

M-quantile Area-level Models for Robust Small Area Estimation without Reliance on Design-based Variances

María Bugallo Porto

In collaboration with:

María José Lombardía and Alexandro Aneiros-Batista

MODES research groups & CITIC, University of A Coruña, Spain

✿ Problem of interest

What is our motivation for this research?

Main goal

Answer the requests of the research project: “Small Area Estimation for the use of auxiliary information in survey data disaggregation”.

Funding institution

Spanish National Institute of Statistics (INE)

Work plan (2025–2026)

- Estimate labour and wage indicators in unplanned domains with small sample sizes
- Perform analysis for specific regions and subpopulations

Methodology

Small Area Estimation (SAE) techniques

Spanish Quarterly Labour Cost Survey

Within a project funded by the Spanish National Institute of Statistics

Estimating the effective hours worked in small areas in Spain, benchmarking results against the following aggregates:

- ✓ Autonomous Communities and sectors of activity.
- ✓ Divisions and company size groups.



- Unplanned domains with no sample.

- In the **2^o quarter of 2024**, domains with data are reduced to $D = 2500$.

Small Area Estimation Problem

Sample information is scarce in many estimation domains.

	0%	20%	40%	60%	80%	100%	0%	20%	40%	60%	80%	100%
n_d	1	3	5	8	11	74	5	25	43	62	79	96
N_d	2	10	29	85	369	37,263	7	54	136	355	1,100	37,604
f_d	0.12	3.03	8.33	16.67	30.77	85.71	0.15	1.90	4.15	7.95	14.51	66.67

(a) NUTS 2 \times NACE \times Size.

(b) NUTS 2 \times NACE.

Table 1: Summary of some descriptive results.

Small Area Estimation Problem

Sample information is scarce in many estimation domains.

	0%	20%	40%	60%	80%	100%	0%	20%	40%	60%	80%	100%
n_d	1	3	5	8	11	74	5	25	43	62	79	96
N_d	2	10	29	85	369	37,263	7	54	136	355	1,100	37,604
f_d	0.12	3.03	8.33	16.67	30.77	85.71	0.15	1.90	4.15	7.95	14.51	66.67

(a) NUTS 2 x NACE x Size.

(b) NUTS 2 x NACE.

Table 1: Summary of some descriptive results.

The **Hájek estimator** of the population mean is given by:

$$\widehat{Y}_d^{dir} = \frac{1}{\widehat{N}_d} \sum_{j \in s_d} \omega_{dj} y_{dj}, \quad \widehat{N}_d = \sum_{j \in s_d} \omega_{dj}.$$

* Background and motivation

Review of existing area-level approaches

Classic: **Fay-Herriot model** (LMM, Fay and Herriot, 1979):

$$\widehat{Y}_d^{dir} = \bar{\mathbf{X}}_d' \boldsymbol{\beta} + u_d + e_d, \quad u_d \sim N(0, \sigma_u^2), \quad d = 1, \dots, D.$$

Review of existing area-level approaches

Classic: **Fay-Herriot model** (LMM, Fay and Herriot, 1979):

$$\widehat{Y}_d^{dir} = \bar{\mathbf{X}}_d' \boldsymbol{\beta} + \mathbf{u}_d + e_d, \quad \mathbf{u}_d \sim N(0, \sigma_u^2), \quad d = 1, \dots, D.$$

Recent: **Fay-Herriot M-quantile model** (MQ, Marchetti et al., 2025):

$$\widehat{Y}_d^{dir} = \bar{\mathbf{X}}_d' \boldsymbol{\beta}_\psi(\mathbf{q}_{d,\sigma_d}) + e_{\psi,d}, \quad \mathbf{q}_{d,\sigma_d} \in (0, 1), \quad d = 1, \dots, D.$$

Review of existing area-level approaches

Classic: **Fay-Herriot model** (LMM, Fay and Herriot, 1979):

$$\widehat{Y}_d^{dir} = \bar{\mathbf{X}}_d' \boldsymbol{\beta} + \mathbf{u}_d + e_d, \quad \mathbf{u}_d \sim N(0, \sigma_u^2), \quad d = 1, \dots, D.$$

