

Catalogue no. 12-539-X

# Statistique Canada : lignes directrices concernant la qualité



Cinquième édition – octobre 2009



Statistique  
Canada

Statistics  
Canada

Canada

Pour toute demande de renseignements au sujet de ce produit ou sur l'ensemble des données et des services de Statistique Canada, visiter notre site Web à [www.statcan.gc.ca](http://www.statcan.gc.ca). Vous pouvez également communiquer avec nous par courriel à [infostats@statcan.gc.ca](mailto:infostats@statcan.gc.ca) ou par téléphone entre 8 h 30 et 16 h 30 du lundi au vendredi aux numéros suivants :

#### **Centre de contact national de Statistique Canada**

Numéros sans frais (Canada et États-Unis) :

Service de renseignements	1-800-263-1136
Service national d'appareils de télécommunications pour les malentendants	1-800-363-7629
Télécopieur	1-877-287-4369

Appels locaux ou internationaux :

Service de renseignements	1-613-951-8116
Télécopieur	1-613-951-0581

#### **Programme des services de dépôt**

Service de renseignements	1-800-635-7943
Télécopieur	1-800-565-7757

#### **Comment accéder à ce produit**

Le produit n° 12-539-X au catalogue est disponible gratuitement sous format électronique. Pour obtenir un exemplaire, il suffit de visiter notre site Web à [www.statcan.gc.ca](http://www.statcan.gc.ca) et de choisir la rubrique « Publications ».

#### **Normes de service à la clientèle**

Statistique Canada s'engage à fournir à ses clients des services rapides, fiables et courtois. À cet égard, notre organisme s'est doté de normes de service à la clientèle que les employés observent. Pour obtenir une copie de ces normes de service, veuillez communiquer avec Statistique Canada au numéro sans frais 1-800-263-1136. Les normes de service sont aussi publiées sur le site [www.statcan.gc.ca](http://www.statcan.gc.ca) sous « À propos de nous » > « Offrir des services aux Canadiens ».

# Statistique Canada : lignes directrices concernant la qualité

Cinquième édition – octobre 2009

Publication autorisée par le ministre responsable de Statistique Canada

© Ministre de l'Industrie, 2009

Tous droits réservés. Le contenu de la présente publication électronique peut être reproduit en tout ou en partie, et par quelque moyen que ce soit, sans autre permission de Statistique Canada, sous réserve que la reproduction soit effectuée uniquement à des fins d'étude privée, de recherche, de critique, de compte rendu ou en vue d'en préparer un résumé destiné aux journaux et/ou à des fins non commerciales. Statistique Canada doit être cité comme suit : Source (ou « Adapté de », s'il y a lieu) : Statistique Canada, année de publication, nom du produit, numéro au catalogue, volume et numéro, période de référence et page(s). Autrement, il est interdit de reproduire le contenu de la présente publication, ou de l'emmagasiner dans un système d'extraction, ou de le transmettre sous quelque forme ou par quelque moyen que ce soit, reproduction électronique, mécanique, photographique, pour quelque fin que ce soit, sans l'autorisation écrite préalable des Services d'octroi de licences, Division des services à la clientèle, Statistique Canada, Ottawa, Ontario, Canada K1A 0T6.

Décembre 2009

N° 12-539-X au catalogue  
ISSN 1708-6264

Périodicité : irrégulier

Ottawa

This publication is available in English upon request (catalogue no. 12-539-XIE)

---

## Note de reconnaissance

Le succès du système statistique du Canada repose sur un partenariat bien établi entre Statistique Canada et la population, les entreprises, les administrations canadiennes et les autres organismes. Sans cette collaboration et cette bonne volonté, il serait impossible de produire des statistiques précises et actuelles.

## Préface

Les Canadiens et leur gouvernement ont créé Statistique Canada pour avoir accès à une source d'information en laquelle ils peuvent avoir confiance. La confiance peut uniquement s'établir si les données que Statistique Canada produit correspondent aux besoins du pays et représentent bien le monde que nous tentons de décrire. En d'autres mots, l'information doit être pertinente et de grande qualité.

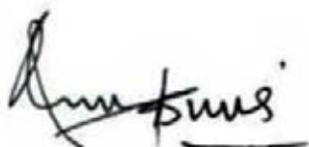
La qualité est donc essentielle à l'exécution du mandat de Statistique Canada, qui consiste à produire de l'information. Elle comporte d'ailleurs une caractéristique fondamentale qu'il est important de bien comprendre : la qualité, définie comme la représentativité de l'univers que nous tentons de cerner, se détériore automatiquement en l'absence de mesures proactives. C'est la raison pour laquelle les méthodes que nous employons pour maintenir la représentativité de nos données doivent évoluer au fur et à mesure que le monde qui nous entoure évolue.

À la lumière de ces réalités avec lesquelles il doit composer, Statistique Canada fournit depuis longtemps une orientation quant aux principes à appliquer dans ses plans d'enquête. Pour ce faire, l'organisme regroupe ses expériences et ses conclusions à l'égard de ce que constituent les « pratiques exemplaires » dans un ensemble de Lignes directrices concernant la qualité. La première édition des Lignes directrices concernant la qualité est parue en 1985. Des éditions révisées ont été publiées en 1987, 1998 et 2003. Compte tenu de la nécessité d'actualiser les lignes directrices régulièrement, le présent document, résultat d'une importante mise à jour de l'édition précédente, a été produit afin de refléter les progrès accomplis en méthodologie d'enquêtes au cours des six dernières années.

Les lignes directrices présentées dans le présent document ne s'appliquent pas toutes également à chacun des processus d'acquisition de données. Il faut donc juger minutieusement leur pertinence et leur importance en fonction des exigences et des contraintes propres à chaque programme. L'utilisation de ce document commande donc beaucoup d'attention professionnelle et de jugement.

Bien que les lignes directrices qui figurent dans ce document ne sauraient remplacer l'expertise et le jugement du personnel chargé de concevoir les enquêtes, le souci de la qualité qui en émane doit se refléter dans l'ensemble de nos activités. Il incombe à tous les employés intervenant dans les activités statistiques de veiller à ce que la qualité soit considérée comme hautement prioritaire lors de la conception et de la mise en œuvre des méthodes et des procédures statistiques qu'ils gèrent.

Je tiens à remercier les nombreux experts de Statistique Canada qui ont collaboré à la préparation des Lignes directrices concernant la qualité au fil des ans. En outre, les conseils du Comité des méthodes et des normes ont permis d'améliorer le document.



Munir A. Sheikh  
Statisticien en chef

# Table des matières

---

---

	Page
Introduction .....	6
Figure 1 Niveaux 1 et 2 du Modèle statistique général du processus opérationnel .....	10
Étapes de l'enquête.....	14
1. Objectifs, utilisations et utilisateurs .....	16
2. Concepts, variables et classifications .....	19
3. Couverture et bases de sondage .....	22
4. Plan d'échantillonnage .....	27
5. Conception du questionnaire .....	32
6. Collecte, saisie et codage des données .....	36
7. Utilisation des données administratives .....	45
8. Réponse et non-réponse .....	52
9. Vérification .....	58
10. Imputation .....	62
11. Pondération et estimation .....	68
12. Désaisonnalisation et estimation de la tendance-cycle .....	72
13. Étalonnage et techniques connexes .....	78
14. Évaluation de la qualité des données .....	81
15. Contrôle de la divulgation .....	85
16. Diffusion et communication des données .....	89
17. Analyse et présentation des données .....	93
18. Documentation .....	98

# Introduction

---

L'information statistique est essentielle au fonctionnement des démocraties modernes. L'absence de données de qualité mettrait gravement en péril les processus décisionnels, la répartition de milliards de dollars en ressources et la capacité des gouvernements, des entreprises, des établissements et du grand public de comprendre la réalité sociale et économique du pays. Un organisme statistique national comme Statistique Canada joue un rôle crucial dans la production et la diffusion de renseignements statistiques.

La crédibilité d'un organisme statistique dans ses efforts pour s'acquitter de ce rôle clé repose sur les piliers suivants : la production de renseignements statistiques de grande qualité, un bon rapport coût-efficacité, la protection des renseignements personnels, la confidentialité et le maintien en poste d'un effectif très compétent et motivé.

La qualité et, plus particulièrement, la pertinence de l'information qu'il produit revêtent une importance fondamentale pour l'organisme statistique. Si ce dernier est incapable de produire des données de grande qualité, les utilisateurs comme les fournisseurs perdraient vite confiance en l'organisme, rendant sa mission impossible à réaliser. Cette section, qui se veut une introduction aux lignes directrices à suivre en la matière, présente les principes de l'assurance-qualité appliqués à Statistique Canada.

## **Les principes de l'assurance-qualité à Statistique Canada**

La structure de gestion, les politiques et les lignes directrices, les mécanismes de consultation, l'approche de réalisation et de gestion des projets et l'environnement du Bureau ont été mis en place pour faciliter et assurer une gestion efficace de la qualité. Le Cadre d'assurance de la qualité de Statistique Canada (Statistique Canada, 2002c) décrit les mécanismes fondamentaux de la gestion de la qualité.

Ce cadre est constitué d'un large éventail de mécanismes et de processus qui influent sur divers niveaux des programmes du Bureau et à l'échelle de l'organisme. L'efficacité de ce cadre ne dépend pas d'un seul processus ou mécanisme, mais bien de l'effet combiné de nombreuses mesures interdépendantes qui s'appuient sur les intérêts professionnels et la motivation du personnel et qui se renforcent mutuellement avec comme objectif la satisfaction des besoins des clients. Ces mesures mettent l'accent sur le professionnalisme du Bureau et illustrent le souci rattaché à la qualité des données. Un des traits dominants de cette stratégie est la synergie attribuable au fait que les nombreux intervenants des programmes du Bureau travaillent dans un cadre où on privilégie la cohérence des processus et l'uniformité des messages. Les Lignes directrices concernant la qualité, qui s'inscrivent à l'intérieur de ce cadre, offrent un document d'accompagnement dans lequel sont décrites les pratiques exemplaires jalonnant toutes les « étapes » d'un programme statistique. Elles sont destinées aux membres de l'équipe de projet chargés de l'élaboration et de la mise en œuvre des programmes statistiques.

Huit principes directeurs orientent l'ensemble de ces mécanismes, processus et pratiques.

## **La qualité est relative, et non un absolu**

Une des caractéristiques importantes de la gestion de la qualité est l'atteinte d'un équilibre entre les objectifs de qualité et les contraintes des ressources financières et humaines, la volonté des répondants de fournir les données de base et les demandes concurrentes pour une plus grande quantité de renseignements plus détaillés. La gestion de la qualité ne signifie pas l'optimisation de la qualité à tout prix, mais l'atteinte d'un juste équilibre entre la quantité et la qualité de l'information produite par les programmes du Bureau et les ressources disponibles. Dans chacun des programmes, le défi consiste à faire des compromis judicieux entre les besoins changeants des clients, les coûts, le fardeau des répondants et les divers éléments ou dimensions de la qualité.

Les données statistiques tirent leur importance des fins auxquelles elles sont destinées. Il s'en suit que les données statistiques ne peuvent être utiles que si elles sont pertinentes et qu'elles représentent adéquatement le monde qu'elle cherche à décrire. Ce concept, adopté par de nombreux organismes statistiques, est désigné sous l'expression « adaptation à l'utilisation » ou le terme « adéquation ». Il suppose de la part de l'organisme statistique une connaissance approfondie des fins auxquelles sont destinées les données, d'où la nécessité d'entretenir des liens constants avec les utilisateurs.

Suivre ce principe, c'est reconnaître que la qualité « parfaite » n'est ni souhaitable ni accessible – en fait, elle est même rarement possible. Les données sont sujettes à de fréquentes erreurs, à l'échantillonnage et à d'autres étapes, et il revient à l'organisme statistique d'atteindre un équilibre entre différents facteurs, tels l'exactitude, le coût et le fardeau du répondant, lors de l'élaboration d'un programme statistique. L'objectif ne consiste pas à réduire l'erreur en soi; chaque programme statistique doit être conçu dans les limites de ce qui est faisable, en tenant compte de l'importance que revêtent les données pour les utilisateurs.

Statistique Canada s'efforce d'intégrer à tous ses programmes et produits les principes de pertinence et de qualité. La qualité de ses statistiques officielles repose sur l'emploi de méthodes scientifiques éprouvées et adaptées progressivement aux besoins changeants des clients, à la réalité en évolution que le Bureau tente de mesurer et à la capacité ou la volonté des répondants de fournir des données fiables et actuelles. Les présentes Lignes directrices concernant la qualité sont un des outils qui aideront le Bureau à intégrer la qualité à la conception de chacun des programmes.

### **La qualité est multidimensionnelle**

Ces vingt dernières années, les organismes statistiques sont parvenus à un consensus : la « qualité » de l'information statistique est multidimensionnelle. Statistique Canada définit la qualité en fonction de six aspects, et d'autres organismes statistiques se sont dotés de cadres similaires. Bien que les définitions établies varient légèrement, elles mettent toutes en relief le fait qu'il n'existe pas de mesure unique de la qualité des données.

À Statistique Canada, les différents aspects de la qualité se définissent comme suit :

**La pertinence de l'information statistique** reflète la mesure dans laquelle cette information répond aux besoins réels des clients. C'est en examinant cet élément qu'on détermine si l'information disponible permet de mieux comprendre les enjeux qui sont importants pour les utilisateurs. Aussi la pertinence apparaît-elle comme l'aspect le plus important de la qualité, si bien qu'on pourrait la considérer comme l'un des piliers de l'organisme statistique. Cet aspect relève du domaine des utilisateurs de l'information : il ne peut être déterminé par l'organisme statistique lui-même. Par comparaison, les autres aspects de la qualité sont davantage du ressort de l'organisme statistique.

**L'exactitude de l'information statistique** est la mesure dans laquelle l'information décrit correctement le phénomène qu'elle devait évaluer. Généralement, elle est caractérisée par l'erreur dans les estimations statistiques et est décomposée en composantes de biais (erreur systématique) et de variance (erreur aléatoire). L'exactitude peut également être décrite en fonction des sources d'erreur majeures qui peuvent mener à l'inexactitude (p. ex., couverture, échantillonnage, non-réponse, réponse).

**L'actualité de l'information statistique** renvoie à l'intervalle entre le point de référence (ou la fin de la période de référence) auquel se rapporte l'information et la date à laquelle l'information est diffusée. Habituellement, l'actualité se trouve en relation d'équilibre avec l'exactitude. L'actualité de l'information influera sur sa pertinence.

**L'accessibilité de l'information statistique** renvoie à la facilité avec laquelle on peut obtenir l'information auprès du Bureau. Cela comprend la facilité avec laquelle on peut certifier l'existence de l'information, ainsi que l'à-propos de la forme ou du médium par le biais duquel on peut accéder à l'information. En outre, le coût de l'information peut représenter un facteur d'accessibilité pour certains utilisateurs.

**L'intelligibilité de l'information statistique** reflète la disponibilité de l'information et des métadonnées supplémentaires nécessaires à l'interprétation et à l'utilisation appropriées des renseignements. Normalement, cette information comprend les variables, les classifications et les concepts sous-jacents utilisés, la méthode de collecte des données et le traitement, ainsi que les indications ou les mesures de l'exactitude de l'information statistique.

**La cohérence** reflète la mesure dans laquelle on peut réussir à regrouper cette information avec d'autres renseignements statistiques dans un cadre analytique général et au fil du temps. L'utilisation de classifications, des populations cibles et de concepts normalisés favorise la cohérence, tout comme l'utilisation d'une méthode commune d'une enquête à l'autre. La cohérence ne sous-tend pas nécessairement l'uniformisation numérique complète.

Ces dimensions de qualité se chevauchent et se recoupent. La gestion de la qualité oblige donc à les considérer toutes. La négligence d'un seul aspect peut entraîner la faillite d'un programme statistique entier.

### **Chaque employé doit intervenir dans l'assurance de la qualité**

Au moyen de ses politiques, de ses lignes directrices et de ses communications internes, Statistique Canada insiste auprès de ses employés sur le fait que chacun a un rôle à jouer en matière d'assurance de la qualité, des personnes affectées aux tâches de production quotidiennes aux cadres des échelons les plus élevés. Cette approche s'inspire de la philosophie de Deming, soit que la qualité n'est pas un élément à inspecter en cours de processus, mais à intégrer dès le départ. C'est la raison pour laquelle aucune entité n'est explicitement responsable de l'assurance de la qualité à Statistique Canada.

Les pratiques exercées par le Bureau en matière de ressources humaines témoignent également du principe selon lequel la qualité est l'affaire de tous. Les programmes de recrutement, de formation et de perfectionnement mettent grandement l'accent sur les compétences techniques ainsi que sur la connaissance des caractéristiques de données de grande qualité.

### **Pour équilibrer tous les aspects de la qualité, la gestion par équipe de projet est l'approche indiquée**

Comme la qualité est multidimensionnelle, les différents aspects de la qualité tendent à être les domaines d'expertise de différents groupes au sein du Bureau. Au fil des ans, Statistique Canada a constaté que l'approche de gestion par équipe de projet offre le meilleur moyen d'atteindre l'équilibre entre les divers aspects de la qualité.

À Statistique Canada, la gestion de la qualité est assurée dans le contexte d'un cadre de gestion matricielle, c'est-à-dire que la gestion des projets relève d'une organisation fonctionnelle. À cet égard, le Bureau est divisé en six secteurs. Trois de ces secteurs sont principalement responsables des programmes statistiques de production et d'analyse des données dans différents domaines spécialisés (p. ex. statistiques sociales, statistiques des entreprises, comptes nationaux). Les trois autres secteurs sont avant tout chargés du maintien d'une infrastructure et de la prestation de services pour les programmes statistiques (p. ex. méthodologie, informatique, activités de collecte, diffusion, systèmes de gestion). Un programme statistique est généralement géré par une division spécialisée, et ses responsables font un usage intensif des ressources de l'infrastructure et des secteurs de service.

Le recours à une approche interdisciplinaire de gestion par équipe de projet pour la conception ou la restructuration d'un programme statistique est important. Il permet de s'assurer qu'on accorde suffisamment d'attention à la qualité de toutes les composantes et étapes du programme au cours de sa conception, de sa mise en œuvre et de son évaluation. Il incombe aux organisations fonctionnelles de s'assurer que les équipes de projet sont composées de personnes possédant la crédibilité et les compétences requises pour représenter leur secteur fonctionnel. Le personnel spécialisé met à contribution sa connaissance du contenu, des besoins des clients et de la pertinence des données. Pour leur part, les méthodologistes apportent leurs connaissances relativement à l'équilibre à atteindre entre les méthodes statistiques et la qualité des données, particulièrement en ce qui concerne l'exactitude, l'actualité et le coût. Les experts des opérations ont quant à eux l'expérience voulue des méthodes opérationnelles et s'attardent aux questions d'ordre pratique, à l'efficacité, au personnel sur le terrain et aux répondants. Enfin, les experts des systèmes connaissent bien les normes technologiques et les outils à utiliser pour concevoir les projets. Ils apportent aussi une dimension systémique aux projets.

C'est au sein d'une telle équipe que sont prises les nombreuses décisions et que sont faits les compromis nécessaires à l'atteinte d'un juste équilibre entre la qualité, d'une part, et le coût et le fardeau de réponse, d'autre part. Ensemble, les membres de l'équipe doivent trouver un équilibre entre les pressions contradictoires afin de mettre au point une conception optimale. Le fait que chacun des membres de l'équipe fasse partie d'une organisation fonctionnelle spécialisée, à laquelle ils peuvent avoir recours au besoin pour obtenir des ressources de gestion et diverses ressources plus spécialisées, permet de résoudre les problèmes techniques et les différends qui peuvent survenir au cours de la réalisation d'un projet.

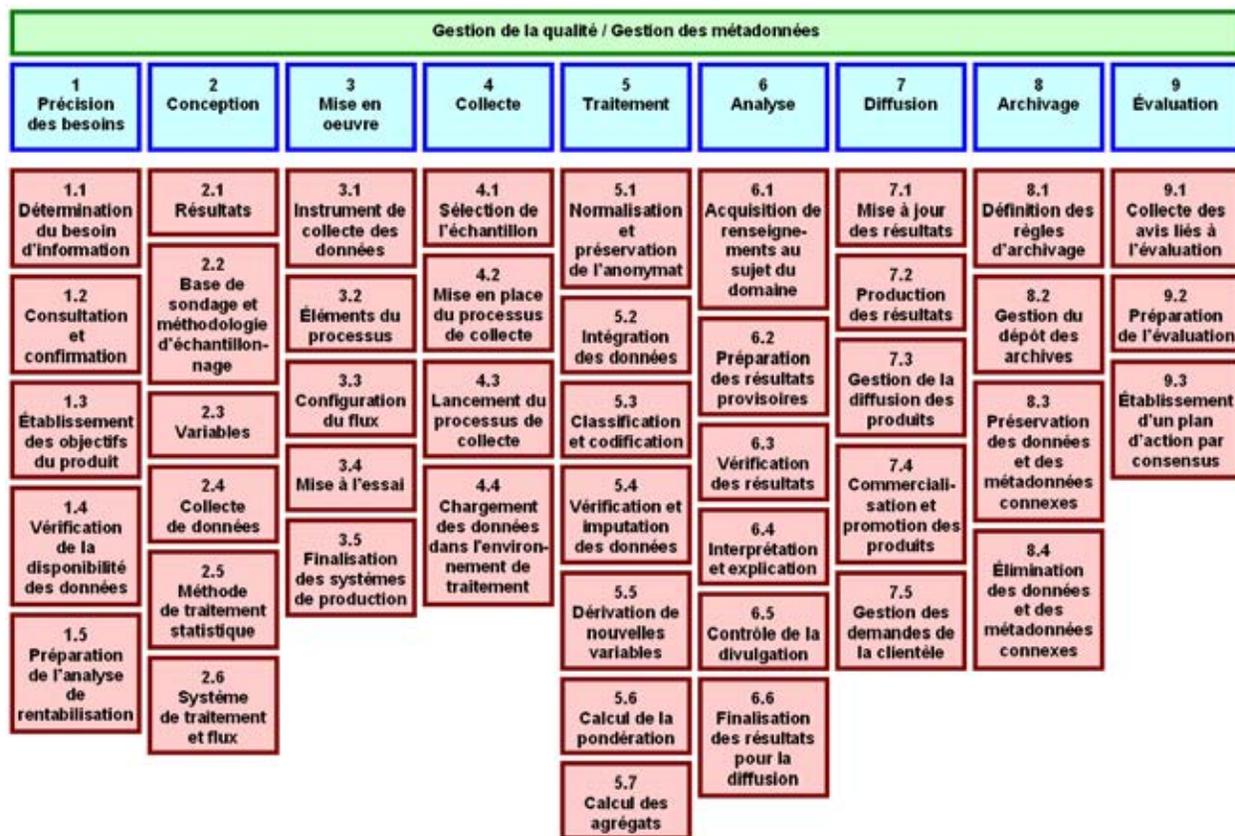
Les projets sont habituellement orientés par un comité directeur formé de cadres supérieurs qui peut réunir des gestionnaires de chacun des principaux secteurs participants. Ce comité, qui fait partie d'un mécanisme d'approbation officielle de la conception et de la mise en œuvre du programme, fournit une orientation d'ensemble ainsi que les grands paramètres relatifs au budget et à la conception. Il veille aussi à ce que les ressources appropriées soient mises à la disposition des responsables du projet et s'emploie à résoudre les problèmes qui ne peuvent pas être réglés de façon satisfaisante au sein de l'équipe de projet.

### **La qualité doit guider chaque étape du processus**

Les organismes statistiques ont souvent modélisé le processus statistique afin de le faciliter. La deuxième édition (avril 1987) des Lignes directrices concernant la qualité de Statistique Canada contenait un diagramme schématique du processus d'enquête statistique. Plus récemment, des organismes comme Statistics New Zealand, l'Australian Bureau of Statistics, Statistique Suède, Statistique Norvège et Statistique Pays-Bas, ainsi que les participants au groupe de travail conjoint de la CEE-ONU, d'Eurostat et de l'OCDE (METIS), ont élaboré différentes versions d'un modèle statistique général du processus opérationnel, le Generic Statistical Business Process Model. Le modèle mis au point par le groupe METIS (Secrétariat de la CEE-ONU, 2008) est représenté à la figure 1. Il est basé sur le modèle de Statistics New Zealand, mais plusieurs autres organismes, y compris Statistique Canada, y ont également contribué.

Les divers modèles proposent tous la division du processus en un certain nombre de phases ou d'étapes. Bien que le contenu de chacun varie, ces modèles comportent des éléments communs : la précision des besoins des utilisateurs, la conception du programme, une phase de mise en œuvre ou d'intégration (caractéristiques, systèmes, guides des opérations, formation, etc.), une phase d'exécution (collecte, vérification, etc.) et une phase d'évaluation. Parmi les principes de base à suivre figure la nécessité de considérer la qualité à toutes les étapes. Si les besoins de l'utilisateur sont mal établis ou incompris, toutes les étapes subséquentes aboutiront à des données qui ne sont pas pertinentes. Des mesures doivent également être prises pour assurer le bon déroulement de la phase de conception; même parfaites, les phases de mise en œuvre et d'exécution ne pourraient compenser cette lacune. Cependant, une bonne conception ne suffit pas, car le processus doit aussi être mis en œuvre de façon adéquate. De plus, sans une évaluation appropriée, l'organisme statistique ignorera si le programme statistique a atteint ses objectifs.

Figure 1 Niveaux 1 et 2 du Modèle statistique général du processus opérationnel



Le principe selon lequel la qualité doit être intégrée à chaque étape, combiné à la notion voulant que la qualité soit de nature multidimensionnelle, mène logiquement à la conceptualisation du processus de gestion de l'assurance de la qualité sous la forme d'un modèle défini par les aspects liés, d'une part, à la qualité en soi (pertinence, exactitude, etc.) et, d'autre part, aux phases de l'enquête (précision des besoins, conception, mise en œuvre, exécution et évaluation). Une approche de gestion de la qualité globale exige de prendre en compte tous les éléments de ce modèle.

### Les mesures de l'assurance de la qualité doivent être adaptées au programme

À Statistique Canada, les gestionnaires de programmes statistiques sont responsables de la mise en œuvre de leurs programmes. Chaque gestionnaire d'activités statistiques doit veiller à ce que les méthodes et procédures du programme statistique témoignent du souci de la qualité du Bureau. Il a toujours été clair que les Lignes directrices concernant la qualité s'inscrivaient dans cette approche : elles visent à aider les gestionnaires de programmes, et non à imposer des règles. On ne s'attend pas à ce que chaque programme soit conforme à la moindre ligne directrice; cela serait à la fois futile et d'un coût exorbitant, compte tenu des écarts entre les différents programmes statistiques quant à leur importance. On s'attend plutôt à ce que les gestionnaires de programmes statistiques et les équipes de projet qui les appuient prennent les décisions nécessaires.

## **Il faut informer les utilisateurs de la qualité des données afin qu'ils puissent mesurer si l'information statistique est adéquate à leur propre utilisation**

Enfin, pour que l'utilisateur puisse recourir de façon éclairée à l'information statistique offerte, il doit être capable d'évaluer si les données sont d'une qualité suffisante. Il peut le faire lui-même pour certains aspects de la qualité, comme l'actualité des données. Toutefois, pour d'autres aspects, comme l'intelligibilité, la cohérence et même la pertinence, cette évaluation n'est pas si simple. Mentionnons plus particulièrement l'aspect de l'exactitude, pour lequel les utilisateurs ne disposent souvent d'aucun moyen d'évaluation, devant alors se fier à l'organisme statistique.

Au fil des ans, le Bureau s'est doté de politiques et d'outils pour aider les utilisateurs. Il a longtemps appliqué la Politique visant à informer les utilisateurs sur la qualité des données et la méthodologie (PIUQDM), qui contient, à l'intention des utilisateurs, certains renseignements de base sur la qualité des données et la méthodologie. Diffusées au moyen du Quotidien, toutes les données sont accompagnées d'un lien vers la Base de métadonnées intégrées (BMDI), qui fournit de l'information sur les concepts, les définitions, les sources de données et la méthodologie utilisés pour chacun des programmes statistiques.

En ce qui concerne l'exactitude, la Politique précise également que toutes les données diffusées doivent être accompagnées de renseignements sur trois sources d'erreurs courantes – la couverture (la différence entre la population cible et l'échantillon utilisé pour mener l'enquête), la non-réponse (la proportion de l'échantillon n'ayant pas répondu) et l'erreur d'échantillonnage (une erreur imputable à l'échantillon et non à l'enquête en soi) – ainsi que sur d'autres sources, également importantes (p. ex. les erreurs touchant les réponses, les erreurs de traitement, les erreurs introduites à l'étape du contrôle de la divulgation).

## **La qualité doit demeurer au premier plan de toutes les activités**

À moins que des mesures proactives ne soient prises, la qualité des données se détériore avec le temps. Par exemple, un « manque de pertinence » se produit souvent en raison du délai entre l'émergence d'un besoin de données et la capacité d'y répondre. Des lacunes touchant d'autres aspects de la qualité peuvent aussi apparaître alors que les taux de réponse diminuent sous l'effet de changements d'attitude dans la société, ou lorsque les systèmes deviennent désuets ou qu'une restructuration des méthodes s'impose.

L'organisme statistique doit veiller à ce que la qualité demeure au premier plan de toutes ses activités pour diminuer le risque de manque de pertinence et empêcher une dégradation de la qualité au fil des ans. Il faut constamment évaluer, vérifier ou examiner la qualité des processus et des résultats des programmes. De tels mécanismes d'examen doivent être intégrés aux processus de planification et de prise de décisions du Bureau. Entre autres exemples de mécanismes d'examen de la qualité de la sorte, mentionnons le Conseil national de la statistique et les arrangements bilatéraux supérieurs avec des ministères et organismes fédéraux clés – pour engager un dialogue efficace avec les intervenants –, la présentation de rapports intégrés sur les programmes – pour assurer l'évaluation systématique des résultats statistiques – et la mise en œuvre d'un programme d'examen de la qualité – pour procéder à l'évaluation officielle des processus statistiques.

## **But et portée des lignes directrices**

La section 2 de ce document réunit des lignes directrices et des listes de contrôle liées à de nombreuses questions dont on doit tenir compte dans la poursuite des objectifs de qualité que sous-tend l'exécution des activités statistiques. Le document s'attarde principalement à la façon d'assurer la qualité grâce à la conception ou à la mise en œuvre efficace et adéquate d'un projet ou d'un programme statistique, des débuts jusqu'à l'évaluation, la diffusion et la documentation des données. Ces lignes directrices sont fondées sur les connaissances et l'expérience collective d'un grand nombre d'employés de Statistique Canada. On espère que les Lignes directrices concernant la qualité seront utiles au personnel chargé de la planification et de la conception des enquêtes et d'autres programmes statistiques, ainsi qu'à ceux qui évaluent et analysent les résultats de ces programmes.

Le principal objectif des Lignes directrices concernant la qualité consiste à fournir une liste exhaustive de principes directeurs et de pratiques exemplaires à appliquer lors de la conception d'enquêtes. Pour mieux apprécier la portée de ces lignes directrices, il importe de définir le sens donné aux termes « enquête » et « conception ».

Le terme « enquête » est un générique utilisé pour désigner toutes les activités visant la collecte ou l'acquisition de données statistiques. Il englobe :

- le recensement, par lequel on tente de recueillir des données sur tous les membres de la population;
- l'enquête par sondage, dans laquelle on recueille des données sur un échantillon (habituellement aléatoire) des membres de la population;
- la collecte des données provenant des dossiers administratifs, où les données sont tirées des documents initialement conservés à des fins non statistiques;
- les activités statistiques dérivées, qui entraînent l'estimation et la modélisation de données qui peuvent être tirées des sources de données statistiques existantes.

Les lignes directrices concernent principalement les recensements et les enquêtes par sondage. Il est manifeste que la qualité des activités statistiques dérivées est grandement déterminée par la qualité des parties des composantes. Par conséquent, le présent document ne s'attarde pas directement à ces activités.

Le terme « conception » s'applique à la délimitation de tous les aspects d'une enquête, de la détermination d'un besoin de données jusqu'à la production des résultats finaux (le fichier de microdonnées, les séries statistiques et l'analyse).

Le corps du présent document (section 2) traite principalement des questions de qualité liées à la conception des enquêtes. Toutefois, il importe de se rappeler que le contexte dans lequel on prépare une enquête impose des contraintes à la conception de cette enquête. Chaque nouvelle enquête, tout en visant à répondre à certains besoins immédiats en matière d'information, ajoute également de l'information à une base de données statistiques qu'on peut utiliser à des fins qui sont beaucoup plus nombreuses que celles établies au moment de la conception de l'enquête. Il est donc important de s'assurer que le résultat de chaque enquête peut, dans la mesure du possible, être intégré à des données sur des sujets connexes tirées d'autres enquêtes et utilisé en conjonction avec ces données. Cela sous-tend la nécessité d'examiner et de respecter les normes statistiques sur le contenu et les domaines qui ont été mises en place pour offrir une cohérence et une harmonie des données dans le système statistique national. Parmi ces normes, on trouve des cadres statistiques (comme le Système de comptabilité nationale), des systèmes de classification statistique (comme ceux rattachés à l'industrie ou à la géographie), ainsi que d'autres concepts et définitions qui précisent les variables statistiques à mesurer. L'utilité des nouvelles données statistiques est accrue jusqu'au point où on peut les utiliser en conjonction avec des données existantes.

## **Bibliographie**

Brackstone, G. (1999). La gestion de la qualité des données dans un bureau de statistique. *Techniques d'enquête*, 25, 157-171.

Deming, W.E. (1982) *Quality, Productivity, and Competitive Position*, Cambridge, MA: Massachusetts Institute of Technology.

Fellegi, I. (1996). Characteristics of an effective statistical system. *International Statistical Review*, 64, 165-197

Statistique Canada (1987). *Lignes directrices concernant la qualité*. Statistique Canada.

Statistique Canada (2000c). Base de métadonnées intégrée: [http://stdsweb/standards/imdb/imdb-menu\\_f.htm](http://stdsweb/standards/imdb/imdb-menu_f.htm)

Statistique Canada (2000d). Politique visant à informer les utilisateurs de la qualité des données et la méthodologie. Manuel des politiques, 2.3, Statistique Canada, Ottawa, Ontario.

Statistique Canada (2002c). Le cadre d'assurance de la qualité de Statistique Canada – 2002. Publication n° 12-586-X au catalogue, Statistique Canada.

Statistique Canada (2003d). Manuel des politiques de Statistique Canada.

Trewin, D. (2002). L'importance d'une culture de la qualité. *Techniques d'enquête*, 28, 135-145.

UNECE Secretariat (2008) « Generic Statistical Business Process Model: Version 3.1 – December 2008 », Joint UNECE/Eurostat/OECD Work Session on Statistical Metadata (METIS).

# Étapes de l'enquête

---

Cette section est divisée en sous-sections qui correspondent aux activités principales d'une enquête type. Les sous-sections présentent toutes la même structure : elles décrivent la portée et l'objet, les principes, les lignes directrices et les Indicateurs de qualité liées à chaque activité. La première sous-section traite de l'étape à laquelle on détermine les objectifs, les utilisations et les utilisateurs. Les sous-sections qui suivent décrivent les autres étapes de l'enquête dans l'ordre chronologique où elles se présentent habituellement. Cependant, il existe des interrelations importantes entre certaines étapes, comme par exemple entre la conception du questionnaire et les opérations de collecte et de saisie des données. Pour cette raison, certaines sous-sections renvoient à d'autres sous-sections. En outre, les étapes d'enquête dont il est question dans ce document ne se déroulent pas nécessairement selon une séquence stricte. Certaines activités peuvent être menées simultanément, par exemple, la constitution de la base de sondage, les plans d'échantillonnage et la conception du questionnaire. D'autres étapes, telles que l'évaluation de la qualité des données et la documentation sont rattachées à la majorité des autres activités et ne constituent pas des activités distinctes en soi.

## Portée et objet

Sous la rubrique Portée et objet, on donne une description de l'activité et on indique son impact potentiel sur la qualité. L'objectif de l'étape et la raison pour laquelle elle est importante est déclaré brièvement.

## Principes

Les principes sont les politiques, les approches et les orientations générales qui régissent la conception d'une activité donnée; on accorde une priorité aux principes qui sont rattachés à la qualité.

## Lignes directrices

Les lignes directrices sont des règles de pratique reconnues qui ont été élaborées lors de la conception et de la mise en œuvre des enquêtes statistiques. Toutes ces lignes directrices ne s'appliquent pas à chaque enquête, mais offrent des listes de contrôle pour faciliter la conception de l'enquête. Il faut toutefois faire preuve de jugement afin d'évaluer les considérations suggérées par ces lignes directrices.

D'un autre côté, Statistique Canada applique des politiques qui ont une incidence sur de nombreux aspects des activités statistiques du Bureau et qui peuvent comporter des exigences liées à la mise en œuvre d'activités particulières. Ces politiques sont décrites dans le Manuel des politiques de Statistique Canada. Lorsqu'une politique a une incidence sur un sujet particulier traité dans les lignes directrices, on indique l'existence et la pertinence de cette politique.

## Indicateurs de qualité

Les mesures de la qualité donnent une mesure directe de la qualité des données mais, en pratique, on peut rarement les calculer de manière explicite. Par exemple, dans le cas de l'exactitude, il est presque impossible de mesurer le biais de non-réponse puisqu'il peut être difficile de déterminer les caractéristiques des personnes qui ne répondent pas. Certaines informations peuvent par contre être fournies pour aider à « indiquer » la qualité. Les indicateurs de qualité se composent habituellement d'information représentant un sous-produit du processus statistique. Ils ne mesurent pas directement la qualité mais fournissent suffisamment d'information pour donner une idée de la qualité. Incluant dans ce section sont les mesures de qualité ou ils existent, ainsi que les indicateurs de qualité.

L'information présentée dans la rubrique Indicateurs de qualité sera utile aux méthodologistes chargés de produire des mesures de la qualité pour accompagner les produits statistiques. Elle intéressera également les gestionnaires d'enquêtes et les utilisateurs de données, qui se serviront des indicateurs pour évaluer et comparer la qualité des produits statistiques à leur utilisation. Par ailleurs, cette section saura intéresser les directeurs de secteurs de programme parce qu'elle servira de fondement à la surveillance du rendement ayant trait à la qualité des processus et des produits dans un secteur de programme.

# 1 Objectifs, utilisations et utilisateurs

## 1.1 Portée et objet

Les objectifs dictent la nature des informations requises, établis en fonction d'un programme, d'une problématique de recherche ou d'une hypothèse justifiant la nécessité d'avoir accès à cette information. Les utilisations viennent circonscrire et préciser, quant à elles, le type d'informations requises, en décrivant, par exemple, les décisions dont l'information recueillie peut engendrer ainsi que leur rôle dans la justification des décisions. Enfin, les utilisateurs sont les organisations, les organismes, les groupes ou les personnes qui sont censés se servir des résultats visés par le projet.

La première étape de planification d'une activité statistique consiste à en déterminer les objectifs. En les énonçant clairement, on oriente les étapes subséquentes du projet. Rien n'empêche de les réviser, autant de fois qu'il est nécessaire, pendant l'élaboration d'une enquête. Pour atteindre les objectifs, il sera parfois nécessaire de concevoir une nouvelle enquête, de remanier une enquête existante, d'utiliser des produits de données existants et des dossiers administratifs, voire une combinaison de ces moyens. Les besoins en information tels qu'énoncés dans les objectifs doivent justifier le fardeau de réponse qu'ils engendreront. En outre, il faut s'assurer que la pertinence des résultats visés par le projet pour la collectivité d'utilisateurs ciblée est clairement établie.

