

## EMPIRICAL BAYES ESTIMATES FOR A 2-WAY CROSS-CLASSIFIED ADDITIVE MODEL\*

BY LAWRENCE D. BROWN

*University of Pennsylvania*

BY GOURAB MUKHERJEE

*University of Southern California*

AND

BY ASAF WEINSTEIN

*Stanford University*

We develop an empirical Bayes procedure for estimating the cell means in an unbalanced, two-way additive model with fixed effects. We employ a hierarchical model, which reflects exchangeability of the effects within treatment and within block but not necessarily between them, as suggested before by [Lindley and Smith \(1972\)](#). The hyperparameters of this hierarchical model, instead of considered fixed, are to be substituted with data-dependent values in such a way that the point risk of the empirical Bayes estimator is small. Our method chooses the hyperparameters by minimizing an unbiased risk estimate and is shown to be asymptotically optimal for the estimation problem defined above. The usual empirical Best Linear Unbiased Predictor (BLUP) is shown to be substantially different from the proposed method in the unbalanced case and therefore performs sub-optimally. Our estimator is implemented through a computationally tractable algorithm that is scalable to work under large designs. The case of missing cell observations is treated as well. We demonstrate the advantages of our method over the BLUP estimator through simulations and in a real data example, where we estimate average Nitrate levels in water sources based on their locations and the time of the day.

**1. Introduction.** Multilevel cross-classified models are pervasive in statistics, with applications ranging from detecting sources of variability in medical research ([Goldstein \*et al.\*, 2002](#)) to understanding micro-macro linkages in social studies ([Mason \*et al.\*, 1983](#); [Zaccarin and Rivellini, 2002](#)). These models offer a natural and flexible approach to specify meaningful

---

\*Supported in part by grants NSF DMS .

*Keywords and phrases:* Shrinkage estimation; Empirical Bayes; Two-way ANOVA; Oracle Optimality; Stein's unbiased risk estimate (SURE); Empirical BLUP.

latent structures and, importantly, a systematic way to use all information for simultaneously analyzing the effects of more than one factor (Rasbash and Goldstein, 1994). Hierarchical cross-classified models have classically been used to decompose the total variability of the response into individual sources and for prediction in random-effects models. Nevertheless, ever since the appearance of the James-Stein estimator (James and Stein, 1961) and its Bayesian interpretation (Stein, 1962; Lindley, 1962), the usefulness of such models in estimation problems involving multiple *nonrandom* effects has been well recognized.

Hierarchical models have been used to facilitate shrinkage estimators in linear regression models since the early 1970s (Efron and Morris, 1972). In both theoretical and more applied work, various authors have employed hierarchical models to produce estimators that shrink *towards* a subspace (e.g., Sclove, 1968; Oman, 1982; Jiang *et al.*, 2011; Tan, 2014) or *within* a subspace (e.g., Lindley and Smith, 1972; Rolph, 1976; Kou and Yang, 2015); see Section 2 of the last reference for a discussion on the difference between the two types of resulting estimators. Cross-classified additive models are in a sense the most immediate extension of Stein’s canonical example. Specifically, unlike in a general linear model, the symmetries of within-batch effects can be regarded as a-priori information, which suggest the use of exchangeable priors, such as those proposed by Lindley and Smith (1972) and Efron and Morris (1973). In the case of balanced design, the properties of resulting shrinkage estimators are by now well understood and have a close relationship to the James-Stein estimator. Indeed, when all cell counts are equal, multiple one-way, homoscedastic estimation problems emerge; for these the James-Stein estimator has optimality properties under many criteria. But in the unbalanced case, the problems of estimating the effects corresponding to different batches are intertwined due to lack of orthogonality in the design matrix; hence, the situation in the case of unbalanced design is substantially different.

This paper deals with empirical Bayes (EB) estimation of the cell means in the two-way fixed effects additive model with unbalanced design. We consider a family of Bayes estimators resulting from a normal hierarchical model, which reflects within-batch exchangeability and is indexed by a set of hyper-parameters that govern the prior. Any corresponding estimator that substitutes *data-dependent* values for the hyper-parameters is referred to as an empirical Bayes estimator. We propose an empirical Bayes procedure that is asymptotically optimal for the estimation of the cell means under squared loss. In our asymptotic analysis, the number of row and column levels tends to infinity. Importantly, the so-called empirical BLUP (Best Linear Unbiased

Predictors) estimators, using the usual maximum-likelihood approach in estimating the hyperparameters, are shown to perform sub-optimally in the unbalanced case. Instead of using the maximum-likelihood criterion, we choose the values for the hyper-parameters by minimizing an unbiased estimate of the risk (URE), which leads to estimates that are different in an essential way. The proposed approach is appealing in the fixed effects case, because it uses a criterion directly related to the risk instead of using the likelihood under the postulated hierarchical model.

Using the URE criterion to calibrate tuning parameters has been proposed in many previous works and in a broad range of parametric and nonparametric estimation problems (Li, 1986; Ghosh *et al.*, 1987; Donoho *et al.*, 1995; Johnstone and Silverman, 2004; Candes *et al.*, 2013, to name a few). Recently, Xie *et al.* (2012) employed URE minimization to construct alternative empirical Bayes estimators to the usual ones in the Gaussian mean problem with known *heteroscedastic* variances and showed that it produces asymptotically uniformly better estimates. Our work can be viewed as a generalization of Xie *et al.* (2012) from the one-way unbalanced layout to the two-way unbalanced layout.

The two-way unbalanced problem presents various new challenges. The basis for the difference, of course, lies in the facts that the two-way case imposes structure on the mean vector, which is nontrivial to handle due to missingness and imbalance in the design. Some of the implications are that the analysis of the performance of EB methods is substantially more involved than in the one-way scenario; in addition, the implementation of the URE estimator, which is trivial in the one-way scenario, becomes a cause of concern, especially with a growing number of factor levels. We offer an implementation of the corresponding URE estimate that in the all-cells-filled case has comparable computational performance to that of the standard empirical BLUP in the popular R package `lme4` of Bates (2010). Our theoretical analysis of the two-way case differs in fundamental aspects from the optimality proof techniques usually used in the one-way normal mean estimation problem. To tackle the difficulties encountered in the two-way problem, where computations involving matrices are generally unavoidable, we developed a flexible approach for proving asymptotic optimality based on efficient pointwise risk estimation; this essentially reduces our task to controlling the moments of Gaussian quadratic forms.

We would also like to point out that the current work is different from the recent extensions of Kou and Yang (2015) of the URE approach to the general Gaussian linear model. While the setup considered in that paper formally includes our setup as a special case, their results have limited

implications for additive cross-classified models; for example, the covariance matrix used in their second level of the hierarchy is not general enough to accommodate the within-batch exchangeable structure we employ and is instead governed by a single hyper-parameter. Moreover, their asymptotic results require keeping the dimension of the linear subspace fixed, whereas the number of factor levels is increasing in our setup.

**Organization of the paper.** In Section 2 we describe our estimation setup for the simplest case when there are no missing observations. In Section 3 we introduce the more general model, which allows missing observations, and describe a unified framework for estimation across all scenarios – missing or non-missing. In Section 4 we show that our proposed estimation methodology is asymptotically optimal and is capable of recovering the directions and magnitude for optimal shrinkage; this is established through the notion of oracle optimality. Section 5 is devoted to the special case of a balanced design. After describing the computation details in Section 6, we report the results from extensive numerical experiments in Section 7. Lastly, in Section 8 we demonstrate the applicability of our proposed method on a real-world problem concerning the estimation of the average nitrate levels in water sources based on location and time of day.

## 2. Model Setup and Estimation Methods.

2.1. *Basic Model and Estimation Setup.* **Additive model with all cells filled.** Consider the following basic two-way cross-classified additive model with fixed effects:

$$(1) \quad \begin{aligned} y_{ij} &= \eta_{ij} + \epsilon_{ij}, & 1 \leq i \leq r \text{ and } 1 \leq j \leq c, \\ \text{where } \eta_{ij} &= \mu + \alpha_i + \beta_j & \text{and } \epsilon_{ij} \sim N(0, \sigma^2 K_{ij}^{-1}). \end{aligned}$$

$K_{ij}$  is the number of observations, or the *count* in the  $(i, j)^{\text{th}}$  cell;  $\sigma^2 > 0$  is assumed to be known; and  $\epsilon_{ij}$  are independent Gaussian noise terms. Model (1) is over-parametrized, hence the parameters  $\mu, \boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_r)^{\top}, \boldsymbol{\beta} = (\beta_1, \dots, \beta_c)^{\top}$  are not identifiable without imposing further side conditions; however, the vector of cell means  $\boldsymbol{\eta} = (\eta_{11}, \eta_{12}, \dots, \eta_{rc})^{\top}$  is always identifiable. Our goal is to estimate  $\boldsymbol{\eta}$  under the sum-of-squares loss

$$(2) \quad L_{r,c}(\boldsymbol{\eta}, \hat{\boldsymbol{\eta}}) = \frac{1}{rc} \|\hat{\boldsymbol{\eta}} - \boldsymbol{\eta}\|^2 = \frac{1}{rc} \sum_{i=1}^r \sum_{j=1}^c (\hat{\eta}_{ij} - \eta_{ij})^2.$$

In model (1) the unknown quantities  $\alpha_i$  and  $\beta_j$  will be referred to as the  $i$ -th “row” (or “treatment”) and the  $j$ -th “column” (or “block”) effects,

respectively. In the *all-cells-filled* model,  $K_{ij} \geq 1$  for  $1 \leq i \leq r$  and  $1 \leq j \leq c$ ; the more general model, which allows some empty cells, is presented in Section 3. We would like to emphasize the focus in this section on the loss (2) rather than the weighted quadratic loss

$$L_{r,c}^{\text{wgt}}(\boldsymbol{\eta}, \hat{\boldsymbol{\eta}}) = \frac{1}{rc} \sum_{i=1}^r \sum_{j=1}^c K_{ij} (\hat{\eta}_{ij} - \eta_{ij})^2 ,$$

which is sometimes called the “prediction” loss, and under which asymptotically optimal estimation has been investigated before (Dicker, 2013). Nevertheless, in later sections results are presented for a general quadratic loss, which includes the weighted loss as a special case.

**Shrinkage estimators for the two-way model.** The usual estimator of  $\boldsymbol{\eta}$  is the weighted least squares (WLS) estimator, which is also maximum-likelihood under (1). The WLS estimator is unbiased and minimax but can be substantially improved on in terms of quadratic loss by shrinkage estimators, particularly when  $r, c \rightarrow \infty$  (Draper and Van Nostrand, 1979). Note that through out this paper we represent vectors in bold and matrices by capital letters. As the starting point for the shrinkage estimators proposed in this paper, we consider a family of Bayes estimators with respect to a conjugate prior on  $(\boldsymbol{\alpha}, \boldsymbol{\beta})$

$$\alpha_1, \dots, \alpha_r \stackrel{\text{i.i.d.}}{\sim} N(0, \sigma_A^2) \quad \text{and} \quad \beta_1, \dots, \beta_c \stackrel{\text{i.i.d.}}{\sim} N(0, \sigma_B^2) ,$$

where  $\sigma_A^2, \sigma_B^2$  are hyper-parameters. This prior extends the conjugate normal prior in the one-way case and was proposed by Lindley and Smith (1972) to reflect exchangeability within rows and columns separately. In vector form, the two-level hierarchical model is:

$$(3) \quad \begin{array}{ll} \text{Level 1:} & \mathbf{y} | \boldsymbol{\eta} \sim N_p(\boldsymbol{\eta}, \sigma^2 M) \quad \boldsymbol{\eta} = \mathbf{1}\mu + Z\boldsymbol{\theta} \quad \boldsymbol{\theta}^\top = (\boldsymbol{\alpha}^\top, \boldsymbol{\beta}^\top) \\ \text{Level 2:} & \boldsymbol{\theta} \sim N_q(0, \sigma^2 \Lambda \Lambda^\top) , \end{array}$$

where  $M = \text{diag}(K_{11}^{-1}, K_{12}^{-1}, \dots, K_{rc}^{-1})$  is an  $rc \times rc$  matrix and  $Z = [Z_A \ Z_B]$  with  $Z_A = I_r \otimes \mathbf{1}_c$  and  $Z_B = \mathbf{1}_r \otimes I_c$ . The  $(r+c) \times (r+c)$  matrix

$$\Lambda = \begin{bmatrix} \sqrt{\lambda_A} I_r & 0 \\ 0 & \sqrt{\lambda_B} I_c \end{bmatrix}$$

is written in terms of the relative variance components  $\lambda_A = \sigma_A^2/\sigma^2$  and  $\lambda_B = \sigma_B^2/\sigma^2$ . Henceforth, for notational simplicity, the dependence of  $\Lambda$  on the model hyper-parameters will be kept implicit. As shown in Lemma S.1.1

of the supplementary materials, the marginal variance of  $\mathbf{y}$  in (5) is given by  $\sigma^2 \Sigma$  where

$$(4) \quad \Sigma = Z\Lambda\Lambda^\top Z^\top + M = \lambda_A Z_A Z_A^\top + \lambda_B Z_B Z_B^\top + M.$$

At this point a comment is in order regarding shrinkage estimators for the general homoscedastic linear model. Note that model (1) could be written for individual, homoscedastic observations (with an additional subscript  $k$ ) instead of for the cell averages. With the corresponding design matrix, the two-way additive model is therefore a special case of the homoscedastic Gaussian linear model,  $\mathbf{y} \sim N_n(X\boldsymbol{\gamma}, \sigma^2 I)$ , where  $X \in \mathbb{R}^{n \times p}$  a known matrix and  $\boldsymbol{\gamma} \in \mathbb{R}^p$  is the unknown parameter. Thus, the various Stein-type shrinkage methods that have been proposed for estimating  $\boldsymbol{\gamma}$  can also be applied to our problem. Specifically, a popular approach is to reduce the problem of estimating  $\boldsymbol{\gamma}$  to the problem of estimating the mean of a  $p$ -dimensional *heteroscedastic* normal vector with known variances (see, e.g., [Johnstone, 2011](#), Section 2.9) by applying orthogonal transformations to the parameter  $\boldsymbol{\gamma}$  and data  $\mathbf{y}$ . Thereafter, Stein-type shrinkage estimators can be constructed as empirical Bayes rules by putting a prior which is either i.i.d. on the transformed coordinates or i.i.d. on the original coordinates of the parameter ([Rolph, 1976](#), referred to priors of the first type as *proportional* priors and to those of the second kind as *constant* priors). In the case of factorial designs, however, neither of these choices is very sensible, because they do not capture the (within-batch) symmetries of cross-classified models. Hence, procedures relying on models that take exchangeability into account can potentially achieve a significant and meaningful reduction in estimation risk. The estimation methodology we develop here incorporates the exchangeable structure of (3).

**Empirical Bayes estimators.** The following is a standard result and is proved in Section S.1.1 of the supplementary materials.

LEMMA 2.1. *For any fixed  $\mu \in \mathbb{R}$  and non-negative  $\lambda_A, \lambda_B$  the Bayes estimate of  $\boldsymbol{\eta}$  in (3) is given by:*

$$(5) \quad E[\boldsymbol{\eta}|\mathbf{y}] = \mathbf{y} - M\Sigma^{-1}(\mathbf{y} - \mathbf{1}\mu),$$

where the hyper-parameters  $\lambda_A, \lambda_B$  are involved in  $\Sigma$  through  $\Lambda$ .

Instead of fixing the values of  $\mu, \lambda_A, \lambda_B$  in advance, we may now return to model (1) and consider the parametric family of estimators

$$\mathcal{S}[\tau] = \left\{ \hat{\boldsymbol{\eta}}^S(\mu, \lambda_A, \lambda_B) = \mathbf{y} - M\Sigma^{-1}(\mathbf{y} - \mathbf{1}\mu) : \mu \in [\hat{a}_\tau(\mathbf{y}), \hat{b}_\tau(\mathbf{y})], \lambda_A \geq 0, \lambda_B \geq 0 \right\}.$$

Above,  $\mu$  is restricted to lie within  $\hat{a}_\tau(\mathbf{y}) = \text{quantile}\{y_{ij} : 1 \leq i \leq r, 1 \leq j \leq c; \tau/2\}$  and  $\hat{b}_\tau(\mathbf{y}) = \text{quantile}\{y_{ij} : 1 \leq i \leq r, 1 \leq j \leq c; 1 - \tau/2\}$ , the  $\tau/2$  and  $(1 - \tau/2)$  quantiles of the observations. The constraint on the location hyper-parameter  $\mu$  is imposed for technical reasons but is moderate enough to be well justified. Indeed, an estimator that shrinks toward a point that lies near the periphery or outside the range of the data is at the risk of being non-robust and seems to be an undesirable choice for a Bayes estimator corresponding to (3), which models  $\boldsymbol{\alpha}$  and  $\boldsymbol{\beta}$  as having zero means. In practice  $\tau$  may be taken to be 1% or 5%.

An empirical Bayes estimator is obtained by selecting for each observed  $\mathbf{y}$  a (possibly different) candidate from the family  $\mathcal{S}[\tau]$  as an estimate for  $\boldsymbol{\eta}$ ; equivalently, an empirical Bayes estimator is any estimator that plugs data-dependent values  $\hat{\mu}, \hat{\lambda}_A, \hat{\lambda}_B$  into (5), with the restriction that  $\hat{\mu}$  is in the allowable range. In the next section, we propose a specific criterion for estimating the hyperparameters.

