

Plans par grappes et à plusieurs degrés

Sommaire

1. Principe, justification et premiers exemples
 - a. Principe
 - b. Justification
 - c. Premiers exemples

2. Plans par grappes
 - a. Cas général : tirage des grappes à probabilités inégales
 - b. Tirage des grappes proportionnellement à leur taille
 - c. Tirage des grappes à probabilités égales
 - d. Comparaison avec un SAS de même taille
 - e. Considérations pratiques

3. Plans à plusieurs degrés
 - a. Principe général
 - b. Tirage SAS à chaque degré
 - c. Tirages autopondérés
 - d. Considérations pratiques

4. Exemples des enquêtes ménages INSEE

1. Principe, justification et premiers exemples

a. Principe

On accède aux unités échantillonnées par l'intermédiaire de regroupements de ces unités.

Exemple d'une enquête auprès des ménages :

- *sélection d'un échantillon de communes,*
- *puis on retient :*
 - *tous les ménages des communes retenues : plans par grappes*
 - *un échantillon de ménages dans les communes retenues : plans à 2 degrés où les communes s'appellent « unités primaires » et les ménages « unités secondaires »*
 - *un échantillon d'individus dans les ménages et communes retenus aux degrés précédents : plans à 3 degrés*

b. Justifications

Avantages :

- Absence de base de sondage ou mauvaise qualité : seule la connaissance exhaustive des unités primaires (par ex, des villes) est nécessaire.
- Economie
- Gain de temps

Inconvénients :

- Perte de précision par rapport à un sondage aléatoire simple sans remise de même taille

Effets de grappes : les unités statistiques regroupées dans une même unité primaire ont souvent tendance à se ressembler.

↳ Il faut donc soigner le premier degré de tirage.

c. Premiers exemples

- En industrie, contrôle par lots : produits de série conditionnés par caisses.
- Etudes médicales : on accède aux patients ou à des prescriptions en interrogeant des médecins ou des laboratoires d'analyses médicales
- Sondages électoraux « sortie des urnes » réalisés auprès d'électeurs à la sortie de certains bureaux de vote.
- Enquête Emploi et enquêtes ménages INSEE

Les notations :

- Population : $U = \{1, \dots, k, \dots, N\}$

- M unités primaires (UP) ou M grappes : $U = \bigcup_{i=1}^M U_i$

Dont m sont échantillonnées dans S_G ou S_{UP}

- N unités secondaires (US) : $N = \sum_{i=1}^M N_i$ où $N_i = \text{card}(U_i)$

Choix de n_i unités parmi N_i dans les m UP i tirées dans S_I

- Total de la variable Y :

$$t_Y = \sum_{i=1}^M t_{Yi} = \sum_{k=1}^N Y_k \quad \text{Où} \quad t_{Yi} = \sum_{k \in U_i} Y_k$$

- Moyenne de la variable Y :

$$\mu_Y = \frac{1}{N} \sum_{k=1}^N Y_k = \frac{1}{N} \sum_{i=1}^M N_i \mu_{Yi} = \frac{t_Y}{N}$$

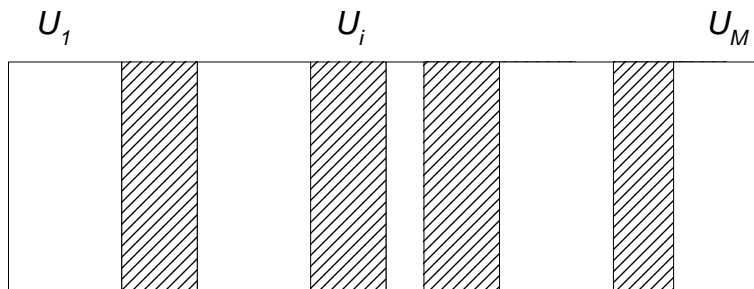
- Dispersion du total de la variable Y entre les grappes ou les UP

$$S^2_T = \frac{1}{M-1} \sum_{i=1}^M \left(t_{Yi} - \frac{t_Y}{M} \right)^2$$

2. Plans par grappes

Principe :

