

Préparation des données

I/collecte des données :

→ Données structurées : un tableau individus X

Variable	x_1	...	x_n
individu			

→ Données non structurées : tout ce qui n'est pas organisé sous forme d'un tableau de données :

la messagerie, les images, les vidéos

II/Nettoyage des données :

détection des données manquantes, des données aberrantes

détection des données manquantes

Résultat : `na.fail()`
↓
summary

détection des données aberrantes

nuage de points
boîte à moustache

Si $\frac{\text{taille de NA}}{\text{dim(data)}} < 1\%$
on supprime

traitement des données aberrantes et des données manquantes

traitement des données manquantes
- calcul des NA

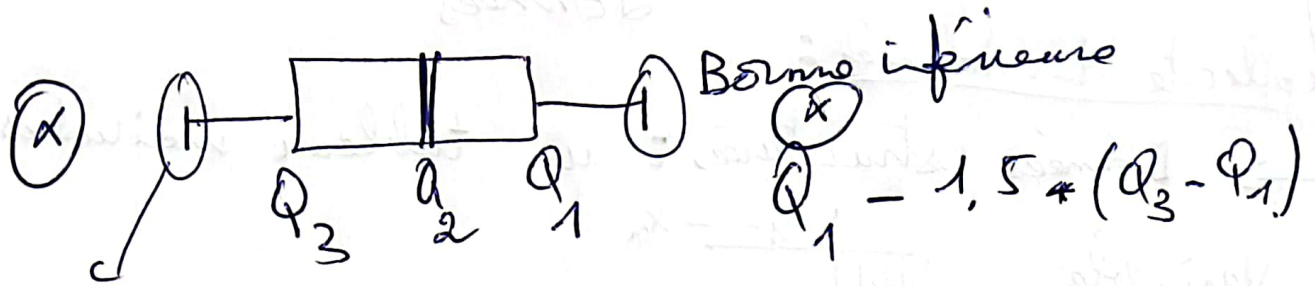
$\frac{\text{sum(is.na(data))}}{\text{Prod(dim(data))}}$

Imputation simple

Generalized imputation

similar case imputation

→ détection des valeurs aberrantes :
boîte à moustache



Bornes
supérieures

$$Q_3 + 1.5 * (Q_3 - Q_1)$$

→ traitement des valeurs manquantes :

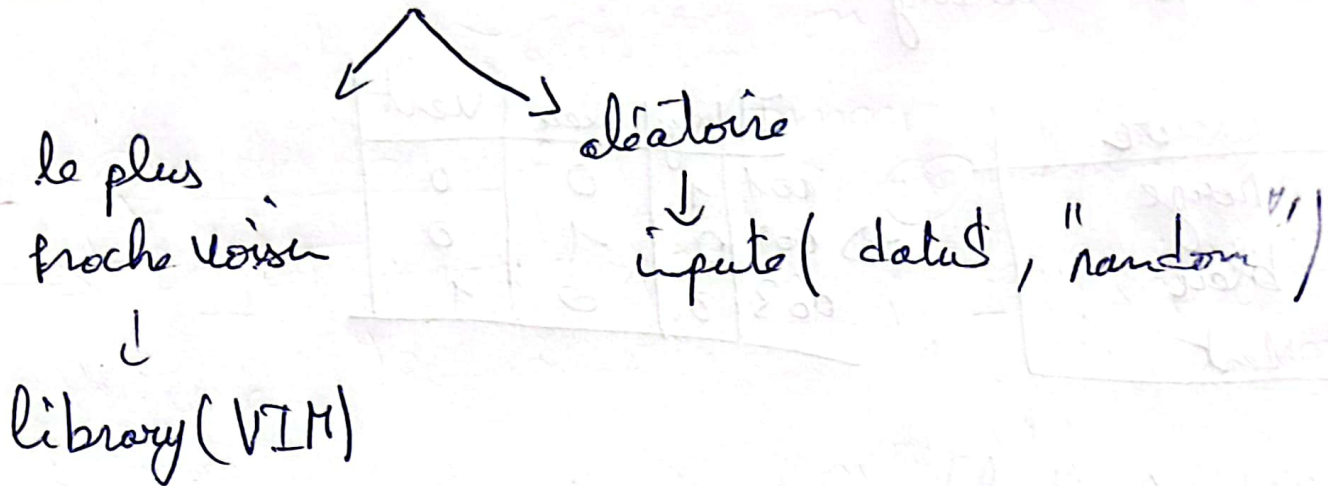
1 Generalized imputation : on remplace les données
manquantes par : la moyenne } variable quantitative
la médiane } variable qualitative

le mode variable qualitative

library(Hmisc)