Recent: **Fay-Herriot M-quantile model** (MQ, Marchetti et al., 2025):

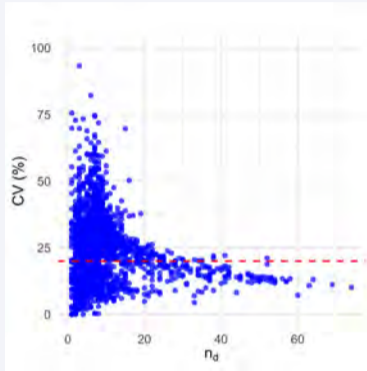
$$\widehat{Y}_d^{dir} = \bar{\mathbf{X}}_d' \boldsymbol{\beta}_\psi(\mathbf{q}_{d,\sigma_d}) + e_{\psi,d}, \quad \mathbf{q}_{d,\sigma_d} \in (0, 1), \quad d = 1, \dots, D.$$

Common features ($d = 1, \dots, D$)

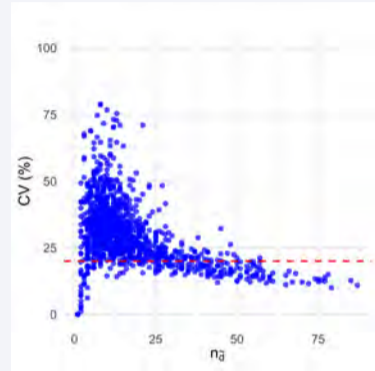
- Model errors: $e_d, e_{\psi,d} \sim N(0, \sigma_d^2)$
- Known sampling variances

$$\widehat{\sigma}_d^2 := \frac{1}{\widehat{N}_d^2} \sum_{j \in s_d} \omega_{dj} (\omega_{dj} - 1) (y_{dj} - \widehat{Y}_d^{dir})^2.$$

Why not rely on the design-based variances?



(a) NUTS 2 x NACE x Size.



(b) NUTS 2 x NACE.

Figure 1: Design-based CV estimates for the Hájek estimator of effective hours versus small-area sample sizes. Red line: CV = 20%.

✿ Methodological proposal

Less restrictive area-level MQ models

Extending the unit-level approach by **Chambers and Tzavidis (2006)**

Novel: For each quantile index $q \in (0, 1)$,

$$\widehat{Y}_d^{dir} = \overline{\mathbf{X}}_d' \boldsymbol{\beta}_\psi(q) + e_{\psi,d}(q), \quad d = 1, \dots, D,$$

where $Q_q(e_{\psi,d}(q); \sigma_q, \psi | \overline{\mathbf{X}}_d) = 0$; and common variance σ_q^2 .

Less restrictive area-level MQ models

Extending the unit-level approach by **Chambers and Tzavidis (2006)**

Novel: For each quantile index $q \in (0, 1)$,

$$\widehat{Y}_d^{dir} = \overline{\mathbf{X}}_d' \boldsymbol{\beta}_\psi(q) + e_{\psi,d}(q), \quad d = 1, \dots, D,$$

where $Q_q(e_{\psi,d}(q); \sigma_q, \psi | \overline{\mathbf{X}}_d) = 0$; and common variance σ_q^2 .

For SAE: A **distribution-free** strategy captures inter-area variability:

$$q_d = \underset{0 < q < 1}{\text{solution}} \left\{ \widehat{Y}_d^{dir} = \overline{\mathbf{X}}_d' \boldsymbol{\beta}_\psi(q) \right\}.$$

After fitting, a grid-search is performed in practice for each domain.