*2.2. Estimation Methods.* The usual empirical Bayes estimators are derived relying on hierarchical model (3). The fixed effect  $\mu$  and the relative variance components  $\lambda_A$  and  $\lambda_B$  are treated as unknown fixed parameters to be estimated based on the marginal distribution of  $\mathbf{y}$  and substituted into (5). For any set of estimates substituted for  $\lambda_A$  and  $\lambda_B$ , the general mean  $\mu$  is customarily estimated by generalized least squares, producing an empirical version of what is known as the Best Linear Unbiased Predictors (BLUP). There is extensive literature on the estimation of the variance components (see chapters 5 and 6 of [Searle \*et al.\*, 2009](#)), with the main methods being maximum-likelihood (ML), restricted maximum-likelihood (REML), and the ANOVA methods (Method-of-Moments), including the three original ANOVA methods of Henderson ([Henderson, 1984](#)). Here we concentrate on the commonly used maximum-likelihood estimates, which are implemented in the popular R package `lme4` ([Bates \*et al.\*, 2014](#)). If  $\mathcal{L}(\mu, \lambda_A, \lambda_B; \mathbf{y})$  denotes the marginal likelihood of  $\mathbf{y}$  according to (3), then the maximum-likelihood (ML) estimates are

$$(7) \quad (\hat{\mu}^{\text{ML}}, \hat{\lambda}_A^{\text{ML}}, \hat{\lambda}_B^{\text{ML}}) = \underset{\mu \in [\hat{a}_\tau, \hat{b}_\tau], \lambda_A \geq 0, \lambda_B \geq 0}{\text{arg max}} \quad \mathcal{L}(\mu, \lambda_A, \lambda_B; \mathbf{y}).$$

The corresponding empirical Bayes estimator is  $\hat{\boldsymbol{\eta}}^{\text{ML}} = \hat{\boldsymbol{\eta}}^{\text{S}}(\hat{\mu}^{\text{ML}}, \hat{\lambda}_A^{\text{ML}}, \hat{\lambda}_B^{\text{ML}})$  and will be referred to as EBMLE (for Empirical Bayes Maximum-Likelihood).

LEMMA 2.2. *The ML estimates defined in (7) satisfy the following*

equations:

(8)

$$\text{I. } \hat{\boldsymbol{\mu}} = \hat{\boldsymbol{\mu}}_1 \cdot I\{\hat{\boldsymbol{\mu}}_1 \in [\hat{a}_\tau, \hat{b}_\tau]\} + \hat{a}_\tau \cdot I\{\hat{\boldsymbol{\mu}}_1 < \hat{a}_\tau\} + \hat{b}_\tau \cdot I\{\hat{\boldsymbol{\mu}}_1 > \hat{b}_\tau\}$$

$$\text{where, } \hat{\boldsymbol{\mu}}_1 = (\mathbf{1}^\top \hat{\Sigma}^{-1} \mathbf{y}) / (\mathbf{1}^\top \hat{\Sigma}^{-1} \mathbf{1}).$$

If  $\hat{\boldsymbol{\mu}}_1 \in [\hat{a}_\tau, \hat{b}_\tau]$  and  $\hat{\lambda}_a, \hat{\lambda}_b$  are both strictly positive, they satisfy

$$\text{II. } \text{tr}(\hat{\Sigma}^{-1} Z_A Z_A^\top) - \sigma^{-2} \mathbf{y}^\top (I - \hat{P})^\top \hat{\Sigma}^{-1} Z_A Z_A^\top \hat{\Sigma}^{-1} (I - \hat{P}) \mathbf{y} = 0$$

$$\text{III. } \text{tr}(\hat{\Sigma}^{-1} Z_B Z_B^\top) - \sigma^{-2} \mathbf{y}^\top (I - \hat{P})^\top \hat{\Sigma}^{-1} Z_B Z_B^\top \hat{\Sigma}^{-1} (I - \hat{P}) \mathbf{y} = 0,$$

$$\text{where } \hat{P} = \mathbf{1}(\mathbf{1}^\top \hat{\Sigma}^{-1} \mathbf{1})^{-1} \mathbf{1}^\top \hat{\Sigma}^{-1}.$$

The derivation is standard and provided in Section S.1.1.1 of the supplements, which also contain the estimating equation for the case when  $\hat{\boldsymbol{\mu}}_1 \notin [\hat{a}_\tau, \hat{b}_\tau]$ . If the solution to the estimating equations (9) includes a negative component, adjustments are needed in order produce the maximum-likelihood estimates of the scale hyper-parameters (see Searle and McCulloch, 2001, Section 2.2b-iii for a discussion of the one-way case).

**Estimation of hyper-parameters.** We propose an alternative method for estimating the shrinkage parameters. Following the approach of Xie *et al.* (2012), for fixed  $\tau \in (0, 1]$  we choose the shrinkage parameters by minimizing unbiased risk estimate (URE) over estimators  $\hat{\boldsymbol{\eta}}^S$  in  $\mathcal{S}[\tau]$ . By Lemma S.1.2 of the supplements, an unbiased estimate of the risk of  $\hat{\boldsymbol{\eta}}^S$ ,

$$R_{r,c}(\boldsymbol{\eta}, \hat{\boldsymbol{\eta}}^S(\mu, \lambda_A, \lambda_B)) \triangleq \frac{1}{rc} \mathbb{E} \|\hat{\boldsymbol{\eta}}^S(\mu, \lambda_A, \lambda_B) - \boldsymbol{\eta}\|^2,$$

is given by

$$\widehat{\text{URE}}(\mu, \lambda_A, \lambda_B) = \frac{1}{rc} \{\sigma^2 \text{tr}(M) - 2\sigma^2 \text{tr}(\Sigma^{-1} M^2) + (\mathbf{y} - \mathbf{1}\mu)^\top [\Sigma^{-1} M^2 \Sigma^{-1}] (\mathbf{y} - \mathbf{1}\mu)\}.$$

Hence we propose to estimate the tuning parameters of the class  $\mathcal{S}[\tau]$  by

$$(11) \quad (\hat{\boldsymbol{\mu}}^U, \hat{\lambda}_A^U, \hat{\lambda}_B^U) = \arg \min_{\mu \in [\hat{a}_\tau, \hat{b}_\tau], \lambda_A \geq 0, \lambda_B \geq 0} \widehat{\text{URE}}(\mu, \lambda_A, \lambda_B).$$

The corresponding empirical Bayes estimator is  $\hat{\boldsymbol{\eta}}^{\text{URE}} = \hat{\boldsymbol{\eta}}^S(\hat{\boldsymbol{\mu}}^U, \hat{\lambda}_A^U, \hat{\lambda}_B^U)$ . As in the case of maximum likelihood estimation, there is no closed-form solution to (11), but we can characterize the solutions by the corresponding estimating equations.



LEMMA 2.3. *The URE estimates of (11) satisfy the following estimating equations:*

(12)

$$\text{I. } \hat{\mu} = \hat{\mu}_1 \cdot I\{\hat{\mu}_1 \in [\hat{a}_\tau, \hat{b}_\tau]\} + \hat{a}_\tau \cdot I\{\hat{\mu}_1 < \hat{a}_\tau\} + \hat{b}_\tau \cdot I\{\hat{\mu}_1 > \hat{b}_\tau\}$$

$$\text{where, } \hat{\mu}_1 = (\mathbf{1}^\top [\hat{\Sigma}^{-1} M^2 \hat{\Sigma}^{-1}] \mathbf{y}) / (\mathbf{1}^\top [\hat{\Sigma}^{-1} M^2 \hat{\Sigma}^{-1}] \mathbf{1}) .$$

If  $\hat{\mu}_1 \in [\hat{a}_\tau, \hat{b}_\tau]$  and  $\hat{\lambda}_a, \hat{\lambda}_b$  are both strictly positive, they satisfy:

$$\text{II. } \text{tr}(\hat{\Sigma}^{-1} Z_A Z_A^\top \hat{\Sigma}^{-1} M^2) - \sigma^{-2} \mathbf{y}^\top (I - \hat{P})^\top \hat{\Sigma}^{-1} Z_A Z_A^\top \hat{\Sigma}^{-1} M^2 \hat{\Sigma}^{-1} (I - \hat{P}) \mathbf{y} = 0$$

$$\text{III. } \text{tr}(\hat{\Sigma}^{-1} Z_B Z_B^\top \hat{\Sigma}^{-1} M^2) - \sigma^{-2} \mathbf{y}^\top (I - \hat{P})^\top \hat{\Sigma}^{-1} Z_B Z_B^\top \hat{\Sigma}^{-1} M^2 \hat{\Sigma}^{-1} (I - \hat{P}) \mathbf{y} = 0 ,$$

where  $\hat{P} = \mathbf{1}(\mathbf{1}^\top [\hat{\Sigma}^{-1} M^2 \hat{\Sigma}^{-1}] \mathbf{1})^{-1} \mathbf{1}^\top \hat{\Sigma}^{-1} M^2 \hat{\Sigma}^{-1}$ .

The derivation is provided in Section S.1.1.1 of the supplementary materials. Comparing the two systems of equations (9) and (13) without substituting the value of  $\mu$ , it can be seen that the URE equation involves an extra term  $\hat{\Sigma}^{-1} M^2$  in both summands of the left-hand side, as compared to the ML equation. The estimating equations therefore imply that the ML and URE solutions may differ when the design is unbalanced. In Section 4, we show that the URE estimate  $\hat{\boldsymbol{\eta}}^{\text{URE}}$  is asymptotically optimal as  $r, c \rightarrow \infty$ , and the numerical simulations in Section 7 demonstrate that in certain situations EBMLE performs significantly worse.

**3. Estimation in Model with Missing Cells.** A more general model than (1) allows some cells to be empty. Hence, consider

$$(14) \quad \begin{aligned} y_{ij} &= \eta_{ij} + \epsilon_{ij} \quad \text{for } (i, j) \in \mathcal{E} \\ \eta_{ij} &= \mu + \alpha_i + \beta_j \quad \text{and } \epsilon_{ij} \sim N(0, \sigma^2 K_{ij}^{-1}) , \end{aligned}$$

where  $\mathcal{E} = \{(i, j) : K_{ij} \geq 1\} \subseteq \{1, \dots, r\} \otimes \{1, \dots, c\}$  is the set of indices corresponding to the nonempty cells. As before,  $\sigma^2 > 0$  is assumed to be known. Our goal is in general to estimate all cell means that are *estimable* under (14) rather than only the means of observed cells. For ease of presentation and without loss of generality, from here on we assume that  $\mathcal{E}$  is a connected design (Dey, 1986) so that all  $rc$  cell means are estimable.

We will need some new notation to distinguish between  $\mathbb{E}[\mathbf{y}] \in \mathbb{R}^{|\mathcal{E}|}$  and the  $rc$  vector consisting of all cell means. In general, the notation in (3) is reserved for quantities associated with the observed variables. As before,  $\boldsymbol{\theta} = (\boldsymbol{\alpha}^\top, \boldsymbol{\beta}^\top)^\top$ . The matrix  $M = \text{diag}(K_{ij}^{-1} : (i, j) \in \mathcal{E})$ , where the indices of diagonal elements are in lexicographical order. Let  $Z_c = [\mathbf{1}_{rc} \ I_R \otimes \mathbf{1}_C \ \mathbf{1}_R \otimes I_C]$

be the  $rc \times (r+c+1)$  design matrix associated with the unobserved complete model. The  $|\mathcal{E}| \times (r+c+1)$  ‘‘observed’’ design matrix  $Z$  is obtained from  $Z_c$  by deleting the subset of rows corresponding to  $\mathcal{E}^c$ . With the new definitions for  $Z$  and  $M$ , we define  $\Sigma$  by (4). Finally, let  $\boldsymbol{\eta}_c = Z_c \boldsymbol{\theta} \in \mathbb{R}^{rc}$  be the vector of all estimable cell means and  $\boldsymbol{\eta} = Z \boldsymbol{\theta} \in \mathbb{R}^{|\mathcal{E}|}$  be the vector of cell means for only the observed cells of (14). Hence, assuming  $\mathcal{E}$  corresponds to connected design, we consider estimating  $\boldsymbol{\eta}_c$  under the normalized sum-of-squares loss.

Note that since  $\boldsymbol{\eta}_c$  is estimable, it must be a linear function of  $\boldsymbol{\eta}$ . The following lemma is an application of the basic theory of estimable functions and is proved in the Section S.1.2 of the supplementary materials.

LEMMA 3.1. *If  $\boldsymbol{\eta}_c$  is estimable, then  $\boldsymbol{\eta}_c = Z_c(Z^\top Z)^- Z^\top \boldsymbol{\eta}$ , where  $(Z^\top Z)^-$  is any generalized inverse of  $Z^\top Z$ .*

In particular, writing  $Z^\dagger$  for the Moore-Penrose pseudo-inverse of  $Z$ , we therefore have  $\boldsymbol{\eta}_c = Z_c Z^\dagger \boldsymbol{\eta}$ . Thus, we can rewrite the loss function as

$$(15) L_{r,c}(\boldsymbol{\eta}_c, \hat{\boldsymbol{\eta}}_c) \triangleq \frac{1}{rc} \|\hat{\boldsymbol{\eta}}_c - \boldsymbol{\eta}_c\|^2 = \frac{1}{rc} (\hat{\boldsymbol{\eta}} - \boldsymbol{\eta})^\top Q (\hat{\boldsymbol{\eta}} - \boldsymbol{\eta}) = L_{r,c}^Q(\boldsymbol{\eta}, \hat{\boldsymbol{\eta}}),$$

where

$$(16) \quad Q = (Z_c Z^\dagger)^\top Z_c Z^\dagger.$$

In other words, the problem of estimating  $\boldsymbol{\eta}_c$  under sum-of-squares loss can be recast as the problem of estimating  $\boldsymbol{\eta} = \mathbb{E}[\mathbf{y}]$  under appropriate quadratic loss. This allows us to build on the techniques developed in the previous section and extend their applicability to the loss in (15). The standard unbiased estimator of  $\boldsymbol{\eta}_c$  is the weighted least squares estimator. The form of the Bayes estimator for  $\boldsymbol{\eta}$  under (3) is not affected by the generalized quadratic loss  $L_{r,c}^Q$  and is still given by (5), with  $M, \Sigma^{-1}$  as defined in the current section. As before, for any pre-specified  $\tau \in (0, 1]$  we consider the class of estimators  $\mathcal{S}[\tau]$  defined in (6). The EBMLE estimates the hyper-parameters  $\mu, \lambda_A, \lambda_B$  based on the marginal likelihood  $\mathbf{y}$  according to (3), where  $M, \Sigma^{-1}$  are as defined in the current section. As shown in Lemma S.1.3 of the supplements, an unbiased estimator of the point risk corresponding to (15),

$$R_{r,c}^Q(\boldsymbol{\eta}, \hat{\boldsymbol{\eta}}^S(\mu, \lambda_A, \lambda_B)) \triangleq \mathbb{E}\{L_{r,c}^Q(\boldsymbol{\eta}, \hat{\boldsymbol{\eta}}^S(\mu, \lambda_A, \lambda_B))\},$$

is given by

$$(17) \quad \widehat{\text{URE}}^Q(\mu, \lambda_A, \lambda_B) = (rc)^{-1} [\sigma^2 \text{tr}(QM) - 2\sigma^2 \text{tr}(\Sigma^{-1}MQM) + (\mathbf{y} - \mu\mathbf{1})^\top [\Sigma^{-1}MQM\Sigma^{-1}](\mathbf{y} - \mu\mathbf{1})].$$

The URE estimates of the tuning parameters are

$$(18) \quad (\hat{\mu}^{U_Q}, \hat{\lambda}_A^{U_Q}, \hat{\lambda}_B^{U_Q}) = \arg \min_{\mu \in [\hat{a}_\tau, \hat{b}_\tau], \lambda_A \geq 0, \lambda_B \geq 0} \widehat{\text{URE}}^Q(\mu, \lambda_A, \lambda_B),$$

and the corresponding EB estimate is  $\hat{\eta}^{\text{URE}} = \hat{\eta}^S(\hat{\mu}^{U_Q}, \hat{\lambda}_A^{U_Q}, \hat{\lambda}_B^{U_Q})$ . Equivalently, the estimate for  $\eta_c$  is  $\hat{\eta}_c^{\text{URE}} = Z_c Z_c^\dagger \hat{\eta}^S(\hat{\mu}^{U_Q}, \hat{\lambda}_A^{U_Q}, \hat{\lambda}_B^{U_Q})$ . The estimating equations for the URE as well as ML estimates of  $\mu, \lambda_A, \lambda_B$  can be derived similarly to those in the all-cells-filled model.

**4. Risk Properties and Asymptotic Optimality of the URE Estimator.** We now present the results that establish the optimality properties of our proposed URE-based estimator. We present the result for the quadratic loss  $L_{r,c}^Q$  of the previous section with the matrix  $Q$  defined in (16). Substituting  $Q$  with  $I_{rc}$  will give us the results for the fully-observed model (1), which are also explained. In proving our theoretical results we make the following assumptions:

**A1. On the parameter space:** We assume that the parameter  $\eta_c$  in the complete model is estimable and satisfies the following second order moment condition:

$$(19) \quad (A1) \quad \lim_{r,c \rightarrow \infty} \frac{1}{rc} \sum_{i=1}^r \sum_{j=1}^c \eta_{i,j}^2 < \infty.$$

This assumption is very mild, and similar versions are widely used in the EB literature (see Assumption  $C'$  of Xie *et al.*, 2012). It mainly facilitates a shorter technical proof and can be avoided by considering separate analyses of the extreme cases.

