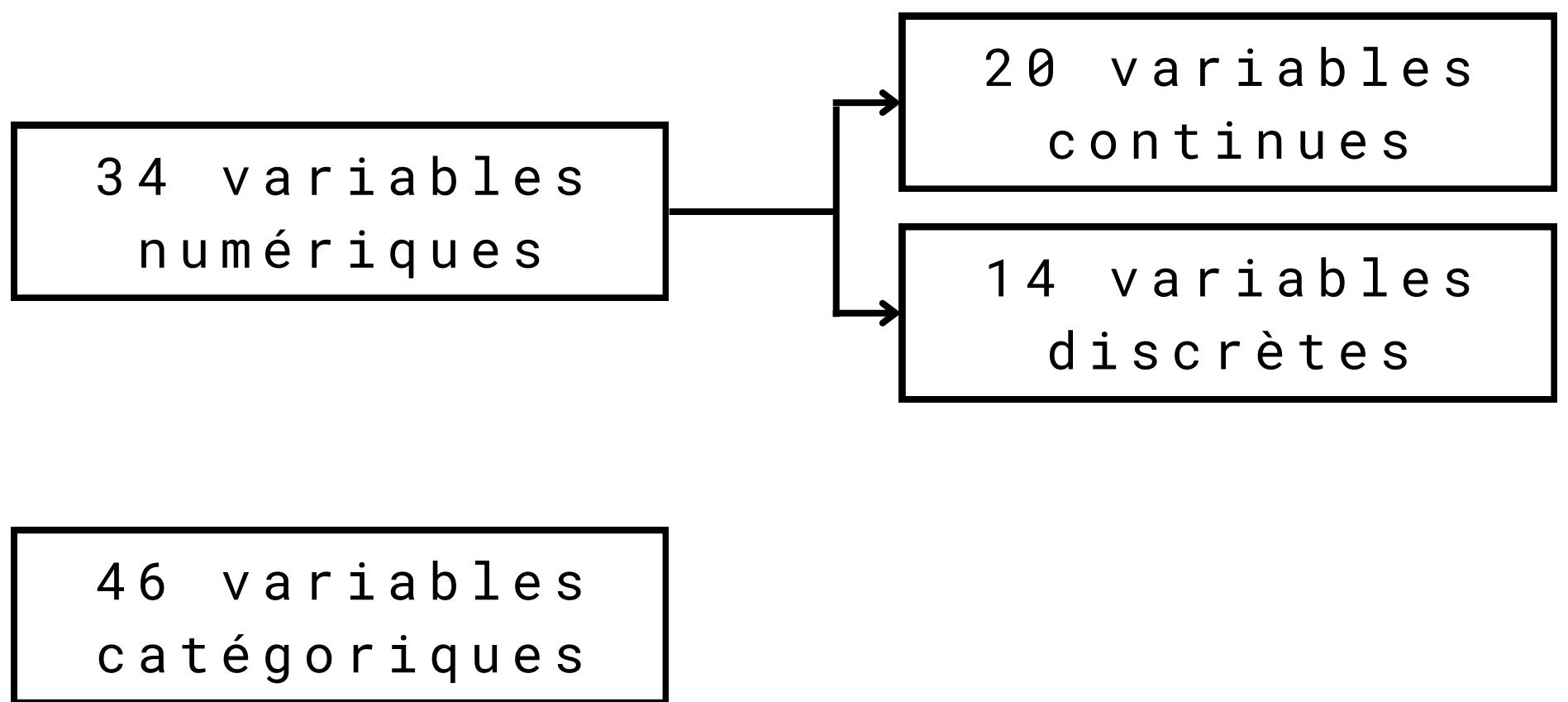
34 variables
numériques

46 variables
catégoriques

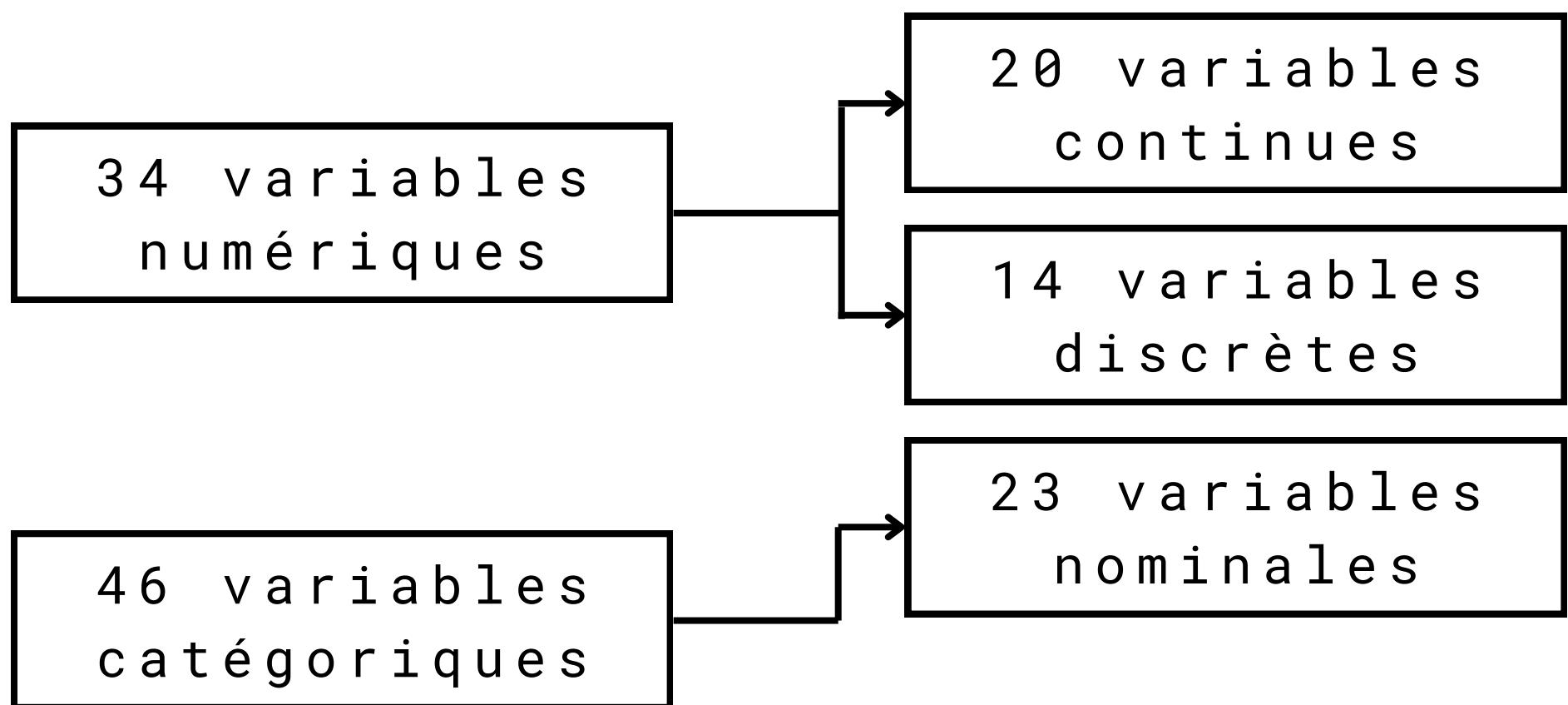
Le dataset



