

APPLICATION OF THE FAY-HERRIOT MODEL WITH CORRELATED DESIGN ERRORS: THE LABOUR FORCE SURVEY EXPERIENCE

François Verret¹ and Braedan Walker²

ABSTRACT

The Fay-Herriot model is an area-level model commonly used to perform small area estimation since it efficiently integrates the model and design random processes. It assumes independence between the area estimates. However, survey weights are usually calibrated at a more aggregated level than the small area, which can introduce non-negligible correlation between area estimates when the characteristic of interest is closely related to a calibration variable. This is the case of the Labour Force Survey (LFS), where monthly estimates of total employment are required for a detailed mapping of the ten provinces and the survey weights are calibrated to provincial population totals. This paper describes the successful experience of the LFS in using the Fay-Herriot model extended for correlated design errors. It covers smoothing of the design variances, derivation of a working design variance matrix, point estimation and mean squared prediction error estimation, raking of the point estimates and generalization of key diagnostics.

KEY WORDS: Small area estimation, working variance matrix, raking, model diagnostics.

RÉSUMÉ

Le modèle de Fay-Herriot est un modèle au niveau des domaines communément utilisé pour l'estimation pour petits domaines parce qu'il intègre de façon efficace la randomisation due au modèle et celle due au plan de sondage. Il suppose l'indépendance entre les estimations des petits domaines. Cependant, les poids de sondage sont habituellement calés à un niveau plus agrégé que le domaine, ce qui peut introduire des corrélations non-négligeables entre les estimations des domaines lorsque la caractéristique d'intérêt est fortement liée à une variable de calage. C'est le cas de l'Enquête sur la population active (EPA), où des estimations mensuelles du total de l'emploi sont requises pour un découpage détaillé des dix provinces et où le calage est fait à des totaux de population provinciaux. Cet article illustre l'expérience fructueuse de l'EPA dans l'utilisation du modèle de Fay-Herriot étendu pour des erreurs sous le plan de sondage corrélées. Il couvre le lissage des variances sous le plan de sondage, la dérivation d'une matrice de travail de variance sous le plan, l'estimation ponctuelle et l'estimation de l'erreur quadratique moyenne de prédiction, la réconciliation des estimations ponctuelles et la généralisation de diagnostics clés.

MOTS CLÉS : Estimation pour de petits domaines, matrice de travail de variance, réconciliation, diagnostics du modèle.

1. INTRODUCTION

Statistical agencies are often tasked with the challenge of delivering estimates at a level that is a lot less aggregated than the levels used to define the original objectives of their surveys. This is done in large part to better represent small or minority population groups. To meet the demand, small area estimation (SAE) methodologies are often considered. In common practice, the Fay-Herriot (FH) model (Fay & Herriot, 1979), an area-level model, is used to perform SAE. Area-level models such as the FH model have an advantage over unit-level models of only requiring the data to be available at the aggregated area-level rather than at the unit-level. This avoids the requirement for unit-level record linkages between survey records (which provide the dependent variable of the model) and auxiliary data sources records (which provide the independent variables of the model). It is also notably more difficult to account for the potentially informative sampling design with unit-level models than with area-level models like the FH model, which integrates in a straightforward manner the model and design random processes.

The Canadian Labour Force Survey (LFS) has been using the FH model in recent years to produce small area estimates of total employment, unemployment rates and wage percentiles for a variety of domains. For the estimation of total

¹ François Verret, Statistics Canada, 100 Tunney's Pasture Driveway, Ottawa, Ontario, Canada, K1A 0T6

² Braedan Walker, Statistics Canada, 100 Tunney's Pasture Driveway, Ottawa, Ontario, Canada, K1A 0T6

employment, estimates are produced and disseminated monthly for a mapping of the ten provinces defined by Census metropolitan areas (CMA), Census agglomerations (CA) and Self-contained labour areas (SLA). This represents over 600 areas in total. CMAs contain larger cities, have a population of at least 100,000 people and have at least 50,000 people in their core area. CAs contain smaller towns and have a population of at least 10,000 people. SLAs cover the rest of the provinces. They are formed by combining Census subdivisions (CSD) based on place of work patterns obtained from the latest Census long form data (OECD, 2020). In general, CMAs and CAs correspond to urban areas, and SLAs to rural areas.

