

STRATÉGIE DE CALAGE DE L'ENQUÊTE SUR LA DYNAMIQUE DU TRAVAIL ET DU REVENU

Sylvie LaRoche¹

RÉSUMÉ

L'enquête sur la dynamique du travail et du revenu (EDTR) est une enquête longitudinale par panel menée auprès des individus. L'échantillon est composé de deux panels qui se chevauchent. Dans le cadre de la stratégie de calage harmonisée pour les statistiques du revenu, l'EDTR a mis au point un estimateur par régression combinant des comptes démographiques et fiscaux pour ajuster certaines estimations d'enquête à des totaux connus de la population. Ce document décrit le nouvel estimateur et les variables auxiliaires utilisées ainsi que l'impact de la nouvelle stratégie sur les estimations tirées de l'enquête.

MOTS CLÉS : Calage; enquête longitudinale; estimateur par régression; variables auxiliaires.

ABSTRACT

The Survey of Labour and Income Dynamics (SLID) is a longitudinal panel survey of individuals. The sample is composed of two overlapping panels. In the context of the harmonised calibration strategy for income statistics, SLID developed a regression estimator that uses a combination of demographic and fiscal control counts to adjust survey estimates to known population counts. This paper presents the estimator and the auxiliary variables used, along with the impact of the new calibration strategy on the survey estimates.

KEY WORDS: Auxiliary variables, Calibration, Longitudinal Survey, Regression estimator.

1. INTRODUCTION

1.1 Description de l'enquête

L'enquête sur la dynamique du travail et du revenu (EDTR) est une enquête longitudinale par panel menée auprès des individus. Elle vise à mesurer les changements au niveau du bien-être économique des individus ainsi que les facteurs pouvant influencer ces changements, plus particulièrement les facteurs déterminants au niveau des caractéristiques démographiques, familiales et au niveau de l'activité. L'EDTR produit des estimations longitudinales ainsi que transversales. Depuis 1996, elle est devenue la principale source d'information transversale sur le revenu des ménages et des individus, remplaçant son prédécesseur, l'enquête sur les finances des consommateurs (EFC). Chaque panel de l'EDTR a une durée de 6 ans et un nouveau panel est introduit à chaque trois ans. Ainsi, deux panels se chevauchent toujours. Chaque panel couvre près de 17 000 ménages et environ 35 000 adultes. Dans ce document, il sera question de la composante transversale de l'enquête.

Dans le processus de pondération de l'EDTR, la dernière étape, celle du calage aux marges, est utilisée pour ajuster les poids d'enquête afin que les estimations tirées de l'enquête pour certaines caractéristiques-clés de la population respectent les totaux de population fiables (appelés totaux de contrôle) provenant d'autres sources externes de données. Avant l'année de référence 2000, l'EDTR n'utilisait que des totaux de contrôle par province, groupes d'âge et sexe dérivés à partir des projections démographiques basées sur le recensement. Comme ces projections démographiques sont révisées suite à la

¹ Sylvie LaRoche, Division des méthodes d'enquêtes auprès des ménages, Statistique Canada, 18^e étage, Immeuble R.H. Coats, Ottawa, Ontario, K1A 0T6, Canada, sylvie.laroche@statcan.ca

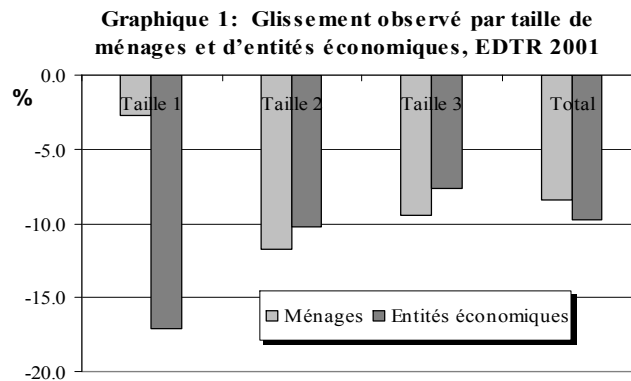
tenu d'un nouveau recensement, l'EDTR procède donc à une révision historique des poids d'enquête environ tous les cinq ans.

Cet article discute de la nouvelle stratégie de calage de l'EDTR. On énonce d'abord les problèmes ayant conduit à un changement de stratégie pour le calage de l'EDTR, de l'évolution de celle-ci au cours des années, menant finalement à la stratégie de calage nouvellement implantée. Par la suite, on montre les impacts de la nouvelle stratégie de calage sur les estimations de l'enquête et leur précision puis on conclut en discutant entre autres de considérations importantes découlant de l'implantation de cette stratégie.

