

# ÉCHANTILLONNAGE DES AGGLOMÉRATIONS DE L'IPC POUR LA BASE 2015\*

20 mars 2015

*Patrick SILLARD*<sup>1(\*\*)</sup>, *Laurence JALUZOT*<sup>2(\*\*)</sup>,

*(\*\*)* Insee, Direction des statistiques démographiques et sociales

## Introduction

L'indice des prix à la consommation (IPC) est un indice de Laspeyres à panier fixe annuel, indice chaîné annuellement. En d'autres termes, l'indice est calculé, dans le courant de l'année, à partir de l'observation des prix des produits constituant le panier de biens et services suivis pour l'année considérée, en référence aux prix des mêmes produits observés au mois de décembre de l'année précédente. Puis cet indice-base décembre est chaîné pour constituer l'IPC proprement dit, en référence à l'année de base (actuellement 1998). L'Insee a décidé de réaliser, sur l'année 2015, un changement de base de l'IPC. Ce changement de base est justifié par une triple motivation : l'introduction de la nouvelle nomenclature européenne de diffusion des indices de prix à la consommation (COICOP), la rénovation du suivi des produits frais destinée à aligner les concepts suivis sur les règlements européens et, la rénovation de l'échantillon géographiques du sondage IPC.

Le présent article vise à présenter les principes retenus pour la sélection des agglomérations de collecte de la base 2015, ainsi que la détermination, par optimisation de variance, du nombre de relevés à réaliser par type de produit et agglomération. Ce texte traite la question de la sélection des points de vente de métropole<sup>3</sup>. Un état des lieux de l'échantillon actuel (base 1998) est proposé en première section.

---

\*Nous remercions Jérôme Accardo, Pascal Ardilly, Sébastien Faivre, Olivier Sautory ainsi que les rapporteurs du Comité du label pour les discussions fructueuses que nous avons eues avec eux dans le cadre de la préparation de cet article.

1. [patrick.sillard@insee.fr](mailto:patrick.sillard@insee.fr)

2. [laurence.jaluzot@insee.fr](mailto:laurence.jaluzot@insee.fr)

3. Dans le cadre de la collecte IPC, des observations sont réalisées dans les DOMs. Du point de vue de l'échantillonnage, chaque DOM est considéré globalement comme un agrégat unique et la collecte est théoriquement réalisée sur l'ensemble du territoire. L'annexe E donne quelques éléments à propos du nombre de relevés, en relation avec la précision des indices calculés.

# 1 État des lieux : localisation des actuels points de vente de l'IPC

Les points de vente de l'IPC sont rattachés, en base 1998, à 96 agglomérations (au sens de l'IPC) métropolitaines. Au mois d'octobre 2013, il y avait, dans l'échantillon IPC, 26 104 points de vente actifs, c'est-à-dire des points de vente dans lesquels il existe un produit qui figure soit dans l'échantillon 2012, soit dans les ordres de recherche<sup>4</sup> des produits pour 2013 et dont le prix a été observé en octobre ou auparavant.

La notion d'agglomération IPC<sup>5</sup> n'est pas confondue avec celle des unités urbaines : elle est même protéiforme puisque dans certains cas, des unités urbaines disjointes appartiennent à la même agglomération IPC. À l'opposé, la collecte de Lyon ne concerne aujourd'hui qu'une petite partie de l'unité urbaine, de sorte que l'on peut considérer que l'agglomération IPC de Lyon est une fraction de l'unité urbaine de Lyon. Pourtant, dans l'esprit des concepteurs de l'échantillon IPC de 1990 (Ardilly & Guglielmetti 1992, Ardilly & Guglielmetti 1993), l'unité urbaine est l'individu statistique sélectionné dans le sondage géographique. Naturellement, cet individu était celui défini à l'époque du tirage, c'est-à-dire tel qu'il était défini dans le recensement de 1990. Quoiqu'il en soit, il est raisonnable de revenir à la notion d'unité urbaine qui est une notion précisément définie et pour laquelle on dispose d'un univers de sondage sur la base du dernier recensement de la population.

La première étape d'identification de l'existant nécessite la localisation des points de vente actifs de métropole. Pour cela, on considère les 26 104 points de vente actifs de métropole en octobre 2013. 95,6% de ces points de vente<sup>6</sup> comportent un numéro SIRET identifié, suite à l'opération de siretisation des points de vente conduite en 2011-2012 par les enquêteurs et les gestionnaires de sites-prix. Depuis lors, ce champ est tenu à jour lors des ouvertures/fermetures de points de vente.

Les points de vente qui comportent un numéro SIRET appariable avec la base SIRENE le sont dans une seconde étape. Ainsi, des 26 104 points de ventes actifs, 86,8% comportent

---

4. i.e. produits nouvellement intégrés en octobre ou novembre dans le champ de la collecte IPC et qui ont vocation à intégrer le panier de biens et services suivis pour l'année suivante.

5. une autre façon de définir l'"agglomération IPC" est de considérer qu'elle correspond à l'agrégat géographique élémentaire, c'est-à-dire celui pour lequel on procède, à partir des prix des produits, à un calcul d'indice élémentaire. Au-delà de ce grain, l'agrégation est réalisée par agrégation de Laspeyres. Afin d'éviter l'ambiguïté du terme d'agglomération IPC, on parle d'*agglos* pour désigner un grain d'agrégat géographique élémentaire à partir de la section 2.

6. Certains points de vente, comme les photomaton, ne sont pas enregistrés en tant que point de vente dans la base SIRENE. Ils ne comportent donc pas d'identifiant SIRET. Afin de distinguer ces cas, le champ SIRET de la base des points de vente de l'IPC se voit attribuer une valeur fictive normalisée (14 caractère "9" successifs). Ces cas sont, dans cette statistique, considérés comme non siretisés. Il y a 2,7% des points de vente actifs en octobre 2013 qui sont concernés par un SIRET fictif. Le complément (1,7%) comporte un champ SIRET vide.

un écho dans la base SIRENE. Ces points de vente sont immédiatement localisés, puisque la base SIRENE dispose pour chaque point de vente, de l'adresse et de l'identifiant de commune INSEE.

Pour les autres, on utilise un algorithme de reconnaissance de chaînes de caractères qui permet d'identifier la commune d'implantation du point de vente. Pour le reste (environ 400 points de vente), un codage manuel a été réalisé. Au final, l'ensemble des 26 104 points de vente actifs fin 2013 sont localisés à la commune. Il est dès lors possible de localiser la collecte de l'IPC et d'examiner sa répartition sur les unités urbaines métropolitaines. Cette analyse est réalisée par Jaluzot (2014). Elle est reproduite en annexe F.

## 2 Les principes de base de l'échantillon IPC

Faivre (2012) a proposé un schéma renouvelé dans la perspective de l'introduction des données de caisses dans l'IPC. Même si la base 2015 ne verra pas, *stricto sensu*, la substitution de données de caisses à de la collecte terrain, ces travaux permettent d'orienter la réflexion.

On limite ici le raisonnement au champ de la consommation localisée, c'est-à-dire celle pour laquelle une répartition géographique est connue et les prix varient spatialement. Le complément relève du champ des tarifs, c'est-à-dire de la consommation dont les prix ne varient pas spatialement ou pour laquelle on dispose *a priori* d'une information précise sur la variabilité spatiale des prix (exemple : le transport ferroviaire de voyageurs), ou celle pour laquelle la répartition spatiale est très concentrée par rapport à celle de la population générale (par exemple les achats de planches à voile). La définition précise du champ de la consommation localisée est un exercice, partie intégrante de la définition des variétés. Préalablement à la définition de l'échantillon de l'IPC, il convient donc de séparer les variétés dont la logique de collecte correspond à une consommation localisée de celles dont on considère qu'elles correspondent à une consommation non localisée. Les prix correspondant à une consommation localisée sont en principe collectés sur le terrain par enquête, tandis que la consommation non localisée l'est de manière centralisée.

Pour le tirage de l'échantillon de produits dont les prix sont collectés sur le terrain (consommation localisée), les opérations suivantes doivent être réalisées :

1. Fixer les  $H$  strates parmi lesquelles on va sélectionner les agglos. Ces strates sont des croisements de zones géographiques (ZEAT<sup>7</sup> pour l'échantillon IPC base 1998) et

---

7. Les ZEAT sont des regroupements de régions administratives (code entre parenthèses) : REGION PARISIENNE (1) - Ile de France ; BASSIN PARISIEN (2) - Bourgogne, Centre, Champagne-Ardenne, Basse et Haute Normandie, Picardie ; NORD (3) - Nord Pas-de-Calais ; EST (4) - Alsace, Franche-Comté, Lorraine ; OUEST (5) - Bretagne, Pays de la Loire, Poitou-Charentes ; SUD-OUEST (7) - Aquitaine, Limousin, Midi-Pyrénées ; CENTRE-EST (8) - Auvergne, Rhône-Alpes ; MEDITERRANEE (9) - Languedoc-Roussillon, Provence-Alpes-Côte d'Azur, Corse

de types d'agglomérations ( pour l'échantillon IPC bases 1990 et 1998 – voir Ardilly & Guglielmetti (1992) et Faivre (2012)).

2. Déterminer le nombre  $n_h$  d'agglos retenues pour chaque strate  $h$ ,  $h \in \{1, \dots, H\}$ .

3. Pour une agglo  $a$  retenue, fixer le nombre de relevés  $n_{a,v}$  à réaliser pour la variété  $v$ .

En première approximation, on souhaite que le nombre total de relevés  $\sum_{v,a} n_{a,v}$  et que le nombre d'agglos  $\sum_h n_h$  restent les mêmes qu'en base 1998. Dans un second temps, on cherchera à optimiser ces nombres selon un critère qui sera précisé par la suite.

Le plan de sondage est le suivant :

- i. stratifié avec détermination du nombre  $n_h$  d'agglos pour chaque strate  $h$  ;
- ii. au sein de chaque strate  $h \in \{1, \dots, H\}$ , on procède à un sondage à deux degrés :
  - a- on sélectionne  $n_h$  agglos selon un plan à préciser ;
  - b- au sein de l'agglo  $a$  sélectionnée, on détermine un nombre de relevés  $n_{a,v}$  à réaliser pour une variété  $v$  donnée et on sélectionne les relevés en question selon un plan à préciser.

### 3 Variable à estimer et estimateur

*Notations : Dans toute cette partie, les indices se réfèrent à une variété donnée. L'indice de variété  $v$  est donc omis dans cette section.*

L'indice d'ensemble est le fruit d'une agrégation de Laspeyres. En d'autres termes, quelle que soit la partition  $J$  de la consommation considérée, l'indice d'ensemble  $I$ , étant donnés les indices élémentaires  $I_j$  calculés pour les ensembles  $j \in J$ , vaut :

$$I = \sum_{j \in J} w_j I_j$$

où  $w_j$  est le poids de l'ensemble  $j$  dans la dépense de consommation des ménages. Par définition  $\sum_{j \in J} w_j = 1$ .

Les indices élémentaires calculés dans l'IPC sont des indices calculés sur le couple  $(v, a)$ . L'indice d'ensemble vaut donc, si on imagine une partition  $A$  du territoire en somme d'agglos (on omet le  $v$  de la variété),

$$I = \sum_{a \in A} w_a I_a \tag{1}$$

où  $w_a$  est le poids de l'agglo  $a$  dans la consommation des ménages (*i.e.*  $\sum_a w_a = 1$ ).

Cet indice peut être estimé par sondage, d'une part en sélectionnant des agglos et d'autre part, en déterminant un estimateur de  $I_a$  à partir de l'observation de prix de produits appartenant à la variété d'intérêt  $v$  au sein de l'agglo  $a$ . Si la probabilité d'inclusion

des produits contribuant à l'estimateur de  $I_a$  est indépendante de celle de l'agglomération  $a$  dans l'échantillon (cas d'un sondage à deux degrés), alors la variance d'ensemble est la somme, d'une part, d'une variance intra spécifique à  $a$  et ne dépendant que du nombre de produits  $n_a$  sélectionnés dans l'échantillon sur lequel se fonde l'estimateur de  $I_a$  et d'autre part, d'une variance extra ne dépendant que des indices estimés  $\hat{I}_a$  et du nombre  $n_a$  d'agglomérations sélectionnées. En effet, notons :

- $\mathbf{1}_{\mathcal{A}}(a)$  la variable de Cornfield caractérisant l'appartenance de l'agglomération  $a$  à l'échantillon sélectionné  $\mathcal{A}$ ,  $\pi_a = \Pr(\mathbf{1}_{\mathcal{A}}(a) = 1)$  et  $\pi_{a,a'} = \Pr(\{\mathbf{1}_{\mathcal{A}}(a) = 1\} \cap \{\mathbf{1}_{\mathcal{A}}(a') = 1\})$ ;
- $\mathbf{1}_{\mathcal{S}_a}(i)$  la variable de Cornfield caractérisant l'appartenance du produit  $i$  à l'échantillon  $\mathcal{S}_a$  de produits sélectionnés dans l'agglomération  $a$ .

L'estimateur de l'indice d'ensemble  $I$  est estimé par  $\hat{I}^{\mathcal{A}}$ , fondé sur l'échantillon  $\left(\mathcal{A}, \bigcup_{a \in \mathcal{A}} \mathcal{S}_a\right)$  et de la forme :

$$\hat{I}^{\mathcal{A}} = \sum_{a \in \mathcal{A}} \omega_a \hat{I}_a \quad (2)$$

où

$$\begin{cases} \hat{I}_a &= \sum_{i \in \mathcal{S}_a} \omega_i \frac{p_i}{p_i^0} \\ \hat{I}_a &= \prod_{i \in \mathcal{S}_a} \left(\frac{p_i}{p_i^0}\right)^{\omega_i} \end{cases} \quad \text{et} \quad \sum_{a \in \mathcal{A}} \omega_a = 1$$

Naturellement, pour que l'indice (2) soit sans biais, il est nécessaire en outre que  $\mathbb{E}(\hat{I}_a) = I_a$ .

Il est possible d'écrire  $\hat{I}^{\mathcal{A}}$  à l'aide des variables de Cornfield :

$$\hat{I}^{\mathcal{A}} = \sum_{a \in \mathcal{A}} \omega_a \hat{I}_a \times \mathbf{1}_{\mathcal{A}}(a)$$

Moyennant quoi,  $\mathbb{E}(\hat{I}^{\mathcal{A}}) = I$ , où  $I$  est donné par la formule (1) lorsque :

$$\omega_a = \frac{w_a}{\pi_a}$$

L'estimateur  $\hat{I}^{\mathcal{A}}$  de  $I$  correspond donc à la formule :

$$\hat{I}^{\mathcal{A}} = \sum_{a \in \mathcal{A}} \frac{w_a \hat{I}_a}{\pi_a} \quad (3)$$

les poids  $\frac{w_a}{\pi_a}$  étant de somme unitaire sur<sup>8</sup>  $\mathcal{A}$ .