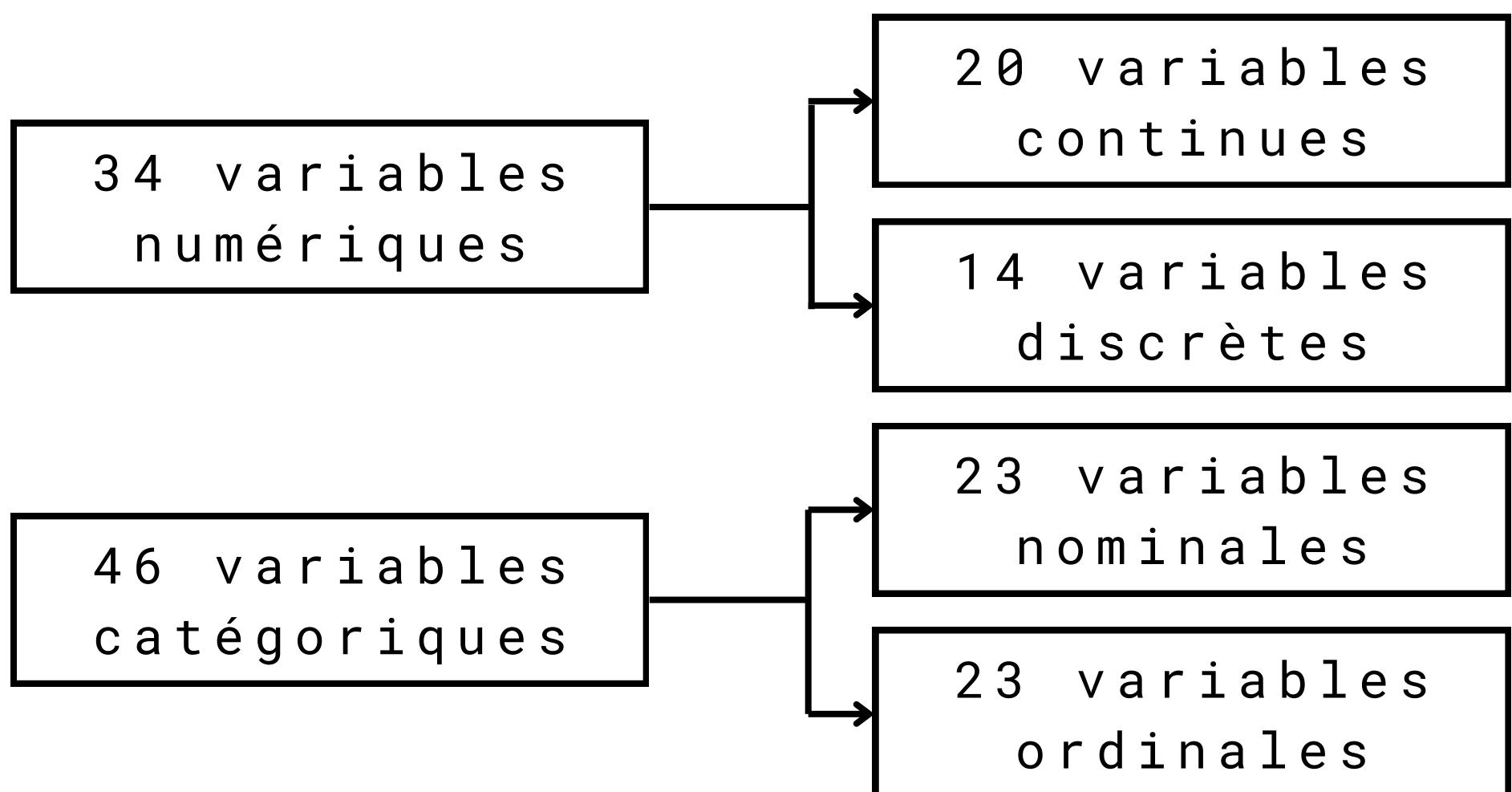
Le dataset



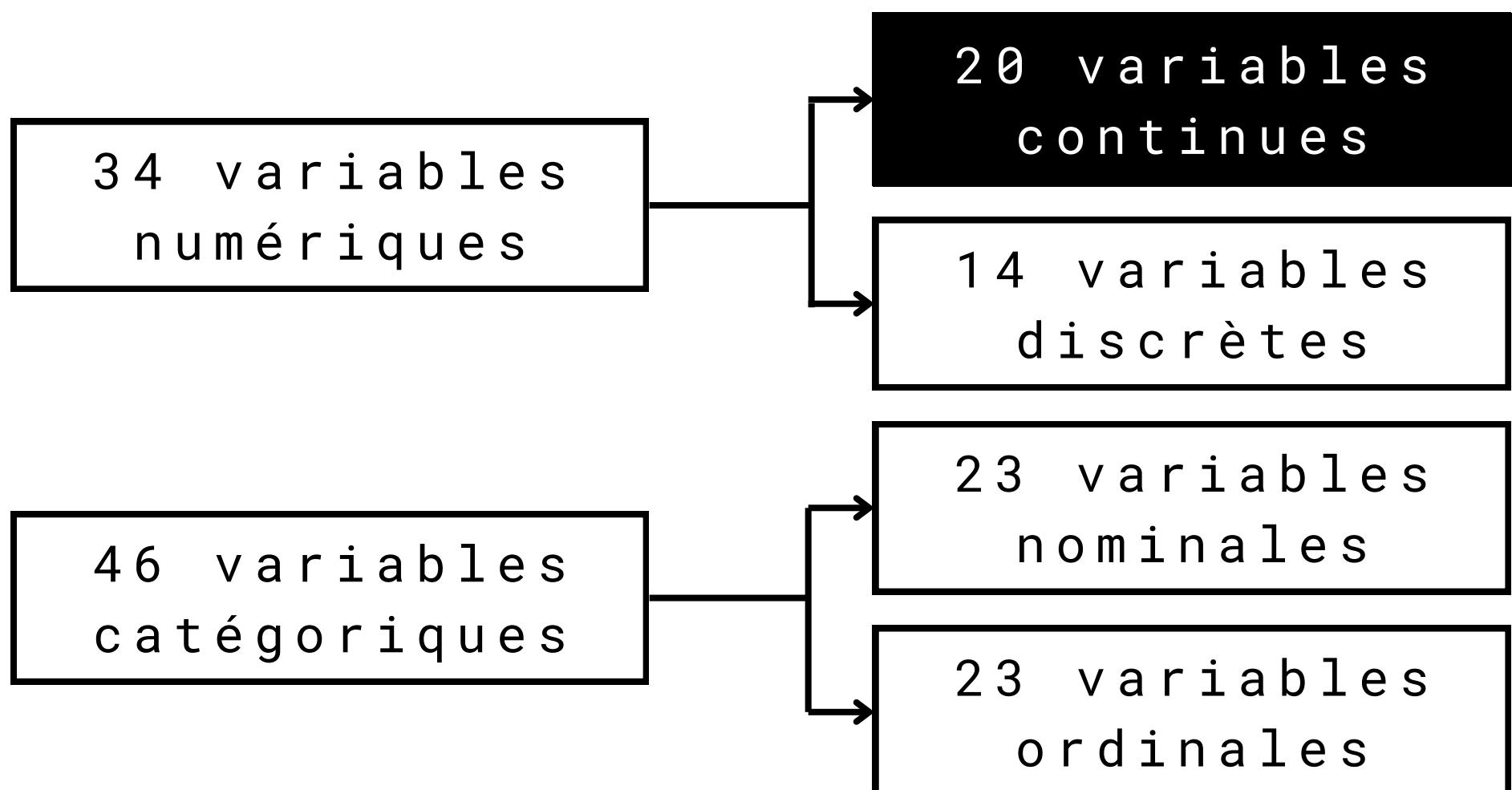
Le dataset



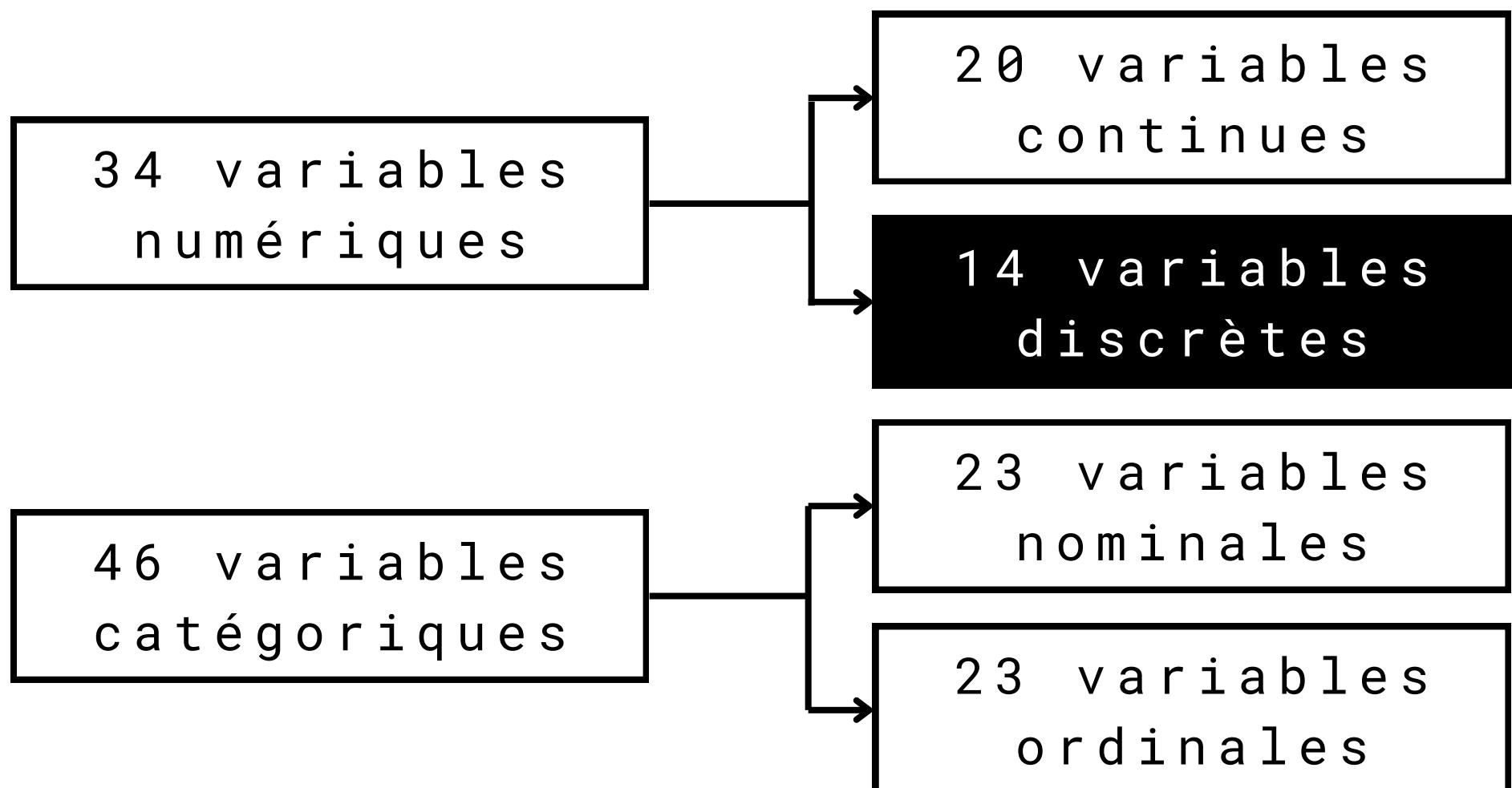