## 1.2 Principes

L'énonciation des objectifs devrait inclure les éléments suivants : hypothèses à vérifier, exigences particulières à l'égard des données et utilisation de celles-ci, attentes en matière de qualité des données, contraintes budgétaires et dates de livraison prévues. C'est aussi à cette étape que le concept, la définition, l'unité d'analyse et la population visée, lesquels seront abordés dans les sections suivantes, sont définis. Ainsi, les utilisateurs visés – ou même les utilisateurs potentiels – pourront déterminer si, et dans quelle mesure, les résultats visés par le projet répondent à leurs besoins.

## 1.3 Lignes directrices

### 1.3.1 Planification

- Établir les objectifs et les contraintes de l'enquête en collaboration avec les utilisateurs importants et les parties concernées. Établir et maintenir des liens avec les utilisateurs de l'information des secteurs public et privé, ainsi qu'avec le public général, afin de rehausser, d'une part, la pertinence de l'information produite et d'améliorer, d'autre part, la commercialisation des produits et des services. Parmi les utilisateurs importants, mentionnons les représentants des marchés potentiels, les décideurs et les agents, qui se servent de l'information à des fins législatives. Avant de se lancer dans des plans ou des restructurations d'envergure, envisager à mener une étude de faisabilité et/ou un essai pilote. Mener systématiquement de vastes consultations axées sur les utilisateurs, afin de déterminer ses options, sur le plan du contenu, et la pertinence de l'enquête. Cette démarche permet de statuer sur la nécessité d'une enquête transversale ou longitudinale, en plus de susciter l'intérêt du public à l'égard du programme et de l'encourager à participer lors de la collecte des données.
- Orienter l'analyse des besoins des utilisateurs et des données requises sur la recherche des solutions les plus rentables à court et à long terme dans le contexte de programme statistique. À priori, demander aux utilisateurs de préciser leur plan d'analyse ou les tableaux de diffusion qu'ils proposent, afin de mieux cerner leurs besoins. Avant de concevoir une nouvelle activité statistique (ou d'en restructurer une existante, pour une enquête longitudinale comme transversale), analyser les statistiques accessibles dans le secteur en question en termes de sources, fréquence, qualité, actualité, etc. Déterminer le meilleur équilibre entre l'efficacité à répondre aux besoins des clients à partir de statistiques disponibles versus le coût et le temps requis pour entreprendre une nouvelle activité afin produire des nouvelles statistiques.

- Lorsqu'il existe des objectifs sur la qualité des données explicites, il faut les intégrer à l'énoncé des objectifs de l'enquête en indiquant les aspects mesurables de la qualité pour toute la population ou pour des domaines particuliers. On peut établir des objectifs pour mesurer, par exemple, l'erreur d'échantillonnage, les taux de couverture, les taux de réponse et l'actualité. En ce qui concerne les données administratives et les activités statistiques dérivées, la qualité des données sera directement liée à la qualité des sources de données d'entrée.
- Pendant la planification, définir les contraintes opérationnelles telles que la période de référence, les coûts, les ressources et les méthodes de collecte de données. Tenir compte également des facteurs suivants : utilisation des réponses par personne interposée, rappel des répondants et nécessité de mesurer la variation des données saisonnières. À chaque étape du processus, le projet devrait se détailler et se préciser. On formule d'abord des hypothèses générales que les estimations sont basées sur la méthodologie qu'il convient d'adopter; puis, ces dernières se précisent à mesure que progresse la planification. En outre, ces estimations doivent s'appuyer sur des données historiques, lorsque ces dernières sont accessibles, et nécessitent de constantes mises à jour, c'est-à-dire chaque fois qu'on révisé les objectifs.
- Lorsque vous tentez de déterminer la stratégie convenant le mieux aux besoins des utilisateurs, viser le juste équilibre entre ces mêmes besoins, d'une part, puis les objectifs de l'enquête et les questions liées au budget, au fardeau de réponse et à la protection des renseignements personnels, d'autre part. On devrait prévoir des approches de rechange (méthodologie, moyen et mode de collecte des données, fréquence, degré de détail géographique, etc.) pour s'assurer de parvenir à la solution optimale. Évidemment, la marge de manœuvre pour ses ajustements peut être très restreinte si l'enquête est soumise à des obligations légales. Il est parfois nécessaire de mener une enquête pilote ou une étude de faisabilité pour parvenir à cette solution.
- Prenez soin de passer en revue les activités statistiques en cours à intervalles réguliers. Les programmes statistiques doivent évoluer, s'adapter et innover pour satisfaire les exigences changeantes des utilisateurs qu'ils desservent ou encore celles de nouveaux utilisateurs. Il faut donc réexaminer le but de toute activité statistique, de même que l'énoncé de ses objectifs, au fur et à mesure qu'évoluent ou changent les besoins et les contraintes (budget, période de référence, ressources, etc.) des utilisateurs, afin de renforcer la pertinence du produit statistique qui en découle. Il est parfois souhaitable de restructurer des enquêtes existantes de fond en comble pour préserver la fiabilité de séries statistiques principales, surtout si nos sources d'information ont changé ou si la façon de les diffuser a été repensée ou restructurée.

## 1.4 Indicateurs de qualité

Principal élément de la qualité : pertinence

Procéder à la description et la classification des principaux utilisateurs des résultats visés par le projet.

Décrire les besoins des mêmes utilisateurs ainsi que l'utilisation qu'ils comptent faire des produits de données, en dressant des plans d'analyse et en préparant des tableaux de diffusion. Prendre en compte tout écart existant entre les besoins des utilisateurs et les résultats visés.

Si des changements sont apportés à un programme d'enquête dont les estimations sont ajustées à une série chronologique, évaluer l'incidence des changements en question sur la série chronologique.

### References

BLANC, M., W. RADERMACHER et T. KORMER. 2001. « Quality and users », *International Conference on Quality in Official Statistics*, Stockholm, Suède.

BRACKSTONE, G.J. 1993. « Data Relevance: Keeping Pace with User Needs », *Journal of Official Statistics*, vol. 9, n° 1, 1993, p. 49 à 56.

STATISTIQUE CANADA. 2003. *Méthodes et pratiques d'enquête*, produit n° 12-587-XIF au catalogue de Statistique Canada, Ottawa, 396 p.

## 2 Concepts, variables and classifications

### 2.1 Portée et objet

Les concepts sont des idées générales ou abstraites qui traduisent le phénomène social et/ou économique à l'étude. Autrement dit, ce sont les sujets d'enquête et d'analyse qui intéressent les utilisateurs.

Une variable repose sur deux éléments, soit une unité statistique et sa propriété. L'unité statistique est l'unité d'observation ou de mesure pour laquelle les données sont recueillies ou dérivées (p. ex. les personnes ou les ménages, dans les enquêtes sociales, les entreprises ou les établissements, dans les enquêtes-entreprises) (Statistique Canada, 2008). La propriété est une caractéristique ou un attribut de cette unité statistique. Les variables doivent être définies de façon claire et précise, sans aucune équivoque, pour satisfaire les besoins de l'analyse en prévision de laquelle on recueille les données (Statistique Canada, 2004).

La classification est le regroupement systématique des valeurs qu'une variable peut prendre, ce qui inclut les classes mutuellement exclusives et l'ensemble des valeurs. En outre, la classification présente souvent une structure hiérarchique qui permet d'agréger les données pour en faciliter l'analyse et l'interprétation. Plusieurs types de classifications peuvent servir à répertorier les données d'une variable donnée (Statistique Canada, 2004).

### 2.2 Principes

Statistique Canada veut s'assurer que l'information produite par le Bureau trace un portrait logique et cohérent de l'économie, de la société et de l'environnement du Canada, en plus d'offrir des ensembles de données qui puissent s'analyser conjointement et en conjonction avec des informations issues d'autres sources. Afin d'atteindre ces objectifs, le Bureau adopte certains cadres conceptuels, recourt à des désignations et des définitions normalisées, lorsqu'il réfère aux populations, unités statistiques, concepts, variables et classifications de ses programmes, et applique, pour finir, des méthodes de collecte et de traitement uniformes dans la production des données statistiques de toutes ses enquêtes. Sur ce plan, la marche à suivre est dictée par la Politique concernant les normes de Statistique Canada (Statistique Canada, 2004) qui respecte les normes suivantes, par ordre décroissant d'importance : les normes ministérielles, les normes recommandées et les normes particulières des programmes (Statistique Canada, 2004).

### 2.3 Lignes directrices

#### 2.3.1 Observation des normes

- Énoncer clairement les concepts et les variables du projet, puis ce à quoi ils sont censés servir. Définir les concepts, variables, classifications, unités statistiques et populations en recourant aux définitions normalisées énoncées par la Politique concernant les normes de Statistique Canada (Statistique Canada, 2004). Tenir compte, dans le choix des règles d'affectation des noms, des écarts entre la norme et l'usage. Autrement dit, réserver les titres normalisés aux éléments définis dans les normes.
- Recourir à des définitions normalisées permet de comparer les données provenant de sources diverses et d'intégrer ces mêmes données à plusieurs sources (Statistique Canada, 2004). Statistique Canada dispose de classifications types destinées aux industries, aux produits, aux programmes d'enseignement, aux professions, à la comptabilité et à la géographie (Statistique Canada 2007a (SCIAN), 2007b (SPAN), 2000 (CPE), 2006a (CNP-S), 2006b (PC) et 2007c (CGT), ainsi que pour plusieurs autres domaines pertinents pour les statistiques sociales et économiques.
- Outre les classifications de Statistique Canada, il existe des classifications internationales normalisées, qui sont produites par le Bureau de statistique des Nations Unies, le Bureau international du Travail, Eurostat et d'autres organismes régionaux et internationaux. La Division des normes de Statistique Canada a dressé une liste officielle des concordances entre certaines de ces classifications et celles du Bureau. Pour transmettre des données à des organismes internationaux, privilégier d'abord cette liste, à moins qu'elle soit exempte du terme recherché.

- Choisir des unités d'observation normalisées pour faciliter la comparaison des données. Les classifications sont généralement conçues en fonction d'unités d'observation précises. Par exemple, le Système de classification des industries de l'Amérique du Nord (SCIAN) vise principalement des établissements.
- Il faut connaître les activités statistiques dérivées ou les cadres statistiques (p. ex. le Système de comptabilité nationale) dont les définitions des concepts et des variables peuvent avoir une incidence considérable sur des activités de collecte de données précises (Statistique Canada, 1989).
- Il y a parfois plusieurs façons d'évaluer un concept. Lorsqu'on choisit des variables et des classifications, en vue de l'évaluation, il faut tenir compte de plusieurs facteurs tels que l'accessibilité des informations, le fardeau de réponse, la méthode de collecte, le contexte de réponse (contexte idéal pour poser la/question(s)), le traitement des données (plus particulièrement les techniques de révision, d'imputation et de pondération), les dossiers administratifs (présence de l'information requise), ainsi que les coûts associés à la collecte et au traitement. En fait, lorsqu'on interprète un concept, le succès de l'entreprise dépend de l'approche retenue. Par ailleurs, il faut savoir qu'une variable peut s'avérer désuète, bien que pertinente à l'origine, lorsqu'elle est soumise à de nouveaux facteurs. Il faut alors la modifier ou la remplacer. Il faut donc s'assurer de recourir à la version la plus récente de toute variable approuvée. À ce chapitre, les plus récentes mises à jour sont accessibles sur le site Web de Statistique Canada.
- En l'absence de normes officielles, étudier les concepts, variables et classifications employés dans le cadre de programmes statistiques connexes. Consulter la Division des normes, au besoin

### **2.3.2 Utilisation des classifications**

- Chaque classification doit demeurer la plus flexible possible. Afin de s'en assurer, coder les microdonnées et garder les fichiers au niveau le plus bas de la classification appropriée dans ce contexte. Il est parfois nécessaire d'agréger à un niveau supérieur pour répondre aux besoins particuliers de l'analyse ou respecter des contraintes en matière de confidentialité ou de fiabilité des données. Lorsque c'est le cas, adopter, si possible, les classes ou les agrégations du niveau supérieur dictées par la norme. Sinon, opter pour une stratégie de regroupement commune, puis documenter les écarts entre la norme et les niveaux de classifications/agrégations choisis. Recourir à des classifications qui reflètent à la fois les niveaux détaillés et agrégés. Expliquer aux utilisateurs de quelle façon ces classifications s'intègrent au niveau supérieur (c.-à-d. moins détaillé).

## **2.4 Indicateurs de qualité**

Principaux éléments de la qualité : cohérence, intelligibilité, pertinence.

Décrire les principaux concepts statistiques, notamment les mesures statistiques, la population, les variables, les unités, les domaines et la période de référence. Cette information permet aux utilisateurs de constater la pertinence du produit par rapport à leurs besoins.

Fournir des références exactes lorsqu'on adopte des concepts, des variables et des classifications normalisés.

Décrire, justifier et, dans la mesure du possible, mesurer (qualitativement, sinon quantitativement) tout écart par rapport aux normes. Cette démarche permet aux utilisateurs de mesurer la pertinence des données et augmente leur intelligibilité.

## Bibliographie

STATISTIQUE CANADA. 1989. *Guide de l'utilisateur pour le Système de comptabilité nationale du Canada*, publication n° 13-589-XPB au catalogue de Statistique Canada, Ottawa, 106 p.

STATISTIQUE CANADA. 2000. *Classification des programmes d'enseignement – CPE* (en ligne), <http://www.statcan.gc.ca/concepts/classification-fra.htm>.

STATISTIQUE CANADA. 2004. « Politique concernant les normes », *Manuel des politiques de Statistique Canada*.

STATISTIQUE CANADA. 2006a. *Classification nationale des professions pour statistiques (CNP-S) 2006*, publication n° 12-583-X au catalogue de Statistique Canada, Ottawa, 648 p.

STATISTIQUE CANADA. 2006b. *Plan comptable : situation financière et les résultats financiers des entreprises privées* (en ligne), [http://stdsweb/english/Subjects/Standard/coa-standdraft\\_f.htm](http://stdsweb/english/Subjects/Standard/coa-standdraft_f.htm).

STATISTIQUE CANADA. 2007. *Base de métadonnées intégrée* (en ligne), [http://stdsweb/standards/imdb/imdb-menu\\_f.htm](http://stdsweb/standards/imdb/imdb-menu_f.htm).

STATISTIQUE CANADA. 2007. *Système de classification des industries de l'Amérique du Nord (SCIAN) 2007* (en ligne), <http://www.statcan.gc.ca/subjects-sujets/standard-norme/naics-scian/2007/index-indexe-fra.htm>.

STATISTIQUE CANADA. 2007a. *Système de codage pour la classification des industries (SCCI) version 1.4*, publication n° 12F0074XCB au catalogue de Statistique Canada, <http://www.statcan.gc.ca/subjects-sujets/standard-norme/napcs-scpn/napcs-scpn-fra.htm>.

STATISTIQUE CANADA. 2007b. *Système de classification des produits de l'Amérique du Nord (SPAN) 2007* (en ligne), <http://www.statcan.gc.ca/subjects-sujets/standard-norme/napcs-scpn/napcs-scpn-fra.htm>.

STATISTIQUE CANADA. 2007c. *Classifications géographiques type, CGT 2006*, produits n°s 12-571-X et 12-572 au catalogue de Statistique Canada, Ottawa.

STATISTIQUE CANADA. 2007c. *Classifications géographiques type, CGT 2006* (en ligne), <http://www.statcan.gc.ca/subjects-sujets/standard-norme/sgc-cgt/geography-geographie-fra.htm>.

STATISTIQUE CANADA. 2008. *Unités statistiques normalisées* (en ligne), <http://www.statcan.gc.ca/concepts/units-unites-fra.htm>.

### **3 Couverture et bases de sondage**

#### **3.1 Portée et objet**

La population cible est l'ensemble des unités sur lesquelles on désire obtenir des renseignements et pour lesquelles il faut produire des estimations. Il se peut que certains facteurs d'ordre pratique obligent le chercheur à définir la population observée en excluant des unités de la population cible; la population observée peut aussi comporter des unités définies différemment de manière à permettre l'accès à la population cible.

La base de sondage est une liste, un document ou un dispositif qui délimite et identifie les éléments de la population observée, en plus d'y donner accès. De façon générale, les bases de sondage entrent dans deux grandes catégories : les bases aréolaires et les bases listes. La base liste est composée d'une liste d'unités faisant partie de la population observée. La base aréolaire est habituellement composée d'une hiérarchie d'unités géographiques parmi lesquelles se trouvent des unités de la population observée; en d'autres termes, il est possible de subdiviser les unités qui se situent à un niveau de cette base pour former les unités du niveau suivant. Lorsqu'ils sont conjugués, les éléments de la base aréolaire forment la population de la base de sondage. Dans bien des cas, les bases de sondage sont bien plus que de simples listes d'unités ou que des cartes indiquant les limites des unités géographiques. Le plus souvent, elles fournissent d'autres renseignements (p. ex., identification, personnes-ressources, classification, adresse, taille, cartes relatives aux unités géographiques) qui serviront aux fins de l'enquête.

La couverture correspond au degré d'exhaustivité des renseignements relatifs à la population cible. Il serait possible de dériver ces renseignements si l'enquête portait sur l'ensemble des unités de la base de sondage. Les erreurs de couverture sont des divergences entre les statistiques sur la population cible et celles sur la population de la base de sondage. Ces erreurs sont attribuables au sous-dénombrement (surdénombrement) de la population cible dans la base de sondage, ainsi qu'aux erreurs survenant dans le cadre des opérations de l'enquête. Les erreurs de couverture se traduisent par des écarts entre les estimations des éléments réellement couverts par l'enquête et ceux que l'on prévoyait estimer. Elles touchent autant la dimension spatiale que temporelle.

#### **3.2 Principes**

Il doit exister une concordance raisonnable entre la population observée et la population cible pour que les résultats de l'enquête soient pertinents

De même, la base de sondage devrait se conformer à la population observée. Toute erreur de couverture dans la base de sondage (p. ex., omissions d'unités visées par l'enquête, inclusion d'unités hors champ, erreurs de classification des unités et enregistrement en double de certaines unités) risque de compliquer le processus d'enquête et d'entraîner, par le fait même, une hausse de ses coûts, ainsi qu'une diminution de l'actualité et de l'exactitude de ses estimations (à cause de biais et de variances).

Les données de la base de sondage doivent être exactes et à jour, car elles servent à la stratification, à la sélection des échantillons, à la collecte, au suivi, au traitement des données, à l'imputation, aux estimations, au couplage des enregistrements, à l'évaluation de la qualité et aux analyses. Si ces données présentent des erreurs, cela entraînera probablement un biais ou une diminution de la fiabilité des estimations de l'enquête ainsi qu'une augmentation des coûts de collecte des données.

La conception et les opérations des enquêtes doivent comporter des procédures visant à réduire le plus possible les erreurs de couverture et leurs répercussions.

### 3.3 Lignes directrices

#### 3.3.1 Conception

- À l'étape de la planification de l'enquête, mettre des bases de sondage potentielles à l'essai, afin d'en évaluer la pertinence et la qualité. Évaluer la couverture de la base de sondage et des unités de collecte cibles.
- Si aucune base de sondage ne peut assurer, à elle seule, la couverture de la population cible telle que requise, se tourner vers une approche méthodologique qui repose sur des bases de sondage multiples (combinaison de deux bases de sondage ou plus). Il peut s'agir, par exemple, de jumeler une base liste et une base aréolaire ou encore deux bases listes ou plus. Avant de recourir à cette approche, il faut toutefois s'assurer qu'aucune base de sondage ne peut assurer, à elle seule, une couverture suffisante. Par ailleurs, il est permis d'y recourir lorsqu'on ne peut se servir des bases de sondages existantes, soit parce qu'elles sont incomplètes, bien qu'abordables sur le plan des dépenses, ou, au contraire, qu'elles sont trop onéreuses bien qu'exhaustives.
- Lors des enquêtes téléphoniques, envisager la méthode de composition aléatoire (CA), qu'elle soit employée seule ou en conjonction avec d'autres bases aréolaires ou bases listes.
- Il arrive parfois qu'on ne dispose d'aucune base de sondage qui réponde à la fois aux exigences économique et statistique de l'enquête, c'est-à-dire qui permette de traiter la population d'unités auxquelles s'intéresse l'enquête sans engendrer des coûts trop élevés. Il faut alors envisager de recourir à des méthodes d'échantillonnage à plusieurs degrés ou d'échantillonnage indirect.
- Plusieurs listes sont tenues à jour à Statistique Canada afin de pouvoir servir de bases de sondage pour les enquêtes de l'organisme. Le Registre des entreprises peut servir aux enquêtes-entreprises et aux enquêtes institutionnelles. Pour les enquêtes agricoles, le Registre des fermes constitue la base de sondage habituelle. Dans le cas des enquêtes-ménages, le Registre des adresses, la base de sondage de l'Enquête sur la population active (base aréolaire) et les unités géographiques du Recensement de la population sont des options à envisager. Si, dans une situation donnée, aucune de ces bases de sondage ne constitue le meilleur choix par rapport à la population cible d'une enquête, on se penchera sur d'autres bases de sondage possibles (p. ex., listes d'immigrants, bases de données sur les importateurs ou les exportateurs).
- S'assurer que la base de sondage correspond le plus possible, du point de vue de ses mises à jour, à la période de référence de l'enquête.
- Conserver et stocker l'information sur l'échantillonnage, la rotation et la collecte des données afin que l'on puisse coordonner les enquêtes et mieux gérer les relations avec les répondants ainsi que le fardeau de réponse. Indiquer, par exemple, la fréquence à laquelle chaque unité est sélectionnée par chacune des enquêtes utilisant une même base de sondage.
- Pour les activités statistiques rattachées à des sources administratives ou pour les activités statistiques dérivées, où les changements de couvertures peuvent échapper au contrôle du gestionnaire immédiat, déterminer et surveiller la couverture, et négocier les changements requis avec le gestionnaire de la source.
- Pour pallier l'erreur de couverture d'une base de sondage, ajuster les données ou se servir de données complémentaires tirées d'autres sources.
- Lorsque des enquêtes portent sur une même population cible, utiliser la même base de sondage, si possible, afin d'augmenter la cohérence, d'éviter les contradictions, de faciliter la combinaison des estimations tirées des enquêtes et de réduire les coûts relatifs à la mise à jour et à l'évaluation des bases de sondage.
- Comme le veut la pratique courante, à Statistique Canada, utiliser le Système automatisé de regroupement des territoires (SARTE) lorsque vient le temps de créer des sous-ensembles géographiques dans la base de sondage; il s'agit d'un système automatisé de délimitation et de vérification spatiales partiellement générique servant à créer des unités géographiques contigus et non chevauchants dans une base de sondage.

- Lorsqu'il existe des bases de sondage multiples, il est possible de se servir de celles-ci pour évaluer l'exhaustivité d'une base de sondage.
- Implanter des procédures afin de repérer et de corriger d'éventuelles erreurs de couverture dans la base de sondage. Fournir de la rétroaction pour mettre et maintenir à jour la base de sondage.
- Établir des procédures et développer de la formation à l'intention du personnel chargé de la collecte et du traitement des données, afin de réduire les erreurs de couverture au minimum (p. ex., une procédure permettant d'obtenir la confirmation exacte des listes de logements pour les unités de la base aréolaire composant l'échantillon).
- Élaborer les questionnaires d'enquête et la documentation connexe de manière à réduire au minimum les erreurs de couverture imputables aux répondants (p. ex., mentions erronées de personnes relativement à un logement dans un questionnaire, omission d'endroits faisant partie du champ d'une enquête axée sur les établissements).

### **3.3.2 Mise à jour**

- Afin de rehausser et/ou de maintenir le niveau de qualité de la base de sondage, établir des procédures permettant d'éliminer les enregistrements en double et de faire les mises à jour concernant les naissances, les décès, les unités hors du champ de l'enquête et les changements de caractéristiques.
- Intégrer ces mises à jour dans la base de sondage le plus tôt possible.
- Réduire au minimum les erreurs dans la base de sondage en misant sur la formation adéquate du personnel, en insistant sur l'importance de la couverture et en implantant des procédures d'assurance de la qualité pour tout ce qui touche les activités associées à la base de sondage.
- En ce qui concerne les bases aréolaires, implanter une procédure de vérification des cartes, afin que la délimitation des régions géographiques qui servent à développer le plan d'échantillonnage soit évidente et exemptes de tout chevauchement (p. ex., vérifications sur le terrain ou comparaison avec des cartes provenant d'autres sources). Si cela convient, recourir au Registre des adresses pour vérifier le listage des adresses résidentielles effectué sur le terrain, ce qui permettra de repérer les sous-dénombrements ou les surdénombrements.
- Examiner et améliorer la procédure d'identification des unités cibles qui ont été omises ou mal codées, et établir une procédure pour atténuer ce problème.
- Implanter des procédures permettant de repérer et de réduire au minimum les omissions et les classifications erronées – qui peuvent engendrer le sous-dénombrement – de même que les inclusions erronées et les enregistrements en double – qui peuvent engendrer le surdénombrement.

### **3.3.3 Documentation**

- Dans la documentation relative à l'enquête, donner la définition de la population cible et de la population observée, puis noter d'éventuelles différences entre elles. Décrire la base de sondage et faire état des erreurs de couverture qu'elle comporte.
- Faire état des écarts répertoriés entre la couverture de l'enquête et les principaux besoins des utilisateurs.

## **3.4 Indicateurs de qualité**

Principaux éléments de la qualité : exactitude, pertinence.

Si l'on fait abstraction des différences entre la population cible et la population observée, les erreurs de couverture sont attribuables soit au sous-dénombrement, soit au surdénombrement. Le sous-dénombrement survient lorsqu'on omet, par erreur, d'inclure certaines unités dans le fichier de la base de sondage; le surdénombrement survient lorsque des unités sont incluses par erreur dans ce même fichier (p. ex., des décès). Les erreurs de classification – par exemple au niveau d'un secteur industriel – donnent lieu à une

erreur de couverture de la base de sondage, du fait d'un sous-dénombrement au niveau de la classification « correcte » et d'un surdénombrement imputable à une classification erronée. Plusieurs processus – sélection des échantillons, traitement des données, imputation, estimation, couplage d'enregistrements, évaluation de la qualité et analyse – s'appuient sur les données de classification. La collecte des données et le suivi peuvent tirer profit, quant à eux, des coordonnées de personnes-ressources. Les imperfections de la base de sondage, comme les erreurs de couverture et les caractéristiques désuètes, sont susceptibles de biaiser les estimations de l'enquête ou d'en réduire la fiabilité, ainsi que de faire augmenter les coûts de collecte des données.

Les erreurs de couverture (sous-dénombrements/surdénombrements) peuvent miner les résultats d'une enquête et remettre en cause leur pertinence et leur exactitude.

- Contrôler la qualité de la base de sondage en évaluant périodiquement sa couverture ainsi que la qualité des informations dont on dispose sur les caractéristiques des unités. Il existe plusieurs techniques, pour ce faire :
  - apparier la base de sondage ou un échantillon de celle-ci avec des sources comparables – on trouve souvent ce genre de données dans les dossiers administratifs – pour évaluer la population observée ou des sous-ensembles de cette dernière;
  - analyser les résultats de l'enquête pour identifier les enregistrements en double, les décès, les unités hors champ et les changements relatifs aux caractéristiques;
  - utiliser des questions précises du questionnaire pour faciliter le suivi de l'information sur la couverture et la classification. Vérifier auprès des autorités locales (p. ex., les bureaux régionaux, le personnel des enquêtes sur le terrain, les unités d'enquête elles-mêmes);
  - vérifier la base de sondage ou certains de ses sous-ensembles sur le terrain (ce qui peut inclure la vérification d'unités hors champ);
  - comparer la base de sondage avec un échantillon d'unités issues d'une base aréolaire correspondante;
  - faire la mise à jour de la base de sondage afin de mettre en lumière les changements qu'elle aurait subis;
  - vérifier la cohérence des chiffres par rapport à d'autres sources ou par rapport à des données issues de répétitions conçues spécialement à cette fin;
  - recourir à des informations tirées d'autres enquêtes qui emploient la même base de sondage aux fins de l'évaluation (Lessler et Kalsbeek, 1992).
- Contrôler la base de sondage entre le moment où l'échantillon est sélectionné et celui auquel la période de référence de l'enquête correspond.
- Définir et comparer la population cible et la population observée.
- Dans le cas de recensements, mesurer l'erreur de couverture au moyen d'enquêtes postcensitaires (Hogan, 2003) ou d'enquêtes de la contre-vérification des dossiers (Statistique Canada, 2004) ainsi que d'études connexes portant sur le surdénombrement. Dans le cas de recensements de la population et des logements, on peut aussi évaluer les erreurs de couverture en comparant chiffres de dénombrement aux estimations démographiques. Fournir une estimation non seulement de l'erreur nette, mais également de ses composantes.
- Fournir des estimations de l'erreur de couverture ou du taux de glissement autant pour les utilisateurs actuels des estimations de cette enquête que pour ceux qui concevront ses prochaines éditions.

## Bibliographie

- ARCHER, D. 1995. « Maintenance of Business Registers », *Business Survey Methods*, B.G. Cox et coll., New York, Wiley-Interscience, p. 85 à 100.
- HARTLEY, H.O. 1962. « Multiple Frame Surveys », *Proceedings of the Social Statistics Section*, American Statistical Association, p. 203 à 206.
- HOGAN, H. 2003. « The Accuracy and Coverage Evaluation: Theory and Design », *Techniques d'enquête*, vol. 29, n° 2, p. 129 à 138.
- KOTT, P.S. et F.A. VOGEL. 1995. « Multiple-Frame Business Surveys », *Business Survey Methods*, B.G. Cox et coll., New York, Wiley-Interscience, p. 185 à 203.
- LANIEL, N. et H. FINLAY. 1991. « Data Quality Concerns with Sub-annual Business Survey Frames », *Proceedings of the Section on Survey Research Methods*, American Statistical Association, p. 202 à 207.
- LESSLER, J.T. et W.D. KALSBECK. 1992. *Nonsampling Errors in Surveys*, New York, Wiley, 432 p.
- MASSEY, J.T. 1988. « An Overview of Telephone Coverage », *Telephone Survey Methodology*, R.M. Groves et coll., New York, Wiley, p. 3 à 8.
- STATISTIQUE CANADA. 2008. *Méthodologie de l'Enquête sur la population active du Canada*, publication no71-526XIF au catalogue de Statistique Canada, Ottawa, 116 p.
- STATISTIQUE CANADA. 2003. *Méthodes et pratiques d'enquête*, publication n° 12-587-XIF au catalogue de Statistique Canada, Ottawa, 396 p.
- STATISTIQUE CANADA. 2004. *Couverture : rapport technique du Recensement de 2001*, publication n° 92-394-XIF au catalogue de Statistique Canada, Ottawa, 85 p.
- SWAIN, L., J.D. DREW, B. LAFRANCE et K. LANCE. 1992. « La création d'un registre des adresses résidentielles pour améliorer la couverture du Recensement du Canada de 1991 », *Techniques d'enquête*, vol. 18, n° 1, p. 139 à 155.

## 4 Plan d'échantillonnage

### 4.1 Portée et objet

L'échantillonnage est un moyen de sélectionner un sous-ensemble d'unités d'une population cible dans le but de recueillir des renseignements. Ces renseignements sont utilisés pour tirer des conclusions au sujet de la population en général. Le sous-ensemble d'unités sélectionnées à l'échantillonnage est appelé échantillon. Le plan d'échantillonnage englobe tout ce qui concerne la manière de regrouper les unités dans la base, de déterminer la taille de l'échantillon, de répartir l'échantillon dans les diverses classifications des sous-sections de la base de sondage et de sélectionner l'échantillon. Les choix relatifs au plan d'échantillonnage sont influencés par de nombreux facteurs comme le degré de précision et de détail visé pour les informations à livrer, l'existence de bases de sondage appropriées, la disponibilité de variables auxiliaires permettant la stratification et la sélection de l'échantillon, les méthodes d'estimation qui seront appliquées et le budget alloué, du point de vue du « temps » et des « ressources ».

### 4.2 Principes

Il existe deux types d'échantillonnage : l'échantillonnage non probabiliste et l'échantillonnage probabiliste. L'échantillonnage non probabiliste repose sur la sélection subjective d'unités au sein d'une population. Il est généralement rapide, simple et abordable. Étant donné ses caractéristiques, ce type d'échantillonnage est parfois utile pour mener des études préliminaires, tenir des groupes de discussion et faire des études de suivi. Or, pour pouvoir tirer des conclusions sur la population entière, il n'en demeure pas moins qu'il faut supposer, souvent à tort, que l'échantillon est représentatif. L'échantillonnage probabiliste repose quant à lui sur trois principes généraux qui tracent les limites de son cadre statistique. Le premier principe est la randomisation, soit la sélection aléatoire des unités de l'échantillon. Selon le second principe, toutes les unités de la population observée ont une probabilité positive connue d'être sélectionnées dans l'échantillon. Le troisième est le calcul de cette probabilité, qui permet ensuite d'établir des estimations générales et des estimations de l'erreur d'échantillonnage. L'échantillonnage probabiliste reste le meilleur choix, pour la plupart des programmes statistiques, car il permet de tirer des conclusions fiables sur l'ensemble de la population et de quantifier l'erreur dans les estimations.

Le plan d'échantillonnage devrait être aussi simple que possible. Il a pour objectif de livrer des estimations exactes et suffisamment précises pour répondre aux exigences de l'enquête. La précision d'une estimation est mesurée selon sa variance. Le manque d'exactitude est révélé par les biais, qui sont souvent attribuables à des facteurs indépendants de l'échantillonnage, comme les erreurs de déclaration et de mesure, l'inexactitude du traitement, ainsi que les erreurs liées à la non-réponse et aux déclarations incomplètes.

### 4.3 Lignes directrices

#### 4.3.1 Plan

- Pour déterminer la taille d'un échantillon, il faut tenir compte des niveaux de précision nécessaires à la production des estimations de l'enquête, du type de plan (p. ex., échantillonnage en grappes, stratification) et d'estimateur utilisés, de l'accessibilité des informations auxiliaires et des coordonnées des personnes-ressources, des contraintes budgétaires, ainsi que de certains facteurs, comme la non-réponse, la présence d'unités hors champ, l'attrition dans les enquêtes longitudinales, etc. Pour les enquêtes périodiques, il faut tenir compte des additions et des suppressions d'unités prévisibles dans la population observée, qui est en constante évolution. Il importe de souligner que la précision des estimations d'une enquête tient généralement davantage à la taille de l'échantillon total qu'au taux d'échantillonnage (ratio de la taille de l'échantillon par rapport à la taille de la population).
- Il faut se rappeler que la plupart des enquêtes génèrent des estimations pour plusieurs variables et que le fait d'optimiser l'échantillon pour une variable en particulier peut avoir des effets négatifs sur d'autres variables importantes. Il faut gérer ce problème en déterminant d'abord les variables les plus importantes; on obtient alors un sous-ensemble de variables permettant de déterminer quelle stratégie adopter pour l'échantillonnage. Cette stratégie sous-tend souvent un compromis entre les stratégies optimales s'appliquant à chacune des variables de ce sous-ensemble. Consulter Bethel (1989).

- La stratification consiste à diviser la population en sous-ensembles qui sont appelés strates. Chaque strate fournit un échantillon indépendant. Le choix des strates est dicté par les objectifs de l'enquête, la disponibilité des variables de la base de sondage, la distribution de la variable d'intérêt et le niveau de précision visé pour les estimations. La majorité des enquêtes produisent des estimations sur divers domaines d'intérêt (p. ex., les provinces). Il faut en tenir compte dans le plan de l'enquête – si possible – en stratifiant la population de manière appropriée (p. ex., par province). Si tel n'est pas le cas, il faudra envisager de recourir à des méthodes spéciales, à l'étape de l'estimation, pour produire les estimations de ces domaines (voir Imputation). Afin d'être efficace sur le plan statistique, il importe de s'assurer que chacune des strates contienne des unités aussi homogènes que possible par rapport aux informations recueillies par l'enquête. Pour les enquêtes longitudinales, choisir des variables de stratification qui correspondent à des caractéristiques reconnues pour leur stabilité au fil du temps.
- Mener des études pour évaluer plusieurs options par rapport aux méthodes d'échantillonnage, à la stratification et à la répartition. L'utilité de ces études dépend de la disponibilité et de l'actualité des données qui les alimentent – qu'il s'agisse de données administratives, d'enquêtes ou de recensements antérieurs – et de la relation qu'entretiennent ces dernières avec les variables importantes de l'enquête. Consulter Kish (1988).
- Déterminer le taux de réponse attendu au moyen d'un prétest ou de données tirées d'éditions précédentes de la même enquête ou d'enquêtes similaires. Ce taux peut servir à déterminer la taille de l'échantillon. L'échantillon peut être divisé en vagues successives qui seront relâchées au besoin selon la taille de l'échantillon obtenu par strate. Pour les enquêtes longitudinales, il faut utiliser l'attrition cumulée prévue pour un nombre de cycles donné.

#### 4.3.2 Méthodes

- Pour les populations qui sont hautement asymétriques, il faut créer une strate de grandes unités dont l'inclusion dans l'enquête est certaine (la strate à tirage complet). En général, ces grandes unités représenteront une part substantielle des totaux de population. Afin de réduire le fardeau du répondant, il se peut qu'il faille créer une strate de très petites unités à exclusion de la population observée. Consulter Baillargeon et coll. (2007). Il importe de bien distinguer la portion non sondée de la population observée (strate à tirage nul), qui appartient à la population observée sans toutefois faire partie de l'échantillon, et les unités hors du champ de l'enquête, qui n'appartiennent pas à la population observée. La contribution de la strate à tirage nul peut être estimée au moyen de modèles.
- Il arrive que l'information nécessaire à la stratification de la population ne soit pas accessible dans la base de sondage. Lorsque c'est le cas, on peut se servir d'un plan d'échantillonnage à deux phases (ou double), qui sélectionne un grand échantillon lors de sa première phase, pour obtenir les informations nécessaires à la stratification. Ce premier échantillon est ensuite stratifié. Pendant la seconde phase, on retient un sous-échantillon pour chacune des strates du premier échantillon. Il importe de se questionner sur le coût de l'échantillonnage pour chaque phase, sur la disponibilité de l'information requise par chaque phase et sur les gains associés à la stratification de l'échantillon de première phase, pour ce qui est de la précision.
- Dans la pratique, il n'est pas toujours possible de sélectionner directement les unités qui fourniront les renseignements nécessaires ou de communiquer directement avec elles. Il arrive que la démarche ne soit pas rentable ou qu'on ne dispose pas d'informations suffisantes pour la mener à bien. En pareil cas, on peut se servir d'un plan d'échantillonnage à deux degrés : on sélectionne d'abord des grappes (appelées les unités primaires d'échantillonnage) d'unités déclarantes, puis un échantillon d'unités déclarantes à l'intérieur de chaque grappe sélectionnée. Il est possible que des contraintes budgétaires ou d'une autre nature nécessitent plus de deux degrés (un plan à plusieurs degrés). Pour l'échantillonnage, déterminer le nombre d'étapes nécessaires et le type d'unité approprié, et ce, à chaque étape. Pour chaque type d'unité, vérifier les éléments suivants : disponibilité d'une base d'unités adéquate – ou possibilité d'en créer une – à chaque étape, facilité de la prise de contact et de la collecte/mesure des données, qualité des données fournies par l'unité et coût de la collecte. Les plans à plusieurs degrés sont, par définition, des plans d'échantillonnage par grappes. Bien qu'elles réduisent le coût de la collecte de données, les grappes peuvent accroître les variances attribuables à la corrélation intragrappe.