La variance de  $\hat{I}^{\mathcal{A}}$  s'écrit alors :

$$\text{var}(\hat{I}^{\mathcal{A}}) = \sum_{a \in \mathcal{A}} \left(\frac{w_a}{\pi_a}\right)^2 \text{var} \left[ \hat{I}_a \times \mathbf{1}_{\mathcal{A}}(a) \right] + \sum_{\substack{a \neq a' \\ (a, a') \in \mathcal{A}^2}} \frac{w_a w_{a'}}{\pi_a \pi_{a'}} \text{cov} \left[ \hat{I}_a \times \mathbf{1}_{\mathcal{A}}(a), \hat{I}_{a'} \times \mathbf{1}_{\mathcal{A}}(a') \right]$$

Les plans de sondage consistant à sélectionner  $\mathcal{A}$  dans  $A$  et  $\mathcal{S}_a$  dans l'ensemble des produits de l'agglomération  $a$  étant indépendants, les variables  $\mathbf{1}_{\mathcal{A}}(a)$  et  $\hat{I}_a$  sont indépendantes.

8. et sont, le cas échéant, renormalisés pour ce faire.

En revanche, les variables  $\mathbf{1}_{\mathcal{A}}(a)$  et  $\mathbf{1}_{\mathcal{A}}(a')$  sont, sauf hypothèse complémentaire, liées. Il découle de ce constat et de l'expression précédente que :

$$\begin{aligned} \text{var}(\hat{I}^{\mathcal{A}}) &= \overbrace{\sum_{a \in A} \left(\frac{w_a}{\pi_a}\right)^2 \pi_a I_a^2 (1 - \pi_a) + \sum_{\substack{a \neq a' \\ (a, a') \in A^2}} \frac{w_a w_{a'}}{\pi_a \pi_{a'}} I_a I_{a'} (\pi_{aa'} - \pi_a \pi_{a'})}^{V1} \\ &\quad + \underbrace{\sum_{a \in A} \left(\frac{w_a}{\pi_a}\right)^2 \pi_a \text{var}(\hat{I}_a)}_{V2} \end{aligned} \quad (4)$$

On constate que la variance précédente se décompose en une somme de deux termes : le terme  $V1$ , appelé variance de premier degré, ne dépend que du tirage d'agglos, c'est la variance inter-agglos ; le terme  $V2$ , appelé variance de second degré, ne dépend que du tirage des produits dans les agglos sélectionnées, c'est la variance intra-agglos. Cette décomposition est une propriété classique des sondages à deux degrés (Tillé 2001, Ardilly 2006). S'agissant de  $V1$ , tout se passe comme s'il s'agissait de la variance du  $\pi$ -estimateur du total (pour un échantillon  $\mathcal{A}$  donné) :

$$\hat{y}^{\mathcal{A}} = \sum_{a \in \mathcal{A}} \frac{w_a I_a}{\pi_a} \quad (5)$$

où la caractéristique individuelle sommée est  $y_a = w_a I_a$  (voir Tillé *op. cit.*, page 35). Ainsi, pour cette partie de l'estimateur  $\hat{I}^{\mathcal{A}}$ , on peut appliquer les résultats de théorie des sondages sur le  $\pi$ -estimateur du total en se rappelant que la variable sommée n'est pas  $I_a$  mais  $w_a I_a$ .

## 4 Optimisation de l'échantillon

*Notations* : À la différence du paragraphe 3, dans ce paragraphe, lorsqu'une quantité se réfère à une variété, ceci est explicitement indiqué en indice ( $v$  pour variété).

L'idée de base de l'optimisation consiste à déterminer les probabilités d'inclusion des agglos  $\pi_a$  et le nombre d'observations  $n_{a,v}$ , pour une variété donnée, de sorte que  $\text{var}(\hat{I}^{\mathcal{A}})$  donnée par l'expression (4) soit minimale. Naturellement, il faut procéder à quelques approximations pour pouvoir accéder à une expression de la variance qui soit utilisable analytiquement. Nous effectuons donc les hypothèses suivantes :

— **H1** – On procède à l'optimisation en deux étapes distinctes :

1. Les probabilités d'inclusion  $\pi_a$  des agglos sont déterminées *ex-ante* sur la base d'un critère global portant sur l'ensemble des variétés, comme proposé par Ardilly & Guglielmetti (1993).
2. Conditionnellement aux  $\pi_a$  de l'étape 1, on détermine pour la variété  $v$ , le nombre d'observations  $n_{a,v}$  dans l'agglomération  $a$  de façon à ce que le nombre total

d'observations de la variété  $v$  soit  $n_v$  et que la variance, caractérisée par le facteur  $V^2$  dans l'expression (4) soit minimale. Le nombre d'observations  $n_{a,v}$  est donc à ce stade conditionnel aux  $\pi_a$  et à  $n_v$ .

- **H2** – On suppose que le plan de sondage de l'échantillon d'observations pour une variété  $v$  dans une agglo  $a$  donnée est un sondage aléatoire simple avec remise de taille fixée égale à  $n_{a,v}$ . Cette hypothèse est acceptable dans la mesure où le taux de sondage est en général inconnu (on ne connaît pas le nombre de produits rattachés à une variété vendus dans une agglo donnée) et supposé très faible. Moyennant quoi,  $\text{var}(\hat{I}_{a,v}) = \frac{\sigma_{a,v}^2}{n_{a,v}}$ , où  $\sigma_{a,v}^2$  est un paramètre de variance<sup>9</sup> que l'on supposera indépendant de  $n_{a,v}$ .
- **H3** –  $w_v$  représente le poids de la variété  $v$  dans la dépense de consommation des ménages.  $\sum_v w_v = 1$ .  $w_{a,v}$  correspond au poids de l'agglo  $a$  dans la dépense de consommation des ménages. Moyennant quoi,  $\sum_{a \in A} w_{a,v} = w_v$ ,  $\sum_v w_{a,v} = w_a$  où  $w_a$  est le poids de l'agglo  $a$  dans la dépense de consommation des ménages, et  $w_{a,v}/w_v$  est le poids de l'agglo  $a$  dans la dépense de consommation des ménages de la variété  $v$ .

Sous les hypothèses **H1**, **H2** et **H3** précédentes, il est possible de déterminer le nombre d'observations requis pour le couple  $(a, v)$ . En effet, la variance à minimiser de l'estimateur  $\hat{I}_v$  de l'indice de la variété  $v$ , issu de l'agrégation des indices de var-agglo, est<sup>10</sup> :

$$\varphi(\mathbf{n}_v) = \frac{1}{w_v^2} \sum_{a \in A} \left( \frac{w_{a,v}}{\pi_a} \right)^2 \pi_a \frac{\sigma_{a,v}^2}{n_{a,v}} \quad (6)$$

où  $\mathbf{n}_v$  est le vecteur de composantes  $(n_{a,v})_{a \in A}$ . On peut donc la minimiser sous la condition que, en moyenne (i.e. en tenant compte de la probabilité de sélection de l'agglo  $a$ ), le nombre d'observations pour la variété  $v$  soit égal à une quantité exogène  $n_v$ , c'est-à-dire que  $\sum_{a \in A} \pi_a n_{a,v} = n_v$ . Pour résumer, on définit  $\mathbf{n}_v$  comme<sup>11</sup> :

$$\mathbf{n}_v = \underset{\mathbf{v}}{\text{argmin}} \left\{ \sum_{a \in A} \left( \frac{w_{a,v}}{\pi_a} \right)^2 \pi_a \frac{\sigma_{a,v}^2}{\nu_a} \left| \sum_{a \in A} \pi_a \nu_a = n_v \right. \right\} \quad (7)$$

Finalement, on obtient<sup>12</sup> :

$$n_{a,v} = \frac{w_{a,v} \sigma_{a,v} / \pi_a}{\sum_{a \in A} w_{a,v} \sigma_{a,v}} \times n_v \quad (8)$$

9. Par exemple, pour une variété homogène, un estimateur de  $\sigma_{a,v}^2$  est  $\hat{\sigma}_{a,v}^2 = \frac{1}{n_{a,v}-1} \sum_{i \in (a,v)} \left( \frac{p_i}{p_i^0} - \hat{I}_{a,v} \right)^2$  où  $\hat{I}_{a,v} = \frac{1}{n_{a,v}} \sum_{i \in (a,v)} \frac{p_i}{p_i^0}$ ,  $p_i$  étant le prix courant du produit  $i$  et  $p_i^0$  son prix au mois de base.

10. En effet, dans la formule (4), l'indice d'ensemble (i.e. pour une variété donnée) est la somme pondérée d'indices élémentaires d'agglos, les poids utilisés dans cette somme étant de somme unitaire. Par rapport aux notations introduites à l'hypothèse **H3**, il convient dès lors de considérer des poids normalisés, ce qui est le cas des  $w_{a,v}/w_v$  dont la somme sur  $A$  est unitaire.

11. Le facteur  $1/w_v^2$  est omis à ce stade du programme d'optimisation car sa présence ne modifie pas la solution.

12. voir annexe B avec  $a_i \rightarrow w_{a,v} \sigma_{a,v} / \sqrt{\pi_a}$ ,  $\alpha_i \rightarrow \sqrt{\pi_a}$  et  $K \rightarrow n_v$

Ce résultat conduit à un nombre d'observations par couple  $(a, v)$  inconditionnel, c'est-à-dire qu'il ne dépend pas du choix de l'échantillon  $\mathcal{A}$ . En contrepartie, le nombre  $n_{a,v}$  est aléatoire. On pourrait à l'inverse calculer un nombre optimal  $n_{a,v}^{\mathcal{A}}$  conditionnellement à l'échantillon  $\mathcal{A}$ . Dans ce cas,  $n_{a,v}^{\mathcal{A}}$  est défini par<sup>13</sup> :

$$\mathbf{n}_v^{\mathcal{A}} = \underset{\nu}{\operatorname{argmin}} \left\{ \sum_{a \in \mathcal{A}} \left( \frac{w_{a,v}}{\pi_a} \right)^2 \frac{\sigma_{a,v}^2}{\nu_a} \mid \sum_{a \in \mathcal{A}} \nu_a = n_v \right\} \quad (9)$$

L'intérêt de cette seconde approche est qu'il n'y a plus d'aléa sur le nombre d'observations par couple  $(a, v)$  dès que  $\mathcal{A}$  est sélectionné : le nombre total d'observations est fixé et égal à la somme des observations réalisées sur l'ensemble des agglos de l'échantillon  $\mathcal{A}$ .

On montre<sup>14</sup> que :

$$n_{a,v}^{\mathcal{A}} = \frac{w_{a,v} \sigma_{a,v} / \pi_a}{\sum_{a \in \mathcal{A}} w_{a,v} \sigma_{a,v} / \pi_a} \times n_v \quad (10)$$

Les deux estimateurs (8) et (10) sont en tout point comparables. En effet, dans le premier, la présence de  $\pi_a$  au numérateur permet de dilater le poids de l'agгло  $a$  en tenant compte du fait que si elle est sélectionnée, elle "représente" davantage qu'elle-même (car  $\pi_a \leq 1$ ) ; au dénominateur la somme des poids est réalisée sur tout l'univers, donc la probabilité d'inclusion n'apparaît pas. À l'inverse, pour le second, la somme au dénominateur étant réalisée sur l'échantillon sélectionné  $\mathcal{A}$ , il convient là-aussi de dilater les poids des agglos sélectionnées pour tenir compte du fait qu'elles représentent davantage qu'elles-mêmes. En espérance, les deux estimateurs sont cohérents.

À ce stade, il nous reste à préciser les quantités  $\pi_a$  (et en particulier le nombre d'agglos retenues) et le nombre  $n_v$  d'observations pour chaque variété. Pour préciser le nombre total d'observations  $n = \sum_v n_v$ , nous effectuons les hypothèses suivantes, complémentaires aux précédentes :

- **H4** – le nombre d'observations total (pour le champ de la consommation localisée) résulte d'une optimisation de la variance de second degré sous contrainte de coût. Le coût d'ensemble est  $C$  et le coût d'un relevé élémentaire pour la variété  $v$  est  $c_v$ . Ces paramètres sont exogènes. Les indices de variété obtenus par agrégation des indices pour les couples  $(a, v)$  sont indépendants. Le nombre d'observations effectué pour le couple  $(a, v)$  est celui résultant de l'optimisation précédente, et donné par la relation (8).

On note  $c_v$  le coût élémentaire d'un relevé de la variété  $v$ ,  $C$  le coût total de l'opération de collecte. L'indice d'ensemble résulte d'une agrégation des indices élémentaires de variétés  $I_v$  conformément à la relation<sup>15</sup> :

$$I = \sum_v w_v I_v$$

---

13. Voir aussi note N°11.

14. voir annexe B avec  $a_i \rightarrow w_{a,v} \sigma_{a,v} / \pi_a$ ,  $\alpha_i \rightarrow 1$  et  $K \rightarrow n_v$

15. Pour être complet, avec ces notations, nous avons en particulier  $\sum_v w_v = 1$  et  $\sum_{a \in A} w_{a,v} = w_v$ . Même si dans le cas général,  $w_{a,v}$  est quelconque, il pourra, par exemple par secteur de variété ou strate



où  $w_v$  est le poids de la variété  $v$  dans la dépense de consommation des ménages. Il en découle que la variance de second degré pour l'indice d'ensemble vaut<sup>16</sup> ( $\widehat{V2}^{\mathcal{A}} = \text{var}^{2d}(\widehat{I}^{\mathcal{A}})$ ) :

$$\widehat{V2}^{\mathcal{A}} = \sum_v w_v^2 \times \left\{ \frac{1}{w_v^2} \sum_{a \in \mathcal{A}} \left( \frac{w_{av}}{\pi_a} \right)^2 \frac{1}{n_v} \sigma_{a,v}^2 \left[ \frac{w_{av} \sigma_{av} / \pi_a}{\sum_{a \in \mathcal{A}} w_{av} \sigma_{av} / \pi_a} \right]^{-1} \right\}$$

Après simplification, cette variance s'écrit :

$$\widehat{V2}^{\mathcal{A}} = \sum_v \frac{(\kappa_v^{\mathcal{A}})^2}{n_v} \quad (11)$$

où

$$\kappa_v^{\mathcal{A}} = \sum_{a \in \mathcal{A}} w_{av} \sigma_{av} / \pi_a$$

On détermine le vecteur  $\mathbf{n} = (n_1, \dots, n_V)$  en minimisant la variance sous contrainte de coût global de collecte :

$$\mathbf{n} = \underset{\mathbf{v}}{\text{argmin}} \left\{ \sum_v \frac{(\kappa_v^{\mathcal{A}})^2}{n_v} \mid \sum_v c_v n_v = C \right\} \quad (12)$$

On montre<sup>17</sup>, pour tout  $v \in \{1, \dots, V\}$ , que l'estimateur du nombre total d'observations par variété  $v$ , conditionnel à l'échantillon d'agglos sélectionnées  $\mathcal{A}$  est :

$$n_v^{\mathcal{A}} = \frac{C}{c_v} \times \frac{\kappa_v^{\mathcal{A}} \sqrt{c_v}}{\sum_v \kappa_v^{\mathcal{A}} \sqrt{c_v}} \quad (13)$$

On en déduit que la variance totale optimisée vaut :

$$\widehat{V2}^{\mathcal{A}} = \frac{1}{C} \left( \sum_v \kappa_v^{\mathcal{A}} \sqrt{c_v} \right)^2 \quad (14)$$

C'est la borne inférieure de la variance que l'on peut obtenir avec l'échantillon  $\mathcal{A}$  et pour un coût de collecte  $C$ .

