

SAE - Echantillonnage et Estimation

La Bretagne



Partie 1 : Estimation du nombre d'habitants d'une région de France : La Bretagne :

Partie 1.1 : Echantillonnage aléatoire simple :

Ce travail, mené à l'aide du logiciel R, vise à développer une compréhension approfondie de l'incertitude et de la précision associées à l'estimation d'une variable mesurable au sein d'une population, et ce par le biais de la construction d'intervalles de confiance fondés sur des techniques d'échantillonnage. Dans un premier temps, l'approche adoptée repose sur un échantillonnage aléatoire simple à probabilités égales, dans lequel chaque individu de la population dispose de la même chance d'être sélectionné. Dans un second temps, l'étude se poursuit par l'application d'un échantillonnage stratifié, méthode permettant de segmenter la population en sous-groupes homogènes afin d'améliorer la précision des estimations.

Pour commencer, nous avons filtré les régions afin de ne sélectionner que celle sur laquelle nous allons travailler. Ainsi, nous allons présenter notre recherche sur la Bretagne.

Nous avons débuté par l'importation du jeu de données qui regroupe toutes les populations par départements et communes de la France. Nous avons donc stocké toutes les données du fichier dans un data frame. Nous avons également intégré la librairie « sampling » qui nous a permis de procéder à des tirages aléatoires sur les communes.

```
1 # Importation de librairie
2 library(sampling)
3
4 # Importation des données du fichier CSV
5 data <- read.csv2("communes.csv", sep=";", dec=".", header=TRUE, fileEncoding = "latin1")
6 Estimations <- data.frame(Total= numeric(), Total_estime=numeric(), IDCinf=numeric(), IDCsup=numeric(), marge= numeric())
7
```

Nous avons choisi de travailler sur la région de la Bretagne, par conséquent nous avons extrait toutes les données relatives à cette région dans une table nommée données. Puis, nous avons choisi de garder uniquement trois variables (colonnes) : le code département, les communes et la population totale.

```
8 ## Partie 1.1
9 # Filtrer pour avoir que les Code.département, Commune et Population.totale de la Bretagne
10 donnees <- data[data$Nom.de.la.région == "Bretagne", c("Code.département", "Commune", "Population.totale")]
11 donnees$Population.totale <- as.numeric(gsub(" ", "", donnees$Population.totale))
12 # Affichage des 6 premières lignes
13 head(donnees)
14
```

L'utilisation du code département garantit que, même si certaines communes portent le même nom, elles seront identifiées distinctement. Ensuite nous avons converti la colonne "Population.totale" en une variable numérique, et en enlevant les espaces (1 200 devient 1200).

Par la suite, nous avons défini une variable appelée U qui représente toutes les communes de la région de Bretagne. Cette variable répertorie tous les noms des communes de cette région, qui est la plus peuplée du pays. Cette région possède 1 207 communes. Nous avons aussi calculé le nombre total d'habitants dans la région : la Bretagne abrite 12 384 734 résidents. Nous avons opté pour cette démarche car, ultérieurement, nous prévoyons d'évaluer la population de cette zone en utilisant un échantillon.

```
15 # Création d'une table contenant l'ensemble des communes et compte du nombre total
16 U = donnees$Commune
17 head(U)
18 N = length(U)
19
20 # Calcul du nombre total d'habitants
21 T = sum(donnees$Population.totale)
```