Direct estimation of total employment for those geographies is performed using standard design-based methodologies. A direct estimate for a given domain is obtained using survey data and weights of that domain only. The monthly LFS sample is large (more than 53,000 households are sampled; Statistics Canada, 2020), and the survey produces good quality estimates for most CMAs. However, the LFS is not designed to produce precise estimates for all areas each month. Therefore, many small areas will have little to no sample size (about a third of the areas have no sample size each month). For these areas, the estimate of the total can be zero and will be of poor quality. The corresponding variance estimate of the total estimate will be of poor quality for the same reason. SAE is thus used to produce improved estimates to disseminate. On top of this, good predictive auxiliary data is available and this makes gains in precision of SAE substantial.

This paper describes the experience of the LFS using the FH model to produce monthly estimates of total employment for CMAs, CAs and SLAs. Section 2 outlines the initial use of the basic FH model and why the model needed to be extended to dependent sampling errors. Section 3 details the model fitting for the extended model. It covers the determination of the between-area design covariances, and the extensions of the various diagnostics used for the basic model. Section 4 outlines the advantages of using the extended model over the basic model. A conclusion follows in Section 5.

2. SMALL AREA ESTIMATION OF TOTAL EMPLOYMENT USING THE BASIC FH MODEL

2.1 The Basic FH Model and its Application to Estimate Total Employment

SAE of total employment for CMAs, CAs and SLAs was originally performed using the basic FH model. This linear mixed model is given by

$$\hat{\theta}_m = z_m^T \beta + b_m v_m + e_m,$$

where $m = 1, \dots, M$ represents the area, $\hat{\theta}_m$ is the direct survey estimate, z_m is a vector of auxiliary information, b_m is a predetermined constant, v_m is the model error and e_m is the error produced by survey sampling. For this application, the auxiliary information is a function of population and employment insurance (EI) beneficiaries counts and $b_m = N_m$, the number of people aged 15 and over in the area. The model can be decomposed into the sampling model and the linking model, which are respectively

$$\hat{\theta}_m = \theta_m + e_m$$

and

$$\theta_m = z_m^T \beta + b_m v_m,$$

where θ_m is the target parameter of interest one would obtain if a census of the area was performed. For the application this parameter is the total population employed of area m . For the sampling model, it is assumed that $e_m \sim_{\text{ind}}(0, \text{var}_p(e_m | z_m) = \psi_m)$, and for the linking model that $v_m \sim_{\text{iid}} N(0, \text{var}_m(v_m | z_m) = \sigma^2)$. When combining the two models, we are interested in performing model and design based inference, we assume normality of the sampling errors e_m and we get

$$\hat{\theta}_m = z_m^T \beta + b_m v_m + e_m,$$

where $e_m \sim_{\text{ind}} N(0, \text{var}_{mp}(e_m | z_m) = \tilde{\psi}_m)$ are assumed independent of $v_m \sim_{\text{iid}} N(0, \text{var}_{mp}(v_m | z_m) = \sigma^2)$.

For estimation purposes, the sampling variance of the direct estimator, $\tilde{\psi}_m$, is assumed to be known, although in practice it is not, as is normally the case for SAE applications. From an application of the Rao-Wu bootstrap (Rao & Wu, 1988), we

have for the LFS variance estimates $\hat{\psi}_m$. For areas with large sample size, $\hat{\psi}_m$ is used as the value of $\tilde{\psi}_m$. For other areas, a prediction from a generalized variance function model (Beaumont & Bocci, 2016) $\tilde{\psi}_m$ is used. Given the “known” sampling variances, the model variance σ^2 and the vector β are estimated using restricted maximum likelihood (REML).

Fitting this model, one can obtain from the fixed part of the model the synthetic estimator $z_m^T \hat{\beta}$. However, the empirical best linear unbiased predictor (EBLUP) is given by the composite estimator

$$\hat{\theta}_m^{\text{SAE}} = \hat{\gamma}_m \hat{\theta}_m + (1 - \hat{\gamma}_m) z_m^T \hat{\beta},$$

where $\hat{\gamma}_m = \frac{b_m^2 \hat{\sigma}^2}{b_m^2 \hat{\sigma}^2 + \tilde{\psi}_m}$. The composite estimator is thus a convex combination of the direct estimator $\hat{\theta}_m$ and of the synthetic estimator $z_m^T \hat{\beta}$. Its proximity to one estimator or the other depends on the relative model and design variance for the given area through $\hat{\gamma}_m$. For example, if the model predicts very precisely, then the model variance will be close to 0, $\hat{\gamma}_m$ will be close to 0, and the small area estimate will be close to the synthetic estimate. If in an area a census was performed, then there will be no sampling variance in that area, $\hat{\gamma}_m$ will be equal to 1 and the small area estimate will be the direct estimate, which will be equal to the target parameter/census value θ_m .