## 5 Les formules de variance des micro-indices

Ces formules ont été examinées par Ardilly & Guglielmetti (1992) dans un contexte où les formules élémentaires étaient légèrement différentes de celles pratiquées aujourd'hui.

---

d'agglos  $\mathcal{S}(a, v)$ , être compris comme le produit de poids marginaux :  $w_{a,v} = w_a^{S(a,v)} \times w_v$ . Moyennant quoi, on a alors  $\sum_{a \in A} w_a^{S(a,v)} = 1$ .

16. On repart ici de la formule (6) et on somme en pondérant chaque indice de variété  $\hat{I}_v$  par son poids  $w_v$  (mis au carré dans la formule de variance). On s'appuie ici sur l'expression de l'estimation de la variance de  $\hat{I}_v$  fondée sur l'échantillon  $\mathcal{A}$ , soit les relations (9) et (10).

17. voir annexe B avec  $a_i \rightarrow \kappa_v^{\mathcal{A}}$ ,  $\alpha_i \rightarrow \sqrt{c_v}$  et  $K \rightarrow C$

En particulier, à l'époque, la formule appliquée pour les variétés hétérogènes consistait en une moyenne arithmétique de rapports de prix ; la formule pratiquée aujourd'hui est une moyenne géométrique de rapports de prix. Les formules ont été revues par Petit (2014) et Balcone (2014). Les formules proposées ici sont démontrées en annexe C. Elles correspondent aux formules établies pour la variance intra au niveau des micro-indices par Petit (2014) et Balcone (2014).

## 5.1 Variétés homogènes

Un indice de prix d'agglomération pour une variété homogène donnée (on omet ici l'indice de variété et celui d'agglomération) est calculé à l'aide d'une formule d'indice de Laspeyres : soient  $(p_i^t)_{i \in \{1, \dots, n\}}$  les prix des  $n$  biens composant le panier d'intérêt à l'instant  $t$  et  $(p_i^0)_{i \in \{1, \dots, n\}}$  les prix de ces mêmes biens à l'instant de base, alors l'indice estimé est<sup>18</sup> :

$$\hat{I} = \frac{\frac{1}{n} \sum_{i=1}^n p_i^t}{\frac{1}{n} \sum_{i=1}^n p_i^0}$$

Sous l'hypothèse où les  $n$  biens ont été sélectionnés par sondage aléatoire simple sans remise dont le taux de sondage est négligeable, alors on montre (voir annexe C) qu'un estimateur de la variance (intra) est :

$$\widehat{\text{var}}(\hat{I}) = \left( \frac{1}{\hat{P}^0} \right)^2 \frac{1}{n(n-1)} \sum_{i \in \mathcal{S}} \left( p_i^t - \hat{I} p_i^0 \right)^2 \quad (15)$$

où  $\hat{P}^0$  est la moyenne empirique des prix de base sur l'échantillon de produits retenus.

Par rapport au paramètre  $\sigma_{av}$  évoqué au paragraphe 4 (hypothèse **H2**), il suffit de considérer l'expression précédente multipliée par le nombre d'observations qui concourent à la somme<sup>19</sup> et d'en prendre la racine carrée.

## 5.2 Variétés hétérogènes

Un indice de prix d'agglomération pour une variété hétérogène donnée (on omet ici l'indice de variété et celui d'agglomération) est calculé à l'aide d'une formule d'indice de Laspeyres géométrique. Avec les mêmes notations qu'au paragraphe précédent, l'indice estimé est :

$$\hat{I} = \left( \prod_{i=1}^n \frac{p_i^t}{p_i^0} \right)^{1/n}$$

18. Dans tout ce texte, on omet l'éventuel facteur multiplicatif 100 qui permet, en l'absence de variation de prix, de fixer le niveau de base à 100. En l'absence de ce facteur, le niveau de base est égal à 1.

19. En effet, on définit  $\sigma_{a,v}$  de sorte que  $\text{var}(\hat{I}_{a,v}) = \frac{\sigma_{a,v}^2}{n_{a,v}}$ , où  $n_{a,v}$  est le nombre d'observations de l'échantillon pour la variété  $v$  dans l'agglomération  $a$ .

Sous l’hypothèse où les  $n$  biens ont été sélectionnés par sondage aléatoire simple sans remise dont le taux de sondage est négligeable, alors on montre (voir annexe C) qu’un estimateur de la variance (intra) est :

$$\widehat{\text{var}}(\hat{I}) = \frac{\hat{I}^2}{n(n-1)} \sum_{i \in \mathcal{I}} \left[ \ln \left( \frac{p_i^t}{p_i^0} \right) - \ln \hat{I} \right]^2 \quad (16)$$

Comme précédemment, s’agissant du paramètre  $\sigma_{av}$  évoqué au paragraphe 4 (hypothèse **H2**), il suffit de considérer l’expression précédente multipliée par le nombre d’observations qui concourent à la somme et d’en prendre la racine carrée.

### 5.3 Variétés produits frais

Le calcul des indices de produits frais répond à une mécanique assez différente de celle utilisée pour les variétés ordinaires, puisque l’agrégation est réalisée au niveau poste-strate, et non au niveau var-agglo. Par ailleurs, les produits frais pèsent peu dans l’échantillon et sont observés dans la plupart des agglos IPC. L’enjeu à les inclure dans l’analyse, notamment pour la sélection des agglos, est donc relativement faible tandis que leur intégration contribuerait à accroître la complexité de l’algorithme. Comme Ardilly & Guglielmetti (1992), nous ne prenons pas en compte les produits frais dans cette analyse.

## 6 Le plan du sondage de premier degré

Comme évoqué ci-dessus (paragraphe 3), la variance du sondage de premier degré s’apparente à celle d’un  $\pi$ -estimateur du total de la variable d’agglo  $w_{a,v}I_{a,v}$  où  $w_{a,v}$  est le poids du couple  $(a, v)$  dans la dépense de consommation des ménages. Il est clair que cette variable est corrélée à toute variable de taille d’agglo<sup>20</sup>.

En conséquence, il serait justifié (voir Tillé (2001)) d’adopter un plan de sondage des agglos avec probabilités d’inclusion inégales proportionnelles à une variable de taille. En pratique, ce constat donne un objectif à suivre pour le choix des probabilités d’inclusion des agglos dans l’échantillon. Un tel plan de sondage permet d’assurer une précision optimisée du  $\pi$ -estimateur (Tillé 2001) par rapport par exemple à un sondage à probabilités d’inclusion égales. Cependant, le choix des agglos est aussi dicté par l’échantillon antérieur d’agglos. Donc le plan de sondage ne peut pas suivre *stricto-sensu* un plan qui serait défini de manière purement *ad hoc*, indépendamment de l’échantillon actuel d’agglos.

---

20. De ce point de vue, différentes variables sont susceptible de correspondre à  $w_{a,v}$  et le choix, sur le plan des principes du sondage, n’est pas essentiel pourvu que la variable soit bien liée à la taille des agglos : on peut penser à une variable de taille liée à la démographie ou à une variable prenant en compte la géographie de la consommation des ménages. Des éléments de cadrage sont proposés par Jaluzot (2014) – voir annexe F.

Le sondage destiné à la sélection des agglos est un sondage stratifié par zone géographique (ZG) croisées avec les classes de tailles d’agglos (TA). Ces strates définissent une partition de l’espace géographique. On note les strates  $h \in \{1, \dots, H\}$ .

Pour cet exercice, on distingue deux étapes :

1. détermination du nombre  $n_h$  d’agglos par strate en minimisant un critère une fonction de variance<sup>21</sup> d’ensemble (de premier degré) sous contrainte de coût ;
2. sondage au sein de chaque strate avec probabilité d’inclusion approximativement proportionnelle à une variable de taille.

Ces deux étapes sont précisées ci-après.

## 6.1 Le nombre d’agglos par strate

Le nombre d’agglos par strate doit être déterminé de sorte que la variance de premier niveau du sondage à deux degrés soit minimale. Au sein de chaque strate, les agglos sont sélectionnés par sondage aléatoire avec probabilité d’inclusion proportionnelle à une variable de taille. On notera  $\pi_a$  la probabilité d’inclusion dans l’échantillon d’une agglo  $a$ . Elle appartient à une strate  $h$  donnée. Dans la mesure où le nombre d’agglos par strate, noté  $n_h$  pour la strate  $h$ , est la variable d’intérêt à ce stade, il convient de déterminer ce nombre de sorte que la variance d’ensemble soit minimale, sous contrainte d’un coût lié au nombre d’agglos retenu pour chaque strate. Dans les étapes exposées au paragraphe 4, la liste des agglos est fixée et les coûts appliqués pour l’optimisation sont ceux de la collecte de relevés, conditionnellement à la liste des agglos et au nombre total de relevés à réaliser sur l’ensemble des agglos. Les coûts qui doivent être pris en compte dans l’étape présente, qui consistent à ventiler un nombre total de relevés fixé sur des agglos en plus ou moins grand nombre, s’apparentent à des coûts d’approche, pour les enquêtes ménages. Faute d’information particulière, il n’y a pas lieu de considérer que ces coûts diffèrent d’une agglo à l’autre. Nous supposons donc que la fonction de coût est simplement la somme des nombres d’agglos obtenus sur chacune des strates (soit un coût unitaire par agglo incluse dans l’échantillon). Cette hypothèse est également celle retenue par Ardilly & Guglielmetti (1992) et Faivre (2012) dans leurs travaux.

Comme vu à la relation<sup>22</sup> (5), on estime la variance du  $\pi$ -estimateur du total de la variable  $w_{a,v}J_{a,v}$  sur la strate. Ardilly (2014) propose d’utiliser, pour le calcul de variance,

---

21. On parle ici de fonction de variance d’ensemble car celle-ci est une forme approximée de ce que serait la variance d’ensemble si un tirage sans remise avec probabilité d’inclusion proportionnelle à une variable de taille était réalisé.

22. Dans la somme correspondante,  $A$  désigne l’ensemble des agglos du territoire. En d’autres termes,  

$$A = \bigcup_{\{1, \dots, H\}} h.$$

une formule découlant de l'approximation proposée par Ardilly (2006) – page 156 :

$$V1_{h,v} = \sum_{a \in h} \pi_{a,h} (1 - \pi_{a,h}) \left( \frac{w_{a,v} I_{a,v}}{\pi_{a,h}} - \frac{1}{N_h} \sum_{j \in h} \frac{w_{j,v} I_{j,v}}{\pi_{j,h}} \right)^2$$

où  $N_h$  est le nombre total d'agglos dans la strate  $h$ . On suppose donc que  $\pi_{a,h} = w_a \times n_h / W_h$  où  $n_h$  est le nombre (inconnu) d'agglos à sélectionner dans la strate  $h$ ,  $w_a$  est défini conformément à la relation (1) et  $W_h = \sum_{a \in h} w_a$ . Puis, comme on ne connaît pas le poids en dépense de consommation sur une maille élémentaire  $(a, v)$  mais que les poids marginaux sont connus, nous posons que<sup>23</sup>  $w_{a,v} = w_a \times w_v$ . Il en découle que :

$$V1_{h,v}(n_h) = w_v^2 W_h \sum_{a \in h} w_a \left( \frac{1}{n_h} - \frac{w_a}{W_h} \right) \left( I_{a,v} - \frac{1}{N_h} \sum_{j \in h} I_{j,v} \right)^2 \quad (17)$$

Formellement, il convient de retenir le vecteur  $\mathbf{n} = (n_1, \dots, n_H)$  comme solution du problème de minimisation de la variance sous contrainte que le nombre total d'agglos soit égal à  $\mathcal{N}$  :

$$\mathbf{n} = \underset{\nu}{\operatorname{argmin}} \left\{ \sum_{h=1}^H \sum_v V1_{h,v}(\nu_h) \mid \sum_{h=1}^H \nu_h = \mathcal{N} \right\}$$

Compte tenu du très faible nombre d'agglos par strate figurant dans l'échantillon actuel de l'IPC, il n'est pas possible d'estimer convenablement une variance  $\left( I_{a,v} - \frac{1}{N_h} \sum_{j \in h} I_{j,v} \right)^2$ . Par conséquent, nous choisissons de considérer que la variance précédente est identique pour toutes les strates et est donc égale à un paramètre  $\sigma_v^2$  ne dépendant pas de  $h$ . Finalement, le problème d'optimisation revient à<sup>24</sup> :

$$\mathbf{n} = \underset{\nu}{\operatorname{argmin}} \left\{ \sum_{h=1}^H \frac{W_h^2}{\nu_h} \mid \sum_{h=1}^H \nu_h = \mathcal{N} \right\} \quad (18)$$

Finalement<sup>25</sup> (voir annexe B),

$$n_h = \mathcal{N} \times W_h \quad (19)$$

23. Avec ces notations  $\sum_a w_a = 1$ ,  $\sum_v w_v = 1$  et  $\sum_{a,v} w_{a,v} = 1$ .

24. Les multiplicateurs positifs et autres constantes qui s'ajoutent à la fonction objectif sont simplifiés dans cette expression.