- Sélection d'un échantillon S_G de grappes
- Recensement dans chaque grappe retenue



- Probabilités d'inclusion des individus et des grappes:

$$\pi_k = \pi_i, \quad \forall k \in U_i$$
$$\pi_{kl} = \begin{cases} \pi_i, & \text{si } k, l \in U_i \\ \pi_{ij}, & \text{si } k \in U_i \text{ et } l \in U_j \end{cases}$$

- Taille d'échantillon aléatoire :

$$n_s = \sum_{k \in S} 1 = \sum_{i \in S_G} N_i$$

a. Cas général : tirage de grappes à probabilités inégales

Estimateur de Horvitz-Thompson du total

$$\hat{t}_y = \sum_{i \in S_G} \frac{t_{yi}}{\pi_i}$$

- Estimateur sans biais
- De variance :

$$Var(\hat{t}_y) = \sum_{i,j \in U_I} \Delta_{ij} \frac{t_{yi}}{\pi_i} \frac{t_{yj}}{\pi_j}$$

ou Sen-Yates-Grundy

- Estimateur sans biais de la variance si $\pi_{ij} > 0$ pour toutes les grappes i et j :

$$\hat{Var}(\hat{t}_y) = \sum_{i,j \in S_G} \frac{\Delta_{ij}}{\pi_{ji}} \frac{t_{yi}}{\pi_i} \frac{t_{yj}}{\pi_j}$$

ou Sen-Yates-Grundy

Estimation de la moyenne

Si N connu ,

Estimateur de Horvitz-Thompson

$$\hat{u}_y = \frac{\hat{t}_y}{N} = \frac{1}{N} \sum_{i \in S_G} \frac{t_{yi}}{\pi_i}$$

- Estimateur sans biais
- De variance : $\frac{1}{N^2} Var(\hat{t}_y)$
- estimée par : $\frac{1}{N^2} \hat{V}ar(\hat{t}_y)$

Si N inconnu ,

Estimateur de Hajek :

$$\hat{u}_{y,Hajek} = \frac{\hat{t}_y}{\hat{N}} = \frac{\sum_{i \in S_G} \frac{t_{yi}}{\pi_i}}{\sum_{i \in S_G} \frac{N_i}{\pi_i}}$$

- Estimateur biaisé
- Plus précis en présence d'un effet taille :

$$Var(\hat{u}_{y,Hajek}) \cong \frac{1}{N^2} \sum_{i,j=1,\dots,M} \Delta_{ij} \frac{N_i(\mu_{Yi} - \mu_Y)}{\pi_i} \times \frac{N_j(\mu_{Yj} - \mu_Y)}{\pi_{lj}}$$

$$Var(\hat{u}_{y,Hajek}) \cong \frac{1}{N^2} \sum_{i,j \in S_G} \frac{\Delta_{ij}}{\pi_{ij}} \frac{N_i(\mu_{Yi} - \hat{u}_{y,Hajek})}{\pi_i} \times \frac{N_j(\mu_{Yj} - \hat{u}_{y,Hajek})}{\pi_{lj}}$$

b. Tirage des grappes proportionnellement à leur taille

$$\pi_i = m \frac{N_i}{N} = \pi_k, \quad \forall k \in U_i$$

En moyenne, la taille d'échantillon vaut :

$$E(n_s) = \sum_{i \in U_I} N_i m \frac{N_i}{N} = \frac{m}{N} \sum_{i \in U_I} N_i^2$$

