

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25

Smoothing and Benchmarking for Small Area Estimation

Rebecca C. Steorts*, Timo Schmid[†] and Nikos Tzavidis[‡]

February 2, 2020

Abstract

Small area estimation is concerned with methodology for estimating population parameters associated with a geographic area defined by a cross-classification that may also include non-geographic dimensions. In this paper, we develop constrained estimation methods for small area problems: those requiring smoothness with respect to similarity across areas, such as geographic proximity or clustering by covariates; and benchmarking constraints, requiring weighted means of estimates to agree across levels of aggregation. We develop methods for constrained estimation decision-theoretically and discuss their geometric interpretation. The constrained estimators are the solutions to tractable optimization problems and have closed-form solutions. Mean squared errors of the constrained estimators are calculated via bootstrapping. Our approach assumes the Bayes estimator exists and is applicable to any proposed model. In addition, we give special cases of our techniques under certain distributional assumptions. We illustrate the proposed methodology using web-scraped data on Berlin rents aggregated over areas to ensure privacy.

1 Introduction

Small area estimation (SAE) deals with estimating many parameters, each associated with an “area”—a geographic domain, a demographic group, an

*Department of Statistical Science and Computer Science, Duke University, Durham, North Carolina, Email: beka@stat.duke.edu

[†]Freie Universität Berlin Institute of Statistics and Econometrics, Berlin, Germany, Email: timo.schmid@fu-berlin.de

[‡]University of Southampton Southampton Statistical Sciences Research Institute, Southampton, United Kingdom, Email: n.tzavidis@soton.ac.uk

26 experimental condition, etc. Areas are “small” since there is little or no
27 information about any one area. Estimates of a parameter based only on
28 observations from the associated area, called direct estimates, can be im-
29 precise. To increase precision, one tries to “borrow strength” from related
30 areas, and hierarchical and empirical Bayesian models are one way to do
31 so. Since the pioneering work of [Fay and Herriot \(1979\)](#) and [Battese et al.
32 \(1988\)](#), such models have dominated SAE, with many successful applica-
33 tions in official statistics, sociology, epidemiology, political science, business,
34 etc. ([Rao and Molina 2015](#)). Recently, SAE has been applied in other fields,
35 such as neuroscience, and performs as well as common approaches such as
36 smoothed ridge regression and elastic net ([Wehbe et al. 2015](#)).

37 We extend these classical approaches in two directions, both of which
38 have been the subject of recent interest in the SAE literature. One direc-
39 tion is to take direct account of information about the proximity of areas
40 in space or time. In many applications, it is reasonable to expect that the
41 parameters will be smooth, so that nearby areas will have similar parame-
42 ters, but this is not altogether standard within SAE ([Rao and Molina 2015](#)).
43 Incorporating spatial dependence directly into Bayesian models leads to sta-
44 tistical and computational difficulties, yet it seems misguided to discard such
45 information. The other direction is “benchmarking,” the imposition of con-
46 sistency constraints on (weighted) averages of the parameter estimates. A
47 simple form of benchmarking is when the average of the parameter estimates
48 must match a known global average. When there are multiple levels of ag-
49 gregation for the estimates, there can be issues of internal consistency as
50 well.

51 We provide a unified approach to smoothing and benchmarking by re-
52 garding them both as *constraints* on Bayes estimates. Benchmarking corre-
53 sponds to equality constraints on global averages and variances. Similarly,
54 smoothing corresponds to an inequality constraint on the “roughness” of
55 estimates (how much the parameter estimates of nearby areas differ). The
56 motivation of this smoothing is based upon manifold learning and frequen-
57 tist non-parametrics, where loss functions are augmented by a penalty. Such
58 a penalty term is in the spirit of ridge regression, where a transformation
59 of the parameters is performed and additional shrinkage is carried out. Our
60 penalty corresponds to how much estimates at nearby points in the domain
61 should differ.

62 Decision-theoretically, we obtain smoothed, benchmarked estimates by
63 minimizing the Bayes risk subject to these constraints, extending the ap-
64 proaches of [Datta et al. \(2011\)](#) and [Ghosh and Steorts \(2013\)](#) (themselves in
65 the spirit of [Louis \(1984\)](#) and [Ghosh \(1992\)](#)). Geometrically, the constrained

66 Bayes estimates are found by projecting the unconstrained estimates into
67 the feasible set. If the constraints are linear, then the resulting optimization
68 can be solved in closed form, requiring nothing more than basic matrix oper-
69 ations on the unconstrained Bayes estimates. Another strong advantage of
70 our decision-theoretic and geometric approach is its generality. We require
71 no distributional assumptions on the data or on the unconstrained Bayes es-
72 timator. Our results apply whether the unconstrained estimator is linear or
73 non-linear. The relevant notion of proximity between areas may be spatial
74 or more abstract. It can also include clustering on covariates not directly
75 included in the model. Finally, we are able to prove known cases under our
76 proposed approach, where the Bayes and frequentist estimates are in fact
77 the same. Finally, we illustrate our proposed methodology on rental prices
78 in Berlin.

79 The rest of the paper proceeds as follows. Section 2 describes related
80 work. Section 3 provides the proposed general framework for smoothing
81 in small area estimation. Section 3.1, introduces notation used throughout
82 the paper. Section 3.2 proposes a general result for SAE in the context
83 of smoothing. Section 3.3 proposes special cases of our general framework
84 under the area-level Fay-Herriot model. Section 4 extends our generalized
85 result in Section 3.2 to benchmarking constraints. Section 4.2 derives special
86 cases under the area-level Fay-Herriot model. Section 4.3 discusses choices
87 of the smoothness penalties. Section 5 proposes a non-parametric boot-
88 strap for mean squared error (MSE) estimation. Section 6 applies the pro-
89 posed methodology to web-scraped data for estimating average rent prices
90 in Berlin. Section 7 concludes with a discussion and future work.

91 2 Related Work

92 The proposed methodology for SAE with benchmarking and smoothing gen-
93 eralizes the work of [Datta et al. \(2011\)](#) and [Ghosh and Steorts \(2013\)](#), which
94 take a decision theoretic approach to SAE. However, this literature does not
95 allow for spatial smoothing. The approach proposed in this paper is also
96 similar to that of [Wehbe et al. \(2015\)](#) in the sense that spatial smooth-
97 ing is considered; however, these authors do not consider benchmarking.
98 Moreover, the authors focus more on a neuroscience application and less
99 on developing a general methodology for SAE methodology. Other relevant
100 literature includes [Pratesi and Salvati \(2008\)](#), who proposed a spatial em-
101 pirical best linear unbiased predictor under the Fay-Herriot model with a
102 simultaneous autoregressive (SAR) structure for the random effects and an

103 analytic based MSE. [Souza et al. \(2009\)](#) account for spatial relationships
104 when fitting hierarchical Bayesian exponential growth models. [Rao and Yu](#)
105 [\(1994\)](#) proposed a linking model that does not include area-specific random
106 effects in essence to avoid over-smoothing, which is a valid concern when
107 proposing any type of smoothing constraint.

108 Previous efforts at smoothing in SAE problems have smoothed either
109 the raw data or direct estimates. In contrast, we smooth estimates based on
110 models which do not themselves include spatial structure. Computationally,
111 this is much easier than expanding the models. Our optimization problems
112 can be solved in closed form and retain the advantages of model-based es-
113 timation. This approach to smoothing also combines naturally with the
114 imposition of benchmarking constraints, which has never been handled to
115 our knowledge in the literature before.

116 Our proposed methodology employs ideas about smoothing on graphs
117 and manifolds from frequentist non-parametrics and machine learning. In
118 particular, we take advantage of “Laplacian” regularization ideas ([Belkin](#)
119 [et al. 2006](#); [Corona et al. 2008](#); [Lee and Wasserman 2010](#)), where the loss
120 function is augmented by a penalty term which reflects how much estimates
121 at nearby points in the domain differ. Such regularization is designed to
122 ensure that estimates vary smoothly with respect to the intrinsic geometry
123 of some underlying graph or manifold. (Smoothness on a domain is rep-
124 resented mathematically by the domain’s Laplacian operator, which is the
125 generator for diffusion processes.) This generalizes the roughness or curva-
126 ture penalties from spline smoothing ([Wahba 1990](#)) to domains geometrically
127 more complicated than \mathbb{R}^d . We are unaware of any previous application of
128 Laplacian regularization to SAE problems, though spline smoothing is often
129 used in spatial statistics, including traditional SAE applications to disease
130 mapping ([Kafadar 1996](#)).

131 **3 Smoothing for Small Area Estimation**

132 In this section, we provide a generalized approach for SAE. First, we intro-
133 duce notation used throughout the paper (Section [3.1](#)). Second, we provide
134 a general result for SAE in the context of smoothing (Section [3.2](#)). Third,
135 we provide special cases of this result under the area-level Fay-Herriot model
136 (Section [3.3](#)).

137 **3.1 Notation and Terminology**

138 In this section, we present general notation that is used throughout the
 139 remainder of the paper. We assume m areas, and for each area i , we estimate
 140 an associated scalar quantity θ_i , collectively $\boldsymbol{\theta} = (\theta_1, \theta_2, \dots, \theta_m)$. We denote
 141 the response (direct estimator) by $\mathbf{y} = (y_1, y_2, \dots, y_m)$. “Areas” are often
 142 spatial regions, but they might be different demographic groups. Our goal is
 143 to estimate the unknown parameter $\boldsymbol{\theta}$ by some estimator $\boldsymbol{\delta} = (\delta_1, \delta_2, \dots, \delta_m)$,
 144 where we denote the optimal estimator by $\hat{\boldsymbol{\delta}} = (\hat{\delta}_1, \hat{\delta}_2, \dots, \hat{\delta}_m)$.