**A2. On the design matrix:** Denoting the largest eigenvalue of a matrix  $A$  by  $\lambda_1(A)$ , the matrix  $Q$  in (16) is assumed to satisfy

$$(20) \quad (A2) \quad \lim_{r,c \rightarrow \infty} (rc)^{-1/8} (\log(rc))^2 \nu_{r,c} \lambda_1(Q) = 0,$$

where  $\nu_{r,c} = \max\{K_{ij} : (i,j) \in \mathcal{E}\} / \min\{K_{ij} : (i,j) \in \mathcal{E}\}$ . As shown in Lemma A.6 in the Appendix,  $\lambda_1(Q)$  equals the largest eigenvalue of  $(Z_c' Z_c)(Z' Z)^\dagger$ . Intuitively, it represents the difference in information between the observed data matrix and the complete data matrix  $Z_c$ . If there are many empty cells,  $\lambda_1((Z_c' Z_c)(Z' Z)^\dagger)$  will be large and may violate the above condition. On the contrary, in the case of the completely observed data we have  $\lambda_1(Q) = 1$  (see Lemma A.6). Thus, in that case the assumption reduces

to  $\lim_{r,c \rightarrow \infty} (rc)^{-1/8} (\log(rc))^2 \nu_{r,c} = 0$ . This condition amounts to controlling in some sense the extent of imbalance in the number of observations procured per cell. Here, we are allowing the imbalance in the design to asymptotically grow to infinity but at a lower rate than  $(rc)^{1/8} / (\log(rc))^2$ . This assumption on the design matrix is essential for our asymptotic optimality proofs. Section A of the Appendix shows its role in our proofs and a detailed discussion about it is provided in the supplementary materials.

**Asymptotic optimality results.** The following theorem forms the basis for the results presented in this section:

**THEOREM 4.1.** *Under Assumptions A1-A2, with  $d_{r,c} = m_{r,c}^7 \nu_{r,c}^3 \lambda_1^3(Q)$  and  $m_{r,c} = \log(rc)$  we have*

$$\lim_{\substack{r \rightarrow \infty \\ c \rightarrow \infty}} d_{r,c} \cdot \left\{ \sup_{\substack{|\mu| \leq m_{r,c} \\ \lambda_A, \lambda_B \geq 0}} \mathbb{E} \left| \widehat{\text{URE}}_{r,c}^Q(\mu, \lambda_A, \lambda_B) - L_{r,c}^Q(\boldsymbol{\eta}, \widehat{\boldsymbol{\eta}}^S(\mu, \lambda_A, \lambda_B)) \right| \right\} = 0.$$

Theorem (4.1) shows that the unbiased risk estimator approximates the true loss pointwise uniformly well over a set of hyper-parameters where  $\lambda_A, \lambda_B$  can take any non-negative value and the location hyper-parameter  $\mu$  is restricted to the set  $[-m_{r,c}, m_{r,c}]$ , which grows as  $r, c$  increases. The set of all hyper-parameters considered in  $\mathcal{S}[\tau]$  differs from the aforementioned set, as there  $\mu$  was restricted to be in the data-dependent set  $[\hat{a}_\tau, \hat{b}_\tau]$ . However, as  $r, c \rightarrow \infty$ ,  $[\hat{a}_\tau, \hat{b}_\tau]$  is asymptotically contained in  $[-m_{r,c}, m_{r,c}]$  (see Lemma A.4), so Theorem 4.1 asymptotically covers all hyper-parameters considered in  $\mathcal{S}[\tau]$  for any  $\tau \in (0, 1]$ . This explains intuitively why in choosing the hyper-parameters by minimizing an unbiased risk estimate as in (17), we can expect the resulting estimate  $\widehat{\boldsymbol{\eta}}^S(\widehat{\mu}^{\text{UQ}}, \widehat{\lambda}_A^{\text{UQ}}, \widehat{\lambda}_B^{\text{UQ}})$  to have competitive performance. To compare the asymptotic performance of our proposed estimate, we define the oracle loss (OL) hyper-parameter as:

$$(\tilde{\mu}^{\text{OL}}, \tilde{\lambda}_A^{\text{OL}}, \tilde{\lambda}_B^{\text{OL}}) = \arg \min_{\mu \in [\hat{a}_\tau, \hat{b}_\tau]; \lambda_A, \lambda_B \geq 0} L^Q(\boldsymbol{\eta}, \widehat{\boldsymbol{\eta}}^S(\mu, \lambda_A, \lambda_B))$$

and the corresponding oracle rule

$$(21) \quad \tilde{\boldsymbol{\eta}}_c^{\text{OL}} = Z_c Z_c^\dagger \widehat{\boldsymbol{\eta}}^S(\tilde{\mu}^{\text{OL}}, \tilde{\lambda}_A^{\text{OL}}, \tilde{\lambda}_B^{\text{OL}}).$$

Note that the oracle rule depends on the unknown cell means  $\boldsymbol{\eta}_c$  and is therefore not a “legal” estimator. It serves as the theoretical benchmark for the minimum attainable error by any possible estimator: by its definition,

no EB estimator in our class can have smaller risk than  $\boldsymbol{\eta}_c^{\text{OL}}$ . The following two theorems show that our proposed URE-based estimator performs asymptotically nearly as well as the oracle loss estimator. The results hold for any class  $\mathcal{S}[\tau]$  where  $\tau \in (0, 1]$ . These results are in terms of the usual quadratic loss on the vector of all cell-means. Note that, based on our formulation of the problem in sections 2 and 3, both theorems 4.2 and 4.3 simultaneously cover the missing and fully-observed model.

**THEOREM 4.2.** *Under Assumptions A1-A2, for any  $\epsilon > 0$  we have*

$$\lim_{\substack{r \rightarrow \infty \\ c \rightarrow \infty}} P\{L_{r,c}(\boldsymbol{\eta}_c, \hat{\boldsymbol{\eta}}_c^{\text{URE}}) \geq L_{r,c}(\boldsymbol{\eta}_c, \tilde{\boldsymbol{\eta}}_c^{\text{OL}}) + \epsilon\} = 0 .$$

The next theorem asserts that under the same conditions, the URE-based estimator is asymptotically as good as the oracle estimator in terms of risk.

**THEOREM 4.3.** *Under Assumptions A1-A2, the following holds:*

$$\lim_{\substack{r \rightarrow \infty \\ c \rightarrow \infty}} R_{r,c}(\boldsymbol{\eta}_c, \hat{\boldsymbol{\eta}}_c^{\text{URE}}) - \mathbb{E}[L_{r,c}(\boldsymbol{\eta}_c, \tilde{\boldsymbol{\eta}}_c^{\text{OL}})] = 0 .$$

Finally, as the oracle performs better than any empirical Bayes estimator associated with  $\mathcal{S}[\tau]$ , a consequence of the above two theorems is that the URE-based estimator cannot be improved by any other such empirical Bayes estimator.

**COROLLARY 4.1.** *Under Assumptions A1-A2, it holds that for any estimator  $\hat{\boldsymbol{\eta}}^{\text{S}}(\hat{\mu}, \hat{\lambda}_A, \hat{\lambda}_B)$  corresponding to the class  $\mathcal{S}[\tau]$  we have*

$$(a) \quad \lim_{r \rightarrow \infty, c \rightarrow \infty} P\{L_{r,c}^{\text{Q}}(\boldsymbol{\eta}, \hat{\boldsymbol{\eta}}^{\text{URE}}) \geq L_{r,c}^{\text{Q}}(\boldsymbol{\eta}, \hat{\boldsymbol{\eta}}^{\text{S}}(\hat{\mu}, \hat{\lambda}_A, \hat{\lambda}_B)) + \epsilon\} = 0 .$$

$$(b) \quad \limsup_{r \rightarrow \infty, c \rightarrow \infty} R_{r,c}^{\text{Q}}(\boldsymbol{\eta}, \hat{\boldsymbol{\eta}}^{\text{URE}}) - R_{r,c}^{\text{Q}}(\boldsymbol{\eta}, \hat{\boldsymbol{\eta}}^{\text{S}}(\hat{\mu}, \hat{\lambda}_A, \hat{\lambda}_B)) \leq 0 .$$

Unlike the above two theorems, this corollary is based on the quadratic loss  $L^{\text{Q}}$ . It emphasizes the nature of our estimating class  $\mathcal{S}[\tau]$ . In Section 2 we saw that the EBMLE and URE generally produce different solutions in unbalanced designs; combined with Corollary (4.1), this implies that, asymptotically EBMLE generally does not achieve the optimal risk of an EB estimator corresponding to the class  $\mathcal{S}[\tau]$  (otherwise the EBML estimate for  $\boldsymbol{\eta}$  would have to be very close to the URE estimate).

The proofs of theorems 4.2 and 4.3 and that of Corollary 4.1 is left to Section A of the Appendix. The proofs rely heavily on the asymptotic risk estimation result of Theorem 4.1, which in turn uses the asymptotic risk

properties of estimators in  $\mathcal{S}[\tau]$ . Below, we sketch its proof by describing the interesting risk properties of these estimations.

To conclude this section, we would like to point out the qualitative differences between the type of results included in the current section and the results for the one-way normal mean estimation problem exemplified in [Xie \*et al.\* \(2012\)](#) and especially point out the differences in the proof techniques. In estimation theory, the optimality of shrinkage estimators in one-way problems is usually studied through a sequence model (see Ch. 2 of [Johnstone, 2011](#)), where there is a natural indexing on the dimensions in the parametric spaces. In unbalanced designs, the cell mean estimation problem in 2-way layouts cannot be reduced to estimating independent multiple vectors, and so there is no indexing on the parametric space under which the “row” effects and the “column” effects can be decoupled. Thus, the approach of [Xie \*et al.\* \(2012\)](#), which would require showing uniform convergence of the difference between the URE and the loss over the hyper-parametric space, i.e., showing  $L_1$  convergence of  $\sup_{\mu \in [\hat{a}_\tau, \hat{b}_\tau]; \lambda_A, \lambda_B \geq 0} |\widehat{\text{URE}}_{r,c}^Q(\mu, \lambda_A, \lambda_B) - L_{r,c}^Q(\boldsymbol{\eta}, \hat{\boldsymbol{\eta}}^S(\mu, \lambda_A, \lambda_B))|$  to 0, cannot be trivially adapted to the two-way layout. Instead, in [Theorem 4.1](#) we show the pointwise convergence of the expected absolute difference between the URE and the loss. Specifically, we show that as  $r, c \rightarrow 0$ , it converges at a rate  $d_{r,c}$  uniformly over the essential support of the hyper-parameters. Using this pointwise convergence, its rate and the properties of the loss function (see [Section A.2](#)), we prove the optimality results of [theorems 4.2, 4.3](#), which are of the same flavor as those in [Xie \*et al.\* \(2012\)](#) for the one-way case. Our pointwise convergence approach greatly helps to tackle the difficulties encountered when passing to the two-way problem, where computations involving matrices are generally unavoidable. Our pointwise convergence result is proved by a moment-based concentration approach, which translates the problem into bounding moments of Gaussian quadratic forms involving matrices with possibly dependent rows and columns. The following two lemmas, which are used in proving [Theorem 4.1](#), display our moment-based convergence approach, where the concentration of relevant quantities about their respective mean is proved. To prove [Theorem 4.1](#) we first show [Lemma 4.1](#), in which the URE methodology estimates the risk in  $L_2$  norm pointwise uniformly well for all estimators in  $\mathcal{S}[\tau]$  that shrink towards the origin (i.e., with  $\mu$  set at 0). Thereafter, in [Lemma 4.2](#) we prove that the loss of those estimators concentrate around their expected values (risk) when we have large number of row and column effects.

LEMMA 4.1. *Under Assumptions A1-A2, with  $d_{r,c} = m_{r,c}^7 \nu_{r,c}^3 \lambda_1^3(Q)$ ,*

$$m_{r,c} = \log(rc),$$

$$\lim_{\substack{r \rightarrow \infty \\ c \rightarrow \infty}} d_{r,c}^2 \cdot \left\{ \sup_{\lambda_A, \lambda_B \geq 0} \mathbb{E} \left[ \widehat{\text{URE}}_{r,c}^Q(0, \lambda_A, \lambda_B) - R_{r,c}^Q(\boldsymbol{\eta}, \widehat{\boldsymbol{\eta}}^S(0, \lambda_A, \lambda_B)) \right]^2 \right\} = 0.$$

LEMMA 4.2. *Under Assumptions A1-A2, with  $d_{r,c} = m_{r,c}^7 \nu_{r,c}^3 \lambda_1^3(Q)$ ,  $m_{r,c} = \log(rc)$ ,*

$$\lim_{\substack{r \rightarrow \infty \\ c \rightarrow \infty}} d_{r,c}^2 \cdot \left\{ \sup_{\lambda_A, \lambda_B \geq 0} \mathbb{E} \left[ L_{r,c}^Q(\boldsymbol{\eta}, \widehat{\boldsymbol{\eta}}^S(0, \lambda_A, \lambda_B)) - R_{r,c}^Q(\boldsymbol{\eta}, \widehat{\boldsymbol{\eta}}^S(0, \lambda_A, \lambda_B)) \right]^2 \right\} = 0.$$

If we restrict ourselves to only estimators in  $\mathcal{S}[\tau]$  that shrink towards the origin, then Theorem 4.1 follows directly from the above two lemmas. As such, for this subset of estimators, the lemmas prove a stronger version of the theorem with convergence in  $L_2$  norm. The proof is extended to general shrinkage estimators by controlling the  $L_1$  deviation between the true loss and its URE-based approximation through the nontrivial use of the location invariance structure of the problem. The proofs of all these results are provided in Section A of the Appendix. The results for the weighted loss  $L_{r,c}^{\text{wgt}}$  (defined in Section 2) are discussed in Section S.1.3.1 of the supplements.

**5. URE in Balanced Designs.** In this section we inspect the case of a balanced design,  $K_{ij} = k$ ,  $1 \leq i \leq r, 1 \leq j \leq c$ . We show that under a balanced design the problem essentially decouples into two independent one-way problems, in which case the URE and EBMLE estimates coincide (see also Xie *et al.*, 2012, Section 2). As a bonus, the analysis will suggest another class of shrinkage estimators for the general, unbalanced two-way problem by utilizing the one-way estimates of Xie *et al.* (2012).

To carry out the analysis, suppose without loss of generality that  $K = 1$ . Let the grand mean and the row and column main effects be

$$(22) \quad m = \mu + \alpha + \beta, \quad a_i = \alpha_i - \alpha., \quad b_j = \beta_j - \beta.$$

and let  $\mathbf{a} = (a_1, \dots, a_r)^\top$ ,  $\mathbf{b} = (b_1, \dots, b_c)^\top$ . Then, in the balanced case, the Bayes estimator  $\widehat{\boldsymbol{\eta}}^S(y_{..}, \lambda_A, \lambda_B)$ , obtained by substituting the mean of  $\mathbf{y}$  for  $\mu$  in (5), is

$$(23) \quad \left\{ \widehat{\boldsymbol{\eta}}_{ij}^S(y_{..}, \lambda_A, \lambda_B) \right\} = \widehat{m}^{\text{LS}} + c_\alpha(\lambda_A) \widehat{\mathbf{a}}_i^{\text{LS}} + c_\beta(\lambda_B) \widehat{\mathbf{b}}_j^{\text{LS}},$$

$$(24) \quad \text{where } \widehat{m}^{\text{LS}} = y_{..}, \quad \widehat{\mathbf{a}}_i^{\text{LS}} = y_{i.} - y_{..}, \quad \widehat{\mathbf{b}}_j^{\text{LS}} = y_{.j} - y_{..}$$

are the least squares estimators, and  $c_\alpha := c_\alpha(\lambda_A) = \lambda_A/(\lambda_A + \sigma^2/c)$  and  $c_\beta := c_\beta(\lambda_B) = \lambda_B/(\lambda_B + \sigma^2/r)$  are functions involving, respectively, only  $\lambda_A$  or only  $\lambda_B$ . Its risk  $R(\boldsymbol{\eta}, \widehat{\boldsymbol{\eta}}^S(y_{..}, \lambda_A, \lambda_B))$  decomposes as

$$(25) \quad \mathbb{E}\left\{(\widehat{m}^{\text{LS}} - m)^2\right\} + \frac{1}{r}\mathbb{E}\left\{\sum_{i=1}^r (c_\alpha \widehat{a}_i^{\text{LS}} - a_i)^2\right\} + \frac{1}{c}\mathbb{E}\left\{\sum_{j=1}^c (c_\beta \widehat{b}_j^{\text{LS}} - b_j)^2\right\}$$

due to the orthogonality of the vectors corresponding to the three sums-of-squares (detailed derivation is provided in the supplements). Consequently, one obtains URE by writing URE for each of summands above. Moreover, since

$$(26) \quad \widehat{\boldsymbol{m}}^{\text{LS}} \sim N(m, \sigma^2 \lambda_m^2), \quad \widehat{\boldsymbol{a}}^{\text{LS}} \sim N_r(a, \sigma^2 \Lambda_a), \quad \widehat{\boldsymbol{b}}^{\text{LS}} \sim N_c(b, \sigma^2 \Lambda_b),$$

minimizing URE jointly over  $(c_\alpha, c_\beta)$  therefore consists of minimizing separately the “row” term over  $c_\alpha$  and the “column” term over  $c_\beta$ . Each of these is a “one-way” Gaussian homoscedastic problem, except that the covariance matrices  $\Lambda_\alpha, \Lambda_\beta$  are singular because the main effects are centered. The unbiased risk estimator will naturally take this into account and will possess the “correct” degrees-of-freedom.