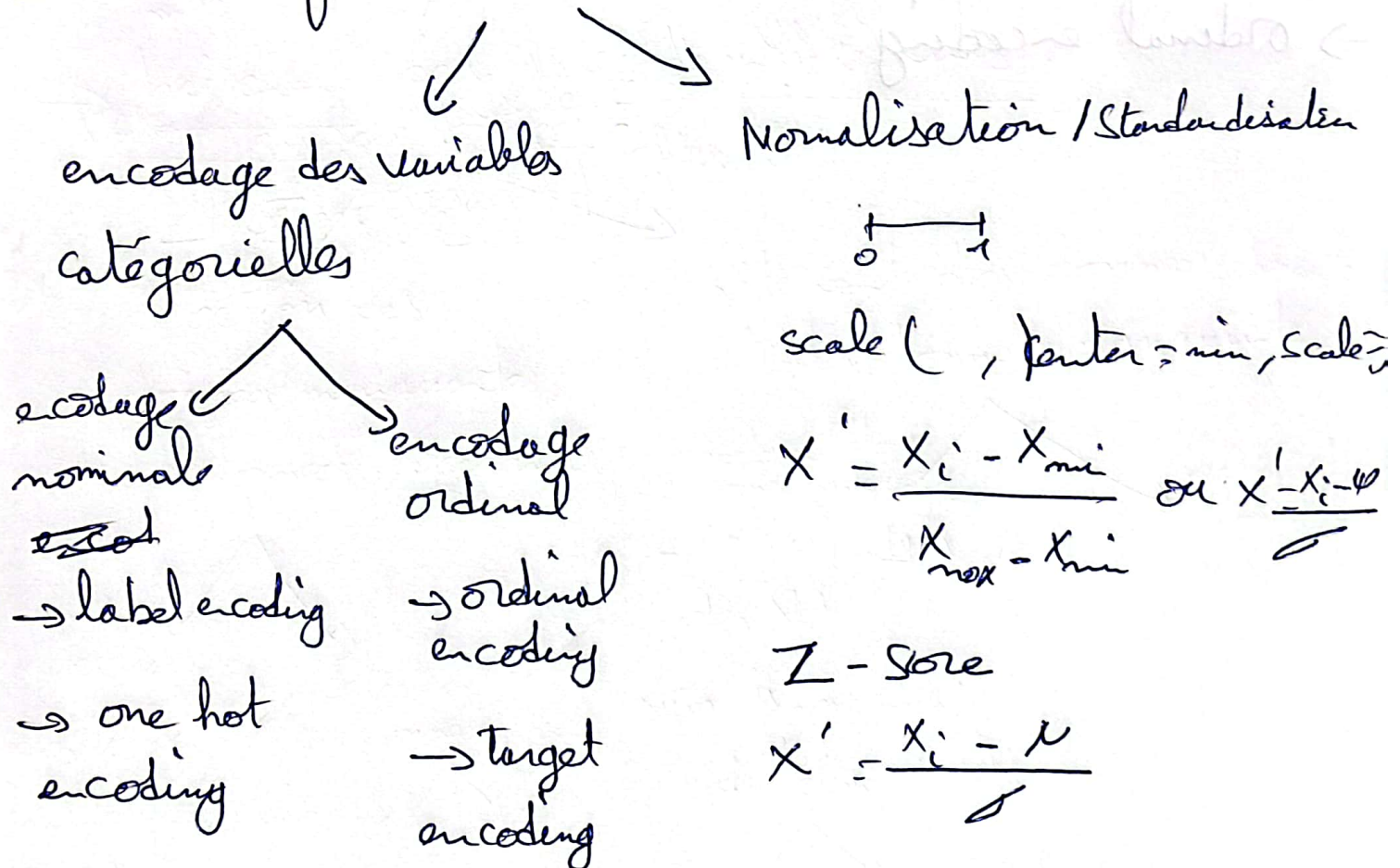
> impute (data\$, fun=medef)

2 Similar case imputation



$KNN(data[, c(" ", " ")])$

III transformation des données



transformation des données

1 encodage des variables catégorielles

l'objectif de l'encodage: l'encodage est la conversion des variables catégorielles en format numérique.

