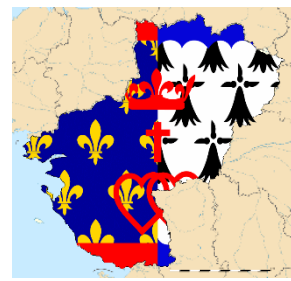


# SAE – Échantillonnage et Estimation

Population des Pays de la Loire



## Partie 1 : Estimation du nombre d'habitants

Dans le cadre de cette SAE, nous avons tenté d'estimer la population de la région Pays de la Loire en utilisant deux méthodologies de sondage aléatoire : le sondage aléatoire simple et le sondage stratifié.

Les résultats obtenus via le sondage aléatoire simple sont présentés sur un graphique en fin de rapport, en comparaison avec la population réelle de la région, estimée à 3,9 millions d'habitants. On observe que la population estimée varie fortement selon les tirages, avec un minimum de 1,3 million et un maximum de 5,5 millions d'habitants. La marge d'erreur est très élevée, atteignant parfois un niveau proche de la population réelle, ce qui limite la fiabilité de cette méthode.

Cependant, lors de deux tirages, l'estimation s'est révélée très proche de la population réelle. De manière générale, la majorité des estimations présentent un écart de 500 000 à 1 000 000 d'habitants par rapport à la valeur réelle.

```

1 # -----
2 # Initialisation
3 # -----
4
5 library(sampling) # Chargement du package pour les fonctions d'échantillonnage
6
7 # Chargement du fichier CSV contenant les données sur les communes françaises
8 Commune_fr = read.csv2("population_francaise_communes.csv", sep=";", dec=".", header=TRUE)
9
10 # -----
11 # Conception et nettoyage
12 # -----
13
14 Commune_fr <- Commune_fr[1:6, ] # Suppression des 6 premières lignes (souvent des en-têtes ou commentaires)
15 Commune_fr <- Commune_fr[, -9] # Suppression de la colonne 9 (inutile pour l'analyse)
16 Commune_fr <- Commune_fr[, -10] # Suppression de la colonne 10
17
18 colnames(Commune_fr) <- Commune_fr[1, ] # Redéfinition des noms de colonnes à partir de la 1ère ligne restante
19 Commune_fr <- Commune_fr[-1, ] # Suppression de cette ligne qui sert maintenant d'en-tête
20
21 # -----
22 # Partie 1.1 : Échantillonnage aléatoire simple
23 # -----
24
25 # Filtrer les données pour ne garder que la région "Pays de la Loire"
26 donnees <- Commune_fr[Commune_fr$Nom de la région == "Pays de la Loire", ]
27 donnees <- donnees[, c(3, 7, 9)] # Sélection des colonnes d'intérêt
28
29 # Nettoyage de la colonne "Population totale"
30 donnees$Population totale <- gsub(" ", "", donnees$Population totale) # Supprimer les espaces
31 donnees$Population totale <- as.numeric(donnees$Population totale) # Conversion en numérique
32
33 head(donnees) # Affichage des premières lignes
34
35 u = donnees$Commune # Liste des communes
36 head(u)
37
38 N = length(u) # Taille de la population
39 print(N)
40
41 T = sum(donnees$Population totale) # Somme totale des populations
42 print(T)
43
44 n = 100 # Taille de l'échantillon
45 E = sample(u, n) # Tirage aléatoire simple
46 head(E)
47
48 donnees1 = donnees[donnees$Commune %in% E, ] # Sous-échantillon
49 head(donnees1)
50
51 xbar = mean(donnees1$Population totale) # Moyenne estimée
52 print(xbar)
53
54 idcmoy = t.test(donnees1$Population totale)$conf.int # Intervalle de confiance
55 print(idcmoy)
56
57 T_test = N * xbar # Estimation du total
58 print(T_test)
59
60 idcT = idcmoy * N # Intervalle de confiance pour le total
61 print(idcT)
62
63 marge = (idcT[2] - idcT[1]) / 2 # Marge d'erreur
64 print(marge)

```

Voici le code utilisé pour installer la librairie sampling ainsi que pour lire et trier les données nécessaires à l'analyse.