The maximum-likelihood estimates for the two-way random-effects additive model do not have a closed-form solution even for balanced data (Searle *et al.*, 2009, Ch. 4.7 d.), so it is not possible that they always produce the same estimates as discussed above. On the other hand, the REML estimates coincide with the positive-part Moments method estimates (Searle *et al.*, 2009, Ch. 4.8), which, in turn, reduce (for known  $\sigma^2$ ) to solving separately two one-way problems involving  $\widehat{\boldsymbol{a}}^{\text{LS}}$  for the rows and  $\widehat{\boldsymbol{b}}^{\text{LS}}$  for the columns. These have closed-form solutions and are easily seen to coincide with the URE solutions.

In the unbalanced case, (23) no longer holds, and so the Bayes estimates for  $\boldsymbol{a}$  and  $\boldsymbol{b}$  are each functions of both  $\widehat{\boldsymbol{a}}^{\text{LS}}$  and  $\widehat{\boldsymbol{b}}^{\text{LS}}$ . We can nevertheless use shrinkage estimators of the form (23) and look for “optimal” constants  $c_\alpha = c_\alpha(\lambda_A)$  and  $c_\beta = c_\beta(\lambda_B)$ . Appealing to the asymptotically optimal



one-way methods of [Xie \*et al.\* \(2012\)](#), we consider the estimator

(27)

$$\widehat{\eta}_{ij}^{\text{XKB}} = \widehat{m}^{\text{LS}} + \widehat{c}_\alpha^{\text{XKB}} \widehat{a}_i^{\text{LS}} + \widehat{c}_\beta^{\text{XKB}} \widehat{b}_j^{\text{LS}}, \quad 1 \leq i \leq r, 1 \leq j \leq c,$$

(28)

$$\text{where, } \widehat{c}_\alpha^{\text{XKB}} = \arg \min_{c_\alpha \in [0,1]} \widehat{\text{URE}} \left\{ \sum_{i=1}^r (c_\alpha \widehat{a}_i^{\text{LS}} - a_i)^2 \right\},$$

(29)

$$\widehat{c}_\beta^{\text{XKB}} = \arg \min_{c_\beta \in [0,1]} \widehat{\text{URE}} \left\{ \sum_{j=1}^c (c_\beta \widehat{b}_j^{\text{LS}} - b_j)^2 \right\}.$$

A slight modification of the parametric SURE estimate of [Xie \*et al.\* \(2012\)](#) that shrinks towards 0 is required to accommodate the covariance structure of the centered random vectors  $\widehat{\mathbf{a}}^{\text{LS}}, \widehat{\mathbf{b}}^{\text{LS}}$ . Contrasting the performance of the optimal empirical Bayes estimators corresponding to this class of shrinkage estimators with that corresponding to the class  $\mathcal{S}[\tau]$  of EB estimators can be taken to quantify the relative efficiency of using one-way methods in the two-way problem.

**6. Computation of the URE Estimator.** To compute the hyperparameter estimates by the URE method, one could attempt to solve the estimating equations in (13), which have no closed-form solution. For example, one could fix the value of  $\lambda_A$  to some initial positive value and solve the first equation in  $\lambda_B$ . Then, plug the solution into the second equation and solve for  $\lambda_A$ , and keep iterating between the two equations until convergence. If this approach is taken, a non-trivial issue to overcome will be obtaining the actual minimizing values  $\lambda_A$  and  $\lambda_B$  when one of the solutions to (13) is negative. Another issue will be ascertaining the global optimality of the solutions, as  $\widehat{\text{URE}}$  is not necessarily convex in  $(\mu, \lambda_A, \lambda_B)$ . To bypass these issues, we minimize  $\widehat{\text{URE}}$  by conducting a grid search on the scale hyperparameters, and  $\mu$  is subsequently estimated by (12).

A major hindrance for computations in large designs is the occurrence of the  $(rc) \times (rc)$  matrix  $\Sigma^{-1}$ , which depends on  $\lambda_A$  and  $\lambda_B$ . Inverting it can be a prohibitive task for even moderately large values of  $r$  and  $c$ , and it would need inversion at every point along the grid for a naive implementation. In our implementation, we adopt some of the key computational elements from the `lme4` package [Sec. 5.4 [Bates, 2010](#)] and produce an algorithm that works as fast as the computation of the EBMLE estimate with the `lme4` R-package.

For the case of no empty cells, the pivotal step in our implementation is the representation of the  $\widehat{\text{URE}}$  criterion by the following expression:

$$(30) \quad \widehat{\text{URE}} = (rc)^{-1} \left[ -\sigma^2 \text{tr}(M) + 2\sigma^2 \text{tr}\{(\Lambda^\top Z^\top M^{-1} Z \Lambda + I_q)^{-1} (\Lambda^\top Z^\top Z \Lambda)\} \right. \\ (31) \quad \left. + \|MV^{-1}(\mathbf{y} - \mathbf{1}\mu)\|^2 \right].$$

The detailed steps for deriving (30) are provided in Section B of the appendix where the computation of each of the above terms is also elaborately explained. (30) is numerically minimized jointly over  $(\lambda_A, \lambda_B)$ , where the key step in evaluating it for a particular pair  $(\lambda_A, \lambda_B)$  is employing a sparse Cholesky decomposition for the matrix  $\Lambda^\top Z^\top M^{-1} Z \Lambda + I_q$ . This decomposition takes advantage of the high sparsity of  $\Lambda^\top Z^\top M^{-1} Z \Lambda + I_q$ . It first determines the locations of non-zero elements in the Cholesky factor, which do not depend on the values of  $(\lambda_A, \lambda_B)$  and hence this stage is needed only once during the numerical optimization. This is the only costly stage of the decomposition and determining the values of the non-zero components is repeated during the numerical optimization. For the empty-cells case, the implementation is very similar after using the reduction to quadratic loss in  $L^Q$  described in Section 3.

**7. Simulation Study.** We carry out numerical experiments to compare the performance of the URE based estimator to that of different cell means estimators discussed in the previous sections. As the standard technique we consider the weighted Least Squares estimator  $\widehat{\boldsymbol{\eta}}^{\text{LS}} = \widehat{\boldsymbol{\mu}}^{\text{LS}} \mathbf{1} + Z \widehat{\boldsymbol{\theta}}^{\text{LS}}$ , where  $(\widehat{\boldsymbol{\mu}}^{\text{LS}}, \widehat{\boldsymbol{\theta}}^{\text{LS}})$  is any pair that minimizes

$$(\mathbf{y} - \boldsymbol{\mu} \cdot \mathbf{1} - Z\boldsymbol{\theta})^\top M^{-1} (\mathbf{y} - \boldsymbol{\mu} \cdot \mathbf{1} - Z\boldsymbol{\theta}).$$

The two-way shrinkage estimators reported are the maximum-likelihood empirical Bayes (EBML) estimator  $\widehat{\boldsymbol{\eta}}^{\text{ML}}$  and the URE based estimator  $\widehat{\boldsymbol{\eta}}^{\text{URE}}$ , as well as versions of these two estimators which shrink towards the origin (i.e., with  $\boldsymbol{\mu}$  fixed at 0); these are designated in Table 1 as “EBMLE (origin)” and “URE (origin)”. We also consider the generalized version of  $\widehat{\boldsymbol{\eta}}^{\text{XKB}}$  discussed in Section 5 which shrinks towards a general data-driven location and estimates the scale hyper-parameters based on two independent one-way shrinkage problems. For a benchmark we consider the oracle rule  $\widehat{\boldsymbol{\eta}}^{\text{OL}} = \widehat{\boldsymbol{\eta}}^{\text{S}}(\widetilde{\boldsymbol{\mu}}, \widetilde{\lambda}_A, \widetilde{\lambda}_B)$  where,

$$(32) \quad (\widetilde{\boldsymbol{\mu}}, \widetilde{\lambda}_A, \widetilde{\lambda}_B) = \arg \min_{\boldsymbol{\mu}, \lambda_A \geq 0, \lambda_B \geq 0} \|\mathbf{y} - M \boldsymbol{\Sigma}^{-1} (\mathbf{y} - \boldsymbol{\mu} \cdot \mathbf{1}) - \boldsymbol{\eta}\|^2.$$

Since for any  $\mathbf{y}$  the oracle rule minimizes the loss over all members of the parametric family (6), its expected loss lower bounds the risk achievable by any empirical Bayes estimator of the form (5).

**Simulation setup.** We report results across 6 simulation scenarios. For each of them, we draw  $(\boldsymbol{\alpha}, \boldsymbol{\beta}, M^{-1} = \text{diag}(K_{11}, K_{12}, \dots, K_{rc}))$  jointly from some distribution such that the cell counts  $K_{ij}$  are i.i.d. and  $(\boldsymbol{\alpha}, \boldsymbol{\beta})$  are drawn from some conditional distribution given the  $K_{ij}$ s. We then draw  $y_{ij} \sim N(\mu + \alpha_i + \beta_j, \sigma^2 K_{ij})$  independently, fixing  $\mu = 0$  throughout and setting  $\sigma^2$  to some (known) constant value. This process is repeated for  $N = 100$  time for each pair  $(r, c)$  in a range of values, and the average squared loss over the  $N$  rounds is computed for each of the estimators mentioned above. With 100 repetitions, the standard error of the average loss for each estimator is at least one order-of-magnitude smaller than the estimated differences between the risks; hence, the differences can be safely considered significant. The URE estimate is computed using the implementation described in Section 6, and the oracle “estimate” is computed employing a similar technique. The EBMLE estimate is computed using the R package `lme4` (Bates *et al.*, 2014).

Table 1 shows the estimation errors of different estimators as a fraction of the estimated risk of the Least Squares (LS) estimator. We have equal number of row and column levels for all experiments except for scenario (c). In Figure 1, we have the plot of the mean square errors (MSE) of the URE, EBMLE, LS and the Oracle loss (OL) estimators across the six experiments as the number of levels in the design varies. It shows how the estimation errors of the different estimators compare with the minimum achievable (oracle) error rates as the number of levels in the designs increases. The general pattern reflected in the subplots shows an initial sharp decline with a gradual flattening-out of the error rates as the number of levels exceeds 100, suggesting a setting within the asymptotic regime. In all the examples, the performance of our proposed URE based method is close to that of the oracle when the number of levels is large; for levels greater than 60, there is no other estimator which is much better at any instance than the URE. On the contrary, in all examples except scenario (a) the EBMLE performs quite bad, and gets outperformed even by the “one-way” XKB estimator. In cases with dependency between the effects and the cell counts, even the LS estimator can be preferable to the EBMLE (experiments (b) and (d)).

(a) *Hierarchical Gaussian Model.* For  $L \in \{20, 60, \dots, 180\}$  we set  $r = c = L$  and  $\sigma^2 = 25$ .  $K_{ij}$  are independent such that  $P(K_{ij} = 1) = 0.9$  and  $P(K_{ij} = 9) = 0.1$ . For  $1 \leq i, j \leq L$ ,  $\alpha_i, \beta_j$  are drawn from a  $N(0, \sigma^2/(4L))$  distribution independently of the  $K_{ij}$ s. The joint distribution of the row

	(a)	(b)	(c)	(d)	(e)	(f)
LS	1.00	1.00	1.00	1.00	1.00	1.00
EBMLE	0.31	1.79	0.48	1.37	0.21	0.96
URE	<b>0.31</b>	<b>0.45</b>	<b>0.19</b>	<b>0.21</b>	<b>0.18</b>	<b>0.58</b>
EBMLE (origin)	0.31	0.69	0.45	1.42	0.58	0.95
URE (origin)	0.31	0.46	0.20	0.53	0.57	0.63
XKB	0.31	0.58	0.28	0.44	0.20	-
Oracle	0.30	0.42	0.16	0.20	0.17	0.56

TABLE 1

Estimation errors relative to the Least Squares (LS) estimator. The columns in the table correspond to the six simulation examples described in section 7.

effects, column effects and the  $K_{ij}$ s in this example obeys the Bayesian model under which the parametric estimator (5) is derived. Hence the true Bayes rule is of that form, and the EBMLE is expected to perform well estimating the hyperparameters from the marginal distribution of  $\mathbf{y}$ . Indeed, the risk curve of the EBMLE approaches that of the oracle rule and seems to perform best for relatively small value of  $L$ . The MSE of the URE estimator, however converges to the oracle risk as  $L$  increases. Interestingly, the performance of the XKB estimator seems to be comparable to that of URE and EBMLE for large values of  $L$ .

(b) *Gaussian model with dependency between effects and cell counts.* For  $L \in \{20, 60, \dots, 180\}$  we set  $r = c = L$  and  $\sigma^2 = 25$ . In this example the  $K_{ij}$  are no longer independent of the random effects. We take  $K_{ij} = 1 \cdot (1 - Z_i) + 25 \cdot Z_i$  where  $Z_i \sim \text{Bin}(1, 0.5)$  independently, so that the cell frequencies are constant in each row. If  $Z_i = 1$ ,  $\alpha_i$  is drawn from a  $N(1, \sigma^2/(100 \cdot 2L))$  distribution, and otherwise from a  $N(0, \sigma^2/(2 \cdot L))$  distribution.  $\beta_j$  are drawn independently from a  $N(0, \sigma^2/(2L))$  distribution. The advantage of our URE method over the EBMLE is clear in Figure 1; in fact, even the LS estimator seems to do better than the EBMLE for the values of  $L$  considered here, a consequence of the strong dependency between the cell frequencies and the random effects. Again the XKB estimator performs surprisingly well.

(c) *Scenario (b) for different number of row and column effects.* This example is the same as example (b), except that we fix  $c = 40$  throughout and study the performance of the different estimators as number of row levels  $r = L \in \{20, 60, \dots, 180\}$  varies. The performance of the LS estimator relative to the other methods is much worse than in the previous examples. The performance the URE estimator gets closer to that of the oracle as  $r = L$  increases. The MSE of the XKB is significantly higher than that of the URE but much lower than that of the EBMLE.

(d) *Non-Gaussian row effects.* For  $L \in \{20, 60, \dots, 180\}$  we set  $r = c = L$  and  $\sigma^2 = 25$ . In this example the row effects are *determined* by the  $K_{ij}$ . We take

$K_{ij} = 1 \cdot (1 - Z_i) + 25 \cdot Z_i$  where  $Z_i \sim \text{Bin}(1, 0.5)$  independently, and set  $\alpha_i = 1 \cdot (1 - Z_i) + (1/25) \cdot Z_i$ .  $\beta_j$  are drawn independently from a  $N(0, \sigma^2/(2L))$  distribution. The URE estimator performs significantly better than the other estimators for large values of  $L$ , with about 50% smaller estimated risk for  $L = 180$  than that of the XKB estimator, and even much better compared to the other methods.

(e) *Correlated Main Effects.* For  $L \in \{20, 60, \dots, 180\}$  we set  $r = c = L$  and  $\sigma^2 = 25$ . In this example both the row and the column effects are determined by the  $K_{ij}$ . The cell frequencies  $K_{lj} = \max(T_l, 1)$ ,  $1 \leq l \leq L, 1 \leq j \leq L$ , where  $T_l$ ,  $1 \leq l \leq L$ , are drawn independently from a mixture of a  $\text{Poisson}(1)$  and  $\text{Poisson}(5)$  distributions with weights 0.9 and 0.1, respectively. The row and column effects are  $\alpha_l, \beta_l = 1/T_l$ ,  $1 \leq l \leq L$ . The MSE of the URE estimator is smaller than that of EBMLE by 14.7% ( $\widehat{\text{sd}}(\text{diff}) < 4 \cdot 10^{-5}$ ) for  $L = 200$ , but difference is not as big as in previous examples. The LS estimator performs considerably worse than the rest.

(f) *Missing Cells.* In the last example we study the performance of the estimators when some cells are empty. The setting is exactly as in example (b), except that after the  $K_{ij}$  are drawn, each  $K_{ij}$  is independently set to 0 (corresponding to an empty cell) with probability 0.2. In accordance with the theory, the performance of the URE estimator approaches the oracle loss, and for  $L = 180$  achieves significantly smaller risk than that of the EBMLE, although not as significantly smaller as in example (b) with all cells filled (40% vs 75% smaller than EBMLE for examples (f) and (b), respectively). The performance of the LS estimator is comparable to that of the EBMLE. The XKB estimator is not considered here as it is not applicable when some data are missing.