Le dataset



Le dataset



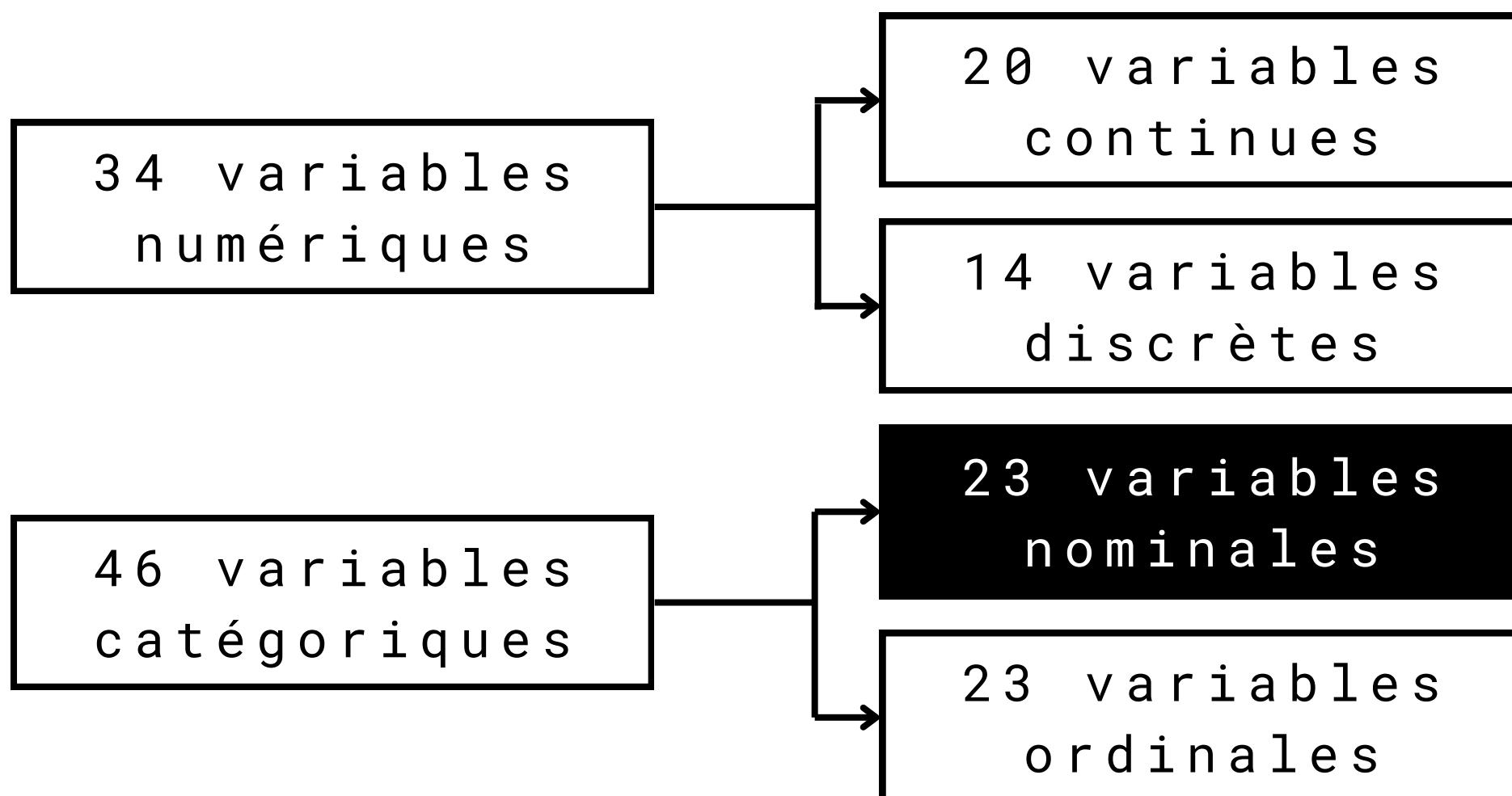
Le dataset



GarageCars

Capacité du garage en nombre de voitures.