Linear prediction and MSE estimation

Plug-in: Based on the new AMQ models for $q = \hat{q}_d$,

$$\hat{Y}_d^{amq} = \bar{\mathbf{X}}_d' \hat{\beta}_\psi(\hat{q}_d), \quad d = 1, \dots, D.$$

Linear prediction and MSE estimation

Plug-in: Based on the new AMQ models for $q = \hat{q}_d$,

$$\hat{Y}_d^{amq} = \bar{\mathbf{X}}_d' \hat{\beta}_\psi(\hat{q}_d), \quad d = 1, \dots, D.$$

Two options for the MSE estimation of the predictors:

Analytical: Based on a Taylor series expansion:

$$\begin{aligned} mse_d^{anlt} &= (\mathbf{a}'_d - \mathbf{1}'_d) \operatorname{diag}(\hat{\sigma}_g^2)_{1 \leq g \leq D} (\mathbf{a}'_d - \mathbf{1}'_d)' + (2a_{dd} - 1) \hat{\sigma}_d^2 \\ &+ \left(\sum_{g=1}^D a_{dg} \bar{\mathbf{X}}_g' \hat{\beta}_\psi(\hat{q}_g) - \hat{Y}_d^{dir} \right)^2. \end{aligned}$$

Linear prediction and MSE estimation

Plug-in: Based on the new AMQ models for $q = \hat{q}_d$,

$$\hat{Y}_d^{amq} = \bar{\mathbf{X}}_d' \hat{\beta}_\psi(\hat{q}_d), \quad d = 1, \dots, D.$$

Two options for the MSE estimation of the predictors:

Analytical: Based on a Taylor series expansion:

$$\begin{aligned} mse_d^{anlt} &= (\mathbf{a}'_d - \mathbf{1}'_d) \operatorname{diag}(\hat{\sigma}_g^2)_{1 \leq g \leq D} (\mathbf{a}'_d - \mathbf{1}'_d)' + (2a_{dd} - 1) \hat{\sigma}_d^2 \\ &+ \left(\sum_{g=1}^D a_{dg} \bar{\mathbf{X}}_g' \hat{\beta}_\psi(\hat{q}_g) - \hat{Y}_d^{dir} \right)^2. \end{aligned}$$

Bootstrap: More flexible semiparametric approach (...):

$$mse_d^{boot} = \frac{1}{B} \sum_{b=1}^B \left(\hat{Y}_d^{amq*(b)} - \bar{Y}_d^{*(b)} \right)^2.$$

Data Requirements

Linear Prediction and MSE estimation

	Area-level				Unit-level	
	Hájek	FH	FHMQ	AMQ	MQ	NER
y_{dj}	Green				Green	
ω_{dj}	Green					
x'_{dj}					Green	
$\overline{\mathbf{X}}'_d$		Green	Green	Green	Green	
\widehat{Y}_d^{dir}		Green	Green	Green		
$\widehat{\sigma}_d^2$		Green	Green	Yellow		
	Minimal	Strong	Moderate	Minimal	Moderate	Strong

Color coding: green = required, yellow = optional, white = not required.

✿ Empirical results

Summary of the simulation results (I)

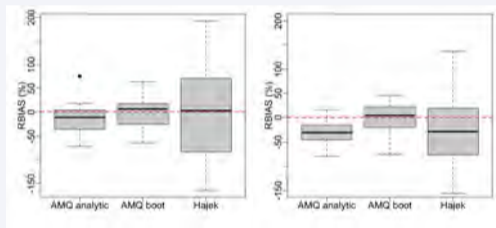
Simulation 1: Performance of population mean predictors

	No outliers	Outliers	
	[0, 0]	[3% e, 0]	[3% e, 10% u]
FH model	✓✓		
FHMQ model		✓✓	✓
AMQ model		✓	✓✓

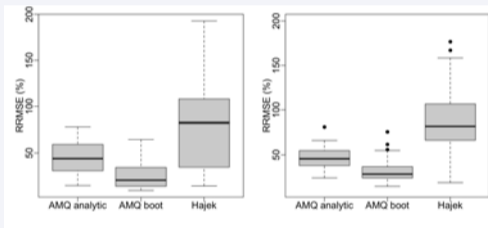
Table 2: Prediction performance in model-based simulations; $n_d = 5$.

Summary of the simulation results (II)

Simulation 2: Reliability of the variance and MSE estimators



(a) Relative bias: $[0,0]$ & $[e,u]$.



(b) Relative root-MSE: $[0,0]$ & $[e,u]$.

Figure 2: Performance of the variance and MSE estimators; $n_d = 5$.

