

```
1 #####
2 #
3 # Strategies de recherche en sciences sociales
4 #
5 # Introduction au logiciel R - Partie I
6 #
7 # Professeur: Dominique Joye
8 # Assistants: Francesco Laganà, Julien Chevillard, Julie Falcon
9 #
10 # Index
11 # Introduction informatique:
12 #
13 # Introduction a R
14 # Operation elementaires avec R
15 # Management des fichiers avec R
16 #
17 #
18 #
19
20 # Ou telecharger R?
21 #
22 # Pour Windows et Mac: http://stat.ethz.ch/CRAN/
23 #
24 # Pour Linux (Ubuntu): dans le terminal tapez: sudo apt-get install r-base
25 #
26 # PRECISATION: R64 bit soit R32bit?
27 #
28 # Vous pouvez installer le R64 sur des ordinateurs dont le processeur
29 # supporte le 64bit.
30 # Autrement vous devriez installer le R32 bit.
31 # Pour le travail qu'on fera dans le seminaire on utilisera R32
32 #
33 # Quand on ecrit une synthaxe la symbole "#" represente le "comment": R
34 # interpret tout ce qui suit comme n'etant pas une commande.
35
36 # Seance organisee a l'aide de syntaxes afin que vous puissiez
37 # tout faire ou reproduire vous-memes.
38 # Les syntaxes repondent en effet au principe de transparence d'une
39 # demarche scientifique : tout est evaluable et reproductible
40
41 # Travailler avec des synatxes vous permettra par ailleurs de reproduire
42 # facilement ce que vous avez fait, ou de modifier vos analyses rapidement
43 #
44 # R comme outil de la production de science est public et son code
45 # est libre.
46 # Chacun peut lire et modifier le code pour les analyses.
47 # Il existe differents bulletins:
48 #
49 # http://journal.r-project.org/
50 #
51 # http://www.jstatsoft.org/
52 #
53 # Info sur R
54 #
55 # http://www.r-project.org/
56 #
57 # http://www.statmethods.net/
58 #
59 # Une importante caracteristique de R est que chacun peut contribuer au
60 # projet avec
61 # des extensions, qui permettent d'effectuer differents types d'analyse.
```

```

62 #
63 #
64 # Installer des extensions:
65 #
66 # install.packages("nom_de_l'extension", dep=T)
67
68 install.packages("rgrs", dep=T)
69
70 # Charger les extensions
71 # (Il faut faire ca avant toute analyse).
72 # Faire reference a cette page web:
73 # http://www.statmethods.net/interface/packages.html
74 #
75 # commande (exemple): library(extension)
76
77 library(foreign)
78 library(rgrs)
79 library(Hmisc)
80
81 #
82 # R est une collection d'extensions qui permettent de faire de
83 # l'analyse statistique.
84 # Le premier pas consiste a assigner a une variable une certaine valeur:
85 # Cela est permis par le signe d'assignation "<-"
86
87 # Dans le cas d'un scalaire:
88
89 x <- 4
90
91 # ou dans le cas d'un vecteur:
92
93 y.vector <- c(2,6,8,10,20)
94
95 # Position d'elements dans un vecteur
96 # Les paranteses carrees donnent l'element de position 3 dans le vecteur
97
98 y.vector[3]
99
100
101 #####
102 # R comme calculateur
103 #####
104
105 # Vous pouvez utiliser R comme calculateur
106
107 4+5
108
109 # or
110
111 x <- 4
112
113 y <- 5
114
115 x+y
116
117 y.vec2 <- y.vector^2
118
119 logyvec <- log(y.vector)
120
121
122 # Combien d'element il y a dans notre vecteur?

```

```

123
124 length(y.vector)
125
126
127 # R peut travailler aussi avec des matrices:
128
129 m1<-matrix(1:12,nrow=3)
130
131 m1
132
133 #####
134 # Object et typologies de donnees
135 #####
136
137 # Differement des autres logiciels statistique R peut manager different
138 # types de donnees
139 # La chose la plus importante est de considerer que tout dans R est
140 # considere comme un objet
141
142 # IMPORTANT: Il existe des caracteres reserves qu'on ne peut pas utiliser
143 # dans R:
144
145 # 1) JAMAIS nommer une variable en commençant par un chiffre
146 # Exemple:
147
148 0a <- 56
149
150 # Mais le contraire (utiliser un chiffre qui ne figure pas en premiere
151 # position dans le nom de la variable ou de l'objet) est possible:
152
153 a0 <- 56
154
155 a0
156
157 # 2) En outre il y a des caracteres reserves:
158
159 # FALSE TRUE Inf NA NaN NULL
160 # break else for function if in next repeat while
161
162 # 3) Il n'est pas conseille de nommer des objets avec le nom d'une fonction:
163 # Exemples: Eviter de nommer vos objet: 1) c (parce que dans le langage R
164 # cela signifie "combine")
165 #
166 #
167 #
168 #
169 #
170 # Types d'objets dans R:
171
172 # character
173
174 str <- c("Dominique","Julie","Julien","Francesco")
175
176 str
177
178 is.character(str)
179
180 # numeric: les elements appartiennent aux nombres Reels
181 # Ils representent une quantite sur un espace continu: Ex. Revenu
182
183 sqz <- seq(1.575, 5.125, by=0.05)