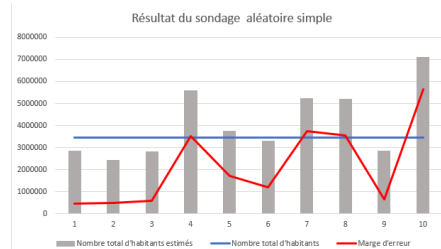
Par la suite, nous avons procédé à l'estimation de la population de la Bretagne en nous basant sur un échantillon aléatoire de 100 communes. Ce tirage nous a permis de constituer un data frame, contenant les informations de chaque commune tirée. Par la suite, nous avons déterminé le nombre moyen d'habitants par commune dans cet échantillon de 100 communes sélectionnées. Cela nous a donné la possibilité d'évaluer la dimension de la population dans cette région. Effectivement, étant donné que nous connaissons le nombre moyen d'habitants, le nombre total de communes et le nombre de communes incluses dans l'échantillon, nous sommes alors en mesure d'estimer le nombre d'habitants en Bretagne. Nous avons également effectué le calcul de l'IDC ainsi que de la marge d'erreur. En créant une boucle, nous avons réalisé ces calculs en utilisant 10 tirages aléatoires distincts.

```
23 # Boucle pour réaliser les 10 estimations et les ajouter dans un data frame
24 for (i in 1:10) {
25   # Tirage d'échantillon selon un sondage aléatoire simple
26   n=100
27   E=sample(U,n)
28   length(E)
29   head(E)
30
31   # Nouvelle table contenant seulement les communes tirés
32   donnees1 <- donnees[donnees$Commune %in% E,]
33
34   # Nombre moyen d'habitants de l'échantillon E
35   xbar=mean(donnees1$Population.totale)
36   xbar
37
38   #IDC à 95%
39   idcmoy = t.test(donnees1$Population.totale)$conf.int
40   idcmoy
41
42   # Nombre d'habitants estimé
43   T_est = N*xbar
44   T_est
45
46   # IDC de T
47   idcT = idcmoy*N
48   idcT
49
50   # marge d'erreur
51   marge = (idcT[2]-idcT[1])/2
52   marge
53
54   # Ajouter une nouvelle ligne au data frame
55   Estimations <- rbind(Estimations, data.frame(Total= T,Total_estime=T_est,
56                                                IDCinf=idcT[1], IDCsup=idcT[2], marge= marge))
57 }
```

Les différents résultats montrent qu'il y a tout de même des écarts conséquents entre la valeur estimée et la valeur réelle, cela est dû à la taille des communes sélectionnées dans l'échantillon de 100

communes. Effectivement, on observe d'énormes disparités en termes de population entre les différentes communes.

	Nombre total d'habitants	Nombre total d'habitants estimés	IDC binf	IDC bsup	Marge d'erreur
1	3463439	2829317	2062076	2989125	463524
2	3463439	2411670	1920810	2902531	490860
3	3463439	2819422	2236355	3402489	583067
4	3463439	5582303	2066940	9097665	3515363
5	3463439	3736244	2037083	5435405	1699161
6	3463439	3286898	2085054	4488742	1201844
7	3463439	5221267	1483746	8958788	3737521
8	3463439	5194807	1634205	8755410	3560602
9	3463439	2830415	2176412	3484418	654003
10	3463439	7108544	1479841	12737246	5628703



Néanmoins, nous pouvons remarquer grâce au graphique que nous avons 3 estimations qui sont très proche de la réalité. Mais cela reste une petite part, par rapport aux 7 autres essais effectués.

Partie 1.2 : Echantillonnage aléatoire stratifié :

Par la suite, afin de minimiser autant que possible ces disparités, nous avons opté pour la création d'un échantillon de 100 communes en recourant à la méthode des quotas. Autrement dit, en organisant le tirage selon des normes qui reflètent les mêmes proportions que la population visée, i-e en créant des strates.

```

22 # Affichage des quartiles
23 summary(donnees$Population.totale)
24
25 # strates :
26 # On transfère le dataframe 'donnees' dans un autre dataframe appelé 'datastrat' pour préciser
27 # que l'on va manipuler les données avec les différentes strates de communes
28 datastrat = donnees
29 datastrat$strate = cut(donnees$Population.totale, breaks = c(0,716,1390,2762,230000), labels = c(1,2,3,4))
30 # Affichage des 6 premières lignes du nouveau data frame
31 head(datastrat)
32
33 # Sondage stratifié
34 datastrat = datastrat[order(datastrat$strate), ]

```

Nous avons ainsi débuté par la détermination des différents quartiles de notre groupe de population. Ces derniers nous ont permis de concevoir des strates entre les diverses municipalités de Bretagne, dans l'objectif d'estimer la population de manière plus précise. Nous avons donc 4 strates différentes qui sont, les communes de moins de 716 habitants, les communes entre 716 et 1390 habitants, les communes entre 1390 et 2762 habitants et pour finir toutes celles qui ont plus de 2762 habitants. Nous avons par la suite réalisé un data frame pour ordonner les communes selon leur strat respective.

