

IMPACT D'ERREURS D'APPARIEMENT SUR L'INFÉRENCE STATISTIQUE LORS D'ÉTUDES DE MORTALITÉ PAR COHORTES

D. Krewski^{1, 2, 4}, Y. Wang³, S. Bartlett³, J.M. Zielinski³ et R. Mallick^{1,2}

RÉSUMÉ

Grâce aux méthodes de couplage d'enregistrements, il est désormais plus facile d'effectuer des études de mortalité par cohortes où il y a couplage électronique des données d'exposition d'une base d'information et des données de mortalité d'une autre base. Le présent article est consacré à l'incidence des erreurs de couplage sur les estimations d'indicateurs épidémiologiques de risque comme les taux comparatifs de mortalité et les paramètres de modèles de régression de risques relatifs. Il démontre que ces indicateurs peuvent être entachés d'un biais et d'un surcroît de variabilité à cause d'erreurs de couplage, les faux liens et les non-liens introduisant un biais respectivement positif et négatif dans les estimations de taux normaliser de mortalité. Ces erreurs accroissent toujours l'incertitude des estimations, mais on peut réussir à éliminer le biais dans le cas particulier d'une égalité des faux positifs et des faux négatifs pour des états homogènes définis par un croisement des covariables d'intérêt.

MOTS CLÉS : Étude par cohortes; Couplage informatisé d'enregistrements; Erreurs de couplage; Valeurs seuils de pondération de couplage; Régression de Poisson; Régression de risques relatifs; Taux comparatif de mortalité

1. INTRODUCTION

Ces dernières années, on a effectué diverses études chronologiques par cohortes en épidémiologie environnementale en exploitant les bases de données administratives en place comme sources d'information (Howe et Spasoff, 1986; Carpenter et Fair, 1990). En général, il s'agit de coupler les enregistrements sur l'exposition humaine aux risques du milieu et les enregistrements sur l'état de santé, ce qui souvent se fait à l'aide de méthodes informatisées de mise en correspondance individuelle d'enregistrements appartenant à des bases de données différentes. Dans une étude de mortalité par cohortes, on détermine l'état civil de chaque membre d'une cohorte par couplage avec les dossiers de mortalité que tiennent les organismes publics. Un excès de mortalité dans une cohorte par rapport à l'ensemble de la population peut s'expliquer par le degré d'exposition de ses membres.

Plus précisément, il y a couplage d'enregistrements lorsqu'on met en correspondance deux éléments d'information ou plus pour une même entité (Bartlett et coll., 1993). Les méthodes de couplage informatique d'enregistrements (CIE) sont aujourd'hui des plus raffinées. Elles font appel à des algorithmes complexes pour évaluer les probabilités de juste appariement de deux enregistrements (Hill, 1988; Newcombe, 1988).

¹Centre R. Samuel McLaughlin d'évaluation des risques pour la santé des populations, Université Carleton, Ottawa, Ontario, Canada, K1N 6N5.

²School of Mathematics and Statistics, Université Carleton, Ottawa, Ontario, Canada, K1S 5B6.

³Direction générale de la santé environnementale et sécurité des consommateurs, Ottawa, Ontario, Canada, K1A 0L2.

⁴C'est la personne avec qui on doit correspondre.

Statistique Canada a mis au point un système CIE appelé CANLINK qui peut traiter les mises en correspondance intrafichiers et interfichiers (Howe et Lindsay, 1981; Smith et Silins, 1981). Dans ce système, on attribue des valeurs de pondération à des paires d'enregistrements en fonction des probabilités d'appariement. On fixe deux seuils. Les appariements ayant une valeur de pondération plus grande que le seuil supérieur sont considérés comme de bons liens. Ceux dont la valeur de pondération se situe en deçà du seuil inférieur sont considérés comme des non-liens. On résout les cas d'appariement possible dont la valeur de pondération est comprise entre les seuils supérieure et inférieure en employant tout complément d'information si on en a, sinon on fixe une valeur seuil unique qui permet de distinguer les liens des non-liens.

On veille jalousement sur la confidentialité des enregistrements protégés par la *Loi sur la statistique* dans toute étude où il y a couplage d'enregistrements. Toutes les études prévoyant une mise en correspondance avec des bases de données protégées doivent satisfaire à de rigoureux critères d'examen et d'approbation avant de se réaliser. Tous les fichiers en couplage renfermant des données d'identification restent sous la garde de Statistique Canada (Labossière, 1986).

On a appliqué des méthodes informatisées de couplage d'enregistrements pour mettre en correspondance les données d'exposition au milieu et la Base canadienne de données sur la mortalité (BCDM). Ainsi, on a entrepris une étude sur les exploitants agricoles canadiens qui vise à dégager les éventuels rapports entre les causes de décès et diverses variables sociodémographiques et agricoles, et plus particulièrement l'utilisation de pesticides, chez plus de 326 000 de ces exploitants (Jordan-Simpson et coll., 1990). Dans cet exercice, on a mis en couplage la BCDM et les recensements de la population et de l'agriculture de 1971. Une autre étude permanente à grande échelle s'appuie sur le Registre national des doses (RND) du Canada (Ashmore et Grogan, 1985; Ashmore et Davies, 1989). Ce registre renferme des données sur l'exposition professionnelle aux rayonnements ionisants chez plus de 400 000 Canadiens depuis 1950. On a récemment effectué le couplage RND-BCDM afin d'examiner les liens entre l'excès de mortalité par cancer et l'exposition professionnelle à de faibles rayonnements ionisants (Ashmore et coll., 1997, 1998). Plus récemment encore, il y a eu raccordement du RND et de la Base canadienne de données sur le cancer (Sont et coll., 2001). Fair (1989) a dressé une liste complète des autres études sanitaires où il y a mise en correspondance de données d'exposition et de l'information de la BCDM.