25. Par rapport aux travaux antérieurs, ce développement appelle la remarque suivante. Le nombre d'agglos enquêtées par strate est ici le fruit d'une minimisation de variance sous la contrainte que la somme des coûts marginaux de collecte d'une agglo soient les mêmes qu'actuellement. Cette approche est différente de celle proposée par Ardilly & Guglielmetti (1992) et reprise par Faivre (2012) qui minimisent la somme des coûts (avec une hypothèse de coûts marginaux par strate unitaire) sous la contrainte que la variance est la même qu'auparavant. La forme proposée ici est plus simple à mettre en œuvre et théoriquement duale, sous la réserve que le coût et la variance de référence des relations de contrainte associées aux deux problèmes soient déduits l'un de l'autre, ce qui n'est vraisemblablement pas le cas. Pour le calcul d'optimisation, voir annexe B avec  $a_i \rightarrow W_h$ ,  $\alpha_i \rightarrow 1$  et  $K \rightarrow \mathcal{N}$  où  $\sum_{h=1}^H W_h = 1$ .

## 6.2 La sélection dans chaque strate conditionnellement à $n_h$

Comme on l'a vu, il est judicieux d'adopter un plan de sondage tel que les probabilités d'inclusion s'approchent d'une proportionnalité à une variable de taille d'agglomération. On peut faire l'hypothèse que dans l'échantillon actuel, c'est bien le cas. En effet, cette propriété est mentionnée par Ardilly & Guglielmetti (1992) – page 75 : *“La taille d'échantillon par groupe<sup>26</sup> est à peu près proportionnelle à la taille moyenne des agglomérations du groupe auquel elles appartiennent.”*

La sélection des agglomérations s'opère conditionnellement à leur appartenance à l'échantillon d'agglomération de l'IPC-base 1998. Nous notons  $\mathcal{I}$  l'échantillon IPC-base 1998 et  $\mathcal{X}$  l'échantillon nouvelle base. On raisonne ci-après sur des ensembles d'agglomérations sur lesquels un sondage aléatoire est réalisé.

On cherche à faire en sorte que la probabilité d'inclusion d'une agglomération  $a$  dans l'échantillon  $\mathcal{X}$  soit égale à une valeur  $\pi_a^{\mathcal{X}}$  connue, tout en faisant dépendre la probabilité de sélection dans l'échantillon  $\mathcal{X}$  de l'appartenance de l'agglomération dans l'échantillon  $\mathcal{I}$ . En effet, l'idée est de faire en sorte, autant que possible, qu'une agglomération présente dans l'échantillon  $\mathcal{I}$  ait un maximum de chance de figurer dans  $\mathcal{X}$  et, en particulier, davantage qu'une agglomération qui ne figure pas dans  $\mathcal{I}$ .

On a classiquement :

$$\Pr(a \in \mathcal{X}) = \Pr(a \in \mathcal{X} | a \in \mathcal{I}) \times \Pr(a \in \mathcal{I}) + \Pr(a \in \mathcal{X} | a \notin \mathcal{I}) \times \Pr(a \notin \mathcal{I})$$

On fait l'hypothèse ici qu'on connaît  $\Pr(a \in \mathcal{I}) = \pi_a^{\mathcal{I}}$  et  $\Pr(a \notin \mathcal{I}) = 1 - \pi_a^{\mathcal{I}}$ . Moyennant quoi,

$$\pi_a^{\mathcal{X}} = \pi_a^{\mathcal{I}} \times \Pr(a \in \mathcal{X} | a \in \mathcal{I}) + (1 - \pi_a^{\mathcal{I}}) \times \Pr(a \in \mathcal{X} | a \notin \mathcal{I}) \quad (20)$$

En pratique, on souhaite évoluer le moins possible par rapport à l'échantillon d'origine. En d'autres termes, on cherche à maximiser la probabilité d'inclusion dans le nouvel échantillon d'une agglomération qui figurerait dans l'ancien (i.e. la valeur de  $\Pr(a \in \mathcal{X} | a \in \mathcal{I})$ ) tandis que la relation (20) serait vérifiée.

On montre (cf. annexe D) que deux situations peuvent survenir :

- soit  $\pi_a^{\mathcal{X}} \leq \pi_a^{\mathcal{I}}$  et dans ce cas il convient de retenir les probabilités conditionnelles d'inclusion suivantes :

$$\begin{cases} \Pr(a \in \mathcal{X} | a \in \mathcal{I}) = \pi_a^{\mathcal{X}} / \pi_a^{\mathcal{I}} \\ \Pr(a \in \mathcal{X} | a \notin \mathcal{I}) = 0 \end{cases}$$

---

26. NDR: comprendre strate.

- soit  $\pi_a^{\mathcal{X}} > \pi_a^{\mathcal{J}}$  et dans ce cas il convient de retenir les probabilités conditionnelles d'inclusion suivantes :

$$\begin{cases} \Pr(a \in \mathcal{X} | a \in \mathcal{J}) = 1 \\ \Pr(a \in \mathcal{X} | a \notin \mathcal{J}) = (\pi_a^{\mathcal{X}} - \pi_a^{\mathcal{J}})/(1 - \pi_a^{\mathcal{J}}) \end{cases}$$

En faisant l'hypothèse, pour un groupe<sup>27</sup>  $G$  d'agglos sur lequel on souhaite sélectionner un échantillon aléatoire d'agglos, que la probabilité d'inclusion est proportionnelle à une variable  $T_a$  de taille d'agflo, alors

$$\pi_a^G = n_G \times \frac{T_a}{\sum_{a \in G} T_a}$$

où  $n_G$  est le nombre d'agglos sélectionnées dans le groupe  $G$  (ce nombre est exogène à ce stade de la procédure). On peut construire de la sorte les  $\pi_a^{\mathcal{J}}$  et les  $\pi_a^{\mathcal{X}}$ . Puis, sur  $G$  on dispose d'une liste  $G^{\mathcal{J}}$  d'agglos appartenant à l'échantillon  $\mathcal{J}$ . Conformément à ce qui précède, sur  $G$ , on applique l'algorithme de sélection suivant :

1. Si  $a \in G^{\mathcal{J}}$ , deux cas peuvent se produire :
  - (a) si  $\pi_a^{\mathcal{X}} \leq \pi_a^{\mathcal{J}}$ , alors  $a$  est retenue dans  $\mathcal{X}$  avec une probabilité  $\pi_a^{\mathcal{X}}/\pi_a^{\mathcal{J}}$  ;
  - (b) si  $\pi_a^{\mathcal{X}} > \pi_a^{\mathcal{J}}$ , alors  $a$  est retenue dans  $\mathcal{X}$ .
2. Si  $a \notin G^{\mathcal{J}}$ , deux cas peuvent se produire :
  - (a) si  $\pi_a^{\mathcal{X}} \leq \pi_a^{\mathcal{J}}$ , alors  $a$  n'est pas retenue dans  $\mathcal{X}$  ;
  - (b) si  $\pi_a^{\mathcal{X}} > \pi_a^{\mathcal{J}}$ , alors  $a$  est retenue dans  $\mathcal{X}$  avec une probabilité  $(\pi_a^{\mathcal{X}} - \pi_a^{\mathcal{J}})/(1 - \pi_a^{\mathcal{J}})$ .

Considérons un groupe  $G$  dans lequel on doit sélectionner  $n_G$  agglos. D'un point de vue pratique, dans le cas 1b précédent, les agglos sont retenues. On suppose qu'il y en a  $n_G^{(1)}$  dans ce cas. Puis, pour les autres agglos, on procède par tirage systématique proportionnel à la variable de probabilité d'inclusion conditionnelle indiquée aux 1a et 2b précédents en sélectionnant  $n_G - n_G^{(1)}$  agglos. Cette sélection est implémentée sous SAS comme indiqué par Ardilly (2006) – page 157.

## 7 Application numérique

### 7.1 Détermination du nombre $n_h$

#### 7.1.1 Choix des strates

Dans un premier temps, on fixe la nature des strates retenues. Pour l'échantillon IPC, on procède comme Ardilly & Guglielmetti (1992) en réunissant les agglos par taille d'unité

---

<sup>27</sup>. Typiquement, une strate.

urbaine<sup>28</sup> et en divisant le territoire métropolitain en ZISP<sup>29</sup> (zone d'intervention d'un site prix).

### 7.1.2 Choix des poids $w_{a,v}$

Comme vu au paragraphe 6, le poids  $w_{a,v}$  se décompose en un produit  $w_a \times w_v$ . Le poids des variétés est celui issu des chiffres de la Comptabilité nationale étant données la liste des variétés à la fin de l'année 2013. S'agissant des poids d'agglos  $w_a$ , on utilise une proxy de la dépense de consommation des ménages au lieu d'achat. Cette variable, construite par Jaluzot (2014), est fondée sur<sup>30</sup> les éléments suivants :

1. La distribution de la dépense de consommation alimentaire des ménages selon la taille d'agglomération est connue par l'enquête Budget des familles 2011. La variable de taille d'unité urbaine utilisée ici est celle des données du recensement codée en 9 catégories<sup>31</sup>.
2. Les poids précédents sont ensuite ventilés, au sein des classes de taille d'unités urbaines, sur l'ensemble des unités urbaines de la classe, proportionnellement à la population municipale 2010.

### 7.1.3 Le calcul de $n_h$

Avec les éléments précédents, on applique la formule (19). Pour  $\mathcal{N}$ , on prend le nombre total d'agglos devant figurer dans l'échantillon (fixé de manière exogène à 100 pour la métropole). Il arrive que pour une strate donnée  $h$ , la formule aboutisse à un nombre d'agglos sélectionnées supérieur à celui qui figure dans la strate. Dans ce cas, il convient de saturer le nombre d'agglos sélectionnées pour la strate au nombre disponible :  $n_h = N_h$ . Puis de relancer l'algorithme de calcul en réduisant l'échantillon à sélectionner à un échantillon de taille  $\mathcal{N} - N_h$ . Les poids  $W_h$  doivent être renormalisés sur le nouvel univers du sondage (l'ensemble d'origine moins la strate saturée  $h$ ) de façon à être de somme unitaire sur le nouvel univers. On procède de la même façon itérativement jusqu'à ce que tous les  $n_h$  soient déterminés.

---

28. On distingue 5 classes : A=UU de Paris ; B2=UU de plus de 200 000 habitants ; B1=UU de 100 000 à 200 000 habitants ; C=UU de 20 000 à 100 000 habitants ; D=UU de 2 000 à 20 000 habitants. Les communes rurales sont dans une classe H complémentaire.

29. Une zone d'intervention des sites prix (ZISP) est une partition de l'espace métropolitain en 7 zones. C'est une somme de régions administratives, à l'exception du département de l'Yonne rattaché à la ZISP du site-prix de Saint-Quentin-en-Yvelines et non à celle de Lyon et du département des Ardennes rattaché à la ZISP du site-prix de Lille et non à celle de Nancy.

30. Voir annexe F.

31. Code "taille de l'unité urbaine" (basé sur la population municipale au recensement 2007 pour les UU 2010) : 0-Rural ; 1-Unités urbaines de 2 000 à 4 999 habitants ; 2-Unités urbaines de 5 000 à 9 999 habitants ; 3-Unités urbaines de 10 000 à 19 999 habitants ; 4-Unités urbaines de 20 000 à 49 999 habitants ; 5-Unités urbaines de 50 000 à 99 999 habitants ; 6-Unités urbaines de 100 000 à 199 999 habitants ; 7-Unités urbaines de 200 000 à 1 999 999 habitants ; 8-Agglomération de Paris.



## 7.2 La sélection des agglomérations

Pour chaque strate, nous procédons conformément à l’algorithme développé au paragraphe 6.2. Nous raisonnons, pour les agglos, sur l’enveloppe des unités urbaines au sens du RP de 2010. Nous supposons que la probabilité d’inclusion d’une telle agglo dans l’échantillon de l’IPC-base 98 est proportionnelle au nombre d’habitants compris dans cette agglo au sens de la population sans double compte de 1990. Le coefficient de proportionnalité est tel que la somme des probabilités d’inclusion sur une strate soit égale au nombre d’unités urbaines couvertes par la collecte IPC actuelle au seuil de 100 relevés minimum par mois<sup>32</sup>. Les agglos sont sélectionnées de sorte que leur probabilité d’inclusion dans l’échantillon-base 2015 soit proportionnelle à une variable de pondération, fondée sur une proxy de la dépense de consommation des ménages localisée au lieu d’achat<sup>33</sup>.

Le tableau 1 donne le nombre d’agglos par strate selon la variable de pondération utilisée, conditionnellement au fait que l’échantillon sélectionné comprend 100 unités urbaines. La dernière colonne du tableau indique le nombre d’unités urbaines sélectionnées dans le tirage systématique.

## 7.3 Détermination du nombre de relevés par variété et par agglomération

Pour le calcul du nombre de relevés  $n_v$  par variété, on utilise pour les coûts  $c_v$  les temps utilisés pour calculer les temps de travail des enquêteurs (Sillard 2012a, Sillard 2012b) et reportés dans la table 3.

À l’aide de ces éléments et des probabilités d’inclusion qui découlent du mode de sélection des agglos décrit au paragraphe 7.2, on en déduit le nombre de relevés requis par variétés. Enfin, en appliquant la formule (10), on en déduit le nombre d’observations par variété et agglo.

L’optimisation du nombre de relevés par unité urbaine et par variété conduit à réviser également les cibles de nombre de relevés. Le graphique 1 donne un tracé du nombre optimal de relevés par agglomération à coût de collecte d’ensemble inchangé, pour l’échantillon sélectionné. Le tableau 4 donne le nombre moyen de relevés par catégorie d’unité urbaine dans l’échantillon actuel et dans l’échantillon-type optimisé en base 2015. Par rapport à la pratique actuelle, le nombre de relevés<sup>34</sup> devrait augmenter sur les petites agglos, passant

---

32. Si la probabilité d’inclusion ainsi obtenue est supérieure à 1, elle est saturée à 1 et le calcul des probabilités d’inclusion des autres agglos est relancé avec un nouveau coefficient de proportionnalité tel que la somme soit unitaire sur les probabilités itérées. Les itérations se poursuivent ainsi jusqu’à ce qu’aucune probabilité d’inclusion calculée ne soit strictement supérieure à 1.

33. Cette proxy est fondée sur les poids en dépense de consommation alimentaire des ménages par type d’agglo (Jaluzot 2014), répartis par agglo, au sein des classes de type d’agglos, en proportion de la variable de population des ménages au dernier RP (population municipale 2010).