Summary of the simulation results (III)

Effect of the sample sizes on the MSE estimation

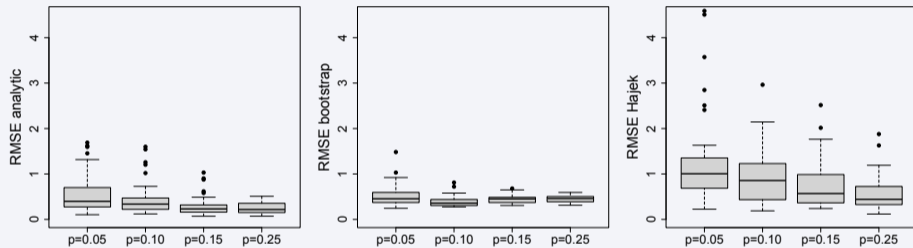


Figure 3: Performance of MSE (analytical & bootstrap) and design-based variance estimates for $D = 40$, $N = 200$ and $n = Np \in \{10, 20, 30, 50\}$.

* Application to real data

Spanish Quarterly Labour Cost Survey

Within a project funded by the Spanish National Institute of Statistics

Estimating the effective hours worked in small areas in Spain, benchmarking results against the following aggregates:

- ✓ Autonomous Communities and sectors of activity.
- ✓ Divisions and company size groups.



- Unplanned domains with no sample.

- In the **2^o quarter of 2024**, domains with data are reduced to $D = 2500$.

Prediction and root-MSE estimation

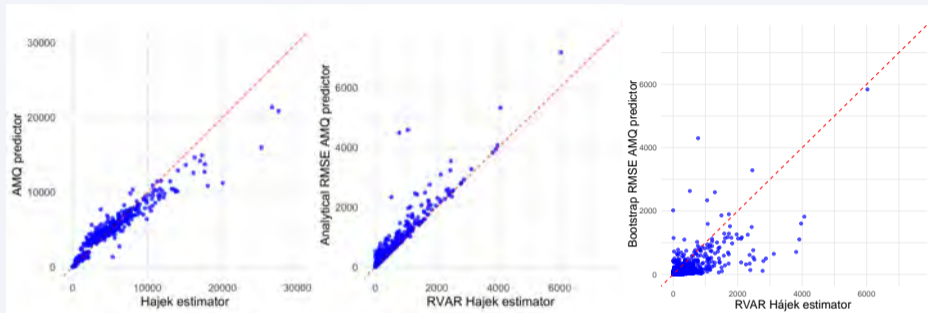


Figure 4: Left: AMQ predictor vs Hájek estimator for population means of effective hours; Center & right: analytical and semiparametric bootstrap RMSEs of AMQ ($B = 100$) vs design-based RVARs.

- Aligned AMQ predictions and Hájek estimates.
- Reduced root-MSE estimates for the AMQ predictor.

Bootstrap & design-based coefficients of variation

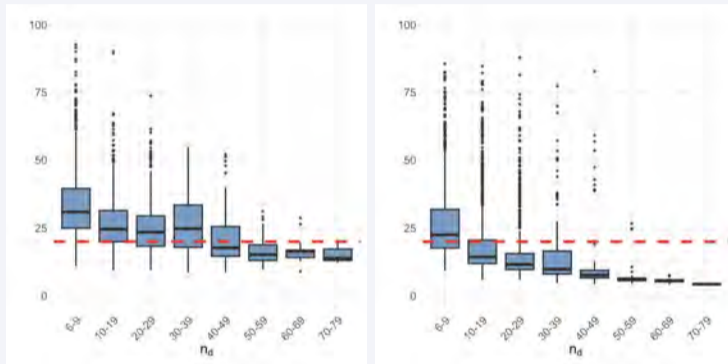


Figure 5: Left: CV bootstrap of the AMQ predictor; Right: Design-based CV of the Hájek estimator. Red line: CV=20%.

Main contribution

This work extends Small Area Estimation methods to business surveys using flexible M-quantile area models.

- ✓ **More reliable inference:** Hájek design-based variances are poorly estimated.
- ✗ **Benchmarking warning:** Not recommended under high heterogeneity across domains.
- 👉 **Next steps:** Temporal extensions and deeper analysis of benchmarking.

Improving robustness and flexibility in business survey estimation

THANK YOU