Les études par couplage d'enregistrements sont avantageuses à plusieurs égards par rapport aux études épidémiologiques classiques. Le recours aux bases de données administratives en place obvie à la nécessité de recueillir de nouvelles données aux fins des études sur la santé, et il est souvent possible d'obtenir assez aisément de grandes tailles d'échantillon. Selon la nature des bases d'information exploitées, le couplage peut être un moyen économique d'examen de nombreux rapports possibles dans des études épidémiologiques. Le couplage d'enregistrements présente aussi certains inconvénients. Ainsi, on est généralement peu maître des données recueillies et le suivi peut être impossible dans un nombre appréciable de cas. Un autre inconvénient est l'apparition d'erreurs de couplage, qui est ici notre propos même. Inévitablement, des enregistrements qui se correspondent ne seront pas couplés et d'autres qui ne se correspondent pas seront couplés par erreur.

On s'est relativement peu employé à établir l'incidence de ces erreurs sur l'inférence statistique. Neter et coll. (1965) se sont servis d'un simple modèle de régression linéaire pour analyser l'incidence des erreurs introduites par le couplage. Leurs résultats indiquent que ces erreurs viennent gonfler la variance résiduelle et introduire un biais dans le paramètre de pente estimé. Winkler et Scheuren (1991) ont exprimé le biais en question dans des estimations de coefficients de régression linéaire. Les progrès de l'estimation des taux d'erreur de couplage grâce aux travaux de Belin et Rubin (1991) ont permis à Scheuren et Winkler (1993) de mettre au point une méthode améliorée de correction de biais.

Notre propos sera de cerner l'incidence des erreurs de couplage sur l'inférence statistique dans des études de mortalité par cohortes. Nous décrirons à la section 2 les modèles de régression de risques relatifs ayant servi à l'analyse des données et établirons des expressions pour le nombre observé et prévu de décès par ces modèles. À la section 3, nous examinerons l'effet des erreurs de couplage sur le dénombrement observé et prévu de décès et d'années-personnes d'exposition. À la section 4, nous analyserons l'incidence de ces mêmes erreurs sur les

estimations de taux normalisés de mortalité (TNM) et de paramètres de régression de risques relatifs. L'un et l'autre de ces types d'erreurs peuvent introduire un biais et un surcroît de variabilité dans les estimations de paramètres. Enfin, nous livrerons nos conclusions à la section 5.

2. MODÈLES DE RÉGRESSION DE RISQUES RELATIFS

Les méthodes statistiques d'analyse sont bien établies dans le cadre des études de mortalité par cohortes (Breslow et Day, 1987). Une telle analyse vise principalement à déterminer si l'exposition à l'agent d'intérêt accroît le taux de mortalité chez les membres d'une cohorte. La mortalité est caractérisée par la fonction de risque, qui spécifie le taux de mortalité comme fonction du temps. Soit T le moment du décès. La fonction de risque au moment u est alors formellement définie de la manière suivante :

$$\lambda(u) = \lim_{\Delta u \downarrow 0} \Pr \{u \leq T < u + \Delta u | T \geq u\}. \quad (1)$$

Soit $\lambda_i(u)$ la fonction de risque pour une cause particulière de décès au moment u du membre $i=1, \dots, N$ d'une cohorte de taille N . Soit $z_i(u)$ un vecteur correspondant de covariables particulières à ce membre. Nous supposons que ces covariables ont pour effet de modifier le risque de référence $\lambda^*(u)$ suivant le modèle de régression de risques relatifs

$$\lambda_i(u) = \lambda^*(u) \gamma \{\beta' z_i(u)\}, \quad (2)$$

où γ est une fonction positive des covariables et β , un vecteur de paramètres de régression.

Les modèles multiplicatif et additif sont deux cas d'espèce du modèle général de régression de risques relatifs. On peut définir la fonction γ dans (2) par

$$\log \gamma(z) = \frac{(1+z)^\rho - 1}{\rho}. \quad (3)$$

Lorsque $D=1$, le modèle général se ramène au modèle multiplicatif

$$\lambda_i(u) = \lambda^*(u) \exp\{\beta' z_i(u)\}. \quad (4)$$

Le modèle des risques proportionnels, qui a été introduit par Cox (1972), sert largement à l'analyse des données de mortalité (Kalbfleisch et Prentice, 1980). Le modèle additif

$$\lambda_i(u) = \lambda^*(u) + \beta' z_i(u) \quad (5)$$

se présente comme un cas limite lorsque $D=0$.

Soit t_i^0 et t_i^1 les âges respectifs à l'entrée dans l'étude et au moment de l'impossibilité du suivi (par fin de l'étude ou encore retrait ou décès) du i^{e} membre de la cohorte. Soit $\delta_i = 1$ ou 0 selon que le i^{e} membre était

décédé ou non au moment de l'impossibilité du suivi. On peut ainsi écrire la fonction de vraisemblance en expression logarithmique selon le modèle des risques relatifs (2) :

$$\log L = \sum_{i=1}^N \left\{ \delta_i \log(\gamma\{\beta' z_i(t_i^1)\}) - \int_{t_i^0}^{t_i^1} \gamma\{\beta' z_i(u)\} \lambda^*(u) du \right\}. \quad (6)$$

Si on a une seule covariable $z_i(u) \equiv 1$, l'estimation de maximum de vraisemblance de $\theta = \exp\{\beta\}$ se ramène au taux normalisé de mortalité $TNM = OBS/EXP$, où $OBS = \sum_{i=1}^N \delta_i$ et $EXP = \sum_{i=1}^N e_i$ sont respectivement le nombre observé et le nombre prévu de décès avec $e_i = \int_{t_i^0}^{t_i^1} \lambda^*(u) du$.