The quality of the estimates obtained from the FH model and the validity of the estimated mean squared prediction error of those estimates directly depend on the validity of the model hypotheses. For that reason, a backward selection of the regressors is done, and numerous diagnostics are performed in practice to validate the model. Furthermore, an outlier detection algorithm is applied using the standardized residuals on both the FH model and the variance smoothing model to ensure erroneously influential areas are excluded from the model fit. Finally, for each domain, it is determined if the final estimate to publish should be the composite estimate or the original direct estimate using a diagnostic developed by Lesage, Beaumont & Bocci (2021). This diagnostic estimates the probability that the composite estimate is the most precise. The direct estimator is chosen for a given domain if the probability for that domain is less than 25%.

2.2 Initial Signs of Failure of the Independence of Sampling Errors Assumption

Since totals are estimated for a geographical mapping of the provinces, the sum of the small area estimates by province can be compared to the original direct provincial estimate of the LFS, which is of good quality by design. Table 1 compares those two provincial estimates for the month of November 2021. The relative difference can be large, especially for the smallest provinces such as Prince Edward Island (PEI) with a relative difference of 9.7%. Those relative differences are consistent for all months.

Table 1 – Provincial discrepancies between direct estimate of total employment and the aggregate of small area estimates using the basic Fay-Herriot model for the month of November 2021

Province	Direct estimate	Small area estimate	Difference	Relative difference (%)
NL	229,601	224,525	5,076	2.2
PEI	81,098	73,215	7,883	9.7
NS	477,326	483,855	-6,529	-1.4
NB	360,082	360,306	-224	-0.1
QUE	4,380,216	4,361,616	18,600	0.4
ONT	7,660,158	7,778,329	-118,171	-1.5
MAN	669,344	652,554	16,790	2.5
SASK	556,206	509,706	46,500	8.4
ALTA	2,298,973	2,277,357	21,616	0.9
BC	2,708,386	2,701,059	7,327	0.3

A raking process (Dagum & Cholette, 2006) was applied to the estimates to eliminate the differences. This process minimizes the following distance function

$$\sum_{m=1}^M \frac{(\hat{\theta}_m^{\text{SAE}} - \hat{\theta}_m^{\text{SAE, raked}})^2}{|c_m \hat{\theta}_m^{\text{SAE}}|},$$

under the constraint that the sum of the resulting estimates $\hat{\theta}_m^{\text{SAE, raked}}$ is equal to the sum of the direct estimates $\hat{\theta}_m$ by province. Constants c_m are chosen “alterability coefficients” which control relative changes to the original estimates. For this application they were chosen to be the relative root mean square prediction error (RRMSPE) of the small area estimates, so that the most precise estimates are affected the least by raking.

A parametric bootstrap using the assumptions of the basic FH model was then used to estimate the mean squared prediction error (MSPE) of the resulting estimators. Since raking was performed and under the independence assumption, one would expect the estimated MSPE of the sum of the raked small area estimates at the provincial level to be close to the estimated variance of the provincial direct estimates. Table 2 compares the direct estimates CVs to the RRMSPE of the aggregates of raked estimates at the provincial level. The two quantities can be quite different, especially in the smallest provinces. This indicates that the independence assumption does not hold and that the direct area estimates must generally be negatively correlated. A more direct way to verify that the independence assumption does not hold is to compare the sum of the variances of the direct area estimates of a province to the variance of the direct estimate of that province. This can be done before any small area model is fitted.