Par la suite, nous avons calculé l'effectif et le poids des strates. Nous avons fait une boucle pour procédé à une sélection aléatoire de 100 communes issues des 4 strates distinctes. Nous avons sélectionné 25 municipalités pour chaque stratification. Le tirage aléatoire étant sans remise, nous avons calculé le taux de sondage (environ $0.08 < 0.1$), et ces tirages peuvent alors être assimilé à des tirages avec remise. Ainsi, chaque commune est indépendante l'une de l'autre.

```

44 # Boucle pour réaliser les 10 estimations et les ajouter dans un data frame
45 for (i in 1:10) {
46   # Tirage d'un échantillon stratifié de taille n = 100
47   n = 100
48   nh = round(c(n*Nh[1]/N, n*Nh[2]/N, n*Nh[3]/N, n*Nh[4]/N))
49
50   # taux de sondage dans les strates
51   fh = nh/Nh
52   fh

```

Ensuite, nous avons réalisé des tirages aléatoires simples dans chacune des strates puis calculer leurs moyennes et leurs variances respectives. En utilisant ces moyennes de différents échantillons, nous pourrions alors établir des estimations plus précises de la population de Bretagne par rapport aux premières estimations. En utilisant chaque moyenne des échantillons, nous avons été en mesure de déterminer la moyenne globale et la variance générale.

```

54 # Sondage strat (sans remise dans les strates)
55 st = strata(datastrat, stratanames = c("strate"),size = nh, method = "srswr")
56 data1 = getdata(datastrat, st)
57 head(data1)
58 length((data1$Commune))
59
60 # Sous-échantillons
61 ech1 = data1[data1$strate == 1, ]
62 ech1
63 ech2 = data1[data1$strate == 2, ]
64 ech2
65 ech3 = data1[data1$strate == 3, ]
66 ech3
67 ech4 = data1[data1$strate == 4, ]
68 ech4
69
70 # Moyenne des sous-échantillons
71 m1 = mean(ech1$Population.totale)
72 m2 = mean(ech2$Population.totale)
73 m3 = mean(ech3$Population.totale)
74 m4 = mean(ech4$Population.totale)
75
76 # Moyenne des 4 sous échantillons réunis
77 Xbarst = (Nh[1]*m1 + Nh[2]*m2 + Nh[3]*m3 + Nh[4]*m4)/N
78
79 # Variances des 4 sous-échantillons
80 var1 = var(ech1$Population.totale)
81 var2 = var(ech2$Population.totale)
82 var3 = var(ech3$Population.totale)
83 var4 = var(ech4$Population.totale)
84
# Moyenne des 4 sous échantillons réunis
Xbarst = (Nh[1]*m1 + Nh[2]*m2 + Nh[3]*m3 + Nh[4]*m4)/N
# Estimation de la variance de Xbarst
varXbarst = ((gh[1])^2)*(1-fh[1])*var1/(nh[1]) + ((gh[2])^2)*(1-fh[2])*var2/(nh[2]) +
((gh[3])^2)*(1-fh[3])*var3/(nh[3]) + ((gh[4])^2)*(1-fh[4])*var4/(nh[4])

```

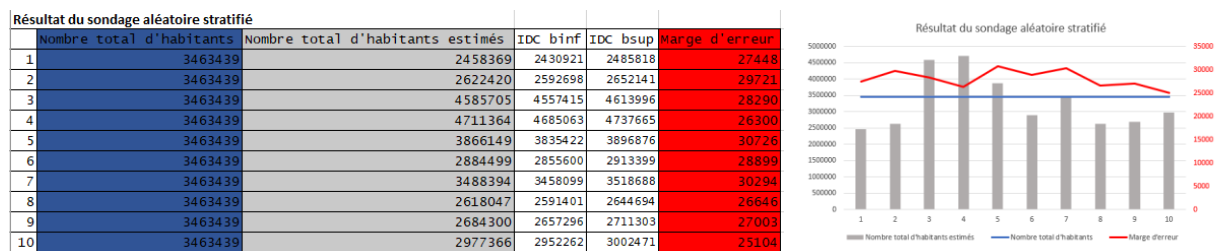
En utilisant chaque moyenne et variance des échantillons, nous avons été en mesure de calculer l'IDC pour μ à 95%. Pour finir, nous avons estimé la population de Bretagne, son IDC et la marge d'erreur à dix reprises avec une boucle.

```

89 # IDC pour  $\mu$  à 95%
90 alpha = 0.05
91 binf = Xbarst - qnorm(1-alpha/2)*sqrt(varXbarst)
92 bsup = Xbarst + qnorm(1-alpha/2)*sqrt(varXbarst)
93 idcmoy = c(binf, bsup)
94
95 # Estimation du total de la population (T)
96 Tstr = N * Xbarst
97 Tstr
98
99 # Estimation de l'IDC pour T
100 binf = idcmoy[1] * N
101 bsup = idcmoy[2] * N
102 idcT = c(binf,bsup)
103 idcT
104
105 # marge d'erreur
106 marge = (idcT[2]-idcT[1])/2
107 marge
108
109 # Ajouter une nouvelle ligne au data frame
110 Estimations <- rbind(Estimations, data.frame(Total= T,Total_estime=Tstr,
111 IDCinf=idcT[1], IDCsup=idcT[2], marge= marge))
112 }