Dans le cas de grandes tailles d'échantillon, la maximisation de la fonction de vraisemblance (6) peut être d'un calcul laborieux. Breslow et coll. (1983) simplifient le calcul en posant que les covariables prennent des valeurs constantes dans les états où se trouve successivement un membre de la cohorte pendant que dure l'étude. Les états en question sont définis par classement recoupé des covariables d'intérêt. Plus précisément, on pose qu'il y a J états $\{S_j; j = 1, \dots, J\}$, de sorte que $z_i(u) = z_j$ chaque fois que le i^e membre est dans l'état S_j au moment u . Ces états s'excluent les uns les autres et il n'y en a pas d'autres, aussi chaque membre de la cohorte se trouve-t-il dans un état et un seul à tout moment u . On peut alors écrire la fonction de vraisemblance en expression logarithmique (6) de la manière suivante :

$$\log L = \sum_{j=1}^J \left\{ d_j \log(\gamma\{\beta' z_j\}) - \gamma\{\beta' z_j\} e_j \right\}, \quad (7)$$

où

$$e_j = \sum_{i=1}^N \int_{[z_i(u) \in S_j]} \lambda^*(u) du \quad (8)$$

est la contribution au nombre prévu de décès de toutes les années-personnes d'observation dans l'état S_j , et où d_j désigne le nombre total de décès dans cet état. Pour $\Lambda_j(\beta) = \log(\gamma\{\beta' z_j\})$, l'estimation de maximum de vraisemblance $\hat{\beta}$ de β s'obtient comme la solution de l'équation de caractérisation

$$\frac{\partial \log L}{\partial \beta} = \sum_{j=1}^J \frac{\partial \Lambda_j(\hat{\beta})}{\partial \beta} \left\{ d_j - \exp\{\Lambda_j(\hat{\beta})\} e_j \right\} = 0 \quad (9)$$

3. INCIDENCE DES ERREURS DE COUPLAGE SUR LE NOMBRE OBSERVÉ ET PRÉVU DE DÉCÈS

Deux types d'erreurs peuvent se produire en cas de couplage informatisé d'enregistrements (CIE) de fichiers de données (Fellegi et Sunter, 1969). Il y a faux positif lorsqu'on considère à tort comme décédé un membre de la cohorte qui est vivant. Il y a faux négatif lorsqu'on considère comme toujours vivant un sujet décédé. Dans cette section, nous examinerons l'incidence des erreurs de couplage sur le nombre observé et prévu de décès respectivement. Pour ce faire, nous définissons d'abord des jeux d'indices intra-états qui représenteront des ensembles d'enregistrements qui ont été bien ou mal couplés.

3.1 Erreurs de couplage

Soit C_j l'ensemble de descripteurs des membres de la cohorte qui passent par l'état S_j . Dans l'hypothèse d'absence d'erreurs de couplage, soit A_j et D_j les ensembles respectifs de descripteurs des membres qui restent vivants tout au long de l'état S_j et qui décèdent dans ce même état. À noter que $A_j \cup D_j = C_j$ et $A_j \cap D_j = \{\emptyset\}$.

Soit A_j^L et D_j^L les ensembles de descripteurs A_j et D_j dans l'hypothèse de présence d'erreurs de couplage. Nous définissons en outre D_j^P et A_j^N comme les ensembles d'indices de faux positifs et de faux négatifs dans le j^{e} état S_j . Ces ensembles satisfont aux relations $A_j^L = (A_j - D_j^P) \cup A_j^N$ et $D_j^L = (D_j - A_j^N) \cup D_j^P$, où $A_j \cup D_j = C_j$ et $A_j^L \cap D_j^L = \{\emptyset\}$.

Nous pouvons décrire l'incidence des erreurs de couplage sur la fonction de vraisemblance dans (7). Soit t_{ij}^0 , t_{ij}^1 et t_{ij}^2 les moments où le i^{e} membre entre, meurt (si la mort survient) et/ou sort du $j^{\text{ième}}$ état S_j . Par (8) et les décompositions de A_j^L et D_j^L , on peut formuler de la manière suivante le nombre prévu de décès e_j^L dans le j^{e} état en cas d'erreurs de couplage :

$$\begin{aligned}
 e_j^L &= \sum_{i \in A_j^L} \int_{t_{ij}^0}^{t_{ij}^2} \lambda^*(u) du + \sum_{i \in D_j^L} \int_{t_{ij}^0}^{t_{ij}^1} \lambda^*(u) du \\
 &= \left(\sum_{i \in A_j} \int_{t_{ij}^0}^{t_{ij}^2} + \sum_{i \in A_j^N} \int_{t_{ij}^0}^{t_{ij}^2} - \sum_{i \in D_j^P} \int_{t_{ij}^0}^{t_{ij}^2} \right) \lambda^*(u) du + \\
 &\quad \left(\sum_{i \in D_j} \int_{t_{ij}^0}^{t_{ij}^1} + \sum_{i \in D_j^P} \int_{t_{ij}^0}^{t_{ij}^1} - \sum_{i \in A_j^N} \int_{t_{ij}^0}^{t_{ij}^1} \right) \lambda^*(u) du \tag{10} \\
 &= \sum_{i \in A_j} \int_{t_{ij}^0}^{t_{ij}^2} \lambda^*(u) du + \sum_{i \in D_j} \int_{t_{ij}^0}^{t_{ij}^1} \lambda^*(u) du + \sum_{i \in A_j^N} \int_{t_{ij}^0}^{t_{ij}^2} \lambda^*(u) du - \sum_{i \in D_j^P} \int_{t_{ij}^0}^{t_{ij}^1} \lambda^*(u) du \\
 &= e_j - \Delta e_j,
 \end{aligned}$$

où

$$e_j = \sum_{i \in A_j} \int_{t_{ij}^0}^{t_{ij}^2} \lambda^*(u) du + \sum_{i \in D_j} \int_{t_{ij}^0}^{t_{ij}^1} \lambda^*(u) du \tag{11}$$

et

$$\Delta e_j = e_j^P - e_j^N \tag{12}$$

avec

$$e_j^P = \sum_{i \in D_j^P} \int_{t_{ij}^1}^{t_{ij}^2} \lambda^*(u) du \quad (13)$$

et

$$e_j^N = \sum_{i \in A_j^N} \int_{t_{ij}^1}^{t_{ij}^2} \lambda^*(u) du. \quad (14)$$

Le terme Δe_j représente le biais dont est entaché le nombre prévu de décès dans le j° état à cause d'erreurs de couplage. Il découle de (10) et (12) que les faux positifs et les faux négatifs tendent respectivement à diminuer et à augmenter le nombre prévu de décès.