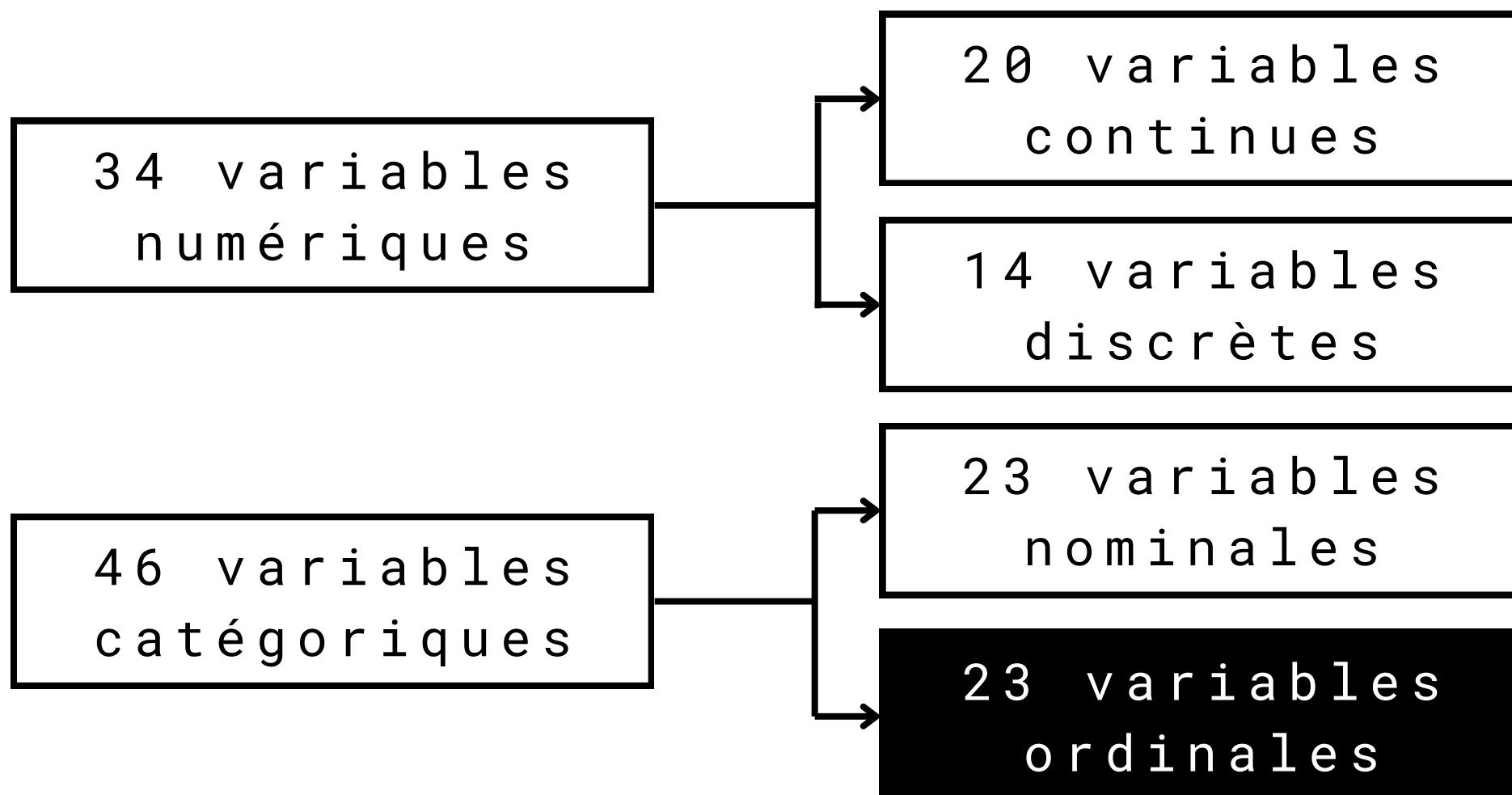
Le dataset



RoofMat1

- *ClyTile* : Argile ou tuile
- *CompShg* : Bardeaux standards (composites)
- *Membran* : Membrane
- *Metal* : Métal
- *Roll* : Rouleau
- *Tar&Grv* : Gravier et goudron
- *WdShake* : Bardeaux de bois (secoués)
- *WdShngl* : Bardeaux de bois (classiques)

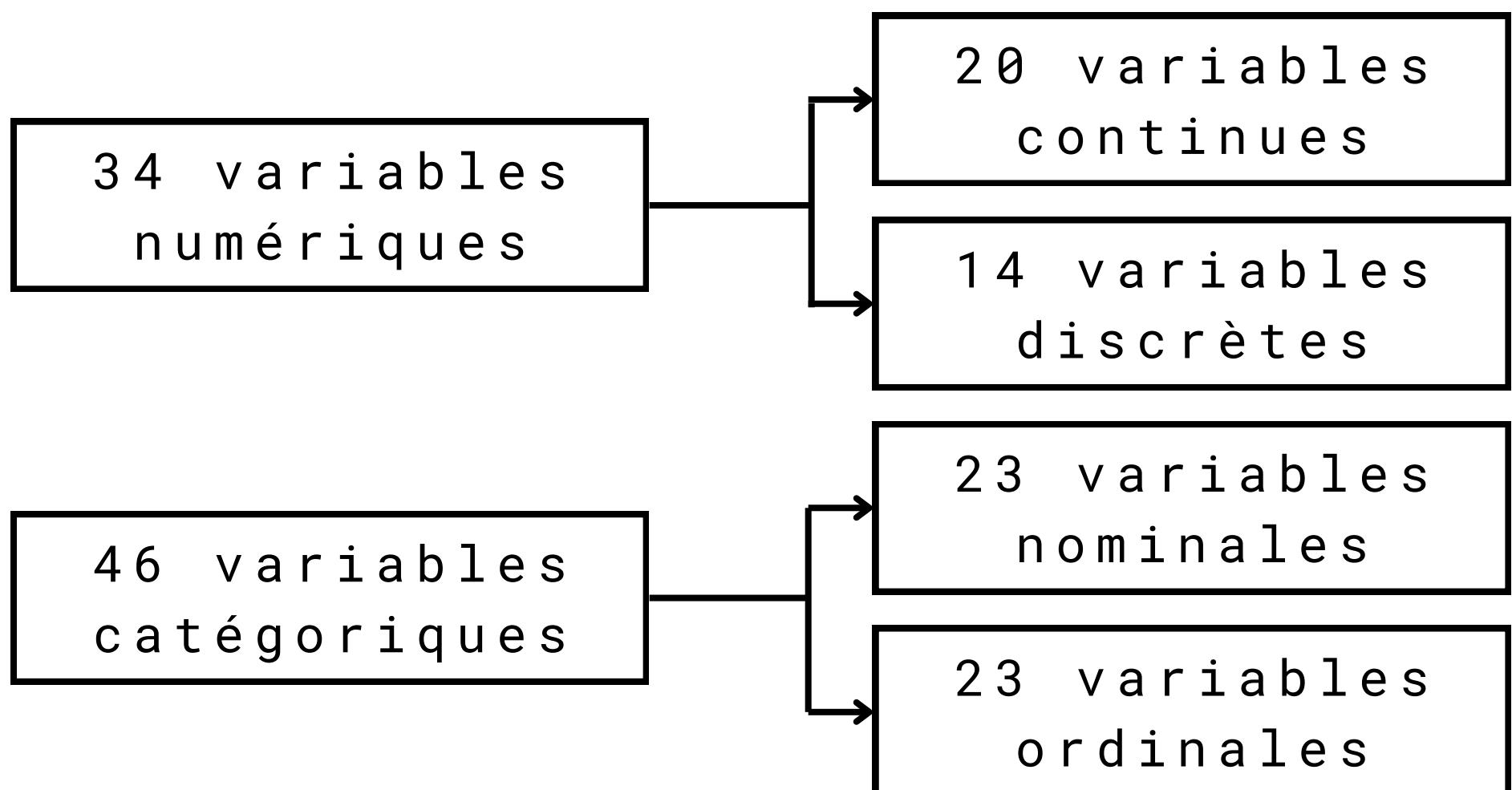
Le dataset



GarageQual

- *Ex* : Excellente
- *Gd* : Bonne
- *TA* : Typique/Moyenne
- *Fa* : Passable
- *Po* : Mauvaise
- *NA* : Pas de garage (non applicable)

