

Partie1_Apurement

KABASSINA Gnimdou Ange & Ibrahima GNING

2022-06-27

Contents

Introduction	1
I. Choix et constitution de la base d'étude	1
I.1 Importation des tables de données des sections 1 & 2	2
I.2 Extraction des variables d'intérêt	2
I.3 Fusion des données	3
II. Recodage ou labélisation des valeurs des variables	4
III. Gestion des valeurs manquantes	8
IV. Gestion des valeurs aberrantes	13
variables portant sur la migration	14

Introduction

Dans les pays en développement, les enquêtes auprès des ménages sont la source privilégiée pour produire des indicateurs de suivi des conditions de vie des populations. c'est dans ce contexte que dans l'espace UEMOA (Union Economique et Monétaire Ouest Africaine) pour le suivi et l'évaluation de la pauvreté et des conditions de vie des ménages dans chacun des pays membres, il a été mis en place l'Enquête Harmonisée sur les Conditions de vie des Ménages. Et c'est l'analyse des données recueillies dans le cadre de cette enquête au Sénégal qui fait l'objet de cette étude. La présente partie du travail consiste à préparer une table de données propre pour l'analyse : l'appurement.

I. Choix et constitution de la base d'étude

Les données recueillies et mises à notre disposition sont réparties en 20 sections portant chacune sur un aspect de la vie des ménages. Le but de notre travail sera d'étudier le volet éducation contenu dans la section 2 ainsi que quelques caractéristiques sociodémographiques des ménages répertoriés dans la section 1. Ainsi notre base d'étude portera sur certaines variables des section 1 et 2.

I.1 Importation des tables de données des sections 1 & 2

Pour constituer notre base nous utiliserons les informations contenues dans les sections 1 et 2. La première étape est donc d'importer les données qui y correspondent.

```
## Appel de la bibliothèque nécessaire pour ouvrir des fichier du format dta
library(haven)

## Définissons le chemin d'accès aux tables de données
path <- "C:\\Users\\HP ProBook\\Documents\\Mes Cours ISEP2\\R\\Projet\\Base de données"

## importation des bases des sections 1 & 2
section_1 <- read_dta(paste(path,"s01_me_SEN_2021.dta",sep="\\")) # tables de la section 1
section_2 <- read_dta(paste(path,"s02_me_SEN_2021.dta",sep="\\")) # tables de la section 2
```

I.2 Extraction des variables d'intérêt

I.2.1 Caractéristiques Sociodémographiques

Ici nous nous intéresseront à l'étude de:

- L'âge de mariage (pour voir entre autres si les mariages précoces perssistent)
- La nationalité des différents enquêtés (pour voir le degré d'ouverture du pays c'est à dire s'il y a assez de diversité)
- Les causes des migrations et les localités d'origines des migrants

Pour atteindre ces objectifs on s'intéressera à certaines variables à savoir :