Table 2 – Provincial discrepancies between the CV of the direct estimate of total employment and the RRMSPE of the aggregate of small area estimates using the basic Fay-Herriot model for the month of November 2021

Province	Direct CV	RRMSPE of aggregates of raked estimates	Ratio
NL	1.5%	5.4%	3.6
PEI	1.5%	8.0%	5.3
NS	1.2%	3.9%	3.3
NB	1.3%	3.8%	2.9
QUE	0.6%	1.4%	2.3
ONT	0.5%	1.0%	2.0
MAN	0.9%	3.7%	4.1
SASK	1.1%	4.8%	4.4
ALTA	0.8%	2.1%	2.6
BC	0.8%	1.8%	2.3

The sampling covariance present in the direct estimates of the LFS is due to estimation procedures. In terms of sampling, LFS strata are small and numerous, which should guarantee that Horvitz-Thompson estimators (i.e. design weighted estimators) of the small area totals are independent. However, final direct LFS estimators at the area level are dependent because the calibration adjustment of the LFS weights is done at the provincial level and some calibration variables (total population by age and gender groups) are strongly correlated with the variable of interest (total employment). When calibrating at the province level, the variance of the total employment estimate at the province level is reduced because a large proportion of the total working age population is employed. However, variance of total employment at the small area level should not be reduced proportionally because of the weaker relationship between total employment at that level and total working age population at the province level. This means the sum of the off-diagonal terms in the provincial variance matrix of area-level direct estimates should decrease to a negative value with provincial weight calibration.

Table 2 indicates that the square root of the sum of the variances of direct area-level estimates of a given province is at least twice the square root of the variance of the provincial direct estimate. That is, we have that for any province p

$$\sqrt{\sum_{m \in \Omega_p} \hat{\psi}_m} \geq 2.0 \sqrt{\sum_{m \in \Omega_p} \hat{\psi}_m + \sum_{m \neq m' \in \Omega_p} \hat{\psi}_{mm'}},$$

where Ω_p is the set of areas in province p . This means the sum of the off-diagonal terms of the provincial variance matrix $\sum_{m \neq m' \in \Omega_p} \hat{\psi}_{mm'}$ is at most $-3/4$ times the sum of the variance of the direct area-level estimates $\sum_{m \in \Omega_p} \hat{\psi}_m$. This is a significant departure from the independence assumption of the basic FH model. For this reason, the FH model extended for dependent sampling errors was considered instead.

3. SMALL AREA ESTIMATION OF TOTAL EMPLOYMENT USING THE FH MODEL EXTENDED FOR DEPENDENT SAMPLING ERRORS

The Fay-Herriot model extended for dependent sampling errors can be written in matrix form as

$$\hat{\theta} = Z\beta + Bv + e,$$

where $v \sim N(0, \sigma^2 I)$ is independent of $e \sim N(0, \Psi)$ and matrix B is a diagonal matrix with values b_m in its diagonal. The difference with the basic model is that Ψ may not be a diagonal matrix. This matrix is assumed to be known for estimation purposes. With the basic model, smoothed variances were obtained to determine the diagonal terms of the Ψ matrix and covariance terms were assumed to be 0. The challenge in practice is thus to obtain the covariance terms of Ψ in the presence of small sample sizes. The properties of the desired $\tilde{\Psi}$ matrix are:

1. The diagonal should correspond to the smoothed direct variance estimates $\tilde{\psi}_m$ obtained for the basic model.
2. The covariance structure should be one that respects the covariance structure of some aggregate totals. In our case, we will respect the variance-covariance structure of the provincial estimates for the LFS since both LFS weight calibration and raking of the FH estimates are done at that level.
3. The matrix should be a proper covariance matrix and must adhere to all properties of a covariance matrix (symmetric, positive semi-definite).

The key idea to derive such a matrix is that variance matrices appear naturally as a result of statistical processes. Thus, it is assumed that direct estimates come from a somewhat simplified statistical process to satisfy the third property. The statistical process is specifically chosen so that the first and second properties hold as well. Inspired by LFS sampling and weight calibration at the provincial level, it was assumed that the direct estimates $\hat{\theta}$ result from independent variables X (to represent independent sampling) that have been raked to provincial estimates Y (to represent weight calibration). The assumed raking thus corresponds to minimizing

$$\sum_{m=1}^M \frac{(X_m - \hat{\theta}_m)^2}{|c_m X_m|}$$

under the constraint $G\hat{\theta} = Y$, where G is a matrix of zeros and ones summing the direct estimates $\hat{\theta}$ to the vector of provincial estimates Y . We use the variance matrix of the ‘‘raked’’ estimates $\hat{\theta}$ to define $\tilde{\Psi}$. Note that the goal here is to obtain the variance matrix of direct estimates $\hat{\theta}$, which comes before applying a FH model or raking small area estimates.