34. On rappelle qu’un jour-plein de collecte représente une petite centaine de relevés.

TABLE 1 – Nombre d’unités urbaines par strate métropolitaine

Strate		$N_h$	IPC	optimi- sation	$n_h$	sélec. auto.
Taille d’UU	ZISP					
A	11	1	1	1	1	1
	08	3	2	3	3	3
	20	6	5	6	6	6
	28	5	5	5	5	5
	36	5	5	5	5	5
	40	5	5	3,62	4	3
	64	7	5	5,25	5	3
B1	08	5	2	2,34	2	1
	20	4	2	1,79	2	1
	28	4	1	1,94	2	1
	36	2	1	1,15	1	1
	40	4	1	1,59	2	1
	64	3	0	1,44	1	0
C	08	28	5	4,18	4	0
	11	14	3	1,65	2	0
	20	29	5	4,22	4	0
	28	32	6	4,77	5	2
	36	18	3	3,03	3	1
	40	29	5	3,69	4	1
	64	38	7	5,32	5	1
D	08	285	10	5,21	5	0
	11	97	3	1,85	2	1
	20	256	4	4,08	4	2
	28	290	5	4,96	5	2
	36	255	6	4,66	5	1
	40	275	5	4,25	4	1
	64	518	15	9,01	9	0
<i>Total</i>		2218	117	100	100	43

**Note** : ZISP 08=Aquitaine, 11=IdF (y.c. Yonne), 20=NPdC (y.c. Ardennes), 28=Rhône-A. (n.c. Yonne), 36=L-R, 40=Lorraine (n.c. Ardennes), 64=Bretagne;  $N_h$  est le nombre d’unités urbaines de la strate; IPC est le nombre d’unités urbaines de la strate couvertes par la collecte actuelle de l’IPC; “optimisation” désigne le nombre d’unités urbaines qu’il conviendrait de retenir dans la strate si le poids de la strate était strictement proportionnel à la variable de dépense de consommation au lieu d’achat;  $n_h$  est le nombre d’unités urbaines (UU) retenu pour la strate en base 2015. La colonne “sélec. auto” indique, parmi les unités urbaines sélectionnées dans l’échantillon base 2015, combien sont sélectionnées avec une probabilité d’inclusion (conditionnelle à leur appartenance à l’échantillon) certaine.

TABLE 2 – Poids des regroupements d’unités urbaines utilisés pour la stratification

Strate géographique	Poids (en %)			
	démographique (1)	alimentaire (2)	(1) normalisé sans le rural	(2) normalisé sans le rural
A	16,7	18,1	21,5	19,5
B2	24,6	25,6	31,7	27,6
B1	5,3	7	6,8	7,6
C	13,6	18,6	17,5	20,1
D	17,4	23,4	22,4	25,2
Rural	22,5	7,4		
<i>Total</i>	100	100	100	100

**Note** : la colonne de poids démographique est issue de la population au sens du RP 2010 et celle de l’alimentaire est issue du traitement de l’enquête BDF en localisation au lieu d’achat (Jaluzot 2014). Les deux dernières colonnes de poids correspondent aux deux premières, respectivement, le poids du rural étant annulé. Ceci revient à faire l’hypothèse que les achats réalisés dans le rural sont effectués dans les unités urbaines, au prorata du poids de ces unités urbaines dans la consommation.

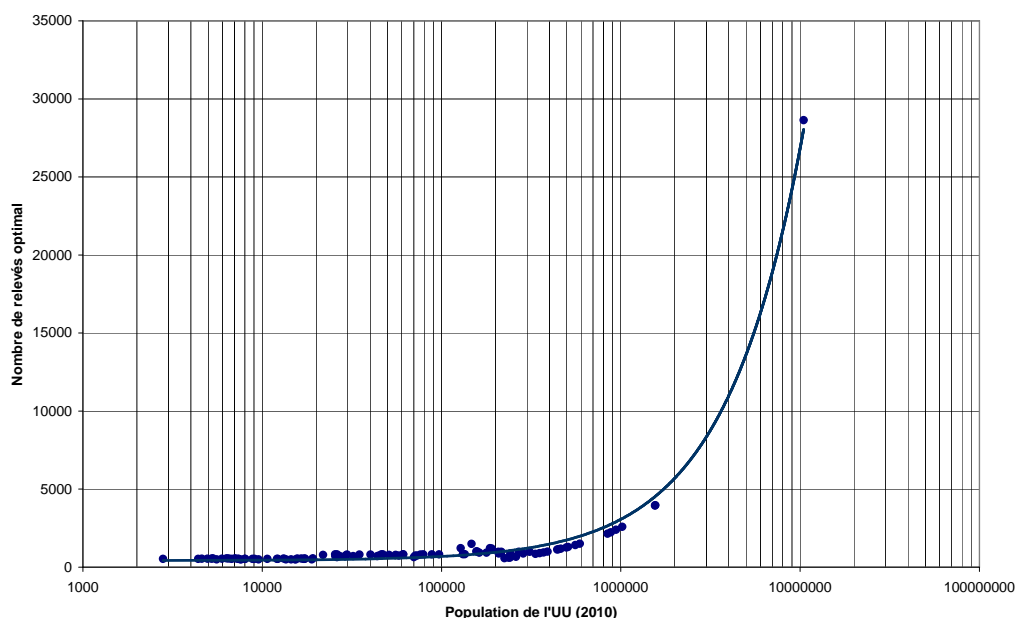
TABLE 3 – Durées moyennes d’observation par secteur de relevé utilisées pour le coût élémentaire de relevé  $c_v$ 

secteur	durée moyenne hors déplacement	normalisée y.c. dépl. d’un PV à l’autre	Nombre de point de vente par relevé
Bien durables	102s	1.30	0.27
Habillement	60s	0.86	0.19
Alimentaire	49s	0.53	0.09
P. manufacturé	78s	1.16	0.27
Services	25s	1.73	0.61

**Note** : la colonne 3 comprend une mesure  $\alpha_i$  normalisée de temps élémentaires sectoriels issus de la colonne 2. La normalisation est telle que  $\sum_{i \in \text{secteur}} \alpha_i \text{Obs}_i = 112000$ , 112000 relevés étant approximativement le nombre de relevés réalisés au mois de décembre pour l’IPC en 2012 et 2013. Cette mesure tient compte des déplacements d’un point de vente à l’autre lors d’une tournée d’enquêteur. Le temps de déplacement est fixé à 7min13s par point de vente fréquenté. Le nombre de points de vente fréquentés par relevé, selon le secteur, est indiqué en colonne 4. Ces calculs sont issus de l’analyse des temps d’observation élémentaire par relevé (voir Sillard (2012b)).

d'une moyenne de 177 relevés actuellement à une moyenne de 529 relevés pour les unités urbaines de moins de 20 000 habitants (D) ; il devrait également augmenter pour les unités urbaines de 20 000 à 100 000 habitants (C), passant d'une moyenne de 464 à 767 relevés en moyenne ; pour les unités urbaines de 100 000 à 200 000 habitants (B1) et pour les unités urbaines de plus de 200 000 habitants (hors Paris – B2), le nombre d'observations devait diminuer, comme l'indiquent les chiffres du tableau 4. A noter que pour ces unités urbaines, le nombre optimal de relevés à réaliser est très dépendant de la taille d'unité urbaine, donc l'indication donnée par l'évolution de la moyenne du nombre de relevés entre l'échantillon actuel et l'échantillon base 2015 est à relativiser en conséquence<sup>35</sup>. Enfin, sur l'unité urbaine de Paris, le nombre de relevés à réaliser devrait passer de 22 900 à 28 600.

FIGURE 1 – Lien entre nombre optimal de relevés et population de l'unité urbaine



**Note :** Tracé des couples de points (population de l'UU, nombre optimal de relevés) pour les 100 agglomérations d'un échantillon type. La courbe est l'approximation linéaire optimale pour ces points. Elle est d'équation  $Nb\_relevés = 0.0026 \times pop\_mun + 424$ . La même relation calculée sur la collecte de l'année 2013 en base 1998 est  $Nb\_relevés = 0.0022 \times pop\_mun + 456$ . *Rem. : Échelle logarithmique pour les abscisses.*

La répartition du nombre de relevés par secteur de produit est également légèrement modifiée par rapport à l'existant. Après optimisation du nombre de relevés par variétés, les principales modifications (cf. table 4) portent sur l'alimentaire, où le nombre de relevés devrait passer de 36 000 à 49 000 et sur l'habillement où le nombre de relevés devrait

<sup>35</sup>. Un calcul de dépendance du nombre optimal de relevés à réaliser en fonction de la taille d'unité urbaine est indiqué en note de la figure 1.

diminuer de 21 000 à 11 000. La masse de relevés consacrés aux autres secteurs est peu modifiée.

TABLE 4 – Nombre de relevés en fonction du type d’agglomération

Type d’agglo	Nb. relevés b. 1998	Nb. relevés b. 2015
<i>A</i>	22862	28633
<i>B2</i>	2105 [365 ; 5048]	1389 [570 ; 3966]
<i>B1</i>	1485 [1263 ; 1909]	1063 [826 ; 1491]
<i>C</i>	464 [137 ; 1145]	767 [637 ; 820]
<i>D</i>	177 [100 ; 381]	529 [489 ; 562]

**Note** : les types d’unité urbaines sont – l’unité urbaine de Paris (A) ; – les unités urbaines de plus de 200 000 habitants (B2) hors Paris ; – les unités urbaines de 100 000 à 200 000 habitants (B1) ; – les unités urbaines de 20 000 à 100 000 habitants (C) ; – les unités urbaines de moins de 20 000 habitants (D, hors zones rurales). Les intervalles correspondent aux valeurs extrêmes observées sur les agglos de la catégorie concernée. Calcul réalisé sur les agglos de la base 2015.

TABLE 5 – Nombre de relevés en fonction du type de secteur

Secteur	Base 1998 (total)	Base 1998 (Agglos conservées)	Base 2015
AL	35 568	31 801	49 380
BD	6 544	6 074	5 652
HA	21 033	20 451	11 449
MA	29 939	27 683	29 019
SE	19 683	18 119	21 375
<b>Total</b>	<b>112 767</b>	<b>104 128</b>	<b>116 875</b>

**Note** : AL=alimentaire, BD=biens durables, HA=habillement, MA=produits manufacturés, SE=services ; “Agglos conservées” désigne les agglos de la base 1998 qui sont conservées en base 2015.

S’agissant des produits frais, ceux-ci sont exclus de l’exercice d’optimisation car leurs relevés ne correspondent pas, pour l’instant, au schéma retenu pour le reste de la collecte IPC. En revanche, le changement de base est l’occasion de faire évoluer la collecte pour le rendre homologué à celle des produits alimentaires. Dans ce contexte, le nombre de relevés de produits frais est déterminé de sorte que l’effort de collecte consenti en base 1998 par l’Insee perdure en base 2015. En effet, les relevés de produits frais sont en particulier utilisés pour la publication de prix moyens par type de produits et l’optimisation du

nombre de relevés aux seules fins de publication d’indice de prix ne prendrait pas en compte cette exigence particulière. Il a donc été décidé de maintenir le dimensionnement actuel de l’échantillon. Ceci définit le nombre total de relevés de produits frais à réaliser en nouvelle base (en l’occurrence 37 000 sur la métropole). Ces relevés seront ventilés entre agglos de collecte proportionnellement au nombre de relevés réalisés au titre de l’alimentaire.

## 7.4 Le nombre de relevés par forme de vente et par point de vente

Pour chaque variété de l’IPC, la distribution marginale des ventes par forme de vente<sup>36</sup> est supposée connue de l’expert sectoriel de la division des prix à la consommation en charge du suivi de la variété concernée. Cette connaissance repose sur différentes sources de documentations, notamment les données issues des panélistes. La répartition marginale est ensuite appliquée, par agglo, au nombre de relevés fixé.

Dans le but de ne pas engendrer d’“effets de grappes” dans la mesure de l’évolution des prix, les enquêteurs sont invités à ne pas multiplier les relevés dans un point de vente donné. Ce nombre dépend de la taille du point de vente et a été fixé sur la base de critères statistiques (Division des prix à la consommation 2006). Ce principe est conservé en base 2015. Le nombre de relevés important à réaliser sur de petites agglomérations en base 2015 peut se heurter au faible nombre de points de vente sur site, compte tenu notamment de la contrainte précédente. Au cas par cas, il est donc nécessaire d’adapter le calcul statistique à la réalité du terrain. Dans un tel cas, tout se passe comme si le poids de l’agglo en dépense de consommation des ménages avait été initialement sur-évalué. Il convient donc de le revoir à la baisse, en gardant en tête que la probabilité d’inclusion de l’agglo dans l’échantillon est inchangée, c’est-à-dire qu’elle reste dictée par le poids  $w_a$  initialement utilisé pour le calcul du plan de sondage. Ainsi, à probabilité d’inclusion de l’agglo  $a$  dans l’échantillon inchangée, les résultats issus des formules (10) et (13) sont à recalculer conformément à ces formules après avoir diminué les  $w_{a,v}$  pour l’agglo  $a$  considérée (et pour des  $\pi_a$  inchangés). Au premier ordre, ceci va se traduire par le maintien du nombre total de relevés à réaliser pour la variété  $v$  concernée (variable  $n_v$  dans les expressions précédentes) et un “transfert” des relevés en moins dans l’agglo  $a$  vers l’ensemble des autres agglos, approximativement en proportion des nombres de relevés initialement calculés pour elles au titre de la variété considérée. Le processus d’adaptation des cibles de nombre de relevés décrit ci-dessus est mis en œuvre quand le terrain l’impose. En pratique, ceci ne concerne qu’un faible nombre de couples (agglos, variété).

---

36. les points de vente de l’IPC sont classés parmi 11 formes de vente (avec leur code entre parenthèses) : Hypermarché (10), Supermarché (20), Maxidiscount (25), Supérette (30), Magasin populaire (40), Grand magasin (50), Grande surface spécialisée (60), Magasin traditionnel (70), Marché (80), Services (90), autres (99).