- les clés d'identification des ménages : nous les utiliserons pour fusionner les deux sous bases que nous aurons créés-
- le sexe : pour des analyses stratifiées
- Anné de naissance : pour calculer l'âge des enquêtés (pour mieux emputer certaines valeurs manquantes et faire des analyses statifiée)
- Localité d'origine
- Causes de l'installation dans la localité de résidence
- Nationalité

```
# choisissons nos variables d'intérêt, tout en les renommant
## section 1: caractéristiques
carac_men <- section_1[1:6] # variables d'identification tel que la clé du ménage
names(carac_men)[6] <- "member_id"
carac_men[c("Sexe", "Situation_Matrimoniale", "Age_au_premier_mariage", "Nationalite", "Ne_dans_sa_localite",
            "Localite_d_origine_des_migrants", "Principale_cause_de_migration",
            , "Annee_de_naissance")] <- section_1[c("s01q01", "s01q07", "s01q10", "s01q15", "s01q17", "s01q18", "s01q19", "s01q20", "s01q21", "s01q22", "s01q23", "s01q24", "s01q25", "s01q26", "s01q27", "s01q28", "s01q29", "s01q30", "s01q31", "s01q32", "s01q33", "s01q34", "s01q35", "s01q36", "s01q37", "s01q38", "s01q39", "s01q40", "s01q41", "s01q42", "s01q43", "s01q44", "s01q45", "s01q46", "s01q47", "s01q48", "s01q49", "s01q50", "s01q51", "s01q52", "s01q53", "s01q54", "s01q55", "s01q56", "s01q57", "s01q58", "s01q59", "s01q60", "s01q61", "s01q62", "s01q63", "s01q64", "s01q65", "s01q66", "s01q67", "s01q68", "s01q69", "s01q70", "s01q71", "s01q72", "s01q73", "s01q74", "s01q75", "s01q76", "s01q77", "s01q78", "s01q79", "s01q80", "s01q81", "s01q82", "s01q83", "s01q84", "s01q85", "s01q86", "s01q87", "s01q88", "s01q89", "s01q90", "s01q91", "s01q92", "s01q93", "s01q94", "s01q95", "s01q96", "s01q97", "s01q98", "s01q99", "s01q100", "s01q101", "s01q102", "s01q103", "s01q104", "s01q105", "s01q106", "s01q107", "s01q108", "s01q109", "s01q110", "s01q111", "s01q112", "s01q113", "s01q114", "s01q115", "s01q116", "s01q117", "s01q118", "s01q119", "s01q120", "s01q121", "s01q122", "s01q123", "s01q124", "s01q125", "s01q126", "s01q127", "s01q128", "s01q129", "s01q130", "s01q131", "s01q132", "s01q133", "s01q134", "s01q135", "s01q136", "s01q137", "s01q138", "s01q139", "s01q140", "s01q141", "s01q142", "s01q143", "s01q144", "s01q145", "s01q146", "s01q147", "s01q148", "s01q149", "s01q150", "s01q151", "s01q152", "s01q153", "s01q154", "s01q155", "s01q156", "s01q157", "s01q158", "s01q159", "s01q160", "s01q161", "s01q162", "s01q163", "s01q164", "s01q165", "s01q166", "s01q167", "s01q168", "s01q169", "s01q170", "s01q171", "s01q172", "s01q173", "s01q174", "s01q175", "s01q176", "s01q177", "s01q178", "s01q179", "s01q180", "s01q181", "s01q182", "s01q183", "s01q184", "s01q185", "s01q186", "s01q187", "s01q188", "s01q189", "s01q190", "s01q191", "s01q192", "s01q193", "s01q194", "s01q195", "s01q196", "s01q197", "s01q198", "s01q199", "s01q200", "s01q201", "s01q202", "s01q203", "s01q204", "s01q205", "s01q206", "s01q207", "s01q208", "s01q209", "s01q210", "s01q211", "s01q212", "s01q213", "s01q214", "s01q215", "s01q216", "s01q217", "s01q218", "s01q219", "s01q220", "s01q221", "s01q222", "s01q223", "s01q224", "s01q225", "s01q226", "s01q227", "s01q228", "s01q229", "s01q230", "s01q231", "s01q232", "s01q233", "s01q234", "s01q235", "s01q236", "s01q237", "s01q238", "s01q239", "s01q240", "s01q241", "s01q242", "s01q243", "s01q244", "s01q245", "s01q246", "s01q247", "s01q248", "s01q249", "s01q250", "s01q251", "s01q252", "s01q253", "s01q254", "s01q255", "s01q256", "s01q257", "s01q258", "s01q259", "s01q260", "s01q261", "s01q262", "s01q263", "s01q264", "s01q265", "s01q266", "s01q267", "s01q268", "s01q269", "s01q270", "s01q271", "s01q272", "s01q273", "s01q274", "s01q275", "s01q276", "s01q277", "s01q278", "s01q279", "s01q280", "s01q281", "s01q282", "s01q283", "s01q284", "s01q285", "s01q286", "s01q287", "s01q288", "s01q289", "s01q290", "s01q291", "s01q292", "s01q293", "s01q294", "s01q295", "s01q296", "s01q297", "s01q298", "s01q299", "s01q300", "s01q301", "s01q302", "s01q303", "s01q304", "s01q305", "s01q306", "s01q307", "s01q308", "s01q309", "s01q310", "s01q311", "s01q312", "s01q313", "s01q314", "s01q315", "s01q316", "s01q317", "s01q318", "s01q319", "s01q320", "s01q321", "s01q322", "s01q323", "s01q324", "s01q325", "s01q326", "s01q327", "s01q328", "s01q329", "s01q330", "s01q331", "s01q332", "s01q333", "s01q334", "s01q335", "s01q336", "s01q337", "s01q338", "s01q339", "s01q340", "s01q341", "s01q342", "s01q343", "s01q344", "s01q345", "s01q346", "s01q347", "s01q348", "s01q349", "s01q350", "s01q351", "s01q352", "s01q353", "s01q354", "s01q355", "s01q356", "s01q357", "s01q358", "s01q359", "s01q360", "s01q361", "s01q362", "s01q363", "s01q364", "s01q365", "s01q366", "s01q367", "s01q368", "s01q369", "s01q370", "s01q371", "s01q372", "s01q373", "s01q374", "s01q375", "s01q376", "s01q377", "s01q378", "s01q379", "s01q380", "s01q381", "s01q382", "s01q383", "s01q384", "s01q385", "s01q386", "s01q387", "s01q388", "s01q389", "s01q390", "s01q391", "s01q392", "s01q393", "s01q394", "s01q395", "s01q396", "s01q397", "s01q398", "s01q399", "s01q400", "s01q401", "s01q402", "s01q403", "s01q404", "s01q405", "s01q406", "s01q407", "s01q408", "s01q409", "s01q410", "s01q411", "s01q412", "s01q413", "s01q414", "s01q415", "s01q416", "s01q417", "s01q418", "s01q419", "s01q420", "s01q421", "s01q422", "s01q423", "s01q424", "s01q425", "s01q426", "s01q427", "s01q428", "s01q429", "s01q430", "s01q431", "s01q432", "s01q433", "s01q434", "s01q435", "s01q436", "s01q437", "s01q438", "s01q439", "s01q440", "s01q441", "s01q442", "s01q443", "s01q444", "s01q445", "s01q446", "s01q447", "s01q448", "s01q449", "s01q450", "s01q451", "s01q452", "s01q453", "s01q454", "s01q455", "s01q456", "s01q457", "s01q458", "s01q459", "s01q460", "s01q461", "s01q462", "s01q463", "s01q464", "s01q465", "s01q466", "s01q467", "s01q468", "s01q469", "s01q470", "s01q471", "s01q472", "s01q473", "s01q474", "s01q475", "s01q476", "s01q477", "s01q478", "s01q479", "s01q480", "s01q481", "s01q482", "s01q483", "s01q484", "s01q485", "s01q486", "s01q487", "s01q488", "s01q489", "s01q490", "s01q491", "s01q492", "s01q493", "s01q494", "s01q495", "s01q496", "s01q497", "s01q498", "s01q499", "s01q500", "s01q501", "s01q502", "s01q503", "s01q504", "s01q505", "s01q506", "s01q507", "s01q508", "s01q509", "s01q510", "s01q511", "s01q512", "s01q513", "s01q514", "s01q515", "s01q516", "s01q517", "s01q518", "s01q519", "s01q520", "s01q521", "s01q522", "s01q523", "s01q524", "s01q525", "s01q526", "s01q527", "s01q528", "s01q529", "s01q530", "s01q531", "s01q532", "s01q533", "s01q534", "s01q535", "s01q536", "s01q537", "s01q538", "s01q539", "s01q540", "s01q541", "s01q542", "s01q543", "s01q544", "s01q545", "s01q546", "s01q547", "s01q548", "s01q549", "s01q550", "s01q551", "s01q552", "s01q553", "s01q554", "s01q555", "s01q556", "s01q557", "s01q558", "s01q559", "s01q560", "s01q561", "s01q562", "s01q563", "s01q564", "s01q565", "s01q566", "s01q567", "s01q568", "s01q569", "s01q570", "s01q571", "s01q572", "s01q573", "s01q574", "s01q575", "s01q576", "s01q577", "s01q578", "s01q579", "s01q580", "s01q581", "s01q582", "s01q583", "s01q584", "s01q585", "s01q586", "s01q587", "s01q588", "s01q589", "s01q590", "s01q591", "s01q592", "s01q593", "s01q594", "s01q595", "s01q596", "s01q597", "s01q598", "s01q599", "s01q600", "s01q601", "s01q602", "s01q603", "s01q604", "s01q605", "s01q606", "s01q607", "s01q608", "s01q609", "s01q610", "s01q611", "s01q612", "s01q613", "s01q614", "s01q615", "s01q616", "s01q617", "s01q618", "s01q619", "s01q620", "s01q621", "s01q622", "s01q623", "s01q624", "s01q625", "s01q626", "s01q627", "s01q628", "s01q629", "s01q630", "s01q631", "s01q632", "s01q633", "s01q634", "s01q635", "s01q636", "s01q637", "s01q638", "s01q639", "s01q640", "s01q641", "s01q642", "s01q643", "s01q644", "s01q645", "s01q646", "s01q647", "s01q648", "s01q649", "s01q650", "s01q651", "s01q652", "s01q653", "s01q654", "s01q655", "s01q656", "s01q657", "s01q658", "s01q659", "s01q660", "s01q661", "s01q662", "s01q663", "s01q664", "s01q665", "s01q666", "s01q667", "s01q668", "s01q669", "s01q670", "s01q671", "s01q672", "s01q673", "s01q674", "s01q675", "s01q676", "s01q677", "s01q678", "s01q679", "s01q680", "s01q681", "s01q682", "s01q683", "s01q684", "s01q685", "s01q686", "s01q687", "s01q688", "s01q689", "s01q690", "s01q691", "s01q692", "s01q693", "s01q694", "s01q695", "s01q696", "s01q697", "s01q698", "s01q699", "s01q700", "s01q701", "s01q702", "s01q703", "s01q704", "s01q705", "s01q706", "s01q707", "s01q708", "s01q709", "s01q710", "s01q711", "s01q712", "s01q713", "s01q714", "s01q715", "s01q716", "s01q717", "s01q718", "s01q719", "s01q720", "s01q721", "s01q722", "s01q723", "s01q724", "s01q725", "s01q726", "s01q727", "s01q728", "s01q729", "s01q730", "s01q731", "s01q732", "s01q733", "s01q734", "s01q735", "s01q736", "s01q737", "s01q738", "s01q739", "s01q740", "s01q741", "s01q742", "s01q743", "s01q744", "s01q745", "s01q746", "s01q747", "s01q748", "s01q749", "s01q750", "s01q751", "s01q752", "s01q753", "s01q754", "s01q755", "s01q756", "s01q757", "s01q758", "s01q759", "s01q760", "s01q761", "s01q762", "s01q763", "s01q764", "s01q765", "s01q766", "s01q767", "s01q768", "s01q769", "s01q770", "s01q771", "s01q772", "s01q773", "s01q774", "s01q775", "s01q776", "s01q777", "s01q778", "s01q779", "s01q780", "s01q781", "s01q782", "s01q783", "s01q784", "s01q785", "s01q786", "s01q787", "s01q788", "s01q789", "s01q790", "s01q791", "s01q792", "s01q793", "s01q794", "s01q795", "s01q796", "s01q797", "s01q798", "s01q799", "s01q800", "s01q801", "s01q802", "s01q803", "s01q804", "s01q805", "s01q806", "s01q807", "s01q808", "s01q809", "s01q810", "s01q811", "s01q812", "s01q813", "s01q814", "s01q815", "s01q816", "s01q817", "s01q818", "s01q819", "s01q820", "s01q821", "s01q822", "s01q823", "s01q824", "s01q825", "s01q826", "s01q827", "s01q828", "s01q829", "s01q830", "s01q831", "s01q832", "s01q833", "s01q834", "s01q835", "s01q836", "s01q837", "s01q838", "s01q839", "s01q840", "s01q841", "s01q842", "s01q843", "s01q844", "s01q845", "s01q846", "s01q847", "s01q848", "s01q849", "s01q850", "s01q851", "s01q852", "s01q853", "s01q854", "s01q855", "s01q856", "s01q857", "s01q858", "s01q859", "s01q860", "s01q861", "s01q862", "s01q863", "s01q864", "s01q865", "s01q866", "s01q867", "s01q868", "s01q869", "s01q870", "s01q871", "s01q872", "s01q873", "s01q874", "s01q875", "s01q876", "s01q877", "s01q878", "s01q879", "s01q880", "s01q881", "s01q882", "s01q883", "s01q884", "s01q885", "s01q886", "s01q887", "s01q888", "s01q889", "s01q890", "s01q891", "s01q892", "s01q893", "s01q894", "s01q895", "s01q896", "s01q897", "s01q898", "s01q899", "s01q900", "s01q901", "s01q902", "s01q903", "s01q904", "s01q905", "s01q906", "s01q907", "s01q908", "s01q909", "s01q910", "s01q911", "s01q912", "s01q913", "s01q914", "s01q915", "s01q916", "s01q917", "s01q918", "s01q919", "s01q920", "s01q921", "s01q922", "s01q923", "s01q924", "s01q925", "s01q926", "s01q927", "s01q928", "s01q929", "s01q930", "s01q931", "s01q932", "s01q933", "s01q934", "s01q935", "s01q936", "s01q937", "s01q938", "s01q939", "s01q940", "s01q941", "s01q942", "s01q943", "s01q944", "s01q945", "s01q946", "s01q947", "s01q948", "s01q949", "s01q950", "s01q951", "s01q952", "s01q953", "s01q954", "s01q955", "s01q956", "s01q957", "s01q958", "s01q959", "s01q960", "s01q961", "s01q962", "s01q963", "s01q964", "s01q965", "s01q966", "s01q967", "s01q968", "s01q969", "s01q970", "s01q971", "s01q972", "s01q973", "s01q974", "s01q975", "s01q976", "s01q977", "s01q978", "s01q979", "s01q980", "s01q981", "s01q982", "s01q983", "s01q984", "s01q985", "s01q986", "s01q987", "s01q988", "s01q989", "s01q990", "s01q991", "s01q992", "s01q993", "s01q994", "s01q995", "s01q996", "s01q997", "s01q998", "s01q999", "s01q1000", "s01q1001", "s01q1002", "s01q1003", "s01q1004", "s01q1005", "s01q1006", "s01q1007", "s01q1008", "s01q1009", "s01q1010", "s01q1011", "s01q1012", "s01q1013", "s01q1014", "s01q1015", "s01q1016", "s01q1017", "s01q1018", "s01q1019", "s01q1020", "s01q1021", "s01q1022", "s01q1023", "s01q1024", "s01q1025", "s01q1026", "s01q1027", "s01q1028", "s01q1029", "s01q1030", "s01q1031", "s01q1032", "s01q1033", "s01q1034", "s01q1035", "s01q1036", "s01q1037", "s01q1038", "s01q1039", "s01q1040", "s01q1041", "s01q1042", "s01q1043", "s01q1044", "s01q1045", "s01q1046", "s01q1047", "s01q1048", "s01q1049", "s01q1050", "s01q1051", "s01q1052", "s01q1053", "s01q1054", "s01q1055", "s01q1056", "s01q1057", "s01q1058", "s01q1059", "s01q1060", "s01q1061", "s01q1062", "s01q1063", "s01q1064", "s01q1065", "s01q1066", "s01q1067", "s01q1068", "s01q1069", "s01q1070", "s01q1071", "s01q1072", "s01q1073", "s01q1074", "s01q1075", "s01q1076", "s01q1077", "s01q1078", "s01q1079", "s01q1080", "s01q1081", "s01q1082", "s01q1083", "s01q1084", "s01q1085", "s01q1086", "s01q1087", "s01q1088", "s01q1089", "s01q1090", "s01q1091", "s01q1092", "s01q1093", "s01q1094", "s01q1095", "s01q1096", "s01q1097", "s01q1098", "s01q1099", "s01q1100", "s01q1101", "s01q1102", "s01q1103", "s01q1104", "s01q1105", "s01q1106", "s01q1107", "s01q1108", "s01q1109", "s01q1110", "s01q1111", "s01q1112", "s01q1113", "s01q1114", "s01q1115", "s01q1116", "s01q1117", "s01q1118", "s01q1119", "s01q1120", "s01q1121", "s01q1122", "s01q1123", "s01q1124", "s01q1125", "s01q1126", "s01q1127", "s01q1128", "s01q1129", "s01q1130", "s01q1131", "s01q1132", "s01q1133", "s01q1134", "s01q1135", "s01q1136", "s01q1137", "s01q1138", "s01q1139", "s01q1140", "s01q1141", "s01q1142", "s01q1143", "s01q1144", "s01q1145", "s01q1146", "s01q1147", "s01q1148", "s01q1149", "s01q1150", "s01q1151", "s01q1152", "s01q1153", "s01q1154", "s01q1155", "s01q1156", "s01q1157", "s01q1158", "s01q1159", "s01q1160", "s01q1161", "s01q1162", "s01q1163", "s01q1164", "s01q1165", "s01q1166", "s01q1167", "s01q1168", "s01q1169", "s01q1170", "s01q1171", "s01q1172", "s01q1173", "s01q1174", "s01q1175", "s01q1176", "s01q1177", "s01q1178", "s01q1179", "s01q1180", "s01q1181", "s01q1182", "s01q1183", "s01q1184", "s01q1185", "s01q1186", "s01q1187", "s01q1188", "s01q1189", "s01q1190", "s01q1191", "s01q1192", "s01q1193", "s01q1194", "s01q1195", "s01q1196", "s01q1197", "s01q1198", "s01q1199", "s01q1200", "s01q1201", "s01q1202", "s01q1203", "s01q1204", "s01q1205", "s01q1206", "s01q1207", "s01q1208", "s01q1209", "s01q1210", "s01q1211", "s01q1212", "s01q1213", "s01q1214", "s01q1215", "s01q1216", "s01q1217", "s01q1218", "s01q1219", "s01q1220", "s01q1221", "s01q1222", "s01q1223", "s01q1224", "s01q1225", "s01q1226", "s01q1227", "s01q1228", "s01q1229", "s01q1230", "s01q1231", "s01q1232", "s01q1233", "s01q1234", "s01q1235", "s01q1236", "s01q1237", "s01q1238", "s01q1239", "s01q1240", "s01q1241", "s01q1242", "s01q1243", "s01q1244", "s01q1245", "s01q1246", "s01q1247", "s01q1248", "s01q1249", "s01q1250", "s01q1251", "s01q1252", "s01q1253", "s01q1254", "s01q1255", "s01q1256", "s01q1257", "s01q1258", "s01q1259", "s01q1260", "s01q1261", "s01q1262", "s01q1263", "s01q1264", "s0
```