Under this hypothetical model and choosing alterability coefficients $c_m = \sqrt{V(\hat{\theta}_m)}/X_m$ for the hypothetical raking, Verret & Walker (2024) showed that $\hat{\theta} = AX + CY$, for fixed matrices A and C and that the working variance matrix sought will be

$$V(\hat{\theta}) = \tilde{\Psi} = \text{Adiag}(\sigma_{X,1}^2, \dots, \sigma_{X,M}^2)A^T + C\Sigma_Y C^T, \quad (1)$$

where $\sigma_{X,1}^2, \dots, \sigma_{X,M}^2$ are the (unknown) variances of the hypothetical independent variables X ; matrices A and C and the diagonal elements of $V(\hat{\theta})$ are known; and Σ_Y is the matrix of known covariances at the aggregated level. The diagonal of equation (1) represents a system of M equations with M unknowns. Using linear regression theory, the values of $\sigma_{X,1}^2, \dots, \sigma_{X,M}^2$ that generate covariances at the aggregated level Σ_Y and that preserves the smoothed variances as much as possible is given by

$$\sigma_X^2 = \left(A^{2T} A^2 \right)^{-} A^{2T} \left[D \left(V(\hat{\theta}) \right) - D(C \Sigma_Y C^T) \right],$$

where A^2 is the matrix of the squared elements of A and operator “D” turns the diagonal of the matrix in its argument into a vector. Using these values in the right-hand side of equation (1), we obtain the complete matrix we are looking for $\tilde{\Psi}$. In essence, this extends the smoothing made for the diagonal of $\tilde{\Psi}$ to the off-diagonal elements while preserving the known covariance structure Σ_Y .

The EBLUP or composite estimator of the extended model is given by

$$\hat{\theta}^{\text{SAE}} = \hat{H}^T Z \hat{\beta} + (I - \hat{H}^T) \hat{\theta},$$

where $\hat{H}^T = I - \hat{\sigma}^2 B^2 \hat{\Sigma}_{zz}^{-1}$ and $\hat{\Sigma}_{zz} = \hat{\sigma}^2 B^2 + \tilde{\Psi}$. Available standard SAE software only fit the basic FH model, including Statistics Canada's generalized system G-Est (Estevao et al., 2023) and the sae package in R. Hence, estimation with the extended model needed to be programmed. Further theoretical developments and their programming were necessary to generalize the steps of the methodology from the original project to the extended FH model. Those are:

1. The backward variable selection
2. The condition index collinearity diagnostic
3. The outlier detection
4. The many graphical diagnostics available for the basic FH model in G-Est
5. The diagnostic of Lesage et al. (2021).

For the first four items, the data had to be made independent by standardizing them beforehand using the following formula:

$$\left(\hat{\sigma}^2 B^2 + \tilde{\Psi} \right)^{-1/2} \hat{\theta} = \left(\hat{\sigma}^2 B^2 + \tilde{\Psi} \right)^{-1/2} (Z\beta + Bv + e).$$

To extend the Lesage et al. (2021) diagnostic, the same steps as in the original paper were followed assuming that the variance matrices and β are known, but for the extended FH model. They are outlined next.

Step 1: Derive the expression of the design MSPE matrices of the direct estimator and EBLUP vectors.

We have

$$\begin{aligned} \text{MSPE}_p(\hat{\theta}) &= V(\hat{\theta} - \theta | v) \\ &= V(Z\beta + Bv + e - Z\beta - Bv | v) = \Psi \end{aligned}$$

and

$$\begin{aligned} \text{MSPE}_p(\hat{\theta}^{\text{SAE}}) &= V(\hat{\theta}^{\text{SAE}} - \theta | v) \\ &= V(H^T Z\beta + (I - H^T)(Z\beta + Bv + e) - Z\beta - Bv | v) \\ &= V((I - H^T)e - H^T Bv | v) \\ &= (I - H^T)\Psi(I - H) + H^T Bv v^T B H, \end{aligned}$$

where $H^T = I - \sigma^2 B^2 \Sigma_{zz}^{-1}$ and $\Sigma_{zz} = \sigma^2 B^2 + \Psi$.

Step 2: Find the expression giving the condition when the MSPE of the EBLUP of a given area is smaller than that of its direct estimator.

$\text{MSPE}_p(\hat{\theta}_i^{\text{SAE}})$ will be smaller than $\text{MSPE}_p(\hat{\theta}_i)$ when

$$1_i^T \Psi 1_i \geq 1_i^T (I - H^T) \Psi (I - H) 1_i + 1_i^T H^T Bv v^T B H 1_i \Leftrightarrow 1_i^T [\Psi - (I - H^T) \Psi (I - H)] 1_i \geq 1_i^T H^T Bv v^T B H 1_i,$$

where 1_i is a vector with a 1 in position i and a 0 elsewhere.