Par la décomposition de D_j^L , on peut formuler de la manière suivante le nombre observé de décès d_j^L en cas d'erreurs de couplage :

$$d_j^L = d_j + \Delta d_j, \quad (15)$$

où

$$\Delta d_j = d_j^P - d_j^N. \quad (16)$$

Ici, $d_j = \dim(D_j)$, $d_j^P = \dim(D_j^P)$ et $d_j^N = \dim(A_j^N)$, $\dim(A)$ désignant le nombre d'éléments de l'ensemble A . Le terme Δd_j représente le biais dont est entaché le nombre observé de décès dans le j° état à cause d'erreurs de couplage. Il découle de (15) et (16) que les faux positifs et les faux négatifs augmenteront et diminueront respectivement le nombre observé de décès.

On établit souvent l'état civil par couplage avec la BCDM, fonds d'information qui est généralement plus vaste que ce qui se rapporte directement à la cohorte d'intérêt. Si on relie à tort les données d'exposition d'une personne vivante à celles d'une personne décédée, cette dernière ne fera habituellement pas partie de la cohorte. Ainsi, les années-personnes d'exposition qu'apporte cette personne se termineront de façon prématurée dans l'année de son présumé décès; les années-personnes d'exposition perdues correspondent au temps s'écoulant entre le présumé décès et la fin du suivi. Par ailleurs, si on relie à tort les données d'exposition d'une personne décédée à celles d'une personne vivante, les années-personnes d'exposition qu'apporte cette dernière augmenteront du nombre d'années comprises entre l'année de son décès effectif et la fin de la période de suivi. Ainsi, les faux positifs et les faux négatifs dégonflent et gonflent respectivement le nombre d'années-personnes d'exposition de la cohorte.

3.2 Espérances et variances des biais du nombre observé et prévu de décès

L'incidence des erreurs de couplage sur le nombre observé et prévu de décès dépend des taux de fausse positivité (falsipositivité) et de fausse négativité (falsinégativité). Soit D_{ij} (\overline{D}_{ij}) l'événement du décès (ou du non-décès) du i° membre de la cohorte dans le j° état et L_{ij} (\overline{L}_{ij}) l'événement de caractérisation de ce membre comme étant décédé (n'étant pas décédé) à la suite du couplage avec la BCDM. Comme Fellegi et Sunter (1969), nous définissons le taux de fausse positivité comme la probabilité conditionnelle $p_{ij}^I = \Pr \{L_{ij} | \overline{D}_{ij}\}$.

De même, le taux de fausse négativité est $p_{ij}^{\text{II}} = \Pr \left\{ \overline{L}_{ij} \mid D_{ij} \right\}$. Par souci de simplicité, nous supposons que les taux de fausse positivité et de fausse négativité sont constants dans chaque état $j = 1, \dots, J$. Plus précisément, $p_{ij}^{\text{I}} = p_j^{\text{I}}$ et $p_{ij}^{\text{II}} = p_j^{\text{II}}$ pour $i \in C_j$. L'hypothèse conviendra chaque fois que les sujets se trouvant dans le même état sont hautement homogènes, plus particulièrement pour des attributs comme la qualité des descripteurs individuels qui influent sur les taux d'erreur de couplage. Il est improbable que l'on satisfasse pleinement dans la pratique à cette hypothèse idéale, mais celle-ci permet de considérablement simplifier l'évaluation subséquente des effets des erreurs de couplage.

Soit ξ_{ij} une variable indicatrice, $\xi_{ij} = 1$ s'il y a faux positif pour le i^{e} sujet dans le j^{e} état et $\xi_{ij} = 0$ dans les autres cas. De même, soit $\psi_{ij} = 1$ dans le cas d'un faux négatif pour le même sujet et $\psi_{ij} = 0$ dans les autres cas. À supposer que les erreurs de couplage relatives aux divers sujets soient indépendantes, l'espérance et la variance du biais du nombre observé de décès dans le j^{e} état Δd_j sont

$$E\{d_j^{\text{P}} - d_j^{\text{N}}\} = E\left\{ \sum_{i \in C_j} \xi_{ij} - \sum_{i \in C_j} \psi_{ij} \right\} = \sum_{i \in C_j} \Delta p_{ij} = \Delta p_j N_j \quad (17)$$

et

$$\begin{aligned} \text{Var}\{d_j^{\text{P}} - d_j^{\text{N}}\} &= \sum_{i \in C_j} \text{Var}\{\xi_{ij}\} + \sum_{i \in C_j} \text{Var}\{\psi_{ij}\} - \sum_{i \in C_j} \text{Cov}\{\xi_{ij}, \psi_{ij}\} \\ &= \sum_{i \in C_j} \{p_{ij}^{\text{I}}(1 - p_{ij}^{\text{I}}) + p_{ij}^{\text{II}}(1 - p_{ij}^{\text{II}})\} \\ &= \{p_j^{\text{I}}(1 - p_j^{\text{I}}) + p_j^{\text{II}}(1 - p_j^{\text{II}})\} N_j, \end{aligned} \quad (18)$$

où $\Delta p_{ij} = p_{ij}^{\text{I}} - p_{ij}^{\text{II}}$ et $\Delta p_j = p_j^{\text{I}} - p_j^{\text{II}}$. À noter que, comme il ne peut y avoir simultanément de faux positif et de faux négatif,

$$\text{Cov}\{\xi_{ij}, \psi_{ij}\} = E\{(\xi_{ij} - \psi_{ij})^2\} - (E\{\xi_{ij} - \psi_{ij}\})^2 = 0. \quad (19)$$