I.2.2 Education

Dans le volet éducation, il sera question d'étudier le niveau d'alphabétisation ; les problèmes généralement rencontrés par les élèves et étudiants sénégalais ; les raisons d'abandons et aussi de non scolarisation (on pourra voir par exemple si les gens continuent à penser que les filles n'ont pas le droit de faire des études formelles); l'âge d'entrée à l'école... On pourra aussi voir l'impact de la COVID-19 en voyant si les écoles ont fermé et si pendant la fermeture les élèves ont pu rester en contact avec l'administration de l'école ou les enseignants. Ainsi quelques unes des variables qui nous intéressent dans la section deux sont:

- Si l'enquêté peut lire, écrire, comprendre, un texte en français, en langue locale ou dans une autre langues;

__ Si l'enquêté fait des études formelles;

- Age d'entrée à l'école
- problèmes rencontrés à l'école
- Moyens utilisés pour rester en contact avec l'école pendant la fermeture

```
# choisissons nos variables d'intérêt, tout en les renommant
## section 2: Education
edu <- section_2[1:6] # variables d'identification tel que la clé du ménage
names(edu)[6] <- "member_id"
edu[c("Peut_lire_français", "Peut_lire_la_langue_locale", "Peut_lire_une_autre_langue", "Peut_ecrire_fr",
edu[c("Peut_comprendre_un_texte_en_français", "Peut_comprendre_un_texte_en_langue_locale", "Peut_compre
edu[c("Fait_une_ecole_formelles", "Raisons_de_non_frequentation_d_ecole_formelle", "Fait_une_ecole_non_fo
edu[c("Fermeture_de_l_ecole", "Reste_en_contact_avec_l_ecole", "Resultat_2019_2020", "Raisons_d_abandon_de
edu[c("Diplome_le_plus_eleve", "Reste_en_contact_avec_l_ecole_par_SMS", "Reste_en_contact_avec_l_ecole_pa
edu[c("Insuffisance_de_livres_et_ou_de_fournitures", "Insuffisance_de_tables_bancs_et_d_equipements", "Abs
```

I.3 Fusion des données

Maintenant que nous avons pu extraire les informations qui nous intéressent, nous allons toutes les regrouper en un seul endroit. En d'autres termes on devra à présent fusionner les deux sous bases constituées grâce à la fonction merge.

```
library(dplyr)

##
## Attachement du package : 'dplyr'

## Les objets suivants sont masqués depuis 'package:stats':
##
##      filter, lag

## Les objets suivants sont masqués depuis 'package:base':
##
##      intersect, setdiff, setequal, union
```

```
# fusion des deux sous tables pour avoir notre table de données d'étude
## création de la clé de fusion
carac_men <- mutate(carac_men, key = paste(interview__key,interview__id,id_menage,member_id))
edu <- mutate(edu, key = paste(interview__key,interview__id,id_menage,member_id))
my_data <- merge(carac_men, edu, by = "key")

dim(my_data)
```

```
## [1] 22017    54
```

Notre base d'étude ainsi constituée comprends 22017 lignes et 58 colonnes. Nous allons à présent y retirer les variables qui se répètent et celle qui ne seront pas trop essentielles à notre analyse. Ensuite nous ajouterons une variable Age dans laquelle nous calculerons l'âge de l'enquêté grâce à son année de naissance. Aussi, les observation en double (doublons) ne nous interesse t-elles pas car elles biaiserons nos analyses ; il faut donc s'en débarrasser.

```
# Supprimons les doublons et les variables en double qui ne nous interessent pas
# c'est à dire les identifiants des ménages
my_data1 <- distinct(my_data) # suppression des doublons
dim(my_data1)
```

```
## [1] 21768    54
```

```
my_data2 <- my_data1[-c(2:7,16:21)] # suppression des colonnes ininteressantes

# calculons l'âge des individus
my_data2 <- mutate(my_data2, Age = 2022 - Annee_de_naissance)
```

Les données recueillies sur l'éducation ne concernent que les individus de plus de 3 ans; de même la majeure partie des variables qui nous intéressent et des analyses que nous souhaitons faire, n'ont pas vraiment d'intérêt pour des personnes aussi jeunes (moins de 3 ans). Du coup nous éliminerons cette tranche de la population dans nos données.

```
# excluons les individus de moins de 3 ans car ils ne nous intéressent pas
my_data3 <- subset(my_data2, Age > 3)
dim(my_data3)
```

```
## [1] 13817    43
```

```
# sauvegardons la base ainsi obtenue dans un fichier csv.
write.csv(my_data3, paste(path,"Base_projet2.csv", sep="//"))
```

Notre table de données finale comporte donc 13817 lignes et 43 colonnes ou variables.