Estimateur de Horvitz-Thompson du total :

$$\hat{t}_y = \frac{N}{mN_i} \sum_{i \in s_G} t_{yi}$$

Estimateur de Horvitz-Thompson de la moyenne :

$$\hat{\mu}_y = \frac{1}{m} \sum_{i \in s_G} \mu_{yi}$$

c. Tirage des grappes à probabilités égales

$$\pi_i = \frac{m}{M} = \pi_k, \quad \forall i = 1, \dots, M, \forall k \in U$$

En moyenne, la taille d'échantillon vaut :

$$E(n_s) = \sum_{i \in U_I} N_i \frac{m}{M} = \frac{Nm}{M}$$

Estimateurs de Horvitz-Thompson du total et de la moyenne :

$$\hat{t}_y = \frac{M}{m} \sum_{i \in s_G} t_{yi}$$

$$\hat{\mu}_y = \frac{M}{Nm} \sum_{i \in s_G} N_i \mu_{yi}$$

- Estimateurs sans biais
- De variance pour le total :

$$\text{Var}(\hat{t}_y) = M^2 \left(1 - \frac{m}{M}\right) \frac{S_T^2}{m} \quad \text{Où} \quad S_T^2 = \frac{1}{M-1} \sum_{i=1}^M \left(t_{Yi} - \frac{t_Y}{M}\right)^2$$

- Estimée sans biais par :

$$\hat{\text{Var}}(\hat{t}_y) = M^2 \left(1 - \frac{m}{M}\right) \frac{s_T^2}{m} \quad \text{où} \quad s_T^2 = \frac{1}{m-1} \sum_{i \in s_G} \left(t_{Yi} - \frac{\hat{t}_Y}{M}\right)^2$$

↪ Bonne précision quand s_T^2 est faible et m grand

d. Comparaison avec un SAS de même taille

Supposons $N_i = N_o = N/M$ pour toutes les grappes

Alors,

$$n = m N_o \quad \text{et} \quad m/M = n/N$$

$$\hat{\mu}_y = \frac{1}{m} \sum_{i \in S_G} \mu_{yi}$$

$$\boxed{\text{Var}(\hat{\mu}_y) = \left(1 - \frac{m}{M}\right) \frac{M}{M-1} \frac{\sigma_{Yinter}^2}{m}} \quad \text{avec} \quad \sigma_{Yinter}^2 = \sum_{i=1}^M \frac{N_o}{N} (\mu_{Yi} - \mu_Y)^2$$

Car :

$$\begin{aligned} \text{Var}(\hat{\mu}_y) &= \left(\frac{M}{N}\right)^2 \left(1 - \frac{m}{M}\right) \frac{1}{m} \frac{1}{M-1} \sum_{i=1}^M \left(N_o \mu_{Yi} - \frac{t_Y}{M}\right)^2 \\ &= M^2 \left(1 - \frac{m}{M}\right) \frac{1}{m} \frac{1}{M-1} \sum_{i=1}^M \left(\frac{N_o}{N} \mu_{Yi} - \frac{N_o}{N} \mu_Y\right)^2 \end{aligned}$$

Avec un SAS sans remise de même taille n , on aurait eu :

$$\boxed{\text{Var}(\hat{\mu}_{Ysas}) = \left(1 - \frac{m}{M}\right) \frac{1}{n} S_y^2}$$

↳ **Perte de précision du plan par grappes due aux « effets de grappes » : « qui se ressemble s'assemble »**

- *clientèles typées selon les médecins, leur lieu d'exercice*
- *hétérogénéité des bureaux de vote, etc...*

Plus précisément, on mesure l'**effet de sondage** du plan par grappes :

$$Deff = \frac{Var(\hat{\mu}_y)}{Var(\hat{\mu}_{y,sas})} \cong 1 + (N_o - 1)\rho_{intra}$$

et $Deff > 1$ si $\rho_{intra} > 0$

$Deff < 1$ si $\rho_{intra} < 0$

qui s'obtient en écrivant :

$$Var(\hat{\mu}_y) \cong \left(1 - \frac{m}{M}\right) \frac{1}{mN_o} S_y^2 [1 + (N_o - 1)\rho_{intra}]$$

avec le *coefficient de corrélation intra-grappe* :

$$\rho_{intra} = \frac{\sum_{i=1}^M \sum_{j \neq j' \in U_k} (Y_{ij} - \mu_Y)(Y_{ij'} - \mu_Y)}{(N_o - 1)(N - 1)S_Y^2} \cong \frac{N_o \frac{\sigma_{Y \text{ inter}}^2}{\sigma_Y^2} - 1}{N_o - 1}$$

e. Considérations pratiques

Comment construire les grappes ?