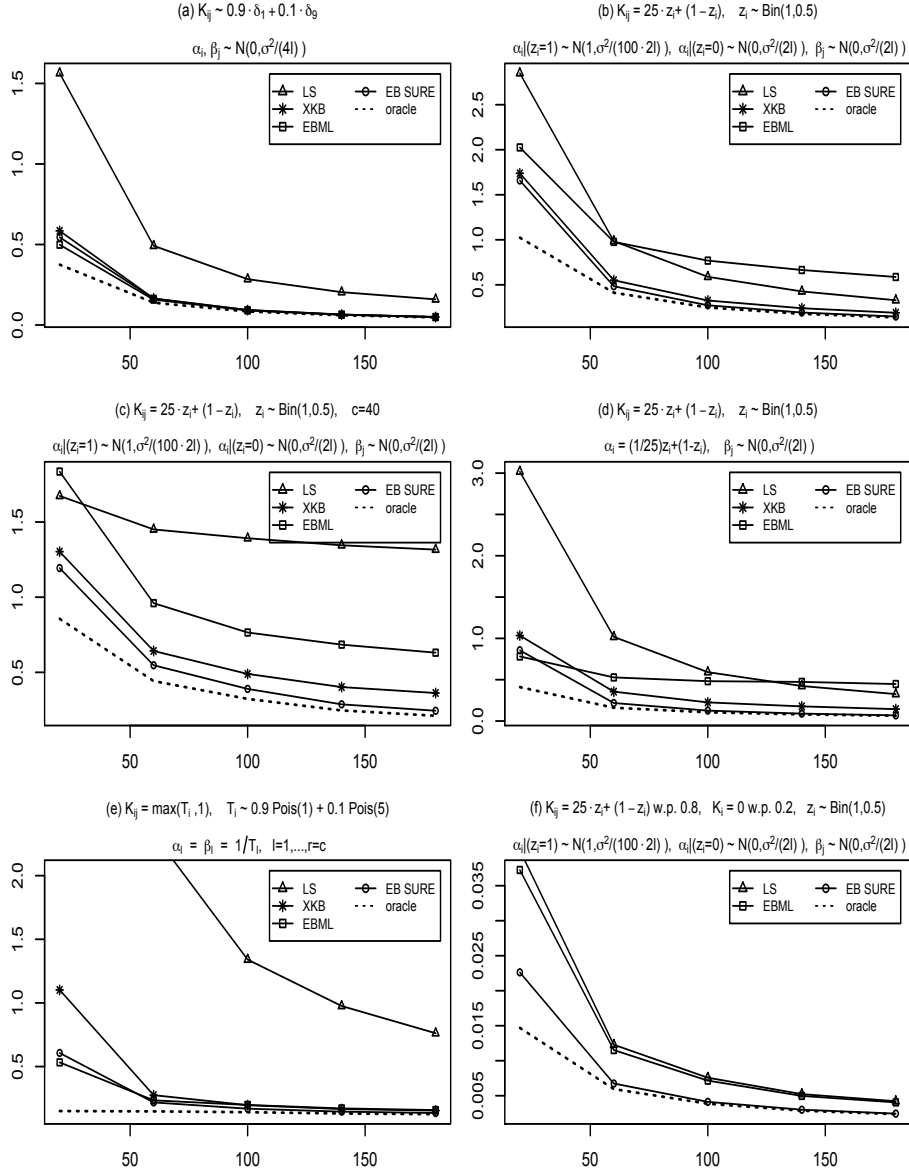


FIG 1. Risk of the various estimators in the six simulation scenarios of Table 1. The ordinate shows the risk of the estimators while we vary  $L$  along the abscissa. Recall,  $L = r = c$  for all experiments in the table except (d) where  $L = r$  and  $c$  was fixed at 40. .

**8. Real Data Example.** We analyze data collected on Nitrate levels measured in water sources across the US. Nitrates are chemical units found in drinking water that may lead to adverse health effects. According to the U.S. Geological Survey (USGS), excessive nitrate levels can result in restriction of oxygen transport in the bloodstream. The data was obtained from the Water Quality Portal cooperative (<http://waterqualitydata.us/>).

We consider estimating the average Nitrate levels based on the location of the water resource and time when the measurement was taken. Specifically, we fit the homoscedastic Gaussian, additive two-way model

$$(33) \quad y_{ijk} = \eta_{ij} + \epsilon_{ijk}, \quad \eta_{ij} = \mu + \alpha_i + \beta_j \quad k = 1, \dots, K_{ij}$$

where  $\alpha_i$  is the effect associated with the  $i$ -th level of a categorical variable indicating the hour of the day when the measurement was taken (by rounding to the past hour, e.g., for 14:47 the hour is 14);  $\beta_j$  is the effect associated with the  $j$ -th US county; and  $y_{ijk}$  is the corresponding log-transformed measurement of Nitrate level (in mg/l). The errors  $\epsilon_{ijk}$  are treated as i.i.d. Gaussian with a fixed (known) variance equal to the the LS estimate  $\hat{\sigma}^2$ . We used records from January and February of 2014, and concentrated on measurements made between 8:00 and 17:00 as those were the most active hours. This yielded a total of 858 observations categorized into 9 different hour-slots (8-16) and 108 counties across the entire country. The data is highly unbalanced: 57% of the cells are empty, and the cell counts among the nonempty cells vary between 1 to 12. Figure 2 (left panel) shows the residuals from the standard LS fit for the data (note that this assumes independence of the noise terms). The alignment with the normal quantiles is better around the center of the distribution.

A two-way Analysis-of-Variance yielded a highly significant p-value for county ( $< 10^{-5}$ ) but not for hour (0.25), for comparing the models with an without each variable (i.e., using Type II sums of squares). For the estimation problem, we considered the two-way shrinkage estimators, EBMLE and URE, as well as the “pre-test” estimator which, failing to reject the null hypothesis for the overall effect of hour, proceeds with fitting the one-way LS estimate by county. We will refer to the latter as the “one-way” estimator or as “LS-county”. As a two-way estimator, it can be interpreted as shrinking all the way to zero on hour, while providing no shrinkage at all for county. The “usual” estimator is the LS estimator based on (33).

Applying the shrinkage estimators to the entire data set, we observe that both shrink the LS estimates, but the shrinkage factors are quite different. Table 2 shows the estimates of the relative variance components  $\lambda_A$  and  $\lambda_B$ , corresponding to hour and county, respectively, as well as the

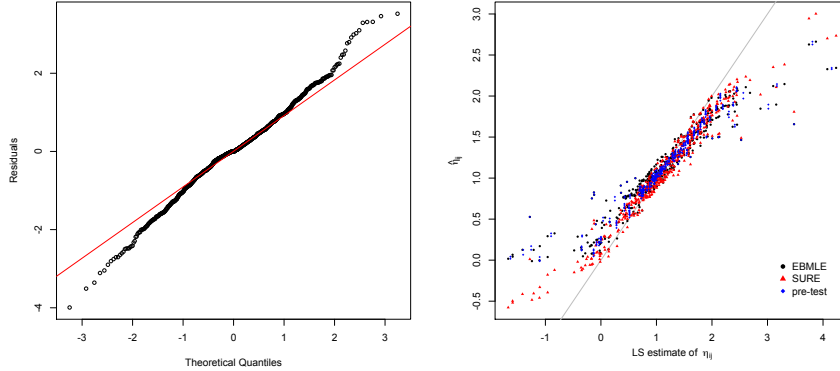


FIG 2. *Left: Normal Q-Q plot for the residuals of the LS fit to the two-way model for water data. Right: Plot of Shrinkage estimates vs. LS estimates of the cell means. The horizontal coordinate is the LS estimate and the vertical coordinate is an alternative estimate: EBMLE, URE or LS-county. EBMLE exhibits most shrinkage. The gray line is the identity line.*

estimates of the fixed term  $\mu$ , for each of the shrinkage estimators. There is a marked difference between the two methods in the estimates of the two variance components. Figure 2 displays fitted values based on the two competing methods, as well as the one-way estimator (LS-county), against the corresponding LS estimate. In terms of shrinkage magnitude, it seems that EBMLE exhibits the most shrinkage among the three, and URE the least among the three, although the differences are not very big. Note that the individual shrinkage patterns could not be immediately anticipated from the values in Table 2 because of the imbalance in the data.

	$\mu$	county	hour
EBMLE	1.10	0.57	0.05
URE	0.78	0.07	0.80

TABLE 2

*Estimated fixed effect ( $\mu$ ) and relative shrinkage factors.*

To compare the performance of the different estimators we carried out two separate analyses. In the first one, we split the data evenly and used the first portion for estimation and the second portion for validation. The second analysis is a data-informed simulation intended to compare performance of the estimators when the additive model (33) is correctly specified.

We begin with comparing the predictive performance against a holdout set. Recall that in the case of missing cells our aim is to estimate the vector



$\boldsymbol{\eta}_c$  of *all estimable* cell means. For a random even split of the data into two subsets  $\mathbf{y}^{(1)}, \mathbf{y}^{(2)}$ , denote by  $\hat{\boldsymbol{\eta}}_c^{(1)}$  an estimate of  $\boldsymbol{\eta}_c$  based on  $\mathbf{y}^{(1)}$  and denote by  $\hat{\boldsymbol{\eta}}_c^{\text{LS}(2)}$  the *Least Squares* estimate of  $\boldsymbol{\eta}_c$  based on  $\mathbf{y}^{(2)}$ . As reflected in notation, we assume that the set of estimable cells is the same for the two portions. Then under (33),  $\hat{\boldsymbol{\eta}}_c^{\text{LS}(2)}$  is an unbiased estimator of  $\boldsymbol{\eta}_c$  and

$$(34) \quad \text{SSPE}[\hat{\boldsymbol{\eta}}_c^{(1)}] = \|\hat{\boldsymbol{\eta}}_c^{(1)} - \hat{\boldsymbol{\eta}}_c^{\text{LS}(2)}\|^2$$

is the Sum of Squared Prediction Error of  $\hat{\boldsymbol{\eta}}_c^{(1)}$ . Instead of averaging (34) directly over random splits, we could use the average of the estimated Total Squared Error

$$\widehat{\text{TSE}}[\hat{\boldsymbol{\eta}}_c^{(1)}] = \text{SSPE}[\hat{\boldsymbol{\eta}}_c^{(1)}] - R(\boldsymbol{\eta}_c, \hat{\boldsymbol{\eta}}_c^{\text{LS}(2)})$$

where for any fixed split  $R(\boldsymbol{\eta}_c, \hat{\boldsymbol{\eta}}_c^{\text{LS}(2)}) = \text{tr}[\text{Cov}(Z_c Z_c^\dagger \hat{\boldsymbol{\eta}}_c^{\text{LS}(2)})]$  and is as an unbiased estimator of the expected risk of  $\hat{\boldsymbol{\eta}}_c^{(2)}$  under a random even split (assuming that  $\hat{\sigma}^2$  is the true variance). Unlike in the other sections we use the un-normalized sum-of-squares loss here, but this will not make any difference because *relative* estimated risks are compared. Note that under (33) the average of  $\|\hat{\boldsymbol{\eta}}_c^{\text{LS}(1)} - \hat{\boldsymbol{\eta}}_c^{\text{LS}(2)}\|^2/2$  over random splits of the data is an unbiased estimator of the expected risk of  $\hat{\boldsymbol{\eta}}_c^{\text{LS}(2)}$  under a random even split; we use it for our calculations in place of  $R(\boldsymbol{\eta}_c, \hat{\boldsymbol{\eta}}_c^{\text{LS}(2)})$  to allow more flexibility in case of departures from the assumed model.

The first row of Table 3 shows the average  $\widehat{\text{TSE}}$  for the two shrinkage estimators and the one-way estimator, as fraction of  $\widehat{\text{TSE}}_{\text{LS}}$ , the average  $\widehat{\text{TSE}}$  for the LS estimator  $\hat{\boldsymbol{\eta}}_c^{\text{LS}(1)}$ , over  $N = 1000$  random splits of the data. We removed from the analysis all counties for which there was a total of less than 8 observations, and recorded the estimated TSEs for each of the  $N$  rounds where the random split resulted in the same set of estimable cells for the two portions of the split. Hence the averages (and standard errors) are based on a slightly smaller effective number of simulation rounds,  $N' = 927$ .

Both shrinkage estimators show significant improvement over LS in terms of estimating the cell means. The EBMLE performs slightly better, with TSE 16% smaller than URE. The estimated relative risk of the one-way estimator is smaller than LS but bigger than the two (empirical) linear shrinkage methods. The pre-test estimator is known to be dominated by a positive-part James-Stein estimator, and, for small values of the parameter, to perform better than the standard (LS) estimator (Sclove *et al.*, 1972); this assumes balanced design, a correctly-specified model, and would entail testing the ‘preliminary’ hypothesis at each round to decide whether to use

the one- or two-way LS; none of these is exactly true of the current analysis, but the outcome of our analysis (also of the simulation analysis, reported next, in which at least misspecification is not a concern) is still in some informal sense consistent with the theoretical results.

	EBMLE	URE	LS-county
validation	<b>0.42</b>	0.5	0.72
simulation	0.81	<b>0.72</b>	0.98

TABLE 3

*Estimated relative TSE for various estimators. The first row of the table corresponds to analysis with validation. The second row corresponds to the data-informed simulation, in which data was simulated according to the additive model. Standard errors are  $< 0.005$ .*

*The URE method seems to perform better under the assumed additive model.*

As the estimators discussed in this paper are designed for the additive model (33), for our second analysis we compare the performance of the different methods (LS, LS-county, EBMLE and URE) when the data is actually generated from the additive model. We set the LS estimate  $\boldsymbol{\eta}^{\text{LS}}$  for the model (33) and the corresponding  $\hat{\sigma}^2$ -based on all 858 observations from all 108 counties – as the “truth”, then draw an independent vector  $\mathbf{y}^* \sim N_n(\boldsymbol{\eta}^{\text{LS}}, \hat{\sigma}^2 I)$ ,  $n = \sum_{i,j} K_{ij}$ , and compute the sum of squared loss  $\|\hat{\boldsymbol{\eta}}_c^* - \boldsymbol{\eta}_c\|^2$  for each estimator  $\hat{\boldsymbol{\eta}}_c^*$ , where the asterisk indicates that the estimate is based on  $\mathbf{y}^*$  only. This process was repeated  $N = 500$  times. The second row of Table 3 shows the estimated risk of the two shrinkage estimators and the one-way estimator as a fraction of the risk of LS. All three estimators have higher risks (relative to LS) compared to the previous analysis, and the URE now has estimated relative risk about 10% smaller than EBMLE. The one-way estimator now barely improves over the standard LS estimator. As both EBMLE and the URE estimators (as well as the pre-test estimator) are designed for the additive model, the results from this analysis might be considered a better basis for comparison between the methods.

**9. Discussion.** We considered estimation under sum-of-squares loss of the cell means in a two-way linear model with additive fixed effects, where the focus was on the unbalanced case. Minimax shrinkage estimators exist which differ from, and hence dominate, the Least Squares estimator for the more general linear regression setup (Rolph, 1976). However, such estimators do not exploit the special structure of the two-factor additive model, and might lead to undesirable shrinkage patterns which are difficult to interpret. Instead, we considered a parametric class of Bayes estimators corresponding to a prior motivated from exchangeability considerations.

The resulting estimates exhibit meaningful shrinkage patterns and, when appropriately calibrated, achieve significant risk reduction as compared to the Least Squares estimator in practical situations.

To calibrate the Bayes estimator we considered substituting the hyperparameters governing the prior with data-dependent values, and proposed a method which chooses these values in an asymptotically optimal way. We contrasted the proposed estimator with the traditional likelihood based empirical BLUP estimator, which was shown to generally produce asymptotically sub-optimal estimates of the cell means. Since it relies on the postulated two-level model, the likelihood based empirical BLUP estimator might be led astray when there is dependency between the cell counts and the true cell means; this was clearly shown in our simulation examples.

The theory developed here employs proof techniques that differ in fundamental aspects from those commonly used to prove asymptotic optimality in the one-way normal mean estimation problem. We offered a flexible approach for proving asymptotic optimality by showing efficient point-wise risk estimation. It greatly helped to tackle the difficulties encountered in two-way problem, where computations involving matrices are generally unavoidable. Our proof techniques can be extended to  $k$ -way additive models, although computational difficulty might become a problem when  $k$  is even moderately large. It would be interesting to investigate whether computationally efficient methods can be developed for the higher-way unbalanced layout.

**10. Acknowledgement.** The authors would like to thank Tony Cai, Samuel Kou and Art Owen for helpful discussions.

## Appendix.

### APPENDIX A: PROOFS OF THE ASYMPTOTIC OPTIMALITY RESULTS OF SECTION ??

Throughout this section we present our proofs assuming  $\sigma = 1$ . It is done mainly for the ease of presentation, and the proofs can easily be modified for any arbitrary but known value of  $\sigma$ . Next we introduce some notation. We denote by  $\sigma_k(A)$  the  $k$ -th largest singular value of a matrix  $A$ . We denote by  $\lambda_k(B)$  the  $k$ -th largest eigenvalue of a symmetric matrix  $B$ . Also, we denote  $G \doteq M\Sigma^{-1}$  and  $H \doteq G^T Q G = \Sigma^{-1} M Q M \Sigma^{-1}$  where  $M$ ,  $\Sigma^{-1}$  and  $Q$  are defined in Section 3. We define  $W = M^{\frac{1}{2}} \Sigma^{-1} M^{\frac{1}{2}}$ . As  $0 \prec M \preceq \Sigma$  we have  $W \preceq I$ , and also  $W^2 \preceq I$ . We will use the following result of [Searle \*et al.\*, 2009](#) (Theorem S4, Page 467).