II. Recodage ou labélisation des valeurs des variables

Lorsque les données sont importées depuis SAS, SPSS ou Stata, cela permet notamment de conserver le codage original du fichier importé. Mais il faut noter que ces étiquettes de valeur n'indique pas pour autant de manière systématique le type de variable (catégorielle ou continue). Les vecteurs labellisés n'ont donc pas

vocation à être utilisés pour l'analyse, notamment le calcul de modèles statistiques. Ils doivent être convertis en facteurs (pour les variables catégorielles) ou en vecteurs numériques (pour les variables continues). Pour cela nous utiliserons la fonction `to_factor` de la bibliothèque "labelled"

```
# Recodons les valeurs des variables
# Méthode 1 (fastidieuse)
## sexe: 1 = masculin & 2 = féminin
library(dplyr)
my_data4 <- mutate(my_data3, Sexe = if_else(Sexe == 1, "Homme", "Femme", missing = NULL))

# Méthode 2 : transformer en type factor moins fastidieux
# Recodons les valeurs des variables grâce à la fonction to_factor de la librairie "labelled"
library(labelled)
str(my_data3) # pour détecter les variables catégorielles à recoder

## 'data.frame': 13817 obs. of 43 variables:
## $ key : chr "00-01-63-80 1df091103eb748ac859408a..."
## $ Sexe : dbl+lbl [1:13817] 1, 2, 1, 1, 1, 1, ...
## ..@ label : chr "1.01. Quel est le sexe de %roster..."
## ..@ format.stata: chr "%10.0g"
## ..@ labels : Named num 1 2
## .. ..- attr(*, "names")= chr [1:2] "Masculin" "Féminin"
## $ Situation_Matrimoniale : dbl+lbl [1:13817] 2, 2, 1, 1, 1, 1, ...
## ..@ label : chr "1.07. Quelle est la situation matrimoniale de %roster..."
## ..@ format.stata: chr "%10.0g"
## ..@ labels : Named num 1 2 3 4 5 6 7
## .. ..- attr(*, "names")= chr [1:7] "Célibataire" "Marié(e) monogame" "Marié(e) polygame" "Union l..."
## $ Age_au_premier_mariage : num 25 NA NA NA NA NA 19 17 NA NA ...
## $ Nationalite : dbl+lbl [1:13817] 13, 13, 13, 13, 13, 13, ...
## ..@ label : chr "1.15. De quelle nationalité est %roster..."
## ..@ format.stata: chr "%10.0g"
## ..@ labels : Named num 1 2 3 4 5 6 7 8 9 10 ...
## .. ..- attr(*, "names")= chr [1:17] "Bénin" "Burkina Faso" "Cape-vert" "Cote d'ivoire" ...
## $ Ne_dans_sa_localite_de_residence : dbl+lbl [1:13817] 1, NA, 1, 1, 1, 1, ...
## ..@ label : chr "1.17. %roster... est-il né au %s01q15?"
## ..@ format.stata: chr "%10.0g"
## ..@ labels : Named num 1 2
## .. ..- attr(*, "names")= chr [1:2] "Oui" "Non"
## $ Localite_d_origine_des_migrants : dbl+lbl [1:13817] NA, 1, NA, NA, NA, NA, ...
## ..@ label : chr "1.19. Quelle est la dernière localité où %roster... a vécu avant de v..."
## ..@ format.stata: chr "%10.0g"
## ..@ labels : Named num 1 4 5 6 7 8 9 10 11 12 ...
## .. ..- attr(*, "names")= chr [1:20] "Capitale" "Bénin" "Burkina Faso" "Cape-vert" ...
## $ Principale_cause_de_migration : dbl+lbl [1:13817] NA, 9, NA, NA, NA, NA, ...
## ..@ label : chr "1.21. Quelle était la raison principale pour laquelle %roster... est v..."
## ..@ format.stata: chr "%10.0g"
## ..@ labels : Named num 1 2 3 4 5 6 7 8 9 10 ...
## .. ..- attr(*, "names")= chr [1:20] "Envoyé par sa famille pour travailler" "Est venu avec son pa..."
## $ Annee_de_naissance : num 1983 1997 2001 1991 1991 ...
## $ Peut_lire_français : dbl+lbl [1:13817] 1, 2, 2, NA, 1, 2, ...
## ..@ format.stata: chr "%10.0g"
## ..@ labels : Named num 1 2
## .. ..- attr(*, "names")= chr [1:2] "Oui" "Non"
## $ Peut_lire_la_langue_locale : dbl+lbl [1:13817] 2, 2, 2, NA, 2, 2, ...
```

```

## ..@ format.stata: chr "%10.0g"
## ..@ labels : Named num 1 2
## ..- attr(*, "names")= chr [1:2] "Oui" "Non"
## $ Peut_lire_une_autre_langue : dbl+lbl [1:13817] 2, NA, 2, NA, 2, 2,
## ..@ format.stata: chr "%10.0g"
## ..@ labels : Named num 1 2
## ..- attr(*, "names")= chr [1:2] "Oui" "Non"
## $ Peut_ecrire_français : num 1 0 0 NA 1 0 0 0 0 0 ...
## $ Peut_ecrire_la_langue_locale : num NA 0 0 NA 0 0 0 0 0 0 ...
## $ Peut_ecrire_une_autre_langue : num 0 0 0 NA 0 0 0 0 0 NA ...
## $ Peut_comprendre_un_texte_en_français : num 1 0 0 NA 1 0 0 0 0 0 ...
## $ Peut_comprendre_un_texte_en_langue_locale : num 0 0 0 NA 0 0 0 NA NA 0 ...
## $ Peut_comprendre_un_texte_en_autre_langue : num 0 0 0 NA 0 0 0 0 NA 0 ...
## $ Fait_une_ecole_formelles : dbl+lbl [1:13817] 1, 2, NA, NA, 1, NA,
## ..@ label : chr "2.03. %rosteritle% a-t-il fait ou fait-il des études actuellement dans un
## ..@ format.stata: chr "%10.0g"
## ..@ labels : Named num 1 2
## ..- attr(*, "names")= chr [1:2] "Oui" "Non"
## $ Raisons_de_non_frequentation_d_ecole_formelle : dbl+lbl [1:13817] NA, 14, 14, NA, NA, 14,
## ..@ label : chr "2.04. Pour quelle raison principale %rosteritle% n'a-t-il pas fait des ét
## ..@ format.stata: chr "%10.0g"
## ..@ labels : Named num 1 2 3 4 5 6 7 8 9 10 ...
## ..- attr(*, "names")= chr [1:15] "Trop jeune" "Pas d'école, école trop éloignée" "Refus de la
## $ Fait_une_ecole_non_formelle : dbl+lbl [1:13817] NA, 1, 1, NA, NA, 1,
## ..@ label : chr "2.05. Est ce que %rosteritle% a suivi une école non formelle ou une form
## ..@ format.stata: chr "%10.0g"
## ..@ labels : Named num 1 2
## ..- attr(*, "names")= chr [1:2] "Oui" "Non"
## $ Age_d_entre_a_l_ecole : num 6 NA NA NA 7 NA NA NA NA NA ...
## $ Gerant_de_l_ecole : dbl+lbl [1:13817] NA, NA, NA, NA, NA, NA,
## ..@ label : chr "2.09. Qui gère l'école fréquentée par %rosteritle% au cours de l'année %a
## ..@ format.stata: chr "%10.0g"
## ..@ labels : Named num 1 2 3 4 5 6
## ..- attr(*, "names")= chr [1:6] "Gouvernement" "Privé religieux" "Privé non religieux" "Privé
## $ Fermeture_de_l_ecole : dbl+lbl [1:13817] NA, NA, NA, NA, NA, NA,
## ..@ label : chr "2.09a. L'école de %rosteritle% a-t-elle été fermée momentanément à caus
## ..@ format.stata: chr "%10.0g"
## ..@ labels : Named num 1 2
## ..- attr(*, "names")= chr [1:2] "Oui" "Non"
## $ Reste_en_contact_avec_l_ecole : dbl+lbl [1:13817] NA, NA, NA, NA, NA, NA,
## ..@ label : chr "2.09b. %rosteritle% était-il/elle en contact avec les enseignants ou l'
## ..@ format.stata: chr "%10.0g"
## ..@ labels : Named num 1 2
## ..- attr(*, "names")= chr [1:2] "Oui" "Non"
## $ Resultat_2019_2020 : dbl+lbl [1:13817] NA, NA, NA, NA, NA, NA,
## ..@ label : chr "2.10. Quel résultat %rosteritle% a-t-il obtenu au cours de l'année %annee
## ..@ format.stata: chr "%10.0g"
## ..@ labels : Named num 1 2 3 4 5
## ..- attr(*, "names")= chr [1:5] "Diplômé, études achevées" "Passe en classe supérieure" "Echec
## $ Raisons_d_abandon_de_l_ecole : dbl+lbl [1:13817] NA, NA, NA, NA, NA, NA,
## ..@ label : chr "2.11. Pour quelle raison %rosteritle% a-t-il abandonné l'école en cours d
## ..@ format.stata: chr "%10.0g"
## ..@ labels : Named num 1 2 3 4 5 6 7 8 9 10 ...
## ..- attr(*, "names")= chr [1:16] "A obtenu un emploi" "S'est marié" "C'est une fille" "Grosses