```

```

184
185 sqz
186
187 is.numeric(sqz)
188
189 # integer
190 # Il representent une quantite sur un espace discret (nombres entiers): Ex.
191 # Household size
192
193 sqint <- seq(1,8, by=1)
194
195 sqint
196
197 sqint <- as.integer(sqint)
198
199 # logical
200 # Il s'agit d'operateurs booleens : ils peuvent assumer les valeurs VRAI ou
201 # FAUX
202
203 log <- c(TRUE, FALSE, TRUE, TRUE, TRUE, FALSE)
204
205 # complex, nombres complexe
206 # Ils appartiennent pas aux nombres reels.
207
208 cmp <- complex(real=1:10, imaginary=-1:9)
209
210 ##### Differents objets peuvent etre regroupes dans une liste:
211
212 elle <- list(STR=str, SQZ = sqz, SQINT= sqint, BOOL=log, CMP=cmp)
213 elle
214
215 #####
216 # Managing your analyses:
217 #
218 # Setting your working directory / Assigner votre dossier/emplacement de
219 # travail
220 #
221 # IMPORTANT
222 # Sur ce site vous pourrez trouver beaucoup de documentation:
223 # http://www.statmethods.net/interface/workspace.html
224
225 setwd("/Users/yourlogin/Desktop")
226
227 # knowing your directory:
228
229 getwd()
230
231 # Which files are stored in your working directory? / quels sont
232 # les fichiers figurant dans votre dossier de travail ?
233
234 dir()
235
236 # Il est parfois utile de sauvegarder le contenu de vos analyses
237 # Pour ce faire, il faut creer un fichier log de vos analyses
238 # On fait ca dans R avec la commande sink()
239
240
241 # output directed to output.txt in c:\projects directory.
242 # output overwrites existing file. no output to terminal.
243 # Attention, la commande ci-dessous va effacer le document appele
244 # "Seance6.txt"

```

```

245 # veuillez donc a changer le nom du fichier dans la commande si vous
246 # souhaitez conserver le fichier original
247 # sink("Seance6.txt")
248 # output directed to myfile.txt in cwd. output is appended
249 # to existing file. output also send to terminal.
250 # vous pouvez ajouter comme ci-dessous "append" la commande initiale
251 # Dans ce cas, le fichier original ne sera pas efface, mais les
252 # modifications et nouvelles commandes/analyses que vous avez effectuees y
253 # seront ajoutees
254
255 # sink("Seance6.txt", append=TRUE, split=TRUE)
256
257 # Knowing the past history of your command:
258
259 history(max.show=100)
260
261 # R est case sensitive: majuscules et minuscules font la difference:
262
263 x <- 46
264
265 x
266
267 X
268
269 #####
270 #
271 # DATA FRAMES:
272 #
273 #####
274
275 # Le dataframe est l'objet de R que nous interesse le plus:
276 # Il est comme une matrice mais, a la difference d'une matrice
277 # Il accepte aussi des valeurs non numeriques.
278 # Par consequent chaque matrice peut etre un dataframe mais
279 # un dataframe n'est pas forcement une matrice.
280
281 # En R on convertit une matrice ou un objet en dataframe avec
282 # la commande: data.frame
283
284 mdat <- data.frame(m1)
285
286 is.data.frame(m1)
287
288 is.data.frame(mdat)
289
290
291 # Il y a differents moyens de lire les dataframes dans R:
292 # 1) Manuellement avec la fonction scan()
293
294 x<- scan()
295 1 2 4 5 6 7 9 10 11 23 9 32 4 5 34 23
296
297 x
298
299 # 2) Lire un fichier deja organise comme data frame:
300
301
302 # Depuis un format proprietare:
303 # La livrairie foreign permet de lire des fichier crees avec d'autres
304 # logiciels
305 # comme SPSS, SAS, Stata

```

```

306 # Pour cela vous devriez utiliser l'Extension foreign
307 #
308
309 # Nous allons travailler avec un format de type text file / txt, les donnees
310 sont separees
311 # par une virgule (*.csv pour comma separated value :fichier ou les donnees
312 sont separees par des virgules ).
313 # Pour l'analyser, vous devez pouvoir lire ce type de fichier dans R
314 # Pour cela, nous utiliserons la fonction read.table()
315 #
316
317 ess08csv <- read.table("ess08_Methode_31032011.csv", sep=",", header=T)
318
319 # Explication de la commande:
320 #     1) ess08csv represente le nouvel objet ou vos donnees
321 #     sont enregistrees
322 #     2) entre parentheses, vous avez le fichier .csv ou .dat que vous
323 #     voulez lire ("ess08_Methode_31032011.csv")
324 #     3) l'option "header" dit a R que la premiere ligne contient les noms
325 #     des variables
326 #     4) l'option "sep" dit a R que les colonnes (donc les variables) sont
327 #     separees par une virgule.
328
329 d <- ess08csv
330
331
332 # Pour savoir quels objets sont presents dans votre workspace il faut
333 # utiliser la fonction ls()
334
335 ls()
336
337 #####
338 #
339 # Trois facons de sauver votre output
340 # Sauver le workspace
341
342 save(d, file="ess08_Methode19Apr2011.rda")
343
344 # Sauver l'image de votre travail:
345 save.image(file = "Methode_Objjet31032011.RData")
346
347 # Sauver vos données
348 write.table(d, file="ess08_Methode_31032011.csv", sep=",")
349
350 #
351 # Obtenir de l'aide:
352 # la commande "help(fonction)" nous permet de demander a R de l'aide sur une
353 fonction:
354
355 help(read.table)
356
357 # Ou de facon equivalente:
358
359 ?read.table
360
361
362
363
364