La méthode utilisée consiste à effectuer une estimation de la population à partir d'un échantillon de 100 villes tirées aléatoirement dans la région des Pays de la Loire. Pour chaque tirage, la population estimée, l'intervalle de confiance (IDC) et la marge d'erreur ont ensuite été calculés.

Les résultats obtenus via le sondage aléatoire stratifié sont présentés dans le graphique en fin de rapport, montrant la population estimée, la population réelle (3,9 millions d'habitants) et la marge d'erreur associée.

La population estimée maximale atteint 4,3 millions d'habitants, avec seulement deux sondages dépassant les 4 millions. Le minimum observé est de 2,8 millions d'habitants. On constate que la marge d'erreur est, dans l'ensemble, nettement plus faible que celle obtenue avec le sondage aléatoire simple.

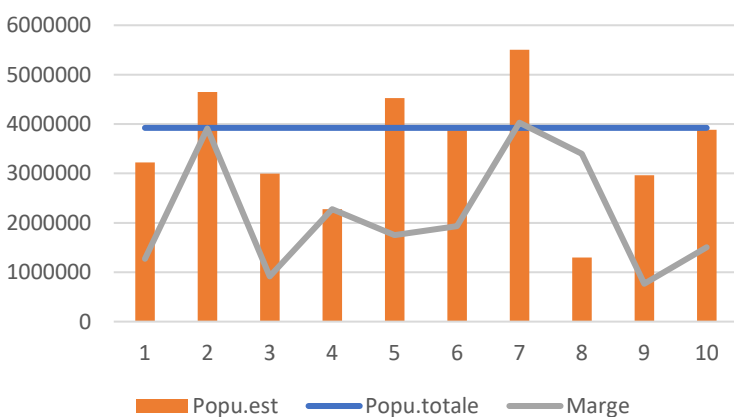
En moyenne, l'écart entre la population réelle et la population estimée est significativement réduit. La plupart des estimations se situent entre 200 000 et 600 000 habitants d'écart par rapport à la valeur réelle.

```
66 # -----
67 # Partie 1.2 : Échantillonnage aléatoire stratifié
68 # -----
69
70 summary(donnees$`Population totale`) # Statistiques descriptives
71
72 # Création des strates
73 donnees$strate = cut(donnees$`Population totale`, breaks=c(0, 525, 1173, 2719, 325857), labels=c(1, 2, 3, 4))
74 datastrat = donnees[, c("Commune", "Population totale", "strate")]
75 head(datastrat)
76
77 data = datastrat[order(donnees$strate), ]
78 head(data)
79
80 Nh = table(data$strate) # Taille de chaque strate
81 print(Nh)
82
83 N = sum(Nh)
84 print(N)
85
86 gh = Nh / N # Proportions de chaque strate
87 print(gh)
88
89 n = 100
90 nh = c(25, 25, 25, 25) # Taille d'échantillon égale par strate
91 print(nh)
92
93 fh = nh / Nh # Taux de sondage
94 print(fh)
95
96 # Tirage aléatoire stratifié
97 st = strata(data, stratanames = c("strate"), size = nh, method = "srswr")
98 data1 = getdata(data, st)
99 head(data1)
100 print(length(data1$Commune)) # Taille totale de l'échantillon
101
102 # Extraction des sous-échantillons
103 ech1 = data1[data1$strate == 1, ]
104 ech2 = data1[data1$strate == 2, ]
105 ech3 = data1[data1$strate == 3, ]
106 ech4 = data1[data1$strate == 4, ]
```