145 Denote the i th area by a (column) vector of covariates

$$\mathbf{x}_i = \begin{pmatrix} x_{i1} \\ x_{i2} \\ \vdots \\ x_{ip} \end{pmatrix},$$

which may include spatial coordinates. We can represent the covariates as
 a design matrix in the following way:

$$X_{m \times p} = \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1p} \\ x_{21} & x_{22} & \dots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ x_{m1} & x_{m2} & \dots & x_{mp} \end{pmatrix}.$$

146 One Bayesian treatment of this problem of finding an optimal estimator
 147 is to define a loss function, and then minimize the posterior risk. That
 148 is, under a defined loss function $L(\boldsymbol{\theta}, \boldsymbol{\delta})$, one goal will be to minimize the
 149 posterior risk with respect to the estimator $\boldsymbol{\delta}$.

Turning to the loss function, we will assume for convenience and for
 the desirability of tractable solutions that our loss function is a weighted
 squared error, where the weight for area i is $\phi_i > 0$, which can be denoted
 by a matrix of weights, Φ . The total loss is denoted by

$$L(\boldsymbol{\theta}, \boldsymbol{\delta}) = \sum_{i=1}^m \phi_i (\theta_i - \delta_i)^2 = (\boldsymbol{\theta} - \boldsymbol{\delta})^T \Phi (\boldsymbol{\theta} - \boldsymbol{\delta}).$$

150 In many SAE applications, the weights Φ reflect variations in measure-
 151 ment precision and can be estimated from the survey design (Pfeffermann
 152 2013; Rao and Molina 2015; Tzavidis et al. 2018). There exists a large lit-
 153 erature regarding proposals for estimating loss function weights. Isaki et al.

154 (2004) proposed taking each weight as the reciprocal of the posterior vari-
 155 ance of the Bayes estimator. For a full review of such choices for the loss
 156 function weights, we refer to Datta et al. (2011), and stress that the choice
 157 of the loss function weights is application specific.

158 Under our proposed methodology, we simply assume that a Bayes esti-
 159 mator exists, and under this framework, the modeling structure can be set
 160 by the user. Of course, in certain situations, the Bayes estimator is the same
 161 as other estimators in the frequentist literature, such as the Best Linear Un-
 162 biased Predictor (BLUP). A full review of such cases can be found in Molina
 163 and Rao (2010) and Ghosh et al. (1994).

164 3.2 General Result

In this section, we propose our general framework for smoothing in SAE. Before doing so, we introduce new terminology that is needed for the remainder of the paper. Consider two different areas i and i' , where $i \neq i'$. We define a symmetric matrix, Q , with elements $q_{ii'} \geq 0$, to control how important it is that the estimate of θ_i is close to the estimate of $\theta_{i'}$. It may often be the case that $q_{ii'} = q(\mathbf{x}_i, \mathbf{x}_{i'})$; i.e., the degree of smoothing of δ_i and $\delta_{i'}$ is a function of the covariates \mathbf{x}_i and $\mathbf{x}_{i'}$. Note also that the $q_{ii'}$ may be discrete-valued, corresponding to clustering of areas, or continuous-valued, corresponding to a metric space of areas. Writing Q in matrix form, we see that

$$Q = \begin{pmatrix} q_{11} & q_{12} & \cdots & q_{1m} \\ q_{21} & q_{22} & \cdots & q_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ q_{m1} & q_{m2} & \cdots & q_{mm} \end{pmatrix}.$$

165

A natural measure of the smoothness of δ_i is the Q -weighted sum of squared differences between elements, $\sum_{i=1}^m \sum_{i'=1}^m (\delta_i - \delta_{i'})^2 q_{ii'}$, where for the remainder of the paper we denote $\sum_{i=1}^m \sum_{i'=1}^m$ as $\sum_{i,i'}$. We add a penalty term

$$\gamma \sum_{i,i'} (\delta_i - \delta_{i'})^2 q_{ii'}$$

166 to our objective function, with the penalty factor $\gamma \geq 0$ chosen to specify
 167 the overall importance of smoothness. (We address the choice of Q below
 168 and of γ in Section 4.3.)

169 Therefore, we seek to minimize the posterior risk of the loss function

$$L(\boldsymbol{\theta}, \boldsymbol{\delta}) = \sum_i \phi_i(\theta_i - \delta_i)^2 + \gamma \sum_{i,i'} (\delta_i - \delta_{i'})^2 q_{ii'}. \quad (3.1)$$

170 Minimizing the posterior expectation of equation 3.1 is equivalent to mini-
 171 mizing

$$\sum_i \phi_i E[(\theta_i - \delta_i)^2 | \mathbf{y}] + \gamma \sum_{i,i'} (\delta_i - \delta_{i'})^2 q_{ii'}. \quad (3.2)$$

172 Finally, we define Ω to be a matrix such that $\sum_{i,i'} (\delta_i - \delta_{i'})^2 q_{ii'} = \boldsymbol{\delta}^T \Omega \boldsymbol{\delta}$,
 173 where Ω is a semi-positive definite matrix. See Lemma A in the Supplemen-
 174 tary Material for the above equivalence.

175 Then minimizing equation 3.2 is equivalent to minimizing

$$(\boldsymbol{\delta} - \hat{\boldsymbol{\theta}}^B)^T \Phi (\boldsymbol{\delta} - \hat{\boldsymbol{\theta}}^B) + \gamma \boldsymbol{\delta}^T \Omega \boldsymbol{\delta}, \quad (3.3)$$

176 where $\hat{\boldsymbol{\theta}}^B = (\hat{\theta}_1^B, \dots, \hat{\theta}_m^B)$. See Datta et al. (2011) and Ghosh and Steorts
 177 (2013) for details on this equivalence. Then we have the following result.

178 **Theorem 3.1.** *The smoothed Bayes estimator is*

$$\tilde{\boldsymbol{\theta}}^S = (I_m + \gamma \Phi^{-1} \Omega)^{-1} \hat{\boldsymbol{\theta}}^B.$$

179 *Proof.* Differentiating equation 3.3 with respect to $\boldsymbol{\delta}$ and setting the gradient
 180 to zero at $\tilde{\boldsymbol{\theta}}^S$ yields $\Phi(\tilde{\boldsymbol{\theta}}^S - \hat{\boldsymbol{\theta}}^B) + \gamma \Omega \tilde{\boldsymbol{\theta}}^S = \mathbf{0}$. Then

$$(\Phi + \gamma \Omega) \tilde{\boldsymbol{\theta}}^S = \Phi \hat{\boldsymbol{\theta}}^B \implies \tilde{\boldsymbol{\theta}}^S = (I_m + \gamma \Phi^{-1} \Omega)^{-1} \hat{\boldsymbol{\theta}}^B.$$

181 Since equation 3.3 is a positive-definite quadratic form in $\boldsymbol{\delta}$, the solution is
 182 unique. \square

183 3.3 Area-Level Fay-Herriot Model

184 In this section, we consider a special case of our general result in Section
 185 3.2, where our only assumption was that the Bayes estimate exists. In this
 186 section, we consider a special case of Theorem 3.1 in Section 3.2, where we
 187 assume the standard Fay-Herriot model (Fay and Herriot 1979).

Before proceeding, we review the Fay-Herriot model and a few standard results that follow from assuming this model, which is a special case of the general framework proposed in Section 3.2. More specifically, we consider the area-level Fay-Herriot model

$$\begin{aligned} \mathbf{y}_{m \times 1} &= \boldsymbol{\theta}_{m \times 1} + \mathbf{e}_{m \times 1} \\ \boldsymbol{\theta}_{m \times 1} &= X_{m \times p} \boldsymbol{\beta}_{p \times 1} + \mathbf{u}_{m \times 1}, \end{aligned} \quad (3.4)$$

where

$$\mathbf{e}_{m \times 1} \stackrel{ind}{\sim} \text{MVN}(\mathbf{0}, D) \quad \text{and} \quad \mathbf{u}_{m \times 1} \stackrel{ind}{\sim} \text{MVN}(\mathbf{0}, \sigma_u^2 I_m), \quad (3.5)$$

and

$$D_{m \times m} = \text{Diag}(D_1, \dots, D_m),$$

$$B_{m \times m} = \text{Diag}(D_1(\sigma_u^2 + D_1)^{-1}, \dots, D_m(\sigma_u^2 + D_m)^{-1}).$$

188 Note that $\boldsymbol{\beta}$ is a $p \times 1$ vector of unknown regression coefficients and the
 189 $\text{rank}(X) = p (< m)$. Equation 3.4 was first considered in the context of
 190 estimating income for small areas by [Fay and Herriot \(1979\)](#).

191 If both σ_u^2 or D are unknown, then the model is not identifiable. This
 192 is seen by writing the marginal distribution of \mathbf{y} , which can be shown to be
 193 $\mathbf{y} \sim \text{MVN}(X\boldsymbol{\beta}, V)$, where $V = \text{Diag}(\sigma_u^2 + D_1, \dots, \sigma_u^2 + D_m)$. There is clearly
 194 an identifiability issue when both σ_u^2 or D are unknown. (In a Bayesian
 195 setting, this is the marginal distribution of \mathbf{y} after integrating out $\boldsymbol{\theta}$).