```

```

1 #####
2 #
3 # Strategies de recherche en sciences sociales
4 #
5 # Introduction au logiciel R - Partie II
6 #
7 # Professeur: Dominique Joye
8 # Assistants: Francesco Laganà, Julien Chevillard, Julie Falcon
9 #
10 # 1) Exploration de la base des données
11 # 2) Tableaux des frequences
12 # 3) Manipulation des données: Selectionner des observations
13 # et des variables et Types de recodage
14 # 4) Analyzing relations between two variables:
15 # 4.1) Variables nominales;
16 # 4.2) Variable interval ~ variable nominale
17 # 4.3) Variable interval ~ variable interval
18 # 4.4) Variable interval ~ variable nominale + variable interval
19 # 5) Appendix :
20 # 5.1) Elementary database management
21 # 5.2) Missing data recoding
22 # 5.3) Selectionner variables par indexation
23 #####
24
25
26 # Loading libraries into R:
27 library(rgrs)
28 library(Hmisc)
29 library(car)
30 library(psych)
31
32 #####
33 # 1) Exploration de la base des données
34 #####
35
36 # Loading data into R
37 load("ess08_Methode7Apr2011.Rda")
38
39 # Exploration du databse:
40
41 # Which are the variables you have?
42
43 names(d)
44
45 # On accede aux descriptives de la variable avec
46
47 describe(d$txearn)
48
49
50 # What kind of variables you have? factors? Numeric?
51 # La fonction str() nous permet de lire savoir le nom des variables et
52 # leur type: facteur, numerique, integer, etc...
53
54 str(d, list.len=200)
55
56
57 # Dans chaque objet il y a different elements. Exemple: dans un data frame il
58 # y a differents variables, comme le genre, le revenu, l'occupation etc...
59 # Pour se referer a une variable presente dans le dataframe d'interet on
60 # utilise
61 # le signe "Dollar" ($)

```

```

62
63 d$gndr
64
65 # si vous n'indiquez pas le dataframe auquel la variable appartient vous
66 # obtiendrez un message d'erreur
67
68 gndr
69
70 #####
71 # 2) Tableaux des frequences
72 #
73 # Variables facteurs (Nominales ou Ordinales)
74 #
75 #####
76
77 table(d$gndr)
78
79 table(d$txearn)
80
81 # Pourcentages: la fonction prop.table:
82
83 prop.table(table(d$gndr))
84
85 # Inspection visuelle des frequences
86 # Pour des variable nominales on peut utiliser le graphic a bars
87 # ou encore le dotplot.
88
89 # Valeurs absolus
90 taprt <- table(d$prtclbch)
91 dotchart(taprt)
92
93 # Pourcentage
94 dotchart(prop.table(table(d$prtclbch))*100, xlim=c(1,100))
95
96 plot(table(d$mar.reg))
97
98 barplot(prop.table(table(d$mar.reg))*100, ylim=c(0,100))
99
100 # Or la meme chose:
101
102 tamar <- prop.table(table(d$mar.reg))*100
103 barplot(tamar, ylim=c(0,100))
104
105 # Numeric variables
106 min(d$gvjbevn.num, na.rm=T) # La plus petite valeur de la variable
107 max(d$gvjbevn.num, na.rm=T) # La plus grande valeur de la variable
108 mean(d$gvjbevn.num, na.rm=T) # La moyenne de la variable
109 sd(d$gvjbevn.num, na.rm=T) # L'ecart type
110 var(d$gvjbevn.num, na.rm=T) # La variance
111 summary(d$gvjbevn.num) # Minimum; 1st Quartile; Mediane, Moyenne, 3rd
112 Quartile, Maximum, Valeurs manquantes
113
114 # Inspection visuelle des variables interval:
115
116 hist(d$gvjbevn.num)
117 abline(v=mean(d$gvjbevn.num, na.rm=T), col="red", lwd=1)
118 abline(v=median(d$gvjbevn.num, na.rm=T), col="red", lwd=1)
119
120 # Approfondissement: Pour verifier si la distribution de la variable est une
121 # distribution normale
122

```



```

123 # hist(d$gvjbevn.num, probability=T)
124 # lines(density(d$gvjbevn.num, na.rm=T), lwd=2, col="red")
125
126
127 #####
128 # 3) Manipulation des données: Selectionner des observations #
129 # et des variables et Types de recodage #
130 #####
131
132 # Renommer une variable
133
134 d$txe <- d$txearn
135
136 # Dans R il y a deux moyens de selectionner observations et variables:
137 # 1) indexation directe (Barnier "R pour sociologues", p. 43-51)
138 # [[[[VOIR APPENDIX]]]].
139
140 # 2) Fonction subset:
141 # Le premier argument indique le dataframe duquel il faut "subsetter";
142 # Le deuxième argument indique la variable et la modalité par laquelle
143 # il est possible de selectionner les obseravtions (gndr== "Male")
144 # le troisième argument indique les variables qu'on veut selectionner
145 # (ne pas indiquer si on veut selectionner toutes les variables)
146
147 # Assignation des conditions:
148 # == egale à
149 # != different de
150 # > strictement superieur à
151 # < strinctement inferieur à
152 # >= superieur ou egal à
153 # <= inferieur ou egal à
154
155 # Ex. 1 selectionner des reponse des hommes en ce qui concerne le vote,
156 # la'acitivté principale
157 # et l'opinion sur la taxation.
158
159 dh <- subset(d, gndr== "Male", select=c("ident", "vote","mnactic","txearn"))
160
161 # Ex. 2: selectionner la sous-population des femmes.
162
163 dfe <- subset(d, gndr== "Female")
164
165 # Ex. 3 On peut aussi selectionner observations par deux ou plusieurs
166 # conditions, par exemple les homme qui peuvent voter (eligible au vote):
167
168 dhvo <- subset(d, gndr== "Male" & vote != "Not eligible to vote",
169 select=c("ident", "vote","prtclbch","mnactic","txearn"))
170
171 table(dhvo$vote)
172
173 table(dh$vote)
174
175 #####
176 # Recoding variables: #
177 #####
178
179 # Ex. Recoder Main activity
180
181 # Creer une nouvelle variable en regroupant les categories
182
183 table(d$mnactic)