De même, l'espérance et la variance du biais du nombre prévu de décès dans le j^{e} état Δe_j peuvent s'exprimer de la manière suivante :

$$E\{e_j^{\text{P}} - e_j^{\text{N}}\} = \Delta p_j \sum_{i \in C_j} \int_{\min(t_{ij}^1, t_{ij}^2)}^{t_{ij}^2} \lambda^*(u) du \quad (20)$$

et

$$\text{Var}\{e_j^{\text{P}} - e_j^{\text{N}}\} = \{p_j^{\text{I}}(1 - p_j^{\text{I}}) + p_j^{\text{II}}(1 - p_j^{\text{II}})\} \sum_{i \in C_j} \left\{ \int_{\min(t_{ij}^1, t_{ij}^2)}^{t_{ij}^2} \lambda^*(u) du \right\}^2, \quad (21)$$

Les résultats indiquent que, si les taux de fausse positivité et de fausse négativité sont égaux dans chaque état, le nombre observé de décès sera égal en moyenne au nombre effectif de décès, auquel cas le biais dont est

entaché le nombre prévu de décès par erreurs de couplage disparaîtra. Il reste que les deux types d'erreurs de couplage introduiront un surcroît de variation dans le nombre observé et prévu de décès.

Il est difficile de minimiser la variance du biais dans (21), puisque p_j^I et p_j^{II} ne sont pas fonctionnellement indépendants. En règle générale, décroître p_j^I , c'est accroître p_j^{II} et vice versa. On pourrait néanmoins minimiser le facteur multiplicatif $\{p_j^I(1-p_j^I)+p_j^{II}(1-p_j^{II})\}$ dans (21) compte tenu de la relation fonctionnelle entre p_j^I et p_j^{II} . Bien que ce facteur soit indépendant du modèle sous-jacent de régression de risques relatifs γ dans (2), l'erreur quadratique moyenne obtenue par une combinaison de (20) et (21) ne saurait être minimisée sans une spécification du risque de référence $\lambda^*(u)$.

4. INCIDENCE DES ERREURS DE COUPLAGE SUR LES ESTIMATIONS DES TAUX NORMALISÉS DE MORTALITÉ ET DES COEFFICIENTS DE RÉGRESSION

4.1 Taux normalisés de mortalité

Pour déterminer l'incidence des erreurs de couplage sur les taux normalisés de mortalité (TNM), nous remplaçons les nombres et observé et prévu de décès d_j et e_j dans une hypothèse d'absence d'erreurs de couplage par les nombres et observé et prévu de décès d_j^L et e_j^L dans une hypothèse de présence de telles erreurs dans l'expression $TNM = \sum d_j / \sum e_j$. Soit TNM_L les taux normalisés de mortalité en cas d'erreurs de couplage. Nous avons alors

$$TNM_L = TNM \left[1 + \frac{\sum \Delta d_j}{\sum d_j} \right] / \left[1 - \frac{\sum \Delta e_j}{\sum e_j} \right]. \quad (22)$$

Il s'ensuit que les faux positifs et les faux négatifs augmenteront et diminueront respectivement les TNM.

On peut ainsi exprimer la différence $\Delta TNM = TNM_L - TNM$:

$$\frac{\Delta TNM}{TNM} = \frac{\sum \Delta d_j}{\sum d_j} + \frac{\sum \Delta e_j}{\sum e_j} + o\left(\frac{\sum \Delta e_j}{\sum e_j}\right). \quad (23)$$

Dans la plupart des études par cohortes, le nombre de décès est bien moindre que la taille de l'ensemble de la cohorte, auquel cas le terme d'erreur est petit dans (23). On peut approcher la moyenne et la variance du biais relatif du TNM par

$$E \left\{ \frac{\Delta TNM}{TNM} \right\} \approx \sum \Delta p_j \left\{ \left(\sum d_j \right)^{-1} N_j + \left(\sum e_j \right)^{-1} \sum_{i \in C_j} \int_{\min(t_{ij}^1, t_{ij}^2)}^{t_{ij}^2} \lambda^*(u) du \right\} \quad (24)$$

et

$$\begin{aligned} \text{Var} \left\{ \frac{\Delta TNM}{TNM} \right\} &\approx \sum \left\{ p_j' (1 - p_j') + p_j'' (1 - p_j'') \right\} \left\{ \left(\sum d_j \right)^{-2} N_j + \right. \\ &\left. + 2 \left(\sum d_j \right)^{-1} \left(\sum e_j \right)^{-1} \sum_{i \in C_j} \int_{\min(t_{ij}^1, t_{ij}^2)}^{t_{ij}^2} \lambda^*(u) du + \left(\sum e_j \right)^{-2} \sum_{i \in C_j} \left(\int_{\min(t_{ij}^1, t_{ij}^2)}^{t_{ij}^2} \lambda^*(u) du \right)^2 \right\} \end{aligned} \quad (25)$$

respectivement.

On peut tirer trois conclusions de (24) et (25). Premièrement, si le taux de fausse positivité est supérieur (inférieur) au taux de fausse négativité, il y aura en moyenne surestimation (sous-estimation) du TNM. Deuxièmement, si les taux de fausse positivité et de fausse négativité sont égaux dans chaque état, les biais introduits dans le TNM par les erreurs de couplage s'annuleront en réalité. Troisièmement, les deux types d'erreurs de couplage introduisent un surcroît de variation dans les estimations du TNM.