Combien de grappes faut-il construire ?

Quelques aspects importants :

- i. Grappes homogènes en **inter** et hétérogènes en **intra** (*logique contraire à celle de la stratification, analogue à celle du systématique*).
- ii. Grappes de tailles voisines ou faible dispersion de tailles des grappes.
- iii. Grappes de faibles tailles (ou beaucoup de grappes constituées).
- iv. Un maximum de grappes tirées

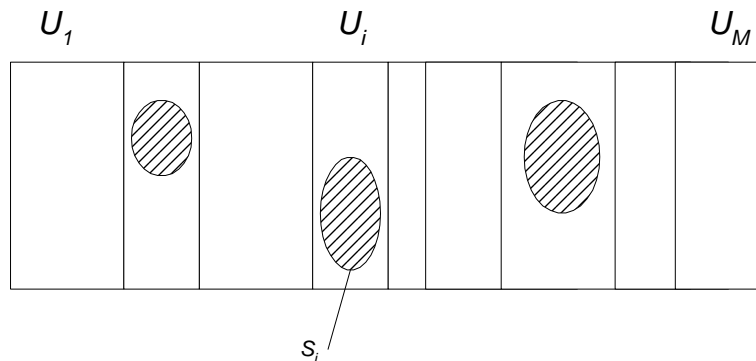
Remarque :

Le tirage systématique est un cas particulier du plan par grappes ($m=1$).

Impossible d'estimer sans biais la variance de l'estimateur du total.

3. Plans à plusieurs degrés

a. Principe général



- i. Un échantillon d'UP, S_{UP} , est tiré dans la population des UP selon un plan de sondage P_{UP}
- ii. Dans chaque U_i sélectionnée ($i \in S_{UP}$), un échantillon S_i d'US est tiré selon un plan de sondage P_{U_i} , indépendamment d'une UP à l'autre.
- iii. L'échantillon complet est : $S = \bigcup_{i \in S_{UP}} S_i$
Sa taille est aléatoire.

Généralisation à plusieurs degrés

- Probabilités d'inclusion des UP dans S_{UP} :

$$\pi_{Ii}, \pi_{Iij}$$

- Probabilités d'inclusion des individus k et l de l'UP i dans S_i , sachant l'UP i tirée :

$$\pi_{k/i}, \pi_{kl/i}$$

- Probabilités d'inclusion des individus k et l dans l'échantillon complet S :

$$\pi_k = \pi_{Ii} \pi_{k/i}, \text{ si } k \in U_i$$

$$\pi_{kl} = \begin{cases} \pi_{Ii} \pi_{kl/i}, & \text{si } k \text{ et } l \in U_i \\ \pi_{Iij} \pi_{k/i} \pi_{l/j}, & \text{si } k \in U_i \text{ et } l \in U_j \end{cases}$$

Estimateur d'Horvitz-Thompson d'un total

$$\hat{t}_y = \sum_{i \in S_{UP}} \frac{\hat{t}_{yi}}{\pi_{Ii}} \quad \text{avec} \quad \hat{t}_{yi} = \sum_{k \in S_i} \frac{y_k}{\pi_{k/i}}$$

$$\text{soit : } \hat{t}_y = \sum_{i \in S_{UP}} \frac{1}{\pi_{Ii}} \sum_{k \in S_i} \frac{y_k}{\pi_{k/i}}$$

- Estimateur sans biais

Estimation d'une moyenne

Si N connu ,

Estimateur de Horvitz-Thompson

$$\hat{u}_y = \frac{\hat{t}_y}{N} = \frac{1}{N} \sum_{i \in S_{UP}} \frac{\hat{t}_{yi}}{\pi_{Ii}}$$