```

```

184 d$mac.r <- as.character(d$mnactic)
185 d$mac.r[d$mnactic=="Paid work" ] <- "Employment"
186 d$mac.r[d$mnactic=="Unemployed, looking for job" |
187 d$mnactic=="Unemployed, not looking for job" ] <- "Unemployment"
188 d$mac.r[d$mnactic=="Education" |
189 d$mnactic=="Permanently sick or disabled" |
190 d$mnactic=="Retired" |
191 d$mnactic=="Housework, looking after children, others" |
192 d$mnactic=="Other"] <- "Not active"
193
194 d$mac.r <- factor(d$mac.r)
195
196 # Check:
197 table( d$mac.r,d$mnactic)
198
199
200 # Découper une variable interval en classes:
201 d$age.rec <- cut(d$agea, c(15, 24, 34, 44, 54, 64, 100 ), include.lowest=T,
202 labels=c("15-24", "25-34", "35-44", "45-54", "55-64", "65>"))
203
204 table(d$age.rec)
205
206
207 # Croiser variables:
208 # Variable factor (nominale ou ordinale) + Variable factor (nominale ou
209 # ordinale)
210 # Ex. gender et employment status:
211
212 # Generer une nouvelle variable où memoriser les resultats:
213 d$gen.act <- NA
214
215 # Generer les consitions:
216 d$gen.act[d$mac.r=="Employment" & d$gndr=="Female" ] <- "Employed Women"
217 d$gen.act[d$mac.r=="Employment" & d$gndr=="Male" ] <- "Employed Male"
218
219 d$gen.act[d$mac.r=="Unemployment" & d$gndr=="Female" ] <- "Unemployment
220 Women"
221 d$gen.act[d$mac.r=="Unemployment" & d$gndr=="Male" ] <- "Unemployment Male"
222
223 d$gen.act[d$mac.r=="Not active" & d$gndr=="Female" ] <- "Not active Women"
224 d$gen.act[d$mac.r=="Not active" & d$gndr=="Male" ] <- "Not active Male"
225
226 d$gen.act <- factor(d$gen.act)
227
228
229 # Variable categorielle + variable interval
230 # Ex. Age et gender:
231 d$age.gndr <- NA
232
233 # Generer les consitions:
234 d$age.gndr[d$agea<=45 & d$gndr=="Female" ] <- "Women <= 45"
235 d$age.gndr[d$agea>45 & d$gndr=="Female" ] <- "Women > 45"
236
237 d$age.gndr[d$agea<=45 & d$gndr=="Male" ] <- "Male <= 45"
238 d$age.gndr[d$agea>45 & d$gndr=="Male" ] <- "Male > 45"
239
240 table(d$age.gndr)
241
242 # Transformation des variables:
243 # Ex. cstandardiser les variables
244 d$z.gvjbevn.num <- (d$gvjbevn.num-mean(d$gvjbevn.num,

```

```

245 na.rm=T))/sd(d$gvjbevn.num, na.rm=T)
246
247 # De cette façon la nouvelle variable est calculé en unités d'ecart type
248 # (ecart type=1) et centrée sur 0 (mean=0)
249
250 mean(d$z.gvjbevn.num, na.rm=T)
251
252 sd(d$z.gvjbevn.num, na.rm=T)
253
254 # Variable interval + variable interval: Computing scores
255
256 # 1. Extraire les variables dont on veut estimer le score:
257 pro <- d[c("gvjbevn.num", "gvhlthc.num", "gvslvol.num" , "gvslvue.num" ,
258 "gvcldcr.num" ,"gvpdlwk.num")]
259
260 # 2. Calculer la moyenne:
261 d$state.resp <- rowMeans(pro, na.rm=T)
262
263 #####
264 # Reliability analysis: #
265 # When we want to compute the score for a given set of variables #
266 # it is possible to assess the reliability of the scale that #
267 # we constructed #
268 # The most used test is the Cronach's Alpha. #
269 # #
270 # 1) Le Cronbach's Alphas peut assumer des valeurs compris parmi 0 et 1 et #
271 # nous donne des indications sul le degre de inter-correlation des #
272 # differentes items qu'on veut examiner. #
273 # #
274 # 2) Si les items sont très correlés parmi eux (consistence interne #
275 # maximale), il est possible de conclure que chaque item donne une #
276 # contribution réelle à la construction de la mesure et que, #
277 # dans l'ensemble tous les items se réfèrent au meme concept. #
278 # #
279 # 3) Si les items ne sont pas très correlés parmi eux, ça signifie que ceux#
280 # là ne constituent pas une mesure satisfaisante du meme construct. #
281 # #
282 # Exemple: si une echelle sur l'anxieté alpha a une valeur de 0.80 #
283 # cela signifie que l'80% du score obtenu par le sujet est du à l'anxieté #
284 # Si la consistence interne de l'echelle est de 0.40, cela signifie que le #
285 # 40% de l'echelle mesure l'anxieté et le restant 60% mesure autre chose. #
286 # #
287 # Valeurs conventionels du alpha de cronbach. #
288 # 0.60 inacceptable; #
289 # 0.60-0.65 pas desirable; #
290 # 0.65-0.70 Acceptable; #
291 # 0.70-0.80 Bon; #
292 # 0.80 > Optimale #
293 # #
294 # Dans R on peut calculer le alpha de Cronbach par la fonction alpha #
295 # contenue dans l'extension psych #
296
297 alpha(pro)
298
299
300
301
302
303
304
305