- Si les échantillons sont sélectionnés dans diverses bases (deux ou plus), il faut se montrer prudent avec les unités appartenant à plus d'une base. Il faut savoir à quelle base chacune de ces unités appartient. Par ailleurs, il faut privilégier des plans d'échantillonnage qui simplifient les procédures d'estimation. Le principe selon lequel le plan doit rester simple est d'autant plus vrai lorsqu'on a recourt à de multiples bases.
- Lorsque la répartition et la taille des échantillons stratifiés sont déterminées, il faut tenir compte des taux de classification erronée prévus dans les unités et de toute autre lacune de la base, sans quoi les estimations de l'enquête seront moins précises que prévu. Il faudra donc s'attaquer à ce problème à l'étape de l'estimation (voir la section 2.10).
- Certains plans d'échantillonnage complexes demandent de calculer un effet de plan de sondage (EPS) pour déterminer la taille de l'échantillon. Pour calculer l'EPS d'une enquête, il faut se servir de résultats d'enquêtes antérieures ou d'enquêtes similaires. Consulter Gambino (2001), Kish (1965) et Gabler et coll. (2006).
- Certaines situations plus complexes – étude de populations rares ou mobiles, échantillonnage à partir d'une liste d'unités qui sont liées aux unités de la population visée, sans pour autant y correspondre directement – peuvent requérir un plan particulier. Il peut s'agir de recourir à certaines techniques comme l'échantillonnage indirect, le sondage par réseaux ou l'échantillonnage par grappes adapté, pour ne nommer que celles-là. Consulter Lavallée (2007) et Thompson et Seber (1996).
- La méthode de composition aléatoire (CA) est très populaire dans certains types d'enquêtes-ménage. Les plans qui recourent à la CA comportent des risques de biais, car ce ne sont pas tous les ménages qui ont des téléphones conventionnels (à fil). Étant donné la prévalence croissante du nombre de ménages ayant uniquement des téléphones cellulaires, le problème va même s'aggraver, sauf si la CA est appliquée aux numéros de téléphone cellulaire. Avant de choisir la méthode du CA pour une enquête, il faut donc mesurer avec soin les risques de biais.

#### 4.3.3 Enquêtes périodiques

- Pour les enquêtes périodiques fondées sur un plan d'échantillonnage où la taille d'échantillon augmente en même temps que la population, il est souvent nécessaire de développer une méthode pour stabiliser la taille de l'échantillon et, par le fait même, les coûts de collecte. Il peut s'agir, par exemple, de la suppression aléatoire, qui permet de stabiliser la taille de l'échantillon, au fil du temps.
- S'assurer que le plan des enquêtes périodiques soit aussi souple que possible, pour pouvoir faire face aux changements futurs, comme l'augmentation ou la réduction de la taille de l'échantillon, la restratification, le rééchantillonnage et l'actualisation des probabilités de sélection. Si des estimations sont requises pour des domaines précis (p. ex., estimations infraprovinciales), former les strates nécessaires à ce calcul en combinant de petites unités stables liées aux domaines concernés (p. ex., petites régions géographiques), si possible. Il sera plus facile, ainsi, de s'adapter à d'éventuels changements dans la définition des strates.
- Si des estimations efficaces du changement sont requises ou si le fardeau de réponse pose problème dans le cadre d'enquêtes périodiques, utiliser un plan d'échantillonnage avec rotation, qui remplace une partie de l'échantillon à chaque période. Le choix du taux de rotation visera le juste équilibre entre la précision nécessaire à l'estimation du changement et le fardeau de réponse des unités déclarantes. Un faible taux de rotation augmente la précision des estimations du changement, bien qu'il risque de diminuer le taux de réponse, au fil du temps, étant donné qu'il accroît le fardeau de réponse. Il a aussi l'avantage de réduire les coûts, lorsque le premier contact est beaucoup plus coûteux que les contacts subséquents.
- Élaborer des procédures visant à surveiller la qualité du plan d'échantillonnage au fil du temps. Mettre en place une stratégie d'actualisation pour le remaniement sélectif des strates gravement altérées par des fluctuations de croissance.

#### 4.3.4 Enquêtes longitudinales

- Pour les enquêtes longitudinales par panel, déterminer la durée du panel (sa durée dans l'échantillon) en tentant de maintenir un juste équilibre entre la satisfaction des besoins de l'enquête (données relatives à la durée), d'une part, et les effets d'attrition et de conditionnement de l'échantillon, d'autre part. Adopter un plan par panels chevauchants (c.-à-d. chevauchement temporel) lorsqu'il faut produire des estimations transversales parallèlement aux estimations longitudinales.
- Il importe grandement de choisir un plan d'échantillonnage dont les caractéristiques sont simples (c.-à-d. base de sondage unique, réduction du nombre de degrés et de phases au minimum), car les procédures d'estimation deviennent extrêmement complexes avec l'augmentation du nombre de vagues.
- Il est recommandé de réserver – principalement – les enquêtes longitudinales à la production d'estimations longitudinales. En tentant de satisfaire à la fois des exigences transversales et longitudinales, on risque de développer un plan et des procédures d'estimation très complexes. S'il faut obtenir des estimations transversales, il est préférable de se servir d'un échantillon « complémentaire » pour tenir compte des naissances et des nouveaux immigrants.

#### 4.3.5 Mise en œuvre

- À l'étape de la mise en œuvre, comparer l'échantillon réel, du point de vue de la taille et des caractéristiques, à l'échantillon attendu. Comparer la précision des estimations aux objectifs sur ce plan. Réévaluer les hypothèses formulées pendant la conception du plan. Par exemple, évaluer la non-réponse (contacts échoués, refus, etc.) et calculer les effets de plan de sondage.
- Préférer les logiciels de sélection d'échantillon généralisés aux systèmes personnalisés. Il peut s'agir du Système généralisé d'échantillonnage (SGECH) mis au point par Statistique Canada. Le SGECH s'avère très utile pour gérer la sélection et la rotation des échantillons, dans le cadre des enquêtes périodiques. En recourant aux systèmes généralisés, on peut s'attendre à réduire les erreurs de programmation et, dans une certaine mesure, les coûts et la durée du développement.

#### 4.3.6 Documentation

- Préparer des documents détaillés et exhaustifs pour chaque aspect du plan d'échantillonnage. Ces documents vont répondre, notamment, aux questions suivantes : quelles bases ont été retenues et pourquoi, comment les unités ont-elles été formées et stratifiées, comment la taille de l'échantillon a-t-elle été déterminée, comment les degrés ou les phases ont-ils été choisis, quels plans de sondage ont été retenus et pourquoi, etc.

### 4.4 Indicateurs de qualité

Principal élément de la qualité : exactitude

En plus de tenir compte des éléments suivants, le lecteur devrait consulter la Politique visant à informer les utilisateurs sur la qualité des données et la méthodologie de Statistique Canada, qui contient des renseignements pertinents, surtout au paragraphe 2.3 de la section E.1.

- Fournir des mesures de la représentativité de l'échantillon : surdénombrement et sous-dénombrement, exclusions, comparaisons avec des sources externes (p. ex., comparer des totaux démographiques externes à ceux obtenus dans le cadre de l'enquête).
- Comparer la taille de l'échantillon observé et celle de l'échantillon prévu. Cette comparaison est d'autant plus importante, pour les enquêtes à plusieurs degrés, car il peut s'avérer difficile, pour les degrés supérieurs à un, de prévoir avec précision la taille de l'échantillon.
- Comparer les taux de réponse, les taux d'attrition et les taux d'unités hors champ à ceux prévus lors de la planification.

- Fournir des mesures de l'erreur d'échantillonnage : produire des variances ou des coefficients de variation (c.v.) et les comparer aux valeurs prévues lors de la planification. Si l'on a recouru à des effets de plan, lors de la planification, il faut les comparer aux effets de plan réels.
- Comparer les c.v. réels des variables employées pour stratifier la base ou répartir l'échantillon aux valeurs cibles établies lors de la conception.
- Si possible, comparer l'homogénéité des strates à celle affichée au moment de leur création. Dans le cas des enquêtes répétées, étudier la détérioration des strates, au fil du temps. Mesurer la fréquence des unités migrantes et des erreurs de classification.

## Bibliographie

BETHEL, J. 1989. « Sample Allocation in Multivariate Surveys », *Survey Methodology*, vol. 15, n° 1, p. 47 à 57.

COCHRAN, W.G. 1977. *Sampling Techniques*, New York, Wiley, 428 p.

GAMBINO, J. 2001. « Design Effect Caveats », Statistique Canada. Document interne.

GABLER, S., S. HADER et P. LYNN. 2006. « Design Effects for Multiple Design Samples », *Survey Methodology*, vol. 2, n° 1, p. 115 à 120.

HIRIDIGLOU, M.A. 1994. « Sampling and Estimation for Establishment Surveys: Stumbling Blocks and Progress », *Proceedings of the Section on Survey Research Methods*, American Statistical Association, p. 153 à 162.

HIRIDIGLOU, M.A. et K.P. SRINATH. 1993. « Problems Associated with Designing Sub-annual Business Surveys », *Journal of Business and Economic Statistics*, n° 11, p. 397 à 405.

KALTON, G. et C.F. CITRO. 1993. « Panel Surveys: Adding the Fourth Dimension », *Survey Methodology*, vol. 19, n° 2, p. 205 à 215.

KISH, L. 1965. *Survey Sampling*, New York, Wiley, 664 p.

KISH, L. 1988. « Multi-purpose Sample Designs », *Survey Methodology*, vol. 14, n° 1, p. 19 à 32.

LAVALLÉE, P. 2007. *Indirect Sampling*, New York, Springer, 256 p.

LOHR, S. 1999. *Sampling. Design and Analysis*, Californie, Duxbury Press, 512 p.

SARNDAL, C. E., B. SWENSONN et J. WRETMAN. 1992. *Model Assisted Survey Sampling*, New York, Springer-Verlag, 694 p.

STATISTIQUE CANADA. 2008. *Méthodologie de l'Enquête sur la population active du Canada*, publication n° 71-526-X au catalogue de Statistique Canada, Ottawa, 116 p.

STATISTIQUE CANADA. 2003. *Méthodes et pratiques d'enquête*, publication n° 12-587-X au catalogue de Statistique Canada, Ottawa, 396 p.

THOMPSON, S.K. et G.A. SEBER. 1996. *Adaptive Sampling*, New York, John Wiley and Sons, 288 p.

TILLÉ, Y. 2001. *Théorie des sondages – Échantillonnage et estimation en populations finies*, Paris, Dunod.

## 5 Conception du questionnaire

### 5.1 Portée et objet

Le questionnaire est composé d'une série de questions permettant de recueillir des informations auprès d'un répondant. Puisqu'il permet le lien entre ce répondant et le chercheur, il joue un rôle de premier plan dans le processus de collecte des données. Les réponses aux questionnaires sont recueillies soit par un intervieweur, soit par l'intermédiaire de méthodes grâce auxquelles le répondant donne seul les informations.

### 5.2 Principes

Les questionnaires jouent un rôle de premier plan dans le processus de collecte des données, en plus d'influencer l'image de l'organisme statistique qui s'en sert. De même, ils influent considérablement sur le comportement du répondant, le rendement de l'intervieweur, le coût de la collecte et les relations avec le répondant; ils ont donc un impact considérable sur la qualité des données.

S'il est bien conçu, le questionnaire devrait recueillir des données correspondant à l'Énoncé des objectifs de l'enquête. Il devrait aussi tenir compte des besoins des utilisateurs en matière de statistique, des exigences administratives et des exigences relatives au traitement des données, ainsi que de la nature et des caractéristiques de la population de répondants. De même, tout bon questionnaire impose un faible fardeau de réponse et s'avère convivial tant pour son répondant que pour l'intervieweur. Enfin, la conception et la formulation de ses questions doivent inciter celui qu'il interroge à donner des réponses les plus exactes possible.

Pour ce faire, il faut qu'il se concentre sur le sujet de l'enquête, s'avère le plus bref possible, comporte des questions qui s'enchaînent bien et facilite le rappel des répondants. S'il est bien conçu, il facilitera, en outre, le codage et la saisie des données. Par ailleurs, il devrait réduire les tâches de vérification et d'imputation au minimum, en plus d'entraîner une réduction globale du coût et du temps consacrés à la collecte et au traitement des données. Pour plus de renseignements, à ce sujet, se reporter à la « Politique concernant l'examen et la mise à l'essai des questionnaires » (Statistique Canada, 2002a).

### 5.3 Lignes directrices

#### 5.3.1 Information aux répondants

- Statistique Canada a pour politique d'informer les répondants sur tout ce qui suit : le but de l'enquête (incluant la description des utilisateurs des statistiques qu'elle vise à produire et l'usage qu'ils comptent en faire), l'autorité qui la régit, les détails concernant son enregistrement, puis tout ce qui touche la participation des répondants (obligatoire ou volontaire), la protection de la confidentialité, les plans de couplage des enregistrements et les ententes relatives au partage des renseignements des répondants (incluant l'identité de ceux qui la concluent). Pour plus de renseignements, à ce sujet, se reporter à la « Politique d'information des répondants aux enquêtes » (Statistique Canada, 1998a).

#### 5.3.2 Pertinence

- Consulter les utilisateurs des données pendant la conception du questionnaire permet de s'assurer qu'on saisit bien la façon dont ils comptent s'en servir. Il faut aussi consulter les textes dédiés au sujet de l'enquête, de même que les enquêtes préalablement menées dans ce domaine, sur les plans national et international, avant de concevoir un nouveau questionnaire. Ainsi, on devrait concevoir un questionnaire efficace répondant aux besoins des utilisateurs.

#### 5.3.3 Contenu et formulation

- Les premières questions doivent pouvoir s'appliquer à tous les répondants, être faciles à comprendre et susciter l'intérêt. Il faut aussi qu'elles montrent que le répondant fait partie de la population visée par l'enquête.

- Formuler les questions en utilisant des mots et en référant à des concepts ayant la même signification pour les répondants et les concepteurs. Dans le cas des entreprises, choisir des questions, des périodes de référence et des catégories de réponse correspondant aux pratiques de l'établissement en matière de tenue de registres.
- Choisir un modèle de questionnaire et des termes qui incitent ceux qu'on interroge à donner des réponses les plus exactes possible. Il faut que le questionnaire se concentre sur le sujet de l'enquête, s'avère le plus bref possible, comporte des questions qui s'enchaînent bien, facilite le rappel des répondants et les oriente vers une source d'information appropriée (Converse et Presser, 1986 et Fowler, 1995).

#### **5.3.4 Cohérence**

- Dans la mesure du possible, harmoniser les concepts et la terminologie employés avec ceux qui prévalent dans l'usage. S'il y a lieu, reprendre des questions tirées d'autres enquêtes.
- Vérifier la concordance des versions française et anglaise du questionnaire.
- Tous les membres de l'équipe participant au projet devraient prendre part à l'examen du questionnaire. Chacun étant susceptible d'exprimer un point de vue différent, cette démarche promet d'enrichir la réflexion à ce sujet. Le débat peut porter sur la capacité du questionnaire à produire des données d'enquête de bonne qualité, sur la programmation simple (dans un environnement assisté par ordinateur) ou sur le choix d'un traitement post-collecte efficace. Plus particulièrement, les membres de l'équipe peuvent évaluer la complexité des questions et leur enchaînement, l'effet de leur structure sur le répondant et le rapport entre leur niveau de détail et la taille de l'échantillon et du plan d'analyse.

#### **5.3.5 Présentation du questionnaire**

- Concevoir les questionnaires que les répondants remplissent de manière autonome afin qu'ils soient intéressants et faciles à remplir. Dans le même ordre d'idées, soigner la lettre d'accompagnement et la page couverture du questionnaire, afin de faire bonne impression. Développer un questionnaire qui porte la marque du professionnalisme. S'il est lu par un intervieweur, s'assurer que le questionnaire est convivial pour celui qui se livre à cette tâche.
- Afin de réduire au minimum le risque d'éventuelles erreurs de déclaration, s'assurer que les instructions à l'intention des répondants et/ou des intervieweurs sont courtes, précises et faciles à repérer. Fournir les définitions requises au début du questionnaire ou en même temps que la question à laquelle elles s'appliquent. S'assurer que les périodes de référence et les unités de réponse sont claires et évidentes pour le répondant; utiliser des caractères gras pour souligner les éléments importants; préciser « incluez » ou « excluez » au sein des questions (et non au sein de directives distinctes); veiller, enfin, à ce que les catégories de réponses soient mutuellement exclusives et à ce qu'elles soient exhaustives.
- En ce qui concerne la présentation du questionnaire, donner des titres ou des en-têtes à chaque section du questionnaire. Intégrer des directives et des zones de réponse qui favorisent des réponses exactes. Se servir de couleurs, d'ombrages, d'illustrations et de symboles pour retenir l'attention des répondants ou des intervieweurs, afin de leur montrer les parties du questionnaire qu'il faut lire et les espaces où inscrire ses réponses. À la fin du questionnaire, aménager un espace pour recevoir les commentaires des répondants et adresser une formule de reconnaissance aux répondants (Converse et Presser, 1986, Fowler, 1995).

#### **5.3.6 Collecte des données**

- Lors de la conception du questionnaire, étudier et évaluer soigneusement différents modes de collecte des données. Peser les « pour » et les « contre » des nouvelles méthodes de collecte électroniques ou par Internet. Le choix d'un mode de collecte exerce des répercussions sur le niveau de détail et la complexité des questions, de même que sur le nombre qu'il est possible de poser. Le sujet de l'enquête et la nature parfois délicate des questions posées doivent aussi être considérées lors du choix du mode de collecte.

- Sensibiliser les concepteurs d'enquêtes et les analystes de données au fait que le mode de collecte influe sur la qualité et la mesure des renseignements recueillis.
- Lors de la conception du questionnaire, se fonder sur les règles optimales de chaque mode de collecte des données. Se rappeler que les aspects suivants peuvent avoir des répercussions considérables sur le comportement des répondants : choix de questions ouvertes et fermées, latitude du répondant face aux choix de réponse (possibilité de cocher un seul ou tous ceux qui s'appliquent), attribution d'un classement et d'une note, ordre des questions et des réponses. (De Leeuw, 2005, et Dillman et Christian, 2003).

### 5.3.7 Essais et évaluation

- Explorer une vaste gamme de méthodes pour évaluer le questionnaire et le mettre à l'essai. Il pourrait s'agir d'essais qualitatifs auprès de groupes de discussion ou encore d'essais cognitifs, préliminaires ou pilotes. La pertinence de telles méthodes et leur fréquence d'utilisation dépendent de plusieurs facteurs et circonstances, comme le type et la taille de l'enquête, son contenu, le choix des questions (questions d'enquêtes antérieures ou normalisées), le statut de la collecte (permanente ou non), la méthode de collecte, le calendrier du projet, le budget et l'accessibilité des ressources. De multiples révisions peuvent s'avérer nécessaires, qui auront un impact sur le coût et le calendrier du projet (Couper, Lessler, Martin, Martin, Presser, Rothgeb, et Singer, 2004).

## 5.4 Indicateurs de qualité

Principaux éléments de la qualité : exactitude, pertinence, cohérence.

- L'erreur de mesure survient lorsqu'il y a un écart entre les valeurs mesurées et les valeurs réelles. Il s'agit d'un biais (erreur systématique relevant d'un instrument de mesure inexact – elle demeure constante dans toutes les répétitions de l'enquête) et d'une variance (fluctuations aléatoires entre les mesures qui s'annulent les unes les autres, dans le cas d'échantillonnages répétés). Les sources d'erreur de mesure sont les suivantes : instrument de mesure, méthode de collecte des données, système d'information du répondant, répondant et intervieweur.
- Il faudrait que les utilisateurs puissent accéder à la description des processus destinés à réduire les erreurs de mesure imputables à l'instrument d'enquête et à optimiser, ce faisant, la comparabilité des données recueillies. Ils pourraient ainsi évaluer l'exactitude et la fiabilité des mesures ainsi que la cohérence des données recueillies par rapport à d'autres renseignements statistiques. Ces processus comprennent l'élaboration du questionnaire, les études pilotes, la mise à l'essai du questionnaire, la formation des intervieweurs, etc.
- L'exactitude des données et leur cohérence peuvent être influencées de plusieurs façons. Il peut s'agir, par exemple :
  - de la nature délicate des renseignements recherchés;
  - d'un écart, entre les répondants et le personnel d'enquête, dans l'interprétation de la terminologie et des concepts de l'enquête;
  - d'un écart entre les concepts et la terminologie, d'une part, et les termes en usage, particulièrement dans le cas des enquêtes-entreprises.

Le cas échéant, il est conseillé d'étudier l'orientation et la magnitude du biais imputable à ces éléments afin d'en évaluer l'incidence sur la qualité des données.

- Il faudrait mettre les questionnaires à la disposition des utilisateurs; cela les aidera à évaluer la pertinence et la cohérence des données en fonction de leurs propres besoins et des autres sources de données auxquelles ils ont accès.
- Tout au long du projet, il faut tenir les utilisateurs informés de l'ensemble des modifications apportées au questionnaire; il faudrait aussi évaluer l'incidence de ces modifications sur la comparabilité des données.
- Les documents portant sur la qualité des données doivent en outre contenir des renseignements sur les problèmes relatifs à la formulation des questions, sur le fardeau de réponse, sur les taux de refus ou sur tout autre renseignement pertinent.

### **Bibliographie**

CONVERSE J.M. et S. PRESSER. 1986. « Survey Questions: Handcrafting the Standardized Questionnaire », *Sage University Paper Series on Quantitative Applications in the Social Sciences*, no07-063, Thousand Oaks, Californie, Sage Publications, 80 p.

COUPER, M. P., J. T. LESSLER, E. A. MARTIN, J. MARTIN, J.M. ROTHGEB et E. SINGER. 2004. *Methods for Testing and evaluating survey questionnaires*, Wiley Series in Survey Methodology, Hoboken, New Jersey, John Wiley and Sons.

DE LEEUW E. 2005. « To Mix or not to Mix Data Collection Modes in Surveys », *Journal of Official Statistics*, vol. 21 n° 2, p. 233 à 255.

DILLMAN, D.A. 2000. *Mail and Internet Surveys. The Tailored Design Method, 2e édition*, John Wiley and Sons, Toronto.

DILLMAN, Don A. et M. CHRISTIAN LEAH. 2003. « Survey Mode as a source of instability in responses across surveys », conférence présentée au Workshop on Stability of Methods for Collecting, Analyzing and Managing Panel Data, March 2003, American Academy of Arts and Science, Cambridge, Massachusetts.

FOWLER, F.J. Jr. 1995. « Improving Survey Questions: Design and Evaluation », *Applied Social Research Methods Series*, n° 38, Thousand Oaks, Californie, Sage Publications, 200 p.

STATISTIQUE CANADA. 1998a. « Politique d'information des répondants aux enquêtes », *Manuel des politiques de Statistique Canada*. (en ligne), [http://icn-rci.statcan.ca/10/10c/10c\\_001\\_f.htm](http://icn-rci.statcan.ca/10/10c/10c_001_f.htm).

## **6 Collecte, saisie et codage des données**

### **6.1 Portée et objet**

La collecte des données réfère à tout processus dont l'objectif est d'acquérir ou de faciliter l'acquisition des données. On procède à cette collecte en demandant et en obtenant des données pertinentes auprès de personnes ou d'organismes au moyen de procédés adéquats. Ces données sont envoyées à l'organisme statistique par le répondant (autodénombrement) ou par l'intervieweur. La collecte comprend également les opérations visant à puiser des renseignements dans des sources administratives qui peuvent, dans certaines circonstances, requérir la permission du répondant afin d'accéder à ses données administratives.

La saisie des données réfère à tout processus inscrivant l'information communiquée par le répondant sur support électronique. Lorsqu'il n'est pas automatisé, ce processus nécessite l'intervention d'employés, qui saisissent les données recueillies au clavier (commis à la saisie des données). Le codage des données réfère à tout processus qui attribue une valeur numérique à une réponse. Bien que fréquemment automatisée, cette étape requiert parfois l'intervention humaine (commis au codage des données), car elle peut demander de prendre des décisions complexes.

Grâce à son degré d'automatisation élevé, les opérations de collecte permettent d'accéder à des par données, c'est-à-dire des informations liées au processus d'enquête. Il peut s'agir d'indicateurs attestant de la présence ou l'absence d'une unité dans un échantillon, de l'historique des appels et des visites, de l'historique des frappes (piste de vérification), du mode de collecte, d'informations administratives (p. ex. profil de l'intervieweur) et d'informations relatives au coût des opérations de collecte.

La collecte des données est plus qu'une source d'information; c'est aussi le principal lien entre l'organisme responsable de l'enquête et le grand public, qu'il faut convaincre d'y participer. La saisie et le codage des données permettent de produire des données formatées qui seront utilisées par tous les processus subséquents de l'enquête. Les opérations de collecte, de saisie et de codage des données accaparent une part importante du budget de l'enquête, car en plus de prendre du temps, elles requièrent des ressources humaines et matérielles considérables.

### **6.2 Principes**

Les répondants représentent la ressource la plus précieuse d'un organisme de réalisation d'enquêtes. Chaque variable qui ne peut être obtenue d'autres sources existantes devient un fardeau pour le répondant. Il importe de réduire au minimum le temps et l'énergie qu'un répondant doit consacrer à fournir des données. De même, il est impératif de satisfaire les exigences relatives à la protection des renseignements personnels et à la sécurité à toutes les étapes de la collecte et du traitement des données. Compte tenu de leur forte incidence sur l'exactitude des données, les opérations de collecte, de saisie et de codage devraient se faire à l'aide d'outils qui permettent de mesurer la qualité et le rendement.

### **6.3 Lignes directrices**

#### **6.3.1 Collecte des données**

- Pour bien planifier le processus de collecte, il faut déterminer les rôles et les responsabilités vis-à-vis tous les aspects de la collecte, incluant la stratégie de communication, la mise en œuvre, l'évaluation, la surveillance, la planification de mesures d'urgence et la sécurité.
- Concevoir le processus de collecte de façon à alléger le fardeau du répondant, à réduire les coûts et à accélérer l'obtention des données les plus exactes possible. La collecte peut se faire au moyen d'autodénombrements, d'interviews téléphoniques ou d'interviews sur place, basées sur un questionnaire papier ou effectuées par voie électronique (c.-à-d. déclaration électronique des données (DED), Internet et interviews assistées par ordinateur). Pour réaliser plus facilement les objectifs établis en matière de conception, le recours à plus d'une méthode pendant le cycle de collecte peut s'avérer

utile. Par exemple, la collecte peut débuter avec un autodénombrement basé sur un questionnaire papier ou accessible sur Internet, et se terminer avec une interview sur place. Dans le cas d'enquêtes à autodénombrement, il est possible de mettre en œuvre plusieurs stratégies au cours de la période de collecte afin d'encourager le retour des questionnaires (p. ex. envoi d'une fiche publicitaire avant la collecte de données, d'une lettre de présentation accompagnant le questionnaire ou d'une carte de rappel, rappel téléphonique ou visite sur place). Vérifier si certaines informations peuvent être acquises en consultant des dossiers administratifs plutôt qu'en recourant aux méthodes de collecte traditionnelles, plus coûteuses et parfois moins exactes. Envisager de recueillir des données dans le contexte d'une enquête supplémentaire à une enquête d'envergure. Cette approche permettrait non seulement de réduire les coûts d'enquête et le fardeau du répondant, mais aussi d'avoir accès à une mine de renseignements utiles à l'ajustement pour la non-réponse. Si possible, réaliser des études pilotes ou procéder à des essais pour évaluer et améliorer les opérations de collecte.

- Établir des procédures et des mesures de contrôle des échantillons pour chaque étape de la collecte de données (p. ex. livraison et renvoi des questionnaires papier, suivi des lacunes ou incohérences et suivi des cas de non-réponse). Grâce à ces procédés, les gestionnaires de la collecte et les intervieweurs seront en mesure d'évaluer l'état d'avancement des opérations, et ce, à tout moment, car ils seront renseignés sur le statut des unités échantillonnées, du début à la fin de la collecte. Ces mesures de contrôle sont particulièrement importantes dans le cadre des enquêtes dont les modes de collecte sont multiples et qui passent d'un mode à l'autre (ou d'un centre de collecte à un autre). Des procédures de contrôle de l'échantillon sont aussi utilisées pour s'assurer que chaque unité échantillonnée franchit toutes les étapes de traitement requises (c.-à-d. la saisie et le codage) et que son statut final est répertorié. Il est possible d'évaluer l'efficacité de ces procédures grâce à des mesures de contrôle des échantillons.
- Établir et entretenir de bonnes relations avec les répondants, afin d'obtenir des taux de réponse satisfaisants. Les mesures utiles à cet égard incluent, notamment, la promotion de l'enquête, l'envoi de lettres informant les répondants qu'ils ont été sélectionnés pour prendre part à l'enquête, la publication de statistiques clés pouvant inciter la participation de répondants (surtout pour les enquêtes longitudinales), l'adoption de stratégies facilitant la communication des informations destinées au public – sites Web, guide d'utilisation du questionnaire ou ligne d'information (surtout pour les enquêtes à autodénombrement) ou l'envoi d'une lettre remerciant les répondants de leur participation. Il s'agit là de mesures qui peuvent contribuer à sensibiliser les unités (individus ou organisations) sélectionnées dans l'échantillon à participer à l'enquête.
- Lors de la collecte des données, veiller à choisir un moment opportun pour communiquer avec le répondant ou la personne désignée du ménage ou de l'organisme répondant. Laisser le répondant communiquer les données selon la méthode et dans le format qui lui conviennent ou qui conviennent à son organisme. Ce faisant, on contribue à l'augmentation des taux de réponse et de la qualité des renseignements obtenus auprès des répondants. Il faut parfois assouplir les modalités de déclaration pour alléger le fardeau du répondant et faciliter la collecte de données. Par exemple, on pourrait suggérer des modalités particulières aux entreprises qui participent à plusieurs enquêtes à la fois. Pour les ménages, lorsque ceci peut être considéré comme une option, il serait utile d'établir des règles permettant de choisir un répondant substitut lorsque le répondant ciblé n'est pas disponible.
- Pour les interviews, choisir le moment opportun, pour appeler ou visiter les unités d'enquête, en se référant aux parodonnées produites lors des cycles précédents de l'enquête ou à celles d'une enquête similaire. Communiquer avec les répondants au moment qui leur convient le mieux et veiller à ce que le nombre d'appels ou de visites n'excède pas les limites acceptables. Établir un ordre de priorité des unités sélectionnées pour établir le contact et faire les entrevues. Ces priorités devraient permettre d'atteindre les tailles d'échantillon visées afin d'obtenir des estimations fiables pour chaque domaine d'intérêt. Ce procédé permettra de formuler des estimations suffisamment exactes (ayant un biais et

une variance faibles) pour être diffusées. Concrètement, dans le cas des enquêtes-entreprises, cela signifie que la priorité serait accordée aux unités les plus importantes ou les plus influentes – au risque de négliger les petites. Pour les enquêtes-ménages, une priorité plus élevée devrait être assignée aux unités les moins susceptibles de répondre à un questionnaire. Il est possible, également, de recourir à la fonction de pointage pour établir des priorités. En ce qui concerne les interviews téléphoniques, se servir d'un système automatisé pour gérer l'ordonnanceur d'appels. Ces systèmes devraient aussi établir les priorités.

- Les intervieweurs sont essentiels au succès de la collecte de données. Il faut donc développer avec soin les manuels et les activités de formation qui leur sont destinés. Cette pratique s'avère la meilleure approche pour obtenir des données de qualité (c.-à-d. un taux de réponse élevé et des réponses exactes), car on s'assure que les concepts et les sujets de l'enquête sont bien compris de tous et que les réponses aux questions sont, par conséquent, appropriées. Plusieurs approches sont possibles en matière de formation, de l'étude à domicile à l'apprentissage en classe, en passant par les interviews fictives ou sur place. Pendant la formation, il est important d'évaluer les compétences des intervieweurs afin de s'assurer qu'ils respectent une liste de critères préétablis (p. ex. la capacité de lire les questions de la façon dont elles sont écrites dans le questionnaire). Ce procédé permettra, en plus, de cerner les forces et les faiblesses de l'intervieweur, du point de vue de ses compétences, de lui faire part de commentaires et d'axer sa formation sur ses points faibles. La surveillance peut se faire sur place ou au moyen d'enregistrements, selon les ressources et le mode d'interview retenu. Consulter les intervieweurs et le personnel directement responsable des opérations de collecte afin de développer de meilleurs outils de formation. Il pourrait aussi s'avérer utile de faire un suivi auprès des répondants pour connaître leur point de vue sur le déroulement de l'interview.
- Mener des recherches pour dépister les répondants et communiquer avec eux lorsque les coordonnées associées à leur nom, dans les unités d'enquête, ne semblent pas à jour. Le dépistage permet d'accroître les taux de réponse et de déterminer si l'unité échantillonnée est toujours admissible à l'enquête. Avant et pendant la collecte, procéder à la mise à jour des coordonnées en se servant des informations figurant dans les sources administratives (c.-à-d. fichiers téléphoniques, autres bases de sondage). Afin d'effectuer un dépistage de qualité, recueillir des renseignements supplémentaires au sujet de l'unité d'échantillonnage (p. ex. les noms d'autres membres de la famille, la relation, l'âge, etc.) pendant la collecte. Des informations supplémentaires provenant de connaissances locales peuvent également s'avérer utiles. Il est recommandé de former une équipe d'experts en dépistage lorsqu'il s'agit d'une enquête répétée ou d'une enquête dont la période de collecte s'étend sur plusieurs mois. Établir des mécanismes permettant aux répondants de faire la mise à jour de leurs coordonnées entre les cycles de l'enquête. Il peut s'agir, par exemple, de leur envoyer une carte d'« avis de changement d'adresse » en les priant de la faire parvenir au Bureau en cas de déménagement. Il est important, par ailleurs, de recueillir l'information obtenue lors du dépistage (p. ex. adresse électronique et numéro de cellulaire) car elle peut servir dans des cycles d'enquête subséquents.
- Pour les enquêtes à autodénombrement, dès la réception des données, vérifier s'il y a des lacunes ou des incohérences en ce qui a trait à l'exactitude de l'information de couverture et à la qualité des données fournies. Dans certains cas, il peut être nécessaire d'effectuer des interviews de suivi (p. ex. lorsque beaucoup d'éléments sont manquants). Fixer l'ordre de telles interviews en fonction de l'importance statistique des unités touchées et des éléments manquants.
- Compte tenu du fait que les enquêtes par autodénombrement tendent à engendrer des taux de réponse inférieurs, envisager de faire un suivi téléphonique auprès des non-répondants – ou de leur rendre visite – pour les convaincre d'y participer, ou en procéder à un interview. S'assurer d'informer rapidement le personnel de collecte de l'enregistrement des questionnaires retournés, afin d'éviter les suivis inutiles. Les suivis sont essentiels dans le cadre des enquêtes longitudinales, qui visent davantage le long terme, car leur échantillon risque une attrition cumulative (et éventuellement des biais) à chacun des cycles d'enquête, en cas de non-réponse. Le suivi appliqué aux unités non répondantes devrait respecter la démarche décrite au point précédent (enquêtes par interviews) pour la gestion des priorités. Il peut s'avérer utile de consulter des paradonnées (p. ex. nombre d'appels ou de visites) pour établir cet ordre de priorité.

- À la fin de la collecte, communiquer avec un sous-échantillon ou l'ensemble des unités non répondantes (y compris les cas non résolus) afin de vérifier qu'elles sont bien admissibles à l'enquête (entreprise active ou non, logement occupé ou non, etc.). Si tel est le cas, il est recommandé d'obtenir une donnée essentielle – comme la taille de l'unité (revenu total de l'entreprise, taille du ménage, etc.) – pour l'ajustement de la non-réponse. Dans certains cas, il est possible de consulter des données administratives courantes sur l'ensemble des unités non répondantes pour obtenir cette information ou pour l'estimer.
- Fournir des plans et des outils permettant de gérer activement la collecte de données en même temps qu'elle progresse. Il peut s'agir de comparer les efforts encourus pour la collecte aux résultats obtenus en utilisant des mesures de productivité (p. ex. nombre d'unités résolues, quotidiennement et au total) et des indicateurs des coûts (p. ex. heures de travail et dépenses liées aux déplacements de l'intervieweur, quotidiennement et au total). Lorsque comparés aux valeurs prévues, ces indicateurs aident également les gestionnaires à prendre de meilleures décisions pendant la période de collecte. Les indicateurs de productivité et de coût (par unité sélectionnée ou questionnaire complété) permettent aussi d'évaluer les coûts et l'effort additionnels à l'accroissement des taux de réponse plus particulièrement vers la fin de la collecte.
- Il faut tout mettre en oeuvre pour garantir la confidentialité des données. Les employés étant en contact avec des données confidentielles doivent connaître les pratiques exemplaires liées à l'impression, à la manipulation et au classement des documents papier, ainsi qu'au traitement des fichiers électroniques. De même, ils doivent respecter les règles relatives à la diffusion de l'information.
- Considérer la possibilité de lancer un programme visant la répétition des interviews afin d'évaluer l'exactitude des opérations de manière globale.
- Étudier les paradonnées pour déterminer comment accroître l'efficacité et la rentabilité des opérations (p. ex. ordre des appels, moment idéal pour les rappels, nombre optimal d'appels ou de visites, etc.), afin d'améliorer les processus et pratiques de collecte actuels et futurs. Par exemple, on peut planifier le prochain cycle d'enquête en se fondant sur la répartition et la durée moyenne des interviews. La durée de l'interview permet d'évaluer, en partie, le fardeau du répondant, ainsi que la nécessité d'offrir de la formation additionnelle à un intervieweur (lorsque la durée de ses interviews s'éloigne considérablement de la moyenne).

### 6.3.2 Saisie de données

- Concevoir le processus de collecte de façon à réduire les coûts, à accélérer l'obtention des données et à en garantir l'exactitude. Des données sont parfois saisies directement par le répondant (par exemple lors de l'utilisation d'Internet ou de la DED), ou par l'interviewer (au moyen d'ITAO ou d'IPAO). En procédant ainsi, on diminue considérablement les coûts liés à la saisie en même temps qu'on augmente sa rapidité, sans compter que l'exactitude des données peut s'accroître par l'intégration de règles de contrôle à l'application sur ordinateur. S'il est impossible de jumeler saisie et collecte, on peut confier la saisie à des préposés (entrée manuelle des données) ou l'automatiser (numérisation suivie d'une reconnaissance intelligente de caractères). Cette dernière option est préférable, car elle réduit les coûts et permet souvent d'améliorer l'exactitude des données.
- Ceux qui mènent des ITAO ou IPAO, et qui saisissent et codent fréquemment des données pendant la collecte, devraient employer des outils et des méthodes de collecte normalisés (p. ex. écrans standard et questions normalisées), afin de faciliter leur travail et de limiter les risques d'erreurs de saisie. Il est possible de valider les entrées de données élémentaires et de corriger des erreurs potentielles, lors de la collecte, en intégrant des règles de contrôle au système de collecte (p. ex. erreurs de frappe, erreurs de réponse et éléments manquants).

- Les préposés à la saisie des données sont essentiels au succès des opérations de saisie. Il faut s'assurer qu'ils possèdent la formation et les outils adéquats, d'où l'importance de préparer des documents et des activités de formation à leur intention, et de leur offrir des séances de formation. En plus d'améliorer les compétences des employés, ces mesures garantissent une saisie exacte des données. L'utilisation de méthodes de contrôle de la qualité est recommandée pour vérifier si le niveau d'exactitude des données saisies par les préposés correspond aux critères préétablis, et pour faire part de commentaires aux préposés à la saisie afin qu'ils améliorent leur rendement.
- Que ce soit à partir de questionnaires papier ou d'images numérisées, la saisie manuelle des données entraîne parfois des erreurs de frappe. Procéder à des vérifications en ligne pour repérer les erreurs que l'opérateur de saisie des données peut corriger (c.-à-d. vérifications qui pointeront les erreurs de frappe). Enregistrer ces cas pour analyse et examen ultérieurs. Lorsque les circonstances le permettent, tester la saisie manuelle avant la réalisation de l'enquête.
- Dans le cas d'une saisie de données automatisée, s'assurer que le questionnaire est conçu de façon à faciliter la numérisation et la reconnaissance intelligente de caractères.
- Lorsqu'on recourt à la saisie automatisée, il se peut que des questionnaires ne puissent pas être numérisés, ou qu'ils le soient sans qu'on puisse en reconnaître les caractères. Lorsque des questionnaires sont endommagés ou mal numérisés, il est suggéré de recourir aux services d'une équipe de préposés à la saisie des données.
- Il faut mettre à l'essai les systèmes automatisés de saisie des données qui sont basés sur la reconnaissance intelligente des caractères à partir d'images numérisées avant de les utiliser. Ces systèmes peuvent être à l'origine d'erreurs systématiques très élevées pour certaines données élémentaires. Envisager d'améliorer les algorithmes et leurs paramètres afin de réduire les taux d'erreur. Dans le cas des données élémentaires pour lesquelles les risques d'erreurs systématiques sont élevés, le recours aux services de préposés à la saisie des données est à considérer.
- Il est également souhaitable de réquisitionner des préposés à la saisie des données afin qu'ils vérifient l'exactitude de la saisie automatique à partir d'un échantillon de l'enquête. Les résultats de cette évaluation permettront d'améliorer le processus.
- Il est recommandé d'adopter des mesures de contrôle efficaces des systèmes pour garantir la sécurité des données saisies, leur transmission et leur manipulation, surtout lorsqu'on recourt aux nouvelles technologies (par ex., collecte des données par cellulaire ou Internet). Prévenir les pertes de données attribuables aux erreurs humaines ou aux défaillances de système, ainsi que la diminution de qualité, voire de crédibilité. Développer des procédures encadrant le processus de destruction des données qui ne sont plus nécessaires.