Pourquoi encoder les variables catégorielles?

les algorithmes de machine learning nécessitent des entrées numériques

Meilleure performance des modèles :

Amélioration de la performance des modèles en facilitant leur capacité à identifier les relations dans les données

Méthodes d'encodage

encodage nominal
(variable nominale)

One hot encoding

Target encoding (encodage par la moyenne)

encodage ordinal
(variable ordinales)

label encoding (encodage des étiquettes)

ordinal guidé par la cible ①

Exemples

* Encodage nominal : transforme chaque modalité en une colonne qui contient des valeurs binaires

ID	couleur
1	Rouge
2	Vert
3	Bleu
4	Rouge
5	Bleu

sous R

> data = data.frame (ID=1:5,
couleur = c("Rouge", "Vert", "Bleu", "Rouge",
"Bleu"))

> one-hot-encoded =
model.matrix (~ couleur - 1, data=data)

	Bleu	Rouge	Vert	(- 1) retire la colonne (ID)
0	0	1	0	
0	0	0	1	
1	0	0	0	
0	1	1	0	
1	0	0	0	

* Encodage ordinal : label encoding : transforme les catégories en valeurs numériques selon leur ordre

ID	Niveau
1	faible
2	Moyen
3	élevé
4	Moyen
5	faible
6	élevé

Sous R

> data = data.frame (ID = 1:6,
Niveau = c("Faible", "Moyen", "élevé",
"Moyen", "faible", "élevé"))

> data \$ Niveau_encode =
as.numeric (factor (data \$ Niveau,
levels = c("Faible", "Moyen", "Elevé")))

ID	Niveau	Niveau_encode
1	faible	1
2	Moyen	2
3	élevé	3
4	Moyen	2
5	faible	1
6	élevé	3

[2] la normalisation et la standardisation :

la normalisation et la standardisation sont deux techniques pour transformer les valeurs des données afin de les rendre comparables et améliorer les performances des modèles d'apprentissage automatique.

(2)

→ Normalisation : la normalisation consiste à transformer les données afin qu'elles se situent dans une plage spécifique, généralement entre 0 et 1. elle est utilisée lorsque vous voulez mettre toutes les caractéristiques sur une même échelle.

Min - Max

$$X_{\text{normalisé}} = \frac{X - X_{\min}}{X_{\max} - X_{\min}}$$

Scalr = scale()

utilisée pour les algorithmes basés sur des distances tels que : les k-plus proches voisins (KNN) ou les réseaux neuronaux, où les différences d'échelle peuvent influencer les résultats.

→ standardisation : la standardisation consiste à transformer les données pour qu'elles aient une moyenne 0 et un écart-type 1, autrement dit pour que les données suivent une distribution normale centrée réduite.

Z - Score

$$X_{\text{standardisé}} = \frac{X - \mu}{\sigma}$$

utilisé avec des algorithmes qui supposent une distribution normale :

→ Normalisation : la normalisation consiste à transformer les données afin qu'elles se situent dans une plage spécifique, généralement entre 0 et 1. elle est utilisée lorsque vous voulez mettre toutes les caractéristiques sur une même échelle.

Min - Max

$$X_{\text{normalisé}} = \frac{X - X_{\min}}{\underbrace{X_{\max} - X_{\min}}_{\text{étendue}}}$$

Scalr R : `scale(data, center = min, scale = max)`

utilisée pour les algorithmes basés sur des distances tels que : les k-plus proches voisins (KNN) ou les réseaux neuronaux, où les différences d'échelle peuvent influencer les résultats.

→ standardisation : la standardisation consiste à transformer les données pour qu'elles aient une moyenne 0 et un écart-type 1, autrement dit pour que les données suivent une distribution normale centrée réduite.

Z - Score

`scale(data, center = mean, scale = 1)`

$$X_{\text{standardisé}} = \frac{X - \mu}{\sigma}$$

utilisé avec des algorithmes qui supposent une distribution normale :