Si N inconnu ,

Estimateur de Hajek

$$\hat{u}_{y,Hajek} = \frac{\hat{t}_y}{\hat{N}} = \frac{\sum_{i \in S_{UP}} \frac{\hat{t}_{yi}}{\pi_{Ii}}}{\sum_{i \in S_{UP}} \frac{\hat{N}_i}{\pi_{Ii}}}$$

b. Tirage SAS à chaque degré

$$\hat{t}_y = \frac{M}{m} \sum_{i \in S_{UP}} \hat{t}_{yi}$$

avec

$$\hat{t}_{yi} = \frac{N_i}{n_i} \sum_{k \in S_i} y_k$$

soit :

$$\hat{t}_y = \frac{M}{m} \sum_{i \in S_{UP}} \frac{N_i}{n_i} \sum_{k \in S_i} y_k$$

Remarque : inutile de connaître N pour estimer le total

Propriétés :

- Estimateur sans biais
- De variance :

$$\text{Var}(\hat{t}_y) = M^2 \left(1 - \frac{m}{M}\right) \frac{S_T^2}{m} + \frac{M}{m} \sum_{i=1}^M N_i^2 \left(1 - \frac{n_i}{N_i}\right) \frac{S_{Yi}^2}{n_i}$$

variance inter-UP + variance intra-UP

$$\text{Où } S_T^2 = \frac{1}{M-1} \sum_{i=1}^M \left(t_{Yi} - \frac{t_Y}{M}\right)^2 \text{ et } S_{Yi}^2 = \frac{1}{N_i-1} \sum_{k=1}^{N_i} (y_k - \mu_{Yi})^2$$

- Estimée sans biais par :

$$\hat{\text{Var}}(\hat{t}_y) = M^2 \left(1 - \frac{m}{M}\right) \frac{s_T^2}{m} + \frac{M}{m} \sum_{i \in S_{UP}} N_i^2 \left(1 - \frac{n_i}{N_i}\right) \frac{s_{Yi}^2}{n_i}$$

$$\text{où } s_T^2 = \frac{1}{m-1} \sum_{i \in S_{UP}} \left(\hat{t}_{Yi} - \frac{\hat{t}_Y}{M}\right)^2 \text{ et } s_{Yi}^2 = \frac{1}{n_i-1} \sum_{k \in S_i} (y_k - \hat{\mu}_{Yi})^2$$

c. Tirages autopondérés

Toutes les unités statistiques de l'échantillon (US) possèdent la même pondération, ou ce qui revient au même, les probabilités d'inclusion des US sont constantes

- i.** Tirage à probabilités égales des UP et des US

- ii.** Tirage des UP proportionnel à leur taille suivi d'un tirage d'un nombre fixe d'US dans chaque UP retenue

Dans ce cas, l'estimateur de la moyenne est la moyenne simple calculée sur l'ensemble des US tirées.

Exemple : enquêtes ménages à l'INSEE

d. Considérations pratiques

Quand utiliser des sondages à plusieurs degrés ?

- Absence de base de sondage ou de mauvaise qualité
- Motifs économiques

Comment améliorer la précision ?

- Avant tout, construire des UP le plus ressemblantes possible entre elles pour limiter les effets de grappes.

Exemple : l'unité ménage est intéressante pour estimer des variables comme le sexe, l'activité, l'âge, etc, mais elle est moins efficace pour étudier le niveau d'instruction, la CS, etc.

- Privilégier le nombre d'UP enquêtées plutôt que le nombre d'US
- Tirer les UP à probabilités inégales
- Stratifier au niveau des UP

4. Exemples des enquêtes ménages INSEE

A- Les grands principes des enquêtes ménages

A.1. Objectif des enquêtes ménages

- Une des missions de l'INSEE
- Rôle du CNIS

Exemples :

Logement, Budget, Revenus et Conditions de Vie, Santé,...