Step 3: Find the conditional distribution of $u = Bv$ given $\hat{\theta}$.

We have $\hat{\theta} = Z\beta + Bv + e = Z\beta + u + e$, with $\hat{\theta}|u \sim N(Z\beta + u, \Psi)$ and $u \sim N(0, \sigma^2 B^2)$. Calculating the joint density and then conditioning on $\hat{\theta}$ we get $u|\hat{\theta} \sim N(\sigma^2 B^2 H \Psi^{-1}(\hat{\theta} - Z\beta), \sigma^2 B^2 H)$.

Step 4: Find the expression of the probability of the event of step 2 under the conditional distribution of step 3.

$$\begin{aligned} P\{1_i^T [\Psi - (I - H^T)\Psi(I - H)]1_i \geq 1_i^T H^T Bv v^T B H 1_i | \hat{\theta}\} &= P\{1_i^T [\Psi - (I - H^T)\Psi(I - H)]1_i \geq 1_i^T H^T u u^T H 1_i | \hat{\theta}\} \\ &= P\{1_i^T [\Psi - (I - H^T)\Psi(I - H)]1_i \geq w_i^2 | \hat{\theta}\}, \end{aligned}$$

where $w_i = 1_i^T H^T u = u^T H 1_i$. This probability is equal to $P(-w_{L,i} \leq w_i \leq w_{L,i} | \hat{\theta})$, where $w_{L,i} = \sqrt{1_i^T [\Psi - (I - H^T)\Psi(I - H)]1_i}$ and is a constant. Since $w_i = 1_i^T H^T u$, we have

$$w_i | \hat{\theta} \sim N\left(\mu_{w_i | \hat{\theta}} = 1_i^T H^T (\sigma^2 B^2 H) \Psi^{-1}(\hat{\theta} - Z\beta), \sigma_{w_i | \hat{\theta}}^2 = 1_i^T H^T (\sigma^2 B^2 H) H 1_i\right).$$

The diagnostic statistic is thus given by

$$D_{1i} = P(-w_{L,i} \leq w_i \leq w_{L,i} | \hat{\theta}) = \Phi\left(\frac{w_{L,i} - \mu_{w_i | \hat{\theta}}}{\sigma_{w_i | \hat{\theta}}}\right) - \Phi\left(\frac{-w_{L,i} - \mu_{w_i | \hat{\theta}}}{\sigma_{w_i | \hat{\theta}}}\right).$$

Step 5: Replace the unknown quantities in step 4 by their estimates.

The final diagnostic consists of using $\hat{D}_{1i} = \Phi\left(\frac{\hat{w}_{L,i} - \hat{\mu}_{w_i | \hat{\theta}}}{\hat{\sigma}_{w_i | \hat{\theta}}}\right) - \Phi\left(\frac{-\hat{w}_{L,i} - \hat{\mu}_{w_i | \hat{\theta}}}{\hat{\sigma}_{w_i | \hat{\theta}}}\right)$, where quantities β , σ^2 , Ψ and Σ_{zz}^{-1} are replaced by their estimated values in D_{1i} . The direct estimator is preferred over the EBLUP if the estimated probability is smaller than 25%.

4. ADVANTAGES OF USING THE EXTENDED FH MODEL OVER THE BASIC FH MODEL

Two major differences have been noticed that suggest using the correct covariance structure is very advantageous. Firstly, variable selection tests can only be trusted with the right variance structure, and using the wrong structure might lead to a suboptimal model. Variables available for backward selection are the working age (15 to 64 years of age) population count, auxiliary variables leading to a quadratic spline of that variable (i.e. the square and other quadratic spline terms of that variable) and the EI beneficiary count. The latter variable is always selected whether using the basic or the extended FH model assumptions. However, with the extended model the linear population count term is always chosen, whereas with the basic model it is not, and a single quadratic spline term is usually chosen instead. The variables selected for the extended model are more sensible since the linear term of a variable should normally be chosen before a spline term. This pattern is what is observed consistently for every month.