### **6.3.3 Codage des données**

- Concevoir le processus de codage de façon à réduire les coûts, à accélérer l'obtention des données et à en garantir l'exactitude. Des données élémentaires sont souvent précodées lors des collectes incluant des questions fermées. De toute évidence, cette approche diminue les coûts associés au codage et peut aussi augmenter l'exactitude. S'il est impossible d'y recourir, parce que la collecte inclut des questions ouvertes, il faut coder les données après la collecte, manuellement (préposés) ou automatiquement (en utilisant un système de codage automatisé par reconnaissance de textes). Cette dernière option est préférable, car elle permet souvent de réduire les coûts et d'améliorer l'exactitude des résultats.
- Lors des opérations de codage manuelles, s'assurer de suivre les procédures de manière constante pour toutes les unités d'enquête, afin d'éviter, le plus possible, de commettre des erreurs. Il est souhaitable que ces opérations soient assistées par ordinateur. Autoriser le personnel ou les systèmes à soumettre les cas difficiles à un petit groupe d'experts en la matière. Centraliser le traitement, afin de réduire les coûts et de pouvoir bénéficier plus facilement des connaissances des experts. Étant donné que la collecte implique son lot d'imprévu, privilégier des processus flexibles pour pouvoir faire des changements, si les normes en matière d'efficacité le requièrent. Si possible, mettre le codage manuel à l'essai avant la réalisation de l'enquête.

- Les préposés au codage des données sont essentiels au succès des opérations de codage. Il faut s'assurer qu'ils possèdent la formation et les outils adéquats, d'où l'importance de préparer des documents et des activités de formation à leur intention, et de leur offrir des séances de formation. En plus d'améliorer les compétences des employés, ces mesures garantissent un codage exact des données. L'utilisation de méthodes de contrôle de la qualité est recommandée pour vérifier si le niveau d'exactitude des données codées par les préposés correspond aux critères préétablis, et pour faire part de commentaires aux préposés au codage afin qu'ils améliorent leur rendement.
- Lors du codage automatisé, créer des fichiers de référence et en faire la mise à jour, afin de maximiser le nombre de phrases reconnues par le système et de limiter les erreurs. Lorsqu'on opte pour le codage automatisé, il arrive souvent que des cas échappent au codage. Il est donc recommandé de recourir aux services d'une équipe de préposés au codage des données pour traiter ces cas.
- Lors d'une enquête, les experts du codage des données devraient aussi procéder à l'évaluation d'un échantillon et vérifier l'exactitude des données codées automatiquement. Les résultats de cette évaluation peuvent servir à augmenter et améliorer le contenu des fichiers de référence utilisés pour le codage des données.

#### **6.3.4 Contrôle de la qualité**

- Employer des méthodes de contrôle de la qualité statistique pour évaluer et améliorer la qualité des opérations de collecte, de saisie et de codage. Analyser ces méthodes, ainsi que leurs résultats, afin d'identifier les principales causes d'erreurs. Transmettre des rapports de rétroaction à ce sujet aux gestionnaires, au personnel, aux sujets matières spécialisés et aux méthodologistes. Utiliser les mesures de la qualité et de la productivité pour orienter la rétroaction destinée aux intervieweurs ou aux opérateurs, de même que pour repérer les éléments engendrant des erreurs lors de la conception des opérations de la collecte ou dans ses procédures de traitement. Les rapports doivent contenir des informations sur la fréquence et sur les sources des erreurs (voir Mudryk et coll., 1994, 1996 et 2002; Mudryk et Xiao, 1996). Il existe divers logiciels facilitant ce genre d'opérations, comme le Système d'analyse des données de contrôle de la qualité (SADCQ) et >NWA Quality Analyst (voir Mudryk, Bougie et Xie, 2002).

#### **6.3.5 Analyse rétrospective**

- Procéder à l'évaluation rétrospective des opérations de collecte, de saisie et de codage des données, puis consigner les résultats en vue d'usages futurs. Évaluer les méthodes employées et tirer des leçons afin d'améliorer chacune des composantes. Les études réalisées après l'enquête s'avèrent souvent utiles à cet égard.
- Lors des processus d'enquête suivant la collecte, veiller à recueillir des informations indiquant si les outils et les procédures de collecte, de saisie et de codage nécessitent des améliorations, sur le plan de la qualité, en prévision des cycles d'enquête futurs. Par exemple, il se peut qu'on observe des indices suggérant l'existence d'un biais dans la réponse ou tout autre problème lié à la collecte, lors de la vérification ou de l'analyse des données.

### **6.4 Indicateurs de qualité**

Principaux élément de la qualité : exactitude

L'incidence des opérations de collecte et de saisie des données (y compris le codage) sur la qualité et les coûts est directe et cruciale, car ces données sont les principaux intrants de l'organisme responsable de la réalisation de l'enquête, et sont souvent à l'origine des dépenses les plus importantes dans le cadre de l'enquête. Par conséquent, la qualité de ces opérations influe grandement sur celle du produit final, plus particulièrement sur le plan de l'exactitude.

Les mesures de la qualité mises en place pendant la collecte des données permettent au gestionnaire de l'enquête de prendre des décisions concernant d'éventuelles modifications ou la restructuration du processus. Les mesures de la qualité les plus importantes sont les taux de réponse, les taux d'erreur de traitement, les taux de suivi et les taux de la non-réponse répartis selon la cause. Lorsqu'elles sont disponibles à tous les niveaux pour lesquels des estimations sont produites et à toutes les étapes du processus, ces mesures permettent d'évaluer le rendement et la qualité des données.

#### **6.4.1 Taux de déclaration par personne interposée**

Signaler les taux de déclaration par personne interposée (c.-à-d. le pourcentage de réponses fournies par un répondant autre que l'unité d'enquête sélectionnée) car ce sont des indicateurs d'erreurs de réponse potentielles.

#### **6.4.2 Taux de non-réponse**

Signaler les taux de non-réponse, car ce sont des indicateurs de biais de non-réponse. Les cas de non-réponse se divisent en plusieurs catégories, par exemple, l'impossibilité d'établir le contact avec le répondant, de son refus de répondre, de son absence temporaire, de problèmes techniques, d'obstacles linguistiques ou de l'état mental ou physique du répondant. Pour refléter l'incertitude liée à la couverture, la non-réponse d'une unité peut aussi se répartir entre les cas résolus (l'unité étant admissible à l'enquête) et les cas non résolus (l'admissibilité n'étant pas déterminée). Il faut indiquer les cas de non-réponse à certains éléments (p. ex. refus et « ne sait pas ») figurant dans les questions principales. Les taux de non-réponse relatifs à des éléments peuvent varier si l'information est recueillie au début ou vers la fin de la période de collecte (certains répondants exigeant davantage d'appels ou de visites). Qu'ils concernent une unité ou certains éléments, les taux de non-réponse peuvent être transmis en fonction du domaine d'intérêt à publier (ce qui en fait potentiellement aussi des critères de publication) et de la sous-population (c.-à-d. grandes et petites entreprises, adultes jeunes ou plus âgés, etc.) pour indiquer dans quelle mesure l'échantillon réel est représentatif de la population. Il est également possible de combiner le taux de non-réponse des unités avec celui enregistré pour des éléments particuliers afin de fournir le taux de non-réponse global par élément. Le taux de conversion des refus et le taux de conversion des cas dépistés (liés à des coordonnées erronées ou devenues désuètes) sont également des indicateurs utiles. Enfin, le taux de refus au premier contact peut être indiqué pour les enquêtes de nature plus délicate.

#### **6.4.3 Erreurs imputables à l'admissibilité ou à l'inadmissibilité à l'enquête**

Lorsqu'on mène une analyse approfondie afin d'évaluer l'exactitude du classement des unités non répondantes, durant la collecte, selon les critères d'admissibilité ou d'inadmissibilité (entreprises : actives ou inactives; logements : occupés ou non occupés), indiquer le taux d'unités classées comme étant admissibles alors qu'elles ne l'étaient pas, ainsi que le taux d'unités classées comme étant inadmissibles alors qu'elles étaient admissibles. Il peut être utile de produire ces taux en fonction des domaines d'intérêt mentionnés plus tôt.

#### **6.4.4 Répartition des interviews selon leur durée moyenne**

Il faut répertorier la durée moyenne des interviews et leur répartition en fonction de la durée. Plus précisément, il importe de révéler le pourcentage d'interviews particulièrement brèves, car elles peuvent indiquer des problèmes afférents aux données déclarées. L'analyse de la durée d'une interview permet également d'évaluer, en partie, le fardeau du répondant.

#### **6.4.5 Incidence du mode de collecte suggestion : effet de mode**

L'« effet de mode » est un biais de mesure attribuable à la méthode de collecte des données. En principe, on peut le mesurer grâce à des modèles expérimentaux qui répartissent les unités d'échantillonnage de façon aléatoire entre deux groupes ou plus. Chaque groupe est enquêté en utilisant un mode de collecte différent. Tous les autres aspects touchant la conception de l'enquête sont contrôlés. Les différences obtenues dans la répartition des réponses entre les différents groupes peuvent alors être comparées et évaluées. D'autres méthodes, comme celles basées sur des scores de propension ou sur une analyse de régression, peuvent être utilisées pour évaluer l'incidence du mode de collecte lorsqu'il est impossible de recourir à des modèles expérimentaux.

#### 6.4.6 Taux de rejet à la vérification

Indiquer le taux de rejet après vérification, ainsi que le nombre et le type de corrections apportées en fonction du domaine, du mode de collecte, du type de traitement, de la présence de données élémentaires et de la langue utilisée lors de la collecte. Ces informations permettent de mieux évaluer la qualité des données ainsi que l'efficacité de la fonction de vérification utilisée lors des opérations de collecte et de saisie. Les taux de rejet après vérification peuvent être répartis en fonction du motif du rejet : élément manquant, incohérence de l'élément déclaré par rapport aux valeurs normales de cet élément ou par rapport à d'autres éléments déclarés). Cette dernière composante est un indicateur d'erreur de mesure (c.-à-d. erreur de réponse + erreur de saisie).

#### 6.4.7 Taux d'erreurs imputables à la saisie ou au codage

Faire état des taux d'erreurs de saisie ou de codage associés aux opérations manuelles ou automatisées. Lorsque l'enquête conjugue des opérations automatisées et manuelles, indiquer un taux composite. Il est possible de calculer les taux globaux ainsi que les taux obtenus en fonction du domaine, du mode de collecte, du type de traitement, de la présence de données élémentaires et de la langue utilisée lors de la collecte.

### Bibliographie

BETHLEHEM, J., F. COBBEN et B. SCHOUTEN. 2008, « Indicateurs de la représentativité des réponses aux enquêtes », *Recueil du Symposium international sur les questions de méthodologie 2008*, Ottawa, Statistique Canada.

COUPER, M.P., R.P. BAKER, J. BETHLEHEM, C.Z.F. CLARK, J. MARTIN, W.L. NICHOLLS II et J. O'REILLY. 1998. *Computer Assisted Survey Information Collection*, New York, Wiley-Interscience, 653 p.

DIELMAN, L. et M.P. COUPER. 1995. « Data Quality in a CAPI Survey : Keying Errors », *Journal of Official Statistics*, vol. 11, n° 2, p. 141 à 146.

DILLMAN, D. A. 2007. *Mail and Internet Surveys. The Tailored Design Method*, New York, Wiley, 554 p.

GROVES, R.M. 1989. *Survey Errors and Survey Costs*, New York, John Wiley and Sons, 620 p.

GROVES, R.M., P. BIEMER, L. LYBERG, J. MASSEY, W.L. NICHOLLS et J. WAKSBERG. 1988. *Telephone Survey Methodology*, New York, Wiley-Interscience, 608 p.

GROVES, R.M. et S.G. HEERINGA. 2006. « Responsive Design for Household Surveys : Tools for Actively Controlling Survey Errors and Costs », *Journal of the Royal Statistical Society, série A*, vol. 169, n° 3, p. 439 à 357.

HUNTER, L. et J.-F. CARBONNEAU. 2005. « An Active Management Approach to Survey Collection », *Recueil du Symposium international sur les questions de méthodologie 2005*, Ottawa, Statistique Canada.

LAFLAMME, F., M. MAYDAN et A. MILLER. 2008. « Using Paradata to Actively Manage Data Collection », *Proceedings of the Section on Survey Research Methods*, American Statistical Association.

LAFLAMME, F. et C. MOHL. 2007. « Research and Responsive Design Options for Survey Data Collection at Statistics Canada », *Proceedings of the Section on Survey Research Methods*, American Statistical Association.

LAFLAMME, F. 2008, « Recherche sur la collecte des données à l'aide de paradonnées à Statistique Canada », *Recueil du Symposium international sur les questions de méthodologie 2008*, Ottawa, Statistique Canada.

LAFLAMME, F. 2008. « Understanding Survey Data Collection Through the Analysis of Paradata at Statistics Canada », American Association for Public Opinion Research 63rd Annual Conference 2008. *Proceedings of the Section on Survey Research Methods*, American Statistical Association.

LEPKOWSKI J.M. et coll. 2007. « Advances in Telephone Survey Methodology », *Second International Conference on Telephone Survey Methodology 2006*, Miami, Wiley Series in Survey Methodology Section, p. 363 à 367.

LYBERG, L., P. BIEMER, M. COLLINS, E. DE LEEUW, C. DIPPO, N. SCHWARZ et D. TREWIN. 1997. *Survey Measurement and Process Quality*, New York, Wiley-Interscience, 808 p.

MUDRYK, W., M.J. BURGESS, P. XIAO, 1996. « Quality control of CATI Operations in Statistics Canada », *Proceedings of the Section on Survey Research Methods*, American Statistical Association, p. 150 à 159.

MUDRYK, W., B. JOYCE, H. XIE, H. 2004. « Generalized Quality Control Approach For ICR Data Capture in Statistics Canada's Centralized Operations », *European Conference on Quality and Methodology in Official Statistics*, Federal Statistical Office Germany.

ROSENBAUM P.R. et D.B. RUBIN. 1983. « The Central Role of the Propensity Score in Observational Studies for Causal Effects », *Biometrika*, vol. 70, n° 1, p. 41 à 45.

STATISTIQUE CANADA. 1998a. « Politique d'information des répondants aux enquêtes », *Manuel des politiques de Statistique Canada* (en ligne), [http://icn-rci.statcan.ca/10/10c/10c\\_001\\_f.htm](http://icn-rci.statcan.ca/10/10c/10c_001_f.htm).

STATISTIQUE CANADA. 2001d. *Normes et lignes directrices pour la déclaration des taux de non-réponse*. Rapport technique de Statistique Canada.

STATISTIQUE CANADA. 2003. *Méthodes et pratiques d'enquête*, produit n° 12-587-XIF au catalogue de Statistique Canada, Ottawa, 396 p.

WILLIAMS, K., C. DENYES, M. MARCH et W. MUDRYK. 1996. « Quality Measurement in Survey Processing », *Recueil du Symposium international sur les questions de méthodologie 1996*, Ottawa, Statistique Canada.

## 7 Utilisation des données administratives

### 7.1 Portée et objet

Les dossiers administratifs sont constitués aux fins de l'administration de divers programmes non statistiques. Par exemple, on conserve des dossiers administratifs pour régulariser le mouvement transfrontalier des biens et des personnes, pour satisfaire aux exigences légales de l'enregistrement de certains événements, comme les naissances et les décès, et pour administrer les avantages sociaux (comme les pensions) ou les obligations (comme les impôts pour les particuliers ou les entreprises). Leur raison d'être est liée à la prise de certaines décisions; l'identité de l'unité correspondant à un dossier donné est donc capitale. Par contre, dans le cas des dossiers statistiques qui, eux, ne servent pas et ne peuvent pas servir de fondement à une mesure visant un individu ou une entreprise, l'identité des individus ne présente aucun intérêt une fois que la base de données est complètement constituée.

Le recours aux dossiers administratifs présente bon nombre d'avantages pour un organisme statistique et ses analystes. Les demandes en matière de statistiques liées à tous les aspects de notre vie, de notre société et de notre économie continuent de croître. Ces demandes se présentent souvent dans un contexte de contraintes budgétaires importantes. En outre, les organismes statistiques se soucient, tout comme plusieurs répondants, de l'accroissement du fardeau de réponse associé aux enquêtes. Les répondants peuvent également réagir de façon négative s'ils ont l'impression d'avoir déjà communiqué des renseignements semblables (p. ex. en ce qui concerne leur revenu) à des enquêtes et à des programmes administratifs. Comme ils sont déjà constitués, les dossiers administratifs n'augmentent ni le coût de la collecte de données, ni le fardeau imposé aux répondants. Les progrès technologiques ont également permis aux organismes statistiques de surmonter nombre d'obstacles associés au traitement des ensembles de données volumineux. Pour toutes ces raisons, les dossiers administratifs sont de plus en plus utilisés à des fins statistiques.

En matière de statistique, les dossiers administratifs servent (i) aux bases de sondage, à titre de base directement ou à titre de supplément d'une base existante, (ii) au remplacement de la collecte de données (p. ex. utilisation des données fiscales pour les petites entreprises au lieu de chercher à obtenir des données d'enquête à leur sujet), (iii) à la vérification et à l'imputation, (iv) à la tabulation directe, (v) à l'estimation de façon indirecte (p. ex. comme information auxiliaire lors de l'estimation par calage, de l'étalonnage ou de la calendarisation) et (vi) à l'évaluation de l'enquête, ce qui comprend la confrontation des données (p. ex. comparaison des estimations de l'enquête avec des estimations issues d'un programme administratif connexe).

D'autre part, il faut se montrer prudent lorsqu'on utilise des données administratives en raison des limites dont on doit tenir compte, notamment (i) le niveau ou l'absence de contrôle de la qualité des données, (ii) la possibilité d'enregistrements partiels ou d'enregistrements complètement manquants (fichier incomplet), (iii) une différence conceptuelle susceptible d'occasionner des problèmes de biais et de couverture, (iv) l'actualité des données (la collecte des données échappant au contrôle de l'organisme statistique, il est possible qu'en raison de circonstances externes l'on ne reçoive pas une partie ou la totalité des données en temps opportun). Il faut également se rappeler que des coûts sont rattachés aux données administratives. Par exemple, des systèmes informatiques sont nécessaires pour nettoyer et compléter les données pour qu'elles soient utiles. Pour un examen des avantages et des inconvénients liés à l'utilisation de données administratives, voir Lavallée (2000).

### 7.2 Principes

Statistique Canada a pour politique d'utiliser des dossiers administratifs chaque fois qu'un tel recours constitue une solution de rechange rentable à la collecte directe de données. Tout comme pour n'importe quel programme d'acquisition de données, il convient de soupeser les coûts et les avantages liés à l'utilisation de dossiers administratifs à des fins statistiques; dans certains cas, cette solution évite les coûts

inhérents à la collecte de données et n'augmente pas le fardeau de réponse, pourvu que la couverture et le cadre conceptuel des données administratives soient compatibles avec la population cible. Dans d'autres circonstances, des coûts peuvent s'appliquer à la saisie des données ou un service peut être demandé en échange de cette utilisation. Selon l'usage qu'on prévoit en faire, il est souvent avantageux de combiner des données administratives avec des données provenant d'une autre source.

Le recours aux dossiers administratifs peut soulever des inquiétudes en ce qui concerne la protection de la confidentialité des renseignements issus du domaine public. Ces inquiétudes prennent de l'importance lorsque les dossiers administratifs sont appariés avec d'autres sources d'information. La Politique d'information des répondants aux enquêtes (Statistique Canada, 1998) exige que Statistique Canada informe tous ses répondants au sujet de renseignements tels que l'objet de l'enquête, les mesures de protection de la confidentialité, les plans de couplage des enregistrements et l'identité des parties à toute entente visant à partager les renseignements transmis par les répondants. Le couplage des enregistrements doit être conforme à la Politique relative au couplage d'enregistrements du Bureau (Statistique Canada, 2008). En particulier, toutes les demandes de couplage d'enregistrements doivent être présentées au Comité de la confidentialité et des mesures législatives et approuvées par le Comité des politiques. Les demandes sont normalement approuvées seulement pour des usages spécifiques. Cependant, dans certains cas, les demandes de données sont approuvées pour un usage récurrent ou continu.

Le recours aux données administratives peut nécessiter la mise en œuvre de certaines étapes – généralement un sous-ensemble – du processus d'enquête par le Bureau, étapes que nous avons décrites dans les sections précédentes. Cela s'explique du fait que plusieurs étapes du processus d'enquête (p. ex. la collecte directe et la saisie des données) sont réalisées par l'organisation responsable des données administratives. Par conséquent, il faut ajouter d'autres lignes directrices à celles qui ont été présentées afin de proposer des façons de compenser les différences au chapitre des objectifs de qualité de l'organisme responsable des données. Par exemple, il peut s'avérer nécessaire d'élaborer un programme complexe de vérification et d'imputation afin d'assurer un certain degré de qualité exigé pour l'utilisation des données.

Il ne faut pas oublier la raison fondamentale qui justifie l'existence de ces dossiers administratifs, qui sont le résultat d'un programme administratif mis en place pour des raisons administratives. Bien souvent, les utilisations statistiques de ces dossiers étaient inconnues lorsque le programme a été mis en œuvre et l'organisme statistique a invariablement une influence limitée sur l'élaboration du programme. Pour cette raison, toute décision relative à l'utilisation des dossiers administratifs doit être précédée d'une évaluation de ces dossiers sur le plan de la couverture, du contenu, des concepts et des définitions, des procédures d'assurance et de contrôle de la qualité mises en place par le programme administratif pour en assurer la qualité, de la fréquence des données, de la rapidité de l'organisme statistique à recevoir les données et de la stabilité du programme au fil du temps. Il va de soi que le coût associé à l'obtention des dossiers administratifs est également un facteur déterminant dans la décision d'y recourir ou non.

## **7.3 Lignes directrices**

### **7.3.1 Le programme administratif**

- Entretien des liens avec le fournisseur des dossiers administratifs. Il faut communiquer avec lui dès qu'on commence à les utiliser. Cependant, il est encore plus important de demeurer constamment en relation étroite avec le fournisseur afin que l'organisme statistique ne soit pas pris au dépourvu par les changements et puissent même les influencer. Les commentaires sur les données statistiques et sur leurs lacunes peuvent être utiles au fournisseur et auront pour effet de renforcer la source de données administratives.

- Il faut connaître le contexte dans lequel l'organisme administratif a créé le programme administratif (p. ex. législation, objectifs et besoins). Cela a une influence très importante sur i) la population couverte, ii) le contenu, iii) les concepts et les définitions, iv) la fréquence et l'actualité, v) la qualité de l'information recueillie et vi) la stabilité au fil du temps. Une attention particulière doit être accordée à la cohérence des concepts et à la qualité des données lorsqu'il existe des sources multiples de données administratives, par exemple lorsque chaque province gère son propre programme.
- Garder à l'esprit que si les renseignements fournis à la source administrative peuvent causer des gains ou des pertes à des personnes ou à des entreprises, il est possible qu'ils soient biaisés et entraînent des problèmes de couverture et de biais imprévus. Il pourra être nécessaire de mener des études spéciales pour évaluer et comprendre ces sources d'erreur.

### 7.3.2 Évaluer la qualité

- Bon nombre des lignes directrices présentées dans les sections précédentes s'appliquent aux dossiers administratifs. Les lignes directrices relatives à l'échantillonnage et à la saisie des données sont pertinentes si les dossiers administratifs existent uniquement sur papier et doivent être codés et saisis. Ces lignes directrices seront également précieuses pour les données administratives accessibles en format électronique, y compris la CED (collecte électronique de données). Soulignons que puisque ces données sont disponibles en format électronique, elles peuvent être implicitement instables et sujettes à d'autres erreurs causées par les processus de traitement et de transmission des données à la source. Les lignes directrices se rapportant à la vérification et à la diffusion doivent être respectées lorsqu'on obtient ou crée un fichier de dossiers administratifs d'individus pour analyse et traitement ultérieurs.
- Collaborer avec les concepteurs chargés de remanier les systèmes administratifs ou d'en concevoir des nouveaux. Cette approche favorise l'intégration des exigences statistiques aux systèmes dès le début du projet. De telles possibilités sont rares; cependant, lorsqu'elles se présentent, les avantages éventuels de la participation de l'organisme statistique valent largement le temps et le travail que requiert cette participation.
- Examiner chaque donnée des dossiers administratifs qu'on prévoit utiliser à des fins statistiques. Évaluer la qualité des données. Comprendre les concepts, les définitions et les procédures qui sous-tendent la collecte et le traitement des données par l'organisme administratif. Certains éléments peuvent être de très mauvaise qualité et donc inutilisables. Par exemple, la qualité du codage d'une classification (comme la profession, l'activité industrielle, la géographie) peut être inadéquate d'un point de vue statistique ou en limiter l'utilisation.
- Garder à l'esprit que la longévité de la source des données administratives et sa portée sont, de façon générale, entièrement entre les mains de l'organisme administratif. Les éléments administratifs qui ont initialement dicté les concepts, les définitions, la couverture, la fréquence, l'actualité et les autres attributs du programme administratif peuvent, au fil du temps, subir des changements qui faussent les séries chronologiques dérivées de la source administrative. Il faut se tenir au courant de ces changements et gérer leur incidence sur le programme statistique.
- Effectuer une évaluation permanente ou périodique de la qualité des données transmises. L'assurance que la qualité des données est conservée est importante, car l'organisme statistique ne contrôle pas le processus de collecte des données. Cette évaluation peut consister en la mise en oeuvre de mesures de protection et de contrôle supplémentaires (p. ex. l'utilisation de méthodes et de procédures statistiques de contrôle de la qualité, les règles de vérification) au moment de la réception des données, de comparaisons avec d'autres sources ou d'études sur le suivi d'un échantillon. Une bonne habitude à prendre est de donner de la rétroaction à ses sources administratives afin de les aider à améliorer la qualité de leurs données.

### 7.3.3 Confidentialité

- Tenir compte des répercussions de la publication de données tirées des dossiers administratifs sur la protection des renseignements personnels. Bien que, en vertu de la Loi sur la statistique, Statistique Canada ait le droit d'accéder aux dossiers administratifs à des fins statistiques, il se peut que ceux qui ont fourni les renseignements à l'origine n'aient pas prévu que ces renseignements seraient utilisés de la sorte. (Statistique Canada, 2005). Par conséquent, les responsables de programme devraient être prêts à justifier cette utilisation secondaire et à expliquer qu'elle est sans gravité et qu'elle dessert l'intérêt public.
- On a parfois recours aux données administratives pour remplacer une série de questions au répondant. En pareil cas, il peut être nécessaire d'obtenir la permission du répondant; il faut alors respecter la Politique d'information des répondants aux enquêtes (Statistique Canada, 1998). En l'absence de consentement, il faut mettre en place des mécanismes de collecte afin de poser les questions d'enquête équivalentes aux répondants.
- Les données administratives renferment souvent des renseignements concernant des personnes ou des entreprises en particulier. Toutes les données diffusées par Statistique Canada sont assujetties aux dispositions - en matière de confidentialité - de la Loi sur la statistique, même lorsque les données sont déjà du domaine public. Par conséquent, il faut tenir compte des lignes directrices sur le contrôle de la divulgation lorsqu'on prépare toutes analyses de données en vue de leur diffusion, y compris la diffusion de données administratives.

### 7.3.4 Non-réponse

- Tout comme les données d'enquête, les données administratives ne sont pas à l'abri de la non-réponse, qu'elle soit partielle ou totale. Dans certains cas, le manque de rapidité dans l'obtention de toutes les données administratives donne lieu à un taux de non-réponse plus élevé. Par conséquent, certaines des lignes directrices sur la non-réponse s'appliqueront. À moins de pouvoir effectuer un suivi des non-répondants et obtenir les réponses voulues, il faut élaborer une procédure d'imputation ou de rajustement des poids pour composer avec la non-réponse. Les sources administratives sont parfois désuètes. Ainsi, dans le cadre du processus d'imputation, il faut accorder une attention spéciale à l'identification des unités actives et/ou inactives. Il peut également être nécessaire de recourir à l'imputation ou à la transformation (p. ex. calendarisation) lorsque certaines unités transmettent leurs données à une fréquence différente (p. ex. hebdomadaire ou trimestrielle) de la fréquence souhaitée (p. ex. mensuelle).

### 7.3.5 Couplage d'enregistrements

- Lorsqu'on doit coupler des dossiers administratifs (p. ex. pour le dépistage de répondants, pour compléter des données d'enquête ou pour analyser des données), on doit respecter la Politique relative au couplage d'enregistrements (Statistique Canada, 2008). L'utilisation d'une seule source de données administratives peut susciter de l'appréhension au plan de la protection des renseignements personnels; cependant, l'appréhension est multipliée lorsque la source administrative est couplée à d'autres sources. En pareil cas, il se peut que les sujets ne sachent pas que les renseignements fournis en deux occasions distinctes sont combinés. La Politique relative au couplage d'enregistrements vise à assurer que l'intérêt public de chaque couplage l'emporte largement sur les atteintes à la vie privée qu'il pourrait occasionner.
- Il n'est pas toujours facile de combiner une source de données administratives à une autre source d'information. Cette tâche est particulièrement ardue lorsqu'il n'y a pas de clé d'appariement commune aux deux sources et que des techniques d'appariement doivent être utilisées. En pareil cas, le type de méthode d'appariement (c.-à-d. l'appariement exact ou statistique) doit être choisi en fonction des objectifs du programme statistique. Lorsque le programme a pour but la création et la mise à jour d'une base de sondage, ou la vérification des données, il faut utiliser un appariement exact. Pour l'imputation ou la pondération, l'appariement exact est préférable, bien qu'un appariement statistique puisse suffire. Lorsqu'on couple les sources afin d'effectuer des analyses de données qui, autrement, ne pourraient pas être réalisées, l'appariement statistique (c.-à-d. l'appariement d'enregistrements ayant des propriétés statistiques similaires) peut s'avérer un choix judicieux (voir Cox et Boruch, 1988, Kovacevic, 1999).

- Lorsqu'on doit procéder à un couplage d'enregistrements, il convient de faire bon usage des logiciels existants. Il existe un certain nombre de logiciels bien documentés, par exemple le Système généralisé de couplage d'enregistrements de Statistique Canada.
- Lorsque les données de plusieurs sources administratives sont combinées, il faut accorder davantage d'attention à la réconciliation des différences potentielles dans les concepts, les définitions, les dates de référence, la couverture et les normes de qualité appliquées à chaque source de données. Parmi les exemples, mentionnons les sources de données sur l'éducation, les rapports sur la santé et le crime, ainsi que les registres des naissances, des mariages, des immatriculations et des véhicules enregistrés, qui sont fournis par diverses organisations et divers organismes gouvernementaux.
- Certaines données administratives sont de nature longitudinale (p. ex. l'impôt sur le revenu et la taxe sur les produits et services). Lorsque des enregistrements de périodes de référence différentes sont combinés, ils constituent des mines de données très riches pour les chercheurs. Il faut demeurer particulièrement vigilant lorsqu'on crée des bases de données longitudinales axées sur des personnes, car leur utilisation soulève des inquiétudes très sérieuses en matière de protection des renseignements personnels. L'identificateur doit être utilisé avec soin, car une unité peut changer d'identificateur avec le temps. Faire le suivi de tels changements afin que l'analyse temporelle des données soit adéquate. Dans certains cas, la même unité peut avoir deux identificateurs ou plus pour la même période de référence, ce qui engendre un dédoublement dans le fichier administratif. Il faut alors élaborer un mécanisme d'élimination du dédoublement.

#### **7.3.6 Documentation**

- Documenter la nature et la qualité des données administratives dès leur évaluation. Ce genre de documents aide les statisticiens à déterminer quels usages conviennent le mieux aux données administratives. Choisir des méthodes adéquates pour le programme statistique en fonction des données administratives et informer les utilisateurs de la méthodologie utilisée et de la qualité des données.

### **7.4 Indicateurs de qualité**

Principaux éléments de la qualité : pertinence, exactitude, actualité, cohérence.

#### **7.4.1 Pertinence**

Les éléments d'information saisis par le système administratif sont-ils le reflet des concepts et des définitions de l'utilisateur des données? Bien qu'il soit souvent moins onéreux d'extraire des données administratives que de les recueillir dans le cadre d'une enquête, les buts de l'analyse doivent être atteints au moyen des données administratives pour que l'opération en vaille la peine. Indiquer la source, la date de référence et la mesure dans laquelle les définitions et les classifications correspondent aux données de l'enquête et aux besoins des utilisateurs des données.

#### **7.4.2 Exactitude**

Il arrive souvent que les données administratives ne soient pas visées par les mêmes procédures de vérification que les données d'enquête. Certaines vérifications sont normalement effectuées par l'organisation administrative, mais leur nature et leurs objets sont habituellement différents de ceux de l'organisme statistique. Il s'ensuit que la qualité des données peut soulever des inquiétudes lorsqu'on utilise des sources administratives à des fins statistiques, particulièrement dans les cas où la possibilité de communiquer de nouveau avec le responsable de l'information est limitée. En outre, les données administratives échantillonnées peuvent ne pas adhérer à aucun plan d'échantillonnage standard, ce qui risque d'introduire des biais et compliquer le calcul des erreurs d'échantillonnage. Enfin, si l'on utilise des données administratives comme base de sondage en plus ou au lieu d'une base de sondage créée grâce à la collecte de données, il pourrait être impossible d'analyser les problèmes de couverture et de non-réponse. D'un point de vue positif, précisons qu'un bon nombre de sources de données administratives

sont des recensements, ce qui signifie qu'il n'y aura pas d'erreur d'échantillonnage dans les estimations qu'on en obtient. Il faut indiquer la contribution des données administratives aux estimations les plus importantes. Si elles servent de base de sondage, il faut déclarer le taux d'imputation pour la non-réponse partielle ou totale et expliquer comment l'imputation a été effectuée. Si on ne fait qu'additionner les données administratives de manière à produire une estimation, inclure une estimation de la perte de précision résultant de l'imputation. Si des données administratives constituent une partie de l'estimation, le reste étant pris en compte par des données d'enquête, déclarer la portion de la base de sondage couverte par les données administratives de même que la portion estimée. Calculer un taux de réponse en combinant la portion de données administratives et celle des données d'enquête selon les explications données par Trépanier et al. (2005).

#### **7.4.3 Actualité**

On doit considérer sérieusement l'actualité des données administratives. Il est fréquent que ce genre de données ne soient disponibles que longtemps après la période de référence. Dans le cas où l'on utilise des données administratives comme base de sondage, celles-ci risquent d'être désuètes au moment où elles pourraient être utilisées. De plus, si les données administratives sont intégrées aux données de l'enquête, il importe qu'elles soient aussi récentes que les données d'enquête, à défaut de quoi tout le processus risque d'être compromis. En revanche, il existe des cas où les systèmes administratifs sont maintenus en temps réel, et les données qu'on en extrait sont beaucoup plus à jour que ne le seraient celles d'une enquête distincte. Indiquer la date de mise à jour de toutes les données administratives utilisées. Expliquer les hypothèses formulées quant à l'utilisation de données administratives désuètes.

#### **7.4.4 Cohérence**

La cohérence est une autre composante importante des données administratives. Ces données sont normalement saisies en vue d'autres utilisations et il s'ensuit qu'elles ne s'intègrent pas forcément à des concepts susceptibles d'avoir déjà été définis par des intérêts en données statistiques. Cela peut se produire dans le cas des concepts et des définitions, mais aussi dans le cas de la couverture et du plan de sondage. Les données administratives peuvent ne couvrir qu'une portion de la population cible, ce qui en rend l'utilisation problématique, ou une stratégie d'échantillonnage pourrait avoir été utilisée, ce qui complique le calcul des poids d'échantillonnage. Dans certains cas, les concepteurs de l'enquête devraient prendre part à la conception des systèmes administratifs, ce qui aurait pour effet d'accroître grandement la cohérence des données. Dresser la liste de toutes les exclusions susceptibles de compliquer les comparaisons avec d'autres données. Les indicateurs peuvent inclure une mesure de la population cible qui n'a pas été couverte.

### **Bibliographie**

BABYAK, C. 2007. « Challenges in Collecting Police-Reported Crime Data », *ICES-III, Proceedings of the Third International Conference on Establishment Surveys, Survey Methods for Businesses, Farms and Institutions*, Montréal, 18 au 21 juin 2007, p. 959 à 966.

BRACKSTONE, G.J. 1987. « Utilisation des dossiers administratifs à des fins statistiques », *Techniques d'enquête*, n° 13, p. 35 à 51.

BRION, P. 2007. « Redesigning French Structural Business Statistics, Using More Administrative Data », *ICES-III, Proceedings of the Third International Conference on Establishment Surveys, Survey Methods for Businesses, Farms, and Institutions*, Montréal, 18 au 21 juin 2007.

COX, L.H. et R.F. BORUCH. 1988. « Record Linkage, Privacy and Statistical Policy », *Journal of Official Statistics*, n° 4, p. 3 à 16.

HAZIZA, D., G. KUROMI et J. BÉRUBÉ. 2007. « Sampling and Estimation in the Presence of Tax Data in Business Surveys at Statistics Canada », *ICES-III, Proceedings of the Third International Conference on Establishment Surveys, Survey Methods for Businesses, Farms and Institutions*, Montréal, 18 au 21 juin 2007.

KOVACEVIC, M. 1999. « Record Linkage and Statistical Matching – They Aren't the Same! », *SSC Liaison*, vol. 13, n° 3, p. 24 à 29.

LAVALLÉE, P. 2000. « Combining Survey and Administrative Data : Discussion Paper », *ICES-II, Proceedings of the Second International Conference on Establishment Surveys, Survey Methods for Businesses, Farms and Institutions*, Buffalo, New York, 17 au 21 juin 2000, p. 841 à 844.

LAVALLÉE, P. 2005. « Indicateurs de la qualité : combinaison des données d'enquêtes et des données administratives », *Recueil du Symposium international sur les questions de méthodologie 2005*, Statistique Canada, Ottawa.

MCKENZIE, R. 2007. « A Statistical Architecture for Economic Statistics », *ICES-III, Proceedings of the Third International Conference on Establishment Surveys, Survey Methods for Businesses, Farms, and Institutions*, Montréal, 18 au 21 juin 2007.

MICHAUD, S., D. Dolson, D. Adams et M. Renaud. 1995. « Combining Administrative and Survey Data to Reduce Respondent Burden in Longitudinal Surveys », *Proceedings of the Section on Survey Research Methods*, American Statistical Association , p. 11 à 20.

PENNECK, S. 2007. « The Future of Using Administrative Data Sources for Statistical Purposes », *ICES-III, Proceedings of the Third International Conference on Establishment Surveys, Survey Methods for Businesses, Farms, and Institutions*, Montréal, Québec, 18 au 21 juin 2007.