When the variance component σ_u^2 is known and $\boldsymbol{\beta}$ has a uniform prior on \mathbb{R}^p , then the Bayes estimator of $\boldsymbol{\theta}$ is

$$\hat{\boldsymbol{\theta}}^B = \hat{\boldsymbol{\theta}}^{\text{BLUP}} = (I_m - B_{m \times m})\mathbf{y}_{m \times 1} + B_{m \times m}X_{m \times p}\tilde{\boldsymbol{\beta}}_{p \times 1}, \quad (3.6)$$

196 where $\tilde{\boldsymbol{\beta}} \equiv \tilde{\boldsymbol{\beta}}(\sigma_u^2) = (X'V^{-1}X)^{-1}X'V^{-1}\mathbf{y}$. It is well known that under
 197 the conditions given above, the Bayes estimator is also the best unbiased
 198 predictor of $\boldsymbol{\theta}$ ([Datta and Ghosh 1991](#); [Datta et al. 2011](#); [Rao and Molina](#)
 199 [2015](#)).

In more realistic settings, σ_u^2 is unknown and must be estimated. Thus, the empirical Bayes estimator becomes

$$\hat{\boldsymbol{\theta}}^{EB} = (I_m - \hat{B}_{m \times m})\mathbf{y}_{m \times 1} + \hat{B}_{m \times m}X_{m \times p}\tilde{\boldsymbol{\beta}}(\hat{\sigma}_u^2)_{p \times 1}. \quad (3.7)$$

200 Note that $\hat{\sigma}_u^2$ in equation 3.7 can be estimated many different ways. For
 201 example, many common estimators are moment estimators or maximum
 202 likelihood estimation. Moment estimation is quite convenient in particu-
 203 lar situations, as shown in [Prasad and Rao \(1990\)](#) and [Steorts and Ghosh](#)
 204 [\(2013\)](#). However, more general techniques can be found in [Rao and Molina](#)
 205 [\(2015\)](#). Thus, in the standard Fay-Herriot situation, where one considers
 206 this very specialized situation, the Bayes and frequentist estimates are the
 207 same, and one may find these estimates without resorting to Markov chain
 208 Monte Carlo (MCMC).

209 Assuming the standard area-level Fay-Herriot model in equation 3.4, we
 210 prove a special case of Theorem 3.1 in Lemma 3.1.

211 **Lemma 3.1.** *Assume the loss function in equation 3.3 and assume the area-*
 212 *level Fay-Herriot model in equation 3.4. Consider two choices of the Bayes*
 213 *estimate, which are $\hat{\theta}^{BLUP}$ and $\hat{\theta}^{EBLUP}$ and are given in equations 3.6 and*
 214 *3.7. The smoothed best linear unbiased Bayes (SBLUP) estimator is*

$$\tilde{\theta}^{SBLUP} = (I_m + \gamma\Phi^{-1}\Omega)^{-1}\hat{\theta}^{BLUP}.$$

215 *and the smoothed empirical best linear unbiased Bayesian (SEBLUP)*
 216 *estimator is*

$$\tilde{\theta}^{SEBLUP} = (I_m + \gamma\Phi^{-1}\Omega)^{-1}\hat{\theta}^{EBLUP}.$$

217 *Proof.* Under the assumption of the Fay-Herriot model and by equations 3.6
 218 and 3.7, the results follows by direction substitution into Theorem 3.1. \square

219 As already mentioned in Section 3.1, there are many ways to choose the
 220 loss function weights, and this is typically application specific. We define
 221 the loss function weights used in our application in Section 6.

222 4 Benchmarking and Smoothing

223 We now turn to situations where our estimates should not just be smooth,
 224 minimizing equation 3.3, but also obey benchmarking constraints. As the
 225 benchmarking constraints are relaxed, we should recover the results of Sec-
 226 tion 3.

227 **Definition 4.1** (Benchmarking constraints, benchmarked Bayes estimator).
 228 *Benchmarking constraints are equality constraints on the weighted means or*
 229 *weighted variances of subsets (possibly all) of the estimates. The bench-*
 230 *marked Bayes estimator is the minimizer of the posterior risk subject to the*
 231 *benchmarking constraints.*

232 The levels to which we benchmark, i.e., the values of the equality con-
 233 straints, are assumed to be given *externally* from some other data source.
 234 For internal benchmarking, we refer to Bell et al. (2013). Our methods
 235 address linear, weighted mean constraints, in a similar manner to that of
 236 Datta et al. (2011) and Ghosh and Steorts (2013); however, our results are
 237 more general.

238 4.1 General Linear Benchmarking Constraints

239 We now return to our general assumptions in Section 3.2. We first con-
 240 sider benchmarking constraints which are linear in the estimate δ , such as

241 means or totals. The general problem is now to minimize the posterior risk
 242 in equation 3.3 subject to the constraints

$$M\boldsymbol{\delta} = \mathbf{t}, \quad (4.1)$$

243 where \mathbf{t} is a given k -dimensional vector and M is a $k \times m$ matrix. This
 244 is equivalent to introducing a k -dimensional vector of Lagrange multipliers
 245 $\boldsymbol{\lambda}$ and minimizing $(\boldsymbol{\delta} - \hat{\boldsymbol{\theta}}^B)^T \Phi(\boldsymbol{\delta} - \hat{\boldsymbol{\theta}}^B) + \gamma \boldsymbol{\delta}^T \Omega \boldsymbol{\delta} - 2\boldsymbol{\lambda}^T (M\boldsymbol{\delta} - \mathbf{t})$. The full
 246 details on this equivalence can be found in Datta et al. (2011).

247 **Theorem 4.1.** *Suppose that equation 4.1 has solutions. Then the con-*
 248 *strained Bayes estimator under the constraint in equation 4.1 is*

$$\tilde{\boldsymbol{\theta}}^{BM} = \Sigma^{-1} \left[\Phi \hat{\boldsymbol{\theta}}^B + M^T (M \Sigma^{-1} M^T)^{-1} (\mathbf{t} - M \Sigma^{-1} \Phi \hat{\boldsymbol{\theta}}^B) \right],$$

249 where $\Sigma = \Phi + \gamma \Omega$.

250 **Remark 4.1.** *Note that the Theorem 4.1 estimator $\tilde{\boldsymbol{\theta}}^{BM}$ can be expressed*
 251 *in terms of the Theorem 3.1 estimator $\tilde{\boldsymbol{\theta}}^S$ as*

$$\tilde{\boldsymbol{\theta}}^{BM} = \tilde{\boldsymbol{\theta}}^S + \Sigma^{-1} M^T (M \Sigma^{-1} M^T)^{-1} (\mathbf{t} - M \tilde{\boldsymbol{\theta}}^S).$$

252 Thus, it can be seen that the benchmarking essentially “adjusts” the estima-
 253 tor $\tilde{\boldsymbol{\theta}}^S$ based on the discrepancy between $M \tilde{\boldsymbol{\theta}}^S$ and the target \mathbf{t} .

Proof of Theorem 4.1. Differentiating with respect to $\boldsymbol{\delta}$ and setting the re-
 sult equal to zero at $\tilde{\boldsymbol{\theta}}^{BM}$ yields

$$\begin{aligned} M^T \boldsymbol{\lambda} &= \Phi(\tilde{\boldsymbol{\theta}}^{BM} - \hat{\boldsymbol{\theta}}^B) + \gamma \Omega \tilde{\boldsymbol{\theta}}^{BM} \\ \implies \tilde{\boldsymbol{\theta}}^{BM} &= \Sigma^{-1} (\Phi \hat{\boldsymbol{\theta}}^B + M^T \boldsymbol{\lambda}). \end{aligned}$$

Then by the constraint,

$$\begin{aligned} \mathbf{t} &= M \Sigma^{-1} (\Phi \hat{\boldsymbol{\theta}}^B + M^T \boldsymbol{\lambda}) \\ &= M \Sigma^{-1} \Phi \hat{\boldsymbol{\theta}}^B + M \Sigma^{-1} M^T \boldsymbol{\lambda}, \end{aligned} \quad (4.2)$$

254 so $\boldsymbol{\lambda} = (M \Sigma^{-1} M^T)^{-1} (\mathbf{t} - M \Sigma^{-1} \Phi \hat{\boldsymbol{\theta}}^B)$. The result follows immediately. \square

Often there is only one linear constraint of the form $\sum_i w_i \delta_i = t$, or
 equivalently $\mathbf{w}^T \boldsymbol{\delta} = t$, for some nonnegative weights w_i and some $t \in \mathbb{R}$.
 This is simply a special case of Theorem 4.1 with $k = 1$ and $M = \mathbf{w}^T$, in
 which case the result simplifies to

$$\tilde{\boldsymbol{\theta}}^{BM} = \tilde{\boldsymbol{\theta}}^S + (t - \mathbf{w}^T \tilde{\boldsymbol{\theta}}^S) (\mathbf{w}^T \Sigma^{-1} \mathbf{w})^{-1} \Sigma^{-1} \mathbf{w}.$$

255 **Geometric Interpretation:** Our formulation of benchmarking and
 256 smoothing as constrained optimization problems has a geometric interpre-
 257 tation. It is well known that the Bayes estimate is the minimizer of the
 258 conditional expectation of the MSE. Since the minimization is taken over
 259 *all* possible values of θ , the Bayes estimate will not respect any constraints
 260 we might wish to impose (except by chance) or unless these constraints are
 261 included in the specification of the prior. Instead, we minimize the MSE
 262 within the feasible set of the constraints. We find the point in the feasible
 263 set which is as close (in the sense of expected weighted squared error) to the
 264 Bayes estimate as possible. That is, we *project* the Bayes estimate into the
 265 feasible set.