The second major difference between the extended and basic FH model is that the quality indicators are negatively impacted by the independence assumption when the estimates are raked. Table 2 showed that under the basic model the provincial RRMSPE of small area estimates was largely overestimated. By construction of the $\tilde{\Psi}$ matrix, the corresponding RRMSPE under the extended model is equal to the CV of the provincial estimate. Table 3 shows the average of the small area RRMSPE by province under the two models. The RRMSPEs of the raked extended FH estimates are not artificially inflated contrarily to those of the basic model, especially for the small provinces. The table also shows the very large average CVs of the direct estimates SAE aims to improve on.

On top of those two major advantages, we do not observe consistency of the β and σ^2 estimates under the basic model when considering only the variables selected under the assumption of the extended model (i.e. these estimates are very different from those of the extended model). Given the strong correlations, it appears that more domains are needed to observe convergence.

Table 3 – Small area level comparison of the CV of the direct estimate of total employment and of the RRMSPE of the small area estimates using the basic and extended Fay-Herriot models for the month of November 2021

Province	Average CV $\left(\sqrt{\tilde{\psi}_m/z_m^T\hat{\beta}}\right)$ of direct estimates	Average RRMSPE of raked basic estimates	Average RRMSPE of FH raked estimates	Average RRMSPE of extended FH estimates
NL	109.5%		14.5%	9.4%
PEI	85.0%		16.0%	7.0%
NS	30.8%		8.8%	6.3%
NB	78.3%		10.2%	7.3%
QUE	97.8%		8.9%	7.1%
ONT	60.9%		7.0%	6.1%
MAN	167.4%		12.7%	7.0%
SASK	483.4%		12.7%	6.8%
ALTA	73.3%		7.8%	6.5%
BC	160.5%		9.2%	7.3%

Table 4 shows that with the extended FH model, raking of SAE estimates is less necessary because the provincial differences are relatively smaller than with the basic model. Moreover, for a given region, the stability over time of the estimates is most often greater with the extended model. Stability over two years of monthly data was studied by measuring for each area the variance of the series of 24 monthly estimates, where smaller variances indicate greater stability. Whether the variable selection is set to that of the extended model or done independently for the basic model, the series are typically more stable (less volatile) with the extended model. With a common selection, for unraked data, 577 of the 623 areas are more stable with the extended model. For the raked data, it is 559 of the 623 regions. For the most populous small areas (the largest CMAs, such as Toronto, Montreal and Vancouver), the results are not as good for the extended model: 28 out of 60 and 18 out of 60 regions respectively do better with the extended model. This is an indication that the working variance matrix needs to be improved for those areas.

Table 4 – Provincial discrepancies between direct estimate of total employment and aggregate of small area estimates using the extended Fay-Herriot model for the month of November 2021

Province	Direct estimate	Small area estimate	Difference	Relative difference (%)	Relative difference under the basic model from Table 1 (%)
NL	229,601	231,945	-2,344	-1.0	2.2
PEI	81,098	78,809	2,289	2.8	9.7
NS	477,326	482,921	-5,595	-1.2	-1.4
NB	360,082	363,074	-2,992	-0.8	-0.1
QUE	4,380,216	4,396,197	-15,981	-0.4	0.4
ONT	7,660,158	7,722,312	-62,154	-0.8	-1.5
MAN	669,344	667,764	1,580	0.2	2.5
SASK	556,206	546,322	9,884	1.8	8.4
ALTA	2,298,973	2,295,764	3,209	0.1	0.9
BC	2,708,386	2,722,514	-14,128	-0.5	0.3

Finally, we compared the May 2016 direct and FH estimates with those of the 2016 Census using the average relative difference (ARD):

$$ARD(\hat{\theta}) = \frac{|\hat{\theta} - \hat{\theta}_{\text{Census}}|}{\hat{\theta}_{\text{Census}}}$$