STATISTIQUE CANADA. 1998. « Politique d'information des répondants aux enquêtes », *Manuel des politiques de Statistique Canada*. (en ligne), [http://icn-rci.statcan.ca/10/10c/10c\\_001\\_f.htm](http://icn-rci.statcan.ca/10/10c/10c_001_f.htm).

STATISTIQUE CANADA. 2005. « Loi sur la Statistique », Ottawa, <http://www.statcan.gc.ca/about-aperçu/act-loi-fra.htm>.

STATISTIQUE CANADA. 2008. « Politique relative au couplage d'enregistrements », *Manuel des politiques de Statistique Canada* (en ligne), [http://icn-rci.statcan.ca/10/10c/10c\\_025\\_f.htm](http://icn-rci.statcan.ca/10/10c/10c_025_f.htm).

TRÉPANIÉ, J., C. JULIEN et J. KOVAR. 2005. « Reporting Response Rates when Survey and Administrative Data are Combined », *Proceedings of the Federal Committee on Statistical Methodology Research Conference*, Arlington, Virginie, 14 au 16 novembre 2005.

WALLGREN, A. et B. WALLGREN. 2007. *Register-based Statistics : Administrative Data for Statistical Purposes*, New York, John Wiley and Sons, 258 p.

## 8 Réponse et non-réponse

### 8.1 Portée et objet

Malgré les plus grands efforts que fournissent les gestionnaires d'enquête et le personnel des opérations pour optimiser la réponse, la plupart des enquêtes, sinon toutes, doivent faire face au problème de la non-réponse.

La réponse fait référence ici à toute donnée obtenue soit directement auprès d'un répondant, soit par l'utilisation de données administratives. Cette définition au sens large de la réponse est nécessaire pour refléter l'utilisation accrue de différentes stratégies de collecte pour une même enquête, une pratique devenue de plus en plus courante. De plus, tout comme les données d'enquête, les données administratives ne sont pas à l'abri de la non-réponse, qu'elle soit partielle ou totale. Cette non-réponse découle parfois du manque de rapidité à l'obtention de toutes les données administratives.

Pour qu'une unité soit classée comme répondante, le degré de réponse d'item ou de réponse partielle (où on obtient une réponse exacte seulement à l'égard de certaines données exigées du répondant) doit correspondre à un seuil minimal en deçà duquel on considère qu'il y a une non-réponse d'unité. Dans un tel cas, la personne, le ménage, l'entreprise, l'institution, l'exploitation agricole ou toute autre unité échantillonnée est considérée comme n'ayant fourni aucune réponse.

Les mécanismes de réponse classiques sont les suivants : non-réponse uniforme [ou réponse manquant entièrement au hasard] où la probabilité de réponse est complètement indépendante des unités et du processus de mesure, et est constante sur l'ensemble de la population ; non-réponse dépendant d'une variable auxiliaire [ou réponse manquant au hasard] où le mécanisme de réponse dépend de certaines données auxiliaires ou des variables disponibles pour toutes les unités mesurées et non-réponse dépendant de la variable d'intérêt [ou réponse ne manquant pas au hasard] où la probabilité de réponse dépend de la variable d'intérêt.

La non-réponse peut avoir deux effets sur les données : premièrement, elle introduit un biais dans les estimations lorsque les non-répondants diffèrent des répondants par rapport aux caractéristiques mesurées; deuxièmement, elle contribue à faire augmenter la variance totale des estimations, car la taille observée de l'échantillon est réduite par rapport à la taille initialement prévue.

### 8.2 Principes

Le degré des efforts fournis pour obtenir une réponse d'un non-répondant est fonction des contraintes de budget, de temps et du personnel, de son incidence sur la qualité générale et du risque de biais de non-réponse. Si la non-réponse persiste, diverses approches existent pour diminuer l'effet de la non-réponse. Les décisions concernant le degré de recherche convenable à instaurer pour élaborer des techniques de correction de la non-réponse seront fonction des mêmes contraintes mentionnées précédemment.

Lors d'entrevues téléphoniques ou personnelles, ou lors d'un suivi, tenter de recueillir autant que possible des informations de base sur le répondant pour éviter de faire des ajustements basés sur des hypothèses un peu plus tard.

Pour traiter la non-réponse, tirer parti autant que possible de l'information auxiliaire disponible.

Un programme efficace de relations avec les répondants, un questionnaire bien conçu, l'utilisation de la gestion active pour un suivi régulier sur les opérations de collecte et une collecte adaptative de données (Laflamme, 2008), sont des éléments essentiels à l'optimisation de la réponse.

## 8.3 Lignes directrices

### 8.3.1 Établir le taux de réponse anticipé

Un des points pour déterminer la taille d'échantillon et gérer la collecte est l'établissement du taux de réponse anticipé. Pour ce faire, utiliser, entre autres moyens, les résultats des cycles précédents de l'enquête, ceux d'un essai préalable ou ceux d'enquêtes semblables.

### 8.3.2 Réduire la non-réponse

S'assurer d'un degré de qualité acceptable durant toutes les étapes de planification et de mise en œuvre de l'enquête pour l'obtention d'un bon taux de réponse. Pour ce faire, il faut garder à l'esprit les facteurs suivants :

- Lors de la conception de l'enquête : L'expérience antérieure du même type d'enquête, le budget total, et l'affectation du budget entre les diverses opérations;
- La qualité de la base de sondage (en ce qui concerne la couverture de la population et la facilité à établir le contact avec le répondant), la population observée et la méthode d'échantillonnage;
- La méthode de collecte des données (par exemple, par la poste, par une interview personnelle ou par une interview téléphonique assistée par ordinateur, par la collecte électronique de données (CED), par l'Internet, ou par la combinaison de quelques méthodes), la période de l'année et la longueur de la période de collecte;
- La stratégie de communication qui sera utilisée pour informer les répondants de l'importance de l'enquête et pour maintenir les relations avec les répondants;
- L'utilisation et l'efficacité des mesures incitatives pour les répondants;
- Le fardeau de réponse imposé (longueur de l'interview, difficulté du sujet, choix du moment et périodicité de l'enquête); la nature et la délicatesse du sujet; la longueur et la complexité du questionnaire; la langue du questionnaire et les antécédents culturels des répondants;
- L'expérience antérieure et les compétences du personnel de collecte en relations interpersonnelles; leur charge de travail; les facteurs liés aux intervieweurs eux-mêmes, comme leur formation; et le roulement potentiel du personnel;
- L'efficacité et la portée de la méthodologie de suivi, ainsi que les difficultés prévues dans le dépistage des répondants qui ont déménagé.

Instaurer une collecte adaptative permettant à la stratégie de collecte d'évoluer avec le temps. Ceci requiert l'instauration d'une gestion active pour un suivi régulier sur les opérations de collecte et d'une collecte adaptative de données aux quatre phases de la collecte : avant le contact initial, après quelques essais, au milieu de la période de la collecte, et vers la fin de la collecte.

### 8.3.3 Mettre en place des procédures de suivi auprès des non-répondants au cours de la collecte

Effectuer un suivi auprès des non-répondants (tous ou un sous-échantillon de ces derniers). Le suivi auprès des non-répondants augmente le taux de réponse et peut aider à vérifier si les répondants et les non-répondants sont dotés de caractéristiques mesurées semblables. La stratégie d'enquête devrait tenir compte d'emblée de la non-réponse en adoptant une perspective de sélection à deux phases.

Établir la priorité des activités de suivi. Par exemple, dans les enquêtes auprès des entreprises, effectuer d'abord un suivi des grandes unités ou des unités influentes, possiblement au risque de manquer les plus petites unités. De même, accorder une plus grande priorité aux unités non-répondantes dans les domaines connus comportant un fort potentiel de biais de non-réponse. On peut utiliser une fonction de caractérisation pour établir la priorité du suivi.

Un suivi est particulièrement important dans le cadre des enquêtes longitudinales, où l'échantillon est assujéti à l'attrition croissante (et possiblement au biais) en raison de la non-réponse dans chacun des cycles d'enquête. Dans ce cas, il faut faciliter le dépistage de grande qualité; obtenir des données de contact additionnelles pour les unités échantillonnées à chaque cycle d'enquête; fournir une carte de changement d'adresse et demander à l'unité échantillonnée d'informer le Bureau si un déménagement a lieu entre les cycles d'enquête. Cela permettra d'obtenir des données de contact actualisées. En outre, les données administratives, les annuaires municipaux et téléphoniques, et de nombreuses autres sources, dont le savoir local, sont précieux pour le personnel de dépistage.

#### **8.3.4 Évaluer l'existence d'un biais potentiel de non-réponse**

Diverses approches existent pour déterminer s'il y a des différences entre les répondants et les non-répondants et évaluer le biais potentiel de non-réponse : suivi spécifique d'unités, suivi de non-répondants, et analyse des caractéristiques connues des répondants et des non-répondants. Les renseignements sur les non-répondants peuvent soit parvenir des vagues précédentes d'information (dans le cas d'enquêtes longitudinales ou avec groupes de renouvellement), ou en utilisant des sources de données externes (par exemple, les fichiers de données administratives ou de paradonnées).

#### **8.3.5 Déterminer le mécanisme de réponse**

L'analyse des caractéristiques des répondants et des non-répondants servira également à établir un modèle de non-réponse en vue de réduire autant que possible le biais dû à la non-réponse et à guider le choix de la méthode appropriée pour compenser la non-réponse.

Pour les enquêtes longitudinales, il faut tenir compte de la structure de non-réponse dans le temps (Hedeker et Gibbons, 2006).

#### **8.3.6 Déterminer une méthode de traitement de la non-réponse**

Les principales approches servant à composer avec les données manquantes sont l'imputation et la repondération.

Le traitement devrait être choisi en fonction du type de non-réponse (totale ou partielle), de la disponibilité de variables auxiliaires et de la qualité du modèle de réponse. De façon générale, la repondération est utilisée pour le traitement de la non-réponse totale. L'imputation est surtout utilisée pour le traitement de la non-réponse partielle, quoiqu'elle peut l'être pour le traitement de la non-réponse totale si des données auxiliaires sont disponibles (enquêtes répétées, données administratives ou autres)

La repondération vise à éliminer, ou du moins à réduire, le biais de non-réponse totale. On peut considérer la repondération sous deux angles : par modèle de non-réponse ou par calage (Särndal, 2007).

Pour l'approche du modèle de non-réponse, un modèle est développé pour estimer les probabilités de réponse inconnues. Les poids de sondage sont alors ajustés par l'inverse des probabilités de réponse estimées (Oh et Scheuren, 1983; Lynn, P., 1996). Pour obtenir une certaine protection contre l'inadéquation du modèle, il est suggéré de former des groupes de réponses homogènes c'est-à-dire de regrouper les unités ayant les mêmes caractéristiques et la même propension à répondre (Haziza et Beaumont, 2007). Plusieurs méthodes peuvent être utilisées à cet effet : les algorithmes d'arbres de décision, comme CHAID dans le logiciel Knowledge Seeker (Kass, 1980; Angoss Software, 1995), les modèles de régression logistique, la méthode de score, l'utilisation de données auxiliaires telles que les paradonnées (Beaumont, 2005; Eltinge, Yansaneh, 1997), etc.

Des systèmes développés à Statistique Canada permettent d'évaluer et de mesurer l'impact de la non-réponse et de l'imputation : GENESIS (GENERalized SIMulation System) quantifie la performance relative de méthodes d'imputation par l'utilisation d'études de simulation et SEVANI (Système pour l'Estimation de la VARIance due à la Non-réponse et à l'Imputation) calcule la variance due à la non-réponse (Beaumont, 2007). À noter que si la variance due à la non-réponse est importante comparativement à la variance d'échantillonnage pour une région donnée, il pourrait peut-être être indiqué, pour respecter le budget, de réduire la taille souhaitée de l'échantillon afin de consacrer davantage de ressources à la prévention de la non-réponse.

### **8.3.7 Évaluer et publier les taux de non-réponse**

Suivre les Normes et lignes directrices pour la déclaration des taux de non-réponse (Statistique Canada, 2001d) afin de faciliter la comparabilité entre les enquêtes. Ces normes décrivent les exigences pour la déclaration des taux de non-réponse, conformément à la Politique visant à informer les utilisateurs sur la qualité des données, pour les recensements ou enquêtes-échantillons qui sont fondés strictement sur la collecte de données directement auprès des répondants. Parmi les sujets abordés, il y a les taux de non-réponse pondérés ou non, les taux de réponse à la collecte des données et à l'estimation, la non-réponse pour les enquêtes secondaires ou longitudinales, le biais dû à la non-réponse, la surveillance des opérations d'enquête, l'évaluation des méthodes de collecte des données, les mesures de la couverture de la base de sondage, la création d'une base de données longitudinales, la déclaration des cas de non-réponse et les exigences en matière de déclaration dans la base de métadonnées intégrées.

Se référer au besoin à des mises en application particulières des normes en fonction de spécificités d'enquêtes. Par exemple, des articles récents traitent des enquêtes utilisant des données administratives pour certaines unités et des données d'enquête pour d'autres (Trépanier et al. 2005), de celles où un mode de collecte mixte est utilisé pour une même unité (Leon, 2007) ou des enquêtes à composition aléatoire (Marchand, 2008).

### **8.3.8 Déterminer et analyser les raisons de la non-réponse**

Noter les raisons de la non-réponse au moment de la collecte (p. ex., refus, non-contact, absence temporaire, problème technique) puisque le degré de biais de non-réponse peut différer en fonction de la raison.

## **8.4 Indicateurs de qualité**

Principal élément de la qualité : exactitude

### **8.4.1 Évaluer les taux de réponse et de non-réponse**

Rédiger une note sur le taux de réponse. Celui-ci peut être calculé de différentes façons, avec interprétations pour des objectifs différents. Se référer aux Normes et lignes directrices pour la déclaration des taux de non-réponse (Statistique Canada, 2001d). Rapporter les taux de réponse pondérés pour illustrer la contribution aux estimations, et utiliser les taux de réponse non-pondérés pour refléter le taux de participation dans la population d'enquête.

Rapporter le taux de non-réponse ventilé par différents types de non-réponse. Cette information peut être utilisée ultérieurement lors de la conception d'autres enquêtes et est utile pour les utilisateurs des données qui doivent interpréter les données. Les pourcentages des unités échantillonnées ayant refusé de répondre, identifiées comme hors champs, n'ayant pu être contacté pendant la période de collecte et ayant répondu partiellement pourraient également être d'intérêt.

Préciser si les estimations d'enquête sont ajustées ou non en vue de compenser la non-réponse. Si les estimations sont ajustées, une description de la procédure de correction doit être jointe.

### **8.4.2 Évaluer la variance due à la non-réponse**

Rapporter la variance due à la non-réponse. Pour ce faire, utiliser SEVANI lorsque le total ou la moyenne d'un domaine ou des méthodes de rééchantillonnage est estimé.

### **8.4.3 Étudier le biais**

Étudier le biais de non-réponse en fonction du mode de collecte et du type de non-réponse.

Dans le cas des enquêtes périodiques, effectuer des études périodiques du biais de non-réponse. Les résultats de ces études doivent être inclus dans les renseignements déclarés aux utilisateurs conformément à la politique.

S'il y a lieu, tenter d'évaluer la mesure dans laquelle les procédures corrigent le biais potentiel.

## Bibliographie

ANGOSS SOFTWARE. 1995. *Knowledge SEEKER – User's Guide*, ANGOSS Software International LTD.

BEAUMONT, J.-F. et J. BISSONNETTE. 2007. « Variance Estimation Under Composite Imputation Using an Imputation Model », conférence présentée à l'Atelier sur le calage et l'estimation dans les enquêtes, Ottawa.

BEAUMONT, J.-F. 2005. « L'utilisation de renseignements sur le processus de collecte des données pour traiter la non-réponse totale au moyen de l'ajustement de poids », *Techniques d'enquête*, vol. 31, n° 2, p. 249 à 254.

ELTINGE, J. L. et I. S. YANSANEH. 1997. « Méthodes diagnostiques pour la construction de cellules de correction pour la non-réponse avec application à la non-réponse aux questions sur le revenu de la U.S. Consumer Expenditure Survey », *Techniques d'enquête*, vol. 23, n° 1, p. 37 à 45.

FULLER, W.A. 1993. *Measurement Error Models*, New York, Wiley, 440 p.

HEDEKER, D. et R.D. GIBBONS, 2006. *Longitudinal Data Analysis*, New York, Wiley, 360 p.

GROVES, R.M., D.A. DILLMAN, J. L. ELTINGE et R. J. A. LITTLE. 2001. *Survey Nonresponse*, New York, Wiley, 520 p.

HAZIZA, D. et J.-F. BEAUMONT. 2007. « On the Construction of Imputation Classes in Surveys », *International Statistical Review*, vol. 75, n° 1, p. 25 à 43.

KASS, G.V. 1980. « An Exploratory Technique for Investigating Large Quantities of Categorical Data », *Applied Statistics*, vol. 29, n° 2, p. 119 à 127.

LAFLAMME, F. 2008. « Using Paradata to Actively Manage Data Collection Survey Process », *Proceedings from the American Statistical Society 2008 Joint Statistical Methods Conference*, Denver, Colorado, American Statistical Association.

LEON, C. A. 2007. « Reporting Response Rates in Characteristic Surveys », *Proceedings from the Statistical Society of Canada 2007 Conference*, St. John's, Terre-Neuve, Société statistique du Canada.

LYNN, P. 1996. « Weighting for Non-reponse. Survey and Statistical Computing », *Proceedings from the Association for Survey Computing 1996 Survey and Statistical Computing conference*, Londres.

MARCHAND, I., R. CHEPITA, P. ST-CYR et D. WILLIAMS. 2008. « La non-réponse dans le cadre d'une enquête à composition aléatoire : l'expérience du cycle 21 (2007) de l'Enquête sociale générale », conférence présentée au *Symposium international sur les questions de méthodologie de Statistique Canada*, 25 au 28 octobre 2008, Ottawa.

OH, H.L. et F. J. SCHEUREN. 1983. « Weighting Adjustment for unit nonresponse », *Incomplete data in Sample Surveys, Vol. 2. Theory and Bibliographies*, W.G. Madow, I. Olkin et D. B. Rubin, New York, Academic Press, p. 143 à 184.

SÄRNDAL, C.-E. et S. LUNDSTRÖM, 2005. *Estimation in Surveys with Nonresponse*, Wiley, New York, 212 p.

STATISTIQUE CANADA. 2000d. « Politique visant à informer les utilisateurs de la qualité des données et la méthodologie », *Manuel des politiques de Statistique Canada* (en ligne), [http://icn-rci.statcan.ca/10/10c/10c\\_010\\_f.htm](http://icn-rci.statcan.ca/10/10c/10c_010_f.htm).

STATISTIQUE CANADA. 2001d. « Normes et lignes directrices pour la déclaration des taux de non-réponse », Rapport technique de Statistique Canada.

TRÉPANIÉ, J., C. JULIEN et J. KOVAR. 2005. « Reporting response Rates When Survey and Administrative Data are Combined », *Federal Committee on Statistical Methodology Research Conference*.

## 9 Vérification

### 9.1 Portée et objet

La vérification des données est l'application de contrôles visant à détecter les entrées manquantes, invalides ou incohérentes ou à mettre en évidence les enregistrements des données qui sont susceptibles de contenir des erreurs. Certains de ces contrôles sous-tendent des relations logiques qui découlent directement des concepts et des définitions. D'autres sont de nature plus empirique ou sont le résultat de l'application d'essais ou de procédures statistiques (p. ex., des techniques d'analyse des valeurs aberrantes). Les contrôles peuvent être fondés sur des données tirées de collectes antérieures de la même enquête ou d'autres sources.

La vérification englobe une vaste gamme d'activités, dont les vérifications des intervieweurs sur le terrain et les avertissements générés par ordinateur au moment de la collecte ou de la saisie des données. Elle comprend en outre la détermination des unités en prévision du suivi et les vérifications détaillées de microdonnées. Enfin, elle comprend la localisation des erreurs pour les besoins de l'imputation, de même que les vérifications de relations complexes, au niveau des macrodonnées, aux fins de la validation des données.

### 9.2 Principes

Un enregistrement de données qui a été modifié par suite de vérifications devrait être plus près de la valeur réelle qu'avant ces modifications. Nous concevons les vérifications pour déceler et corriger les incohérences, et non pas pour produire un biais par suite de l'imposition de modèles implicites. Lorsqu'une vérification plus poussée a un effet négligeable sur les estimations d'enquête finales, il s'agit d'une survérification et elle devrait être évitée.

L'analyse des taux de rejet à la vérification et l'ampleur des changements découlant des vérifications fournissent des renseignements concernant la qualité des données d'enquête et peuvent aussi suggérer des améliorations à l'outil d'enquête.

### 9.3 Lignes directrices

#### 9.3.1 Conception

- La vérification contribue efficacement à déceler les erreurs fatales (Granquist et Kovar, 1997), puisque le processus peut facilement être informatisé. Exécuter cette activité le plus rapidement possible. Bien qu'une certaine intervention manuelle soit nécessaire, un logiciel généralisé et réutilisable peut être particulièrement utile pour cette tâche. Le système Banff de vérification et d'imputation (Statistique Canada, 2009) et le Système canadien de contrôle et d'imputation du recensement (SCANCIR) (Bankier et coll., 1999), sont des exemples de ce type de logiciel. Certaines applications personnalisées peuvent aussi être développées sur la base d'autres logiciels qui ne visent pas uniquement les processus de vérification. Logiplus, le système de Statistique Canada servant à gérer les tables logiques de décision, est un exemple d'un tel logiciel.
- L'informatisation permet aux gestionnaires d'enquête d'augmenter la portée et le volume des contrôles pouvant être effectués, ce qui est tentant pour eux. Minimiser le nombre de ces augmentations si elles font peu de différence dans les estimations de l'enquête. Plutôt que d'accentuer l'effort de vérification, réorienter les ressources vers des activités plus rentables (p. ex., l'analyse des données, l'analyse des erreurs de réponse).
- Déterminer les valeurs des données extrêmes d'une période d'enquête ou entre les périodes d'enquête (cet exercice est appelé processus de détection des valeurs aberrantes). La présence de ce type de données se démarquant de la distribution est un signe précurseur d'erreurs potentielles. Utiliser des méthodes de détection univariées simples (Hidioglou et Berthelot, 1986) ou des méthodes plus complexes et explicites (de Waal, 2000).

- L'incidence des erreurs s'est avérée très variable, particulièrement dans les enquêtes recueillant des données numériques. Il est fréquent qu'un petit nombre d'erreurs soient à la source de la majorité des changements apportés dans les estimations. Envisager d'effectuer la vérification de façon sélective, afin de réaliser des gains d'efficacité potentiels (Granquist et Kovar, 1997), sans incidence négative sur la qualité des données. Les priorités peuvent être établies en fonction des types ou de la gravité des erreurs ou en fonction de l'importance de la variable ou de l'unité déclarante.
- Les taux de succès des vérifications, soit la proportion des vérifications d'avertissement ou d'interrogation qui mettent en évidence les véritables erreurs, se sont avérés peu efficaces, souvent aussi bas que 20 ou 30 %. Élaborer des vérifications qui sont efficaces et contrôler l'efficacité sur une base régulière.
- Il est possible que les vérifications ne permettent pas de détecter les petites erreurs systématiques introduites constamment dans les enquêtes répétitives, erreurs qui peuvent donner lieu à d'importants biais dans les estimations. Le « resserrement » des vérifications n'est pas la solution. Pour détecter ce genre d'erreurs systématiques, utiliser d'autres méthodes, comme les méthodes classiques de contrôle de la qualité, l'analyse et l'examen approfondis des concepts et des définitions, les études postérieures aux interviews, la validation des données et la confrontation des données avec d'autres sources de données qui peuvent être disponibles pour certaines unités.
- Limiter le recours à la vérification pour résoudre les problèmes déjà survenus, surtout dans le cas des enquêtes répétitives. La contribution de la vérification à la réduction des erreurs est limitée. Bien qu'il soit essentiel d'effectuer un peu de vérification, il faut en réduire la portée et réorienter l'objectif. Attribuer une grande priorité à l'apprentissage tiré du processus de vérification. Pour réduire le nombre d'erreurs, s'attarder aux premières phases de la collecte de données plutôt qu'à l'épuration effectuée à la fin. Pratiquer la prévention plutôt que la correction des erreurs. À cette fin, ramener l'étape de vérification aux premières phases du processus d'enquête, de préférence lorsque le répondant est encore disponible, par exemple, en utilisant les méthodes d'interview téléphonique, d'interview sur place ou d'auto-interview assistées par ordinateur.
- Pendant la conception des processus de collecte des données, et particulièrement pendant la vérification et le codage, s'assurer que les procédures sont appliquées à toutes les unités d'étude le plus uniformément possible et qu'elles comportent le moins d'erreurs possible. L'automatisation est préférable. Permettre au personnel ou aux systèmes de soumettre les cas difficiles à un petit groupe de spécialistes compétents. Centraliser le traitement, afin de réduire les coûts et de faciliter le recours aux connaissances spécialisées disponibles. Comme l'information recueillie peut donner lieu à des résultats imprévus, utiliser des processus adaptables pour apporter les changements qui s'imposent s'il y a lieu de le faire du point de vue de l'efficacité.

### **9.3.2 Collecte des données et suivi des questionnaires rejetés au contrôle**

- La vérification peut être utile pour l'épuration de certaines données, mais son rôle principal est de permettre de fournir de l'information sur le processus d'enquête, soit en donnant des mesures de qualité pour l'enquête en cours, soit en suggérant des améliorations pour les enquêtes futures. Envisager la vérification comme une partie intégrante du processus de collecte des données, du point de vue de la collecte de renseignements sur le processus. Dans cette optique, la vérification peut être précieuse pour préciser les définitions, améliorer un instrument d'enquête, évaluer la qualité des données, déterminer les sources des erreurs non dues à l'échantillonnage, servir de base pour l'amélioration future du processus d'enquête complet et fournir des données valables pour améliorer d'autres processus d'enquête et d'autres enquêtes (Granquist, Kovar et Nordbotten, 2006). Afin d'atteindre cet objectif, superviser le processus et produire des pistes de vérification, des diagnostics et des mesures du rendement, et utiliser ces éléments pour établir les pratiques exemplaires.
- Au cours des suivis, ne pas surestimer la capacité des répondants de corriger les erreurs. Leur agrégation peut être différente, leur mémoire limitée et leur apport négligeable. Limiter l'activité de suivi des répondants.
- Pour les enquêtes-entreprises, élaborer une stratégie de suivi sélectif. L'utilisation d'une fonction de pointage (Latouche et Berthelot, 1992) concentre les ressources sur les unités d'échantillon importantes, les principales variables et les erreurs les plus graves.

### 9.3.3 Assurance de la qualité

- S'assurer que toutes les vérifications sont cohérentes à l'interne (c.-à-d. non contradictoires).
- Il faut se rappeler que l'utilisée de la vérification est limitée et que le processus peut en fait être improductif. Bien souvent, les changements de données fondés sur les vérifications sont considérés à tort comme des corrections de données. On peut prétendre qu'à un certain point du processus de vérification, on introduit autant d'erreurs qu'on en corrige. Identifier et respecter cette fin logique du processus.
- Appliquer de nouveau les vérifications aux unités qui ont subi des corrections pour s'assurer qu'aucune autre erreur n'a été introduite directement ou indirectement par le processus de correction.
- Ne pas sous-estimer la capacité du processus de vérification d'intégrer les données signalées aux modèles qu'imposent les vérifications. Il y a un réel danger de créer de faux changements pour la seule raison de s'assurer que les données ne sont pas rejetées à la vérification. Contrôler le processus!
- Le processus de vérification est souvent très complexe. Lorsque la vérification se fait sous le contrôle de StatCan, communiquer les procédures détaillées et à jour et offrir une formation appropriée à tout le personnel concerné et effectuer un suivi des travaux proprement dits. Envisager d'appliquer les procédures de contrôle de la qualité.
- Effectuer un suivi de la fréquence des rejets à la vérification, du nombre et du type de corrections apportées par strate, du mode de collecte, du type de traitement, des données élémentaires et de la langue utilisée pour la collecte. Cela aidera à évaluer la qualité des données et l'efficacité de la fonction de vérification.

### 9.4 Indicateurs de qualité

Principaux éléments de la qualité : exactitude, actualité

L'erreur de mesure est l'erreur qui se produit dans le cadre du processus de déclaration, tandis que l'erreur de traitement est celle qui se produit au moment du traitement des données. La dernière comprend les erreurs dans la saisie des données, le codage, la vérification et la totalisation des données, ainsi que dans l'affectation des poids d'enquête. Même s'il n'est habituellement pas possible de calculer l'erreur de mesure et l'erreur de traitement individuellement, le taux de rejet au contrôle donne une indication de leur importance combinée. Il s'agit du nombre d'unités rejetées par suite des contrôles de vérification divisé par le nombre total d'unités. Les produits doivent être accompagnés d'une définition des deux types d'erreurs et d'une description des principales sources d'erreurs. Cela éclaire les utilisateurs sur les mécanismes en place pour réduire l'erreur, étant donné que cela permet une collecte, une saisie et un traitement des données au point. Les erreurs de mesure et de traitement ont des répercussions sur le biais et la variance.

Les taux de vérification des variables clés doivent être indiqués. Ils peuvent être plus élevés en raison de l'erreur de mesure (p. ex., en raison du mauvais libellé de la question), ou à cause d'une erreur de traitement (p. ex., des erreurs dans la saisie des données).

La contribution totale des valeurs vérifiées aux estimations clés doit être indiquée. Il s'agit de la mesure dans laquelle les valeurs des estimations clés sont modifiées par les données qui ont été vérifiées. Cela peut donner une indication de l'effet de l'erreur de mesure sur les estimations clés. Cet indicateur s'applique uniquement aux moyennes et aux totaux.

La vérification des données est cruciale pour assurer l'exactitude et la cohérence des données. Toutefois, il peut s'agir d'une initiative coûteuse et longue. Il s'agit probablement de l'activité la plus coûteuse d'une enquête par sondage et d'un cycle de recensement. Lorsque cette vérification laborieuse, et souvent manuelle, a une incidence négligeable sur les estimations finales, on parle de survérification. Outre le fait qu'elle soit coûteuse sur le plan des finances, de l'actualité et de l'augmentation du fardeau de réponse, la survérification peut donner lieu à de graves biais, engendrés par l'intégration de données dans les modèles implicites qu'imposent les vérifications.

## **Bibliographie**

BANKIER, M., M. LACHANCE et P. POIRIER. 1999. « A Generic Implementation of the New Imputation Methodology », *Proceedings of the Survey Research Methods Section*, American Statistical Association, p. 548 à 553.

DE WAAL, T., F. VAN DE POL et R. RENSSSEN. 2000. « Graphical Macro Editing: Possibilities and Pitfalls », *Proceedings of the Second International Conferences on Establishment Surveys*, Buffalo, New York.

GRANQUIST, L. et J.G. KOVAR. 1997. « Editing of Survey data : How Much is Enough? », *Survey Measurement and Process Quality*, New York, Wiley, p. 415 à 435.

GRANQUIST, L., J. KOVAR et S. NORDBOTTEN. 2006. « Improving Surveys – Where Does Editing Fit in? », *Statistical Data Editing – Impact on Data Quality Vol. 3, Conference of European Statisticians*, United Nations Statistical Commission and United Nations Economic Commission for Europe.

HIDIROGLOU, M. A. et J.-M. BERTHELOT. 1986. « Statistical Editing and Imputation for Periodic Business Surveys », *Survey Methodology*, n° 12, p. 73 à 83.

LATOUCHE, M. et J.-M. BERTHELOT. 1992. « Use of a Score Function to Prioritize and Limit Recontacts in Editing Business Surveys », *Journal of Official Statistics*, n° 8, p. 389 à 400.

STATISTIQUE CANADA. 2009. Description des fonctions du système Banff pour la vérification et l'imputation. *Rapport technique de Statistique Canada*.

## 10 Imputation

### 10.1 Portée et objet

L'imputation est le processus qui permet d'attribuer des valeurs de remplacement à des données manquantes, invalides ou incohérentes rejetées à l'étape de vérification des données. Ce processus a lieu après le suivi auprès des répondants (si possible), l'examen manuel et la correction des questionnaires (le cas échéant). L'imputation sert surtout au traitement de la non-réponse partielle, mais aussi, à l'occasion, de la non-réponse totale. La non-réponse totale a lieu lorsqu'on n'obtient aucune information utilisable pour un enregistrement donné tandis que la non-réponse partielle a lieu lorsqu'on n'obtient qu'une partie de l'information recherchée. Après l'imputation, le fichier des données d'enquête devrait uniquement renfermer des enregistrements plausibles et cohérents à l'interne; ces enregistrements peuvent ensuite être utilisés pour l'estimation de quantités d'intérêt de la population.

### 10.2 Principes

Selon le principe de Fellegi-Holt (Fellegi et Holt, 1976), il faut changer le moins de valeurs répondues possible de telle sorte que l'enregistrement complété se conforme à chacune des règles de vérification. La détermination des champs à imputer peut se faire avant l'imputation ou au même moment que l'imputation.

L'imputation est effectuée par les personnes qui ont un accès sans restriction aux microdonnées et qui possèdent ainsi de l'information auxiliaire connue pour les unités avec et sans champs faisant l'objet d'une imputation. Cette information auxiliaire peut servir à prédire les valeurs manquantes au moyen d'un modèle de régression, à repérer des donneurs « proches » de receveurs ou à définir des classes d'imputation (p. ex. Haziza et Beaumont, 2007). Elle peut également servir directement de valeurs de remplacement pour les valeurs manquantes inconnues.

Le principe fondamental de l'imputation réside dans l'utilisation de l'information auxiliaire disponible afin d'obtenir des approximations aussi précises que possible des valeurs manquantes inconnues et ainsi produire des estimations de qualité de caractéristiques de la population. L'application de ce principe devrait donc normalement entraîner la réduction autant du biais que de la variance attribuables à l'impossibilité d'observer toutes les valeurs souhaitées.

Les bons processus d'imputation sont informatisés, objectifs et reproductibles. Ils utilisent efficacement l'information auxiliaire à leur disposition, incluent une piste de vérification pour les besoins de l'évaluation et garantissent que les enregistrements imputés sont cohérents à l'interne.

### 10.3 Lignes directrices

#### 10.3.1 Variables auxiliaires

- Le choix des variables auxiliaires servant à l'imputation – qu'on appelle aussi variables de couplage pour l'imputation par donneur – devrait être dicté principalement par l'ampleur de leur association avec les variables à imputer. Choisir ces variables en considérant l'utilisation de techniques de modélisation et consulter des spécialistes du sujet pour obtenir des informations sur les variables.
- Identifier les variables susceptibles d'agir comme variables auxiliaires pour l'imputation des données manquantes en explorant diverses sources de données (p. ex. données d'enquête en cours, données historiques, données administratives, paradonnées, etc.). Examiner la qualité et la pertinence des variables à sa disposition pour savoir lesquelles utiliser comme variables auxiliaires.
- Évaluer le type de non-réponse. Plus spécifiquement, tenter de déterminer quelles variables auxiliaires peuvent expliquer le(s) mécanisme(s) de non-réponse afin d'utiliser ces variables pour enrichir la méthode d'imputation, surtout si elles sont également associées aux variables à imputer.

- Tenir compte du type de caractéristiques à estimer (niveaux c. changements, agrégats de niveau supérieur c. petits domaines, caractéristique transversale c. caractéristique longitudinale) dans le choix des variables auxiliaires et de la stratégie d'imputation, afin de maintenir les relations d'intérêt. Par exemple, utiliser des variables auxiliaires historiques si vous vous intéressez aux changements ou des variables indiquant l'appartenance aux domaines (si possible) si vous vous intéressez à l'estimation pour des domaines.

### 10.3.2 Méthodes d'imputation et mise en œuvre

- Les méthodes d'imputation appartiennent à l'une ou l'autre de ces catégories, stochastiques ou déterministes, selon qu'elles sont fondées sur un processus aléatoire ou non. (Kalton et Kasprzyk, 1986; Kovar et Whitridge, 1995). Il existe plusieurs méthodes d'imputation déterministes : l'imputation logique, l'imputation historique (par exemple, l'imputation par valeur précédente), l'imputation par la moyenne, l'imputation par le ratio, l'imputation par régression et l'imputation par le plus proche voisin. Bien que toutes déterministes, ces méthodes se divisent encore en deux catégories. D'une part, il y a celles qui déduisent la valeur imputée en se fondant uniquement sur l'information disponible pour le non-répondant et d'autres données auxiliaires (logique et historique); d'autre part, il y a celles qui recourent aux données observées d'unités répondantes de l'enquête. Les données observées des unités répondantes peuvent être utilisées directement en opérant un transfert à partir d'un enregistrement donneur ou en utilisant des modèles paramétriques explicites (ratio et régression). Du côté des méthodes d'imputation stochastique, on trouve la méthode du « hot deck » aléatoire, l'imputation par le plus proche voisin, quand on opère une sélection aléatoire à partir de plusieurs « proches voisins », la régression avec résidus aléatoires, ainsi que toute autre méthode déterministe recourant à des résidus aléatoires.
- Il faut consacrer beaucoup d'efforts à la modélisation pour s'assurer de choisir les variables auxiliaires et le modèle d'imputation appropriés (le modèle d'imputation est un ensemble d'hypothèses relatives aux variables à imputer). Après avoir choisi son modèle d'imputation, définir la stratégie d'imputation afin qu'elle soit en accord, autant que possible, avec ce modèle. Cette démarche devrait permettre de mieux contrôler le biais et la variance dus à la non-réponse, sans compter qu'elle peut s'avérer nécessaire pour estimer correctement la variance.
- Faire en sorte que l'enregistrement imputé soit cohérent à l'interne et qu'il ressemble le plus possible à l'enregistrement rejeté à l'étape de vérification des données. Pour ce faire, imputer le moins de variables possible, afin de conserver le plus grand nombre possible de données attribuables au répondant, selon le principe de Fellegi-Holt. L'hypothèse sous-jacente est qu'un répondant est plus susceptible de commettre seulement une ou deux erreurs que plusieurs, même si cela n'est pas toujours vrai en pratique
- Dans le cas de certaines enquêtes, il faut recourir à plusieurs méthodes d'imputation selon la disponibilité de l'information auxiliaire. Habituellement, on établit une hiérarchie de méthodes en utilisant des règles pré-définies. Élaborer et tester avec soin les méthodes associées à chaque échelon hiérarchique. Limiter, autant que possible, le nombre d'échelons. Dans le même ordre d'idées, lorsqu'il est nécessaire de regrouper des classes d'imputation, élaborer et tester les méthodes d'imputation associées à chaque ensemble de classes.
- Lorsqu'on recourt à l'imputation par enregistrement donneur, tenter d'imputer les données d'un enregistrement en utilisant le moins de donneurs différents possible. Sur le plan opérationnel, cela peut équivaloir à un donneur par section du questionnaire, car il est pratiquement impossible de traiter simultanément toutes les variables d'un grand questionnaire. En outre, veiller à limiter le nombre de fois qu'un donneur particulier est utilisé pour imputer des receveurs; cela permet de contrôler la variance des estimateurs imputés. Selon les donneurs dont on dispose, cela peut signifier de faire en sorte que des actions d'imputation équivalentes aient des chances appropriées d'être retenues, ce qui permet d'éviter l'augmentation artificielle de la taille de certains groupes de la population.

- Dans le cas de grandes enquêtes, il se peut qu'il faille traiter les variables séquentiellement sur deux ou plusieurs cycles – au lieu de traiter toutes les variables au cours d'un même cycle – pour réduire les coûts informatiques de l'enquête. En outre, il se peut qu'un enregistrement contienne un grand nombre d'erreurs de réponse. Il peut donc s'avérer difficile de suivre les lignes directrices à la lettre lorsque de tels scénarios se présentent : certaines cas peuvent nécessiter plus d'un donneur (par section du questionnaire) et plus de valeurs imputées que le nombre minimal.