266 The geometry of the feasible set is itself slightly complicated, because of
 267 the constraints imposed. Note that the smoothness penalty in the loss func-
 268 tion may be reformulated as a smoothness constraint of the form $\delta^T \Omega \delta \leq s$
 269 for some $s > 0$. This constraint defines an ellipsoid centered at the ori-
 270 gin. Constraints on weighted means define linear sub-spaces, e.g., planes,
 271 depending on the number of constraints and the number of variables.

272 **Remark 4.2.** *We do not consider benchmarked constraints of weighted vari-*
 273 *abilities in this paper as the problem is non-convex. Geometrically, con-*
 274 *straints of weighted variabilities define the surfaces of cones. The constrained*
 275 *Bayes estimator is the projection of the unconstrained Bayes estimator onto*
 276 *the intersection of the ellipsoid, the linear sub-space, and the cones. We*
 277 *return to this in our discussion of future work.*

278 4.2 Area-Level Fay-Herriot Model with Benchmarking and 279 Smoothing

280 In this section, we return to the assumptions of the area-level Fay-Herriot
 281 model in Section 3.3 and equation 3.4, which allows us to derive a special case
 282 of our generalized approach for smoothing and benchmarking in Lemma 4.1.

Lemma 4.1. *Let us assume the conditions of Theorem 4.1. In addition,*
assume the area-level Fay-Herriot model in equation 3.4. Finally, let us
assume that the Bayes estimator is either $\hat{\theta}^{BLUP}$ or $\hat{\theta}^{EBLUP}$. It immediately
follows that the smoothed, benchmarked BLUP and EBLUP are

$$\tilde{\theta}^{SB-BLUP} = \tilde{\theta}^{SBLUP} + (t - \mathbf{w}^T \tilde{\theta}^{SBLUP})(\mathbf{w}^T \Sigma^{-1} \mathbf{w})^{-1} \Sigma^{-1} \mathbf{w} \quad (4.3)$$

and

$$\tilde{\theta}^{SB-EBLUP} = \tilde{\theta}^{SEBLUP} + (t - \mathbf{w}^T \tilde{\theta}^{SEBLUP})(\mathbf{w}^T \Sigma^{-1} \mathbf{w})^{-1} \Sigma^{-1} \mathbf{w}, \quad (4.4)$$

283 where

$$\tilde{\boldsymbol{\theta}}^{SBLUP} = (I_m + \gamma\Phi^{-1}\Omega)^{-1}\hat{\boldsymbol{\theta}}^{BLUP}.$$

284 and the smoothed empirical best linear unbiased Bayesian (SEBLUP) esti-
 285 mator is

$$\tilde{\boldsymbol{\theta}}^{SEBLUP} = (I_m + \gamma\Phi^{-1}\Omega)^{-1}\hat{\boldsymbol{\theta}}^{EBLUP}.$$

286 *Proof.* In this situation, the result follows directly from Lemma 3.1. \square

287 4.3 Choice of Smoothing Penalties

288 The choice of γ is assumed to be fixed *a priori*. But knowing γ is equivalent
 289 to knowing how smooth the estimate *ought* to be, and this knowledge is
 290 lacking in most applications. In such situations, we suggest obtaining γ by
 291 leave-one-out cross-validation (Corona et al. 2008; Stone 1974; Wahba 1990).

292 For each value of γ and each area i , define $\boldsymbol{\delta}^{(-i)}(\gamma)$ as the solution of the
 293 corresponding optimization problem with the loss-function term for area i
 294 dropped.¹ The smoothness penalty and any applicable benchmarking con-
 295 straints are calculated over the *whole* of the vector $\boldsymbol{\delta}$, not just the non- i
 296 entries. (This ensures that $\boldsymbol{\delta}^{(-i)}(\gamma)$ does meet all the constraints, while still
 297 making a *prediction* about θ_i .)

The cross-validation score of γ is

$$V(\gamma) = \frac{1}{m} \sum_{i=1}^m \left[\delta_i^{(-i)}(\gamma) - \hat{\boldsymbol{\theta}}_i^B \right]^2 \phi_i,$$

298 where $\delta_i^{(-i)}(\gamma)$ denotes the i th component of $\boldsymbol{\delta}^{(-i)}(\gamma)$, and the minimizer of
 299 the cross-validation scores is $\hat{\gamma} = \operatorname{argmin}_{\gamma \geq 0} V(\gamma)$. Direct evaluation of $V(\gamma)$
 300 can be computationally costly. See Wahba (1990) for faster approximations,
 301 such as “generalized cross-validation.”

302 5 Mean Squared Error Estimation

303 It is traditional in SAE to report approximations to the overall prediction
 304 error. This is generally a challenging undertaking, since methods like cross-
 305 validation can be used to evaluate *prediction* error in a way which is compa-
 306 rable across models, but they do not work for *estimation* error. Thus, one
 307 needs to use more strictly model-based approaches, either analytic or based
 308 on the bootstrap.

¹Instead of the sum of squared errors $\sum_{i'=1}^m \phi_{i'}(\delta_{i'} - \theta_{i'})^2$, we use $\sum_{i' \neq i} \phi_{i'}(\delta_{i'} - \theta_{i'})^2$. This amounts to replacing Φ with a matrix whose i th row and column are both 0.

309 Evaluating the MSE of our estimates is especially difficult, since we
 310 combine a model-based estimate with a non-parametric smoothing term.
 311 A straightforward model-based bootstrap would sample from the posterior
 312 distribution of equation 3.4 to generate a new set of true estimates θ^* and
 313 observations \mathbf{y}^* , re-run the estimation on \mathbf{y}^* , and see how close the result-
 314 ing estimates δ^* came to θ^* . However, this presumes the correctness of the
 315 Fay-Herriot model in equation 3.4, which is precisely what we have chosen
 316 *not* to assume through our imposition of the benchmarking/smoothing con-
 317 straints². Note that such constraints do not fit naturally into the generative
 318 model.

319 We evade this dilemma by using a non-parametric bootstrap, a common
 320 approach when the functional form of a regression is known fairly securely
 321 but the distribution of the error terms is not. The bootstrap assumes that
 322 smoothing is appropriate and that we have chosen the right Ω matrix. We
 323 discuss the choice of Ω and the smoothing penalty in Section 6.3 for the
 324 application on the Empirica database. We assume that the loss function
 325 weight for the i th area (ϕ_i) is the inverse of the estimated MSE.

326 5.1 Non-parametric Bootstrap

327 In this section, we describe the use of a non-parametric bootstrap in order
 328 to estimate the MSE. Assume the area-level Fay-Herriot model in equation
 329 3.4; however, the assumed model here can be a parametric, non-parametric,
 330 or semi-parametric. Assume an estimator $\hat{\theta}$ of θ . For example, one could
 331 consider $\hat{\theta}^{EB}$ in equation 3.7, such as the Bayes estimator or the empirical
 332 Bayes estimator. In addition, assume a constrained Bayes estimator $\hat{\theta}^{BM}$,
 333 such as a benchmarked estimator or a smoothed benchmarked estimator.

For convenience, we can re-write the Fay-Herriot model in equation 3.4
 as the following:

$$\mathbf{y} = X\boldsymbol{\beta} + \mathbf{u} + \mathbf{e}. \quad (5.1)$$

Now, define the residuals as $\mathbf{r} = \mathbf{y} - \hat{\boldsymbol{\theta}}$ as in [Carpenter et al. \(2003\)](#). There
 are other non-parametric bootstraps that have been utilized in the small
 area literature that can be found in the work of [Opsomer et al. \(2008\)](#) and
[Molina et al. \(2009\)](#). Next, center and scale the residuals \mathbf{r} and the estimated
 random effects $\hat{\mathbf{u}}$, where we denote these by \mathbf{r}_e^c and \mathbf{r}_u^c , respectively. These
 are centered at $\mathbf{0}$ and scaled to ensure that they have empirical covariances

²If we follow this procedure nonetheless, we always conclude that benchmarking and
 especially smoothing radically increase the MSE by introducing large biases.

equal to D and $\hat{\sigma}_u I_m$, respectively. Next, we bootstrap and re-sample the centered and scaled random effects \mathbf{r}_u^c , and the residuals \mathbf{r}_e^c , in the following way:

$$\mathbf{u}^* \stackrel{\text{iid}}{\sim} \{\mathbf{r}_u^c\} \quad (5.2)$$

$$\mathbf{e}^* \stackrel{\text{iid}}{\sim} \{\mathbf{r}_e^c\} \quad (5.3)$$

$$\mathbf{y}^* = \mathbf{X}\tilde{\boldsymbol{\beta}}(\hat{\sigma}_u^2) + \mathbf{u}^* + \mathbf{e}^*. \quad (5.4)$$

334 All the residuals are re-sampled independently.

335 In equations 5.2 and 5.3, we draw (with replacement) independent and
 336 identically distributed (iid) random variables u_1^*, \dots, u_m^* and e_1^*, \dots, e_m^* where
 337 each u_i^* is equal to each of r_1^c, \dots, r_m^c and each e_i^* is equal to each of r_1^c, \dots, r_m^c
 338 with probability $1/m$ respectively.