```

```

306 #####
307 # 4. Analyzing relations between two variables:
308 #
309 # De l'analyse d'association aux modèles de regression:
310 #
311 # Different types d'analyse en fonction des types de variables à analyser:
312 #
313 # (1) Variables categorielles
314 # Tableaux croisées + Test du Chi carré + V de Cramer
315 #
316 # (2) Variables variables interval + variable categorielle:
317 # T test (deux groupes) ou Analysis de la variance (quand il s'agit de
318 # plusieurs groupes)
319 #
320 # (3) Variable interval + variable interval
321 # Analysis de la correlation ou regression
322 #
323 #####
324 #####
325 #####
326 #
327 # 4.1) Deux variables categorielles (facteurs)
328 # Tableaux croisées + Test du Chi carré + V de Cramer
329 #
330 #
331 # Valeurs Absolus:
332
333 tab <- table(d$txearn, d$gnldr)
334
335 tab
336
337 # Pourcentages ligne:
338 prop.table(tab, 1)
339
340 # Pourcentages colonne:
341 prop.table(tab, 2)
342
343 # Pourcentages case:
344 prop.table(tab)
345
346 # Quand'il s'agit de calculer les pourcentages pour ligne ou colonne,
347 # jamais oublier la regle de Galtung: il faut percentualiser en relation
348 # au valeurs de celle qui est considerée comme variable independente
349 # dans le modèle de base. (Galtung, 1970)
350
351
352 # Tester l'association parmi deux variables categorielles:
353
354 # Comment tester l'hypothese d'association parmi deux
355 # variables categorielles?
356
357 # Association summary(table) -> chi square
358 # Est-ce que la relation parmi genre et taxation est significative?
359
360 # On repond à cette question par le test du chi carré:
361 # Le test du chi carré teste l'hypothese nulle d'absence de relation
362 # parmi deux variables.
363 # Le chi carré il est estimé par la confrontation parmi le tableau
364 # crée et un tableau hypothetique qui represente la distribution
365 # conjointe des deux variables en cas d'independance.
366

```

```

367 # Les deux commandes donnent des coefficients d'associations parmi deux
368 # variables (Chi square et V de Cramer);
369
370 summary(tab)
371
372 # Le Chi square nous dit qu'il y a de l'association parmi deux
373 # variable et que l'association n'est pas random.
374 #
375 # On peut se demander de la force de l'association:
376 # pour repondre à cette question on utilise le
377 # V de cramer. Il est calculé comme suit:
378 #  $V = \text{SQRT}(\text{chi}2 / (n(k - 1)))$ 
379
380 # Où is chi2 est le chi-square and k is the number of rows or columns in the
381 # table.
382 # Cramer's V varies between 0 and 1. Close to 0 it shows little association
383 # between variables. Close to 1, it
384 # indicates a strong association.
385 # Where the table is 2 x 2, use Phi.
386
387 assocstats(tab)
388
389 # Une autre façon d'obtenir les statistiques du chi2 est par la
390 # commande chisq.test
391
392 chisq.test(tab)
393
394 # On peut aussi avoir un tableau avec les
395 # expected frequences en cas d'independance (sous l'hypothese nulle):
396
397 chisq.test(tab)$expected
398
399 # Et les residuals '(observed - expected)/sqrt(expected)'.:
400
401 chisq.test(tab)$residuals
402
403 # NOTE about the chi2.
404 # 1) It is valid only in case of random samples
405 # 2) Sample size: If a chi square test is conducted on a sample with a
406 # smaller size,
407 # then the chi square test will yield # an inaccurate inference. The
408 # researcher, by using chi square test on small samples,
409 # might end up committing a Type II # error.
410 # 3) Expected Cell Count - Adequate expected cell counts. Some require 5 or
411 # more, and others require 10 or more.
412 # A common rule is 5 or more in all cells of a 2-by-2 table, and 5 or more in
413 # 80% of cells in larger tables,
414 # but no cells with zero expected count. When this assumption is not met,
415 # Yates' correction is applied.
416 # 4) Independence - The observations are always assumed to be independent of
417 # each other.
418 # This means chi-square cannot be used to test correlated data (like: matched
419 # pairs, panel data). In those cases you
420 # might want to turn to McNemar's test.
421 #
422 # In case you have small sample sizes use fisher exact test:
423
424 fisher.test(tab)
425
426
427