### 10.3.3 Incidence sur les estimations

- Il importe de garder l'information relative au processus d'imputation dans les fichiers postimputation, en vue d'évaluer l'incidence de ce processus sur les estimations et les variances. L'information dont il est question inclut des variables indiquant quelles valeurs ont été imputées et par quelle méthode, quels donneurs ont permis d'imputer les données d'un enregistrement et ainsi de suite. Conserver les valeurs non imputées et imputées des variables de l'enregistrement à des fins d'évaluation.
- Tenir compte du degré d'imputation et de son incidence lors de l'analyse des données. Même lorsque le degré d'imputation est faible, les changements apportés à un enregistrement peuvent avoir des effets considérables; c'est le cas lorsque les changements concernent de grandes unités ou lorsqu'ils sont grands et qu'ils touchent un petit nombre d'unités. En général, plus le degré et l'incidence de l'imputation sont importants, plus l'analyste doit être judicieux dans l'utilisation des données. Dans de tels cas, les analyses peuvent être trompeuses si les valeurs imputées sont traitées comme des valeurs observées.
- Il se peut que les méthodes d'imputation ne préservent pas les relations entre les variables et qu'elles exercent une influence considérable sur la distribution des données. Par exemple, les valeurs d'un domaine pourraient systématiquement augmenter pendant que diminueraient celles d'un autre domaine sans qu'aucun changement majeur soit survenu à un niveau agrégé. Cela voudrait vraisemblablement dire qu'il faut tenir compte de la variable indiquant l'appartenance aux domaines dans la stratégie d'imputation.
- Évaluer le degré d'imputation et ses effets en recourant aux outils destinés à cette tâche. Il peut s'agir, par exemple, du Système généralisé de simulation et d'imputation (GENESIS), qui permet d'imputer des données dans un environnement simulé et d'évaluer le biais et la variance d'estimateurs imputés dans des conditions particulières.

### 10.3.4 Systèmes généralisés

- Il existe des systèmes généralisés mettant en œuvre une variété d'algorithmes pour des données continues ou catégoriques. Il faut tenir compte de ces systèmes lorsqu'on élabore une méthodologie d'imputation. Généralement, ils sont conviviaux, du moins lorsque les règles de vérifications sont précisées; ils intègrent également des algorithmes permettant de déterminer quels champs imputer. Ils sont bien documentés et ils conservent des pistes de vérification permettant d'évaluer le processus d'imputation. Statistique Canada a présentement accès à deux systèmes : il s'agit du Système généralisé de vérification et d'imputation (SGVI/BANFF) (Kovar et coll., 1988; Statistique Canada, 2000a), pour l'imputation de variables économiques quantitatives et du Système canadien de contrôle et d'imputation du recensement (SCANCIR) (Bankier et coll., 1999), pour l'imputation de variables qualitatives et quantitatives.

### 10.3.5 Estimation de la variance

- Penser à utiliser des techniques permettant de bien mesurer la variance d'échantillonnage en présence d'imputation de même que la variance additionnelle due à la non-réponse et à l'imputation (Lee et al. 2002; Haziza, 2008; Beaumont et Rancourt, 2005). Il faut disposer de ces informations pour satisfaire les exigences de la Politique visant à informer les utilisateurs sur la qualité des données et la méthodologie (Statistique Canada, 2000d; voir l'annexe 2, qui reproduit cette politique). On peut se servir, à cette fin, du Système d'estimation de la variance due à la non-réponse et à l'imputation (SEVANI) développé à Statistique Canada.

- Le rapport final et les recommandations du Comité sur les mesures de la qualité (Beaumont, Brisebois, Haziza, Lavallée, Mohl, Rancourt et Trépanier, 2008) contiennent des lignes directrices supplémentaires pour l'estimation de la variance en présence d'imputation, dont il serait bon de prendre connaissance et de tenir compte avant de mettre en oeuvre toute nouvelle méthodologie ou tout nouveau logiciel.

### 10.3.6 Ressources

- Différentes ressources sont accessibles pour obtenir une formation générale en matière d'imputation ou pour plus de détails sur certains points spécifiques. Tout d'abord, on suggère de suivre le cours « 0423 : Non-réponse et imputation : Théorie et applications », offert par Statistique Canada. Le bulletin d'imputation est également une source intéressante et utile de renseignements sur le sujet. Enfin, des consultants externes, tels David Haziza et J.N.K. Rao, ainsi qu'un certain nombre de consultants internes, notamment les membres de la Division de la recherche et de l'innovation en statistique, les membres du Comité sur les mesures de la qualité et les membres du Comité sur les pratiques d'imputation sont à votre disposition pour répondre à vos questions.

## 10.4 Indicateurs de qualité

Principaux éléments de la qualité : exactitude, actualité, interprétabilité, cohérence.

Règle générale, les estimations obtenues après que la non-réponse a été observée et que l'imputation a servi à traiter cette non-réponse ne sont pas équivalentes à celles qu'on aurait obtenues si toutes les valeurs voulues avaient été observées sans erreur. La différence entre ces deux types d'estimation est appelée l'erreur de non-réponse. Le biais et la variance dus à la non-réponse (c'est-à-dire le biais et la variance dus à l'impossibilité d'observer toutes les valeurs voulues) sont deux quantités liées à l'erreur de non-réponse qui présentent habituellement un intérêt. Ces quantités inconnues, pour lesquelles nous aimerions normalement obtenir une mesure précise, sont rattachées à l'aspect « exactitude » de la qualité.

En théorie, on élimine le biais de non-réponse si la stratégie d'imputation est fondée sur un modèle d'imputation correctement spécifié possédant une bonne puissance prédictive. Un tel modèle d'imputation conduit également à une réduction de la variance due à la non-réponse. Un modèle d'imputation est correctement spécifié si, étant donné les variables auxiliaires choisies, les hypothèses sous-jacentes à ses premiers moments (habituellement la moyenne et la variance) sont valides. Le modèle est prédictif si les variables auxiliaires choisies sont fortement associées aux variables à imputer. Comme on l'a souligné dans les lignes directrices ci-dessus, les variables utilisées dans la définition de l'estimateur et les variables associées au mécanisme de non-réponse devraient être considérées comme des variables auxiliaires potentielles. L'objectif de ces lignes directrices est de faire en sorte que, étant donné les variables auxiliaires choisies, les répondants et les non-répondants se ressemblent par rapport aux variables mesurées.

Il est difficile de mesurer l'ampleur du biais de non-réponse, mais il est possible de dériver des indicateurs qui lui sont associés. Comme l'ampleur du biais de non-réponse dépend de l'adéquation du modèle d'imputation, des techniques classiques de validation du modèle, que l'on peut trouver dans les manuels conventionnels sur la régression, peuvent servir à dériver des indicateurs utiles. Par exemple, on peut utiliser des graphiques des résidus du modèle par rapport aux différentes variables auxiliaires, notamment les valeurs prédites, pour la détection des erreurs possibles de spécification du modèle. Les résidus peuvent également servir à dériver différentes statistiques. Pour la régression logistique, la statistique de test de Hosmer-Lemeshow peut fournir un indicateur utile. Ces indicateurs peuvent également être utiles pour donner une idée de la façon dont la variance due à la non-réponse a été contrôlée, particulièrement ceux qui donnent de l'information sur la force de la relation entre les variables auxiliaires et les variables à imputer.

Outre les diagnostics du modèle ci-dessus, les estimations de la variance due à la non-réponse ou les estimations de la variance totale peuvent fournir de bonnes mesures de la variabilité accrue découlant de la non-réponse pourvu que l'on puisse poser comme hypothèse que le biais de non-réponse est raisonnablement faible. La variance totale est la variance due à l'échantillonnage à laquelle on ajoute une composante de non-réponse pour refléter l'incertitude supplémentaire due à la non-réponse. Il existe de nombreuses méthodes d'estimation de la variance et certains logiciels qui tiennent compte de la non-réponse et de l'imputation. Par exemple, les estimations de la composante due à la non-réponse ou de la variance totale peuvent être obtenues au moyen du système SEVANI.

On peut utiliser d'autres indicateurs utiles pour obtenir une indication du degré d'imputation, mais ils sont plus difficiles à relier directement au biais et à la variance dus à la non-réponse. L'un de ces indicateurs est le taux d'imputation selon la variable et les domaines importants. Pour des estimations des totaux et des moyennes, un autre indicateur utile est la contribution aux estimations clés qui provient des valeurs imputées. Une contribution importante peut indiquer que le biais et/ou la variance dus à la non-réponse ne sont pas négligeables. On peut déterminer d'autres indicateurs de l'incidence de l'imputation sur les estimations finales, ce qui fournit d'autres informations pour jauger la fiabilité des estimations.

Comme on l'a souligné plus haut, on devrait faire un effort de modélisation sérieux avant de s'arrêter sur une stratégie d'imputation. Cela demande du temps et des ressources. En pratique, un équilibre doit être atteint entre le temps que l'on prend pour produire le fichier de données imputées (actualité) et la qualité du modèle d'imputation sous-jacent si l'on veut éviter de retarder indûment la diffusion des données. Lorsque l'utilisation de systèmes généralisés d'imputation est appropriée, elle est susceptible de contribuer à réduire considérablement le délai de traitement, particulièrement le délai de développement de systèmes, et ainsi faire en sorte que l'on puisse consacrer plus de temps à choisir une stratégie d'imputation appropriée.

Enfin, on devrait clairement décrire et fournir aux utilisateurs la méthodologie d'imputation utilisée ainsi que certains indicateurs et certaines mesures mentionnés plus haut. Cela assure une meilleure interprétabilité des résultats de l'enquête. Si cela est possible et pertinent, on doit envisager d'utiliser des méthodes d'imputation semblables dans les enquêtes qui recueillent le même genre de données pour ainsi assurer la cohérence entre ces enquêtes.

## Bibliographie

BANKIER, M., M. LACHANCE et P. POIRIER. 1999. « A Generic Implementation of the New Imputation Methodology », *Proceedings of the Survey Research Methods Section*, American Statistical Association, p. 548 à 553.

BEAUMONT, J.-F., F. BRISEBOIS, D. HAZIZA, P. LAVALLÉE, C. MOHL, E. RANCOURT et J. TRÉPANIÉ. 2008. *Final Report and Recommendations : Variance Estimation in the Presence of Imputation*. Rapport technique du Comité sur les mesures de la qualité de Statistique Canada.

BEAUMONT, J.-F. et É. RANCOURT, 2005. *Variance Estimation in the Presence of Imputation at Statistics Canada*. Conférence présentée à l'assemblée du Comité consultatif des méthodes statistiques de Statistiques Canada de mai 2005.

FELLEGI, I.P. et D. HOLT. 1976. « A Systematic Approach to Automatic Edit and Imputation », *Journal of the American Statistical Association*, n° 71, p. 17 à 35.

HAZIZA, D. et J.-F. BEAUMONT. 2007. « On the Construction of Imputation Classes in Surveys », *International Statistical Review*, n° 75, p. 25 à 43.

HAZIZA, D. 2008. « Imputation and Inference in the Presence of Missing data », *Handbook of Statistics*, vol. 29, D. Pfeffermann and C.R. Rao, Elsevier, (à paraître).

KALTON, G. et D. KASPRZYK, 1986. « Le traitement des données d'enquête manquantes », *Techniques d'enquête*, n° 12, p. 1 à 17.

KOVAR, J.G. et P. WHITRIDGE. 1995. « Imputation of Business Survey Data », *Business Survey Methods*, B.G. Cox et coll., New York, Wiley, p. 403 à 423.

KOVAR, J.G., J. MACMILLAN et P. WHITRIDGE. 1988. *Overview and Strategy for the Generalized Edit and Imputation System*. Document de travail de la Direction de la méthodologie de Statistique Canada noBSMD 88-007 E/F.

LEE, H., E. RANCOURT et C.-E SÄRNDAL. 2002. « Variance Estimation from Survey Data Under Single Imputation », *Survey Nonresponse*, R.M. Groves et coll., New York, Wiley, p. 315 à 328.

STATISTIQUE CANADA. 2000d. « Politique visant à informer les utilisateurs de la qualité des données et la méthodologie », *Manuel des politiques de Statistique Canada* (en ligne), [www.statcan.gc.ca/about-apercu/policy-politique/info\\_user-usager-fra.htm](http://www.statcan.gc.ca/about-apercu/policy-politique/info_user-usager-fra.htm).

STATISTIQUE CANADA. 2000a. *Description des fonctions du Système généralisé de vérification et d'imputation*. Rapport technique de Statistique Canada.

## 11 Pondération et estimation

### 11.1 Portée et objet

Une enquête vise habituellement à estimer des paramètres descriptifs de population, de même que des paramètres d'analyse, sur la base d'un échantillon sélectionné à partir d'une population d'intérêt. Parmi les exemples de paramètres figurent les statistiques descriptives simples, comme les totaux, les moyennes, les ratios et les centiles. Parmi les exemples de paramètres analytiques figurent les coefficients de régression, les coefficients de corrélation et les mesures de l'inégalité de revenu.

Dans une enquête probabiliste, un poids de sondage est associé à chaque unité échantillonnée. Ce poids peut être interprété comme le nombre d'unités typiques dans la population d'enquête que chaque unité échantillonnée représente. Des estimations peuvent être calculées à partir de ces poids ou des poids d'estimation obtenus en rajustant les poids de sondage. Parmi les rajustements courants figurent ceux qui tiennent compte de la non-réponse et qui intègrent des données auxiliaires. Voir Statistique Canada (2003).

La précision d'une estimation est un aspect important de la qualité. Cet aspect est mesuré au moyen de l'erreur type estimée (racine carrée de la variance estimée). On peut améliorer cette précision en intégrant des données auxiliaires dans le processus d'estimation.

### 11.2 Principes

Dans une enquête probabiliste, tous les éléments de la population possèdent une probabilité connue d'être sélectionnés dans l'échantillon. Ces probabilités d'inclusion tiennent compte des aspects du plan d'échantillonnage, comme la stratification, les grappes et la sélection à plusieurs degrés ou à plusieurs phases. Le poids de sondage est égal à l'inverse de la probabilité d'inclusion dans l'échantillonnage à une phase (un degré). Il s'agit du produit de l'inverse des probabilités de sélection de chaque phase (degré) dans un plan à plusieurs phases (plusieurs degrés).

En cas de non-réponse totale, l'échantillon observé est plus petit que l'échantillon initial sélectionné. Afin de compenser la non-réponse totale, on doit procéder à une repondération en rajustant les poids de sondage. Ces facteurs d'ajustement doivent être fondés sur la probabilité de réponse de chaque unité, qui peut être estimée au moyen de modèles.

Si des données auxiliaires sont disponibles, on peut améliorer la précision des estimations. L'intégration de données auxiliaires dans les processus d'estimation est appelée calage. Le calage consiste à rajuster les poids, afin que les estimations de la ou des variables auxiliaires correspondent à des totaux connus (aussi appelés totaux de contrôle). Le calage comprend des estimateurs bien connus comme l'estimateur par la régression, par le ratio et par le ratissage croisé (Deville et Särndal, 1992). Parmi les propriétés souhaitables du calage figurent les suivantes :

- cohérence des estimations entre les différentes sources;
- améliorations possibles de la précision des estimations;
- réduction possible de l'erreur due à la non-réponse totale et de l'erreur de couverture.

Les estimations sont le résultat de la somme des données multipliées par les poids de sondage ou par le poids d'estimation. Deux types d'erreurs sont associés à ces estimations : erreur d'échantillonnage et erreur non due à l'échantillonnage. L'erreur d'échantillonnage est l'erreur causée par l'observation d'un échantillon plutôt que l'ensemble de la population (Särndal et coll., 1992). Elle est mesurée par la variance d'échantillonnage, qui dépend du plan de sondage et des données auxiliaires qui sont utilisées dans le processus d'estimation. Les erreurs non dues à l'échantillonnage comprennent les erreurs de couverture (base de sondage imparfaite), les erreurs de mesure, les erreurs de traitement et les erreurs liées à la non-réponse.

Une estimation de la variance d'échantillonnage peut être calculée au moyen de méthodes comme la linéarisation de Taylor ou des méthodes de rééchantillonnage, comme le jackknife et le bootstrap. Peu importe la méthode utilisée, elle doit intégrer les propriétés du plan d'échantillonnage, comme la stratification, les grappes ou la sélection à plusieurs degrés ou phases, selon le cas.

Il est plus difficile de mesurer les erreurs non dues à l'échantillonnage. Cela peut nécessiter des données additionnelles qui ne sont généralement pas disponibles. Parmi les exemples figurent les mesures répétées, en vue d'évaluer les erreurs de mesure, et le recontact avec les non-répondants, en vue d'évaluer le biais lié à la non-réponse.

### 11.3 Lignes directrices

#### 11.3.1 Pondération

- Un poids doit être associé à chaque unité échantillonnée. Ce poids peut être le poids de sondage ou le poids d'estimation (par exemple, le poids de calage). Si on utilise uniquement le poids de sondage, l'estimateur qui en résulte s'appelle estimateur d'Horvitz-Thompson. Si des données auxiliaires sont utilisées pour le calage, l'estimateur qui en résulte est appelé estimateur de calage. Le poids lié à cet estimateur est appelé poids d'estimation ou de calage. Un poids d'estimation doit être utilisé chaque fois que le poids de sondage a été rajusté pour tenir compte de la non-réponse ou des données auxiliaires.
- Comme il est peu probable que l'on obtienne une réponse complète, on doit faire des rajustements pour tenir compte de la non-réponse afin de réduire le biais attribuable à la non-réponse. L'application de ces rajustements dans des sous-ensembles de population peut réduire le biais lié à la non-réponse. On suppose que les non-répondants ont un comportement similaire à celui des répondants de ces sous-ensembles. Ceux-ci sont délimités à partir de données auxiliaires (Lundström et Särndal, 2005) ou de modèles de propension (Eltinge et Yansaneh, 1997).
- Si des données auxiliaires sont corrélées aux variables d'intérêt, on doit envisager le calage. Ces données auxiliaires doivent au moins être disponibles pour les unités échantillonnées, et les totaux de population correspondants doivent être connus. L'estimateur de calage résultant comportera habituellement une variance plus faible que l'estimateur d'Horvitz-Thompson. En outre, les données auxiliaires pondérées s'ajouteront aux totaux de population.
- Les poids de calage peuvent être très gros ou même négatifs. Si cela se produit, des méthodes existent pour contrôler la fourchette des poids. Voir Huang et Fuller (1978) ou Deville et Särndal (1992).
- L'estimation composite doit être envisagée pour les enquêtes périodiques comportant un chevauchement d'échantillons importants entre les cycles. Il s'agit d'une méthode de calage qui traite les données des cycles précédents comme des variables auxiliaires. Pour plus de détails, voir Gambino, Kennedy et Singh (2001).
- Deux ensembles de poids peuvent être associés aux enquêtes longitudinales : les poids longitudinaux et les poids transversaux. Les poids longitudinaux se rapportent à la population pour laquelle la sélection initiale de l'échantillon longitudinal a été faite. Dans le cas de l'analyse longitudinale, ces poids doivent être rajustés pour tenir compte de l'érosion de l'échantillon. Les poids transversaux rendent compte de la population à un moment donné. Ils peuvent être utilisés pour produire des estimations ponctuelles ou des différences d'estimations ponctuelles entre les périodes.
- Si un échantillonnage double (à deux phases) a eu lieu, les poids doivent rendre compte du plan de sondage et des données auxiliaires disponibles pour la population ou pour l'échantillon de première phase.

### 11.3.2 Estimation

- Le processus d'estimation doit utiliser des poids d'estimation pour calculer les statistiques descriptives et analytiques des domaines d'intérêt. Les poids d'estimation sont équivalents aux poids de sondage, si aucun rajustement n'a été effectué. Les estimateurs correspondants de la variance doivent rendre compte du plan d'échantillonnage, des rajustements des poids de sondage, de l'imputation, ainsi que de la méthode d'estimation. Les variances peuvent être estimées au moyen de méthodes de linéarisation ou de rééchantillonnage (jackknife, répliques répétées équilibrées et bootstrap). Pour plus de détails, voir Wolter (2007).
- Les petits domaines ont trait à des sous-populations dont l'échantillon n'est pas suffisant (ou pour lesquelles il n'existe pas du tout d'échantillon) pour produire des estimations fiables. Il est par conséquent raisonnable d'intégrer les exigences relatives à ces domaines à l'étape du plan d'échantillonnage (Singh, Gambino et Mantel, 1994). Si cela n'est pas possible à l'étape de l'échantillonnage, ou si les domaines ne sont précisés qu'à une étape ultérieure, il faut envisager des méthodes spéciales d'estimation (estimateurs pour petits domaines) à l'étape de l'estimation. Ces méthodes « prennent appui » sur des régions (ou des domaines) connexes, pour réduire l'erreur quadratique moyenne de l'estimateur résultant (Rao, 2003).
- Lorsque ceci est approprié, on doit utiliser un logiciel généralisé d'estimation (Estevao et coll., 1995).

### 11.4 Indicateurs de qualité

Principal élément de la qualité : exactitude

- La qualité d'une estimation ponctuelle est habituellement décrite en termes d'exactitude et de précision. L'exactitude représente la mesure dans laquelle une valeur mesurée correspond, en moyenne, à la valeur réelle. L'exactitude d'un estimateur est évaluée du point de vue de la proximité de la moyenne de ses valeurs réalisées et du paramètre d'intérêt. À cette fin, on compare son espérance sous le plan avec le paramètre, et la différence est appelée biais. La précision, par ailleurs, rend compte du degré de correspondance entre les différentes mesures. La précision est habituellement mesurée au moyen de l'erreur d'échantillonnage : il s'agit de l'erreur qui découle de l'observation d'un échantillon plutôt que de l'ensemble de la population. Si un estimateur est sans biais, son erreur quadratique moyenne est égale à sa variance d'échantillonnage.
- Lorsqu'il existe des estimateurs sans biais et efficaces du point de vue de la variance, on devrait les utiliser. Des estimateurs légèrement biaisés peuvent être utilisés si leur efficacité, mesurée au moyen de l'erreur quadratique moyenne, est plus faible que la variance des estimateurs non biaisés correspondants.
- Le coefficient de variation est habituellement utilisé pour décrire la précision d'une estimation. Il est défini comme l'erreur type de l'estimation divisée par la valeur réelle du paramètre. Une estimation avec un coefficient de variation donné est moins précise qu'une estimation comportant un coefficient de variation plus faible. En raison de la division possible par zéro, ainsi que de problèmes d'interprétation, l'utilisation des coefficients de variation devrait être limitée aux variables d'intérêt positives. Autrement, on doit utiliser les erreurs types.
- Les estimateurs qui intègrent des données auxiliaires reposent sur le principe que les modèles entre les variables cibles et les données auxiliaires s'appliquent à toutes les unités de la population. En pratique, toutefois, il est difficile de déterminer si les hypothèses des modèles sont valides. Les estimations qui utilisent des données auxiliaires devraient être accompagnées par une description des hypothèses des modèles et par une évaluation de l'effet probable de ces hypothèses sur la qualité des estimations.

## Bibliographie

DEVILLE, J.-C. et C.E. SÄRNDAL. 1992. « Calibration Estimators in Survey Sampling », *Journal of the American Statistical Association*, n° 87, p. 376 à 382.

ELTINGE, J.L. et I.S. YANSANEH. 1997. « Méthodes diagnostiques pour la construction de cellules de correction pour la non-réponse avec application à la non-réponse aux questions sur le revenu de la U.S. Consumer Expenditure Survey », *Techniques d'enquête*, n° 23, p. 37 à 45.

ESTEVAO, V., M.A. HIDIROGLOU et C.E. SÄRNDAL. 1995. « Methodological Principles for a Generalized Estimation System at Statistics Canada », *Journal of Official Statistics*, n° 11, p. 181 à 204.

GAMBINO, J., B. KENNEDY et M.P. SINGH. 2001. « Estimation composite par régression pour l'Enquête sur la population active du Canada : évaluation et application », *Techniques d'enquête*, n° 27, p. 69 à 79.

HUANG, E. T. et W.A. FULLER. 1978. « Nonnegative Regression Estimation for Sample Survey Data », *Proceedings of the Social Statistics Section*, American Statistical Association, p. 300 à 303.

LUNDSTRÖM, S. et C.-E. SÄRNDAL. 2005. *Estimation in Surveys with Nonresponse*, New York, John Wiley and Sons.

RAO, J.N.K. (2003). *Small Area Estimation*, New York, John Wiley and Sons.

SÄRNDAL, C.E., B. SWENSSON, et J.H. WRETMAN. 1992. *Model Assisted Survey Sampling*, New York, Springer-Verlag.

SINGH, M.P., J. GAMBINO et H. MANTEL. 1994. « Les petites régions : problèmes et solutions », *Techniques d'enquête*, n° 20, p. 3 à 23.

STATISTIQUE CANADA. 2003. *Méthodes et pratiques d'enquête*, produit n° 12-587-XIF au catalogue de Statistique Canada, Ottawa.

WOLTER, K. 2007. *Introduction to Variance Estimation*, 2e édition, New York, Springer-Verlag.

## 12 Désaisonnalisation et estimation de la tendance-cycle

En mars 2000, des lignes directrices concernant la désaisonnalisation ont été rédigées à l'intention du Comité des méthodes et des normes de Statistique Canada. Le présent texte propose leur mise à jour, à la suite de l'adoption de X-12-ARIMA (Findley et al., 1998) par Statistique Canada et leur normalisation en vue de les harmoniser avec celles du Census Bureau des États-Unis (McDonal-Johnson et al., 2006a, 2006b) et d'EUROSTAT (Mazzi, 2008).

### 12.1 Portée et objet

#### 12.1.1 Désaisonnalisation

Une série chronologique est une séquence de mesures visant une variable observée au fil du temps. Dans la plupart des cas, ces mesures sont interdépendantes et c'est cette interdépendance qui présente un intérêt dans le cadre de la désaisonnalisation. Pour les besoins de la désaisonnalisation, nous supposons que la série chronologique est observée mensuellement ou trimestriellement et qu'elle est constituée de trois éléments distincts, à savoir la tendance-cycle, les effets saisonniers et de calendrier combinés et l'irrégulier. L'objectif de la désaisonnalisation est de déceler et d'estimer les effets saisonniers et de calendrier combinés pour les éliminer de la série chronologique. La série qui en résulte est alors dite désaisonnalisée et ne comprend plus que la tendance-cycle et l'irrégulier.

Les effets saisonniers sont les fluctuations infra-annuelles (mensuelles, trimestrielles) qui se répètent plus ou moins régulièrement d'année en année. Ils résultent des effets combinés des événements reliés au climat, des décisions institutionnelles ou des modes de fonctionnement qui se reproduisent avec une certaine régularité au cours de l'année. Les effets de calendrier sont liés à la composition du calendrier. Ils comprennent les effets de jours ouvrables associés aux nombres de jours de semaine contenus dans un mois, les effets des congés à occurrence variable associés aux congés à date non fixe tels que Pâques et d'autres événements prévisibles du calendrier. Les effets de jours ouvrables se produisent lorsque le niveau d'activité varie d'un jour à l'autre de la semaine. L'effet de Pâques s'assimile à la variation de niveau causée par le déplacement d'un volume d'activité d'avril à mars lorsque Pâques tombe en mars au lieu d'avril (comme c'est habituellement le cas).

La série désaisonnalisée permet d'évaluer la direction de la tendance-cycle, grâce à des comparaisons de mois en mois ou de trimestre en trimestre.

La tendance est le mouvement sous-jacent à long terme. Ce mouvement dure plusieurs années. Le cycle, qu'on appelle également cycle économique, est une oscillation quasi périodique, dont la durée excède un an, autour de la tendance à long terme. Il est caractérisé par l'alternance de périodes d'expansion et de contraction. Il est difficile d'estimer séparément la tendance et le cycle; on les étudie et on les analyse donc ensemble comme un tout appelé tendance-cycle.

L'irrégulier, ou la composante irrégulière, se manifeste dans des variations aléatoires, qui s'avèrent être des fluctuations imprévisibles causées par des événements (de toutes sortes) indépendants de la tendance-cycle, de la saisonnalité ou des effets de calendrier.

#### 12.1.2 Estimation de la tendance-cycle

La désaisonnalisation d'une série très volatile ne suffit pas nécessairement pour pouvoir tirer des conclusions sur la direction de la tendance-cycle. Le cas échéant, il est bon de poursuivre le lissage de la série désaisonnalisée pour éliminer le plus possible la composante irrégulière. L'estimation de la tendance-cycle qui en résulte doit être considérée comme de l'information auxiliaire sur la série désaisonnalisée.

### 12.1.3 X-12-ARIMA

Le fondement du programme de désaisonnalisation de Statistique Canada est la variante X-11 de la Census Method II de 1967 (Shiskin et al., 1967; Ladiray et Quenneville, 2001). En 1980, Statistique Canada ajoutait les modèles autorégressifs à moyenne mobile intégrée (ARIMA pour autoregressive integrated moving average) (Box et Jenkins, 1976) pour pouvoir projeter et rétropoler les séries avant leur désaisonnalisation. Plusieurs autres changements ont alors été instaurés pour améliorer la variante X-11. Cette nouvelle version a été appelée X-11-ARIMA (Dagum, 1980). En général, l'extension des séries, grâce aux prévisions ARIMA, a permis de réduire la révision des séries désaisonnalisées. En 1988, Statistique Canada a diffusé une version améliorée de X-11-ARIMA (Dagum, 1988) utilisable sur micro ordinateurs.

En 1998, le US Census Bureau a lancé la version X-12-ARIMA (Findley et al., 1998). Cette version s'appuie sur la régression linéaire avec erreurs ARIMA (modélisation regARIMA) pour estimer les effets de calendrier, des valeurs aberrantes additives, des changements de niveau et d'autres variables de régression prédéfinies. En outre, X-12-ARIMA permet à ses utilisateurs de définir une régression pour les effets de calendrier inhabituels ou non standards, et comprend une variante de l'algorithme TRAMO (Gomez et Maravall, 1996), pour la modélisation regARIMA automatique. D'autres fonctionnalités du programme sont décrites dans le guide de l'utilisateur (US Census Bureau, 2008). X-12-ARIMA s'adapte à diverses plateformes informatiques, ce qui inclut l'exécutable FORTRAN brut, une version C pour UNIX, et une pour FAME. Enfin, la procédure X12 (SAS Institute Inc, 2007) en SAS® assure la mise en œuvre des options les plus importantes. Sur ce plan, Statistique Canada recommande la méthode X-12-ARIMA, car X-11-ARIMA est progressivement mise hors service et ne bénéficiera plus d'aucun soutien.

## 12.2 Principes de la désaisonnalisation

La désaisonnalisation vise à éliminer les effets saisonniers et de calendrier combinés. Avant d'appliquer cette méthode, il faut donc s'assurer que de tels effets sont présents et qu'il est possible de les estimer correctement. Lorsqu'il est impossible de repérer des effets saisonniers et/ou de calendrier dans une série chronologique, cette série est réputée désaisonnalisée de facto. Les séries désaisonnalisées ne doivent pas présenter de saisonnalité résiduelle et sont généralement plus lisses que les séries brutes correspondantes.

À mesure que des données nouvelles sont accessibles, il est plus facile de bien estimer les composantes d'une série chronologique : on obtient alors des estimations révisées des valeurs désaisonnalisées antérieures qui sont plus exactes et dont il faut tenir compte. Cependant, il est préférable d'éviter les révisions trop fréquentes, car elles peuvent réduire l'utilité des données désaisonnalisées. Il faut privilégier des options de désaisonnalisation qui réduisent les révisions au minimum sans affecter pour autant la qualité globale de l'ajustement. Il faut également implanter une stratégie permettant de réduire au minimum la fréquence à laquelle les données publiées seront révisées.

Lorsqu'elles sont inappropriées, les options de désaisonnalisation peuvent fausser les résultats. Par conséquent, il faut consacrer suffisamment de temps et d'efforts à l'analyse des séries, au choix des options et à la maintenance de ces dernières. Comme chaque combinaison d'options peut produire des résultats différents d'autres combinaisons, il faut traiter les séries chronologiques comme un tout lorsque ces séries mesurent la même activité économique. Cette approche implique habituellement l'utilisation d'options d'ajustement similaires par les divers secteurs de programme concernés, pour garantir la cohérence des résultats.

## 12.3 Principes de l'estimation de la tendance-cycle

En tant que complément des séries désaisonnalisées, les estimations de la tendance-cycle peuvent révéler la direction de la tendance à court terme (durant l'année courante). Lorsque de nouveaux points de donnée sont ajoutés aux séries, on peut calculer plus précisément les estimations antérieures de la tendance-cycle; ces estimations font donc l'objet de révisions. Les estimations de la tendance-cycle sont sensibles à la phase courante du cycle économique (point de retournement, récession, reprise ou expansion); la fiabilité des estimations courantes de la tendance-cycle dépend donc de la proximité d'un point de retournement,

ainsi que de l'amplitude du cycle.

Les estimations de la tendance-cycle doivent correspondre en tous points à la série désaisonnalisée diffusée. Si les valeurs désaisonnalisées sont gelées dans une base de données, la tendance cycle doit être estimée d'après la série désaisonnalisée telle que présentée dans cette base. De même, si la série désaisonnalisée a subi d'autres ajustements, telles qu'une agrégation ou bien une réconciliation, la tendance-cycle doit être estimée d'après la série agrégée ou réconciliée.

#### 12.4 Lignes directrices concernant la désaisonnalisation

- Avant de désaisonnaliser une série pour la première fois, s'assurer que la composante saisonnière est identifiable et qu'on peut l'estimer correctement.
- Si une série est exempte de saisonnalité ou d'effet de calendrier, n'appliquer aucun traitement; la série est alors réputée désaisonnalisée de facto.
- Une série désaisonnalisée ne doit contenir ni saisonnalité résiduelle, ni effet de calendrier résiduel.
- Afin d'identifier et de bien estimer les effets saisonniers et les effets de calendrier, il est recommandé d'utiliser des données couvrant une période de 10 à 15 ans. La période minimale est de cinq (5) ans pour estimer correctement un profil saisonnier et de sept (7) ans pour les effets de calendrier tels que les effets de jours ouvrables et de jours fériés mobiles (les congés à occurrence variable).
- Il est recommandé d'utiliser la modélisation regARIMA pour calculer les facteurs de correction des effets de calendrier et les corrections temporaires comme celles des valeurs aberrantes additives connues ou des changements de niveau. En général, il faut user du même modèle regARIMA pour extrapoler les séries, afin de réduire les révisions au sein de la série désaisonnalisée.
- Les options de désaisonnalisation sont multiples; les plus importantes concernent la sélection du modèle de décomposition, ainsi que la spécification d'un modèle regARIMA et de la longueur des filtres saisonniers et de la tendance-cycle. Les processus de sélection automatique de X-12-ARIMA peuvent être utilisés pour une initialisation de ces options. Si l'on dispose de suffisamment de temps ou que les attentes en matière de qualité sont grandes, la sélection automatique doit être examinée en s'appuyant sur d'autres statistiques, sur les connaissances préalables des experts et sur l'analyse graphique.
- Pour chaque série, il faut examiner périodiquement les options de désaisonnalisation, afin de s'assurer qu'elles demeurent applicables et appropriées et d'accroître la précision. À moins de motifs valables, il faut éviter de modifier les options de désaisonnalisation retenues entre deux révisions planifiées.
- Si les principales options de désaisonnalisation doivent généralement demeurer fixes, entre les révisions, les facteurs d'ajustement et les paramètres du modèle regARIMA doivent quant à eux être actualisés; autrement dit, il faut les recalculer en se servant de tous les points de données disponibles. Il peut y avoir exception à la règle, lorsqu'on sait que les observations les plus récentes ont déjà été soumises à des révisions majeures. Le cas échéant, les facteurs (prévus) pour l'année à venir peuvent s'avérer plus appropriés.
- La désaisonnalisation des séries agrégées (ou composées) comprenant plusieurs séries-composantes peut être indirecte – les composantes désaisonnalisées sont agrégées pour former la série agrégée désaisonnalisée – ou directe – la série agrégée est désaisonnalisée indépendamment. La méthode directe peut engendrer des écarts entre la série agrégée et les composantes agrégées après la désaisonnalisation. Au besoin, appliquer une méthode de réconciliation afin de faire concorder la série agrégée désaisonnalisée directement avec ses composantes désaisonnalisées, sans modifier les composantes non désaisonnalisées dans la mesure du possible.
- Qu'il résulte de l'approche indirecte ou de l'approche directe, l'agrégat ne devrait contenir aucune saisonnalité résiduelle et devrait être relativement lisse. L'approche directe avec réconciliation est la plus indiquée lorsqu'on accorde plus d'importance à la série agrégée qu'à ses composantes ou bien que les composantes présentent des composantes saisonnières fort similaires. En général, il est approprié de recourir à l'ajustement indirect quand les séries-composantes présentent des profils saisonniers très différents et que plusieurs d'entre elles peuvent être désaisonnalisées individuellement.

- Forcer les totaux annuels des données désaisonnalisées à être égaux à ceux de la série originale (ou de la série originale corrigée des effets de calendrier) est une approche rarement justifiée d'un point de vue théorique, mais qui pourrait être adoptée quand il est nécessaire d'assurer la cohérence avec des étalons externes, telles qu'au Système de comptabilité nationale, ou de réconcilier un agrégat et ses composantes.

#### **12.4.1 Révisions des données désaisonnalisées**

- La publication des données désaisonnalisées révisées doit suivre une politique de révision officielle et s'aligner sur le calendrier de diffusion des données non désaisonnalisées.
- Lorsqu'on emploie un facteur saisonnier actualisé, il n'est pas nécessaire de réviser les estimations désaisonnalisées en remontant en arrière de plus d'une période. Il peut y avoir des exceptions à cette règle, lorsqu'on utilise des observations préliminaires : il est alors recommandé de réviser les facteurs saisonniers chaque fois qu'on révisé les données originales. Chaque année, il faut réviser les valeurs désaisonnalisées des trois dernières années dès que les données du premier mois (trimestre) de l'année suivante sont disponibles. Si les valeurs désaisonnalisées sont générées grâce aux facteurs de désaisonnalisation (prévus) de l'année à venir, cette révision annuelle doit s'appliquer aux quatre dernières années.

#### **12.4.2 Estimation de la tendance-cycle**

- Il faut appliquer la méthode d'estimation de la tendance-cycle à la série désaisonnalisée publiée pour s'assurer que la ligne de tendance soit centrée sur la série désaisonnalisée. La méthode de Dagum (1996) ou une adaptation/variante appropriée de cette méthode est recommandée pour l'estimation de la tendance-cycle.
- Informer les utilisateurs du fait que les dernières estimations de la tendance-cycle (et surtout de la toute dernière) peuvent subir des révisions lorsque s'ajoute un point de donnée supplémentaire. Il est possible de signaler cette variabilité (celle des estimations qui se trouvent vers l'extrémité de la série) en se servant d'une ligne pointillée sur le graphique de tendance, par exemple, ou en publiant une note d'information en même temps que les données.
- Il faut réviser les estimations de la tendance en remontant aussi loin que pour la révision des estimations désaisonnalisées. Par ailleurs, lorsqu'on révisé une série mensuelle type, il faut ajouter trois mois (deux pour une série trimestrielle) durant l'année et six mois (deux pour une série trimestrielle) au moment de la révision annuelle.

#### **12.4.3 Présentation des données et accès aux données**

- Il faut calculer les taux de croissance et les variations – d'un mois à l'autre (ou d'un trimestre à l'autre) – en se servant des données désaisonnalisées et les utiliser avec prudence lorsque les séries chronologiques sont hautement instables. Lorsqu'on compare le même mois, d'une année à l'autre, il faut se servir des données corrigées des effets de calendrier ou, en l'absence d'effets de calendrier, des données brutes.
- Les utilisateurs doivent avoir accès à toute la série historique brute, à la série désaisonnalisée et, sur demande, aux options de désaisonnalisation.