339 Re-sampling-based bootstraps are commonly used in assessing uncer-
 340 tainty for regression models. They presume the correctness of the functional
 341 form of the regression, but not of distributional assumptions about the noise.
 342 To summarize, the resampling procedure proceeds in the following way:

- 343 1. From data (\mathbf{x}, \mathbf{y}) , obtain estimates $\hat{\boldsymbol{\theta}}$ and centered and scaled residuals
 344 \mathbf{r}_u^c and \mathbf{r}_e^c .
- 345 2. Repeat B times:
 - 346 (a) Draw \mathbf{u}^* and \mathbf{e}^* by resampling with replacement from \mathbf{r}_u^c and \mathbf{r}_e^c
 347 respectively.
 - 348 (b) Set $\mathbf{y}^* = \mathbf{X}\tilde{\boldsymbol{\beta}}(\hat{\sigma}_u^2) + \mathbf{u}^* + \mathbf{e}^*$.
 - 349 (c) Refit the model on $(\mathbf{x}, \mathbf{y}^*)$ to obtain $\hat{\boldsymbol{\theta}}^*$.
 - 350 (d) Calculate the constrained Bayes estimate $\hat{\boldsymbol{\theta}}^{BM*}$.
- 351 3. Use the distribution of $\hat{\boldsymbol{\theta}}^*$ and $\hat{\boldsymbol{\theta}}^{BM*}$ in bootstrap calculations to obtain
 352 the estimated MSE.

353 Thus, we have proposed a non-parametric bootstrap, where we define
 354 this using the unconstrained estimates. This is important to ensure that
 355 the bootstrap produces synthetic data closer to the observed data. This can
 356 of course be checked for a few situations under the Fay-Herriot model, and
 357 we do so in our application given previous work done by Prasad and Rao
 358 (1990).

359 **6 Application: Estimating Rental Prices in Berlin**

360 In this section, our goal is to estimate the average rent in each of the 447
361 low geographical areas called *Lebensweltlich orientierte Räume* (LORs) in
362 Berlin in 2015. The Berlin Senate Department for Urban Development and
363 Housing is officially responsible for providing official comparative rents for
364 the consumer price index in Berlin. This official data set is comprised of
365 roughly 2,000 apartments. Furthermore, it is collected by the Federal Sta-
366 tistical Office in Berlin-Brandenburg using a stratified sampling design. The
367 strata are defined by type of the apartment, districts, and type of the land-
368 lord.³ Unfortunately, this survey is confidential due to privacy constraints,
369 and it is not possible to access this database.

370 In order to mimic the official data set collected by the Federal Statistical
371 Office in Berlin-Brandenburg, Empirica provided us with a similar data set
372 for 2,000 apartments available for rent in Berlin in 2015. The Empirica data
373 set was obtained via web-scraping and print media. The Empirica data set
374 creates the following two new challenges: (1) reliable estimates at the LOR
375 level are not available due to the very small or zero sample sizes in some
376 LORs and (2) the sample from the Empirica data set may fail to capture
377 important parts of the Berlin rental market. Therefore, we combine direct
378 estimates from the sample of the Empirica data set with small area models
379 that use area level predictors. The small area model is described in detail
380 in Section 6.2. To match the official rent per square meter in Berlin we
381 incorporate a benchmarking constraint (i.e. the fixed amount of rent set by
382 law for the city of Berlin) and investigate the effect of smoothing in Section
383 6.3. The spatial distribution of rent prices in Berlin is also discussed. Before
384 proceeding, we first further describe the Empirica data set in Section 6.1.

385 **6.1 The Empirica Data Set**

386 The Empirica data set is chosen (in order to mimic the data collected by the
387 Federal Statistical Office in Berlin-Brandenburg) according to a stratified
388 sampling design (strata are defined using the region and the size of the
389 apartment) with a sample size of around 2,000. There are a total of 100
390 covariates such as rental price per square meter (excluding costs for water,
391 sewage, trash collection, etc.), number of bedrooms and bathrooms, year
392 of construction, balcony, and the address (including longitude and latitude
393 coordinates).

³We refer to https://www.statistik-berlin-brandenburg.de/publikationen/aufsaetze/2016/HZ_201602-04.pdf for further information.

394 The 48 strata are defined by the cross-classification of the 12 districts
 395 with the living space categories (four categories: $< 40m^2$, $40 - 60m^2$, $60 -$
 396 $90m^2$, $> 90m^2$). This leads to a sample size of 2,083 apartments with 302
 397 in-sample LORs and 142 out-of-sample LORs. The summary statistics of
 398 the sample sizes by LORs are presented in Table 1.

Table 1: Summary statistics over LOR (Empirica Database)

	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
Sample size	3	4	5	6.90	9	24

399 Direct estimation—using only the sample data—of the average rental
 400 price per square meter is not an option because direct estimates are only
 401 available for 302 out of the 447 LORs. In addition, for some areas where
 402 sample data is available, the small samples sizes lead to direct estimates
 403 with a low precision. As such, we attempt to improve the precision of small
 404 area direct estimates by combining the direct estimates from the sample of
 405 the database with small area models. The small area model from which
 406 we derive our initial estimates is described in Section 6.2. It provides an
 407 estimate of the average rental price per square meter at LOR level in Berlin
 408 based on the Empirica database.

409 In addition, because rentals often directly change hands from outgoing
 410 tenants to incoming tenants in Berlin, a market is not covered by online and
 411 print media sources. More specifically, in this secondary market, it is likely
 412 that the rental price remains constant. For this reason, we may expect some
 413 overestimation of the average rental price per square meter by using the
 414 Empirica database. To adjust for the potential lack of representativeness
 415 of the sample data, we incorporate a benchmarking constraint such that
 416 the weighted mean of the average LOR rental price estimates matches the
 417 official rent per square meter of €8.02 in Berlin as published by the Berlin
 418 Senate Department for Urban Development and Housing. In addition to
 419 benchmarking, we add a spatial smoothness constraint across LORs since we
 420 may expect rental prices to be spatially related. Our choice of the Laplacian
 421 and smoothing penalty is discussed in Section 6.3.

422 6.2 The Fay-Herriot Model applied to Empirica

423 In this section, we describe our methodology that we use for analysis and
 424 estimation. In the context of the Empirica data set, θ denotes the true rental
 425 price per square meter for all the LORs, \mathbf{y} denotes the direct estimates

426 based on the survey from the Empirica data set, $D = \text{Diag}(D_1, \dots, D_m)$
 427 denotes the sampling covariance matrix of the direct estimates \mathbf{y} , σ_u^2 denotes
 428 the area-level variance parameter, \mathbf{x}_i denotes the vector of covariates, and
 429 $\boldsymbol{\beta}$ denotes the unknown regression coefficients. Our initial unconstrained
 430 estimates for the average rental price per square meter are derived from the
 431 area-level Fay-Herriot model in equation 3.4 (Fay and Herriot 1979).

432 The final model is selected following the ideas by Marhuenda et al.
 433 (2014). In particular, we used a Kullback symmetric divergence criterion
 434 with a bootstrap adjustment, KICb2. The final model includes seven ag-
 435 gregated (LOR-level) predictors obtained from the Empirica data set: 1)
 436 the average year of construction, 2) the average floor of the apartment in
 437 the building, and 3-7) share of apartments with an energy performance cer-
 438 tificate available (EPC)/ balcony/ elevator/ fitted kitchen/ open fireplace,
 439 respectively. The distribution of the predictors over LORs is presented in
 440 Table 2. The inclusion of these covariates led to an $R^2 = 52\%$ for the linear
 441 model at the aggregated level.

Table 2: Summary statistics over LOR

	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
Year of construction	1909	1936	1955	1956	1973	2010
Floor	0.44	1.77	2.27	2.50	2.83	7.64
Share of EPC	0.00	0.55	0.66	0.64	0.75	1.00
Share of balcony	0.36	0.62	0.72	0.72	0.83	1.00
Share of elevator	0.00	0.16	0.31	0.36	0.54	0.98
Share of fitted kitchen	0.000	0.33	0.47	0.47	0.60	0.94
Share of open fireplace	0.00	0.00	0.01	0.02	0.02	0.28

442 The Fay-Herriot model is estimated by using the *emdi* package in R
 443 (Kreutzmann et al. 2019). Figure 1 presents the average rental price per
 444 square meter based on the Fay-Herriot estimator. We observe that the most
 445 expensive parts of Berlin are around the city centre and the area in the
 446 south-west (Zehlendorf and Grunewald) of Berlin, which is consistent with
 447 official results by the Berlin Senate.

448 6.3 Benchmarking/Smoothing applied to Empirica

449 Figure 1 offers a first picture about the rent per square meter at the LOR-
 450 level in Berlin using a sample from the Empirica data set. As already men-
 451 tioned, the Empirica data set may exclude certain parts of the rental market.

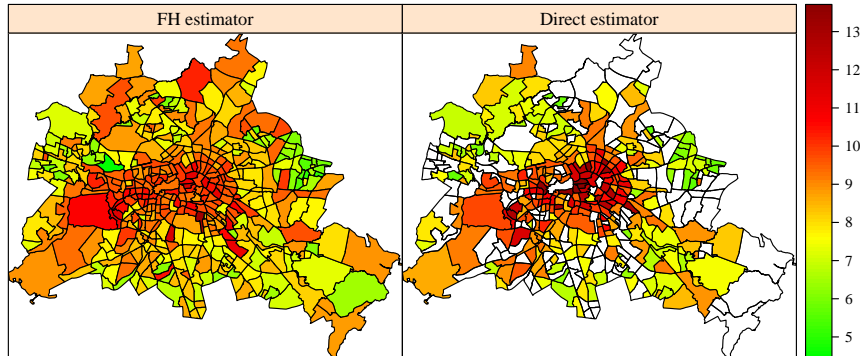


Figure 1: Average rent per square meter in € based on the unconstrained Fay-Herriot estimator (left map) and the direct estimator (right map).