LEMMA A.1. *Central moments of Gaussian Quadratic Forms.* If  $\mathbf{y} \sim N(\boldsymbol{\eta}, V)$  then:

$$\mathbb{E}(\mathbf{y}^\top A \mathbf{y}) = 2\text{tr}[AV] + \boldsymbol{\eta}^\top A \boldsymbol{\eta}, \quad \text{and} \quad \text{Var}(\mathbf{y}^\top A \mathbf{y}) = 2\text{tr}[(AV)^2] + 4\boldsymbol{\eta}^\top AV A \boldsymbol{\eta}.$$

**A.1. Proof of Theorem 4.1, Lemma 4.1 and Lemma 4.2.**

The proof of Theorem 4.1 for the case when the general effect  $\mu = 0$  follows directly from the results of Lemma 4.1 and Lemma 4.2 as  $\mathbb{E}\{\widehat{\text{URE}}_{r,c}^Q(0, \lambda_A, \lambda_B) - L_{r,c}^Q(\boldsymbol{\eta}, \widehat{\boldsymbol{\eta}}^S(0, \lambda_A, \lambda_B))\}^2$  is bounded above by

$$\begin{aligned} & 2 \mathbb{E}\left\{\widehat{\text{URE}}_{r,c}^Q(0, \lambda_A, \lambda_B) - R_{r,c}^Q(\boldsymbol{\eta}, \widehat{\boldsymbol{\eta}}^S(0, \lambda_A, \lambda_B))\right\}^2 \\ & + 2 \mathbb{E}\left\{L_{r,c}^Q(\boldsymbol{\eta}, \widehat{\boldsymbol{\eta}}^S(0, \lambda_A, \lambda_B)) - R_{r,c}^Q(\boldsymbol{\eta}, \widehat{\boldsymbol{\eta}}^S(0, \lambda_A, \lambda_B))\right\}^2. \end{aligned}$$

In fact, in this case we actually prove Theorem 4.1 with the stronger  $L_2$  norm. We now concentrate on proving the lemmas; we will prove the theorem for the general case later by building on the proofs for the  $\mu = 0$  case.

**Proof of Lemma 4.1.** As the URE is an unbiased estimator of the risk of an estimator in  $\mathcal{S}$ , for any fixed  $\lambda_A, \lambda_B \geq 0$ , we have

$$(35) \quad \mathbb{E}[\widehat{\text{URE}}^Q(0, \lambda_A, \lambda_B) - R_{r,c}^Q(\boldsymbol{\eta}; \widehat{\boldsymbol{\eta}}^S(0, \lambda_A, \lambda_B))]^2 = \text{Var}[\widehat{\text{URE}}^Q(0, \lambda_A, \lambda_B)]$$

Based on the expression of the URE estimator in (17) we know that

$$\widehat{\text{URE}}^Q(0, \lambda_A, \lambda_B) = (rc)^{-2} \{ \sigma^2 \text{tr}(QM) - 2\sigma^2 \text{tr}(\Sigma^{-1}MQM) + \mathbf{y}^\top H \mathbf{y} \},$$

and so the RHS of (35) reduces to  $(rc)^{-2} \text{Var}(\mathbf{y}^\top H \mathbf{y})$  which, being the variance of a quadratic form of the Gaussian random vector  $\mathbf{y}$ , can in turn be evaluated by using Lemma A.1 to give

$$(36) \quad \text{Var}[\widehat{\text{URE}}^Q(0, \lambda_A, \lambda_B)] = (rc)^{-2} \{ 2\text{tr}(HMHM) + 4\boldsymbol{\eta}^\top HMH \boldsymbol{\eta} \}.$$

Our goal now is to show that each of the terms on the RHS, after being multiplied by  $d_{r,c}^2$ , uniformly converges to 0 for all choices of  $\lambda_A$  and  $\lambda_B$ . For this purpose, we concentrate on the second term of the RHS first. As  $H$  is p.s.d. by R2 (See Section S.1.5 of Supplement),  $HMH$  is also p.s.d. Thus,  $\boldsymbol{\eta}^\top HMH \boldsymbol{\eta} \leq \lambda_1(HMH) \|\boldsymbol{\eta}\|^2$ . Next, we bound the largest eigen value of  $HMH$  as

$$\begin{aligned} \lambda_1(HMH) &= \lambda_1(\Sigma^{-1}MQM\Sigma^{-1}M\Sigma^{-1}MQM\Sigma^{-1}) \\ &= \lambda_1(\Sigma^{-1}MQM^{\frac{1}{2}}W^2M^{\frac{1}{2}}QM\Sigma^{-1}) \\ &\leq \lambda_1(\Sigma^{-1}MQMQM\Sigma^{-1}). \end{aligned}$$

The last inequality uses  $W^2 \preceq I$ . Again, by R6 of Supplement Section S.1.5, the RHS above equals  $\lambda_1(M^{\frac{1}{2}}QM\Sigma^{-1}\Sigma^{-1}MQM^{\frac{1}{2}})$ . Thus, we have

$$\begin{aligned}\lambda_1(HMH) &= \lambda_1(M^{\frac{1}{2}}QM\Sigma^{-1}\Sigma^{-1}MQM^{\frac{1}{2}}) \\ &= \lambda_1(M^{\frac{1}{2}}QM^{\frac{1}{2}}WM^{-1}WM^{\frac{1}{2}}QM^{\frac{1}{2}}) \\ &\leq \lambda_1(M^{-1})\lambda_1(M^{\frac{1}{2}}QM^{\frac{1}{2}}) \\ &= \lambda_1(M^{-1})\lambda_1^2(M^{\frac{1}{2}}QM^{\frac{1}{2}}).\end{aligned}$$

The inequality follows by using  $WM^{-1}W \preceq \lambda_1(M^{-1})I$ . Thus, we arrive at the following upper bound

$$(rc)^{-2}d_{r,c}^2 \sup_{\lambda_A, \lambda_B \geq 0} \boldsymbol{\eta}^\top H M H \boldsymbol{\eta} \leq (rc)^{-2} d_{r,c}^2 \lambda_1(M^{-1})\lambda_1^2(M^{\frac{1}{2}}QM^{\frac{1}{2}})\|\boldsymbol{\eta}\|^2$$

which, under Assumptions A1 and A2, converges to 0 as  $r, c \rightarrow \infty$ . Now, for the first term in (35) we have

$$\begin{aligned}\text{tr}(H M H M) &= \text{tr}(M^{\frac{1}{2}}H M H M^{\frac{1}{2}}) \\ &= \text{tr}(M^{\frac{1}{2}}\Sigma^{-1}M Q M \Sigma^{-1}M \Sigma^{-1}M Q M \Sigma^{-1}M^{\frac{1}{2}}) \\ &= \text{tr}(W M^{\frac{1}{2}}Q M^{\frac{1}{2}}W^2 M^{\frac{1}{2}}Q M^{\frac{1}{2}}W) \\ &\leq \text{tr}(W M^{\frac{1}{2}}Q M Q M^{\frac{1}{2}}W) && \text{[using } W^2 \preceq I\text{]} \\ &= \text{tr}(M^{\frac{1}{2}}Q M^{\frac{1}{2}}W^2 M^{\frac{1}{2}}Q M^{\frac{1}{2}}) && \text{[we use R6 here]} \\ &\leq \text{tr}(M^{\frac{1}{2}}Q M Q M^{\frac{1}{2}}) && \text{[again using } W^2 \preceq I\text{]} \\ &\leq (rc) \cdot \lambda_1(M^{\frac{1}{2}}Q M Q M^{\frac{1}{2}}) \\ &= (rc) \cdot \lambda_1^2(M^{\frac{1}{2}}Q M^{\frac{1}{2}}),\end{aligned}$$

where the last equation follows by using R6 again. Thus, as  $r, c \rightarrow \infty$ , by Assumption A2 the first term in (35) scaled by  $d_{r,c}^2$  also converges to 0 uniformly over the ranges of  $\lambda_A$  and  $\lambda_B$ . This completes the proof of the lemma.

**Proof of Lemma 4.2** As the risk is the expectation of the loss, to prove the lemma we need to show:

$$d_{r,c}^2 \sup_{\lambda_A, \lambda_B \geq 0} \text{Var}[L^Q(\boldsymbol{\eta}, \hat{\boldsymbol{\eta}}^S(0, \lambda_A, \lambda_B))] \rightarrow 0 \text{ as } r, c \rightarrow \infty.$$

Again, the loss of the estimator  $\widehat{\boldsymbol{\eta}}_0^S = \widehat{\boldsymbol{\eta}}^S(0, \lambda_A, \lambda_B)$  can be decomposed as

$$\begin{aligned} L^Q(\boldsymbol{\eta}, \widehat{\boldsymbol{\eta}}_0^S) &= (rc)^{-1}(\widehat{\boldsymbol{\eta}}_0^S - \boldsymbol{\eta})^\top Q(\widehat{\boldsymbol{\eta}}_0^S - \boldsymbol{\eta}) = (rc)^{-1}(\mathbf{y} - \boldsymbol{\eta} - G\mathbf{y})^\top Q(\mathbf{y} - \boldsymbol{\eta} - G\mathbf{y}) \\ &= (rc)^{-1}\{(\mathbf{y} - \boldsymbol{\eta})^\top Q(\mathbf{y} - \boldsymbol{\eta}) + \mathbf{y}^\top H\mathbf{y} - 2(\mathbf{y} - \boldsymbol{\eta})^\top QG\mathbf{y}\} \\ &= (rc)^{-1}\{L_1 + L_2 - L_3 + L_4\}, \end{aligned}$$

where  $L_1 = (\mathbf{y} - \boldsymbol{\eta})^\top Q(\mathbf{y} - \boldsymbol{\eta})$ ,  $L_2 = \mathbf{y}^\top H\mathbf{y}$ ,  $L_3 = 2\mathbf{y}^\top QG\mathbf{y}$ ,  $L_4 = 2\boldsymbol{\eta}^\top QG\mathbf{y}$ .

Hence, it suffices to show that  $d_{r,c}^2 \sup_{\lambda_A, \lambda_B} \text{Var}((rc)^{-1}L_i) \rightarrow 0$  as  $r, c \rightarrow \infty$  for all  $i = 1, \dots, 4$ . Uniform convergence of the desired scaled variance of  $L_2$  was already shown in the proof of Lemma 4.1.

For the first term  $L_1$  we have:

$$\begin{aligned} \text{Var}[(rc)^{-1}L_1] &= (rc)^{-2}\text{Var}[(\mathbf{y} - \boldsymbol{\eta})^\top Q(\mathbf{y} - \boldsymbol{\eta})] = (rc)^{-1}2 \text{tr}(QMQM) \\ &= (rc)^{-2}2 \text{tr}\{(M^{\frac{1}{2}}QM^{\frac{1}{2}})^2\} \leq (rc)^{-1}2 \lambda_1^2(M^{\frac{1}{2}}QM^{\frac{1}{2}}) \end{aligned}$$

which by Assumption A2 is  $o(d_{r,c}^{-2})$  as  $r, c \rightarrow \infty$  for any value of the hyperparameter. As  $\mathbf{y}$  is normally distributed, the fourth term can be explicitly evaluated as

$$\begin{aligned} 4^{-1}\text{Var}(L_4) &= \text{Var}(\boldsymbol{\eta}^\top QG\mathbf{y}) = \boldsymbol{\eta}^\top QGMG^\top Q\boldsymbol{\eta} \leq \lambda_1(QGMG^\top Q)\|\boldsymbol{\eta}\|^2 \\ &= \lambda_1(QM\Sigma^{-1}M\Sigma^{-1}MQ)\|\boldsymbol{\eta}\|^2 \\ &= \lambda_1(QM^{\frac{1}{2}}W^2M^{\frac{1}{2}}Q)\|\boldsymbol{\eta}\|^2 \\ &\leq \lambda_1(QM^{\frac{1}{2}}M^{\frac{1}{2}}Q)\|\boldsymbol{\eta}\|^2 \\ &\leq \lambda_1(M^{-1})\lambda_1^2(M^{\frac{1}{2}}QM^{\frac{1}{2}})\|\boldsymbol{\eta}\|^2 \end{aligned}$$

which by assumptions A1-A2 is  $o(r^2c^2d_{r,c}^{-2})$ .

The third term requires detailed analysis. First, note that it breaks into two components

$$(37) \quad \text{Var}[(rc)^{-1}L_3] = 4(rc)^{-2}\text{Var}(\mathbf{y}^\top QG\mathbf{y})$$

$$(38) \quad = 8(rc)^{-2}\text{tr}(\widetilde{G}M\widetilde{G}M) + 16(rc)^{-2}\boldsymbol{\eta}^\top \widetilde{G}M\widetilde{G}\boldsymbol{\eta}$$

where,  $\widetilde{G} = QG + G^\top Q$  is a symmetric matrix. We concentrate on the second term of the RHS first. Note that

$$(rc)^{-2}\boldsymbol{\eta}^\top \widetilde{G}M\widetilde{G}\boldsymbol{\eta} \leq (rc)^{-2}\boldsymbol{\eta}^\top \boldsymbol{\eta} \sigma_1(\widetilde{G}M\widetilde{G}).$$

Like before, if we can uniformly bound the largest eigen value of  $\tilde{G}M\tilde{G}$  as  $o(rcd_{r,c}^{-1})$  then the above is  $o(d_{r,c}^{-2})$  as  $r, c \rightarrow \infty$  by Assumption A1. Noting that  $\tilde{G} = QG + G^\top Q = QM\Sigma^{-1} + \Sigma^{-1}MQ$ , we decompose

$$\begin{aligned}\tilde{G}M\tilde{G} &= H_1 + H_1^\top + H_2 + H_3, \text{ where } H_1 = QM\Sigma^{-1}MQM\Sigma^{-1}, \\ H_2 &= QM\Sigma^{-1}M\Sigma^{-1}MQ, \quad H_3 = \Sigma^{-1}MQMQM\Sigma^{-1}.\end{aligned}$$

To uniform bound the eigen values of  $\tilde{G}M\tilde{G}$  we just show that for each of  $i = 1, \dots, 3$ ,  $(rc)^{-1}d_{r,c}^{-2}\sigma_1(H_i) \rightarrow 0$  as  $r, c \rightarrow \infty$ .  $H_1$  is not a symmetric matrix. In this case, note that:

$$\begin{aligned}\sigma_1(H_1) &= \sigma_1(QM^{\frac{1}{2}}WM^{\frac{1}{2}}QM^{\frac{1}{2}}WM^{-\frac{1}{2}}) \\ &= \sigma_1(M^{-\frac{1}{2}}M^{\frac{1}{2}}QM^{\frac{1}{2}}WM^{\frac{1}{2}}QM^{\frac{1}{2}}WM^{-\frac{1}{2}}) \\ &\leq \lambda_1(M^{-\frac{1}{2}}) \cdot \sigma_1(M^{\frac{1}{2}}QM^{\frac{1}{2}}WM^{\frac{1}{2}}QM^{\frac{1}{2}}W) \cdot \lambda_1(M^{-\frac{1}{2}}) \\ &\leq \lambda_1(M^{-\frac{1}{2}}) \cdot \lambda_1(M^{\frac{1}{2}}QM^{\frac{1}{2}}) \cdot \lambda_1(W) \cdot \lambda_1(M^{\frac{1}{2}}QM^{\frac{1}{2}}) \cdot \lambda_1(W) \cdot \lambda_1(M^{-\frac{1}{2}}) \\ &\leq \lambda_1(M^{-1}) \cdot \lambda_1^2(M^{\frac{1}{2}}QM^{\frac{1}{2}})\end{aligned}$$

where the last inequality uses  $W \preceq I$ . For the symmetric matrix  $H_2$  using  $W^2 \preceq I$ , we have

$$\lambda_1(H_2) = \lambda_1(QM^{\frac{1}{2}}W^2M^{\frac{1}{2}}Q) \leq \lambda_1(QM^{\frac{1}{2}}M^{\frac{1}{2}}Q) \leq \lambda_1(M^{-1})\lambda_1^2(M^{\frac{1}{2}}QM^{\frac{1}{2}})$$

which is uniformly controlled at  $o(rcd_{r,c}^{-2})$  by assumption A2. For the other symmetric matrix  $H_3$  we also have

$$\begin{aligned}\lambda_1(H_3) &= \lambda_1(M^{\frac{1}{2}}QM\Sigma^{-1}\Sigma^{-1}MQM^{\frac{1}{2}}) \\ &= \lambda_1(M^{\frac{1}{2}}QM^{\frac{1}{2}}M^{\frac{1}{2}}\Sigma^{-1}M^{\frac{1}{2}}M^{-1}M^{\frac{1}{2}}\Sigma^{-1}M^{\frac{1}{2}}M^{\frac{1}{2}}QM^{\frac{1}{2}}) \\ &= \lambda_1(M^{\frac{1}{2}}QM^{\frac{1}{2}}WM^{-1}WM^{\frac{1}{2}}QM^{\frac{1}{2}}) \\ &\leq \lambda_1(M^{\frac{1}{2}}QM^{\frac{1}{2}})\lambda_1(W)\lambda_1(M^{-1})\lambda_1(W)\lambda_1(M^{\frac{1}{2}}QM^{\frac{1}{2}}) \\ &\leq \lambda_1^2(M^{\frac{1}{2}}QM^{\frac{1}{2}})\lambda_1(M^{-1})\end{aligned}$$

which again is uniformly controlled at  $o(rcd_{r,c}^{-2})$  by assumption A2.