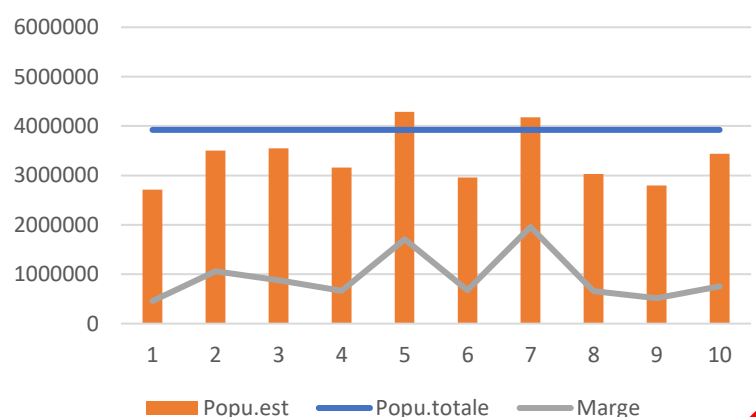
```
108 # Moyennes par strate
109 m1 = mean(ech1$`Population totale`)
110 m2 = mean(ech2$`Population totale`)
111 m3 = mean(ech3$`Population totale`)
112 m4 = mean(ech4$`Population totale`)
113 print(m1); print(m2); print(m3); print(m4)
114
115 # Variances par strate
116 var1 = var(ech1$`Population totale`)
117 var2 = var(ech2$`Population totale`)
118 var3 = var(ech3$`Population totale`)
119 var4 = var(ech4$`Population totale`)
120 print(var1); print(var2); print(var3); print(var4)
121
122 # Moyenne stratifiée estimée
123 xbarst = (Nh[1]*m1 + Nh[2]*m2 + Nh[3]*m3 + Nh[4]*m4) / N
124 print(xbarst)
125
126 # Variance de l'estimateur de moyenne
127 varxbarst = ((gh[1]^2)*(1 - fh[1])*var1/(nh[1]) + ((gh[2]^2)*(1 - fh[2])*var2/(nh[2]) +
128 ((gh[3]^2)*(1 - fh[3])*var3/(nh[3]) + ((gh[4]^2)*(1 - fh[4])*var4/(nh[4])
129 print(varxbarst)
130
131 # Intervalle de confiance pour la moyenne
132 alpha = 0.05
133 binf = xbarst - qnorm(1 - alpha/2) * sqrt(varxbarst)
134 bsup = xbarst + qnorm(1 - alpha/2) * sqrt(varxbarst)
135 idcmoy = c(binf, bsup)
136 print(idcmoy)
137
138 # Estimation du total
139 Tstr = N * xbarst
140 print(Tstr)
141
142 # Intervalle de confiance pour le total
143 idct = c(idcmoy[1]*N, idcmoy[2]*N)
144 print(idct)
145
146 marge = (idct[2] - idct[1]) / 2 # Marge d'erreur
147 print(marge)
```

Cette seconde partie du code permet de définir les différentes strates à partir des quartiles, puis de calculer la population estimée, l'intervalle de confiance (IDC) ainsi que la marge d'erreur pour chaque tirage.

Sondage aléatoire simple

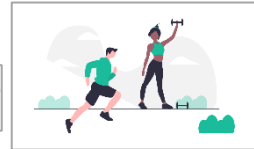


Sondage aléatoire stratifié



Le sondage aléatoire simple ne fonctionne correctement que lorsque la population des villes est homogène. Or, dans notre cas, elle est très hétérogène, ce qui rend cette méthode peu adaptée. Le sondage aléatoire stratifié est donc plus pertinent. Cette SAÉ nous a permis de mieux utiliser le logiciel R pour appliquer des techniques d'échantillonnage, et de comprendre que la méthodologie choisie influence fortement les résultats obtenus.

## Partie 2 : Traitement des données d'une enquête



Dans cette seconde partie, nous travaillons à partir d'un fichier Excel issu d'une enquête sur les habitudes sportives des étudiants.