1.2 Problématique à résoudre

Depuis 1993, des écarts grandissants entre certaines estimations de l'enquête et celles provenant de sources externes ont été observés. Ceci, ainsi que certaines différences entre les estimations des diverses enquêtes sur le revenu et les dépenses ont été à l'origine du lancement, en 1999, du Projet du calage harmonisé de la statistique du revenu. L'objectif de ce projet était d'harmoniser les stratégies de calage des enquêtes sur le revenu et les dépenses afin d'améliorer la comparabilité entre les diverses estimations tirées des enquêtes et celles provenant de sources externes (Tremblay, 2005).

Pour l'EDTR, il y avait deux problèmes importants à résoudre en ce qui concerne les estimations d'enquête. Le premier touche la répartition de l'échantillon selon la taille des ménages² et des entités économiques³. En comparant les estimations tirées de l'EDTR avant calage aux totaux connus de la population (appelé glissement), on observe une sous-estimation des ménages de taille deux et plus supérieure à 8 % alors que le nombre de ménages de taille 1 n'est sous-estimé que d'environ 2% (voir graphique 1). Pour ce qui est des entités économiques, on remarque aussi une sous-estimation du nombre d'entités économiques pour les différentes tailles, la plus importante étant de loin celle des entités économiques de taille 1 qui s'élève à près de 17%. Ce déséquilibre entre les estimations du nombre de ménages et d'entités économiques selon la taille résulte en une sous-estimation très importante du nombre d'entités économiques de taille 1 vivant dans les ménages de taille deux et plus. Lorsqu'aucun contrôle selon la taille de ménage et d'entité économique n'est utilisé dans le calage on observe alors un biais important dans les estimations correspondantes de l'enquête.



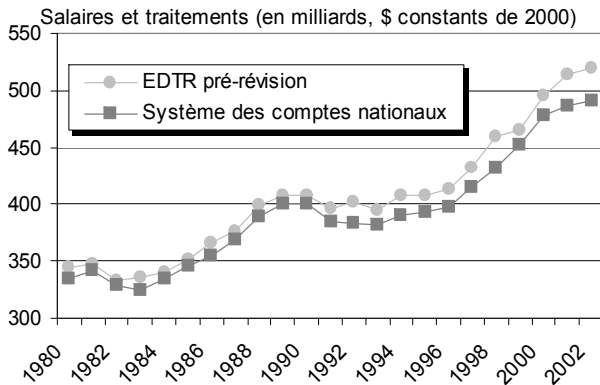
Le deuxième problème observé en terme de représentativité de l'échantillon concerne le revenu. Depuis le début des années 1990, des écarts grandissants ont été observés entre les estimations sur le revenu tirées des enquêtes et celles provenant de sources externes fiables. En effet, lorsqu'on compare la série sur les salaires et traitements agrégés tirées de l'EFC/EDTR à celle provenant du Système de comptabilité nationale présenté dans le graphique 2, on constate qu'entre 1980 et 1991, l'EFC surestimait les salaires agrégés d'environ 1.6% à 3.6% par rapport aux estimations du Système de comptabilité nationale. Puis, à partir de 1992, cet écart s'est agrandi pour fluctuer entre 3.1% et 6.2% (Lathe, 2005).

² Un ménage est constitué d'une personne ou d'un groupe de personnes occupant un même logement.

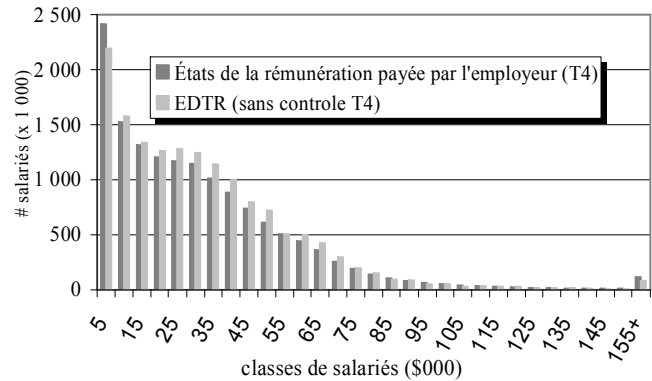
³ Une entité économique (appelé aussi famille économique) est constituée d'un groupe de deux personnes ou plus qui vivent dans le même logement et qui sont apparentées par le sang, par alliance, par union libre ou par adoption. Dans ce document, une entité économique de taille un (appelé aussi personne seule) est une personne vivant seule ou avec une ou plusieurs autres personnes non apparentées par le sang, par alliance, par union libre ou par adoption.