## Références

- Ardilly, P. (2006). *Les techniques de sondage*, 2 edn, Technip.
- Ardilly, P. (2014). Échanges entre la division des prix et M. Pascal Ardilly à propos de la sélection des agglomérations IPC dans la Base 2015, *Communication personnelle, Note interne N° 295/DG75-F320/*, Insee, division des prix à la consommation.
- Ardilly, P. & Guglielmetti, F. (1992). Optimisation de l'échantillon pour le calcul de l'indice de prix à la consommation, *Insee méthodes* **29-30-31**: 71 – 123.
- Ardilly, P. & Guglielmetti, F. (1993). La précision de l'indice des prix : mesure et optimisation, *Économie et Statistique* **267**: 13 – 29.
- Balcone, T. (2014). La précision de l'indice des prix à la consommation, *Note interne N° 1460/DG75-F320/TB*, Insee, division des prix à la consommation.
- Division des prix à la consommation (2006). Guide de l'enquêteur des prix à la consommation, *Manuel de collecte*, INSEE.
- Faivre, S. (2012). Calcul de l'échantillon d'agglomérations optimal pour l'indice des prix suite au passage aux données de caisses, *7<sup>ème</sup> colloque francophone sur les sondages, ENSAI - Rennes, 5 – 7 novembre 2012*.
- Jaluzot, L. (2014). Étude de la géographie de la consommation des ménages, *Note interne N° 172/DG75-F320/LJ*, Insee, division des prix à la consommation.
- Oehlert, G. W. (1992). A Note on the Delta Method, *The American Statistician* **46**(1): 27 – 29.
- Petit, G. (2014). Précision de l'indice des prix à la consommation, *Note interne N° 815/DG75-F320/*, Insee, division des prix à la consommation.
- Sillard, P. (2012a). Proposition de calcul de quotité des enquêteurs-prix métropolitains, *Note interne N° 887/DG75-F320/PS*, Insee, division des prix à la consommation.
- Sillard, P. (2012b). Temps d'observations modélisés pour les enquêteurs prix, *Note interne N° 1523/DG75-F320/PS*, Insee, division des prix à la consommation.
- Tillé, Y. (2001). *Théorie des sondages*, Dunod.

# A Les formules d'agrégation actuellement utilisées dans l'IPC

## A.1 L'agrégation des indices de variétés et au-delà

L'agrégation des indices de variétés  $I_v$  s'opère selon une formule de Laspeyres :

$$I = \sum_v POND_v \times I_v$$

où  $POND_v$  est un poids en valeur. Ce poids est obtenu par la composition d'un poids en valeur du [poste,territoire], issu de la comptabilité nationale et noté  $POND_P$ , et d'un poids de la variété dans le poste  $POSTE$ . Par hypothèse,

$$\sum_{v \in P} POIDS_v = 100$$

et

$$POND_v = \frac{POIDS_v}{100} \times POND_P$$

## A.2 Le calcul de l'indice de variété $I_v$

Par hypothèses, étant donnée une série d'indice  $I_{v,a}$  de varaggio, l'indice de variété se calcule comme suit :

$$I_v = \sum_a POND_{v,a} \times I_{v,a}$$

Le détail des calculs de la variable  $POND_{v,a}$  est donné aux paragraphes suivants.

### A.2.1 Le calcul de $POND_{v,TA}$

On définit d'abord le poids de la variété dans le type d'agglomération ( $TA \in \{A, B, C, D\}$ ) :

$$POND_{v,TA} = \frac{POIDS_{v,TA}}{100} \cdot \frac{POIDS_v}{100} \cdot POND_{P,terri}$$

- les variables  $POIDS_{v,TA}$  et  $POIDS_v$  sont fixés par le sectoriel de la division des prix en charge de suivi de la variété  $v$ , lors de la définition de la variété (opérations de changement d'année). Ainsi défini,  $\sum_{TA \in \{A, B, C, D\}} POIDS_{v,TA} = 100$  caractérise la répartition de la dépense de la consommation de la variété  $v$  sur le territoire considéré selon les 4 types d'agglomérations  $\{A, B, C, D\}$ .

- $POND_{P,terri}$  est calculé à partir des éléments de la Comptabilité nationale.

Avec cette formule, on répartit le poids du [poste×terri] sur chaque couple  $(v, TA)$ .  $POND_{v,TA}$  est donc un poids en dépense, de sorte que :

$$\sum_{\substack{v \in POSTE \\ TA \in \{A, B, C, D\}}} POND_{v,TA} = POND_{P,terri}$$



### A.2.2 Calcul de $POND_{v,a}$ pour le type d'agglomération=A

$$POND_{v,a=PARIS} = POND_{v,TA=A}$$

### A.2.3 Calcul de $POND_{v,a}$ pour le type d'agglomération=B

On répartit  $POND_{v,TA}$  sur l'ensemble des agglomérations  $a \in \{TA = B\}$  selon un coefficient de répartition :

$$POND_{v,a} = R_{v,a} \times POND_{v,TA=B}$$

On définit d'abord le nombre de relevés réalisés dans l'agglomération  $a$  :

$$NBENQ_a = \sum_v NBENQ_{v,a}$$

et

$$R_{v,a} = \frac{POIDS_a \times NBENQ_{v,a}/NBENQ_a}{\sum_{a \in \{TA=B\}} POIDS_a \times NBENQ_{v,a}/NBENQ_a}$$

Le fondement de cette manière de procéder pour le calcul du poids n'apparaît pas clairement. Il conviendra de la ré-examiner pour la base 2015.

## A.3 Calcul de $POND_{v,a}$ pour le type d'agglomération={C,D}

La pondération  $POND_{v,TA}$  sur l'ensemble des agglomérations de type  $TA \in \{C, D\}$  est répartie entre agglomérations selon une formule

$$POND_{v,a} = R_{v,a} \times POND_{v,TA}$$

où

$$R_{v,a} = \frac{NBENQ_{v,a}}{NBENQ_{v,TA}}$$

C'est le poids relatif de l'agglomération  $a$  dans sa taille d'agglomération, en nombre de relevés effectués au titre de la variété considérée.

## B Optimisation du nombre de relevés

On cherche à résoudre le programme d'optimisation suivant :

$$\boldsymbol{\nu}^* = \underset{\boldsymbol{\nu}}{\operatorname{argmin}} \left\{ \sum_i \frac{a_i^2}{\nu_i} \mid \sum_i \alpha_i^2 \nu_i = K \right\}$$

où  $\boldsymbol{\nu}$  est un vecteur d'inconnues de dimension  $M$  ( $\boldsymbol{\nu} = (\nu_1, \dots, \nu_M)$ ) et les  $(a_i)_{1 \leq i \leq M}$ ,  $(\alpha_i)_{1 \leq i \leq M}$  et  $K$  sont des paramètres du programme. Les  $a_i$  sont positifs ou nuls; les  $\alpha_i$  et  $K$  sont strictement positifs.

TABLE 6 – Liste des variables et tables des applications IPC correspondantes

Table IPCNAT	variable	colonne
AGGLO	$POIDS_a$	POIDS
ANVA	$POND_{v,a}$	POND
”	$NBENQU_{v,a}$	NBENQ
VARTA	$POIDS_{v,TA}$	POIDS
”	$NBENQ_{v,TA}$	NBENQ
VARIETE	$POIDS_v$	POIDS
”	$POND_v$	PONDTOT

Le lagrangien du programme s’écrit ( $\lambda$  est le multiplicateur de Lagrange) :

$$\mathcal{L}(\boldsymbol{\nu}, \lambda) = \sum_i \frac{a_i^2}{\nu_i} + \lambda(\sum_i \alpha_i^2 \nu_i - K)$$

Puis,

$$\frac{\partial \mathcal{L}}{\partial \nu_i} = -\frac{a_i^2}{\nu_i^2} + \lambda \alpha_i^2$$

Les conditions de premier ordre  $\frac{\partial \mathcal{L}}{\partial \nu_i} = 0$  conduisent à :

$$\lambda = \frac{a_i^2}{\alpha_i^2 \nu_i^2}$$

Il en découle que

$$\nu_i = \frac{a_i}{\alpha_i \sqrt{\lambda}}$$

Or  $\sum_i \alpha_i^2 \nu_i = K$  donc

$$\sqrt{\lambda} = \frac{1}{K} \sum_i \alpha_i a_i$$

Finalement, pour tout  $i \in \{1, \dots, M\}$ ,

$$\nu_i^* = \frac{K}{\alpha_i^2} \times \frac{\alpha_i a_i}{\sum_i \alpha_i a_i}$$

# C Variance des micro-indices et biais de formule à distance finie

## C.1 Variétés homogènes

On considère un ensemble de bien de taille  $N$ . Chaque bien est indicé par  $i \in \{1, \dots, N\}$ . Le prix des biens est  $p_i^t$  à l'instant  $t$  et  $p_i^0$  à l'instant 0. On cherche à estimer l'indice

$$I = \frac{\frac{1}{N} \sum_{i=1}^N p_i^t}{\frac{1}{N} \sum_{i=1}^N p_i^0}$$

en sélectionnant par sondage aléatoire simple avec remise  $n$  produits. On note  $\mathcal{S}$  cet échantillon et on suppose que le taux de sondage  $n/N$  est très petit devant l'unité, de sorte qu'on peut le négliger. On estime  $I$  par  $\hat{I}$  défini par :

$$\hat{I} = \frac{\frac{1}{n} \sum_{i \in \mathcal{S}} p_i^t}{\frac{1}{n} \sum_{i \in \mathcal{S}} p_i^0}$$

La variance de  $\hat{I}$  provient de la variance des deux variables  $\hat{P}^t = \frac{1}{n} \sum_{i \in \mathcal{S}} p_i^t$  et  $\hat{P}^0 = \frac{1}{n} \sum_{i \in \mathcal{S}} p_i^0$  en tant qu'estimateurs de  $\bar{P}^t = \frac{1}{N} \sum_{i \in \{1, \dots, N\}} p_i^t$  et  $\bar{P}^0 = \frac{1}{N} \sum_{i \in \{1, \dots, N\}} p_i^0$ . On sait calculer  $\text{var}(\hat{P}^t)$  et trouver un estimateur  $\widehat{\text{var}}(\hat{P}^t)$  de cette variance à l'aide des prix observés sur l'échantillon  $\mathcal{S}$ . De même s'agissant de  $\hat{P}^0$ . Il reste donc à calculer  $\text{var}(\hat{I})$  étant données les variance de  $\hat{P}^t$  et de  $\hat{P}^0$ .

Classiquement, on applique le principe de la  $\Delta$ -méthode : si  $\bar{X}_n$ ,  $\bar{Y}_n$  et  $\bar{Z}_n$  sont des variables aléatoires telles que  $\bar{Z}_n = F(\bar{X}_n, \bar{Y}_n)$  où  $F$  est une fonction deux fois dérivable et  $(\bar{X}_n, \bar{Y}_n)$  est un vecteur de variables aléatoire – moyennes arithmétiques de variables iid (de moyenne  $(\mathbb{E}X, \mathbb{E}Y)$ ) sur des échantillons de taille  $n$ , alors, en développant  $F$  au voisinage de  $(\mathbb{E}X, \mathbb{E}Y)$  on a<sup>37</sup> :

$$\begin{aligned} \bar{Z}_n &= F(\mathbb{E}X, \mathbb{E}Y) + \partial_1 F(\mathbb{E}X, \mathbb{E}Y) \cdot (\bar{X}_n - \mathbb{E}X) + \\ &\quad \partial_2 F(\mathbb{E}X, \mathbb{E}Y) \cdot (\bar{Y}_n - \mathbb{E}Y) + \bar{\varepsilon}_n \end{aligned} \quad (21)$$

On montre<sup>38</sup> que  $\text{var}(\bar{\varepsilon}_n) = o(1/n)$ . Il en découle que<sup>39</sup> :

$$\begin{aligned} \text{var}(\bar{Z}_n) &= [\partial_1 F(\mathbb{E}X, \mathbb{E}Y)]^2 \cdot \text{var}(\bar{X}_n) + [\partial_2 F(\mathbb{E}X, \mathbb{E}Y)]^2 \cdot \text{var}(\bar{Y}_n) \\ &\quad + 2 [\partial_1 F(\mathbb{E}X, \mathbb{E}Y)] \cdot [\partial_2 F(\mathbb{E}X, \mathbb{E}Y)] \text{cov}(\bar{X}_n, \bar{Y}_n) + o(1/n) \end{aligned} \quad (22)$$

37.  $\partial_i F$  désigne la dérivée partielle de  $F$  par rapport à sa  $i^{\text{ème}}$  composante.

38. Voir par exemple (Oehlert 1992).

39. Compte tenu des hypothèses,  $\text{var}(\bar{X}_n) = \sigma_X^2/n$  où  $\sigma_X^2$  est la variance de  $X$  ; une propriété similaire est vérifiée par  $Y$ .