452 To correct for this, we incorporate a benchmarking constraint requiring that
 453 the weighted mean of the average rental price estimates matches the official
 454 rent per square meter of €8.02 in Berlin and in addition consider smoothing
 455 the estimates over space.

456 In particular, we consider the following two options for benchmarking
 457 and/or smoothing: (i) benchmarking the mean without spatial smoothing,
 458 and (ii) benchmarking the mean with spatial smoothing. We expect that
 459 smoothing will reduce the variability in the resulting benchmarked estimates.
 460 In addition, in each case (i) and (ii), we choose the benchmarked weights
 461 w_i to be proportional to the number of apartments in each LOR. There is a
 462 rich literature on the choice of such weights. We refer to Ghosh and Steorts
 463 (2013) and Datta et al. (2011) for further details.

464 Figure 2 compares the unconstrained Fay-Herriot estimator to the bench-
 465 marked Fay-Herriot estimator. As expected, the Fay-Herriot estimates of the
 466 average rental price per square meter are higher than the benchmarked esti-
 467 mates for each of the 447 LOR. Observe that the benchmarked Fay-Herriot
 468 estimates are on average around €0.348 lower than the unconstrained Fay-
 469 Herriot estimates. This is expected as the Empirica data set excludes the
 470 secondary rental market, which means that rental prices in the Empirica
 471 data set tend to be lower than those advertised online or in print media.
 472 Intuitively, properties that are advertised in the open market may ignore
 473 the law on rental prices, altogether.

474 Turning now to the use of spatial smoothing, the most important part
 475 of the smoothing procedure is selecting the matrices Ω and Q . Recall that

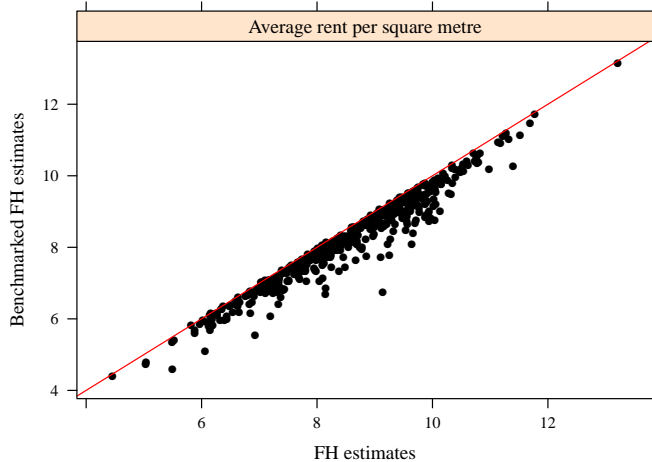


Figure 2: Average rent per square meter in € using (a) the benchmarked Fay-Herriot estimates and (b) the unconstrained Fay-Herriot estimates.

476 Ω is used to measure the smoothness of estimates; and Q shows how similar
 477 the estimates for any two domains should be. This is inevitably application-
 478 specific. In our application, we utilize a simple choice, where where $q_{ii'} = 1$
 479 if the LORs i and i' shared a border, and 0 otherwise. This treats the
 480 LOR as nodes in an unweighted graph, with Q being its adjacency matrix
 481 and Ω its Laplacian. In addition, we considered several alternative ways
 482 of smoothing the Fay-Herriot estimates. One can choose $q_{ii'}$ such that it
 483 decreases with the geographic distance between LORs, regarding the points
 484 at their respective centers. A second approach was to treat the 12 districts
 485 in Berlin as clusters, setting $q_{ii'} = 1$ for LORs within a cluster and $q_{ii'} = 0$
 486 for LORs between them, but neither of these two approaches worked well
 487 under cross-validation. Note that choosing the spatial smoothing parameter
 488 is not an issue in our application as we do not encounter spatial islands,
 489 however, if one does encounter such issues, we would recommend modifying
 490 the definition of a neighbor to be the minimum geographic distance criterion.

491 As described in Section 4.3, the smoothing factor γ was picked by leave-
 492 one-out cross-validation and the final value was $\gamma \approx 0.146$. Figure 3 shows
 493 the smoothed and benchmarked Fay-Herriot estimates versus the uncon-
 494 strained Fay-Herriot estimates. In general, the effect of spatial smoothing
 495 causes an upward adjustment of low values of the unconstrained Fay-Herriot
 496 estimates, and causes a downwards adjustment of higher unconstrained Fay-

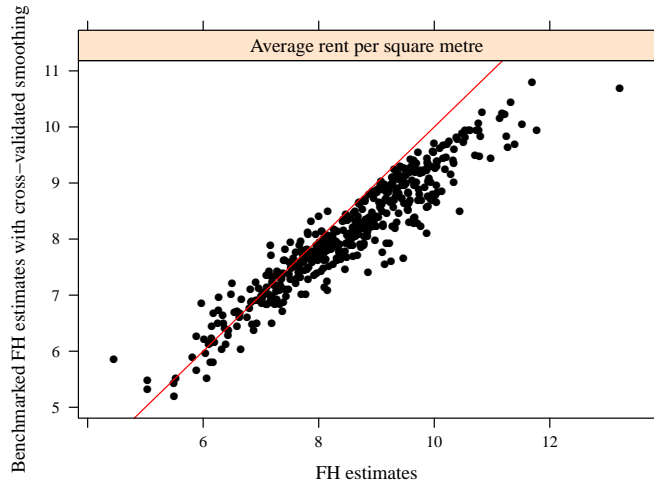


Figure 3: Average rent per square meter in € using (a) the benchmarked Fay-Herriot estimates with cross-validated smoothing and (b) the unconstrained Fay-Herriot estimates.

497 Herriot estimates. However, the majority of the unconstrained Fay-Herriot
 498 estimates are pulled down as a result of both smoothing and benchmarking.
 499 This observation is further confirmed by Figure 4 where the effect of com-
 500 bining smoothing with benchmarking is illustrated this time for the twelve
 501 districts in Berlin. Observe that in each of the twelve districts the smooth
 502 benchmarked Fay-Herriot estimates fall on the line with a slope of less than
 503 1.

504 Table 3 reports the MSEs under the non-parametric bootstrap of Sec-
 505 tion 5 for different combinations of benchmarking and smoothing. In partic-
 506 ular, FH denotes the unconstrained Fay-Herriot estimates, $FH Bench$ the
 507 benchmarked estimates and $FH Bench/Smooth$ the corresponding bench-
 508 marked estimates with cross-validated smoothing. The results are based on
 509 $B = 1000$ bootstrap replications. In addition, we ran a Fay-Herriot model
 510 with spatially correlated random effects, $FH SAR$, using the same adjacency
 511 matrix Q used for the $FH Bench/Smooth$. We followed Pratesi and Salvati
 512 (2009) and used a SAR specifications for the random effects. The MSE of
 513 the $FH SAR$ is estimated by a non-parametric bootstrap ($B = 1000$) as pro-
 514 posed by Molina et al. (2009). Please note that in the case of MSE estimation
 515 under the benchmarked approach with spatially correlated random effects,
 516 $FH SAR Bench$, benchmarking is being implemented with each bootstrap

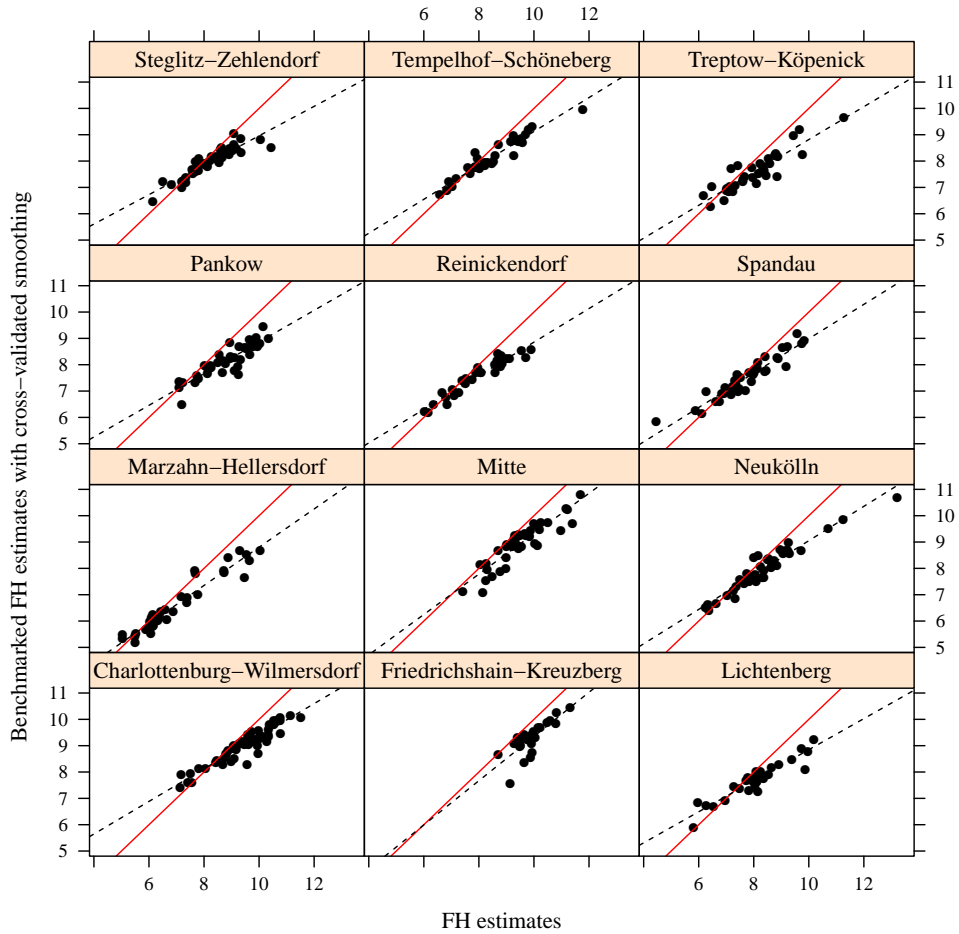


Figure 4: Benchmarked and Spatially Smoothed Fay-Herriot estimates versus unconstrained Fay-Herriot estimates, by region. Large unconstrained Fay-Herriot estimates are adjusted downwards by the benchmarked and spatially smoothed Fay-Herriot estimators, while small unconstrained Fay-Herriot estimators are adjusted upwards. This effect can be seen by the dotted line, denoting the regression line, and the red line, denoting the intersection line.