Par la suite, la répartition des salariés selon les classes de salaires et traitements a été examinée de façon à identifier plus précisément l'origine de cette surestimation des totaux de salaires et traitements agrégés. Tel que le démontre le graphique 3, l'EDTR tend à surestimer le nombre de personnes dans les fourchettes du milieu de la répartition des salaires et traitement alors qu'elle sous-estime le nombre de personnes dont les salaires et traitements se situent aux deux extrémités de la répartition. En regardant ces mêmes distributions pour les années antérieures, il est possible de constater que cette tendance est présente pour l'EDTR mais aussi pour l'EFC. D'autres enquêtes ont aussi observé un problème similaire (Tremblay, 2005).

Graphique 2: Totaux agrégés de salaires et traitements



Graphique 3 : Répartition du nombre de salariés selon la classe de salaires, EDTR 2000



Il était donc important pour l'EDTR de s'attaquer à ces problèmes afin d'améliorer les estimations produites ainsi que la comparabilité de ces dernières avec les estimations provenant d'autres enquêtes et de sources externes.

1.3 La méthodologie de calage pour l'EDTR

L'EDTR utilise un estimateur par régression généralisée (GREG) dans l'étape du calage (Sarndal et coll. 1992). Cet estimateur permet d'ajuster les estimations d'enquêtes à des totaux connus de la population, appelés totaux de contrôle, tout en minimisant la distance entre les poids avant et après calage et en tenant compte de certaines contraintes sur les poids finaux. C'est une méthode d'estimation assistée par modèle qui permet de caler selon plusieurs variables simultanément.

Soit Y , la variable d'intérêt et y , le vecteur $1 \times n$ des valeurs observées dans l'échantillon. La notation utilisée est similaire à celle utilisée par Bankier (2002). Soit \hat{x} la matrice $p \times n$ des valeurs observées dans l'échantillon pour les p variables auxiliaires (ou totaux de contrôle). L'estimateur prend la forme

$$\hat{Y} = y \underset{\sim}{diag}(\underset{\sim}{w}) \underset{\sim}{g}$$

où $\underset{\sim}{w}$ est un vecteur $n \times 1$ des poids avant calage, $\underset{\sim}{diag}(\underset{\sim}{w})$ est une matrice diagonale de dimension $n \times n$ avec les poids avant calage $\underset{\sim}{w}$ sur la diagonale et finalement $\underset{\sim}{g}$, qui représente le vecteur de dimension $n \times 1$ des poids-g. Les poids-g sont calculés de sorte que la fonction de perte L soit minimisée tout en respectant la contrainte $\underset{\sim}{x} \underset{\sim}{diag}(\underset{\sim}{w}) \underset{\sim}{g} = \underset{\sim}{X}$ où $\underset{\sim}{X}$ est le vecteur de dimension $p \times 1$ des p totaux de contrôle. La fonction de perte utilisée dans l'EDTR est

$$L = (\underset{\sim}{g} - \underset{\sim}{1}_n)' \underset{\sim}{I}_n (\underset{\sim}{g} - \underset{\sim}{1}_n)$$

où $\underset{\sim}{I}_n$ est la matrice-identité. Les poids-g résultants sont donnés par

$$\underset{\sim}{g} = \underset{\sim}{1}_n + \underset{\sim}{\hat{x}}' (\underset{\sim}{\hat{x}} \underset{\sim}{\hat{x}}')^{-1} (\underset{\sim}{X} - \underset{\sim}{\hat{x}} \underset{\sim}{1}_n)$$

Pour l'EDTR, les autres restrictions suivantes sont appliquées sur les poids finaux: d'abord, le poids transversal est un poids intégré au niveau du ménage, c'est-à-dire qu'à la suite du calage, tous les membres d'un même ménage auront des poids égaux (Lemaître et Dufour, 1987). L'intégration des poids au niveau du ménage permet entre autres de produire des

estimations au niveau des ménages, des entités économiques et au niveau des personnes. Une autre contrainte utilisée lors du calage est que les poids finaux doivent être supérieurs ou égaux à un.