A.2. L'unité échantillonnée

- On construit des échantillons de logements ordinaires
- On atteint les ménages et/ou les individus par l'intermédiaire de leur logement

A3. Principes d'échantillonnage

- Echantillons probabilistes de logements
- Tous les logements principaux possèdent la même probabilité de participer à l'enquête
- Un même logement ne peut-être interrogé pour des enquêtes distinctes entre deux recensements consécutifs

A.4. Modes de collecte

- Enquêtes généralement en face à face
- Quelques enquêtes par téléphone

Coûts de collecte

Stabilité du réseau d'enquêteurs

↪ *localisation des échantillons*

A.5. Procédé : échantillonnage des enquêtes ménages en 2 temps

- i. 1^{ère} phase de localisation des enquêtes
Tirage de bases de sondage intermédiaires
- ii. 2^{ème} phase de tirage des échantillons des enquêtes :
Choix des logements à interroger dans ces bases
Tirages autopondérés, systématiques, stratifiés

A.6. Sources des bases de sondage

- i. Le dernier recensement de la population (*source : INSEE*)
- ii. Les permis de construire pour les logements construits après le dernier RP (*source : Ministère de l'Équipement*)

En l'absence de liste exhaustive de logements à une date quelconque ou de registres de population

B. La construction de l'Echantillon–Maître (EM99)

B.1. Entre contraintes pratiques et désir de précision

- Etablir une réserve de logements localisée pour alimenter la plupart des enquêtes ménages hors DOM, enquête Emploi, enquêtes locales,...
- Optimiser la précision des résultats nationaux

B.2. Taille de l'EM99 : 2 022 889 logements

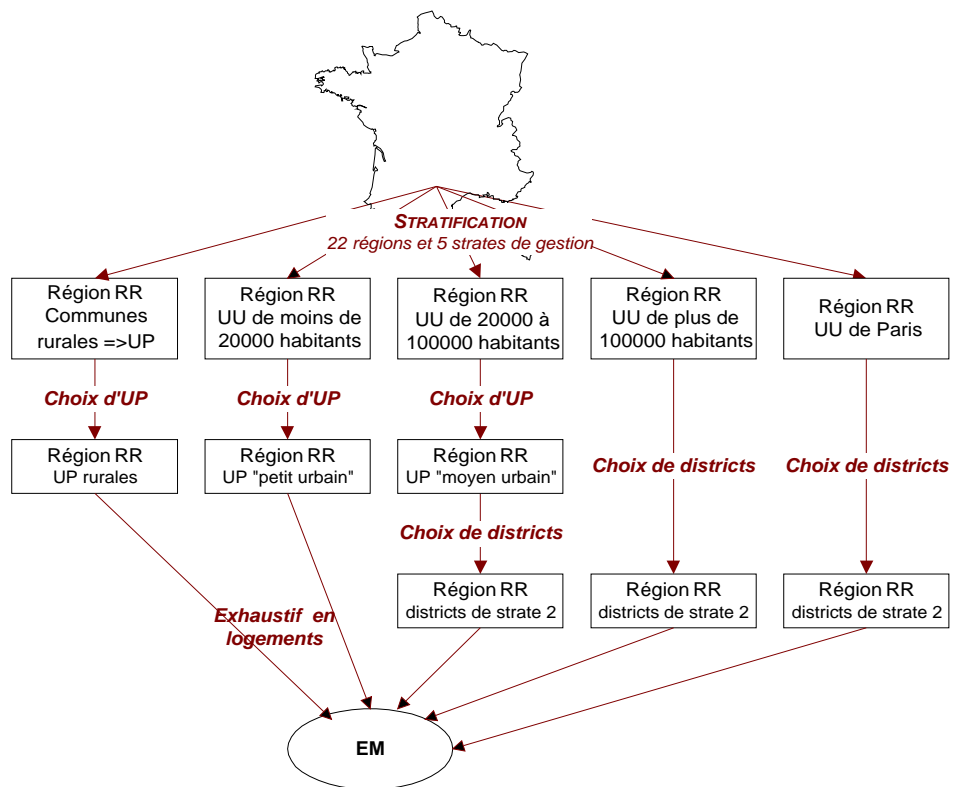
= 7 % des logements recensés en mars 99

Réserve calibrée dans l'optique :

- du nombre moyen d'enquêtes nationales par année,
- de la taille moyenne des enquêtes ménages,
- de la durée de vie de cette réserve

B.3. Constitution de l'EM99

Selon un tirage stratifié à 1 ou 2 degrés selon la strate



□

a. Stratification

Quadrillage du territoire en régions × densité d'habitat

- Collecte des enquêtes en Direction Régionale
- Le degré d'urbanisation explique assez bien le comportement des ménages

↳ *Cette stratification est un élément de qualité.*

b. Les unités primaires

Définition :

- Dans les communes rurales, en petit et moyen urbain, une UP correspond au rayon d'action d'un enquêteur.
- Dans les UU de plus de 100 000 habitants, une UP est une UU.