4.2 Paramètres de régression de risques relatifs

Pour établir l'incidence des erreurs de couplage sur les estimations de paramètres de régression, considérons d'abord le modèle général de régression de risques relatifs (2). Si nous remplaçons les nombres et observé et prévu de décès d_j et e_j dans la fonction de vraisemblance en expression logarithmique (7) par les nombres et observé et prévu de décès en cas d'erreurs de couplage d_j^L et e_j^L , nous obtenons

$$\log L = \sum_{j=1}^J \left\{ d_j^L \log \left(\gamma \left\{ \beta' z_j \right\} \right) - \gamma \left\{ \beta' z_j \right\} e_j^L \right\}. \quad (26)$$

Soit $\hat{\beta}$ et $\tilde{\beta}$ les estimateurs de maximum de vraisemblance (EMV) de β en fonction de $\{d_j, e_j\}$ et $\{d_j^L, e_j^L\}$ respectivement. Si $\hat{\Lambda}_j = \Lambda_j(\hat{\beta})$, la fonction de caractérisation (9) peut s'écrire de la manière suivante :

$$\sum_{j=1}^J \frac{\partial \Lambda_j(\tilde{\beta})}{\partial \beta} \left\{ d_j + \Delta d_j - \exp \left\{ \hat{\Lambda}_j + \Delta \Lambda_j \right\} (e_j - \Delta e_j) \right\} = 0 \quad (27)$$

où $\Delta \Lambda_j = \Lambda_j(\tilde{\beta}) - \Lambda_j(\hat{\beta})$. Si nous supposons que $\Delta \beta = \tilde{\beta} - \hat{\beta}$ est petit, un développement du premier ordre de $\exp \left\{ \hat{\Lambda}_j + \Delta \Lambda_j \right\}$ autour de $\hat{\beta}$ donne

$$\exp \left\{ \hat{\Lambda}_j + \Delta \Lambda_j \right\} \approx \exp \left\{ \hat{\Lambda}_j \right\} + \exp \left\{ \hat{\Lambda}_j \right\} \frac{\partial \hat{\Lambda}_j}{\partial \beta} \Delta \beta \quad (28)$$

Si on substitue (28) dans (27), on obtient :

$$\begin{aligned} & \sum_{j=1}^J \frac{\partial \Lambda_j(\hat{\beta})}{\partial \beta} [d_j - \exp\{\hat{\Lambda}_j\} e_j] + \sum_{j=1}^J \frac{\partial \Lambda_j(\hat{\beta})}{\partial \beta} [\Delta d_j + \gamma \{\hat{\beta}' z_j\} \Delta e_j - \\ & - \gamma \{\hat{\beta}' z_j\} e_j \frac{\partial \Lambda'_j(\hat{\beta})}{\partial \beta} \Delta \beta + \gamma \{\hat{\beta}' z_j\} \Delta e_j \frac{\partial \Lambda'_j(\hat{\beta})}{\partial \beta} \Delta \beta] \approx 0. \end{aligned} \quad (29)$$

Si on utilise (9), la première sommation dans (29) est nulle. Ainsi, comme $\Delta e_j \Delta \beta$ est petit, $\Delta \beta$ peut être approximé par

$$\Delta \beta \approx \left\{ \sum \frac{\partial \hat{\Lambda}_j}{\partial \beta} \gamma \{\hat{\beta}' z_j\} e_j \frac{\partial \hat{\Lambda}'_j}{\partial \beta} \right\}^{-1} \sum \frac{\partial \hat{\Lambda}_j}{\partial \beta} \{ \Delta d_j + \gamma \{\hat{\beta}' z_j\} \Delta e_j \} \quad (30)$$

Il découle de (30) que

$$E \{ \Delta \beta \} \approx \left\{ \sum \frac{\partial \Lambda_j}{\partial \beta} \gamma \{\hat{\beta}' z_j\} e_j \frac{\partial \Lambda'_j}{\partial \beta} \right\}^{-1} \sum \frac{\partial \Lambda_j}{\partial \beta} \Delta p_j \alpha_j, \quad (31)$$

où

$$\alpha_j = N_j + \gamma \{\hat{\beta}' z_j\} \sum_{i \in C_j} \int_{\min(t_{ij}^1, t_{ij}^2)}^{t_{ij}^2} \lambda^*(u) du. \quad (32)$$

De plus,

$$\text{Var} \{ \Delta \beta \} \approx \left\{ \sum \frac{\partial \Lambda_j}{\partial \beta} \gamma \{\hat{\beta}' z_j\} e_j \frac{\partial \Lambda'_j}{\partial \beta} \right\}^{-1} \left[\sum \frac{\partial \Lambda_j}{\partial \beta} \Theta_j \frac{\partial \Lambda'_j}{\partial \beta} \right] \left\{ \sum \frac{\partial \Lambda_j}{\partial \beta} \gamma \{\hat{\beta}' z_j\} e_j \frac{\partial \Lambda'_j}{\partial \beta} \right\}^{-1}, \quad (33)$$

avec

$$\Theta_j = \{ p_j^I (1 - p_j^I) + p_j^{II} (1 - p_j^{II}) \} \sum_{i \in C_j} \left\{ 1 + \gamma \{\hat{\beta}' z_j\} \int_{\min(t_{ij}^1, t_{ij}^2)}^{t_{ij}^2} \lambda^*(u) du \right\}^2. \quad (34)$$

Dans le cas d'espèce du modèle multiplicatif (4), le biais $\Delta \beta$ par erreurs de couplage peut être approché par

$$\Delta \beta \approx (X' W X)^{-1} X' (\Delta D + \Delta W), \quad (35)$$

où $X' = (z'_1, \dots, z'_J)$, $\Delta D' = (\Delta d_1, \dots, \Delta d_J)$, $W = \text{diag}(\exp(z'_1 \hat{\beta}) e_1, \dots, \exp(z'_J \hat{\beta}) e_J)$, and $\Delta W' = (\exp(z'_1 \hat{\beta}) \Delta e_1, \dots, \exp(z'_J \hat{\beta}) \Delta e_J)$. À noter que la matrice de pondération W est la matrice de pondération de Fisher pour $\hat{\beta}$. De (35), il s'ensuit que

$$E \{ \Delta \beta \} \approx (X' W X)^{-1} X' \Pi, \quad (36)$$

où $\Pi' = (\pi_1, \dots, \pi_j)$ avec

$$\pi_j = \Delta p_j \left\{ N_j + \exp(\mathbf{z}'_j \hat{\beta}) \sum_{i \in C_j} \int_{\min(t_{ij}^1, t_{ij}^2)}^{t_{ij}^2} \lambda^*(u) du \right\}. \quad (37)$$