Un grand nombre de totaux de contrôle ont été évalués de façon à sélectionner ceux étant les plus appropriés pour l'enquête. En ce qui a trait aux totaux démographiques, des totaux de contrôle tels que le nombre de personnes selon divers groupes d'âge et de sexe, le nombre de ménages et d'entités économiques selon la taille et selon le type ont été évalués. Pour les totaux de contrôle basés sur le revenu, différentes sources externes ont été considérées, comme par exemple les fichiers T1 et T4 provenant de l'Agence du revenu du Canada. Le fichier T1 contient les données figurant dans la déclaration de revenus et de prestations que produisent les particuliers alors que le fichier T4 contient les données figurant dans l'état de rémunération payée produit par les employeurs. Les nombreuses études effectuées ont démontré que le fichier T4 offrait une meilleure couverture, surtout en ce qui concerne les personnes à faibles revenus (Auger et Tremblay, 2005). Les totaux de contrôle retenus ont donc été ceux du nombre de salariés par classe de salaires et traitements. Les limites des classes sont définies à partir de centiles spécifiques. Pour l'EDTR, les centiles utilisés sont le 10^e, le 25^e, le 50^e, le 65^e, le 75^e et enfin le 95^e, 98^e ou 99^e, selon le nombre d'observations incluses dans la dernière classe. Il est à noter que la limite de classe définie par le 10^e centile a été ajoutée pour l'EDTR en raison de l'importance de la population des faibles revenus pour l'enquête. Le choix des autres centiles a été surtout guidé par le fait qu'une bonne représentativité des hauts revenus était importante pour une meilleure estimation des totaux agrégés.

Dans le cadre du Projet sur le calage harmonisé de la statistique du revenu, des nombreuses études ont été effectuées de façon à identifier la meilleure stratégie de calage, c'est-à-dire celle qui aurait pour effet de réduire les écarts entre certaines estimations-clés provenant de différentes sources tout en limitant l'ampleur de la distorsion dans les poids. Un autre critère important dans le choix de la meilleure stratégie était que celle-ci se devait de minimiser les impacts négatifs sur d'autres estimations d'enquête non contrôlées par le calage. Dans la prochaine section, il est question des changements apportés à la stratégie de calage de l'EDTR au cours du temps et de la nouvelle stratégie implantée cette année.

1.4 Les modifications apportées à la stratégie de calage de l'EDTR

Au cours des dernières années, de nombreuses stratégies de calage ont été analysées et comparées, ce qui a permis d'étudier à fond le processus de calage et ses effets sur les estimations d'enquête. Pour l'EDTR, les améliorations au calage ont été implantées en deux étapes. La première phase du projet a été introduite lors de la diffusion de l'année de référence (AR) 2000, soit en mai 2002. Les changements apportés à la stratégie de calage ne touchaient que les totaux démographiques. Ainsi, des estimations du nombre de ménages et d'entités économiques selon la taille ont été ajoutées aux totaux de contrôle et les groupes d'âge par sexe ont aussi été légèrement modifiés (voir le tableau 1). Seuls les effectifs pour les ménages et entités économiques de taille 1 et 2 ont été retenus dans le calage pour l'EDTR en raison de l'importante distorsion qu'engendrait un plus grand contrôle en terme de nombre de ménages et d'entités économiques selon la taille. Une révision historique a eu lieu et les poids d'enquête ont été révisés rétroactivement jusqu'à l'année 1980. Les totaux de contrôle basés sur les revenus ont été exclus lors de cette première révision car d'autres évaluations de qualité ont été jugées essentielles avant leur introduction dans la stratégie de calage.

Tableau 1 : Évolution des totaux de contrôle de l'EDTR

Totaux de contrôle	Avant révision		Révision historique de 2002 (AR 2000)		Révision historique de 2005 (AR 2003)	
	Groupes d'âge et de sexe	0-6	45-54	0-6	45-54	0-6 (tous)
7-15		55-59	7-15	55-59	7-17	55-64
16-18		60-64	16-17	60-64	18-24	65+
19-24		65-69	18-24	65-69	25-34	
25-34		70+	25-34	70+	35-44	
35-44			35-44			
Nombre de ménages selon la taille	X		1, 2		1, 2	
Nombre d'entités économiques selon la taille	X		1, 2		1, 2	
Nombre de salariés par classes de salaires	X		X		7 classes définies selon les 10 ^e , 25 ^e , 50 ^e , 65 ^e , 75 ^e , 95 ^e ou 98 ^e ou 99 ^e centiles	
Période affectée par la révision	-		1980 à 2000		1990 à 2003	