```

```

## $ Diplome_le_plus_eleve : dbl+lbl [1:13817] 0, NA, NA, NA, 1, NA,
## ..@ label : chr "2.33. Quel est le diplôme le plus élevé obtenu par %rosteritle%?"
## ..@ format.stata: chr "%10.0g"
## ..@ labels : Named num 0 1 2 3 4 5 6 7 8 9 ...
## ..- attr(*, "names")= chr [1:11] "Aucun" "CEPE" "BEPC" "CAP" ...
## $ Reste_en_contact_avec_l_ecole_par_SMS : num NA NA NA NA NA NA NA NA NA NA NA ...
## $ Reste_en_contact_avec_l_ecole_par_applications_mobiles: num NA NA NA NA NA NA NA NA NA NA NA ...
## $ Reste_en_contact_avec_l_ecole_par_Email : num NA NA NA NA NA NA NA NA NA NA NA ...
## $ Reste_en_contact_avec_l_ecole_par_courier : num NA NA NA NA NA NA NA NA NA NA NA ...
## $ Reste_en_contact_avec_l_ecole_Telephone : num NA NA NA NA NA NA NA NA NA NA NA ...
## $ Insuffisance_de_livres_et_ou_de_fournitures : dbl+lbl [1:13817] NA, NA, NA, NA, NA, NA,
## ..@ format.stata: chr "%10.0g"
## ..@ labels : Named num 1 2
## ..- attr(*, "names")= chr [1:2] "Oui" "Non"
## $ Insuffisance_de_tables_bancs_et_d_equipements : dbl+lbl [1:13817] NA, NA, NA, NA, NA, NA,
## ..@ format.stata: chr "%10.0g"
## ..@ labels : Named num 1 2
## ..- attr(*, "names")= chr [1:2] "Oui" "Non"
## $ Absenteisme_des_enseignants_greve : dbl+lbl [1:13817] NA, NA, NA, NA, NA, NA,
## ..@ format.stata: chr "%10.0g"
## ..@ labels : Named num 1 2
## ..- attr(*, "names")= chr [1:2] "Oui" "Non"
## $ Enseignement_pas_satisfaisant : dbl+lbl [1:13817] NA, NA, NA, NA, NA, NA,
## ..@ format.stata: chr "%10.0g"
## ..@ labels : Named num 1 2
## ..- attr(*, "names")= chr [1:2] "Oui" "Non"
## $ Effectifs_plethoriques : dbl+lbl [1:13817] NA, NA, NA, NA, NA, NA,
## ..@ format.stata: chr "%10.0g"
## ..@ labels : Named num 1 2
## ..- attr(*, "names")= chr [1:2] "Oui" "Non"
## $ Insuffisance_d_enseignants : dbl+lbl [1:13817] NA, NA, NA, NA, NA, NA,
## ..@ format.stata: chr "%10.0g"
## ..@ labels : Named num 1 2
## ..- attr(*, "names")= chr [1:2] "Oui" "Non"
## $ Manque_de_toilettes : dbl+lbl [1:13817] NA, NA, NA, NA, NA, NA,
## ..@ format.stata: chr "%10.0g"
## ..@ labels : Named num 1 2
## ..- attr(*, "names")= chr [1:2] "Oui" "Non"
## $ Frequence_des_cotisations : dbl+lbl [1:13817] NA, NA, NA, NA, NA, NA,
## ..@ format.stata: chr "%10.0g"
## ..@ labels : Named num 1 2
## ..- attr(*, "names")= chr [1:2] "Oui" "Non"
## $ Salle_de_classe_en_mauvais_etat : dbl+lbl [1:13817] NA, NA, NA, NA, NA, NA,
## ..@ format.stata: chr "%10.0g"
## ..@ labels : Named num 1 2
## ..- attr(*, "names")= chr [1:2] "Oui" "Non"
## $ Age : num 39 25 21 31 31 26 47 45 6 7 ...

```

variables à recoder

```

variables = c("Sexe", "Situation_Matrimoniale", "Nationalite", "Ne_dans_sa_localite_de_residence",
  "Localite_d_origine_des_migrants", "Principale_cause_de_migration", "Fait_une_ecole_formelle",
  "Raisons_de_non_frequentation_d_ecole_formelle", "Fait_une_ecole_non_formelle", "Gerant_de",
  "Fermeture_de_l_ecole", "Reste_en_contact_avec_l_ecole", "Resultat_2019_2020", "Raisons_d_ab",
  "Diplome_le_plus_eleve", "Insuffisance_de_livres_et_ou_de_fournitures", "Insuffisance_de_ta

```

```
my_data4[variables] <- to_factor(my_data3[variables]) # recodage
```


d'imputer ces NA reviendrait pratiquement à créer nos propres observations ce qui est dangereux. Ainsi nous excluons tout simplement le volet 'Impact de la covid sur le fonctionnement des écoles' dans notre analyse.

```
covid <- c("Fermeture_de_l_ecole", "Reste_en_contact_avec_l_ecole", "Reste_en_contact_avec_l_ecole_par_SMA")
my_data4 <- select(my_data4, -covid)
```

```
## Note: Using an external vector in selections is ambiguous.
## i Use 'all_of(covid)' instead of 'covid' to silence this message.
## i See <https://tidyselect.r-lib.org/reference/faq-external-vector.html>.
## This message is displayed once per session.
```

- Facultés de maîtrise d'une langue :

Une des analyses qui sera portée sur ces données sera de savoir globalement est-ce que les individus maîtrisent le français par exemple. Pour cela les valeurs manquantes de ces variables (catégorielles) seront remplacées par le mode parce que le mode représente la modalité avec la plus grande probabilité d'observation ou d'apparition. En d'autres termes si on pose une question à un individu, il y a une plus grande chance que sa réponse soit le mode ce qui nous permettra d'évaluer la situation la plus probable en terme de maîtrise des langues.

```
## Variables de maîtrise de la langue
langue <- c("Peut_lire_français", "Peut_lire_la_langue_locale", "Peut_lire_une_autre_langue",
            "Peut_ecrire_français", "Peut_ecrire_la_langue_locale", "Peut_ecrire_une_autre_langue",
            "Peut_comprendre_un_texte_en_français", "Peut_comprendre_un_texte_en_langue_locale",
            "Peut_comprendre_un_texte_en_autre_langue")

### cherchons le mode et détection des na et remplacement par le mode
for (x in langue){
  y <- sort(table(my_data4[x])) # on range les modalités en fonction de leurs fréquences dans l'ordre croissant
  print(y) # pour voir les différentes modalités et le mode
  y <- data.frame(y)
  na <- which(is.na(my_data4[x])) # détection des valeurs manquantes
  my_data4[na, x] <- y[2, 1] # remplacement des na
}
```

```
## Peut_lire_français
## Oui Non
## 5678 6139
## Peut_lire_la_langue_locale
## Oui Non
## 574 11336
## Peut_lire_une_autre_langue
## Oui Non
## 2258 9632
## Peut_ecrire_français
## 1 0
## 5460 6410
## Peut_ecrire_la_langue_locale
## 1 0
## 548 11302
## Peut_ecrire_une_autre_langue
## 1 0
## 2195 9623
```

```
## Peut_comprendre_un_texte_en_français
##      1      0
## 5223 6625
## Peut_comprendre_un_texte_en_langue_locale
##      1      0
##      955 10872
## Peut_comprendre_un_texte_en_autre_langue
##      1      0
## 2026 9789
```