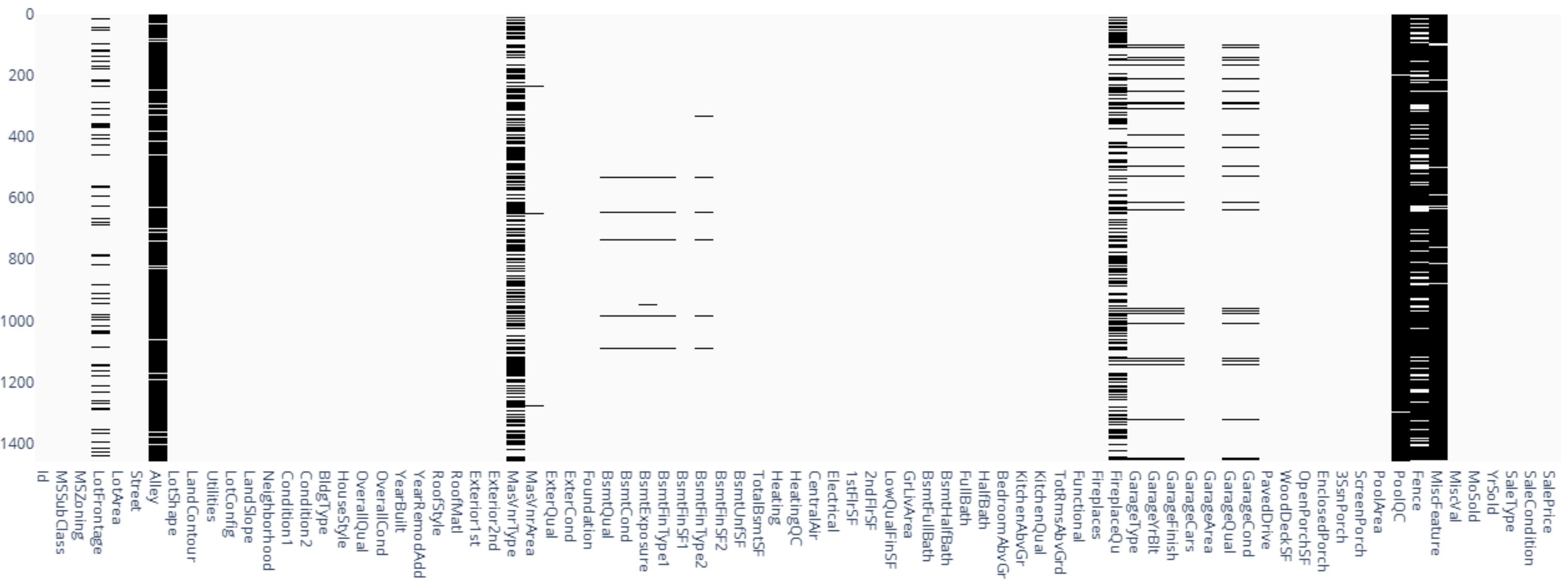
Le dataset



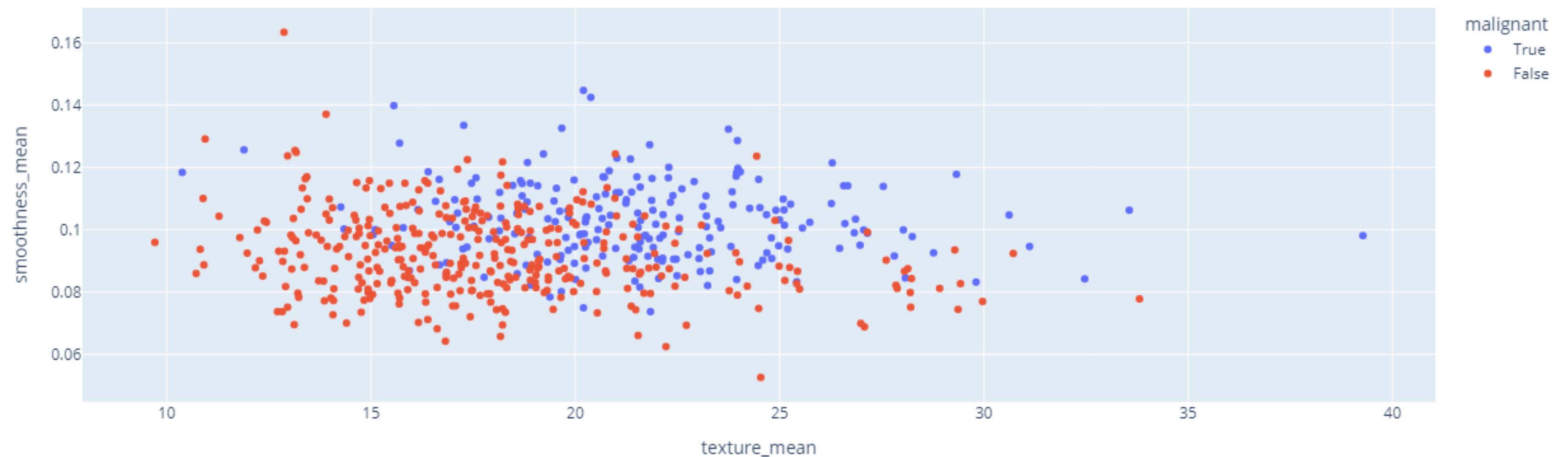
> > > > > > > >

Au travail ! ! !

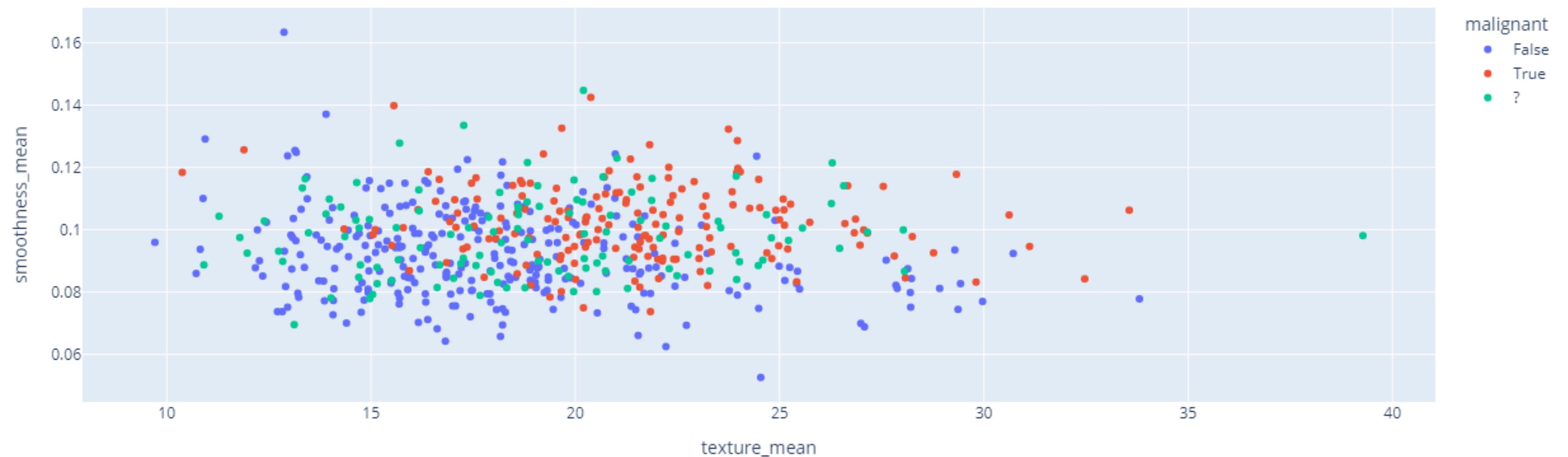
Analyse des données



Rappel du dernier atelier



Rappel du dernier atelier



Métriques



(la dernière fois)
Classification binaire

Métriques



(la dernière fois)
Classification binaire

Exactitude
(Accuracy)

Métriques



(la dernière fois)
Classification binaire

Exactitude
(Accuracy)

Précision
(Precision)

Métriques

(la dernière fois)
Classification binaire

Exactitude
(Accuracy)

Précision
(Precision)

Rappel
(Recall)