Suite aux études effectuées par rapport à la qualité des totaux de contrôle basés sur les revenus et à la mise en place d'une stratégie de vérification de la qualité satisfaisante, la deuxième phase du projet d'harmonisation du calage a été implantée

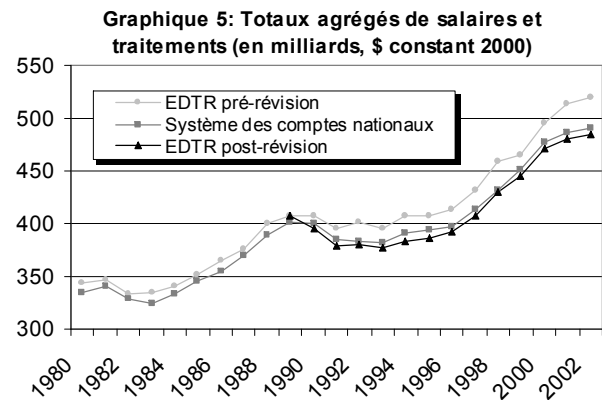
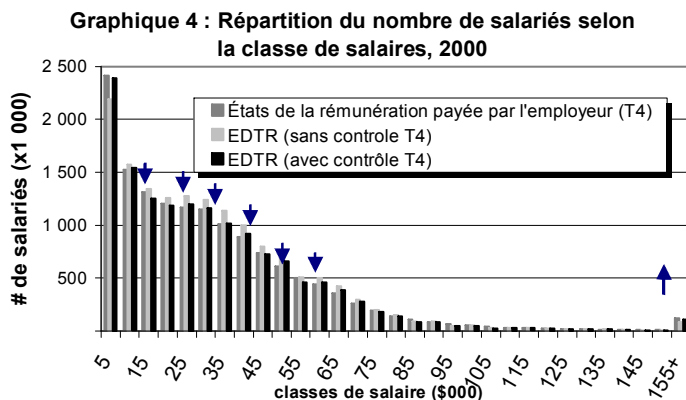
lors de la diffusion des données de l'EDTR pour l'année de référence 2003, soit en mai 2005. La plus importante modification apportée au calage lors de cette révision est celle de l'ajout des totaux de contrôle basés sur le nombre de salariés par classe de salaire. Le nombre de groupes d'âge a aussi été réduit afin de diminuer le nombre total de contrôles utilisés et ainsi limiter la distorsion des poids lors du calage. Étant donnée l'incertitude quant à la qualité des fichiers T4 pour les années antérieures à 1990, cette révision historique des poids n'a été appliquée qu'à partir de l'année 1990 et pour les années ultérieures. Le tableau 1 présente l'évolution des contrôles utilisés selon les différentes révisions qui ont eu lieu pour l'EDTR.

2. IMPACTS SUR LES ESTIMATIONS DE L'EDTR

2.1 Impacts de la nouvelle stratégie de calage sur les estimations d'enquête

Pour analyser les impacts de la nouvelle stratégie de calage, de nombreuses estimations de l'EDTR avant et après la révision ont été calculées et comparées. Dans la mesure du possible, des estimations sur le revenu ainsi que des estimations non reliées au revenu ont été examinées et comparées à des estimations de sources externes. Les résultats obtenus vont dans la direction attendue lorsque les contrôles basés sur les revenus sont utilisés dans le calage.

En ce qui a trait à l'impact de la nouvelle stratégie de calage sur la répartition des salariés dans les différentes classes de salaires et traitements, le graphique 4 montre que les nouvelles estimations sont maintenant plus proches de celles provenant des états de la rémunération payée par les employeurs (Fichier T4). Ainsi, la surestimation des personnes dans les classes du milieu de la répartition des salaires et traitements est réduite ainsi que la sous-estimation du nombre de personnes dont les salaires et traitements se situent aux deux extrémités de la répartition. En terme du total agrégé, les valeurs de salaires et traitements agrégés tirées de l'EFC et de l'EDTR se rapprochent beaucoup plus des estimations provenant du Système de comptabilité nationale tel que montré dans le graphique 5. La surestimation qui existait a été éliminée et fait maintenant place à une légère sous-estimation des totaux agrégés. Il est à noter que celle-ci est due principalement au fait que les estimations de l'EDTR sont tirées d'un échantillon et qu'il est donc très difficile de bien représenter les individus ayant des revenus extrêmement élevés. Or les revenus de ces derniers constituent une part non négligeable des salaires et traitements agrégés (Lathe, 2005).