Dans le cas des variétés homogènes,  $F(X, Y) = \frac{X}{Y}$  donc  $\partial_1 F(X, Y) = \frac{1}{Y}$  et  $\partial_2 F(X, Y) = -\frac{X}{Y^2}$ . Avec les notations précédentes, nous avons :

$$\text{var}(\hat{I}) = \left(\frac{1}{\bar{P}^0}\right)^2 \text{var}(\hat{P}^t) + \left(\frac{\bar{P}^t}{(\bar{P}^0)^2}\right)^2 \text{var}(\hat{P}^0) - 2\frac{1}{\bar{P}^0} \times \frac{\bar{P}^t}{(\bar{P}^0)^2} \text{cov}(\hat{P}^t, \hat{P}^0)$$

On remarque que :

$$\text{var}(\hat{I}) = \left(\frac{1}{\bar{P}^0}\right)^2 \text{var}\left(\hat{P}^t - \frac{\bar{P}^t}{\bar{P}^0} \hat{P}^0\right)$$

Moyennant quoi, il ne reste plus qu'à construire un estimateur de l'expression précédente en remarquant que :

$$\hat{P}^t - \frac{\bar{P}^t}{\bar{P}^0} \hat{P}^0 = \frac{1}{n} \sum_{i \in \mathcal{S}} \left(p_i^t - \frac{\bar{P}^t}{\bar{P}^0} p_i^0\right)$$

Un estimateur de la variance de cette moyenne empirique est classiquement :

$$\widehat{\text{var}}\left(\hat{P}^t - \frac{\bar{P}^t}{\bar{P}^0} \hat{P}^0\right) = \frac{1}{n} \frac{1}{n-1} \sum_{i \in \mathcal{S}} \left[p_i^t - \frac{\hat{P}^t}{\hat{P}^0} p_i^0 - \left(\hat{P}^t - \frac{\hat{P}^t}{\hat{P}^0} \hat{P}^0\right)\right]^2$$

Soit, après simplifications (voir Eq. 15) :

$$\widehat{\text{var}}(\hat{I}) = \left(\frac{1}{\hat{P}^0}\right)^2 \frac{1}{n(n-1)} \sum_{i \in \mathcal{S}} \left(p_i^t - \hat{I} p_i^0\right)^2$$

Il est possible d'évaluer le biais de formule de l'indice, c'est-à-dire l'écart entre  $\mathbb{E}(\hat{I})$  et  $I$ . Ce biais dépend directement de la taille  $n$  de l'échantillon  $\mathcal{S}$ . En effet, si on développe  $\bar{Z}_n$  à l'ordre 2, nous avons :

$$\begin{aligned} \bar{Z}_n &= F(\mathbb{E}X, \mathbb{E}Y) + \partial_1 F(\mathbb{E}X, \mathbb{E}Y) \cdot (\bar{X}_n - \mathbb{E}X) + \partial_2 F(\mathbb{E}X, \mathbb{E}Y) \cdot (\bar{Y}_n - \mathbb{E}Y) + \\ &+ \frac{1}{2} [\partial_{11} F(\mathbb{E}X, \mathbb{E}Y) \cdot (\bar{X}_n - \mathbb{E}X)^2 + \partial_{22} F(\mathbb{E}X, \mathbb{E}Y) \cdot (\bar{Y}_n - \mathbb{E}Y)^2 + \\ &+ 2\partial_{12} F(\mathbb{E}X, \mathbb{E}Y) \cdot (\bar{X}_n - \mathbb{E}X)(\bar{Y}_n - \mathbb{E}Y)] + \bar{\xi}_n \end{aligned} \quad (23)$$

On montre<sup>38</sup> que  $\mathbb{E}(\bar{\xi}_n) = o(1/n)$ . Moyennant quoi, pour les variétés homogènes<sup>40</sup>, le biais relatif s'écrit<sup>41</sup> :

$$\frac{\mathbb{E}(\hat{I}) - I}{I} = \left(\frac{\sigma}{\bar{P}^0}\right)^2 \times \frac{I - \rho}{I} \times \frac{1}{n} + o(1/n) \quad (24)$$

où  $\sigma^2$  est la variance caractérisant la dispersion des  $p_i^0$  (supposés iid) autour de  $\bar{P}^0$  et  $\rho$  est la corrélation des variables de prix  $(p_i^0, p_i^t)_{i \in \{1, \dots, n\}}$ .

40.  $\partial_{11} F(X, Y) = 0$ ,  $\partial_{12} F(X, Y) = -1/Y^2$  et  $\partial_{22} F(X, Y) = 2X/Y^3$ .

41. On suppose que  $\text{cov}(\bar{X}_n, \bar{Y}_n) = 0$ .

## C.2 Variétés hétérogènes

Avec les mêmes notations qu'au paragraphe C.1, on cherche cette fois à estimer l'indice

$$I = \frac{\left( \prod_{i=1}^N p_i^t \right)^{1/N}}{\left( \prod_{i=1}^N p_i^0 \right)^{1/N}}$$

On estime  $I$  par  $\hat{I}$  défini par :

$$\hat{I} = \frac{\left( \prod_{i \in \mathcal{S}} p_i^t \right)^{1/n}}{\left( \prod_{i \in \mathcal{S}} p_i^0 \right)^{1/n}} = \exp \left\{ \frac{1}{n} \sum_{i \in \mathcal{S}} \ln p_i^t - \frac{1}{n} \sum_{i \in \mathcal{S}} \ln p_i^0 \right\}$$

Pareillement,  $\text{var}(\hat{I})$  provient de la variance des composantes aléatoires  $X$  et  $Y$  de la fonction  $F(X, Y) = \exp(X - Y)$ , où  $X = \frac{1}{n} \sum_{i \in \mathcal{S}} \ln p_i^t$  et  $Y = \frac{1}{n} \sum_{i \in \mathcal{S}} \ln p_i^0$  sont des estimateurs de  $\overline{LP}^t = \frac{1}{N} \sum_{i=1}^N \ln p_i^t$  et  $\overline{LP}^0 = \frac{1}{N} \sum_{i=1}^N \ln p_i^0$ . A l'aide d'une  $\Delta$ -méthode, en utilisant le développement (22) nous avons dans le cas présent ( $\partial_1 F(X, Y) = \exp(X - Y)$  et  $\partial_2 F(X, Y) = -\exp(X - Y)$ ) :

$$\begin{aligned} \text{var}(\hat{I}) &= I^2 \text{var}(\widehat{\overline{LP}}^t) + I^2 \text{var}(\widehat{\overline{LP}}^0) - 2I^2 \text{cov}(\widehat{\overline{LP}}^t, \widehat{\overline{LP}}^0) \\ &= I^2 \text{var}(\widehat{\overline{LP}}^t - \widehat{\overline{LP}}^0) \end{aligned}$$

Or

$$\widehat{\overline{LP}}^t - \widehat{\overline{LP}}^0 = \frac{1}{n} \sum_{i \in \mathcal{S}} (\ln p_i^t - \ln p_i^0)$$

donc

$$\widehat{\text{var}} \left( \widehat{\overline{LP}}^t - \widehat{\overline{LP}}^0 \right) = \frac{1}{n} \cdot \frac{1}{n-1} \sum_{i \in \mathcal{S}} \left[ \ln \left( \frac{p_i^t}{p_i^0} \right) - (\widehat{\overline{LP}}^t - \widehat{\overline{LP}}^0) \right]^2$$

Or  $\widehat{\overline{LP}}^t - \widehat{\overline{LP}}^0 = \ln \hat{I}$ . Finalement (voir aussi Eq. (16)),

$$\widehat{\text{var}}(\hat{I}) = \frac{\hat{I}^2}{n(n-1)} \sum_{i \in \mathcal{S}} \left[ \ln \left( \frac{p_i^t}{p_i^0} \right) - \ln \hat{I} \right]^2$$

De même que pour les variétés homogènes, il est possible d'évaluer le biais de formule d'indice pour les variétés hétérogènes, c'est-à-dire l'écart entre  $\mathbb{E}(\hat{I})$  et  $I$ . Ce biais dépend

directement de la taille  $n$  de l'échantillon  $\mathcal{S}$ . On repart du développement (23). Pour les variétés hétérogènes<sup>42</sup>, le biais relatif s'écrit<sup>43</sup> :

$$\frac{\mathbb{E}(\hat{I}) - I}{I} = (1 - \tilde{\rho})\tilde{\sigma}^2 \frac{1}{n} + o(1/n) \quad (25)$$

où  $\tilde{\sigma}^2$  est la variance caractérisant la dispersion des  $\ln p_i$  (supposés iid) autour de  $\bar{P}$  et  $\tilde{\rho}$  est la corrélation des variables de prix  $(\ln p_i^0, \ln p_i^t)_{i \in \{1, \dots, n\}}$ .

### C.3 La pratique du biais de formule

Le terme de biais indiqué aux formules (24) et (25) est généralement assez petit car  $\rho$  est assez proche de 1 (de même que  $I$  pour les variétés homogènes) et le nombre d'observations par agrégat est relativement grand. Pour les variétés de l'IPC, la valeur médiane du coefficient  $(\sigma/\bar{P}^0)^2 \times (I - \rho)/I$  pour les variétés homogènes et du coefficient  $(1 - \tilde{\rho})\tilde{\sigma}^2$  pour les variétés hétérogènes est de 0,06 point d'indice. La pratique actuelle de l'IPC est d'imposer un minimum de 4 observations par agrégat élémentaire var-agglo, à la fois pour limiter l'ampleur du biais de formule mais aussi pour assurer un niveau d'information statistique minimale par agrégat élémentaire. Avec cette disposition, une estimation réalisée sur la collecte actuelle montre qu'un majorant du biais d'indice est de 0,05 point d'indice, s'agissant d'un indice fondé sur la seule collecte terrain. Le terme de biais est donc bien en-deçà de la précision actuelle de la collecte terrain (estimée à 0,1 point). Avec l'optimisation du nombre de relevés, le terme de biais est de l'ordre de l'imprécision d'indice et il convient dès lors d'assurer une vigilance sur chaque micro-agrégat afin d'en limiter le biais, soit en augmentant, le cas échéant, globalement le nombre de relevés d'une variété, soit en répartissant les relevés sur les agglos de manière différente de sorte que les agglos sur lesquelles la variété est observée fasse l'objet d'un nombre de relevés adéquat. Cet arbitrage est à effectuer en prenant aussi en compte les questions de précision particulières à chacune des variétés et générale à l'ensemble de l'échantillon et à son coût de collecte.

## D Calcul de la probabilité d'inclusion optimale pour conserver un maximum d'agglo dans l'échantillon nouveau

On repart de la formule (20) :

$$\pi_a^{\mathcal{X}} = \pi_a^{\mathcal{S}} \times \Pr(a \in \mathcal{X} | a \in \mathcal{S}) + (1 - \pi_a^{\mathcal{S}}) \times \Pr(a \in \mathcal{X} | a \notin \mathcal{S})$$

et on pose :

$$\begin{cases} x &= \Pr(a \in \mathcal{X} | a \notin \mathcal{S}) \\ y &= \Pr(a \in \mathcal{X} | a \in \mathcal{S}) \end{cases}$$

42.  $\partial_{11}F(X, Y) = \exp(X - Y)$ ,  $\partial_{12}F(X, Y) = -\exp(X - Y)$  et  $\partial_{22}F(X, Y) = \exp(X - Y)$ .

43. On suppose que  $\text{cov}(\bar{X}_n, \bar{Y}_n) = 0$  et que  $\tilde{\sigma}^2$  est la même en 0 et en  $t$ .

On cherche  $(x, y)$  se sorte que  $y$  soit maximal (dans l'intervalle  $[0, 1]$ ) et

$$\pi_a^{\mathcal{X}} = \pi_a^{\mathcal{J}} \cdot y + (1 - \pi_a^{\mathcal{J}}) \cdot x$$

tandis que  $x$  est également dans l'intervalle  $[0, 1]$ . Tout ceci revient à maximiser la fonction ( $y \equiv \varphi(x)$  dans l'expression précédente) :

$$\varphi(x) = \frac{\pi_a^{\mathcal{X}}}{\pi_a^{\mathcal{J}}} - \frac{1 - \pi_a^{\mathcal{J}}}{\pi_a^{\mathcal{J}}} x$$

$\varphi$  est décroissante. Pour  $x \in [0, 1]$ , cette fonction est maximale au point  $x = 0$ . Or  $\varphi(0) = \pi_a^{\mathcal{X}} / \pi_a^{\mathcal{J}}$ . Mais  $y \in [0, 1]$ . Par conséquent, deux situations se présentent :

— Si  $\pi_a^{\mathcal{X}} \leq \pi_a^{\mathcal{J}}$  alors la solution

$$(x, y) = \left( 0, \frac{\pi_a^{\mathcal{X}}}{\pi_a^{\mathcal{J}}} \right)$$

est optimale.

— Si  $\pi_a^{\mathcal{X}} > \pi_a^{\mathcal{J}}$ , alors le maximum de  $\varphi(x)$  conditionnellement au fait que  $\varphi(x) \leq 1$  est atteint au point  $\varphi^{-1}(1)$  en lequel  $\varphi$  vaut 1. L'optimum est alors :

$$(x, y) = \left( \frac{\pi_a^{\mathcal{X}} - \pi_a^{\mathcal{J}}}{1 - \pi_a^{\mathcal{J}}}, 1 \right)$$

## E Le cas des DOMs

Chaque territoire DOM est caractérisé, dans l'IPC base 1998, par un unique agrégat géographique. En d'autres termes, partant de prix relevés dans des points de vente du territoire, on ne procède pas au calcul d'agrégat par unité urbaine, mais au calcul d'un agrégat d'ensemble pour le territoire. Il en découle que le nombre de ces agrégat élémentaires, correspondant dans le cas de la métropole aux unités urbaines, est égal à 1 pour chaque territoire DOM. En conséquence, il n'y a pas, dans le cas des DOM, de sondage géographique : par hypothèse, les points de vente retenus pour les observation de l'indice des prix sont représentatifs de l'ensemble des points de vente du territoire et doivent, à ce titre, être répartis sur l'ensemble du territoire.

Ce choix d'un unique agrégat géographique est lié à la nature des biais de formule qui surviennent lorsque peu de produits sont utilisés pour l'estimation des agrégats (micro-indices – voir à ce propos l'annexe C et en particulier la discussion du paragraphe C.3). Pour cette raison, le principe consistant à calculer un seul agrégat par territoire DOM est maintenu en base 2015.

Il est utile de procéder à une optimisation de l'échantillon conformément à ce qui est indiqué au paragraphe 4. Cette optimisation a été réalisée sur les données de l'IPC de 2013. La précision d'un indice calculé sur le champ de la collecte terrain est indiqué, par territoire, dans la table 7.

TABLE 7 – Precision des indices IPC en métropole et dans les DOM

Territoire	Nobs base 1998	$\sigma$	
		IPC base 1998	IPC optimisé
métropole	112 767	0,14	0,022
Guadeloupe	5 666	0,63	0,10
Martinique	5 724	0,67	0,12
Guyane	4 206	0,64	0,10
Réunion	5 826	0,70	0,09

**Note** : écart-type, en points d'indice de décembre de l'année A, base=100 décembre A-1, pour un indice de prix fondé sur le champ restreint à la collecte terrain, hors produits frais. Le calcul correspond à la seule variance intra (ajouter 0,01 d'écart-type inter en métropole). L'optimisation correspond au calcul réalisé pour la base 2015, conformément aux principes exposés aux paragraphes 4, 6 et 7 ; le seul paragraphe 4 s'agissant des DOM.

En optimisant le nombre de relevés par variété par variété, sous contrainte que le coût de collecte d'ensemble, pour chaque territoire, reste identique à son niveau actuel, il serait possible d'atteindre le niveau de précision indiqué dans la table 7.

## F Travaux exploratoires pour le nouveau tirage d'agglomérations dans le cadre du changement de base de l'IPC

Pour le nouveau tirage d'agglomérations (géographie des unités urbaines de 2010) et pour le calcul des pondérations, on voudrait prendre en compte, non seulement le poids démographique des agglomérations (considéré comme un proxy des ventes), mais aussi le poids de la consommation ou des ventes au lieu d'achat dans ces agglomérations. Plutôt que de fixer les probabilités d'inclusion proportionnelles à la population des agglomérations, utilisée comme un proxy des ventes, on pourrait les fixer proportionnellement à la répartition des ventes selon les tailles d'unité urbaine.