#### **12.4.4 Mise en œuvre**

Pour obtenir de l'aide concernant l'interprétation et la mise en œuvre de ces lignes directrices, s'adresser au Centre de recherche et d'analyse en séries chronologiques (CRASC), Division des méthodes d'enquête auprès des entreprises.

## 12.5 Indicateurs de qualité

Les indicateurs suivants peuvent servir à déterminer si une série contient une composante saisonnière :

- Un simple graphique temporel et un graphique année-sur-année. Ces outils permettent d'inspecter la série en vue de déceler des profils saisonniers et de repérer – visuellement – d'autres perturbations;
- Diverses statistiques, telles que les deux tests de Fisher pour la saisonnalité stable et la saisonnalité évolutive, ainsi que les graphiques de spectre décrits dans (Ladiray et Quenneville, 2001; Findley et al., 1998).

La détection de la saisonnalité résiduelle peut se faire grâce aux tests ci-dessus appliqués aux données désaisonnalisées ou au moyen d'autres tests (Ladiray et Quenneville, 2001).

Les statistiques permettant d'évaluer la signification des composantes regARIMA estimées et la qualité globale du modèle ajusté sont décrites dans plusieurs manuels traitant du sujet.

Pour quantifier les révisions, on peut utiliser des statistiques sommaires sur les révisions historiques des niveaux et des changements dans la série désaisonnalisée. Dans le cadre de la désaisonnalisation, les révisions augmentent généralement la précision, car elles reposent sur des observations qui n'étaient pas disponible initialement.

### Bibliographie

BOX, G.E.P. et G.M. JENKINS. 1976. *Time Series Analysis, Forecasting and Control*, San Francisco, Holden Day.

DAGUM, E.B. 1980. *La méthode de désaisonnalisation X-11-ARIMA*, publication n° 12-564F au catalogue de Statistique Canada, Ottawa.

DAGUM, E.B. 1988. *The X11ARIMA/88 Seasonal Adjustment Method - Foundations and User's Manual*, Time Series Research and Analysis Division. Rapport technique de Statistique Canada.

DAGUM, E.B. 1996. « A New Method to Reduce Unwanted Ripples and Revisions in Trend-Cycle Estimates from X-11-ARIMA », *Survey Methodology*, n° 22, p. 77 à 83.

FINDLEY, D.F., B.C. MONSELL, W.R. BELL, M.C. OTTO et B.C. CHEN. 1998. « New Capabilities of the X-12-ARIMA Seasonal Adjustment Program (with Discussion) », *Journal of Business and Economic Statistics*, n° 16, p. 127 à 177.

GOMEZ, V. et A. MARAVALL. 1996. *Programs TRAMO (Time Series Regression with Arima Noise, Missing Observations, and Outliers) and SEATS (Signal Extraction in Arima Time Series). Instructions for the User*. Document de travail 9628 du Département de recherche de la Banque d'Espagne. Voir aussi le document en ligne : <http://www.bde.es/servicio/software/tramo/summprogs.pdf>.

LADIRAY, D. et B. QUENNEVILLE. 2001. « Seasonal Adjustment with the X-11 Method », *Lecture Notes in Statistics*, n° 158, New York, Springer-Verlag.

MAZZI, G.L. 2008. *ESS Guidelines on Seasonal Adjustment* (en ligne), [http://epp.eurostat.ec.europa.eu/pls/portal/docs/PAGE/PGP\\_RESEARCH/PGE\\_RESEARCH\\_04/ESS%20GUIDELINES%20ON%20SA.PDF](http://epp.eurostat.ec.europa.eu/pls/portal/docs/PAGE/PGP_RESEARCH/PGE_RESEARCH_04/ESS%20GUIDELINES%20ON%20SA.PDF).

MCDONAL-JOHNSON, K.M., B. MONSELL, R. FRSCINA et R. FELDPAUSH. 2006. *Seasonal Adjustment Diagnostics, Census Bureau Guideline, Version 1.0.*, Washington, US Census Bureau.

MCDONAL-JOHNSON, K.M., B. MONSELL, R. FRSCINA et R. FELDPAUSH. 2006. *Supporting Document A, Seasonal Adjustment Diagnostics Checklists, Census Bureau Guideline, Version 1.0.*, Washington, US Census Bureau.

SAS INSTITUTE INC. 2007. *SAS/ETS® User's Guide 9.2*, Cary NC, SAS Institute Inc.

SHISKIN, J., A.H. YOUNG et J.C. MUSGRAVE. 1967. *The X-11 Variant of the Census Method II Seasonal Adjustment*. Document technique n° 15 du Bureau of the Census, U.S. Department of Commerce.

US CENSUS BUREAU. 2008. *X-12-ARIMA Reference Manual*, Statistical Research Division, Census Bureau.

## 13 Étalonnage et techniques connexes

### 13.1 Portée et objet

Les programmes statistiques s'appuient souvent sur deux sources de données pour mesurer une même variable cible, à savoir une mesure fréquente destinée à obtenir une estimation précise du changement d'une période à l'autre et une mesure moins fréquente axée sur l'estimation précise du niveau. Sans perte de généralité, à partir d'ici, nous donnerons à la série dont les observations sont les plus fréquentes le nom de série infra-annuelle, tandis que nous utiliserons la série dont la fréquence est plus faible comme série repère (les étalons) et considérerons qu'il s'agit d'une série annuelle.

Par étalonnage, nous entendons les techniques utilisées pour s'assurer de la cohérence entre les séries chronologiques ayant trait à une variable cible mesurée à diverses fréquences, par exemple, infra-annuellement et annuellement. L'étalonnage consiste à imposer le niveau de la série repère, tout en minimisant dans la mesure du possible les révisions aux changements observés dans la série infra-annuelle. Par conséquent, les taux de croissance de la série étalonnée concordent avec ceux des repères (étalons). Dans certaines situations, l'étalonnage peut améliorer la précision et l'actualité du produit statistique.

L'étalonnage non contraignant, l'interpolation, la distribution temporelle, la calendrialisation, le raccordement et la réconciliation sont des techniques connexes qui sont fondées sur des principes et des lignes directrices méthodologiques similaires à ceux de l'étalonnage. L'étalonnage non contraignant est utilisé quand la série des repères peut aussi être révisée. L'interpolation, qui est l'estimation des termes intermédiaires entre des valeurs connues, peut aussi être utilisée pour étalonner les séries sur les stocks. La distribution temporelle est la désagrégation de la série des repères en observations plus fréquentes. La calendrialisation est un cas particulier de la distribution temporelle. Le raccordement est utilisé pour unir différents segments d'une série chronologique en une série chronologique unique cohérente. La réconciliation est utilisé pour imposer des contraintes additives transversales aux composantes d'un système de séries chronologiques. Des renseignements plus détaillés sur ces techniques figurent dans Dagum et Cholette (2006).

L'étalonnage dans le contexte des séries chronologiques ne doit pas être confondu avec les ajustements de la pondération qui peuvent être appliqués à l'étape de l'estimation pour les besoins du calage.

### 13.2 Principes

Il faut s'assurer, au stade de la conception, que les différences d'ordre conceptuel, méthodologique et opérationnel entre les deux sources de données sont aussi peu nombreuses que possible. Les différences entre les séries doivent être examinées minutieusement et bien comprises, après quoi il peut être décidé de manière éclairée si la série doit être publiée telle quelle ou doit être étalonnée afin d'être certain que tous les chiffres concordent. Dans le premier cas, les différences doivent être expliquées aux utilisateurs conformément à la Politique visant à informer les utilisateurs sur la qualité des données et la méthodologie (Statistique Canada, 1998).

Quand les sources de données sont conçues de manière compatible ou que des contraintes externes imposent la cohérence de tous les chiffres, les méthodes d'étalonnage peuvent et — du point de vue statistique — doivent être utilisées. Dans les situations types, il est supposé implicitement que la série infra-annuelle est moins fiable que la série de données annuelles. De par sa nature, le processus d'étalonnage donnera lieu à divers ensembles de révisions des données infra-annuelles et, par conséquent, il est nécessaire d'être disposé à faire les révisions.

Toutes les techniques connexes sont fondées sur les mêmes principes : il faut d'abord comprendre les hypothèses sous-jacentes et confirmer l'applicabilité des méthodes, le plus fréquemment en effectuant une analyse détaillée des données avant et après l'utilisation de la méthode.

### 13.3 Lignes directrices

- Avant d'envisager l'étalonnage, examiner, décrire et quantifier les écarts entre les deux sources de données. Ces écarts doivent être, dans la mesure du possible, réduits au minimum à l'étape de la conception.
- Avant d'envisager l'étalonnage, examiner les différences entre les microdonnées pour les unités d'échantillonnage communes, s'il en existe. Si des corrections sont apportées, elles doivent respecter la nature chronologique des données. Dans le cas de données infra-annuelles, les corrections pourraient viser à améliorer l'exactitude du changement d'une période à l'autre; pour les données annuelles, les corrections doivent tenir compte de l'exactitude du niveau ainsi que de l'exactitude du changement d'une année sur l'autre.
- Ne pas perdre de vue que la conception de la série annuelle pourrait ne pas être compatible avec les objectifs de l'étalonnage. Pour ce dernier, l'enquête annuelle doit fournir à la fois une mesure précise du niveau annuel et une mesure précise du changement annuel, puisqu'ils seront imposés à la série étalonnée.
- Ne pas procéder à l'étalonnage quand les valeurs annuelles sont moins fiables que les sommes annuelles des valeurs de la série infra-annuelle. Le cas échéant, le fait d'imposer les valeurs repères annuelles produira essentiellement une série étalonnée moins fiable.
- Si les sources de données sont conçues différemment, n'envisager l'étalonnage que si de fortes contraintes externes rendent nécessaire la cohérence complète des chiffres. Ne pas oublier que la cohérence résultante pourrait être obtenue au prix d'une réduction de la précision.
- L'étalonnage entraînera des révisions aux données infra-annuelles. N'envisager l'étalonnage que si l'accroissement de la cohérence réduit fortement la confusion parmi les utilisateurs ou que la plus grande précision due à une série annuelle de haute qualité l'emporte sur le fardeau des révisions répétées.
- Procéder à l'étalonnage dans le contexte de la désaisonnalisation lorsqu'il existe des écarts non souhaités entre les totaux annuels de la série brute et les totaux annuels correspondant de la série désaisonnalisée. Au besoin, la série désaisonnalisée peut être étalonnée sur les totaux annuels calculés d'après la série brute.
- Utiliser une méthode d'étalonnage appropriée, telle que les techniques fondées sur la régression décrite dans Dagum et Cholette (2006) ou l'une des diverses méthodes améliorées de Denton décrites dans le Manuel des comptes nationaux trimestriels du Fonds monétaire international (Bloem et coll., 2001). Éviter les simples techniques de répartition proportionnelle, parce que celles-ci introduisent des ruptures entre les années (auxquelles ont donné le nom de « problème de l'escalier » ou step-problem).
- Comprendre les hypothèses sous-jacentes quand les observations les plus récentes de la série infra-annuelle ne possèdent pas de valeur annuelle correspondante — parce que l'année est incomplète ou que la donnée annuelle n'est pas encore disponible. On utilisera dans les méthodes d'étalonnage une projection implicite ou explicite de la valeur annuelle suivante; les projections peuvent être basées sur les données historiques de court terme ou de long terme ou sur des considérations externes.
- Lors de la mise en œuvre de méthodes d'étalonnage ou de méthodes connexes, envisager l'utilisation d'un logiciel généralisé. Cela limite les erreurs de programmation et réduit le coût et le temps de développement. À Statistique Canada, la Méthodologie et l'Informatique offrent un soutien, surtout pour les logiciels suivants : Le logiciel maison SAS Proc Benchmarking pour l'étalonnage, la distribution temporelle et le raccordement; le logiciel maison SAS Proc TSraking – pour la réconciliation; SAS Proc Expand – pour l'interpolation; programme X-12-ARIMA du US Bureau of the Census, SAS Proc X12 ou SAS Proc >ARIMA – pour les méthodes d'inférence statistique appliquées aux séries chronologiques.
- Une aide concernant l'interprétation et la mise en œuvre de ces lignes directrices peut être obtenue auprès du Centre de recherche et d'analyse en séries chronologiques (CRASC), Division des méthodes d'enquête auprès des entreprises.

### 13.4 Indicateurs de qualité

Les indicateurs qui suivent peuvent être utilisés pour décrire et quantifier les écarts :

- énoncés descriptifs des aspects conceptuels : période de déclaration de la source annuelle; définition des variables mesurées et de la population cible, etc.;
- énoncés descriptifs des aspects opérationnels : base de sondage, processus de collecte des données, etc.;
- énoncés descriptifs des aspects méthodologiques : échantillonnage, utilisation et source de données administratives, etc.;
- au besoin, descriptions quantitatives des écarts : dénombrements et dénombrements pondérés par période de déclaration pour estimer l'effet de la non-calendrialisation des données annuelles, différences entre les données administratives, erreurs d'échantillonnage des deux estimations annuelles et variations annuelles correspondantes, etc.

Pour des renseignements plus détaillés et une étude de cas, consulter Yung et coll. (2008).

L'application des méthodes d'étalonnage et l'analyse approfondie des résultats peuvent aussi fournir un bon indicateur du caractère approprié de la méthode. Les différences annuelles, la révision de la série en se basant sur les ratios de la série étalonnée sur la série originale (Bi-ratios en anglais) et la révision des taux de croissance peuvent toutes être étudiées au moyen de graphiques et de statistiques sommaires.

#### Bibliographie

BLOEM, A. M., R. J. DIPPELSMAN et N. Ø. MÆHEL. 2001. *Quarterly National Accounts Manual, Concepts, Data Sources and Compilation*, Washington DC, International Monetary Fund.

DAGUM, E.B. et P.A. CHOLETTE. 2006. « Benchmarking, Temporal Distribution and Reconciliation Methods for Time Series », *Lecture Notes in Statistics*, n° 186, New York Springer. 410 p.

STATISTIQUE CANADA. 1998. « Politique d'information des répondants aux enquêtes », *Manuel des politiques de Statistique Canada* (en ligne), [http://icn-rci.statcan.ca/10/10c/10c\\_001\\_f.htm](http://icn-rci.statcan.ca/10/10c/10c_001_f.htm).

YUNG, W., B. BRISEBOIS, C. TARDIF, G. KUROMI et C. RONDEAU. 2008. *Should Sub-Annual Surveys be Benchmarked to their Annual Counterparts? A Case Study of Manufacturing Surveys*, Ottawa, Statistique Canada, Document de travail BSMD-2008-001. Ottawa, Ontario.

## 14 Évaluation de la qualité des données

### 14.1 Portée et objet

Chaque étape d'une enquête est marquée par le souci d'appliquer des méthodes saines capables de garantir la qualité des données; d'ailleurs, tous les chapitres de ce document en témoignent. L'évaluation de la qualité des données sert à déterminer dans quelle mesure le produit final satisfait aux objectifs initiaux de l'activité statistique, notamment sous l'angle de la fiabilité (exactitude, actualité et cohérence). De plus, elle facilite l'interprétation, pour les utilisateurs, des résultats de l'enquête, et permet à l'organisme statistique d'accroître la qualité de ses enquêtes.

Généralement parlant, il existe deux types d'évaluation de la qualité des données :

- Certification ou validation : les données sont analysées avant leur diffusion officielle, afin d'éviter les erreurs évidentes et d'éliminer les données de piètre qualité; à cette étape de l'enquête, on privilégie la comparaison des données à des sources de données externes ou auxiliaires
- Étude des sources d'erreur : cette démarche permet généralement d'obtenir des renseignements quantitatifs sur les sources précises des erreurs présentes dans les données.

L'évaluation de la qualité des données repose sur des indicateurs établis à chacune des étapes de l'enquête. Ces méthodes d'évaluation sont exposées dans différentes sections du présent document. Plusieurs considèrent qu'il est impossible d'établir un indice de qualité unidimensionnel et unique; or, on peut résumer les différents indicateurs de la qualité et les comparer, sur le plan de leur importance relative et de leurs conséquences.

### 14.2 Principes

Il est essentiel d'évaluer la qualité des données, pour déterminer dans quelle mesure elles sont pertinentes et représentatives. Les utilisateurs sont rarement capables d'exécuter cette tâche. La plupart du temps, elle incombe à l'organisme statistique, qui s'en acquitte et en communique les résultats aux utilisateurs, le plus rapidement possible et sous la forme la plus pratique possible.

Outre la pertinence des données, l'évaluation de la qualité permet de vérifier s'il existe un lien entre certains types d'erreurs et certaines étapes du processus d'enquête. Cette démarche permet – lorsque c'est nécessaire – d'augmenter la qualité d'une prochaine édition de l'enquête ou d'enquêtes similaires.

Les évaluations de la qualité des données menées à Statistique Canada doivent satisfaire aux exigences minimales énoncées dans la Politique visant à informer les utilisateurs sur la qualité des données et la méthodologie (Statistique Canada, 2000). Afin d'y satisfaire, on doit mesurer ou évaluer les erreurs de couverture, les taux de réponse et d'imputation, ainsi que les erreurs d'échantillonnage en fonction des caractéristiques clés (pour une enquête-échantillon).

### 14.3 Lignes directrices

#### 14.3.1 Conception

- Déterminer l'ampleur requise, pour l'évaluation de la qualité des données, en fonction du programme ou du produit soumis à l'évaluation. Pour ce faire, il faut considérer les facteurs suivants : usages et utilisateurs des données; risque d'erreur et incidence des erreurs sur l'usage des données; variation de la qualité, au fil du temps; coût de l'évaluation par rapport au coût total du programme; amélioration de la qualité; augmentation de l'efficacité et de la productivité; degré d'utilité des mesures, pour les utilisateurs, et degré de facilité à les interpréter; possibilités que l'enquête soit répétée.
- L'information requise par les évaluations de la qualité des données est souvent recueillie durant le processus d'enquête. Il faudrait donc inclure le plan d'évaluation dans le plan d'enquête; les rapports sur la qualité des données devraient également figurer dans le calendrier de diffusion de l'enquête.

- S'assurer que les résultats de l'évaluation de la qualité sont valides et qu'ils sont assez récents pour contribuer à l'amélioration des données diffusées. Si cela se révèle impossible, s'assurer, à tout le moins, que les résultats sont assez récents, afin que les utilisateurs puissent analyser les données plus facilement et que les concepteurs des enquêtes puissent améliorer le plan des éditions subséquentes de l'enquête ou d'enquêtes similaires.

#### **14.3.2 Exécution**

- Évaluer la qualité en se fondant sur l'opinion d'experts ou sur une analyse subjective lorsque l'évaluation de la qualité des données ne peut fournir des mesures quantitatives, à cause de la nature du produit, de l'utilisateur, des contraintes de temps, du coût ou de la faisabilité technique.
- En ce qui concerne les enquêtes ou les activités statistiques répétées, il se peut qu'il ne soit pas nécessaire, ou possible, d'évaluer (dans le détail) la qualité de manière constante. Il faut néanmoins le faire périodiquement, pour s'assurer que les activités statistiques atteignent leurs objectifs (sans attendre qu'un problème survienne).
- Impliquer les utilisateurs des résultats d'évaluation – qu'ils viennent d'un organisme statistique ou d'ailleurs – dans la définition des objectifs du programme d'évaluation de la qualité des données. Faire de même pour le processus d'évaluation, lorsque les circonstances le permettent.
- Pour pouvoir mener l'ensemble de ces évaluations, le gestionnaire – ou l'équipe de gestion – de l'enquête doit préalablement identifier les normes qu'il veut respecter et les objectifs qu'il désire atteindre.

#### **14.3.3 Certification ou validation**

- Certifier ou valider l'information statistique chaque fois qu'il est possible ou approprié de le faire.
- Veiller à ce que la certification ou la validation interroge les données au lieu de les rationaliser. Pour ce faire, il est bon d'impliquer des analystes qui n'ont pas pris part à la production des données dans le processus.
- Vérifier la cohérence des données par rapport à des sources de données externes, comme d'autres enquêtes, d'autres éditions de la même enquête ou des données administratives.
- Vérifier la cohérence interne, en calculant des ratios dont les limites probables sont connues (ratio hommes-femmes, valeurs moyennes des biens, etc.), par exemple.
- Analyser la contribution (sur un plan individuel) des grandes unités aux estimations globales (généralement dans le cadre d'enquêtes-entreprises).
- Examiner et interpréter les indicateurs de qualité des données présentés dans les autres sections de ce document et les comparer aux objectifs de production.
- Organiser des rencontres de rétroaction avec le personnel affecté à la collecte et au traitement des données.
- Mandater des spécialistes de l'externe qui sont familiarisés avec l'enquête, pour qu'ils vérifient si ses résultats sont plausibles et qu'ils rédigent un rapport sur les travaux en cours, avant la diffusion des résultats.

#### **14.3.4 Examen des sources d'erreurs**

- Examiner fréquemment les sources d'erreurs, dans le cas des programmes statistiques annuels ou pluriannuels, et occasionnellement, dans le cas des programmes repris à intervalles rapprochés.
- Évaluer, entre autres, les erreurs de couverture et d'échantillonnage, les erreurs attribuables à la non-réponse, ainsi que les erreurs de mesure et de traitement, à la lumière des analyses d'autres étapes de l'enquête.

## 14.4 Indicateurs de qualité

Chacune des sections qui précèdent présente des indicateurs de la qualité adaptés aux caractéristiques des sujets qu'elle aborde. Or, il est bon de recourir à des indicateurs portant sur l'ensemble du projet au lieu de se concentrer sur l'une de ses étapes. Bien souvent, on ne peut pas mesurer ces indicateurs avant la diffusion du produit et il arrive même que l'attente se prolonge au-delà de cette étape. Lorsque c'est le cas, les indicateurs sont exclus de la documentation connexe; ils permettent toutefois de déterminer, approximativement, quelle serait la qualité d'une nouvelle édition du programme ou d'un programme similaire. Voici des exemples de ce genre d'indicateur :

### 14.4.1 Actualité

- Combien de temps le projet a-t-il demandé, du lancement à la clôture? Combien de temps s'est-il écoulé entre la conclusion du projet et la période de référence?
- Combien de temps s'est-il écoulé entre l'étape de la collecte des données et le moment où les estimations relatives aux caractéristiques principales sont devenues accessibles?

### 14.4.2 Pertinence

- Les résultats de l'enquête répondent-ils à ses objectifs et aux besoins analytiques de la collectivité?
- Est-il possible que des populations aient été omises ou que des questions écartées à cause de certaines étapes ou contraintes opérationnelles?
- Comparer les résultats prévus et ceux qui ont été obtenus; justifier d'éventuels écarts.

### 14.4.3 Intelligibilité

- Vérifier que la documentation est complète.
- Calculer le nombre de demandes de renseignements, plus particulièrement celles qui visent à clarifier des informations. Cette mesure s'avère encore plus importante, dans le cas des enquêtes répétées. Tenter de déterminer si ces demandes révèlent une faille dans les fondements du cadre théorique ou de la documentation de l'enquête.

### 14.4.4 Exactitude

- Le projet a-t-il permis de produire des estimations de qualité satisfaisante, et ce, pour chacun des domaines et pour chacune des variables que l'on a prévu étudier? La réponse à cette question peut s'exprimer en pourcentage; par exemple, 86 % des estimations prévues concordent avec les objectifs concernant les c.v.
- Dans le cas d'enquêtes répétées, comparer les estimations principales et leur niveau de qualité (c.v.) aux résultats des versions antérieures. Il faut s'assurer de pouvoir expliquer d'éventuels changements. Traduire les changements relatifs aux c.v. sous forme de pourcentage (plus ou moins élevé que les pourcentages des éditions précédentes). Il est possible de produire des statistiques similaires pour les taux d'imputation, d'erreur, etc.
- Dans le cas d'enquêtes non répétées, envisager de comparer les estimations de l'enquête à des données administratives connexes ou aux estimations d'autres enquêtes. Évidemment, les populations de ces sources peuvent différer les unes des autres; il se peut donc qu'il faille justifier de tels écarts.

### 14.4.5 Cohérence

- Comparer les résultats de l'enquête à ceux d'éditions antérieures; tenter de repérer d'éventuelles différences, d'en identifier les causes et de les quantifier (p. ex. : « l'enquête inclut maintenant les territoires; s'ils avaient été exclus de l'enquête, comme c'était le cas dans les éditions précédentes, les estimations nationales auraient été de 31,4 % plutôt que de 31,5 % »).
- Comparer les résultats de l'enquête aux résultats de sources externes et tenter d'expliquer d'éventuels écarts.

#### 14.4.6 Accessibilité

- Décrire les types et les formats des produits de l'enquête.
- Indiquer le nombre de fois qu'un produit d'enquête a été consulté sur un site Internet accessible au public.
- Mentionner si les données de l'enquête sont stockées dans un fichier de microdonnées à grande diffusion, si certains de leurs produits sont gratuits et si elles sont accessibles dans des centres de données de recherche.

#### Bibliographie

BIEMER, P., R.M. GROVES, N.A. MATHIOWETZ, L. LYBERG et S. SUDMAN. 1991. *Measurement Errors in Surveys*, New York, Wiley.

BIEMER, P. et L. LYBERG. 2003. *Introduction to Survey Quality*, New York, Wiley.

FULLER, W. 1987. *Measurement Error Models*, New York, Wiley.

LESSLER, J.T. et W.D. KALSBECK. 1992. *Nonsampling Errors in Surveys*, New York, Wiley.

LYBERG, L., P. BIEMER, M. COLLINS, E. DE LEEUW, C. DIPPO, N. SCHWARZ, et D. TREWIN. 1997. *Survey Measurement and Process Quality*, New York, Wiley.

STATISTIQUE CANADA. 2000d. « Politique visant à informer les utilisateurs sur la qualité des données et la méthodologie », *Manuel des politiques de Statistique Canada*, [http://www.statcan.gc.ca/about-aperçu/policy-politique/info\\_user-usager-fra.htm](http://www.statcan.gc.ca/about-aperçu/policy-politique/info_user-usager-fra.htm).

STATISTIQUE CANADA. 2002. « Le cadre d'assurance de la qualité de Statistique Canada 2002 », produit n° 12-586-XIF au catalogue de Statistique Canada, Ottawa.

STATISTIQUE CANADA. 2003. *Méthodes et pratiques d'enquête*, produit n° 12-587-XIF, Ottawa.

## 15 Contrôle de la divulgation

### 15.1 Portée et objet

Le contrôle de la divulgation désigne les mesures visant à protéger les données dans le respect des exigences en matière de confidentialité. L'objectif consiste à s'assurer que les dispositions régissant la protection de la confidentialité sont respectées tout en préservant le plus possible l'utilité des données produites. Le programme vigilant de contrôle de la divulgation et de protection de la confidentialité de Statistique Canada contribue grandement à la qualité des données; en effet, les taux de réponse élevés dans les enquêtes du Bureau et la confiance que le public place dans l'organisme en sont tributaires dans une large mesure.

### 15.2 Principes

Les principes qui sous-tendent les activités de contrôle de la divulgation sont presque exclusivement régis par les dispositions de la Loi sur la statistique (1970, S.R.C. 1985, c. S19), plus précisément le paragraphe 17(1) (b) :

aucune personne qui a été assermentée en vertu de l'article 6 ne peut révéler ni sciemment faire révéler, par quelque moyen que ce soit, des renseignements obtenus en vertu de la présente loi de telle manière qu'il soit possible, grâce à ces révélations, de rattacher à un particulier, à une entreprise ou à une organisation identifiables les détails obtenus dans un relevé qui les concerne exclusivement.

Les dispositions de la Loi sur la statistique en matière de confidentialité sont extrêmement rigoureuses. Par conséquent, leur application dans des cas bien précis représente une tâche ardue quoiqu'extrêmement importante. L'objectif premier consiste à s'assurer qu'aucun résultat personnel identifiable ne puisse être inféré dans une fourchette restreinte. De plus, il est nécessaire de protéger l'information, peu importe si le sujet est susceptible d'être considéré confidentiel par les répondants. Enfin, la façon dont le public perçoit la vigilance avec laquelle nous protégeons la confidentialité des statistiques est à tout le moins aussi importante que les mesures réelles que nous prenons pour empêcher la divulgation des données des répondants.

### 15.3 Lignes directrices

#### 15.3.1 Généralités

- Distinguer le type de données à traiter; chaque type ayant des méthodes de contrôle de divulgation qui lui sont propres. Les données tabulaires sont diffusées sous forme de tableaux statistiques comportant souvent de nombreuses dimensions. Elles se divisent davantage en tableaux de fréquences et en tableaux de données quantitatives. Les microdonnées sont des enregistrements anonymisés établis pour les particuliers. Enfin, certaines données de sorties analytiques peuvent également nécessiter un contrôle de la divulgation, surtout si elles ressemblent à des données tabulaires (p.ex., des statistiques ou des histogrammes) ou à des microdonnées (p.ex., des nuages de points ou les valeurs résiduelles d'une régression).
- Consulter les Lignes directrices sur le contrôle de la divulgation (version longue) pour déterminer les méthodes de contrôle les plus appropriées pour vos types de données. Les méthodes d'accès restreint comprennent l'accès aux données à partir de centres de données identifiés, d'avoir un accès à distance sécurisé ou d'avoir un accès limité sous contrats de licence. Les méthodes de diffusion restreinte protègent les données elles-mêmes par réduction ou perturbation de l'information.
- Ne dévoilez pas les paramètres et les règles utilisées pour contrôler la divulgation. La connaissance de ces paramètres peut aider à mieux préciser la valeur de certains répondants.
- Se rappeler en tout temps que l'apparence d'une divulgation peut parfois être aussi néfaste pour l'organisme qu'un cas réel de divulgation.

#### 15.3.2 Divulgation résiduelle

- Tenir compte du risque de divulgation résiduelle. Elle a lieu lorsqu'il est possible d'estimer des données confidentielles par un recoupement de l'information diffusée avec d'autres renseignements accessibles, y compris les diffusions antérieures de l'organisme.

- Dans les tableaux, on doit parfois trouver des cellules complémentaires à supprimer afin de protéger les cellules confidentielles. Les cellules à fréquence zéro peuvent aussi poser un problème de divulgation d'attributs parce qu'elles éliminent certaines possibilités (par exemple, une fréquence zéro pour la catégorie « possède un emploi »). Souvent, il ne suffit pas de supprimer uniquement les cellules confidentielles lorsque la distribution marginale est également diffusée, car il est parfois possible de calculer la valeur exacte des cellules supprimées en résolvant un système d'équations linéaires. Même si cela n'est pas possible, on peut calculer une fourchette de valeurs correspondant à la cellule supprimée en utilisant des méthodes de programmation linéaire, et cette fourchette peut être jugée trop restreinte pour protéger suffisamment la valeur supprimée.
- Vérifier si les catégories et hiérarchies utilisées par les tableaux se chevauchent. Par exemple, des régions publiables peuvent être soustraites de plus grandes régions et entraîner la publication d'une région dont les valeurs seraient confidentielles.
- La divulgation résiduelle a aussi lieu lorsqu'il est possible d'estimer des données confidentielles par un recoupement de l'information diffusée avec d'autres renseignements accessibles, y compris les diffusions antérieures de l'organisme. Il est difficile de formuler des règles afin d'empêcher les divulgations par recoupement lorsque plusieurs produits sont diffusés à partir du même ensemble de données de base, surtout dans les cas de demandes spéciales ou de sorties des centres de données; il faut parfois recourir à l'intervention manuelle. Si des données peuvent être diffusées à partir de plusieurs centres il est nécessaire de coordonner la diffusion ou au minimum d'établir des règles communes pour la diffusion.

### 15.3.3 Microdonnées

- Considérer des méthodes de contrôle de divulgation qui sont appropriées à la diffusion de microdonnées. Les méthodes de réduction des données englobent l'échantillonnage, l'élargissement des catégories de variables (dans le cas de certains groupes identifiables, assurez-vous que la population est assez grande), le regroupement des valeurs extrêmes supérieures et inférieures, la suppression de certaines variables provenant de certains ou de tous les répondants, la suppression de certains répondants du fichier. Les méthodes de modification des données comprennent l'ajout de bruit aléatoire aux microdonnées, la permutation de données, le remplacement de valeurs dans des groupes restreints par des valeurs moyennes ou la suppression de renseignements fournis par certains répondants et leur remplacement par des valeurs imputées.
- Dans les enquêtes longitudinales, déterminer une stratégie convenable avant que l'enquête soit terminée. Les stratégies de diffusion de fichiers de microdonnées provenant d'enquêtes longitudinales posent un problème encore plus épineux. La stratégie doit être élaborée avant que tous les résultats de l'enquête soient disponibles, soit avant la collecte des données pour les prochaines éditions de l'enquête. Comme un des objectifs de la stratégie consiste à définir les variables qui seront diffusées et leur catégorisation, certaines hypothèses doivent être formulées relativement à l'évolution de ces variables dans le temps, notamment à savoir si certaines variables sont susceptibles de devenir des variables clés.
- Dans les cas d'enquêtes de suivi ou de deuxième phase, si l'enquête principale a diffusé ou prévoit diffuser un fichier de microdonnées s'assurer que le fichier de microdonnées ne présente pas de risques additionnels de par le fait qu'on pourrait apparier les microdonnées des deux enquêtes pour créer un fichier composite. Évaluer le taux de succès d'un appariement des deux fichiers et, s'il est important, le risque découlant d'un tel appariement (par exemple, quelles sont les conséquences de l'ajout de variables identificatrices d'une enquête à l'autre).
- En conformité avec la Politique sur la diffusion des microdonnées (Statistique Canada, 1987) s'assurer que le Comité de la diffusion des microdonnées examine tout fichier de microdonnées à grande diffusion.

#### 15.3.4 Divulgence de certains types de renseignements

- Consulter le paragraphe 17(2) de la Loi sur la statistique qui prévoit que certains types de renseignements confidentiels peuvent être diffusés à la discrétion du statisticien en chef et en vertu d'une ordonnance. La diffusion de listes d'entreprises avec adresse et classification industrielle ou la communication de renseignements sur un répondant qui a donné son consentement écrit au préalable (une renonciation) constituent les formes les plus courantes de ce type de divulgation. La diffusion d'information qui s'appuie sur le pouvoir discrétionnaire du statisticien en chef est régie par la Politique relative à la révélation discrétionnaire (Statistique Canada, 2004) et, dans certains cas, par les Lignes directrices relatives à la diffusion de microdonnées non filtrées en vertu d'accords de partage des données prévus par l'article 12 ou en vertu de dispositions de diffusion discrétionnaire des renseignements.

#### 15.3.5 Ressources

- Consulter les ressources disponibles à Statistique Canada en matière de confidentialité :
  - La Division des services d'accès et de contrôle des données offre des avis et conseils à propos des politiques liées à la confidentialité de l'information recueillie par Statistique Canada;
  - Le Comité de la confidentialité et des mesures législatives et ses sous-comités, le Comité d'examen en matière de divulgation et le Comité de la diffusion des microdonnées offrent des stratégies et des pratiques de contrôle de la divulgation;
  - Le Centre de ressources sur le contrôle de la divulgation au sein de la Division des méthodes d'enquêtes auprès des entreprises offre l'aide technique ainsi que l'équipe de soutien des systèmes généralisés pour le logiciel Confid. 169
- Utiliser un logiciel généralisé de contrôle de la divulgation bien établi, tel Confid, plutôt que des systèmes personnalisés. Un tel système réduit le risque d'erreur de mise en œuvre et d'exécution, le risque de divulgation et le risque de « surprotéger » les données, tout en permettant une réduction des coûts et du temps que nécessite la mise au point.

### 15.4 Indicateurs de qualité

Principaux éléments de la qualité : exactitude, accessibilité

En général, les activités de contrôle de la divulgation ont une incidence réductrice sur la qualité des données en cela qu'elles peuvent se traduire par la suppression ou la modification d'un niveau de détail. Le contrôle de la divulgation peut aussi se traduire à limiter l'accès aux données à des groupes de la population tels que les chercheurs. Certaines méthodes telles que la perturbation des données peuvent influencer sur l'exactitude de l'information diffusée. Un biais peut provenir du fait d'arrondir les valeurs ou d'ajouter un bruit aux données.

Il n'est pas réalisable d'offrir une garantie absolue de la confidentialité. Le contrôle de la divulgation est assez complexe et les règles utilisées pour mesurer l'ampleur de la protection offerte sont modérément subjectives. Bien qu'il n'y ait pas de consensus sur les mesures sur la qualité, on retrouve principalement les fonctions de risque et les fonctions de perte.

La fonction de perte mesure l'ampleur de la différence entre les données originales et les données après l'accomplissement de méthodes de contrôle de la divulgation. Pour les données modifiées (p. ex. la perturbation), on mesure la différence relative entre les données avant et après ajustement pour la confidentialité. Dans le cas de données supprimées, on utilise souvent le taux de suppression qui indique la quantité de valeurs qui ont été supprimées par rapport à celles diffusées. Ces indices doivent être produits à différents niveaux de détail et pour divers groupes de répondants (p. ex. pour identifier les groupes industriels les plus touchés par la suppression).

La fonction de risque indique dans une certaine mesure le danger d'identifier un répondant ou une valeur qui lui est rattachée. Généralement pour des données supprimées dans les tableaux, il faut identifier le nombre de cellules supprimées dont la protection est inadéquate c.-à-d. qu'il est possible d'obtenir une approximation trop précise de la valeur supprimée en utilisant l'information provenant des autres cellules. Dans le cas des microdonnées, la plupart des méthodes tendent à mesurer le risque de divulgation en utilisant la méthode de ré-identification pour un ensemble de variables caractéristiques (appelées variables clés) ou en mesurant les tentatives d'appariement avec un fichier externe. Globalement la technique consiste à identifier des combinaisons uniques de la population qui se retrouvent dans l'ensemble de données diffusées.

## Bibliographie

BRACKSTONE, G. et P. WHITE. 2002. « Data Stewardship at Statistics Canada », *Proceedings of the Social Statistics Section*, American Statistical Association, p. 284 à 293.

DOYLE, P., J. LANE, J. THEEUWES J. et L. ZAYATZ. 2001. *Confidentiality, Disclosure and Data Access: Theory and Practical Applications for Statistical Agencies*, North-Holland.

ELLIOT, M., A. HUNDEPOOL, E. SCHULTE NORDHOLT, J.L. TAMBAY et T. WENDE. 2005. *Glossary on Statistical Disclosure Control* (en ligne), <http://neon.vb.cbs.nl/casc/glossary.htm>.

FEDERAL COMMITTEE ON STATISTICAL METHODOLOGY. 2005. *Report on Statistical Disclosure Limitation Methodology. Statistical Policy Working Paper 22. Second version*, Office of Management and Budget, Washington, D.C.

HUNDEPOOL, A. et coll. 2008a.  *$\tau$ -ARGUS version 3.3 User's Manual*, Voorburg, Statistics Netherlands.

HUNDEPOOL, A. et coll. 2008b.  *$\mu$ -ARGUS version 4.2 User's Manual*, Voorburg, Statistics Netherlands.

HUNDEPOOL, A. et coll. 2009. *Handbook on Statistical Disclosure Control, Version 1.1.*, EssNet SDC.

STATISTIQUE CANADA. 1970. *Loi sur la Statistique* (en ligne), [http://icn-rci.statcan.ca/10/10\\_006\\_f.htm](http://icn-rci.statcan.ca/10/10_006_f.htm).

STATISTIQUE CANADA. 1987. « Politique sur la diffusion des microdonnées », *Manuel des politiques de Statistique Canada*, [http://icn-rci.statcan.ca/10/10c/10c\\_026\\_f.htm](http://icn-rci.statcan.ca/10/10c/10c_026_f.htm).

STATISTIQUE CANADA. 2004. « Politique relative à la révélation discrétionnaire », *Manuel des politiques de Statistique Canada*, [http://icn-rci.statcan.ca/10/10c/10c\\_026\\_f.htm](http://icn-rci.statcan.ca/10/10c/10c_026_f.htm).