```

On peut observer à travers le tableau et le graphique suivants que les résultats sont clairement plus précis que ceux de la première partie. Effectivement, l'établissement de strates a facilité l'équilibre entre les grandes et petites communes afin d'éviter toute disproportion.



Partie 2 : Traitement de données d'enquête :

L'objectif de cette partie est de réaliser le test du khi-deux d'indépendance afin de trouver des relations significatives entre la pratique sportive des étudiants et certaines de leurs caractéristiques.

Tout d'abord, nous avons débuté par l'importation du jeu de données qui regroupe l'ensemble des informations sur l'activité sportive des étudiants ainsi que d'autres données les concernant, que nous avons pu récupérer antérieurement grâce à un questionnaire. Ces informations ont été stockées dans un data frame.

```
1 data <- read.csv2("EnqueteSportEtudiant2024.csv", sep=";", dec=".", header=TRUE)
2 head(data)
```

Ce dernier contient 76 variables et 375 observations / réponses, et les individus sont les étudiants de l'IUT de Niort. Quant aux variables, on y retrouve le sexe, le département géographique, le département de formation, le niveau, et d'autres informations recueillies lors du questionnaire. Par ailleurs, les variables sont toutes qualitatives.

Nous avons ensuite décidé de croiser la variable 'sport' avec les variables suivantes : la sante, si l'étudiant est fumeur ou non, le sexe, le département de formation, l'alimentation et la bonification.

```
6 # Tableaux croisés dynamiques de la variable 'sport' avec
7
8 TCD_sante = table(data$sport, data$sante)
9 TCD_sante # Afficher le tableau
10 TCD_fumer = table(data$sport, data$fumer)
11 TCD_fumer # Afficher le tableau
12 TCD_sexe = table(data$sport, data$sexe)
13 TCD_sexe # Afficher le tableau
14 TCD_deptformation = table(data$sport, data$deptformation)
15 TCD_deptformation # Afficher le tableau
16 TCD_alimentation = table(data$sport, data$alimentation)
17 TCD_alimentation # Afficher le tableau
18 TCD_bonif = table(data$sport, data$bonif)
19 TCD_bonif # Afficher le tableau
```


Nous avons ensuite effectué le test d'indépendance du khi-deux entre le sport et les variables choisies.

```
21 # Test d'indépendance du khi-deux entre la variable sport et les autres variables qualitatives
22
23 khideux_sante = chisq.test(TCD_sante)
24 khideux_sante # Afficher le khi-deux
25 khideux_fumer = chisq.test(TCD_fumer)
26 khideux_fumer # Afficher le khi-deux
27 khideux_sexe = chisq.test(TCD_sexe)
28 khideux_sexe # Afficher le khi-deux
29 khideux_deptformation = chisq.test(TCD_deptformation)
30 khideux_deptformation # Afficher le khi-deux
31 khideux_alimentation = chisq.test(TCD_alimentation)
32 khideux_alimentation # Afficher le khi-deux
33 khideux_bonif = chisq.test(TCD_bonif)
34 khideux_bonif # Afficher le khi-deux
```

Précision importante : nous avons choisi arbitrairement de poser le niveau de risque alpha à 0,05.