Comparisons with the 2021 Census would have been timelier, however it was decided not to perform SAE for the month of May 2021 since the relationship between LFS direct estimates and EI data was severely compromised by the COVID-19 pandemic. Table 5 presents the results, which are separated for CMAs/CAs and SLAs. The results are presented, in order, for: direct estimates, estimates from the basic FH model unraked and raked, estimates from the extended FH model unraked and raked, and for estimates from the basic FH model unraked and raked, but with the same variables selected as with the extended model. With the standard direct LFS estimates, we start with large average and median absolute relative differences, especially in SLAs. Note that some CMAs have very good direct estimates. Every SAE method reduces the ARD substantially. Ignoring the results of the last two lines of the table, for CMA/CAs, the basic FH model does best in terms of average ARD and the extended FH in terms of median ARD. For SLAs, the extended FH model is preferable. Raking reduces the ARD in almost every case considered. In this analysis, there are two confounding factors: the variables selected in the model, and the variance matrix used in the estimation processes. The last two lines of Table 5 were included to study the latter factor in isolation. Using the basic model with the variables chosen with the extended model improves the average results of the other (unraked) EBLUPs, even doing better than the extended FH, except for median results in SLAs. It also improves all raked results. In theory, the extended FH without raking would be the EBLUP if the working design variance matrix was perfect. Although using the working variance matrix for variable selection but nothing else is a peculiar and unlikely strategy, those results are another indication that the working covariance matrix might need improvements, for CMAs and CAs in particular.

Table 5 – Mean and median ARD of direct, basic FH and extended FH estimates

	CMA/CA		SLA	
	Mean ARD	Median ARD	Mean ARD	Median ARD
Direct estimate	22.5%	10.3%	68.5%	60.8%
Basic FH unraked	4.4%	3.4%	17.3%	12.4%
Basic FH raked	4.4%	3.4%	13.2%	9.2%
Extended FH unraked	4.7%	3.2%	13.3%	7.8%
Extended FH raked	4.5%	3.3%	12.7%	7.4%
Basic FH unraked (extended FH model variables)	4.1%	3.2%	13.0%	8.4%
Basic FH raked (extended FH model variables)	4.0%	2.8%	11.7%	7.0%

5. CONCLUSION

This paper presented the application of the Fay-Herriot model extended for dependent sampling errors to the Canadian Labour Force Survey. It showed many advantages of the extended model over the basic model for that application. They are clear because of the relationship between the variable of interest (total employment) and variables used in provincial weight calibration of the survey creating significant negative correlations between the direct estimates. This motivates the use of a method such as that of Verret & Walker (2024) to obtain a working design covariance matrix. Although the results are improved compared to those obtained under the independence assumption, they depend on the quality of the approximation of the covariance terms of that matrix. There are indications that further improvements are needed. An easy improvement to consider is including more constraints in the hypothetical raking. Calibration of the LFS weights is also done to the 15 years and over population of the largest regions, so including those constraints in the hypothetical raking should be beneficial. In all cases, the effect of the deviations in the covariances must be studied, and we plan to do so in a simulation context.

REFERENCES

- Beaumont, J.-F., and Bocci, C. (2016). “Small Area Estimation in the Labour Force Survey”, Paper presented at Statistics Canada’s Advisory Committee on Statistical Methods, March 31, 2016.
- Dagum, E.B., and P. Cholette (2006). *Benchmarking, Temporal Distribution, and Reconciliation Methods for Time Series*, New York: Springer-Verlag, Lecture Notes in Statistics 186.

Estevao, V., You, Y., Hidioglou, M., Beaumont, J.-F. and Rubin-Bleuer, S. (2023). “Small Area Estimation-Area Level Model with EBLUP Estimation- Methodology Specifications”, Statistics Canada document.

Fay, R.E., and Herriot, R.A. (1979). “Estimation of Income from Small Places: An Application of James-Stein Procedures to Census Data”. *Journal of the American Statistical Association*, **74**, 269-277.

Lesage, É., Beaumont, J.-F. and Bocci, C. (2021). “Two local diagnostics to evaluate the efficiency of the empirical best predictor under the Fay-Herriot model”. *Survey Methodology*, **47**, 279-297.

OECD (2020). Delineating Functional Areas in All Territories, OECD Territorial Reviews, OECD Publishing, Paris, <https://doi.org/10.1787/07970966-en>.

Rao, J.N.K., and Wu, C.F.J. (1988), “Resampling inference with complex survey data”, *Journal of the American Statistical Association*, **83**, 231-241.

Statistics Canada (2020). Guide to the Labour Force Survey, Catalogue no. 71-543-G. <https://www150.statcan.gc.ca/n1/pub/71-543-g/71-543-g2020001-eng.htm>.

Verret, F., and Walker, B. (2024). “Reverse-engineering a Hypothetical Raking Process for the Estimation of Mean Squared Error of Raked Small Area Estimates”. *Proceedings of Statistics Canada Symposium 2024*, <https://www150.statcan.gc.ca/n1/pub/11-522-x/2025001/article/00006-eng.pdf>.