- Fréquentation d'une école ou non et raisons de non fréquentation

Pour les valeurs manquantes à ce niveau, elle seront remplacées par le mode comme précédemment.

```
formation <- c("Fait_une_ecole_formelles", "Raisons_de_non_frequentation_d_ecole_formelle", "Fait_une_ecole_formelle")
for (x in formation){
  y <- sort(table(my_data4[x])) # on range les modalités en fonction de leurs
  # fréquence dans l'ordre croissant. Ainsi le mode est le premier élément de y
  y <- data.frame(y)
  na <- which(is.na(my_data4[x])) # detection des valeurs manquantes
  my_data4[na,x] <- y[2,1] # remplacement des na
}
```

- Age d'entrée à l'école

La gestion des valeurs manquantes pour cette variable consistera en leur remplacement par la médiane car cette dernière est moins sensible aux valeurs aberrantes.

```
## gestion des NA pour l'âge d'entrée à l'école

x <- which(is.na(my_data4["Age_d_entree_a_l_ecole"])) # detection des Na
my_data4[x, "Age_d_entree_a_l_ecole"] <- median(my_data4[which(is.na(my_data4["Age_d_entree_a_l_ecole"]))])
```

- Principales difficultés rencontrées à l'école :

A ce niveau, nous utiliserons une approche puissante celle des K plus proches voisins. Le principe de cette méthode est simple et assez intuitif. Il se base sur l'hypothèse selon laquelle les individus appartenant à un même groupe ont tendance à avoir des comportements similaires c'est à dire à présenter les mêmes. En d'autres termes pour remplacer la valeur manquante d'un individu, on constitue un groupe des K individus les plus proches de ce individu en termes d'attributs ou caractéristiques et on remplace la valeur manquante par la moyenne des observations pour une variable numérique et par le mode (vote majoritaire) s'il s'agit de variables catégorielles. Cette approche est implémentée dans R grâce à la fonction "KNN" de la bibliothèque "VIM"

```
## gestion des variables relatives aux difficultés rencontrées à l'école

library(VIM)
```

```
## Le chargement a nécessité le package : colorspace
```

```
## Le chargement a nécessité le package : grid
```



```

demographie <- c("Situation_Matrimoniale", "Age_au_premier_mariage",
               "Ne_dans_sa_localite_de_residence", "Localite_d_origine_des_migrants", "Principale_cause_
demographie_na <- c() # vecteur qui recueillera les index de ces lignes
for (x in demographie){
  x <- which(is.na(my_data4[x])) # detection des valeur manquantes
  demographie_na <- c(demographie_na, x) # ajout des index aux autres précédement déterminés
}
y = data_frame(demographie_na)
demographie_na <- demographie_na[!duplicated(y)] # supprimer les doublons car on ne veut pas que les in

library(mice)

```

```

##
## Attachement du package : 'mice'

## L'objet suivant est masqué depuis 'package:stats':
##
##      filter

## Les objets suivants sont masqués depuis 'package:base':
##
##      cbind, rbind

```

```

my_data4[demographie_na,] <- complete(mice(my_data4[demographie_na,],meth = "cart"),1) # imputation

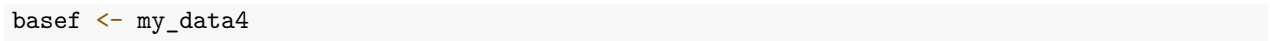
```

```

##
## iter imp variable
## 1 1 Situation_Matrimoniale Age_au_premier_mariage Ne_dans_sa_localite_de_residence Localite_
## 1 2 Situation_Matrimoniale Age_au_premier_mariage Ne_dans_sa_localite_de_residence Localite_
## 1 3 Situation_Matrimoniale Age_au_premier_mariage Ne_dans_sa_localite_de_residence Localite_
## 1 4 Situation_Matrimoniale Age_au_premier_mariage Ne_dans_sa_localite_de_residence Localite_
## 1 5 Situation_Matrimoniale Age_au_premier_mariage Ne_dans_sa_localite_de_residence Localite_
## 2 1 Situation_Matrimoniale Age_au_premier_mariage Ne_dans_sa_localite_de_residence Localite_
## 2 2 Situation_Matrimoniale Age_au_premier_mariage Ne_dans_sa_localite_de_residence Localite_
## 2 3 Situation_Matrimoniale Age_au_premier_mariage Ne_dans_sa_localite_de_residence Localite_
## 2 4 Situation_Matrimoniale Age_au_premier_mariage Ne_dans_sa_localite_de_residence Localite_
## 2 5 Situation_Matrimoniale Age_au_premier_mariage Ne_dans_sa_localite_de_residence Localite_
## 3 1 Situation_Matrimoniale Age_au_premier_mariage Ne_dans_sa_localite_de_residence Localite_
## 3 2 Situation_Matrimoniale Age_au_premier_mariage Ne_dans_sa_localite_de_residence Localite_
## 3 3 Situation_Matrimoniale Age_au_premier_mariage Ne_dans_sa_localite_de_residence Localite_
## 3 4 Situation_Matrimoniale Age_au_premier_mariage Ne_dans_sa_localite_de_residence Localite_
## 3 5 Situation_Matrimoniale Age_au_premier_mariage Ne_dans_sa_localite_de_residence Localite_
## 4 1 Situation_Matrimoniale Age_au_premier_mariage Ne_dans_sa_localite_de_residence Localite_
## 4 2 Situation_Matrimoniale Age_au_premier_mariage Ne_dans_sa_localite_de_residence Localite_
## 4 3 Situation_Matrimoniale Age_au_premier_mariage Ne_dans_sa_localite_de_residence Localite_
## 4 4 Situation_Matrimoniale Age_au_premier_mariage Ne_dans_sa_localite_de_residence Localite_
## 4 5 Situation_Matrimoniale Age_au_premier_mariage Ne_dans_sa_localite_de_residence Localite_
## 5 1 Situation_Matrimoniale Age_au_premier_mariage Ne_dans_sa_localite_de_residence Localite_
## 5 2 Situation_Matrimoniale Age_au_premier_mariage Ne_dans_sa_localite_de_residence Localite_
## 5 3 Situation_Matrimoniale Age_au_premier_mariage Ne_dans_sa_localite_de_residence Localite_
## 5 4 Situation_Matrimoniale Age_au_premier_mariage Ne_dans_sa_localite_de_residence Localite_
## 5 5 Situation_Matrimoniale Age_au_premier_mariage Ne_dans_sa_localite_de_residence Localite_

```

```
vis_miss(my_data4)
```



IV. Gestion des valeurs aberrantes

- Age du premier mariage

13

```
## Age au premier mariage
## les célibataires ne sont pas concernés
y <- which(basef["Situation_Matrimoniale"] == "Célibataire" | basef["Age_au_premier_mariage"] == 9999)
basef[y,"Age_au_premier_mariage"] <- NA
```

variables portant sur la migration

les données sur la migration étant assez semblables, nous conserverons les données d'origines pour pouvoir faire des analyses plus pertinentes

```
## conserver les données d'origines pour faire des analyses plus pertinentes
migr <- c("Localite_d_origine_des_migrants","Principale_cause_de_migration")
basef[migr] <- to_factor(my_data3[migr])
```

- Facultés de maîtrise d'une langue :

ici nous allons juste vérifier les modalités ; a t-on bien 2 modalités ou plus ou moins ? ces dernières sont-elles conformes au questionnaire ?

```
## Variables de maîtrise de la langue
### Gestion valeurs aberrantes
### jettons un coup d'oeil sur les modalités de nos variables
### pour voir si tout va bien

for(x in langue){
  print(table(basef[x]))
}
```

```
## Peut_lire_français
## Oui Non
## 5678 8139
## Peut_lire_la_langue_locale
## Oui Non
## 574 13243
## Peut_lire_une_autre_langue
## Oui Non
## 2258 11559
## Peut_ecrire_français
## 0 1 2
## 6410 5460 1947
## Peut_ecrire_la_langue_locale
## 0 1 2
## 11302 548 1967
## Peut_ecrire_une_autre_langue
## 0 1 2
## 9623 2195 1999
## Peut_comprendre_un_texte_en_français
## 0 1 2
## 6625 5223 1969
## Peut_comprendre_un_texte_en_langue_locale
## 0 1 2
```

```
## 10872 955 1990
## Peut_comprendre_un_texte_en_autre_langue
## 0 1 2
## 9789 2026 2002
```