Métriques

(la dernière fois)
Classification binaire

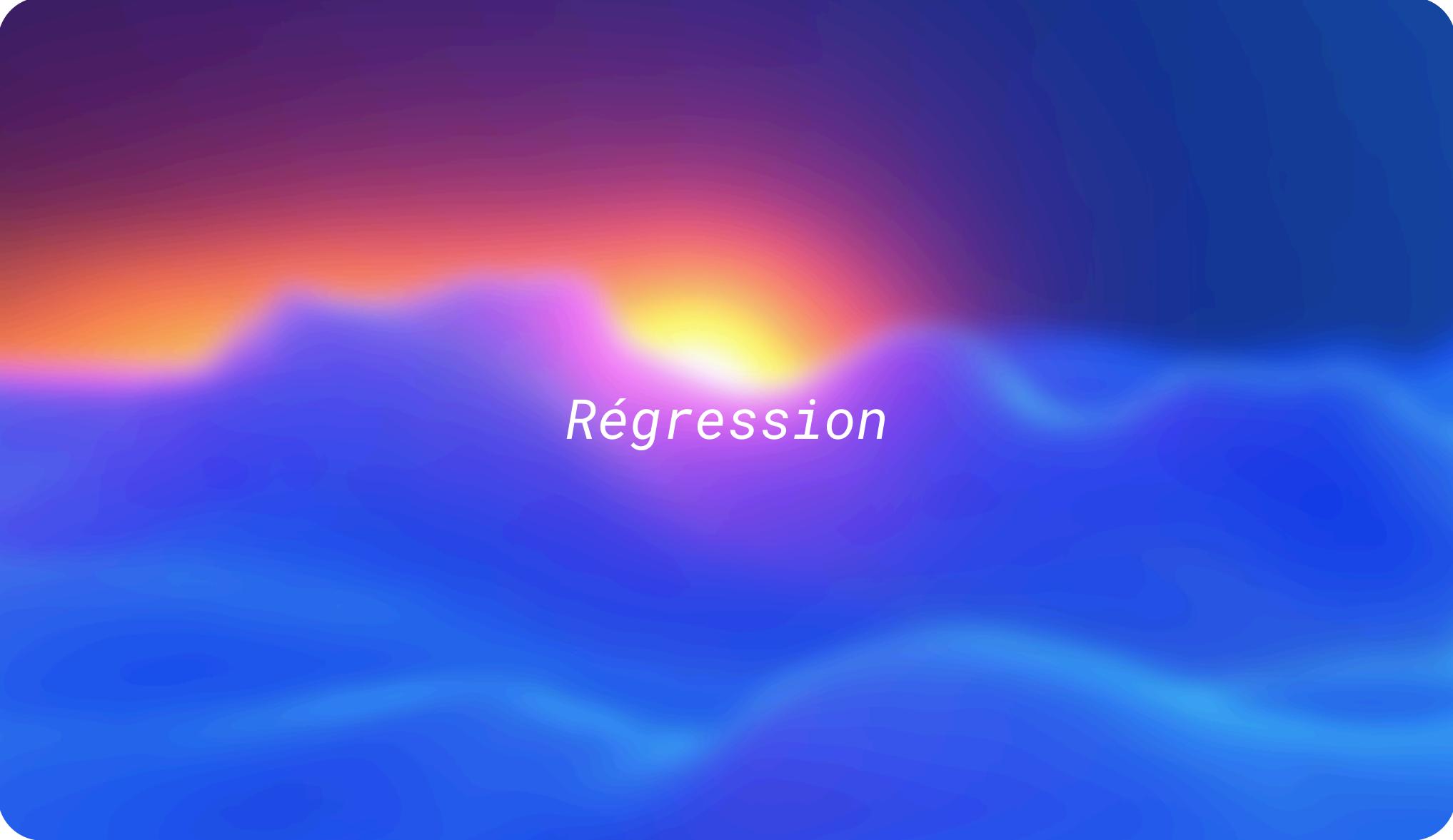
Exactitude
(Accuracy)

Précision
(Precision)

Rappel
(Recall)

Score F1
(F1 Score)

Métriques



Régression

Erreur Absolue Moyenne

$$\text{MAE} = \frac{1}{n} \sum_{i=1}^n |Y_i - \hat{Y}_i|$$

Er r e u r Q u a d r a t i q u e M o y e n n e

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$$

Er r e u r Q u a d r a t i q u e M o y e n n e

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$$

Racine de l'Erreur Quadratique Moyenne

$$\text{RMSE} = \sqrt{\frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2}$$

Validation croisée

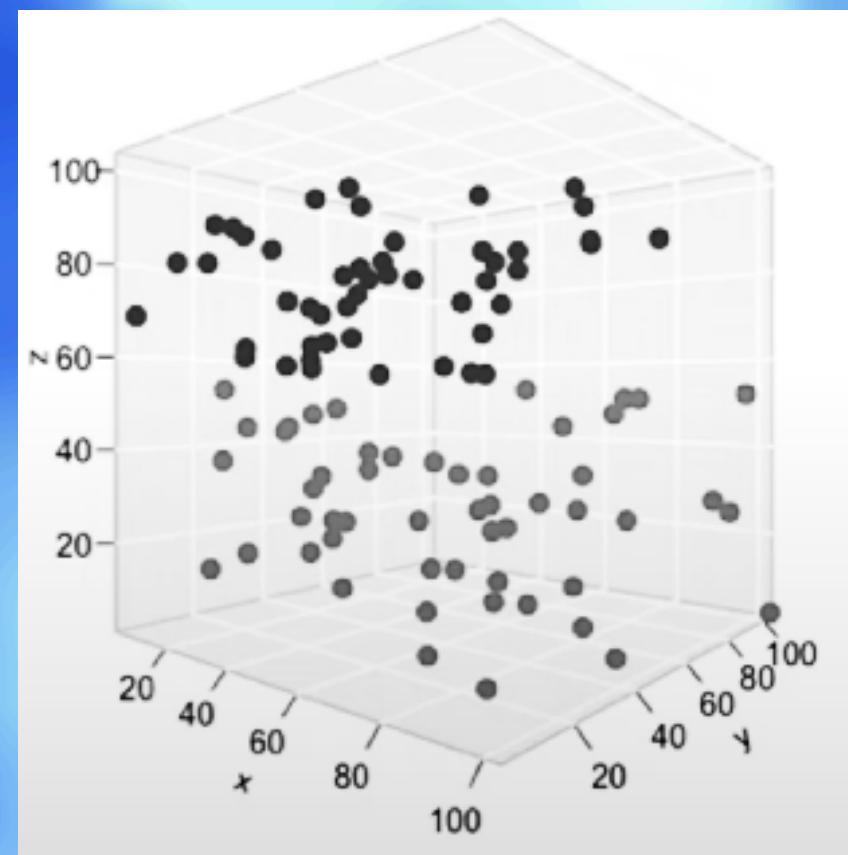
Train
(80 %)

Test
(20 %)