517 sample. First, we observe that the unconstrained Fay-Herriot estimates (*FH*
518 and *FH SAR*) have a smaller MSE compared to the benchmarked estimates.
519 This is expected as the constraint in the estimation introduces additional
520 variability. However, as benchmarking is required in the application, we fo-
521 cus on the three constrained estimates (*FH Bench*, *FH SAR Bench* and *FH*
522 *Bench/Smooth*). Incorporating the spatial effect via smoothing or a SAR
523 structure reduces the variability for most LORs compared to the bench-
524 marked estimates (*Bench*). In addition, the estimated MSEs under the *FH*
525 *SAR Bench* approach and the *FH Bench/Smooth* are comparable. In par-
526 ticular, on average both methods provide similar estimated MSEs with the
527 Fay-Herriot benchmarked and smooth approach also offering less extreme
528 estimated MSEs.

Table 3: Summary statistics of RMSE estimates over LOR

	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
FH	0.07	0.60	0.67	0.63	0.71	0.89
FH SAR	0.10	0.46	0.64	0.61	0.74	2.17
FH SAR Bench	0.10	0.46	0.67	0.71	0.86	2.91
FH Bench/Smooth	0.07	0.67	0.74	0.75	0.81	1.86
FH Bench	0.07	0.68	0.75	0.77	0.85	1.80

529 Having assessed the variability of the three constrained estimates, we
530 have a closer look to the point estimates of the actual rent per square meter
531 in Berlin. Figure 5 presents the benchmarked estimates with and without
532 cross-validated smoothing and the benchmarked Fay-Herriot with spatially
533 correlated random effects. Overall, all maps reflect the current situation
534 of the rental market for apartments in Berlin with higher rents in the city
535 center and in the district Steglitz-Zehlendorf (in the south-west), whereas
536 the districts of Spandau (in the west) and Marzahn-Hellersdorf (in the east)
537 have lower rents compared to other parts in Berlin. For instance, vast hous-
538 ing estates (plattenbau style - large panel system building) were built in
539 the 1980s in Marzahn-Hellersdorf. Most of the plattenbau apartments were
540 built in large settlements on the edge of Berlin making them inconveniently
541 located leading to high vacancy rates and low rent prices. In contrast, the
542 district of Steglitz-Zehlendorf consists of very affluent localities like Dahlem
543 or Zehlendorf. The localities Nikolasee and Wannsee of Steglitz-Zehlendorf
544 are located around the forest of Grunewald and two lakes (Greater and Little
545 Wannsee). These localities are some of the most expensive areas in Berlin
546 for housing.

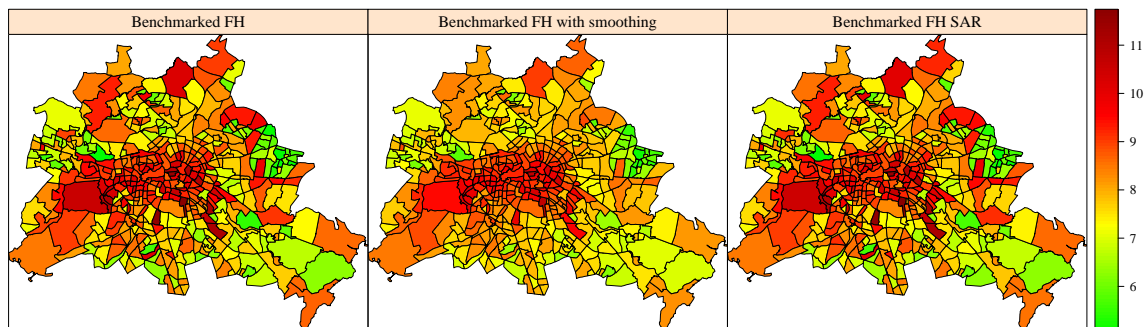


Figure 5: Average rent per square meter in € based on the benchmarked Fay-Herriot estimator with (middle) and without (left map) cross-validated smoothing and the benchmarked Fay-Herriot with spatially correlated random effects (right map).

547 However, we also observe some differences between the three maps. First,
548 some LORs in the city center around the governmental quarter (with the Re-
549 ichstag building, German Chancellery, and Bellevue Palace) have lower than
550 expected rental prices based on the benchmarked estimates (*FH Bench* and
551 *FH SAR Bench*) in Figure 5 and on the unconstrained Fay-Herriot estimates
552 in Figure 1. Additional discrepancies occur for LORs in the suburbs in the
553 north (for instance the locality Blankenfelde in the district Pankow) and
554 south-east (for instance the locality Karolinenhof in the district Treptow-
555 Köpenick). In the latter case LORs have higher than expected rental prices
556 based on the benchmarked estimates (*FH Bench* and *FH SAR Bench*) in
557 Figure 5 and on the unconstrained Fay-Herriot estimates in Figure 1. These
558 localities are bordering the federal state of Brandenburg, are located in ru-
559 ral parts of Berlin and they are the least densely populated areas in Berlin.
560 In addition, Figure 1 reveals that sample information is missing from these
561 LORs, and thus, we heavily rely on the Fay-Herriot equation 3.4. It is pos-
562 sible that especially for LORs with somehow different characteristics (in-
563 frastructure and environment) our estimates based on equation 3.4 suffer
564 from some misspecification. Nevertheless, it appears that our benchmarked
565 estimates with cross-validated smoothing are able to adjust the estimates
566 and protect against potential model misspecification.

567 7 Discussion

568 We have provided a general approach to area-level SAE, where we smooth
569 and benchmark model-based estimates. Our approach yields closed-form
570 solutions without requiring any distributional assumptions. Furthermore,
571 our results apply for linear and non-linear estimators. Finally, we show in
572 the application that smoothing has the potential to improve estimation of
573 rental prices on LOR level in Berlin for most LORs.

574 We now outline some possible extensions, namely extensions to weighted
575 variability constraints, moving beyond squared error loss, and moving from
576 point estimation to full posterior estimates for maximal flexibility. As men-
577 tioned earlier, working beyond a weighted mean constraint and with both a
578 weighted mean and weighted variability would be a more general benchmark-
579 ing framework. The question of how to incorporate variability constraints
580 while maintaining tractability of the model is a potential direction of future
581 research and is beyond the scope of this paper, as the problem may not al-
582 ways be a convex optimization problem. In addition, throughout our paper,
583 we have worked with the squared error loss function. However, it should
584 be possible to replace this with any other loss function. Once the Bayes
585 estimate is obtained, the constrained Bayes estimate would be found by a
586 projection onto the corresponding feasible set.

587 This would presumably mean a need for using numerical optimization
588 when the optimization problem is not tractable. Finally, it may be possible
589 to go beyond point estimates to distributional estimates. Given a sample
590 from the posterior distribution (e.g., from MCMC), it is possible to project
591 each sample point into the feasible set, giving rise to a posterior distribu-
592 tion whose support respects the constraints. This idea is related to that
593 of [Dunson and Neelon \(2003\)](#), however, cannot be directly adapted to our
594 setting. [Dunson and Neelon \(2003\)](#) have proposed constrained Bayes esti-
595 mation through a posterior projection approach, which is appealing in the
596 sense that one fully achieves a Bayesian posterior distribution to the con-
597 strained optimization problem. The constraints considered by the authors
598 are ordered parameters, and do not easily generalize to both weighted means
599 and weighted variabilities in our general framework.

600 Acknowledgements

601 Schmid and Tzavidis are supported by ES/N011619/1 — *Innovations in*
602 *Small Area Estimation Methodologies* from the UK Economic and Social

603 Research Council. Tzavidis is also supported by the InGRID 2 EU-Horizon
604 2020 infrastructure grant (<http://www.inclusivegrowth.eu>). The au-
605 thors thank Empirica-Systeme GmbH (www.empirica-systeme.de) for pro-
606 viding the data set used in the application. The ideas of this paper are
607 of the authors and not of the funding organizations or the data providers.
608 Finally, the authors thank the Editor and the reviewers for comments that
609 significantly improved the paper. The authors also thank David Banks for
610 providing minor comments regarding manuscript.