Now we return to the first term in (37) and upper bound  $\text{tr}(\tilde{G}M\tilde{G}M)$  by  $o(r^2c^2d_{r,c}^{-2})$  when  $r, c \rightarrow \infty$ . Denote  $\dot{G} = Q\dot{G}$  so that  $\tilde{G} = \dot{G} + \dot{G}^\top$ . We have

$$\text{tr}(\tilde{G}M\tilde{G}M) = \text{tr}(\dot{G}M\dot{G}M) + \text{tr}(\dot{G}^\top M\dot{G}^\top M) + 2\text{tr}(\dot{G}M\dot{G}^\top M).$$

Substituting the expression of  $\tilde{G}$  we get

$$\dot{G}M\dot{G}M = QM\Sigma^{-1}MQM\Sigma^{-1}M = QM^{\frac{1}{2}}WM^{\frac{1}{2}}QM^{\frac{1}{2}}QM^{\frac{1}{2}}WM^{\frac{1}{2}},$$

and so we can upper bound its trace as

$$\begin{aligned} \text{tr}(\dot{G}M\dot{G}M) &= \text{tr}(WM^{\frac{1}{2}}QM^{\frac{1}{2}}WM^{\frac{1}{2}}QM^{\frac{1}{2}}) \leq \lambda_1(W)\text{tr}(M^{\frac{1}{2}}QM^{\frac{1}{2}}WM^{\frac{1}{2}}QM^{\frac{1}{2}}) \\ &\leq \text{tr}(M^{\frac{1}{2}}QM^{\frac{1}{2}}M^{\frac{1}{2}}QM^{\frac{1}{2}}) \leq rc \cdot \lambda_1^2(M^{\frac{1}{2}}QM^{\frac{1}{2}}) = o(r^2c^2 d_{r,c}^{-2}) \end{aligned}$$

for any  $\Sigma^{-1}$  and any  $M, Q$  which obeys Assumption A2. Noting that  $\text{tr}(\dot{G}^\top M \dot{G}^\top M) = \text{tr}(M \dot{G} M \dot{G}) = \text{tr}(\dot{G} M \dot{G} M)$ , the second term in (39) is also uniformly bounded by  $o(r^2c^2 d_{r,c}^{-2})$ . Finally, for the third term we have

$$\dot{G}M\dot{G}^\top M = QM\Sigma^{-1}M\Sigma^{-1}MQM = QM^{\frac{1}{2}}W^2M^{\frac{1}{2}}QM \preceq QMQM,$$

and so its trace is upper bounded by

$$\text{tr}(\dot{G}M\dot{G}^\top M) \leq \text{tr}(QMQM) = \text{tr}(M^{\frac{1}{2}}QM^{\frac{1}{2}})^2 = rc \cdot \lambda_1^2(M^{\frac{1}{2}}QM^{\frac{1}{2}}) = o(r^2c^2 d_{r,c}^{-2})$$

by Assumption A2. Thus, we conclude that  $\text{tr}(\tilde{G}M\tilde{G}M)$  is uniformly bounded by  $o(r^2c^2 d_{r,c}^{-2})$  as  $r, c \rightarrow \infty$ . This complete the proof of the lemma.

**Proof of Theorem 4.1 for the general case.** Using the above two lemmas, we now prove our main theorem for the general case. First, note that for arbitrary fixed  $\mu \in \mathbb{R}$ , the loss function decomposes into the following components:

$$\begin{aligned} (\hat{\boldsymbol{\eta}}^S(\mu, \lambda_A, \lambda_B) - \boldsymbol{\eta})^\top Q(\hat{\boldsymbol{\eta}}^S(\mu, \lambda_A, \lambda_B) - \boldsymbol{\eta}) &= (\hat{\boldsymbol{\eta}}^S(0, \lambda_A, \lambda_B) - \boldsymbol{\eta})^\top Q(\hat{\boldsymbol{\eta}}^S(0, \lambda_A, \lambda_B) - \boldsymbol{\eta}) \\ &\quad + \mu^2 \mathbf{1}^\top H \mathbf{1} - 2\mu \mathbf{1}^\top H \boldsymbol{\eta} + 2\mu \mathbf{1}^\top G^\top Q(\boldsymbol{\eta} - \boldsymbol{\eta}). \end{aligned}$$

Comparing it with the definition of  $\widehat{\text{URE}}$  we have:

$$\begin{aligned} \widehat{\text{URE}}_{r,c}^Q(\mu, \lambda_A, \lambda_B) - L_{r,c}^Q(\boldsymbol{\eta}, \hat{\boldsymbol{\eta}}^S(0, \lambda_A, \lambda_B)) &= \widehat{\text{URE}}_{r,c}^Q(0, \lambda_A, \lambda_B) - L_{r,c}^Q(\boldsymbol{\eta}, \hat{\boldsymbol{\eta}}^S(0, \lambda_A, \lambda_B)) \\ &\quad + 2(rc)^{-1} \mu \mathbf{1}^\top G^\top Q(\boldsymbol{\eta} - \boldsymbol{\eta}). \end{aligned}$$

We have already proved the theorem for the case of  $\mu = 0$ ; hence, in light of the above identity, the proof of the general case will follow if we can show:

$$(39) \quad \lim_{\substack{r \rightarrow \infty \\ c \rightarrow \infty}} \sup_{\substack{|\mu| \leq m_{r,c} \\ \lambda_A, \lambda_B \geq 0}} d_{r,c} \cdot (rc)^{-1} \cdot \mathbb{E} |\mu \mathbf{1}^\top G^\top Q(\boldsymbol{\eta} - \boldsymbol{\eta})| = 0.$$

Noting that for any fixed  $\boldsymbol{\eta}$  the random variable  $F = \mathbf{1}^\top G^\top Q(\boldsymbol{\eta} - \boldsymbol{\eta})$  follows a univariate normal distribution with mean 0 and variance  $\mathbf{1}^\top G^\top Q M Q G \mathbf{1}$ , the above holds if

$$(40) \quad \lim_{r \rightarrow \infty, c \rightarrow \infty} m_{r,c} \cdot d_{r,c} \cdot (rc)^{-1} \cdot \sup_{\lambda_A, \lambda_B \geq 0} \{\text{Var}(|\mathbf{1}^\top G^\top Q(\boldsymbol{\eta} - \boldsymbol{\eta})|)\}^{1/2} = 0.$$



Bounding the variance of  $F$  as

$$\begin{aligned}
\text{Var}(\mathbf{1}^\top G^\top Q \mathbf{y}) &\leq \mathbf{1}^\top G^\top Q M Q G \mathbf{1} \\
&= \mathbf{1}^\top \Sigma^{-1} M Q M Q M \Sigma^{-1} \mathbf{1} \\
&\leq \mathbf{1}^\top \mathbf{1} \cdot \lambda_1(\Sigma^{-1} M^{\frac{1}{2}} M^{\frac{1}{2}} Q M^{\frac{1}{2}} M^{\frac{1}{2}} Q M^{\frac{1}{2}} M^{\frac{1}{2}} \Sigma^{-1}) \\
&= rc \cdot \lambda_1(M^{\frac{1}{2}} Q M^{\frac{1}{2}} M^{\frac{1}{2}} \Sigma^{-1} \Sigma^{-1} M^{\frac{1}{2}} M^{\frac{1}{2}} Q M^{\frac{1}{2}}) \\
&\leq rc \cdot \lambda_1(M^{\frac{1}{2}} Q M^{\frac{1}{2}}) \lambda_1(M^{\frac{1}{2}} \Sigma^{-1} \Sigma^{-1} M^{\frac{1}{2}}) \lambda_1(M^{\frac{1}{2}} Q M^{\frac{1}{2}}) \\
&= rc \cdot \lambda_1(M^{\frac{1}{2}} Q M^{\frac{1}{2}}) \lambda_1(W M^{-1} W) \lambda_1(M^{\frac{1}{2}} Q M^{\frac{1}{2}}) \\
&\leq rc \cdot \lambda_1^2(M^{\frac{1}{2}} Q M^{\frac{1}{2}}) \lambda_1(M^{-1}) \leq rc \cdot \lambda_1^2(M^{\frac{1}{2}} Q M^{\frac{1}{2}}) \lambda_1(M^{-1}),
\end{aligned}$$

(40) is proved.

**A.2. Proof of the Decision Theoretic results: Theorems 4.2, 4.3 and Corollary 4.1. Discretization.** In this section, we first define analogous versions of the URE and oracle estimators over a discrete set. Note that in (18) and (21) the URE and oracle estimators involve minimizing the hyper-parameters  $(\mu, \lambda_A, \lambda_B)$  simultaneously over  $\hat{T}_{r,c} = [\hat{a}_\tau, \hat{b}_\tau] \times [0, \infty] \times [0, \infty]$  where the range of the location hyper-parameter  $\mu$  depends on the data. We define a discrete product grid  $\Theta_{r,c} = \Theta_{r,c}^{[1]} \times \Theta_{r,c}^{[2]} \times \Theta_{r,c}^{[3]}$  which only depends on  $r, c$  and not on the data. Details for the construction of  $\Theta_{r,c}$  is provided afterwards. It contains countably infinite grid points as  $r, c \rightarrow \infty$ . We define the discretized version of the oracle estimator where the minimization is conducted over all the points in the discrete grid  $\Theta_{r,c}$  that are contained in  $\hat{T}_{r,c}$ . We define the discretized oracle loss hyper-parameters as

$$(\tilde{\mu}^{\text{OD}}, \tilde{\lambda}_A^{\text{OD}}, \tilde{\lambda}_B^{\text{OD}}) = \arg \min_{(\mu, \lambda_A, \lambda_B) \in \Theta_{r,c} \cap \hat{T}_{r,c}} L_{r,c}^Q(\boldsymbol{\eta}, \hat{\boldsymbol{\eta}}^S(\mu, \lambda_A, \lambda_B)),$$

and the corresponding oracle rule by  $\tilde{\boldsymbol{\eta}}_c^{\text{OD}} = Z_c Z_c^\dagger \hat{\boldsymbol{\eta}}^S(\tilde{\mu}^{\text{OD}}, \tilde{\lambda}_A^{\text{OD}}, \tilde{\lambda}_B^{\text{OD}})$ . We define the URE estimators over the discrete grid by projecting the URE estimates of equation (18) in  $\Theta_{r,c} \cap \hat{T}_{r,c}$ : if the URE hyper-parameters given by Equation (18) are such that:

$$\mu_1 \leq \hat{\mu}^{\text{UQ}} \leq \mu_2, \quad \lambda_1 \leq \hat{\lambda}_A^{\text{UQ}} \leq \lambda_2, \quad \text{and} \quad \lambda_3 \leq \hat{\lambda}_B^{\text{UQ}} \leq \lambda_4$$

where  $\mu_1, \mu_2$  are neighboring points in  $\Theta_{r,c}^{[1]} \cap [\hat{a}_\tau, \hat{b}_\tau]$ ,  $\lambda_1, \lambda_2$  are neighboring points in  $\Theta_{r,c}^{[2]}$  and  $\lambda_3, \lambda_4$  are neighboring points in  $\Theta_{r,c}^{[3]}$ , then the URE

estimates of the tuning parameters over the discrete grid is defined as the minima over the nearest 8-point subset of the grid:

$$(41) \quad (\widehat{\mu}^{\text{UD}}, \widehat{\lambda}_A^{\text{UD}}, \widehat{\lambda}_B^{\text{UD}}) = \underset{(\mu, \lambda_A, \lambda_B) \in \{\mu_1, \mu_2\} \times \{\lambda_1, \lambda_2\} \times \{\lambda_3, \lambda_4\}}{\arg \min} \widehat{\text{URE}}^{\text{Q}}(\mu, \lambda_A, \lambda_B).$$

The corresponding discretized EB estimate is  $\widehat{\boldsymbol{\eta}}^{\text{UD}} = \widehat{\boldsymbol{\eta}}^{\text{S}}(\widehat{\mu}^{\text{UD}}, \widehat{\lambda}_A^{\text{UD}}, \widehat{\lambda}_B^{\text{UD}})$ . The corresponding estimate for  $\boldsymbol{\eta}_c$  is  $\widehat{\boldsymbol{\eta}}_c^{\text{UD}} = Z_c Z_c^\dagger \widehat{\boldsymbol{\eta}}^{\text{S}}(\widehat{\mu}^{\text{D}}, \widehat{\lambda}_A^{\text{D}}, \widehat{\lambda}_B^{\text{D}})$ . If the URE estimators for any of the three hyper-parameters are outside the grid then the nearest boundary of the grid is taken as the UD estimate for that hyper-parameter. We will show afterwards that the probability of such events is negligible. Note that, by construction,  $L(\boldsymbol{\eta}_c, \widehat{\boldsymbol{\eta}}_c^{\text{UD}}) \geq L(\boldsymbol{\eta}_c, \widetilde{\boldsymbol{\eta}}_c^{\text{OD}}) \geq L(\boldsymbol{\eta}_c, \widetilde{\boldsymbol{\eta}}_c^{\text{OL}})$ .

**Construction of the grid  $\Theta_{r,c}$ .** The grid  $\Theta_{r,c}$  is a product grid. The grid  $\Theta_{r,c}^{[1]}$  on the location hyper-parameter  $\mu$  is an equispaced discrete set  $\{-m_{r,c} = \mu[1] < \mu[2] < \dots < \mu[n_1] \leq m_{r,c}\}$  which covers  $[-m_{r,c}, m_{r,c}]$  at a spacing of  $\delta_{r,c}^{[1]}$ . Thus, the cardinality of  $\Theta_{r,c}^{[1]}$ ,  $n_1 = \lceil 2m_{r,c} \{\delta_{r,c}^{[1]}\}^{-1} \rceil$ . We choose the spacing as

$$(42) \quad \delta_{r,c}^{[1]} = \{m_{r,c}^{4/3} \cdot \nu_{r,c} \cdot \lambda_1(Q)\}^{-1}.$$

For constructing the grid  $\Theta_{r,c}^{[2]}$  on the scale hyper-parameter, we consider the following transformation  $\widetilde{\lambda}_A = (1 + \lambda_A)^{-1/2}$ . Note that  $\widetilde{\lambda}_A \in [0, 1]$  as  $\lambda_A$  varies over  $[0, \infty]$ . We construct an equispaced grid on  $\widetilde{\lambda}_A$  between 0 and 1 at a spacing of  $\delta_{r,c}^{[2]}$ :

$$\{0 = \widetilde{\lambda}_A[1] < \widetilde{\lambda}_A[2] < \dots < \widetilde{\lambda}_A[n_2] \leq 1\} \text{ where } \widetilde{\lambda}_A[k] = (k-1)\delta_{r,c}^{[2]} \text{ and } n_2 = \lceil \{\delta_{r,c}^{[2]}\}^{-1} \rceil.$$

The grid on  $\widetilde{\lambda}_A$  is then retransformed to produce the grid  $\Theta_{r,c}^{[2]}$  on the scale hyper-parameter  $\lambda_A$  in the domain  $[0, \infty]$ . The grid  $\Theta_{r,c}^{[3]}$  on  $\lambda_B$  is similarly constructed with  $\delta_{r,c}^{[3]}$  distances between two corresponding grid points in  $\widetilde{\lambda}_B$  scale. The spaces were chosen as:

$$(43) \quad \delta_{r,c}^{[2]} = \delta_{r,c}^{[3]} = \{m_{r,c}^{7/3} \cdot \nu_{r,c} \cdot \lambda_1(Q)\}^{-1}.$$

Now, as  $r, c \rightarrow \infty$ ,  $n_1 = O(m_{r,c}^{7/3} \cdot \nu_{r,c} \cdot \lambda_1(Q))$ ,  $n_2 = O(m_{r,c}^{7/3} \cdot \nu_{r,c} \cdot \lambda_1(Q))$  and thus the cardinality of  $\Theta_{r,c}$  is  $|\Theta_{r,c}| = O(m_{r,c}^7 \nu_{r,c}^3 \lambda_1^3(Q)) = O(d_{r,c})$ .

The following two lemmas enable us to work with the more tractable, discretized versions of the URE and oracle estimators when proving the

decision theoretic results. The first one shows that the difference in the loss between the true estimators and their discretized versions is asymptotically controlled at any prefixed level. The second shows that the URE values for the estimator is also asymptotically close for the discretized version.