Regroupements : UP de taille voisine, de profils moyens voisins

Mode de tirage : à probabilités inégales, proportionnelles à la taille

Tirage équilibré selon la méthode du Cube

c. Les unités secondaires (districts)

Tirage uniquement dans les UU de plus de 20000 habitants

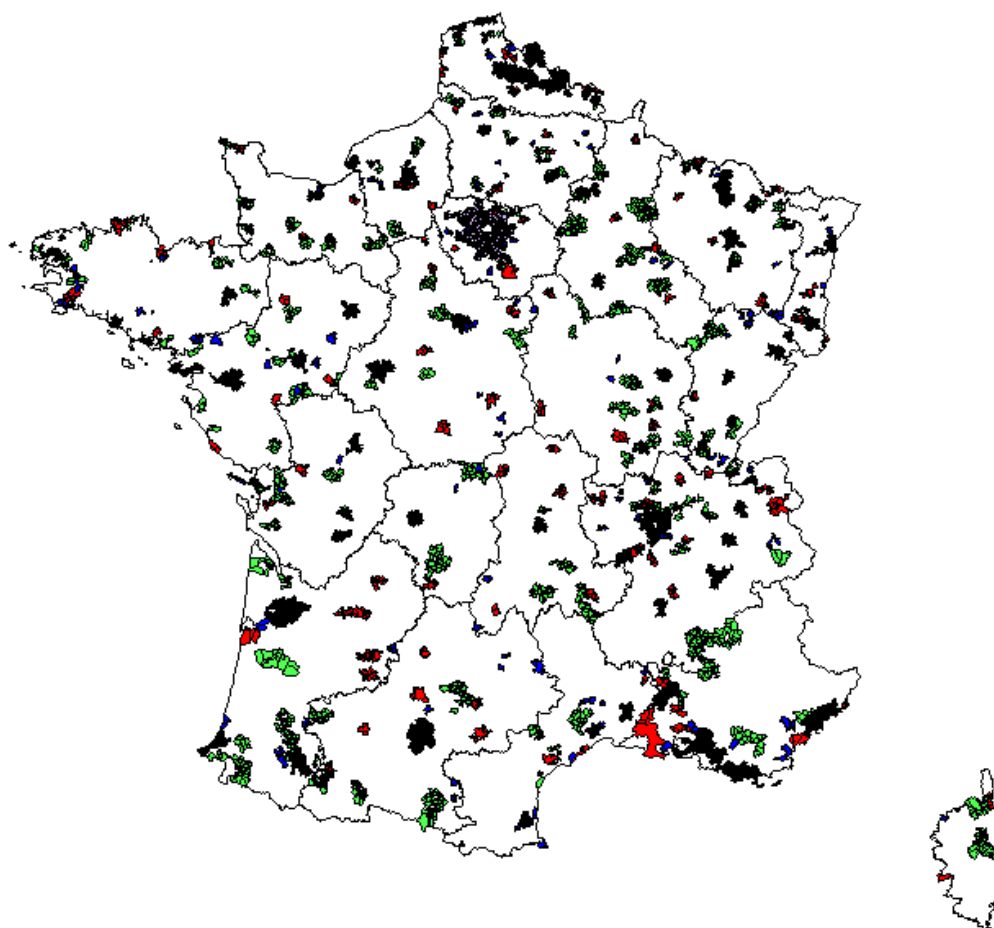
Tirage stratifié par groupes de communes

Tirage à probabilités égales des districts dans l'UU

Tirage équilibré sur l'UU par la méthode du Cube

d. Constitution de l'EM99

Ensemble de tous les logements des UP et des US sélectionnées



Source : INSEE