On voit qu'il y a 3 modalités au lieu de 2 pour certaines variables. De plus nous avons vu depuis l'imputation que ces dernières n'étaient pas conformes au questionnaire. on va donc régler tout ça. la modalité "2" provient de l'imputation et devrait plutôt être égal à 0 le mode (d'après les tableaux de fréquences). il y a donc eu un soucis et 0 doit être remplacé par 2 pour rester conforme au questionnaire.

```
langue2 <- c("Peut_ecrire_français", "Peut_ecrire_la_langue_locale", "Peut_ecrire_une_autre_langue",
             "Peut_comprendre_un_texte_en_français", "Peut_comprendre_un_texte_en_langue_locale",
             "Peut_comprendre_un_texte_en_autre_langue")
for(x in langue2){
  y <- which(basef[x] == 0)
  basef[y,x] <- 2
}
```

- Fréquentation d'une école ou non et raisons de non fréquentation

Cette partie porte sur les variables qui nous renseignent si l'individu fait une école formelle ou non; sinon pourquoi. alors un premier contrôle porte sur les raisons de non fréquentation d'une école formelle car, pour cette question, un individu qui fait déjà une formation formelle n'est pas concerné. Doù une partie des na proviennent du fait que l'individu fait effectivement une formation formelle, ainsi pour de tels enquêtés, la réponse à cette question sera remplacée par NA.

```
## vérification de réponses 2: faire en sorte que la réponse à la question "quels sont les raisons de n
x <- which(basef["Fait_une_ecole_formelles"] == "Oui") # detecter les observations où les individus fon
basef[x,"Raisons_de_non_frequentation_d_ecole_formelle"] <- NA # pour chaque individu (ligne) identifié
```

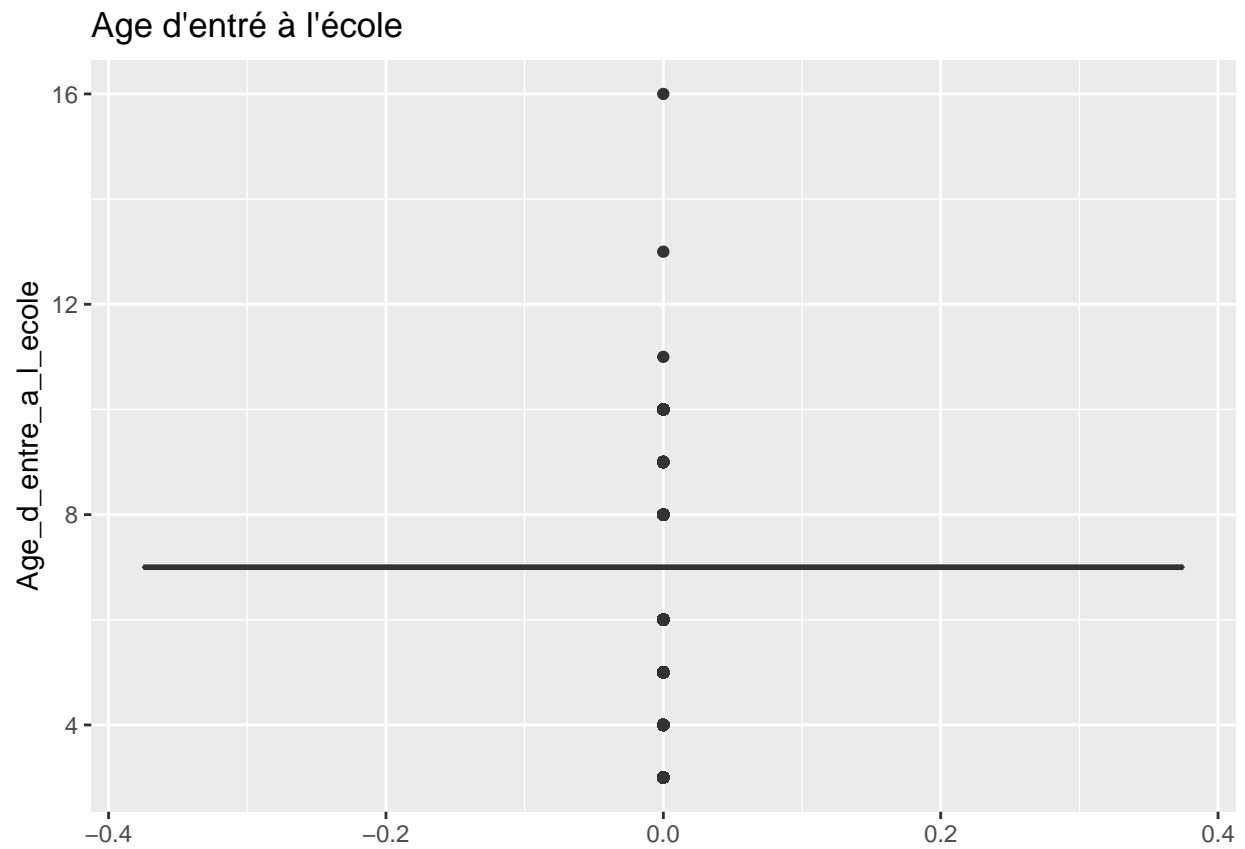
- Age et âge d'entrée à l'école

un premier contrôle à effectuer consiste à s'assurer que tous ce qui n'ont pas fait de formation ne renseignent l'âge d'entrée à l'école.

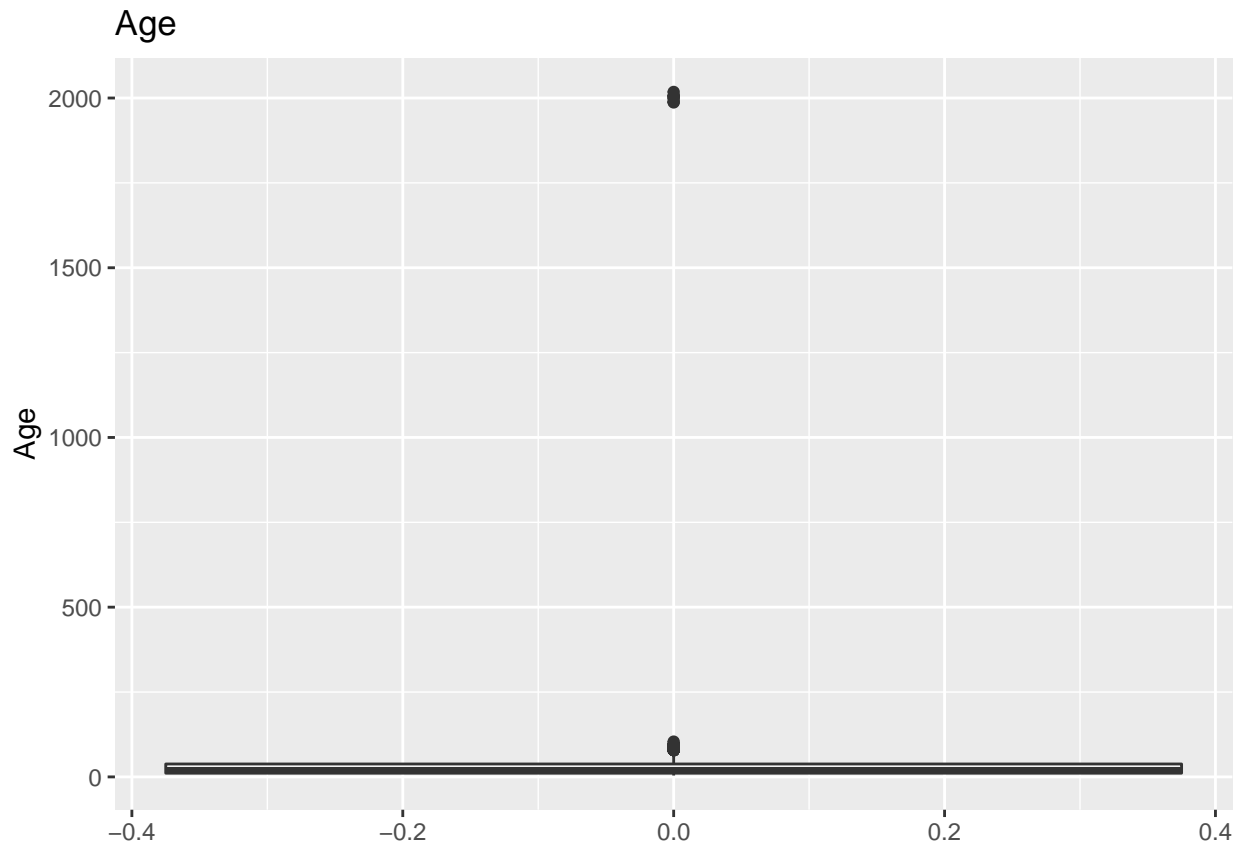
```
## ceux qui n'ont fait aucune formation ne doivent renseigner cette partie
y <- which(basef["Fait_une_ecole_formelles"] == "Non" & basef["Fait_une_ecole_non_formelle"] == "Non")
basef[y,"Age_d_entree_a_l_ecole"] <- NA
```

Pour le second contrôle, consiste à voir s'il y a des outliers. A cet effet, une méthode classiquement employée pour détecter les outliers (ou donnée aberrante), consiste à réaliser un boxplot.

```
library(ggplot2)
ggplot(basef, aes( y=Age_d_entree_a_l_ecole, fill=Age_d_entree_a_l_ecole)) +
  geom_boxplot()+
  ggtitle("Age d'entré à l'école")
```



```
## on fait pareil pour l'age
ggplot(my_data4, aes( y=Age,fill=Age)) +
  geom_boxplot()+
  ggtitle("Age")
```

Sur cette visualisation des données, les outliers sont représentés sous forme de points. Ils correspondent à des observations dont les valeurs sont :

- supérieures à la valeur du 3ème quartile plus 1.5 fois l'intervalle inter-quartile,
- ou inférieures à la valeur du 1er quartile moins 1.5 fois l'intervalle inter-quartile.