De plus,

$$\text{Var} \{ \Delta \beta \} \approx (\mathbf{X}' \mathbf{W} \mathbf{X})^{-1} \mathbf{X}' \Psi \mathbf{X} (\mathbf{X}' \mathbf{W} \mathbf{X})^{-1}, \quad (38)$$

où $\Psi = \text{diag}(\psi_1, \dots, \psi_j)$ avec éléments diagonaux

$$\psi_j = \{ p_j^I (1 - p_j^I) + p_j^{II} (1 - p_j^{II}) \} \sum_{i \in C_j} \left\{ 1 + \exp(\mathbf{z}'_j \hat{\beta}) \int_{\min(t_{ij}^1, t_{ij}^2)}^{t_{ij}^2} \lambda^*(u) du \right\}^2. \quad (39)$$

Avec une covariable unique $z_i = 1$, $\mathbf{X}' \mathbf{W} \mathbf{X} = e^{\hat{\beta}} \Sigma e_j$, $\mathbf{X}' \Delta \mathbf{D} = \Sigma d_j$ et $\mathbf{X}' \Delta \mathbf{W} = e^{\hat{\beta}} \Sigma \Delta e_j$. Dans ce cas,

$$\Delta \beta \approx \left(\Sigma \Delta d_j + e^{\hat{\beta}} \Sigma \Delta e_j \right) / \left(e^{\hat{\beta}} \Sigma e_j \right). \quad (40)$$

Comme

$$TNM = e^{\hat{\beta}} = \frac{\Sigma d_j}{\Sigma e_j}, \quad (41)$$

avec $\Delta \beta = \Delta TNM / TNM$ dans ce cas, nous obtenons

$$\Delta \beta \approx \frac{\Sigma \Delta d_j}{\Sigma d_j} + \frac{\Sigma \Delta e_j}{\Sigma e_j}. \quad (42)$$

Ainsi, (42) peut être considéré comme un cas d'espèce de (23).

Les résultats qui précèdent indiquent que, en cas d'égalité des taux de fausse positivité et de fausse négativité dans chaque état, les biais des estimations des paramètres de régression par erreurs de couplage seront presque éliminés. Il faut cependant préciser que les faux positifs et les faux négatifs introduiront un surcroît de variation dans les estimations des paramètres de régression de risques relatifs.

5. CONCLUSIONS

On utilise les méthodes informatisées de couplage d'enregistrements depuis un certain temps dans les études de mortalité par cohortes, mais on n'a pas encore soumis à un examen détaillé la question de l'incidence des

erreurs de couplage sur la fiabilité des inférences statistiques qui en sont tirées. Les résultats théoriques que nous avons présentés vont en ce sens.

Ils montrent que les faux positifs non seulement gonflent le nombre observé de décès, mais tendent aussi à dégonfler le nombre prévu de décès. À l'inverse, les faux négatifs gonflent le nombre prévu et dégonflent le nombre observé. On a également pu voir que les erreurs de couplage introduisent un biais dans les estimations de TNM, les faux liens et les faux non-liens étant généralement source d'un biais positif et d'un biais négatif. L'estimation des coefficients de régression de risques relatifs est également susceptible d'être entachée d'un biais, mais le sens de ce biais dépendra de la nature des coefficients. Il n'y a pas que les biais que nous venons d'énumérer, puisque les erreurs de couplage ajoutent de l'incertitude aux estimations tant des TNM que des coefficients de régression.

Il est possible de presque éliminer les biais des estimations de TNM si les taux de fausse positivité et de fausse négativité sont égaux dans chaque état défini par un classement recoupé des covariables d'intérêt, à condition que toutes ces covariables prennent des valeurs constantes pour chacun des états en question. Il s'agit de conditions exigeantes auxquelles il est improbable que l'on satisfasse dans la pratique. Il reste que, d'après nos résultats, il serait sans doute souhaitable d'établir des valeurs seuils de couplage qui équilibrent séparément les taux de fausse positivité et de fausse négativité dans des états homogènes.

Les résultats analytiques que nous avons livrés répandent un éclairage considérable sur les effets que peuvent avoir les erreurs de couplage dans les études de mortalité par cohortes, mais il importe d'examiner de tels effets dans des conditions qui soient le plus proches possible de celles qui pourraient se présenter dans la pratique. Voilà pourquoi nous effectuons actuellement, à l'aide de données réelles du Registre national des doses du Canada, une étude par simulation informatique où l'introduction de faux liens et de faux non-liens aux probabilités connues sert à mieux cerner l'incidence des erreurs de couplage sur les estimations de risque de cancer.

Nos résultats peuvent aider à clarifier l'incidence des erreurs de couplage sur l'inférence statistique, mais il reste à développer les méthodes par lesquelles ces erreurs sont prises en compte dans les analyses statistiques. Il peut s'agir de modèles d'erreurs de réponse employés en échantillonnage d'enquête combinées avec les méthodes statistiques classiques d'analyse des données de mortalité par cohortes. Des recherches sont aussi en cours dans ce domaine.