D'autres estimations ont aussi été analysées et comparées. Ainsi, on observe une baisse du revenu total moyen alors que les montants versés pour les prestations d'aide sociale et le taux de chômage sont à la hausse (voir le tableau 2).

En ce qui a trait à la prévalence du faible revenu, les résultats sont aussi tels qu'attendus, c'est-à-dire que la nouvelle stratégie de calage engendre une hausse de la proportion des individus et familles sous le seuil de faible revenu par rapport aux anciennes estimations de l'EDTR, tel que le montre le tableau 2. Il est à noter que les hausses les plus importantes au niveau provincial ont été observées en Colombie-Britannique.

Tableau 2 : Comparaison des estimations de l'EDTR

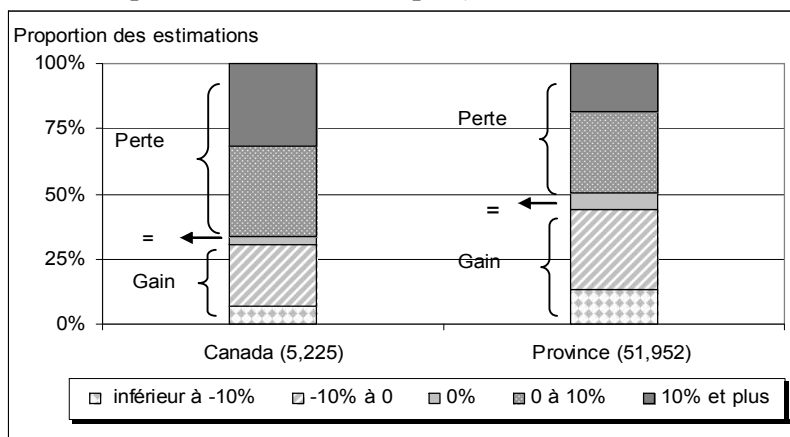
	Variables	Source externe ⁴	EDTR 2000		
			Estimation (écart-type)		
			Pré révision	Post révision	
Canada	Revenu individuel moyen (\$)	22 615	22 842 (174)	22 427 (183)	
	Total agrégé d'aide sociale (M\$)	9 306	7 623 (271)	8 493 (328)	
	Taux de chômage en avril (%)	7,2	6,0 (0,19)	6,4 (0,21)	
	Pourcentage d'immigrants (%)	18,2	19,0 (0,25)	18,9 (0,22)	
	Prévalence du faible revenu (%)	Familles économiques de taille 2 et plus	-	7,9 (0,28)	9,0 (0,28)
Personnes seules		-	28,6 (0,71)	32,9 (0,73)	
Familles monoparentales ayant une femme à la tête		-	33,9 (1,77)	36,3 (1,86)	
Colombie-Britannique	Prévalence du faible revenu (%)	Familles économique de taille 2 et plus	-	9,2 (0,78)	11,3 (0,84)
		Personnes seules	-	28,4 (1,93)	34,4 (2,13)
		Familles monoparentales ayant une femme à la tête	-	27,1 (4,66)	33,3 (5,16)

2.2 Impacts de la nouvelle stratégie de calage sur la précision des estimations

De façon à évaluer si la nouvelle stratégie engendrait un gain de précision des estimations de l'enquête, des comparaisons des erreurs types et des coefficients de variation ont été faites pour la majorité des estimations publiées lors de la diffusion annuelle de l'EDTR. Cette comparaison a été effectuée pour différents types d'estimations sur le revenu tels que des moyennes, totaux, pourcentages, ratios, et autres, et ce pour différents niveaux géographiques. Les résultats obtenus sont plutôt décevants puisque le gain présumé en terme de précision ne s'est pas concrétisé (voir graphique 6). En effet, au niveau national, seulement 30% des estimations ont subi un gain en terme de précision. Pour ce qui est des estimations au niveau provincial, seulement la moitié des estimations ont une meilleure précision malgré le fait que le calage est effectué à ce niveau géographique.

Malgré ces résultats plutôt décevants en ce qui a trait à la précision des estimations, la réduction importante du biais entraîne tout de même une diminution de l'erreur quadratique moyenne, ce qui en soi constitue une nette amélioration.