Voici le tableau présentant les variables qualitatives que nous avons sélectionnées (Sexe, Fumeur, Santé, Fan de sport et Logement), accompagné des résultats des tests statistiques réalisés avec R. Parmi ces variables, seuls les tests portant sur le sexe et sur le fait d'être fan de sport se révèlent statistiquement significatifs.

	Retour des valeurs trouvées sur l'enquête Sport				
	Sexe	Fumer	Sante	Fan	logement
khideux	14.742	0.81111	0.70524	78.823	4.8007
p-valeur	0.0006292	0.6666	0.7028	< 2.2e-16	0.3084
V-Cramer	0.198274			0.4584704	

Voici le script R commenté correspondant à l'ensemble de la deuxième partie. Il comprend la création de tableaux croisés dynamiques, ainsi que les différents calculs liés au test du Khi<sup>2</sup> et au V de Cramer pour analyser les relations entre les variables qualitatives.

```

149 #
150 # Partie 2 : Traitement de données d'enquête
151 #
152
153 Sport <- read.csv2("EnqueteSportEtudiant2024.csv", sep = ";", dec = ",", header = TRUE)
154 head(Sport) # Aperçu
155
156 # Cette base contient :
157 # - Des "individus" : chaque ligne représente un étudiant ayant répondu à l'enquête
158 # - Des "variables" : chaque colonne est une question ou caractéristique (ex : sexe, fumer, santé, logement, sport...)
159 # - Types de variables :
160 # - "Catégorielles" (qualitatives) : sexe, fumer, sante, logement, fan, sport
161 # - "Quantitatives" (quantitatives) : nombre de sport, nombre d'heures de sport, etc.
162 # - "Possiblement numériques" : si d'autres variables (non affichées ici) contiennent des quantités (ex : âge, nombre d'heures de sport, etc.)
163
164 # Création des tableaux croisés
165 TCD_Sexe <- table(Sport$Sport, Sport$sexe)
166 TCD_Fumer <- table(Sport$Sport, Sport$fumer)
167 TCD_Sante <- table(Sport$Sport, Sport$sante)
168 TCD_Fan <- table(Sport$Sport, Sport$fan)
169 TCD_Logement <- table(Sport$Sport, Sport$logement)
170
171 # Affichage des tableaux croisés
172 print(TCD_Sexe)
173 print(TCD_Fumer)
174 print(TCD_Sante)
175 print(TCD_Fan)
176 print(TCD_Logement)
177
178
179 # Tests du khi-deux
180 khideux_Sexe = chisq.test(TCD_Sexe)
181 print(khideux_Sexe)
182
183 khideux_Fumer = chisq.test(TCD_Fumer)
184 print(khideux_Fumer)
185
186 khideux_Sante = chisq.test(TCD_Sante)
187 print(khideux_Sante)
188
189 khideux_Fan = chisq.test(TCD_Fan)
190 print(khideux_Fan)
191
192 khideux_Logement = chisq.test(TCD_Logement)
193 print(khideux_Logement)
194
195 # V de Cramer (Sexe)
196 n <- dim(Sport)[1]
197 p <- nrow(TCD_Sexe)
198 q <- ncol(TCD_Sexe)
199 m <- min(p - 1, q - 1)
200 V_Sexe = sqrt(khideux_Sexe$statistic / (n * m))
201 print(V_Sexe)
202
203 # V de Cramer (Fan de sport)
204 p <- nrow(TCD_Fan)
205 q <- ncol(TCD_Fan)
206 m <- min(p - 1, q - 1)
207 V_Fan = sqrt(khideux_Fan$statistic / (n * m))
208 print(V_Fan)

```

Pour conclure, les résultats montrent qu'il existe un lien significatif entre la pratique sportive, le sexe et le fait d'être fan de sport. Le V de Cramer obtenu pour la relation entre la pratique sportive et le sexe est de 0.198, indiquant une association faible. En revanche, la relation entre la pratique sportive et le fait d'être fan présente un V de Cramer de 0.458, ce qui traduit une association modérée à forte.