REMERCIEMENTS

Ces travaux de recherche ont été en partie financés par une subvention accordée par le Conseil national de recherches en sciences naturelles et en génie du Canada à D. Krewski, titulaire actuel de la chaire CRSNG-CRSH- McLaughlin en évaluation des risques pour la santé des populations à l'Université d'Ottawa. De premières versions du présent document ont été présentées au congrès annuel mixte de l'American Statistical Association à San Francisco du 8 au 12 août 1993 et au congrès annuel de la Société statistique du Canada à Montréal du 10 au 16 juillet 1995. On a eu droit à la version définitive à la séance où on a honoré J.N.K. Rao (18 octobre 2001 à Ottawa) au symposium 2001 de Statistique Canada. Le premier de ses auteurs (D. Krewski) est particulièrement reconnaissant d'avoir été invité à prendre la parole à cette occasion. J.N.K. Rao a été son directeur de thèse de doctorat à une certaine époque.

BIBLIOGRAPHIE

Ashmore, J.P. et Grogan, D. (1985), *The National Dose Registry of Canada@Radiation Protection Dosimetry*, 11, 95-100.

- Ashmore, J.P. et Davies, B.D. (1989), *The National Dose Registry: A Centralized Record Keeping System for Radiation Workers in Canada* in *Applications of Computer Technology to Radiation Protection*. IAEA-SR-136/58, J. Stephan Institute, Ljublyua, 505-520.
- Ashmore, J.P. Krewski, D. et Zielinski, J.M. (1997), *Protocol for a Cohort Mortality Study of Occupational Radiation Exposure Based on the National Dose Registry of Canada* *European Journal of Cancer*, 33, S10-S21.
- Ashmore, J.P., Krewski, D., Zielinski, J.M., Jiang, H., Semenciw, R., et Létourneau E. (1998), *First Analysis of Occupational Radiation Mortality Based on the National Dose Registry of Canada* *American Journal of Epidemiology*, 148, 564-574.
- Bartlett, S., Krewski, D., Wang, Y. et Zielinski, J.M. (1993), *Evaluation des taux d'erreur dans de grandes études par couplage d'enregistrements informatisé* *Technique d'enquête*, 19, 3-13.
- Belin, T.R. et Rubin, D.B. (1991), *Recent Developments in Calibrating Error Rates for Computer Matching* *Proceedings of the 1991 Annual Research Conference*, U.S. Bureau of the Census, 657-668.
- Breslow, N.E., Lubin, J.H. et Langholz, B. (1983), *Multiplicative Models and Cohort Analysis* *Journal of the American Statistical Association*, 78, 1-12.
- Breslow, N.E. et Day, N.E. (1987), *Statistical Methods in Cancer Research, Vol. 2: The Design and Analysis of Cohort Studies*, IARC Scientific Publication No. 82, International Agency for Research on Cancer, Lyon, France.
- Carpenter, M. et Fair, M.E. (Eds.) (1990), *Canadian Epidemiology Research Conference - 1989: Proceedings of Record Linkage Sessions & Workshop*, Ottawa Select Printing, Ottawa.
- Cox, D.R. (1972), *Regression Models and Life Tables (with discussion)* *Journal of Royal Statistical Society*, B34, 187-220.
- Efron, B. (1982), *The Jackknife, the Bootstrap and Other Resampling Plans* *Society for Industrial and Applied Mathematics*, Philadelphia, Pennsylvania, 13-19.
- Fair, M.E. (1989), *Studies and References Relating to Uses of the Canadian Mortality Data Base* *Rapport de la Occupational and Environmental Health Research Unit, Division de la santé, Statistique Canada*, Ottawa.
- Fellegi, I. et Sunter, A. (1969), *A Theory for Record Linkage* *Journal of the American Statistical Association*, 64, 1183-1210.
- Hill, T. (1988), *Generalized Iterative Record Linkage System: GIRLS Strategy (Release 2.7)* *Rapport du Research and General System, Division du développement de systèmes informatiques, Statistique Canada*, Ottawa.
- Howe, G.R. et Lindsay, J. (1981), *A Generalized Iterative Record Linkage Computer System for Use in Medical Follow-up Studies* *Computers and Biomedical Research*, 14, 327-340.
- Howe, G.R. et Spasoff, R.A. (Eds.) (1986), *Proceeding of the Workshop on Computerized Linkage in Health Research*, University of Toronto Press, Toronto.

- Jordan-Simpson, D.A., Fair, M.E., et Poliquin, C. (1990), *A Canadian Farm Operator Study: Methodology* @ *Health Reports*, 2, 141-155.
- Kalbfleish, J.D. et Prentice, R.L. (1980), *The Statistical Analysis of Failure Time Data*, New York, Wiley.
- Labossière, G. (1986), *Confidentiality and Access to Data: The Practice at Statistics Canada* @ *Proceedings of the Workshop on Computerized Record Linkage in Health Research*, University of Toronto Press, Toronto.
- Neter, J., Maynes, E.S., et Ramanathan, R. (1965), *The Effect of Mismatching on the Measurement of Response Errors* @ *Journal of the American Statistical Association*, 60, 1005-1027.
- Newcombe, H.B. (1988), *Handbook of Record Linkage: Methods for Health and Statistical Studies, Administration, and Business*, Oxford Medical Publications, Oxford.
- Scheuren, F. et Winkler, W.E. (1993), "Analyse de régression de fichiers de données couplés par ordinateur" @ *Techniques d'enquête*, 19, 45-65.
- Smith, M.E. et Silins, J. (1981), *Generalized Iterative Record Linkage System* @ *Social Statistics Section, Proceedings of the American Statistical Association*, 128-137.
- Sont, W.N., Zielinski, J.M., Ashmore, J.P., Jiang, H., Krewski, D., Fair, M.E., Band, P. and Létourneau, E. (2001), *First Analysis of Cancer Incidence and Occupational Radiation Exposure Based on the National Dose Registry of Canada* @ *American Journal of Epidemiology*, 153, 309-318.
- Winkler, W.E. et Scheuren, F. (1991), *How Computer Matching Error effect Regression Analysis: Exploratory and Confirmatory Analysis* @ Technical Report, Statistical Research Division, U.S. Bureau of the Census, Washington, D.C.