Graphique 6 : Proportion d'estimation selon la différence relative entre les écart types, avant et après la révision historique (basé sur 57 000 estimations)



3. AUTRES CONSIDÉRATIONS

3.1 Interaction entre les totaux de contrôle de sources différentes

La nouvelle stratégie de calage de l'EDTR utilise pour la première fois une combinaison de totaux de contrôle provenant de sources externes différentes, c'est-à-dire les totaux de contrôle dérivés à partir du recensement de 2001 et ceux provenant du fichier T4. L'utilisation conjointe de ces deux sources fait en sorte que le nombre de non-salariés est indirectement contrôlé

⁴ Fichier T1 provenant de l'Agence du revenu du Canada, Système de comptabilité nationale
Taux de chômage selon l'Enquête sur la population active, taux d'immigration selon le recensement de 2001.

par des contrôles provenant de deux sources. Or, les totaux de contrôle provenant des fichiers T4 sont fixes alors que les totaux démographiques sont sujets à une révision environ tous les cinq ans. Ainsi, lors d'une révision historique et d'une mise à jour des projections démographiques, le nombre de non-salariés peut changer à la hausse ou à la baisse en fonction d'une sous-estimation ou surestimation potentielle de la population par les projections démographiques. Or, les non-salariés possèdent des caractéristiques différentes du reste de la population, par exemple, on y compte une proportion plus élevée de chômeurs ou de personnes récipiendaires de prestations d'aide sociale. Un changement dans le nombre de non-salariés, causé par une mise à jour des totaux démographiques, aura donc des conséquences sur des estimations telles que l'incidence du faible revenu, les totaux agrégés de prestations d'aide sociale, etc. alors qu'auparavant, les impacts dus aux changements dans les projections démographiques étaient peu importants.

3.2 Conclusion

Suite aux modifications apportées à la stratégie de calage de l'EDTR et à l'importante révision historique des données en mai 2005, les estimations tirées de l'EDTR sont maintenant plus proches des estimations provenant des autres enquêtes et de sources externes. Bien que le calage soit un moyen efficace pour corriger des biais dans les échantillons et pouvant être appliqué rétroactivement, il n'en demeure pas moins qu'il est essentiel de continuer à améliorer les méthodes de collecte de façon à réduire les biais dans l'échantillon, sans recourir au calage de façon excessive. Il faut garder en mémoire que le calage doit généralement être utilisé pour des réajustements mineurs de l'échantillon, afin d'éviter d'importantes distorsions de poids ou d'impacts négatifs sur d'autres estimations de l'enquête.

Dans les années qui suivront, il sera aussi important de poursuivre les évaluations de qualité des contrôles utilisés dans le calage ainsi que les impacts de la nouvelle stratégie sur d'autres estimations de l'enquête. Il sera aussi important de considérer d'autres contrôles potentiels pouvant provenir d'autres sources externes et ainsi continuer d'améliorer les estimations de l'EDTR.

REMERCIEMENTS

L'auteur tient à remercier Michel Latouche, Richard Laroche, Mylène Lavigne et Caroline Cauchon pour leurs précieux commentaires.

RÉFÉRENCES

Auger, S. (2005). « Évaluation de la qualité des estimations démographiques et des données fiscales utilisées dans le calage de différentes enquêtes à Statistique Canada ». *Recueil 2005 de la section des méthodes d'enquête, Société statistique du Canada*.

Bankier, M. (2002). « Canadian Census Weighting ». *Rapport technique présenté à la 35^e réunion du Comité consultatif sur les méthodes statistiques, Statistique Canada*. Novembre 2002.

Lathe, H. (2005). « Enquête sur la Dynamique du travail et du revenu : révision historique de 2003 ». *Série de document de recherche – Revenu, Statistique Canada*. Numéro au catalogue 75F0002MIF2005009, juillet 2005.

Lemaître, G. et Dufour, J. (1987). « Une méthode intégrée de pondération des personnes et des familles ». *Techniques d'enquête*, vol. 13, no 2, 211-220.

Sarndal et coll. (1992). *Model Assisted Survey Sampling*. New York: Springer-Verlag

Tremblay, J. (2005). « Aperçu de la stratégie de calage harmonisé des statistiques du revenu de Statistique Canada ». *Recueil 2005 de la section des méthodes d'enquête, Société statistique du Canada*.