Pour récupérer ces valeurs aberrantes, on utilise la fonction `boxplot.stats` et avec la fonction `which` que nous avons déjà beaucoup utilisé, nous pourrions recueillir les index de ces valeurs pour les remplacer par la moyenne.

```
### récupérer les valeurs aberrantes
outlier <- boxplot.stats(basef$Age_d_entre_a_l_ecole)$out
### récupérer les positions des outliers
outlier_ind <- which(basef$Age_d_entre_a_l_ecole %in% c(outlier))
### remplacer les outliers par la moyenne
basef[outlier_ind,"Age_d_entre_a_l_ecole"] <- mean(basef$Age_d_entre_a_l_ecole)

### récupérer les valeurs aberrantes
outlier <- boxplot.stats(basef$Age)$out
### récupérer les positions des outliers
outlier_ind <- which(basef$Age %in% c(outlier))
### remplacer les outliers par la médiane
basef[outlier_ind,"Age"] <- median(basef$Age)
```

- Raison d'abandon des classes et de non fréquentation

Il faut s'assurer pour ces rebrriques que les reponses "c'est une fille" ou "Grossesse" correspondent bien à des individus de sexe féminin.

```
## Raison d'abandon des classes et de non fréquentation d'école formelle
## cherchons les hommes qui ont pour raison d'abandon des classe ou de non fréquentation "c'est une fille"
## remplaçons la raison par NA
### raison de non fréquentation
x1 <- which(basef["Sexe"] == "Masculin" & basef["Raisons_de_non_frequentation_d_ecole_formelle"] == "C'est une fille")
basef[x1,"Raisons_de_non_frequentation_d_ecole_formelle"] <- NA
### raisons d'abandon; on gèrera aussi les cas des homme ayant pour raison "grossesse"
x <- which(basef["Sexe"] == "Masculin" & basef["Raisons_d_abandon_de_l_ecole"] == "C'est une fille")
x2 <- which(basef["Sexe"] == "Masculin" & basef["Raisons_d_abandon_de_l_ecole"] == "Grossesse")
basef[x,"Raisons_d_abandon_de_l_ecole"] <- NA
basef[x2,"Raisons_d_abandon_de_l_ecole"] <- NA
```

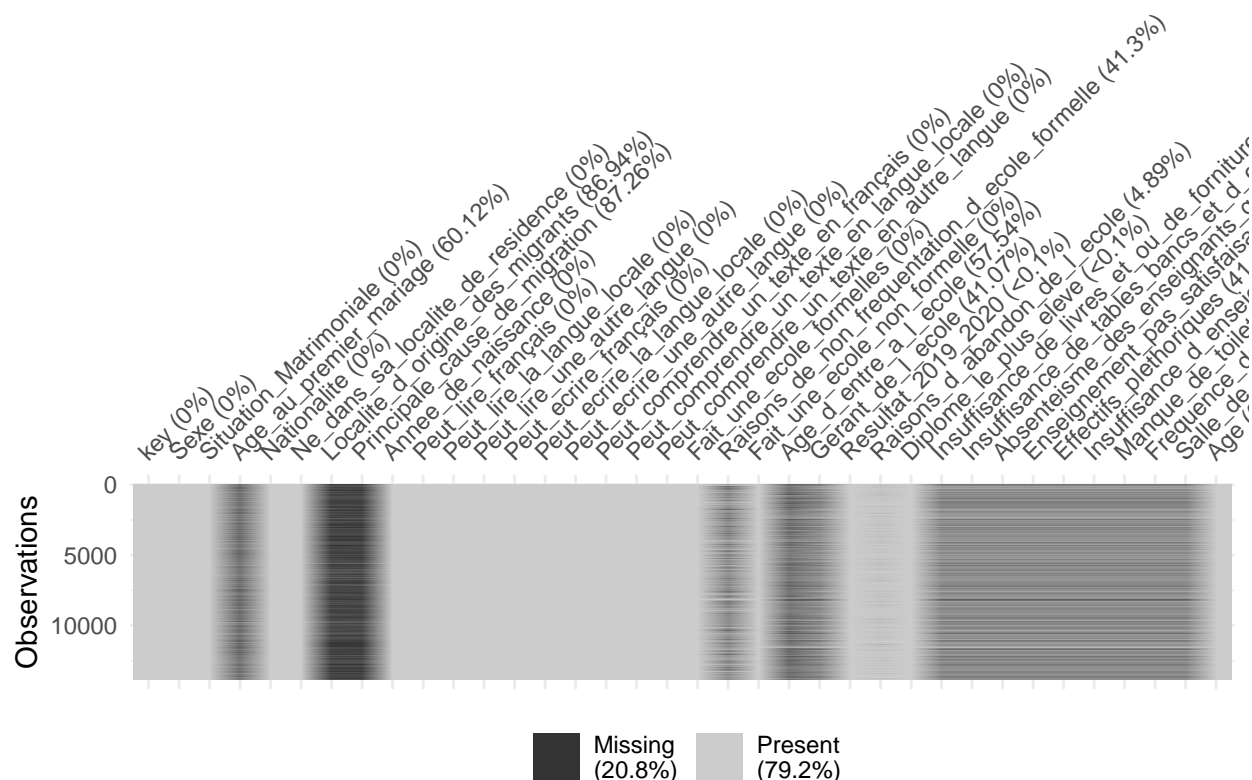
- Gérant de l'école et difficultés rencontrés dans le cadre de la formation

Cette question s'adresse aux individus qui vont actuellement à l'école ainsi un répondant qui ne va ni dans une école formelle ni dans une école informelle ne devrait y répondre. Nous remplacerons ces réponses par des NA. Pareil pour les difficultés rencontrés dans le cadre de la formation.

```
## Gérant de l'école: la réponse de tous ceux qui ne font aucune école doit être NA pareil pour les difficultés
var <- c("Gerant_de_l_ecole",difficulties)
for(x in var){
  y <- which(basef["Fait_une_ecole_formelles"] == "Non" & basef["Fait_une_ecole_non_formelle"] == "Non")
  basef[y,x] <- NA
}
```

Puisque nous avons introduit des NA dernièrement, visualisons les une dernière fois puis exportons notre table de données d'étude finale.

```
# visualisation des NA sur nos données finales
base_finale <- basef
vis_miss(base_finale)
```



```
# sauvegardons la base ainsi obtenue dans un fichier dta et un autre csv.
# mais avant il faut changer les noms qui sont trop longs pour pouvoir les exporter en format .dta
name <- c("Peut_comprendre_un_texte_en_français", "Peut_comprendre_un_texte_en_langue_locale",
          "Peut_comprendre_un_texte_en_autre_langue", "Raisons_de_non_frequentation_d_ecole_formelle",
          "Insuffisance_de_livres_et_ou_de_fournitures", "Insuffisance_de_tables_bancs_et_d_equipements",
          "Absentisme_des_enseignants_greve")
new_names <- c("comprendre_un_texte_en_français", "comprendre_texte_langue_locale",
               "comprendre_texte_en_autre_langue", "Raisons_non_frequentation_EF",
               "Insuffisance_livres_fournitures", "Insuffisance_tablesb_equipements",
               "Absentisme_enseignants_greve")

for(i in 1:length(name)){
  names(base_finale)[names(base_finale) == name[i]] = new_names[i]
}

write.csv(base_finale, paste(path, "Base_projetf.csv", sep = "\\"))
write_dta(base_finale, paste(path, "Base_projetf.dta", sep = "\\"))
```

La première partie prend ainsi fin ! Passons à l'analyse de nos données.