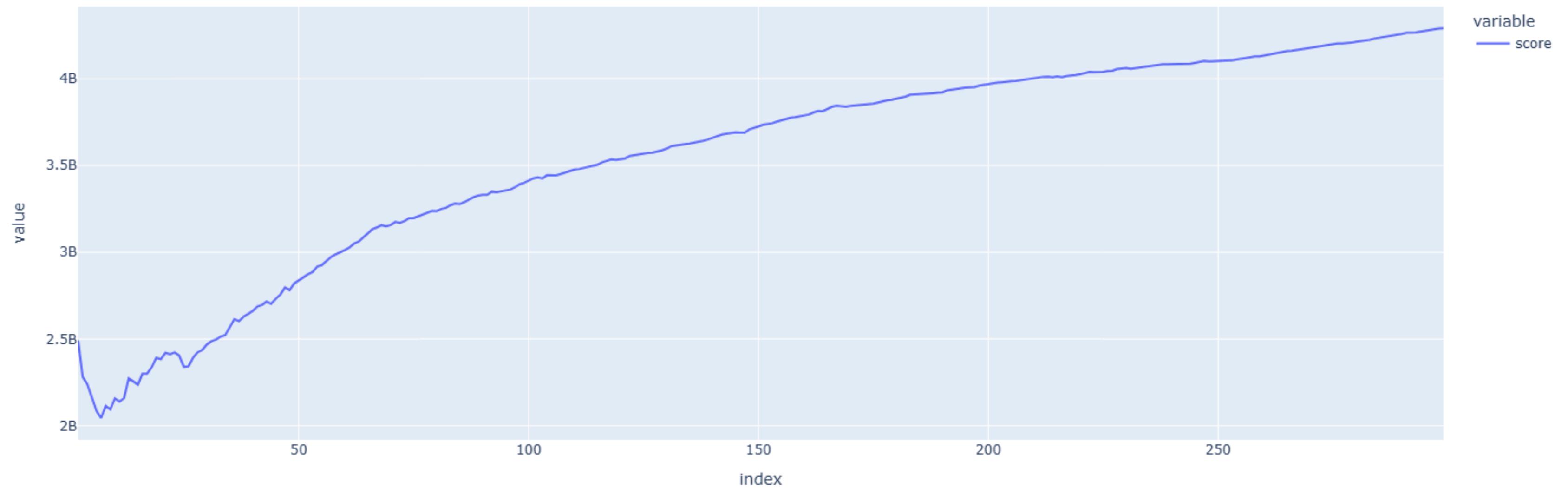
kNN



kNN



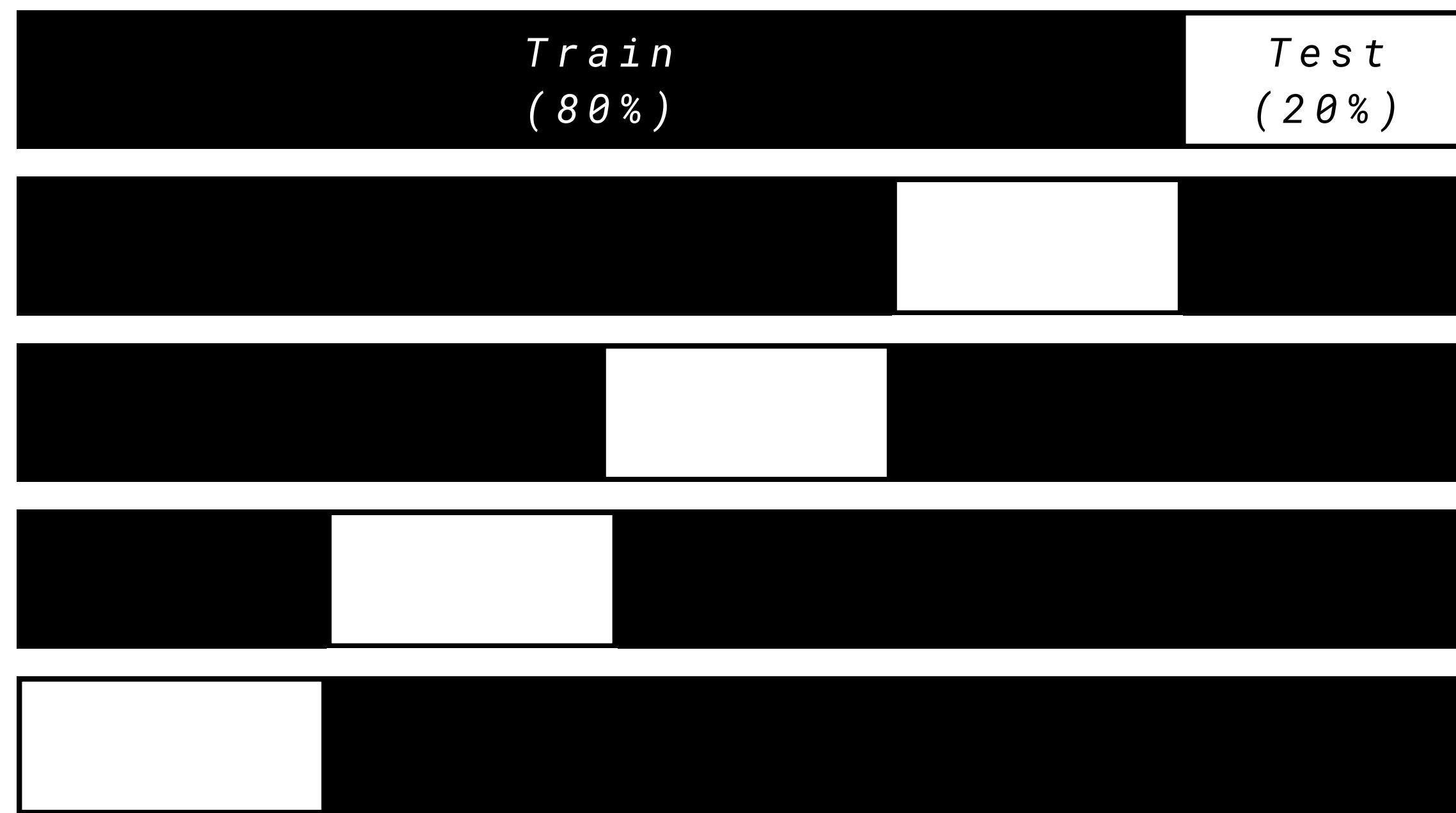
k NN > Hyper - paramètres



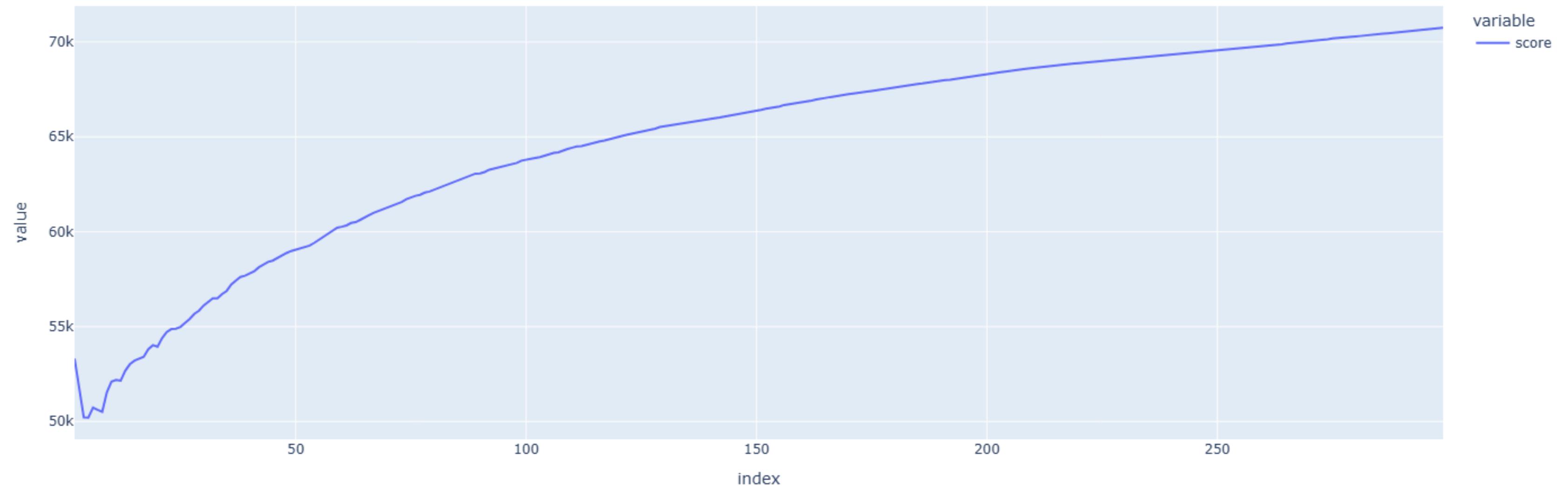
Validation croisée > KFold



Validation croisée > KFold



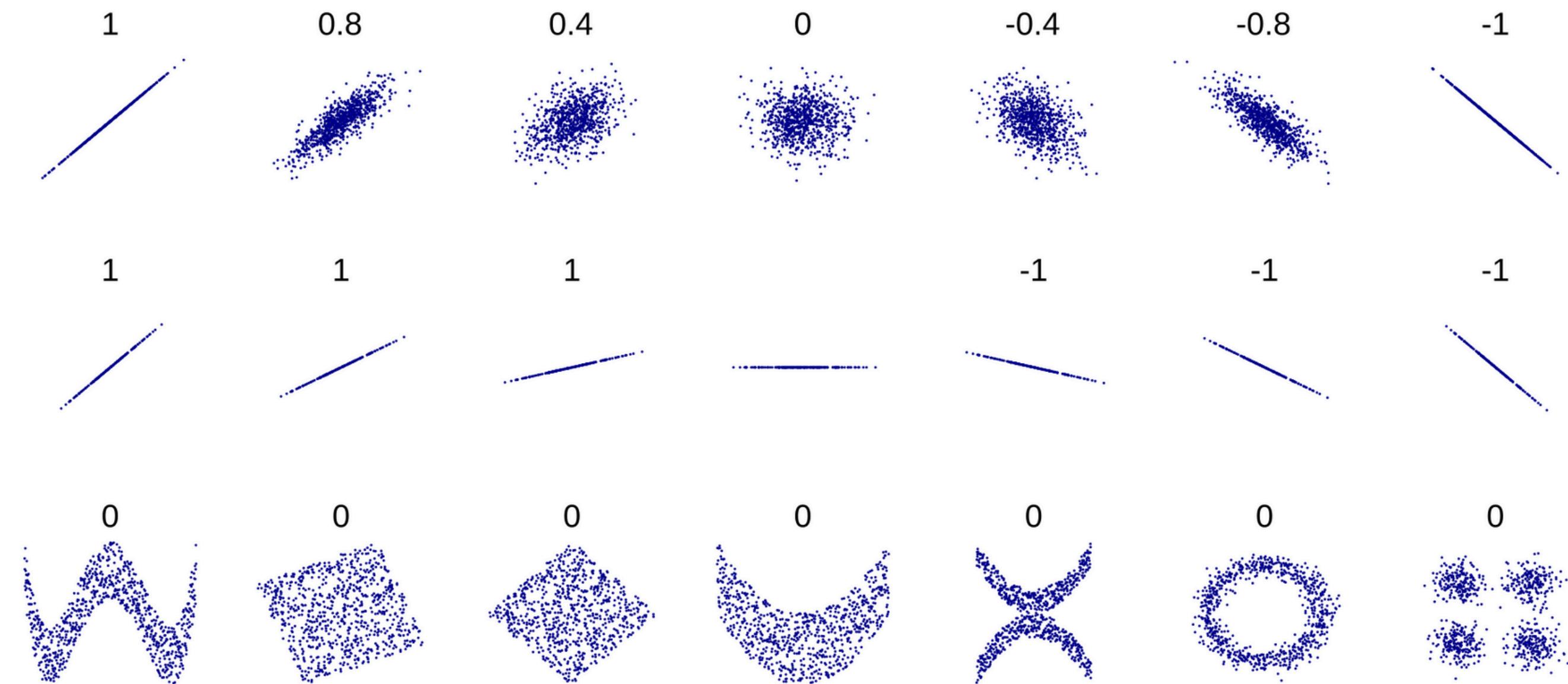
kNN > Hyper-paramètres sous KFold



Corrélation

$$r = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y}$$

Corrélation

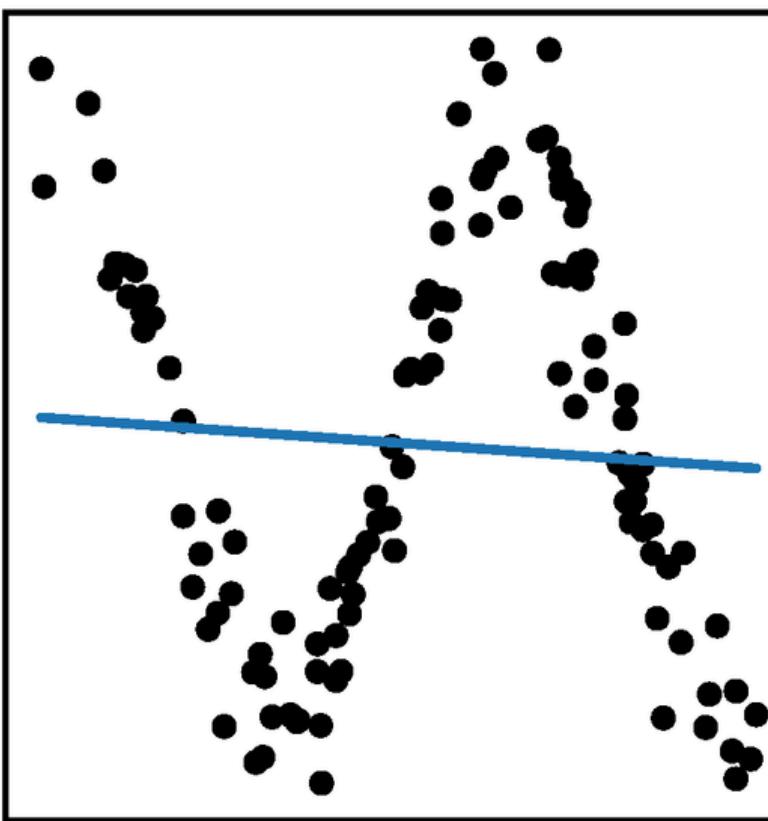


Information mutuelle