Des travaux préliminaires sur les carnets de l'enquête budget des familles (2005 et 2011) et sur l'enquête points de vente 2009 de la division commerce ont été effectués sur le champ de la France métropolitaine. Ils aboutissent à des chiffrages qui pourront servir de référence dans l'établissement des cibles de relevés par groupes de produits et taille d'agglomération en France métropolitaine.

Pour cela plusieurs exploitations ont été faites au niveau France métropolitaine :

- calcul des dépenses faites par les ménages selon le type de produit (regroupement en 12 postes) et la taille d'unité urbaine du lieu d'achat : à partir des carnets de



- dépenses des enquêtes Budget des familles de 2005 et de 2011 ;
- calcul du chiffre des établissements relevant du commerce de détail en magasin et de l'artisanat commercial (NAF 10.13B, 10.71B, 10.71C, 10.71D et division 47 hors groupe 47.8 et 47.9) en France métropolitaine selon la taille d'unité urbaine et selon la forme de vente à partir de l'enquête Points de vente 2009.

## **F.1 Calcul de dépenses au lieu d'achat à partir des carnets de dépenses issus des enquêtes Budget des Familles de 2005 et de 2011**

L'enquête Budget des familles recense les dépenses de consommation par l'intermédiaire d'un questionnaire rétrospectif et d'un carnet de compte, remplis par les individus du ménage de plus de 14 ans. Les carnets de compte remplis par les enquêtés de plus de 14 ans enregistrent les dépenses quotidiennes, ainsi que les petites dépenses irrégulières. La personne interrogée doit noter sur son carnet toutes les dépenses qu'elle effectue durant 14 jours (pour l'enquête de 2005) ou 7 jours (pour l'enquête de 2011). Dans les carnets de dépenses des ménages, les personnes indiquent le lieu d'achat ; on dispose donc des libellés de communes d'achat et du code postal de la commune d'achat. A partir de ces libellés de communes et des codes postaux, on a pu ramener un code géographique de commune en face de chaque ligne de dépense, puis une taille d'unité urbaine. L'exercice a été fait, pour la France métropolitaine, sur les données de BDF 2005, puis sur les données de BDF 2011. Le nombre de lignes de dépenses exploitées est deux fois moins important en 2011 qu'en 2005, car le remplissage du carnet ne concernait plus qu'une semaine de dépenses en 2011 au lieu de deux en 2005 : 495 900 lignes pour 2011 contre 990 861 lignes pour 2005, sur le champ France métropolitaine et fonctions de consommation (codage COICOP – voir table 8) 01 à 12.

L'outil carnet collectant les dépenses sur 2 semaines pour 2005 et sur une semaine pour 2011 ne permet pas de comptabiliser toutes les dépenses des ménages. Les dépenses importantes ou périodiques sont collectées par des questions directes, avec des périodes de référence variables et appropriées à chaque catégorie de dépenses. Les dépenses relevées dans les carnets codés (avec une taille d'unité urbaine du lieu d'achat) ne représentent que 58% des dépenses des ménages en 2005 et 55% en 2011. On ne prend pas en compte le poste 99 (produit non identifié) et le poste 13 (hors champ de la dépense de consommation finale des ménages) pour se ramener à un champ le plus proche possible de celui de l'IPC. Suivant le type de produit, le taux de couverture des carnets dans la dépense globale et le taux de couverture des carnets codés dans la dépense globale est plus ou moins important (table 8) : les carnets codés représentent plus de 90% de la consommation des ménages pour les regroupements 01 et 02 de la COICOP2 : produits alimentaires et boissons non alcoolisées, boissons alcoolisées-tabacs et stupéfiants. Ils pèsent près de 20% des dépenses des ménages en 2011. En 2011, les regroupements couverts entre 70% et 80% sont le "05 :

Ameublement, équipement ménager en entretien courant de la maison” et le “03 : Articles d’habillement et articles chaussants”.

On propose d’utiliser la répartition des dépenses de consommation pour le total 01+02 selon la taille d’unité urbaine, dans le calcul des probabilités d’inclusion des agglomérations selon leur taille (plutôt que d’utiliser uniquement la démographie) et dans le calcul des poids des agglomérations (cf. tableau 9). On pourra également utiliser ces répartitions pour l’établissement des tableaux d’équilibrage par variété, forme de vente et taille d’agglomération.

TABLE 8 – Calcul de taux de couverture des carnets codés (avec une taille d'unité urbaine du lieu d'achat)

COICOP 2 positions	poids dans la dépense totale		taux de couverture carnet codé	
	2005	2011	2005	2011
1	15,5%	16,2%	<b>0,92</b>	<b>0,93</b>
2	2,6%	2,9%	<b>0,91</b>	<b>0,88</b>
3	7,9%	4,9%	0,48	0,76
4	14,9%	15,5%	0,37	0,20
5	7,5%	6,2%	0,59	0,69
6	3,7%	1,8%	0,78	1,63
7	16,0%	17,3%	0,55	0,60
8	3,7%	3,3%	0,27	0,12
9	9,2%	9,6%	0,64	0,52
10	0,7%	0,7%	0,27	0,13
11	5,7%	7,2%	0,61	0,51
12	12,6%	14,5%	0,43	0,30
total	100,0%	100,0%	0,58	0,55

Source : Insee, enquêtes budget des familles 2005 et 2011 ; 01 : Produits alimentaires et boissons non alcoolisées, 02 : Boissons alcoolisées, tabac et stupéfiants, 03 : Articles d'habillement et articles chaussants, 04 : Logement, eau, électricité, gaz et autres combustibles, 05 : Ameublement, équipement ménager en entretien courant de la maison, 06 : Santé (en 2011 le taux de couverture du carnet dépasse 100%, car les dépenses comprennent la part sécurité sociale et la part complémentaire, alors que dans la table de consommation, c'est un reste à charge), 07 : Transports, 08 : Communications, 09 : Loisirs et culture, 10 : Enseignement, 11 : Hôtels, cafés, restaurants, 12 : Autres biens et services.

TABLE 9 – Répartition de la dépense de consommation et de la démographie par taille d'agglomérations

Taille de l'unité urbaine (UU de 2010)	poids démographique	poids dépenses au lieu d'achat (carnet BDF)		poids dépenses '01'+02' au lieu d'achat (carnet BDF)		2011/2005	2011/2005 ('01'+02')
		2005	2011	2005	2011		
<b>France METRO</b>	<b>2007</b>	<b>2005</b>	<b>2011</b>	<b>2005</b>	<b>2011</b>		
0-Rural	22,5	9,0	7,1	8,5	7,4	0,79	0,87
1-UU de 2 000 à 4 999 h	6,7	6,4	6,1	7,8	7,6	0,95	0,98
2-UU de 5 000 à 9 999 h	5,7	6,3	7,0	8,2	9,1	1,10	1,11
3-UU de 10 000 à 19 999 h	5,0	6,6	5,8	8,3	6,7	0,89	0,81
4-UU de 20 000 à 49 999 h	6,3	8,1	9,1	8,9	9,8	1,12	1,10
5-UU de 50 000 à 99 999 h	7,3	9,3	8,5	9,0	8,8	0,92	0,98
6-UU de 100 000 à 199 999 h	5,3	6,5	8,6	5,7	7,0	1,32	1,22
7-UU de 200 000 à 1 999 999 h	24,6	28,5	29,3	26,4	25,6	1,03	0,97
8-Agglomération de Paris	16,7	19,3	18,6	17,3	18,1	0,96	1,04
ensemble (France métropolitaine)	100,0	100,0	100,0	100,0	100,0		

Taille de l'unité urbaine (UU de 2010) - taille aggro IPC	poids démographique	poids dépenses au lieu d'achat (carnet BDF)		poids dépenses '01'+02' au lieu d'achat (carnet BDF)		2011/2005	2011/2005 ('01'+02')
		2005	2011	2005	2011		
<b>France METRO</b>	<b>2007</b>	<b>2005</b>	<b>2011</b>	<b>2005</b>	<b>2011</b>		
Rural (pas dans ipc)	22,5	9,0	7,1	8,5	7,4	0,79	0,87
D - UU de 2 000 à 19 999 h	17,4	19,3	18,9	24,3	23,4	0,98	0,96
C - UU de 20 000 à 99 999 h	13,6	17,4	17,6	17,9	18,6	1,01	1,04
B - UU de 100 000 h ou plus	29,9	35,0	37,9	32,1	32,6	1,08	1,02
A - Agglomération de Paris	16,7	19,3	18,6	17,3	18,1	0,96	1,04
ensemble	100,0	100,0	100,0	100,0	100,0		

01\*\*\*\* PRODUITS ALIMENTAIRES ET BOISSONS NON ALCOOLISEES  
02\*\*\*\* BOISSONS ALCOOLISEES, TABACS ET STUPEFIANTS

Source : Insee, enquêtes budget des familles 2005 et 2011

Le champ actuel de l'indice des prix porte sur l'ensemble des agglomérations de plus de 2000 habitants. L'étude menée à partir des enquêtes budgets des familles permet d'estimer le défaut de couverture géographique à 7% environ. En effet, alors que les habitants des communes du rural représentent 22,5% de la population métropolitaine (population municipale de 2007), 7% de la consommation totale est achetée dans le rural en 2011. Si

on considère uniquement les regroupements 01 et 02, bien couverts, seulement 7,4% des dépenses sont faites dans un lieu d'achat situé dans le rural, en 2011.

## **F.2 Calcul du chiffre des établissements relevant du commerce de détail en magasin et de l'artisanat commercial (NAF 10.13B, 10.71B, 10.71C, 10.71D et division 47 hors groupe 47.8 et 47.9) en France métropolitaine selon la taille d'unité urbaine et la forme de vente : à partir de l'enquête Points vente 2009**

L'enquête Points de vente 2009 a pour objectif de collecter le siret, le chiffre d'affaires, le nombre de personnes occupées et la surface de vente au niveau des établissements commerciaux. Pour ce faire, l'enquête interroge des unités légales du commerce de détail en magasin et de l'artisanat commercial qui répondent pour l'ensemble de ces magasins (NAF 10.13B, 10.71B, 10.71C, 10.71D et division 47 hors groupe 47.8 et 47.9) en France métropolitaine. Dans cette enquête, on se restreint aux unités légales et aux établissements marchands, actifs ou présumés actifs, exploitants, créés strictement avant le 2 janvier 2009, localisés en France métropolitaine et relevant du commerce de détail en magasin. L'échantillon a été stratifié selon trois modalités du zonage en aires urbaines (ZAUER 1999) :

- pôle urbain {1}
- commune mono ou multi-polarisée {2,3}
- espace à dominante rurale {4,5,6}

Ni le code géographique de la commune de localisation de l'établissement, ni la taille d'unité urbaine correspondante ne sont disponibles dans les données de l'enquête. À partir du code Siret et par appariement avec les données par établissement, issues de Sirene, on a pu récupérer le code géographique, puis la taille d'unité urbaine correspondants de l'établissement. A partir de ces données, on reconstitue un équivalent des formes de vente enquêtées dans l'IPC (table 10) et on calcule une répartition du chiffre d'affaire des points de vente d'une taille d'unité urbaine donnée selon les formes de vente reconstituées (table 11).

TABLE 10 – Les formes de ventes et l'enquête "Points de vente"

Codes FV	Signification	Exemple	Enquête PDV 2009
10	Hypermarché	Carrefour	Apet=4711F
20	Supermarché	Intermarché	Apet=4711D
25	Maxidiscount	Dia	
30	Supérette	Casino	Apet=4711C+4711B
40	Magasin populaire	Monoprix	Apet=4711E+4719B
50	Grand magasin	Printemps	Apet=4719A
60	Grande surface spécialisée	Picard, Décathlon	3 premières positions de Apet=474, 475, 476 + Apet=4711A, 4771Z, 4772A, 4772B, 4775Z, 4777Z, 4778A, 4773Z, 4774Z, 4776Z, 4778B, 4778C
70	Magasin traditionnel (-300m <sup>2</sup> )	Boulangerie, Zara	3 premières positions de Apet=472, 107 + Apet=1013B
80	Marché	Marché convention	
90	Services	Café, restaurant, coiffeurs	
99	autres	Services publics et administration, relevés en site prix, achat en bornes	

Source : Insee, enquête Points de vente 2009. A noter : La forme de vente Maxidiscount n'a pas pu être isolée, elle est dans la forme de vente "Supermarché". Les commerces d'alimentation générale d'une surface de vente inférieure à 120 m<sup>2</sup> ont été regroupés avec les supérettes (qui ont une surface comprise entre 120 et 400 m<sup>2</sup>). Les magasins populaires correspondent aux magasins multi-commerces (Apet=4711E). On y a aussi inclus le "non alimentaire non spécialisé - autres commerces (Apet=4719B)" qui sont des commerces de détail non spécialisé sans prédominance alimentaire en magasin d'une surface de vente inférieure à 2500 m<sup>2</sup>. Quand ce type de commerce a une surface de vente égale ou supérieure à 2500 m<sup>2</sup>, il est classé en "grands magasins". Les grandes surfaces spécialisées regroupent les commerces de détail de produits surgelés et commerces de détail dans les secteurs "équipement de l'information et de la communication", "équipement du foyer", "culture, loisirs, sport", "habillement, chaussures", "autres équipements de la personne", "produits pharmaceutiques, articles médicaux et orthopédiques", "autres commerces de détail". L'alimentation spécialisée et l'artisanat commercial sont classés en forme de vente "magasin traditionnel".

TABLE 11 – Répartition du chiffre d'affaires par taille d'unité urbaine et par forme de vente

Taille de l'unité urbaine (UU de 2010) France métro	Hypermarchés	Supermarchés	Supérettes	Magasins multi-commerces populaires	Grands magasins	Grandes surfaces spécialisées	magasins traditionnels	Ensemble
Rural (pas dans ipc)	3,0	33,6	6,8	0,5	0,1	39,9	16,2	100,0
D - UU de 2 000 à 19 999 h	23,3	27,6	1,7	0,8	0,2	37,9	8,5	100,0
C - UU de 20 000 à 99 999 h	35,5	18,1	1,1	1,2	0,2	37,1	6,8	100,0
B - UU de 100 000 h à 199 999 H	26,4	16,0	0,8	0,7	0,6	48,9	6,6	100,0
B - UU de 200 000 h ou plus	30,0	14,2	1,6	1,1	1,2	45,1	6,8	100,0
A - Agglomération de Paris	22,4	10,1	3,9	2,8	3,8	47,8	9,2	100,0

source : insee, enquête Points de vente 2009