```

```

428 #####
429 #
430 # 4.2 Variable interval + variable categorielle:
431 # ANALYSE DE LA VARIANCE:
432 #
433 #####
434
435 # L'analyse de la variance est utilisée pour etudier la relation parmi
436 # une variable (dependante) interval + une variable (independente)
437 # categorielle
438 #
439 # L'Anova permet de comparer la moyenne d'une variable interval en fonction
440 # des categories d'une variable nominale ou ordinale
441 # Tapply estime les moyennes par differents groupes:
442 # Est-ce qu'il y a des differences parmi employés, non-actifs et chomeurs
443 # dans le soutien des actions de redistribution gouvernementales
444 # vers les chomeurs?
445
446 tapply(d$gvslvue, d$mac.r, mean, na.rm=T)
447
448 # To compare two groups you can use t.test
449 t.test(d$gvslvue ~ d$gndr)
450
451 # NOTE: note that the t.test requires normality in the distribution of
452 # the two groups and equal variance (faire reference au texte de Barnier pour
453 # approfondissements.
454 # Les deux tests fondamentaux sont:
455 # 1) shapiro-wilk normality test (fonction: shapiro.test())
456 # 2) Egalité des variances:(fonction var.test())
457 # Au cas ou ces deux assumptions ne sont par rencontrées
458 # il faut utiliser des test non-parametriques:
459 # wilcox.test()
460
461
462 # To compare the significativity in differences of more than two groups: uses
463 # lm() function and anova() functions.
464 # Such functoins allows for an extension of a two-sample t test, testing the
465 # equality of means of more than two groups.
466
467 # Ex. Est-ce que les differences parmi employés, non-actifs et chomeurs
468 # dans le soutien des actions de redistribution gouvernementales
469 # vers les chomeurs sont significatives?
470
471 mod1 <- lm(d$gvslvue ~ d$mac.r)
472
473 summary(mod1)
474
475 anova(mod1)
476
477 # Inspection visuelle de variables interval + categorielle:
478 boxplot(d$gvjbevn.num~d$mac.r)
479
480 # Dans l'Anova le F represente l'equivalent du t dans le t test.
481 # Lien mathematique:  $F=t^2$ 
482
483 # NOTE:
484 # Il y a des assumptions qui doivent etre satisfaites pour l'ANOVA:
485 # Notamment:
486 # 1) Independance des valeurs observés.
487 # 2) Normalité de la distribution de notre variable dependente:
488 # 3) Homogeneité des variances:

```

```

489
490 # La deuxième assumption peut etre évalué de deux façons:
491 # Le test de normalité di Kolmogorov-Sirnov:
492
493 ks.test(d$gvslvue.num, pnorm, mean(d$gvslvue.num, na.rm=T),
494 sd(d$gvslvue.num,na.rm=T))
495
496 # De l'output on doit lire le p-value. Cela doit etre egale à 0.05 ou plus
497 # grand: l'hypothese nulle dit que la distribution est normale.
498 # le ks.test il est sensible aux valeur dupliqués.
499 # Une alternative est le shapiro.test
500
501 shapiro.test(d$gvslvue.num)
502
503 # Egalement au ks.test on doit regarder le p-value. L'ihypothese nulle
504 # dit que notre variable dependente estdistribué selon une distribution
505 # normale.
506
507 # En ce qui concerne l'homogeneité des variances
508
509 bartlett.test(d$gvslvue ~ d$mac.r)
510
511 # De l'output on doit lire le p-value. Il doit etre majeure ou egale
512 # à 0.05. L'hypothese nulle soutient que les variances sont homogenes.
513 # Au cas où cette assumption ne soit pas respectée, il faut calculer l'Anova
514 # avec la correction de Welch.
515
516 #####
517 #
518 # 4.3) Mesures d'association entre deux variables interval:
519 # Coefficient de CORRELATION et Modele de REGRESSION
520 #
521 #####
522
523 # Quand les deux variables sont interval on peut utiliser
524 # analyser la correlation et la regression.
525 # Le coefficient de correlation de Pearson (Pearson's
526 # product-moment correlation)
527 # il est calculé comme le rapport parmi la covariance parmi deux variables (y
528 # and x) et le product de leurs ecart types:
529 #
530 #  $\rho(x,y) = \text{cov}(x,y) / (\text{sigma}(x) * \text{sigma}(y))$ 
531 #
532 # où  $\text{cov}(x,y) = \text{sum}((x-M(x)) * (y-M(y))) / (n-1)$ 
533
534 cor(d$gvjbevn.num, d$iseiego.num, use="complete.obs")
535
536 # Le coefficient de correlation est egale à 0 en cas d'absence de
537 # relation
538 # Il est egale à 1 en cas de très forte association positive
539 # Il est egale à -1 en cas de très forte association negative
540
541 # On peut tester l'hypothese nulle d'absence de relation parmi les deux
542 # variables
543
544 cor.test(d$gvjbevn.num, d$iseiego.num, use="complete.obs")
545
546 # ou encore:
547
548 cor(d$gvjbevn.num, d$gvhlthc.num, use="complete.obs")
549