$$I(X;Y) = D_{\text{KL}}(P_{(X,Y)} \| P_X \otimes P_Y)$$

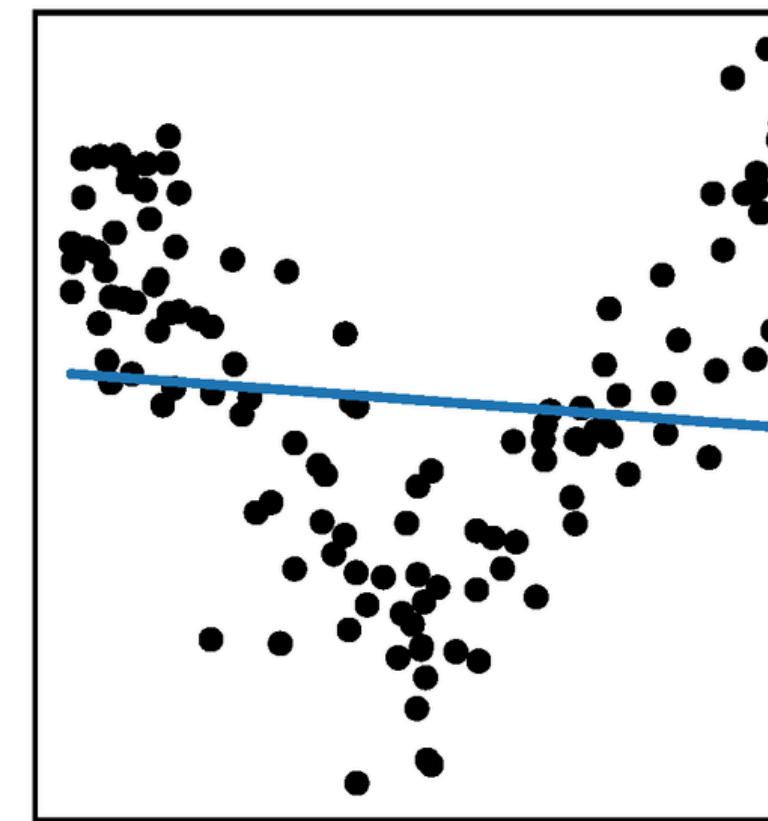
Information mutuelle

(a) $y = \sin(x) + \varepsilon$



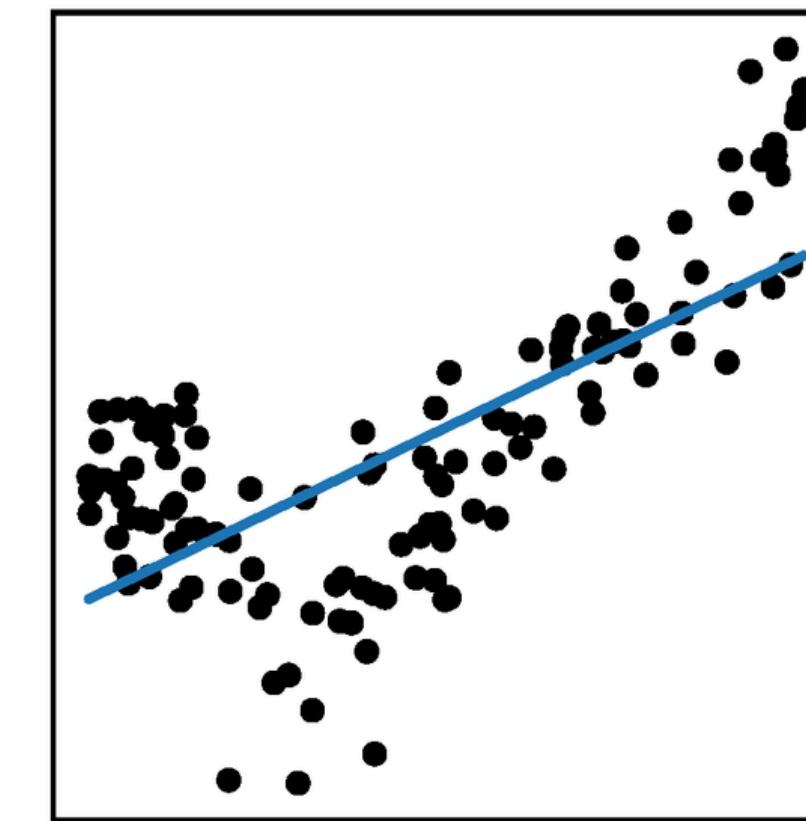
Original: -0.07
Transformed: 0.96
MI: 0.94

(b) $y = |x| + \varepsilon$



Original: -0.10
Transformed: 0.85
MI: 0.85

(c) $y = |x + 0.4| + \varepsilon$



Original: 0.70
Transformed: 0.89
MI: 0.89

[E t v o i l à . . .]

Fin de l'atelier :)

SIAHAAN - GENSOLLEN Rémy , 2025