UN ECONOMIC COMMISSION FOR EUROPE. 2007. *Managing Statistical Confidentiality and Microdata Access – Principles and Guidelines of Good Practice*, Genève, Nations Unies.

WILLENBORG, L. et T. DE WAAL. 1996. « Statistical Disclosure Control in Practice », *Lecture Notes in Statistics*, Springer Verlag.

WILLENBORG, L. et T. DE WAAL. 2000. « Elements of Statistical Disclosure Control », *Lecture Notes in Statistics*, Springer Verlag.

## 16 Diffusion et communication des données

### 16.1 Portée et objet

La diffusion consiste à mettre les données obtenues dans le cadre d'une activité statistique à la disposition des utilisateurs par divers moyens. Lors de chaque diffusion de données, il importe aussi de communiquer efficacement les données à leurs utilisateurs et de faire savoir que les données sont disponibles. Le site Web de Statistique Canada ([www.statcan.ca](http://www.statcan.ca)) est le principal mode de diffusion de l'organisme, mais d'autres modes sont utilisés de façon à répondre aux besoins de certains utilisateurs. Une diffusion de données dans le cadre d'un programme statistique particulier pourra comporter la publication d'un article dans *Le Quotidien*, des tableaux de données (p. ex., tableaux CANSIM ou tableaux sommaires), des publications électroniques et autres ainsi que des métadonnées, toute cette information devenant disponible simultanément. Parmi les autres modes de diffusion, mentionnons les fichiers de microdonnées, les réponses à des demandes spéciales, des allocutions prononcées en public, des exposés ou encore des entrevues données à la télévision ou à la radio.

### 16.2 Principes

Les activités de diffusion et les activités de communications connexes ont pour objectif d'optimiser l'utilisation des données de Statistique Canada et de maintenir la pertinence de Statistique Canada, ce que l'on accomplit :

- en répondant aux besoins des utilisateurs dans le contexte de l'élaboration et de la diffusion de l'information;
- en offrant un accès plus large à l'information (diffusion directe ainsi que par l'entremise d'autres organismes);
- en assurant le plus large accès possible à l'information d'intérêt général, et ce, sans frais, tout en recouvrant les coûts associés à la communication de renseignements spécialisés et au maintien d'une infrastructure de prestation appropriée.

Ce sont là des objectifs fondamentaux aux fins de faire connaître la pertinence des activités de Statistique Canada aux ménages canadiens, aux entreprises, aux établissements, à d'autres organismes statistiques, aux autres ministères fédéraux, aux provinces et aux territoires, et aussi d'obtenir leur soutien à l'égard des activités de collecte de l'organisme. La plupart des objectifs en question sont atteints grâce à la Politique sur les services de diffusion, de communication et de commercialisation de Statistique Canada (Statistique Canada, 1985). Les lignes directrices qui suivent renvoient en outre à plusieurs autres politiques connexes.

### 16.3 Lignes directrices

#### 16.3.1 Diffusion de données statistiques

- Statistique Canada autorisera la diffusion publique de fichiers de microdonnées lorsque les conditions suivantes sont réunies : a) la diffusion donne lieu à une hausse notable de la valeur analytique des données recueillies; b) l'organisme est convaincu que toutes les mesures raisonnables ont été prises pour prévenir l'identification des unités d'enquête (Politique sur la diffusion des microdonnées; Statistique Canada, 1987).
- Ainsi que cela est indiqué dans la politique sur les relations avec les médias (Statistique Canada, 2003a), l'organisme a pour politique d'accepter les demandes d'entrevue présentées par les médias, de fournir des commentaires et d'interpréter les données. Il est interdit en toutes circonstances de tenir des réunions d'information ou des interviews à caractère non officiel.

- Il arrivera à l'occasion que des déclarations erronées concernant Statistique Canada et ses programmes ou politiques, ou encore des erreurs d'interprétation, soient véhiculées par les médias. Lorsque cela se produit, l'organisme évalue aussitôt l'incidence que peut avoir l'erreur et détermine le meilleur moyen d'y donner suite. Les gestionnaires d'enquêtes sont incités à communiquer avec les membres compétents de la Division des communications et des services de bibliothèque afin de suivre une formation en matière de relations avec les médias et d'obtenir de l'aide s'il est nécessaire de fournir des éclaircissements à propos d'une information transmise par les médias.
- Lorsque les données doivent être validées par un organisme externe et que cette validation est censée se traduire, ou s'est traduite auparavant, par une hausse marquée de la qualité des données, des données non diffusées de nature non confidentielle pourront être fournies à un tel organisme pour validation préalablement à la diffusion officielle dans Le Quotidien, en conformité avec les conditions énoncées dans la Politique sur Le Quotidien (diffusion officielle) (Statistique Canada, 2008a).
- En vertu de son mandat, Statistique Canada est autorisé à produire et à publier des estimations (souvent appelées projections ou prévisions) dont les dates de référence sont ultérieures à la date de publication. Leur diffusion doit respecter la Politique sur les estimations ayant les dates de référence futures (Statistique Canada, 2004a).

### 16.3.2 Préparation

- Toutes les publications qui présentent des données statistiques ou des conclusions analytiques doivent contenir une section intitulée « Faits saillants » (Politique sur les faits saillants des publications; Statistique Canada, 2004b).
- La préparation de données provenant du fichier source d'une activité statistique en vue de leur publication comporte généralement de nombreuses étapes. Il faut vérifier ces données pour s'assurer qu'elles correspondent, une fois toutes les étapes de traitement achevées, aux données d'origine. Dans le cas de données regroupées ou de variables dérivées, cela signifie qu'il doit être possible de reproduire les mêmes résultats à partir des données d'origine.
- Examiner à fond toutes les données (y compris les produits sous-jacents) préalablement à la diffusion afin de s'assurer que les données sont exactes, que l'analyse est rigoureuse, que le traitement était adéquat, qu'il est pertinent pour l'organisme de procéder à la publication et que la communication est efficace.
- Recourir s'il y a lieu à des outils automatisés, par exemple Édition électronique intelligente ou un comparateur de texte, pour réduire les risques d'erreur humaine.
- Dans la mesure du possible, éviter de préparer des produits (ébauches préliminaires) lorsque le traitement des données est en cours.
- Le style et la présentation des produits diffusés doivent correspondre à ceux des autres produits de Statistique Canada, ce qui en facilitera l'utilisation. Les articles visant à résumer les principales conclusions, les tendances et les données contextuelles destinées au grand public doivent être rédigés conformément aux lignes directrices sur la rédaction des communiqués pour Le Quotidien.
- Il faut joindre à tous les produits statistiques la documentation relative à la qualité et à la méthodologie, ou citer cette documentation en référence. La documentation en question doit fournir aux utilisateurs des indicateurs de la qualité des données et une description de la méthodologie et des concepts sous-jacents (Politique visant à informer les utilisateurs de la qualité des données et la méthodologie; Statistique Canada, 2000).

### 16.3.3 Vérification

- Veiller à ce que les produits écrits soient examinés par une personne n'ayant pas pris part à l'activité statistique.
- Vérifier soigneusement les chiffres, les périodes de référence (p. ex., au cours du dernier semestre ou du dernier trimestre) et les mots servant à décrire les tendances (p. ex., à la hausse, à la baisse) dans les articles et les publications afin de s'assurer qu'ils sont exacts.
- Éviter de reprendre dans le texte des chiffres fournis dans des tableaux; le cas échéant, veiller à ce que les chiffres concordent.
- Valider les chiffres cités dans les articles et les publications en les comparant aux chiffres d'autres produits tabulaires (p. ex., CANSIM, tableaux sommaires).
- Avant le lancement d'un produit électronique, tester tous ses liens pour s'assurer qu'ils fonctionnent comme prévu.
- Les produits doivent être diffusés simultanément dans les deux langues officielles (Politique sur les langues officielles; Statistique Canada, 2004c). Veiller à ce que le texte soit de grande qualité dans l'une et l'autre langues, et que les deux versions concordent en ce qui a trait aux données et au texte. Un comparateur de texte informatisé est mis à la disposition des auteurs sur le site intranet de Statistique Canada (Statistique Canada, 2008b).
- Tous les produits d'information, et en particulier les produits interprétatifs, analytiques et méthodologiques, dont Statistique Canada est exclusivement ou conjointement responsable, font l'objet d'une évaluation avant d'être diffusés à l'extérieur de l'organisme. L'évaluation doit garantir que le contenu des produits est compatible avec le mandat confié par le gouvernement à Statistique Canada à titre d'organisme statistique gouvernemental, et que les produits sont conformes aux normes d'éthique professionnelle généralement reconnues (Politique concernant l'évaluation des produits d'information : révision institutionnelle et examen par les pairs; Statistique Canada, 2003b).

### 16.4 Indicateurs de qualité

Principaux éléments de la qualité : accessibilité, actualité, pertinence.

- Présenter les produits en précisant les formats, les supports et les divers degrés de détail offerts, afin qu'on puisse déterminer s'ils répondent aux besoins de divers utilisateurs.
- Faire état du délai séparant l'annonce de la date de diffusion et la diffusion du produit. Le fait d'annoncer à l'avance la date de diffusion permet à tous les utilisateurs d'avoir le même accès au produit.
- Faire état du délai séparant la date ou période de référence et la diffusion du produit. Cela permet de déterminer l'actualité du produit dans l'optique des besoins des utilisateurs.
- Faire état du délai séparant la date de diffusion prévue et la date de diffusion réelle. Cela permet de mesurer la ponctualité du produit.
- Documenter les erreurs décelées lors de la vérification, soit la dernière étape avant la diffusion du produit. La détection d'erreurs à un stade aussi tardif de la production de données statistiques augmente le risque de devoir corriger des erreurs après la diffusion.
- Documenter les erreurs détectées après la diffusion.
- Surveiller la fréquence à laquelle les utilisateurs accèdent au produit d'information au fil du temps. Une baisse de cette fréquence pourrait indiquer que le produit perd de sa pertinence.

## Bibliographie

STATISTIQUE CANADA. 1987. « Politique sur la diffusion des microdonnées », *Manuel des politiques de Statistique Canada* (en ligne), [http://icn-rci.statcan.ca/10/10c/10c\\_026\\_f.htm](http://icn-rci.statcan.ca/10/10c/10c_026_f.htm).

STATISTIQUE CANADA. 2000d. « Politique visant à informer les utilisateurs de la qualité des données et la méthodologie », *Manuel des politiques de Statistique Canada* (en ligne), [http://icn-rci.statcan.ca/10/10c/10c\\_010\\_f.htm](http://icn-rci.statcan.ca/10/10c/10c_010_f.htm).

STATISTIQUE CANADA. 2003a. « Politique sur les relations avec les médias : porte-parole et réponse aux médias », *Manuel des politiques de Statistique Canada* (en ligne), [http://icn-rci.statcan.ca/10/10c/10c\\_002\\_f.htm](http://icn-rci.statcan.ca/10/10c/10c_002_f.htm).

STATISTIQUE CANADA. 2003b. « Politique concernant l'évaluation des produits d'information (révision institutionnelle et évaluation par les pairs) », *Manuel des politiques de Statistique Canada* (en ligne), [http://icn-rci.statcan.ca/10/10c/10c\\_011\\_f.htm](http://icn-rci.statcan.ca/10/10c/10c_011_f.htm).

STATISTIQUE CANADA. 2004. « Politique sur les services de diffusion, de communication et de commercialisation », *Manuel des politiques de Statistique Canada* (en ligne), [http://icn-rci.statcan.ca/10/10c/10c\\_015\\_f.htm](http://icn-rci.statcan.ca/10/10c/10c_015_f.htm).

STATISTIQUE CANADA. 2004a. « Politique sur les estimations ayant des dates de référence futures », *Manuel des politiques de Statistique Canada* (en ligne), [http://icn-rci.statcan.ca/10/10c/10c\\_009\\_f.htm](http://icn-rci.statcan.ca/10/10c/10c_009_f.htm).

STATISTIQUE CANADA. 2004b. « Politique sur les faits saillants des publications », *Manuel des politiques de Statistique Canada* (en ligne), [http://icn-rci.statcan.ca/10/10c/10c\\_008\\_f.htm](http://icn-rci.statcan.ca/10/10c/10c_008_f.htm).

STATISTIQUE CANADA. 2004c. « Politique sur les langues officielles », *Manuel des politiques de Statistique Canada* (en ligne), [http://icn-rci.statcan.ca/10/10c/10c\\_034\\_f.htm](http://icn-rci.statcan.ca/10/10c/10c_034_f.htm).

STATISTIQUE CANADA. 2008. « Politique pour Le Quotidien et la diffusion officielle », *Manuel des politiques de Statistique Canada* (en ligne), [http://icn-rci.statcan.ca/10/10c/10c\\_058-fra.htm](http://icn-rci.statcan.ca/10/10c/10c_058-fra.htm).

## 17 Analyse et présentation des données

### 17.1 Portée et objet

L'analyse des données est le processus qui consiste à examiner et à interpréter des données afin d'élaborer des réponses à des questions. Les principales étapes du processus d'analyse consistent à cerner les sujets d'analyse, à déterminer la disponibilité de données appropriées, à décider des méthodes qu'il y a lieu d'utiliser pour répondre aux questions d'intérêt, à appliquer les méthodes et à évaluer, résumer et communiquer les résultats.

Les résultats analytiques soulignent l'utilité des sources de données en jetant de la lumière sur les sujets pertinents. Certains programmes de Statistique Canada dépendent des résultats analytiques à titre de principal produit de données, car, pour des raisons de confidentialité, il est impossible de diffuser les microdonnées. L'analyse des données joue également un rôle clé dans le processus d'évaluation de la qualité des données en indiquant les problèmes liés à la qualité des données dans une enquête particulière. Ainsi, l'analyse peut influencer sur les améliorations futures au processus d'enquête.

L'analyse des données est essentielle pour comprendre les résultats des enquêtes, des sources administratives et des études pilotes, pour obtenir des renseignements sur les lacunes en matière de données, pour concevoir et remanier les enquêtes, pour planifier de nouvelles activités statistiques et pour formuler des objectifs en matière de qualité.

Les résultats de l'analyse des données sont souvent publiés ou résumés dans les diffusions officielles de Statistique Canada.

### 17.2 Principes

Un organisme statistique veille à la pertinence et à l'utilité de l'information que contiennent ses données pour les utilisateurs. L'analyse est le principal outil permettant d'obtenir de l'information à partir des données.

Les données d'une enquête peuvent être utilisées à des fins d'études descriptives ou analytiques. Les études descriptives se réfèrent à l'estimation de mesures agrégées d'une population cible, par exemple les bénéfices moyens des entreprises exploitées par le propriétaire en 2005 ou la proportion de diplômés du secondaire en 2007 qui ont poursuivi des études supérieures au cours des 12 mois suivants. Les études analytiques peuvent servir à expliquer le comportement de caractéristiques ou les relations entre elles; une étude des facteurs de risque d'obésité chez les enfants, par exemple, serait de nature analytique.

Pour être efficace, l'analyste doit comprendre les questions pertinentes (tant celles qui sont actuelles que celles qui sont susceptibles d'émerger à l'avenir) et comment présenter les résultats au public. L'étude du contexte de l'analyse permet à l'analyste de choisir les sources de données et les méthodes statistiques appropriées. Toutes les conclusions présentées dans une analyse, y compris celles qui peuvent avoir une incidence sur les politiques publiques, doivent être appuyées par les données analysées.

### 17.3 Lignes directrices

#### 17.3.1 Préparation initiale

Avant de procéder à une étude analytique, il faut se pencher sur les questions suivantes :

- Objectifs. Quels sont les objectifs de cette analyse? Quel est le sujet abordé? Quelles sont la ou les questions auxquelles il s'agit de trouver une réponse?
- Justification. Pourquoi cette question est-elle intéressante? Comment ces réponses contribueront-elles à la somme des connaissances existantes? Quelle est la pertinence de cette étude?
- Données. Quelles données sont utilisées? Quelle est la meilleure source de données pour cette analyse? Y a-t-il des limites?

- Méthodes d'analyse. Quelles techniques statistiques sont appropriées? Permettront-elles d'atteindre les objectifs?
- Public. Qui s'intéresse à cette question, et pourquoi?

### 17.3.2 Données appropriées

- S'assurer que les données conviennent à l'analyse à effectuer. À cette fin, il faut se pencher sur un grand nombre de détails tels que : la population visée par la source de données est-elle suffisamment reliée à la population cible de l'analyse? Les variables de la source ainsi que les définitions et les concepts sous-jacents sont-ils pertinents dans le cadre de l'étude? La nature longitudinale ou transversale de la source des données convient-elle à l'analyse? La taille de l'échantillon du domaine de l'étude est-elle suffisante pour dégager des résultats convenables? La qualité des données, telle qu'elle est exposée dans la documentation de l'enquête ou évaluée au moyen d'analyse, est-elle suffisante?
- Si plus d'une source de données sert à l'analyse, déterminer si les sources sont cohérentes et comment les intégrer à l'analyse de la manière appropriée.

### 17.3.3 Méthodes et outils appropriés

- Choisir une approche analytique qui convient à la question examinée et aux données à analyser.
- Pour analyser les données d'un échantillon probabiliste, il peut être approprié d'utiliser des méthodes analytiques qui font abstraction du plan d'enquête, si un nombre suffisant des conditions du modèle pour l'analyse sont satisfaites (voir Binder et Roberts, 2003). Toutefois, les méthodes qui intègrent les renseignements sur le plan d'échantillonnage sont généralement efficaces même lorsque certains aspects du modèle sont spécifiés incorrectement.
- Déterminer si l'information sur le plan de sondage peut être intégrée à l'analyse et, le cas échéant, la façon de procéder – par exemple, au moyen de méthodes fondées sur le plan de sondage. Voir Binder et Roberts (2009) et Thompson (1997) pour un examen de diverses approches pour l'inférence sur des données tirées d'un échantillon probabiliste.
  - Voir Chambers et Skinner (2003), Korn et Graubard (1999), Lehtonen et Pahkinen (1995), Lohr (1999) et Skinner, Holt et Smith (1989) pour plusieurs exemples de méthodes analytiques fondées sur le plan de sondage.
  - Pour une analyse fondée sur le plan de sondage, consulter la documentation de l'enquête au sujet de l'approche recommandée pour l'estimation de la variance pour l'enquête. Si l'analyse porte sur les données de plus d'une enquête, déterminer si les différents échantillons ont été sélectionnés indépendamment ou non, et quel en serait l'effet sur l'approche appropriée de l'estimation de la variance.
  - Les fichiers de données pour les enquêtes probabilistes contiennent souvent plus d'une variable de pondération, particulièrement dans le cas d'une enquête longitudinale ou menée dans le but de recueillir des données transversales ainsi que longitudinales. Consulter la documentation de l'enquête et les spécialistes des enquêtes si le choix du meilleur poids à utiliser n'est pas évident pour une analyse fondée sur le plan de sondage dans une enquête particulière.
  - Lorsqu'il s'agit d'analyser des données provenant d'une enquête probabiliste, les renseignements disponibles sur le plan de sondage peuvent être insuffisants pour permettre d'adopter une approche complètement fondée sur le plan de sondage. Évaluer les solutions de rechange qui s'offrent.
- Consulter des spécialistes du sujet à propos de la source des données et les méthodes statistiques si on n'est pas familier avec ces dernières.
- Après avoir déterminé la méthode analytique appropriée aux données, examiner les choix de logiciels qui s'offrent pour l'appliquer. S'il s'agit d'analyser les données provenant d'un échantillon probabiliste au moyen de méthodes fondées sur le plan de sondage, utiliser un logiciel conçu particulièrement pour l'analyse de données d'enquête, puisque les logiciels analytiques standard qui peuvent produire des estimations ponctuelles pondérées ne calculent pas correctement les variances pour des estimations pondérées par les poids de sondage.

- Il est souhaitable d'utiliser un logiciel commercial, s'il convient à la tâche, pour effectuer les analyses choisies, puisque ceux-ci ont généralement été testés davantage que les logiciels non commerciaux.
- Déterminer s'il est nécessaire de reformater les données afin d'utiliser le logiciel choisi.
- Inclure divers diagnostics parmi les méthodes d'analyse utilisées si l'on ajuste des modèles aux données analysées.
- Les sources de données varient beaucoup pour ce qui est des données manquantes. À une extrémité se trouvent les sources qui semblent complètes, dans lesquelles on a tenu compte de toute unité manquante au moyen d'une variable de pondération ayant une composante de non-réponse et toutes les valeurs manquantes des unités déclarantes ont été remplacées par des valeurs imputées. À l'autre extrémité se trouvent les sources de données dans lesquelles aucun traitement n'a été effectué sur les données manquantes. Ainsi, le travail de l'analyste pour traiter les données manquantes peut varier fortement. Il convient de souligner que les mesures à prendre dans le cas de données manquantes dans une analyse sont un sujet de recherche permanent.
  - Se reporter à la documentation au sujet de la source des données pour déterminer la mesure dans laquelle les données manquent, les types de données manquantes et le traitement des données manquantes qui a été effectué. Ces renseignements serviront de point de départ pour déterminer les autres travaux pouvant être requis.
  - Déterminer la façon de traiter la non-réponse totale ou partielle dans l'analyse, en prenant en compte l'importance des données manquantes et les types de données manquantes dans les sources de données utilisées.
  - Déterminer si les valeurs imputées doivent être incluses dans l'analyse et, le cas échéant, la façon dont il convient de les traiter. Si les valeurs imputées ne sont pas utilisées, il faut déterminer quelles autres méthodes peuvent être utilisées pour rendre compte correctement de l'effet de la non-réponse dans l'analyse.
  - Si l'analyse comprend la modélisation, il pourrait être approprié d'inclure certains aspects de la non-réponse dans le modèle analytique.
  - Faire toutes les mises en garde nécessaires sur la façon dont les méthodes utilisées pour traiter les données manquantes peuvent influencer sur les résultats.

#### **17.3.4 Interprétation des résultats**

- Étant donné que la majorité des analyses sont fondées sur des études par observation plutôt que sur les résultats d'une expérience contrôlée, éviter de tirer des conclusions en ce qui concerne la causalité.
- En étudiant les changements survenus au fil du temps, veiller à examiner les tendances à court terme en considérant également les tendances à moyenne et à long terme. Les tendances à court terme ne représentent souvent que de légères fluctuations d'une tendance plus importante à moyen ou à long terme.
- Lorsque possible, éviter les points de référence arbitraires. Privilégier l'utilisation de points de référence comportant une plus grande signification tels que le dernier tournant pour les données économiques, les différences intergénérationnelles pour les statistiques démographiques et les changements législatifs pour les statistiques sociales.

#### **17.3.5 Présentation des résultats**

- Mettre l'accent dans l'article sur les variables et les sujets importants. Lorsque le sujet abordé est trop vaste, l'impact principal du message se trouve souvent atténué.
- Structurer les idées de façon logique, en fonction de leur pertinence ou de leur importance. Recourir à des titres, à des sous-titres et à des encadrés afin de renforcer la structure de l'article.
- Rédiger le texte en langage aussi simple que le sujet le permet. Selon le public cible, il est parfois souhaitable de perdre un peu en précision pour rendre le texte plus lisible.

- Insérer des graphiques en complément du texte et des tableaux pour communiquer le message. Privilégier les titres qui véhiculent un message (p. ex. « Les revenus des femmes demeurent inférieurs à ceux des hommes »), plutôt que des titres de graphique classiques (p. ex. « Revenus selon l'âge et le sexe »). Toujours commenter l'information fournie dans les tableaux et les graphiques afin de permettre au lecteur de mieux la comprendre.
- Lorsque des tableaux sont insérés, la présentation générale doit contribuer à la clarté des données qu'ils contiennent et prévenir les erreurs d'interprétation. Cela comprend l'espacement, la formulation, l'emplacement et l'apparence des titres, les titres de lignes et de colonnes et autre étiquetage.
- Expliquer les pratiques ou les méthodes d'arrondissement. Dans la présentation de données arrondies, le nombre de chiffres significatifs ne doit pas être supérieur à celui qu'exige l'exactitude des données.
- Satisfaire aux exigences en matière de confidentialité (p. ex. taille minimale des cellules) imposées par les enquêtes ou les sources administratives dont les données font l'objet de l'analyse.
- Fournir des renseignements sur les sources de données utilisées ainsi que toutes lacunes dans les données ayant pu avoir une incidence sur l'analyse. Inclure dans le document soit une section sur les données, soit un renvoi indiquant au lecteur où obtenir les détails.
- Fournir des renseignements sur les méthodes analytiques et les outils utilisés. Inclure soit une section portant sur les méthodes, soit un renvoi indiquant au lecteur où obtenir les détails.
- Inclure des renseignements sur la qualité des résultats. Les erreurs types, les intervalles de confiance ou les coefficients de variation fournissent au lecteur des renseignements importants sur la qualité des données. Le choix de l'indicateur peut varier selon l'endroit où l'article est publié.
- S'assurer que toutes les références sont exactes, uniformes et font l'objet de renvois dans le texte.
- S'assurer qu'il n'y a pas d'erreurs dans l'article. Vérifier les détails, par exemple la cohérence des chiffres dans le texte, les tableaux et les graphiques, ainsi que l'exactitude des données externes et des calculs arithmétiques simples.
- S'assurer que ce qui est annoncé dans l'introduction est effectivement exprimé dans le reste de l'article. S'assurer que les conclusions sont cohérentes avec les résultats de l'analyse.
- Faire réviser l'article par d'autres personnes pour en vérifier la pertinence, l'exactitude et l'intelligibilité, peu importe où il doit être diffusé. Comme bonne pratique, demander à quelqu'un de la division qui a fourni les données d'examiner comment ces dernières ont été utilisées. Si l'article doit être diffusé à l'extérieur de Statistique Canada, il doit être soumis à un examen institutionnel ainsi que par les pairs, tel qu'il est précisé dans la Politique concernant l'évaluation des produits d'information (Statistique Canada, 2003).
- Si l'article doit être diffusé dans une publication de Statistique Canada, s'assurer qu'il est conforme aux normes d'édition en vigueur de Statistique Canada. Ces normes sont applicables aux graphiques, aux tableaux et au style, entre autres.
- Comme bonne pratique, envisager de présenter les résultats à des pairs avant de mettre la dernière main au texte. Il s'agit d'un autre type d'examen par les pairs qui peut aider à améliorer l'article. Toujours procéder à une répétition des exposés destinés à des publics externes.
- Consulter les documents disponibles qui pourraient fournir d'autres conseils pour améliorer l'article, comme les Lignes directrices sur la rédaction d'articles d'analyse (Statistique Canada, 2008) et le Guide de rédaction. (Statistique Canada, 2004)

## 17.4 Indicateurs de qualité

Principaux éléments de la qualité : pertinence, intelligibilité, exactitude, accessibilité

Un produit analytique est pertinent s'il y a un public qui s'intéresse (ou qui s'intéressera) aux résultats de l'étude.

Pour que le degré d'intelligibilité d'un article analytique soit élevé, le style de rédaction doit être adapté au public cible. En outre, l'article doit fournir suffisamment de détails pour permettre à une autre personne à laquelle l'accès aux données serait accordé de reproduire les résultats.

Pour qu'un produit analytique soit exact, il faut utiliser les méthodes et les outils appropriés pour produire les résultats.

Pour qu'un produit analytique soit accessible, il doit être mis à la disposition des personnes auxquelles les résultats de la recherche seraient utiles.

### Bibliographie

BINDER, D.A. et G.R. ROBERTS. 2003. « Design-based Methods for Estimating Model Parameters », *Analysis of Survey Data*, R.L. Chambers et C.J. Skinner, Chichester, Wiley, p. 29 à 48.

BINDER, D.A. et G. ROBERTS. 2009. « Design and Model Based Inference for Model Parameters », *Sample Surveys: Inference and Analysis*, D. Pfeffermann et C.R. Rao, Amsterdam, Elsevier.

CHAMBERS, R.L. et C.J. SKINNER. 2003. *Analysis of Survey Data*, Chichester, Wiley.

KORN, E.L. et B.I. GRAUBARD. 1999. *Analysis of Health Surveys*, New York, Wiley.

LEHTONEN, R. et E.J. PAHKINEN. 2004. *Practical Methods for Design and Analysis of Complex Surveys*. 2e édition, Chichester, Wiley.

LOHR, S.L. 1999. *Sampling: Design and Analysis*, Duxbury Press.

SKINNER, C.K., D. HOLT et T.M.F. SMITH. 1989. *Analysis of Complex Surveys*, Chichester, Wiley.

THOMPSON, M.E. 1997. *Theory of Sample Surveys*, Londres, Chapman and Hall.

STATISTIQUE CANADA. 2003. « Politique concernant l'évaluation des produits d'information », *Manuel des politiques de Statistique Canada* (en ligne), [http://icn-rci.statcan.ca/10/10c/10c\\_011\\_f.htm](http://icn-rci.statcan.ca/10/10c/10c_011_f.htm).

STATISTIQUE CANADA. 2004. *Guide de rédaction de Statistique Canada* (en ligne), [http://icn-rci.statcan.ca/10/10d/10d\\_000\\_f.htm](http://icn-rci.statcan.ca/10/10d/10d_000_f.htm).

STATISTIQUE CANADA. 2008. *Lignes directrices sur la rédaction d'articles d'analyse* (en ligne), [http://icn-rci.statcan.ca/10/10g/10g\\_001\\_f.htm](http://icn-rci.statcan.ca/10/10g/10g_001_f.htm)

## 18 Documentation

### 18.1 Portée et objet

La documentation est le compte rendu de l'activité statistique (concepts, définitions et méthodes) qui entoure la collecte, le traitement et l'analyse des données et marque les produits statistiques qui en résultent. Elle vise à favoriser un usage efficace et informé des données. La documentation devrait inclure les indicateurs de qualité générés au cours de l'activité statistique et l'analyse de leur impact sur l'usage des produits résultant de cette activité.

Durant la mise en oeuvre, la documentation assure l'élaboration efficace de l'activité statistique, en plus de faire état des décisions qui lui sont associées et de leur bien-fondé. De plus, les renseignements contenus dans la documentation d'une activité statistique sont utiles à l'élaboration d'activités similaires ou de versions remaniées de cette activité.

### 18.2 Principes

La documentation vise à fournir un compte rendu complet, sans équivoque et polyvalent de l'activité statistique, notamment de ses extraits. Elle s'adresse à divers publics cibles tels que des gestionnaires, du personnel technique, des planificateurs assignés à d'autres enquêtes et des utilisateurs. Il faut qu'elle soit à jour, qu'elle soit aisément et rapidement accessible (dans des délais assurant sa pertinence) et qu'elle soit intelligible pour son public cible. Il est possible d'élaborer la documentation selon une approche multimédia (format papier, format électronique, présentation visuelle). Il faut s'assurer de préserver les documents portant sur l'activité statistique.

### 18.3 Lignes directrices

#### 18.3.1 Garantir la production d'une documentation adaptée au public cible et au contexte général

- Le degré d'érudition des documents doit tenir compte de leur principal public cible. Il faut donc déterminer si les documents doivent être détaillés ou sommaires, techniques ou vulgarisés, etc. Lorsqu'un produit statistique est diffusé par Statistique Canada, la documentation doit respecter les exigences de la Politique visant à informer les utilisateurs de la qualité des données et la méthodologie (Statistique Canada, 2000d).
- L'envergure de la documentation devrait dépendre du statut de l'activité statistique qu'elle cible. Il faut se demander s'il s'agit d'une nouvelle activité ou d'une activité répétée, si cette activité s'apparente à d'autres activités menées du Bureau ou si elle s'en distingue, etc. Il devient alors possible de déterminer s'il faut préparer une nouvelle documentation ou si l'on peut se contenter de références à des documents existants.
- En matière de documentation, les priorités doivent également tenir compte du budget alloué à l'activité statistique, du moment le plus opportun pour diffuser ladite documentation et des bénéfices qu'elle engendre, à court et à long terme. Par ailleurs, il faut éviter les retards, entre la fin de l'activité statistique et la rédaction de la documentation afférente, car ils peuvent nuire à son exactitude, en plus d'affecter son actualité et sa pertinence.

#### 18.3.2 Garantir la production d'une documentation complète et exacte

- Règle générale, la documentation d'une activité statistique se range dans l'une des trois catégories suivantes : 1) documents d'ordre général servant principalement à dresser le portrait actuel de l'activité statistique; 2) documents thématiques fournissant des détails sur sa mise en oeuvre; 3) documents d'évaluation thématiques.

### **Documentation générale**

- Objectifs : inclure des renseignements sur les usages des données et les objectifs qu'elles sont censées satisfaire, sur l'actualité et la fréquence de l'activité statistique et sur les objectifs relatifs à la qualité des données. Les objectifs sont susceptibles de changer, à mesure que l'enquête progresse, et ce, pour plusieurs raisons (contraintes budgétaires, faisabilité anticipée, résultats de nouvelles études pilotes et de nouvelles technologies). Il faut documenter ces changements, car ils ont un impact sur la conception du questionnaire et l'analyse des résultats des mises à l'essai.
- Contenu : inclure les concepts, les définitions et le questionnaire utilisés dans le cadre de l'enquête. Afin de faciliter l'intégration de ce contenu à d'autres sources, signaler le recours à des concepts, des questions, des méthodes et des classifications types et mettre en évidence les différences, s'il y a lieu. Expliquer le rôle des comités consultatifs et des utilisateurs.
- Méthodologie : traiter des questions telles que la population cible, la base de sondage, la couverture, la période de référence, le plan d'échantillonnage, la taille de l'échantillon et sa méthode de sélection, les méthodes de collecte et de suivi en cas de non-réponse, la vérification et l'imputation, l'estimation, l'étalonnage et la révision, la désaisonnalisation et la confidentialité. Fournir un survol méthodologique de l'enquête. Accentuer divers aspects pour répondre aux besoins de divers lecteurs. Regrouper les questions techniques dans un même document rédigé à l'intention du personnel technique.
- Qualité des données : fournir des renseignements d'usage général concernant la couverture, l'erreur d'échantillonnage, l'erreur non due à l'échantillonnage, les taux de réponse, les taux de vérification et d'imputation et leur impact, le facteur d'ajustement dans le temps et la comparabilité générale, les études de validation, les mesures d'assurance de la qualité et toute autre mesure pertinente propre à l'activité statistique en question. Mentionner les facteurs imprévisibles influant sur la qualité des données (inondations, taux de non-réponse élevé, etc.). Inclure la variance totale ou ses composantes, selon la source, puis traiter des biais de réponse et de non-réponse, ainsi que de l'impact et de l'interprétation de la désaisonnalisation, à l'intention des utilisateurs spécialisés.

### **Documentation détaillée sur la mise en œuvre**

- Planification des activités et budget
- Opérations : inclure un manuel des intervieweurs, des guides de formation, des instructions ou un guide destinés aux superviseurs et aux vérificateurs du contrôle de la qualité, des guides destinés au personnel chargé du traitement et de la saisie des données, ainsi que des rapports d'évaluation des opérations et des comptes rendus.
- Pour les interviews assistées par ordinateur, fournir les spécifications de développement de l'application logicielle.
- Systèmes : inclure des renseignements sur les fichiers de données (clichés d'enregistrement, explication des codes, fréquences de base, méthodes de vérification), les systèmes (construction, algorithmes, usages, stockage et extraction) et les rapports de surveillance (temps consacré à des activités précises, sources d'incidents, ordonnancement des essais visant à déterminer si les données sont traitées dans les délais prévus).
- Mise en œuvre : documenter l'ensemble des opérations en spécifiant clairement les données d'entrée et de sortie. Annexer le calendrier de travail de chaque étape de la mise en œuvre à cette documentation.
- Ressources : établir une liste des ressources employées en fonction du temps. Fournir un relevé des dépenses salariales et non salariales (montants et temps). Commenter les dépenses par rapport aux budgets.

## Évaluations

- Produire un rapport général d'évaluation sur l'activité statistique en tant que processus.
- Décrire les tests cognitifs, les tests sur le terrain et les enquêtes pilotes; faire état de leurs résultats et formuler des recommandations par rapport aux spécifications de ces analyses.
- Documenter l'évaluation méthodologique des options de rechange ou du modèle mise en oeuvre (rendement) pour le plan d'échantillonnage.
- Lorsqu'elle est diffusée à l'extérieur du Bureau, la documentation de l'activité statistique doit faire l'objet d'une révision institutionnelle et d'une évaluation par les pairs conformément à la Politique concernant l'évaluation des produits d'information (Statistique Canada, 2003). Il est d'ailleurs préférable qu'elle soit révisée par des gestionnaires, des représentants du public cible ou des pairs même lorsque son usage est restreint à l'interne, afin de garantir sa pertinence, son exactitude et son intelligibilité.
- Garantir l'accessibilité de la documentation
- Afin de faciliter la tâche des utilisateurs, intégrer les éléments de documentation requis à la Base de métadonnées intégrée. (Statistique Canada, 2000c) En tant qu'organe d'archivage des informations sur les enquêtes et programmes de Statistique Canada, la BMDI contient la majeure partie des informations intéressant les utilisateurs, en ce qui concerne la méthodologie et l'exactitude des données. Ces informations sont accessibles sous forme électronique (en empruntant un lien vers la BMDI) et imprimée (texte de la BMDI). Le texte de la BMDI satisfait les exigences de la Politique visant à informer les utilisateurs de la qualité des données et la méthodologie (Statistique Canada, 2000d).
- Opter pour des outils qui centralisent autant que possible les informations sur les activités statistiques, d'une part, et qui structurent l'entreposage et la recherche des documents, d'autre part. Chaque document doit au minimum comporter un titre clair, une date et le nom de ses auteurs (entités ou individus). Les exigences relatives à la préservation des documents de Statistique Canada sont dictées par la Politique concernant la gestion des documents.
- Classer et documenter les références (articles théoriques et généraux, documents associés au projet, bien que produits à l'extérieur de son cadre).

### 18.4 Indicateurs de qualité

Principaux éléments de la qualité : intelligibilité, accessibilité, actualité

- Nombre de documents diffusés à l'externe ayant fait l'objet d'une révision institutionnelle et d'une évaluation par les pairs.
- Produits statistiques diffusés conformes à la Politique visant à informer les utilisateurs de la qualité des données et la méthodologie.

## Bibliographie

NATIONS UNIES. 1983. *Conférence des statisticiens européens*. Ébauche des lignes directrices pour la préparation de présentations sur la couverture et la qualité des statistiques auprès des utilisateurs, Genève, Suisse.

STATISTIQUE CANADA. 2000. « Politique concernant la gestion des documents », *Manuel des politiques de Statistique Canada* (en ligne), [http://icn-rci.statcan.ca/10/10c/10c\\_040\\_f.htm](http://icn-rci.statcan.ca/10/10c/10c_040_f.htm).

STATISTIQUE CANADA. 2000d. « Politique visant à informer les utilisateurs sur la qualité des données et la méthodologie », *Manuel des politiques de Statistique Canada* (en ligne), [http://icn-rci.statcan.ca/10/10c/10c\\_010\\_f.htm](http://icn-rci.statcan.ca/10/10c/10c_010_f.htm).

STATISTIQUE CANADA. 2002c. *Le Cadre d'assurance de la qualité de Statistique Canada – 2002*, publication n° 12-586-XIF au catalogue de Statistique Canada, Ottawa.

STATISTIQUE CANADA. 2003. « Politique concernant l'évaluation des produits d'information », *Manuel des politiques de Statistique Canada* (en ligne), [http://icn-rci.statcan.ca/10/10c/10c\\_011\\_f.htm](http://icn-rci.statcan.ca/10/10c/10c_011_f.htm).

STATISTIQUE CANADA. 2004. « Politique concernant les normes », *Manuel des politiques de Statistique Canada* (en ligne), [http://icn-rci.statcan.ca/10/10c/10c\\_014\\_f.htm](http://icn-rci.statcan.ca/10/10c/10c_014_f.htm).

STATISTIQUE CANADA. 2007. *Base de métadonnées intégrée* (en ligne), [http://stdsweb/standards/imdb/imdb-menu\\_f.htm](http://stdsweb/standards/imdb/imdb-menu_f.htm).