```

```

550 cor.test(d$gvjbevn.num, d$gvhlthc.num, use="complete.obs")
551
552 # NOTE:
553 # Si vous voulez analyser la corrélation parmi différentes variables:
554 # Il faut regrouper ces variables dans un data frame comme on
555 # avait fait avant (ligne 256-258).
556
557 cor(pro, use="complete.obs")
558
559 # Si on trouve un coefficient de corrélation élevé
560 # on peut se poser la question de la meilleure
561 # ligne qui interpole nos données.
562 # Cela signifie de trouver la fonction qui lie
563 # deux phénomènes à travers la formulation
564 # d'un modèle du type:  $f(x) = a + b(x)$ 
565
566 # On répond à cette question par l'analyse de la ligne de
567 # régression:
568
569 # Quel est l'impact du Socioeconomic status sur le soutien au gouvernement?
570
571 mod1 <- lm(d$gvjbevn.num ~ d$iseiego.num)
572 mod1
573
574 # Le coefficient de régression représente le coefficient
575 # angulaire de la ligne. Il nous dit le rapport
576 # parmi variation de notre x et variation de y.
577 # Quel est la variation de y suite à variations de
578 # notre x.
579 # En termes statistiques il représente la partie de
580 # la covariance parmi x et y qui est expliquée par la
581 # variance de x.
582 #
583 #  $\beta = \text{cov}(x, y) / \text{var}(x)$ 
584
585 beta <- cov(d$gvjbevn.num, d$iseiego.num,
586            use="complete.obs") / var(d$iseiego.num, na.rm=T)
587
588 # Visualiser le plot et la ligne de régression
589 plot(jitter(d$iseiego.num, factor=3), jitter(d$gvjbevn.num, factor=2))
590 abline(lm(d$gvjbevn.num ~ d$iseiego.num), col="red")
591
592 # Autre exemple: Corrélation parmi Index of Socio-Economic Status de
593 # l'interviewé et celui de son père.
594 plot(jitter(d$iseifa.num, factor=2), jitter(d$iseiego.num, factor=3))
595 abline(lm(d$iseifa.num ~ d$iseiego.num), col="red")
596
597 #####
598 #
599 # Régression multiple:
600 # Analyses de régression (AnCoVa) :
601 #
602 #####
603
604 # Souvent on veut contrôler si une relation
605 # résiste à des contrôles ultérieures
606 # avec d'autres variables:
607 #
608 # Pour cette raison on utilise des techniques de régression multiple
609 # ou on contrôle notre relation parmi x et y par l'effet d'autres variables:

```



```

611
612 # Reprenons l'exemple du soutien pour les chomeurs (Job for everyone,
613 # governments' responsibility) et supposons que la relation depend de:
614
615 # Position sur l'echelle gauche droite
616 # Statut occupationnel
617 # Socio-economic status
618 # Variables de control:
619 #   Age (continue); Gender.
620
621 # Avant d'estimer une regression:
622
623 # Selectionner les variables qui font partie de votre modele:
624 dm <- d[c("state.resp", "gvjbevn.num", "iseiego.num", "agea.num" ,
625 "lrscale.num", "gndr", "mac.r", "age.rec")]
626
627 # 1) Analyse monovariée de toutes les variables pour detecter outliers...
628 [commande: table(), plot(), summary()]
629
630 # 2) Transformation eventuelle des variables (par exemple centrer par
631 # rapport)
632 # Centrer l'age par rapport à la moyenne (generalement pour toutes les
633 # variables dont le champ de variation
634 # ne comprend pas le zero - cela sert pour
635 # l'interpretation de l'intercepte.
636
637 dm$Cagea <- dm$agea-mean(dm$agea, na.rm=T)
638 dm$Cisei <- dm$iseiego.num-mean(dm$iseiego.num, na.rm=T)
639
640 # 3) Exploration des relations parmi variables [commande: corr()]
641
642 cor(dm[,1:6], use="complete.obs")
643
644 # 4) Creation des contrasts:
645 # Quand dans la regression il y a des variables quantitatives et qualitatives
646 # il faut s'assurer que ces dernières soient traitées comme factors:
647 #
648
649 is.factor(dm$mac.r)
650 is.factor(dm$gndr)
651
652
653 # Visualiser les contrasts:
654 contrasts(dm$mac.r)
655
656 # On peut changer la baseline des contrasts avec l'option suivante:
657
658 contrasts(dm$mac.r) <- contr.treatment(levels(dm$mac.r), base=1)
659
660 mod1 <- lm(d$gvjbevn.num ~ d$iseiego.num + agea.num, data=dm)
661
662 mod1 <- lm(gvjbevn.num ~ Cisei+ Cagea + mac.r+ lrscale.num + gndr, data=dm)
663
664 # Est-ce que l'effet de l'age est lineaire?
665 # On reponds à cette question defragmentant l'age (on utilise age.rec)
666
667 # Creer les contrasts:
668 contrasts(dm$age.rec)
669
670 mod1 <- lm(gvjbevn.num ~ iseiego.num + age.rec + mac.r+ lrscale.num + gndr,
671 data=dm)