611 References

- 612 Battese, G., Harter, R., and Fuller, W. (1988), “An Error-Components
613 Model for Prediction of County Crop Area Using Survey and Satellite
614 Data,” *Journal of the American Statistical Association*, 83, 28–36.
- 615 Belkin, M., Niyogi, P., and Sindhvani, V. (2006), “Manifold Regulariza-
616 tion: A Geometric Framework for Learning from Labeled and Unlabeled
617 Examples,” *Journal of Machine Learning Research*, 7, 239–2434.
- 618 Bell, W., Datta, G., and Ghosh, M. (2013), “Benchmarked Small Area Es-
619 timators,” *Biometrika*, 100, 189–202.
- 620 Carpenter, J. R., Goldstein, H., and Rasbash, J. (2003), “A novel bootstrap
621 procedure for assessing the relationship between class size and achieve-
622 ment,” *Journal of the Royal Statistical Society: Series C (Applied Statis-
623 tics)*, 52, 431–443.
- 624 Corona, E., Lane, T., Storlie, C., and Neil, J. (2008), “Using Laplacian
625 Methods, RKHS Smoothing Splines and Bayesian Estimation as a frame-
626 work for Regression on Graph and Graph Related Domains,” Tech. Rep.
627 TR-CS-2008-06, Department of Computer Science, University of New
628 Mexico.
- 629 Datta, G. and Ghosh, M. (1991), “Bayesian Prediction in Linear Mod-
630 els: Applications to Small Area Estimation,” *The Annal of Statistics*,
631 19, 1748–1770.
- 632 Datta, G. S., Ghosh, M., Steorts, R., and Maples, J. (2011), “Bayesian
633 benchmarking with applications to small area estimation,” *TEST*, 20,
634 574–588.
- 635 Dunson, D. B. and Neelon, B. (2003), “Bayesian inference on order-
636 constrained parameters in generalized linear models,” *Biometrics*, 59,
637 286–295.
- 638 Fay, R. and Herriot, R. (1979), “Estimates of income from small places:
639 an application of James-Stein procedures to census data,” *Journal of the
640 American Stastical Association*, 74, 269–277.
- 641 Ghosh, M. (1992), “Constrained Bayes estimation with applications,” *Jour-
642 nal of the American Statistical Association*, 87, 533–540.

- 643 Ghosh, M., Rao, J., et al. (1994), “Small area estimation: an appraisal,”
644 *Statistical science*, 9, 55–76.
- 645 Ghosh, M. and Steorts, R. C. (2013), “Two-Stage Bayesian Benchmarking
646 as Applied to Small Area Estimation,” *TEST*, 22, 670–687.
- 647 Isaki, C., Tsay, J., and Fuller, W. (2004), “Weighting sample data subject
648 to independent controls,” *Survey Methodology*, 20, 35–44.
- 649 Kafadar, K. (1996), “Smoothing Geographical Data, Particularly Rates of
650 Disease,” *Statistics in Medicine*, 15, 2539–2560.
- 651 Kreutzmann, A.-K., Pannier, S., Rojas-Perilla, N., Schmid, T., Templ, M.,
652 and Tzavidis, N. (2019), “emdi: estimating and mapping disaggregated
653 Indicators,” *Journal of Statistical Software*, 91, 1–33.
- 654 Lee, A. B. and Wasserman, L. (2010), “Spectral Connectivity Analysis,”
655 *Journal of the American Statistical Association*, 105, 1241–1255.
- 656 Louis, T. (1984), “Estimating a population of parameter values using Bayes
657 and empirical Bayes methods,” *Journal of the American Stastical Associ-
658 ation*, 79.
- 659 Marhuenda, Y., Morales, D., and del Carmen Pardo, M. (2014), “Informa-
660 tion criteria for Fay–Herriot model selection,” *Computational Statistics
661 and Data Analysis*, 70, 268 – 280.
- 662 Molina, I. and Rao, J. (2010), “Small area estimation of poverty indicators,”
663 *Canadian Journal of Statistics*, 38, 369–385.
- 664 Molina, I., Salvati, N., and Pratesi, M. (2009), “Bootstrap for estimating
665 the MSE of the Spatial EBLUP,” *Computational Statistics*, 24, 441–458.
- 666 Newman, M. E. J. (2010), *Networks: An Introduction*, Oxford, England:
667 Oxford University Press.
- 668 Opsomer, J. D., Claeskens, G., Ranalli, M. G., Kauermann, G., and Breidt,
669 F. (2008), “Non-parametric small area estimation using penalized spline
670 regression,” *Journal of the Royal Statistical Society: Series B (Statistical
671 Methodology)*, 70, 265–286.
- 672 Pfeiffermann, D. (2013), “New important developments in small area esti-
673 mation,” *Statistical Science*, 28, 40–68.

- 674 Prasad, N. and Rao, J. (1990), “The estimation of the mean squared error
675 of small-area estimators,” *Journal of the American Statistical Association*,
676 85, 163–171.
- 677 Pratesi, M. and Salvati, N. (2008), “Small area estimation: the EBLUP
678 estimator based on spatially correlated random area effects,” *Statistical
679 methods and applications*, 17, 113–141.
- 680 — (2009), “Small Area Estimation in the Presence of Correlated Random
681 Area Effects,” *Journal of Official Statistics*, 25 (1), 37–53.
- 682 Rao, J. and Molina, I. (2015), *Small Area Estimation*, John Wiley & Sons,
683 New York.
- 684 Rao, J. N. K. and Yu, M. (1994), “Small-area estimation by combining
685 time-series and cross-sectional data,” *Canadian Journal of Statistics*, 22,
686 511–528.
- 687 Souza, D. F., Moura, F., and Migon, H. (2009), “Small area population
688 prediction via hierarchical models,” *Catalogue no. 12-001-X*, 203.
- 689 Steorts, R. C. and Ghosh, M. (2013), “On estimation of mean squared errors
690 of benchmarked empirical Bayes estimators,” *Statistica Sinica*, 749–767.
- 691 Stone, M. (1974), “Cross-validatory choice and assessment of statistical pre-
692 dictions,” *Journal of the Royal Statistical Society B*, 36, 111–147.
- 693 Tzavidis, N., Luna, A., Zhang, L. C., Schmid, T., and Rojas-Perilla, N.
694 (2018), “From start to finish: A framework for the production of small
695 area official statistics,” *Journal of the Royal Statistical Society: Series A*,
696 181, 927–979.
- 697 Wahba, G. (1990), *Spline Models for Observational Data*, Philadelphia: So-
698 ciety for Industrial and Applied Mathematics.
- 699 Wehbe, L., Ramdas, A., Steorts, R. C., and Shalizi, C. R. (2015), “Regu-
700 larized Brain Reading with Shrinkage and Smoothing,” *Annals of Applied
701 Statistics*, 9, 1997–2022.

702 **Supplementary Material**

703 **A Lemma on Squared Differences**

704 **Lemma A.1.** *For a suitable matrix Ω ,*

$$\sum_{i,i'} (\delta_i - \delta_{i'})^2 q_{ii'} = \boldsymbol{\delta}^T \Omega \boldsymbol{\delta}.$$

705 *Proof.* Begin by expanding the square and collecting terms:

$$\begin{aligned} & \sum_{i,i'} (\delta_i - \delta_{i'})^2 q_{ii'} \\ &= \sum_{i,i'} \delta_i^2 q_{ii'} + \sum_{i,i'} \delta_{i'}^2 q_{ii'} - 2 \sum_{i,i'} \delta_i \delta_{i'} q_{ii'} \\ &= \sum_i \delta_i^2 \sum_{i'} q_{ii'} + \sum_{i'} \delta_{i'}^2 \sum_i q_{ii'} - 2 \sum_{i,i'} \delta_i \delta_{i'} q_{ii'}. \end{aligned}$$

706 Now define the diagonal matrix $Q^{(r)}$ with elements $q_{ii}^{(r)} = \sum_{i'} q_{ii'}$, and define
707 the diagonal matrix $Q^{(c)}$ with elements $q_{jj}^{(c)} = \sum_i q_{ij}$. Substituting,

$$\begin{aligned} \sum_{i,i'} (\delta_i - \delta_{i'})^2 q_{i,i'} &= \boldsymbol{\delta}^T Q^{(r)} \boldsymbol{\delta} + \boldsymbol{\delta}^T Q^{(c)} \boldsymbol{\delta} - 2 \boldsymbol{\delta}^T Q \boldsymbol{\delta} \\ &= \boldsymbol{\delta}^T \left(Q^{(r)} + Q^{(c)} - 2Q \right) \boldsymbol{\delta}, \end{aligned}$$

708 which defines Ω . □

709 **Remark A.1.** *In an unweighted, undirected graph with adjacency matrix A ,*
710 *the degree matrix D is defined by $D_{ii} = \sum_j A_{ij}$, $D_{ij} = 0$; the graph Laplacian*
711 *in turn is $L = D - A$ (Newman 2010). If Q is an adjacency matrix, then*
712 *$Q^{(r)} = Q^{(c)} = D$, and $\Omega = 2L$.*

713 **Remark A.2.** *By construction, Ω is clearly positive semi-definite. It is not*
714 *positive definite, because $(1 \ 1 \ \dots \ 1)$ is always an eigenvector, of eigenvalue*
715 *zero. This corresponds to the fact that adding the same constant to each δ_i*
716 *does not change $\sum_{i,i'} (\delta_i - \delta_{i'})^2 q_{i,i'}$. (These are of course basic properties of*
717 *graph Laplacians.)*