Pour chaque khi-deux calculé, nous avons obtenu les p-valeurs suivantes :

p-value = 0.70 pour la variable 'sante'

p-value = 0.67 pour la variable 'fumer'

p-value = 0.0006 pour la variable 'sexe'

p-value = 0.0046 pour la variable 'deptformation'

p-value = 0.0002 pour la variable 'alimentation'

p-value = 0.0000000021 pour la variable 'bonif'

Nous avons donc pu déterminer le khi-deux théorique à l'aide de la table du khi-deux, et nous avons pu arriver aux conclusions suivantes :

Concernant la variable santé et la variable fumer, elles sont toutes les deux indépendantes de la variable sport au risque de 5%, car leur khi-deux lue (5,99) est supérieur à leur khi-deux calculé. En revanche, on rejette l'hypothèse d'indépendance pour la variable sexe au risque de 5%, car le khi-deux calculé (14,7) est supérieur au khi-deux lu (5,9). C'est aussi le cas pour le département de formation avec un khi-deux calculé de 18,8 et un khi-deux lue de 12,6. Le test est aussi significatif pour les variables alimentation et bonification, leur khi-deux calculé sont respectivement de 16,7 et de 46,3 tandis que leur khi-deux lu respectif sont 5,9 et 9,5. Ces 4 derniers tests sont donc significatifs, car nous pouvons rejeter l'hypothèse d'indépendance au risque alpha de 5%.

Nous avons ensuite souhaité déterminer le niveau de dépendance de ces 4 relations significatives. Ainsi nous avons calculé le V de Cramer pour la variable 'sexe', 'deptformation', 'alimentation' et 'bonif'.

```

39 # Calcul du V de Cramer de la variable 'sexe', 'deptformation', 'alimentation' et 'bonif'
40
41 # Effectif total
42 n<-nrow(data)
43 |
44 # Nb de lignes
45 p <- nrow(TCD_sexe)
46 # Nb de colonnes
47 q <- ncol(TCD_sexe)
48 # Déterminer le minimum entre le nb de lignes et le nb de colonnes
49 m <- min(p-1,q-1)
50 # V de Cramer
51 V_sexe <- sqrt(khideux_sexe$statistic/(n*m))
52 V_sexe
53
54 p <- nrow(TCD_deptformation)
55 q <- ncol(TCD_deptformation)
56 m <- min(p-1,q-1)
57 # V de Cramer
58 V_deptformation <- sqrt(khideux_deptformation$statistic/(n*m))
59 V_deptformation
60
61 p <- nrow(TCD_alimentation)
62 q <- ncol(TCD_alimentation)
63 m <- min(p-1,q-1)
64 # V de Cramer
65 V_alimentation <- sqrt(khideux_alimentation$statistic/(n*m))
66 V_alimentation
67
68 p <- nrow(TCD_bonif)
69 q <- ncol(TCD_bonif)
70 m <- min(p-1,q-1)
71 # V de Cramer
72 V_bonif <- sqrt(khideux_bonif$statistic/(n*m))
73 V_bonif
74

```

Nous avons réalisé un tableau regroupant tous les V de Cramer de chaque test significatif, et avons rajouté une colonne dans laquelle est inscrit la valeur du V de Cramer le plus élevé.

```

76 # Tableau des V de Cramer
77
78 # Création d'un tableau de données
79 dataCramer = data.frame()
80 # Insertion des valeurs des V de Cramer
81 dataCramer = rbind(dataCramer,c(V_sexe, V_deptformation, V_alimentation, V_bonif))
82 # On renomme les colonnes
83 colnames(dataCramer) <- c("V_sexe", "V_deptformation", "V_alimentation", "V_bonif")
84 # On ajoute une colonne affichant le V de Cramer le plus élevé
85 dataCramer$V_Cramer_Max = max(c(V_sexe, V_deptformation, V_alimentation, V_bonif))
86 # On affiche le tableau de données
87 dataCramer

```

Voici donc le contenu du tableau :

```

> # On affiche le tableau de données
> dataCramer
  V_sexe V_deptformation V_alimentation V_bonif V_Cramer_Max
1 0.198274    0.1582276    0.2107832 0.2485422    0.2485422
> View(dataCramer)

```

Comme nous pouvons le constater, le sexe et le département de formation de l'étudiant ont tous les deux une faible liaison avec sa pratique sportive, puisque leur V de Cramer sont compris entre 0,10 et 0,20. Cependant, le V de Cramer de la variable 'alimentation' et 'bonif' est compris entre 0,20 et 0,25. Cela signifie que la pratique sportive des étudiants est modérément liée à leur alimentation, et au fait qu'ils aient choisi de réaliser une activité sportive en étant bonifié ou non.

Pour conclure, la pratique sportive des étudiants dépend légèrement du sexe de l'étudiant et de son département de formation, mais elle dépend davantage du choix de bonification de l'étudiant ainsi que de son alimentation.