```

```

672
673 summary(mod1)
674
675 # La regression lineaire est seulement un modele de regression qui fait
676 # partie d'une plus grande famille des Generalized Linear Models...
677 # ... mais ça n'est pas objet de notre cours...
678 #
679 # ... FIN
680
681 # A savoir:
682
683 # Trois façons de sauver votre output:
684
685 # Sauver le workspace
686 save(d, file="ess08_Methode19Apr2011.rda")
687
688 # Sauver l'image de votre travail:
689 save.image(file = "Methode_Objjet31032011.RData")
690
691 # Sauver vos données
692 write.table(d, file="ess08_Methode_31032011.csv", sep=",")
693
694
695 #####
696 # APPENDIX
697 # 1) ELEMENTARY DATABASE MANAGEMENT
698 # 2) MISSING DATA RECODING:
699 # 3) SELECTIONNER VARIABLES by indexing
700 #####
701 #####
702 #
703 # 1) ELEMENTARY DATABASE MANAGEMENT
704 #
705 # How many cases are there in your database?
706 # La commande dim() nous permet de savoir combien de cases il y
707 # a dans notre matrice des données.
708
709 dim(d)
710
711 # Finally we want to know which are the attributes of our database more
712 # in detail.
713 # The command attributes will give us a list of informations about our
714 # dataset:
715 # 1) Names of variables;
716 # 2) Class of the object;
717 # 3) The names of the row.
718
719 attributes(d)
720
721 # Differement de Spss R peut manager different type d'objets.
722 # Ceux-là il sont storÈ dans le workspace.
723 # Le workspace est comme bureau de travail, la table sur la quelle vous
724 # mettez ce que vous
725 # produisez dans une session de travail avec R.
726 #
727 # Pour savoir quels object sont present dans votre workspace vous utilisez la
728 # fonction ls()
729
730
731
732

```

```

733 ls()
734
735 # On peut definir et acceder aux Variable labels par les fonctions label et
736 # describe.
737 # Celle ci donne une liste des variables avec des descriptives:
738 # la commande se trouve dans la livrairie Hmisc, pour autant elle ne marche
739 # pas si vous ne lodez sur R cette livrairie.
740
741 label(d$txearn)<-"Taxation for higher versus lower earners"
742
743 # On accede aux descriptives de la variable avec
744
745 describe(d$txearn)
746
747
748 #####
749 #
750 # 2) MISSING DATA RECODING:
751 # Each data contains values that are missing or not applicable.
752 # Il est donc necessaire de dire à R quels valeurs correspondent à missing.
753 # D'habitude c'est des labels comme "Don't know", "Refusal", "No answer", ou
754 # encore 88 ou 99999
755 # La synthax suivante transforme les missing
756 # See the file Dataset setting for the complete synthax:
757
758 d$mnactic[d$mnactic== 88] <- NA
759
760 # Or si vous voulez recoder toutes les valeurs d'un dataframe qui
761 # sont egales à "Don't know" etc... comme missing, il faut user cette
762 # code:
763
764 d[d=="Don't know" | d == "No answer" | d=="Refusal" ] = NA
765
766 # Comme on l'a vu dans la séance 4 il y a different type de variables qui
767 # R peut manager:
768 # Variables Qualitatives:
769 # 1) Echelle nominale:
770
771 d$gndr <- factor(d$gndr)
772 str(d$gndr)
773
774 # 2) Echelle ordinale
775
776 d$edulvl <- ordered(d$edulvl, levels=c("Not Compl.Prim. Educ",
777 "Primary or first stage of basic",
778 "Lower secondary or second stage of basic",
779 "Upper secondary",
780 "Post secondary, non-tertiary",
781 "First stage of tertiary",
782 "Second stage of tertiary"))
783
784 str(d$edulvl)
785
786 # 3) Interval:
787 # Discrete:
788
789 str(d$happy.num)
790
791 # Continues:
792
793 str(d$dweight)

```

```

794
795 # Manipulation sur une variable:
796 # Acceder aux levels d'une variable:
797
798 levels(d$edulvl)
799
800 # De fois les labels des categories d'une variable sont trop longues (surtout
801 quand on fait des graphs).
802 # Par exemple
803
804 table(d$gndr, d$txearn)
805
806 # On peut donner à la variable des labels plus courts des façon a faciliter
807 le plotting
808 # Procedure:
809 # 1. Creer une nouvelle variable:
810 d$txe <- d$txearn
811 levels(d$txe) <- c("HEHSET", "NoneOfThese", "SAMT", "SSET")
812 table(d$gndr, d$txe)
813
814
815 #####
816 #
817 # 3)SELECTIONNER VARIABLES by indexing
818 # Un autre moyen de selectionner variables et observations en R est
819 # par indexation. Dans ce cas on considère le dataframe comme une matrice:
820 # (voir seance 6) de laquelle on slectionne observations et variables.
821
822 # Ex. On commence de la matrice X
823
824 X = matrix(1:100, ncol = 5)
825
826 X
827
828 # De cette matrice on peut extraire lignes, coulones ou cases:
829
830 # Ex. Deuxième coulone:
831 X[,2]
832
833 # Ex. CInquième ligne
834 X[5,]
835
836 # Ex. Element en 5ème ligne et 2ème coulone
837 X[5,2]
838
839 # Quand on travaill avec en dataframe on peut faire la meme chose:
840
841 # Ex. 1 selectioner des reponse en ce qui concerne le vote, la'acitivté
842 principale
843 # et l'opinion sur la taxation.
844
845 dh <- d[,c("ident", "vote", "mnactic", "txearn")]
846
847 # Ex. 2: selectionner la sous-population des femmes.
848
849 dfe <- d[d$gndr=="Female",]
850
851 # Ex. 3 On peut aussi selectionner des souspopulations:
852
853 dhvo <- d[d$gndr== "Male" & d$vote != "Not eligible to vote", c("ident",
854 "vote", "prtclbch", "mnactic", "txearn")]

```

```
855  
856 table(dhvo$vote)  
857  
858 table(dh$vote)
```