LEMMA A.2. *For any fixed  $\epsilon > 0$ , under Assumptions A1-A2,*

- A.  $P\{L_{r,c}(\boldsymbol{\eta}_c, \tilde{\boldsymbol{\eta}}_c^{\text{OD}}) - L_{r,c}(\boldsymbol{\eta}_c, \tilde{\boldsymbol{\eta}}_c^{\text{OL}}) > \epsilon\} \rightarrow 0$  as  $r, c \rightarrow \infty$  and ,
- B.  $\mathbb{E}|L_{r,c}(\boldsymbol{\eta}_c, \tilde{\boldsymbol{\eta}}_c^{\text{OD}}) - L_{r,c}(\boldsymbol{\eta}_c, \tilde{\boldsymbol{\eta}}_c^{\text{OL}})| \rightarrow 0$  as  $r, c \rightarrow \infty$  ,
- C.  $P\{|L_{r,c}(\boldsymbol{\eta}_c, \hat{\boldsymbol{\eta}}_c^{\text{UD}}) - L_{r,c}(\boldsymbol{\eta}_c, \hat{\boldsymbol{\eta}}_c^{\text{URE}})| > \epsilon\} \rightarrow 0$  as  $r, c \rightarrow \infty$  ,
- D.  $\mathbb{E}|L_{r,c}(\boldsymbol{\eta}_c, \hat{\boldsymbol{\eta}}_c^{\text{UD}}) - L_{r,c}(\boldsymbol{\eta}_c, \hat{\boldsymbol{\eta}}_c^{\text{URE}})| \rightarrow 0$  as  $r, c \rightarrow \infty$  .

LEMMA A.3. *For any fixed  $\epsilon > 0$ , under Assumptions A1-A2,*

- A.  $P\{\widehat{\text{URE}}^{\text{Q}}(\hat{\mu}^{\text{UD}}, \hat{\lambda}_A^{\text{UD}}, \hat{\lambda}_B^{\text{UD}}) - \widehat{\text{URE}}^{\text{Q}}(\hat{\mu}^{\text{UQ}}, \hat{\lambda}_A^{\text{UQ}}, \hat{\lambda}_B^{\text{UQ}}) > \epsilon\} \rightarrow 0$  as  $r, c \rightarrow \infty$  and,
- B.  $\mathbb{E}[\widehat{\text{URE}}^{\text{Q}}(\hat{\mu}^{\text{UD}}, \hat{\lambda}_A^{\text{UD}}, \hat{\lambda}_B^{\text{UD}}) - \widehat{\text{URE}}^{\text{Q}}(\hat{\mu}^{\text{UQ}}, \hat{\lambda}_A^{\text{UQ}}, \hat{\lambda}_B^{\text{UQ}})] \rightarrow 0$  as  $r, c \rightarrow \infty$  .

The proof of Lemma A.2 uses the following two lemmas. For shortage of space, the proofs of all these other lemmas (A.2, A.3, A.4 and A.5) is provided in the supplementary materials.

LEMMA A.4. *Under assumption A1 on the parametric space, for any fixed  $\tau \in (0, 1]$ ,  $m_{r,c} = \log(rc)$ , the event  $A_{r,c}(\mathbf{Y}) = \{[\hat{a}_\tau, \hat{b}_\tau] \subseteq [-m_{r,c}, m_{r,c}]\}$  satisfies*

$$P\{A_{r,c}\} \rightarrow 1 \text{ as } n \rightarrow \infty .$$

LEMMA A.5. *Under assumptions A1-A2, for any fixed  $\tau \in (0, 1]$ ,  $m_{r,c} = \log(rc)$ , the event  $A_{r,c}(\mathbf{Y}) = \{[\hat{a}_\tau, \hat{b}_\tau] \subseteq [-m_{r,c}, m_{r,c}]\}$  satisfies:*

- A.  $\mathbb{E}\{|L_{r,c}(\boldsymbol{\eta}_c, \tilde{\boldsymbol{\eta}}_c^{\text{OD}}) - L_{r,c}(\boldsymbol{\eta}_c, \tilde{\boldsymbol{\eta}}_c^{\text{OL}})| \cdot I\{A_{r,c}(\mathbf{Y})\}\} \rightarrow 0$  as  $n \rightarrow \infty$ .
- B.  $\mathbb{E}\{|L_{r,c}(\boldsymbol{\eta}_c, \hat{\boldsymbol{\eta}}_c^{\text{UD}}) - L_{r,c}(\boldsymbol{\eta}_c, \hat{\boldsymbol{\eta}}_c^{\text{URE}})| \cdot I\{A_{r,c}(\mathbf{Y})\}\} \rightarrow 0$  as  $r, c \rightarrow \infty$ .

We next present the proof of the decision theoretic properties where Lemmas A.2, A.3 will be repeatedly used.

**Proof of Theorem 4.2.** We know that

$$\begin{aligned} P\{L(\boldsymbol{\eta}_c, \hat{\boldsymbol{\eta}}_c^{\text{URE}}) \geq L(\boldsymbol{\eta}_c, \tilde{\boldsymbol{\eta}}_c^{\text{OL}}) + \epsilon\} &\leq P\{L(\boldsymbol{\eta}_c, \hat{\boldsymbol{\eta}}_c^{\text{URE}}) \geq L(\boldsymbol{\eta}_c, \tilde{\boldsymbol{\eta}}_c^{\text{OD}}) + \epsilon/2\} \\ &\quad + P\{L(\boldsymbol{\eta}_c, \tilde{\boldsymbol{\eta}}_c^{\text{OD}}) \geq L(\boldsymbol{\eta}_c, \tilde{\boldsymbol{\eta}}_c^{\text{OL}}) + \epsilon/2\}. \end{aligned}$$

The second term converges to 0 by Lemma A.2. The first term is again less than:

$$P\{L(\boldsymbol{\eta}_c, \widehat{\boldsymbol{\eta}}_c^{\text{UD}}) \geq L(\boldsymbol{\eta}, \widetilde{\boldsymbol{\eta}}_c^{\text{OD}}) + \epsilon/4\} + P\{|L(\boldsymbol{\eta}_c, \widehat{\boldsymbol{\eta}}_c^{\text{URE}}) - L(\boldsymbol{\eta}, \widehat{\boldsymbol{\eta}}_c^{\text{UD}})| \leq \epsilon/4\}.$$

The second term in the RHS above converges to 0 as  $r, c \rightarrow \infty$  by Lemma A.2. For the first term note that, by definition,  $\widehat{\text{URE}}^{\text{Q}}(\widehat{\mu}^{\text{UQ}}, \widehat{\lambda}_A^{\text{UQ}}, \widehat{\lambda}_B^{\text{UQ}}) \leq \widehat{\text{URE}}^{\text{Q}}(\widetilde{\mu}^{\text{OD}}, \widetilde{\lambda}_A^{\text{OD}}, \widetilde{\lambda}_B^{\text{OD}})$  which, combined with Lemma A.3, suggests that

$$P\{\widehat{\text{URE}}^{\text{Q}}(\widehat{\mu}^{\text{UD}}, \widehat{\lambda}_A^{\text{UD}}, \widehat{\lambda}_B^{\text{UD}}) \leq \widehat{\text{URE}}^{\text{Q}}(\widetilde{\mu}^{\text{OD}}, \widetilde{\lambda}_A^{\text{OD}}, \widetilde{\lambda}_B^{\text{OD}}) + \epsilon/8\} \rightarrow 0 \text{ as } r, c \rightarrow \infty.$$

Thus, showing  $P\{L(\boldsymbol{\eta}_c, \widehat{\boldsymbol{\eta}}_c^{\text{UD}}) \geq L(\boldsymbol{\eta}, \widetilde{\boldsymbol{\eta}}_c^{\text{OD}}) + \epsilon/4\} \rightarrow 0$  as  $r, c \rightarrow \infty$  can be reduced to showing the following:

$$\lim_{r, c \rightarrow \infty} P\{A(\mathbf{y}; \boldsymbol{\eta}_c) \geq B(\mathbf{y}; \boldsymbol{\eta}_c) + \epsilon/8\} = 0$$

where

$$\begin{aligned} A(\mathbf{y}; \boldsymbol{\eta}_c) &= L(\boldsymbol{\eta}_c, \widehat{\boldsymbol{\eta}}_c^{\text{UD}}) - \widehat{\text{URE}}^{\text{Q}}(\widehat{\mu}^{\text{UD}}, \widehat{\lambda}_A^{\text{UD}}, \widehat{\lambda}_B^{\text{UD}}) \\ B(\mathbf{y}; \boldsymbol{\eta}_c) &= L(\boldsymbol{\eta}_c, \widetilde{\boldsymbol{\eta}}_c^{\text{OD}}) - \widehat{\text{URE}}^{\text{Q}}(\widetilde{\mu}^{\text{OD}}, \widetilde{\lambda}_A^{\text{OD}}, \widetilde{\lambda}_B^{\text{OD}}). \end{aligned}$$

Noting that  $L(\boldsymbol{\eta}_c, \widehat{\boldsymbol{\eta}}_c^{\text{UD}}) = L^{\text{Q}}(\boldsymbol{\eta}, \widehat{\boldsymbol{\eta}}^{\text{UD}})$  and  $L(\boldsymbol{\eta}_c, \widetilde{\boldsymbol{\eta}}_c^{\text{OD}}) = L^{\text{Q}}(\boldsymbol{\eta}, \widehat{\boldsymbol{\eta}}^{\text{S}}(\widetilde{\mu}^{\text{OD}}, \widetilde{\lambda}_A^{\text{OD}}, \widetilde{\lambda}_B^{\text{OD}}))$ , by using Markov's inequality we get

$$P\{A(\mathbf{y}; \boldsymbol{\eta}_c) \geq B(\mathbf{y}; \boldsymbol{\eta}_c) + \epsilon/8\} \leq 8^{-1} \epsilon^{-1} \mathbb{E}\{|A(\mathbf{y}; \boldsymbol{\eta}_c) - B(\mathbf{y}; \boldsymbol{\eta}_c)|\}.$$

By Triangle inequality the RHS above is upper bounded by

$$\begin{aligned} &16 \epsilon^{-1} \mathbb{E}\left\{ \sup_{(\mu, \lambda_A, \lambda_B) \in \Theta_{r,c}} |L^{\text{Q}}(\boldsymbol{\eta}, \widehat{\boldsymbol{\eta}}^{\text{S}}(\mu, \lambda_A, \lambda_B)) - \widehat{\text{URE}}^{\text{Q}}(\mu, \lambda_A, \lambda_B)| \right\} \\ &\leq 16 \epsilon^{-1} \mathbb{E}\left\{ \sum_{(\mu, \lambda_A, \lambda_B) \in \Theta_{r,c}} |L^{\text{Q}}(\boldsymbol{\eta}, \widehat{\boldsymbol{\eta}}^{\text{S}}(\mu, \lambda_A, \lambda_B)) - \widehat{\text{URE}}^{\text{Q}}(\mu, \lambda_A, \lambda_B)| \right\} \\ &\leq 16 \epsilon^{-1} |\Theta_{r,c}| \sup_{\substack{|\mu| \leq m_{r,c} \\ \lambda_A, \lambda_B \geq 0}} \mathbb{E}\left\{ |L^{\text{Q}}(\boldsymbol{\eta}, \widehat{\boldsymbol{\eta}}^{\text{S}}(\mu, \lambda_A, \lambda_B)) - \widehat{\text{URE}}^{\text{Q}}(\mu, \lambda_A, \lambda_B)| \right\}. \end{aligned}$$

As  $|\Theta_{r,c}| = O(d_{r,c})$  by Theorem 4.1, the above expression converges to zero when  $r, c \rightarrow \infty$ . This completes the proof of the theorem.

**Proof of Theorem 4.3.** We decompose the loss into the following three components:

$$\{L(\boldsymbol{\eta}_c, \widehat{\boldsymbol{\eta}}_c^{\text{URE}}) - L(\boldsymbol{\eta}_c, \widehat{\boldsymbol{\eta}}_c^{\text{UD}})\} + \{L(\boldsymbol{\eta}_c, \widetilde{\boldsymbol{\eta}}_c^{\text{OD}}) - L(\boldsymbol{\eta}_c, \widetilde{\boldsymbol{\eta}}_c^{\text{OL}})\} + \{L(\boldsymbol{\eta}_c, \widehat{\boldsymbol{\eta}}_c^{\text{UD}}) - L(\boldsymbol{\eta}_c, \widetilde{\boldsymbol{\eta}}_c^{\text{OD}})\}.$$

By Lemma A.2, the expectation of the absolute value of the first two terms converges to 0 as  $r, c \rightarrow \infty$ . The third term is further decomposed as

$$\begin{aligned} L(\boldsymbol{\eta}_c, \widehat{\boldsymbol{\eta}}_c^{\text{UD}}) - L(\boldsymbol{\eta}_c, \widetilde{\boldsymbol{\eta}}_c^{\text{OD}}) &= \{L(\boldsymbol{\eta}_c, \widehat{\boldsymbol{\eta}}_c^{\text{UD}}) - \widehat{\text{URE}}^Q(\widehat{\boldsymbol{\mu}}^{\text{UD}}, \widehat{\boldsymbol{\lambda}}_A^{\text{UD}}, \widehat{\boldsymbol{\lambda}}_B^{\text{UD}})\} \\ &\quad - \{L(\boldsymbol{\eta}_c, \widetilde{\boldsymbol{\eta}}_c^{\text{OD}}) - \widehat{\text{URE}}^Q(\widetilde{\boldsymbol{\mu}}^{\text{OD}}, \widetilde{\boldsymbol{\lambda}}_A^{\text{OD}}, \widetilde{\boldsymbol{\lambda}}_B^{\text{OD}})\} \\ &\quad + \{\widehat{\text{URE}}^Q(\widehat{\boldsymbol{\mu}}^{\text{UD}}, \widehat{\boldsymbol{\lambda}}_A^{\text{UD}}, \widehat{\boldsymbol{\lambda}}_B^{\text{UD}}) - \widehat{\text{URE}}^Q(\widetilde{\boldsymbol{\mu}}^{\text{OD}}, \widetilde{\boldsymbol{\lambda}}_A^{\text{OD}}, \widetilde{\boldsymbol{\lambda}}_B^{\text{OD}})\}. \end{aligned}$$

By definition  $\widehat{\text{URE}}^Q(\widehat{\boldsymbol{\mu}}^{\text{UD}}, \widehat{\boldsymbol{\lambda}}_A^{\text{UD}}, \widehat{\boldsymbol{\lambda}}_B^{\text{UD}}) \leq \widehat{\text{URE}}^Q(\widetilde{\boldsymbol{\mu}}^{\text{OD}}, \widetilde{\boldsymbol{\lambda}}_A^{\text{OD}}, \widetilde{\boldsymbol{\lambda}}_B^{\text{OD}})$  which, combined with Lemma A.3, suggests that the last term has asymptotically non-positive expectation. Therefore, for all large  $r, c$  values:

$$\begin{aligned} &\mathbb{E}\{L(\boldsymbol{\eta}_c, \widehat{\boldsymbol{\eta}}_c^{\text{UD}}) - L(\boldsymbol{\eta}_c, \widetilde{\boldsymbol{\eta}}_c^{\text{OD}})\} \\ &\leq 2 \mathbb{E}\left\{ \sup_{(\boldsymbol{\mu}, \boldsymbol{\lambda}_A, \boldsymbol{\lambda}_B) \in \Theta_{r,c}} |L^Q(\boldsymbol{\eta}, \widehat{\boldsymbol{\eta}}^S(\boldsymbol{\mu}, \boldsymbol{\lambda}_A, \boldsymbol{\lambda}_B)) - \widehat{\text{URE}}^Q(\boldsymbol{\mu}, \boldsymbol{\lambda}_A, \boldsymbol{\lambda}_B)| \right\} \\ &\leq 2 \mathbb{E}\left\{ \sum_{(\boldsymbol{\mu}, \boldsymbol{\lambda}_A, \boldsymbol{\lambda}_B) \in \Theta_{r,c}} |L^Q(\boldsymbol{\eta}, \widehat{\boldsymbol{\eta}}^S(\boldsymbol{\mu}, \boldsymbol{\lambda}_A, \boldsymbol{\lambda}_B)) - \widehat{\text{URE}}^Q(\boldsymbol{\mu}, \boldsymbol{\lambda}_A, \boldsymbol{\lambda}_B)| \right\} \\ &\leq 2 |\Theta_{r,c}| \sup_{|\boldsymbol{\mu}| \in m_{r,c}; \boldsymbol{\lambda}_A, \boldsymbol{\lambda}_B \geq 0} \mathbb{E}\left\{ |L^Q(\boldsymbol{\eta}, \widehat{\boldsymbol{\eta}}^S(\boldsymbol{\mu}, \boldsymbol{\lambda}_A, \boldsymbol{\lambda}_B)) - \widehat{\text{URE}}^Q(\boldsymbol{\mu}, \boldsymbol{\lambda}_A, \boldsymbol{\lambda}_B)| \right\}. \end{aligned}$$

As  $|\Theta_{r,c}| = O(d_{r,c})$ , the above expression tends to zero when  $r, c \rightarrow \infty$  by Theorem 4.1. This completes the proof of Theorem 4.3.

**Proof of Corollary 4.1.** (a) and (b) are direct consequences, respectively, of Theorems 4.2 and 4.3, since  $L^Q(\boldsymbol{\eta}, \widehat{\boldsymbol{\eta}}^S(\widehat{\boldsymbol{\mu}}, \widehat{\boldsymbol{\lambda}}_A, \widehat{\boldsymbol{\lambda}}_B)) \geq L^Q(\boldsymbol{\eta}, \boldsymbol{\eta}^{\text{OL}})$  and, hence, also  $\mathbb{E}\{L^Q(\boldsymbol{\eta}, \widehat{\boldsymbol{\eta}}^S(\widehat{\boldsymbol{\mu}}, \widehat{\boldsymbol{\lambda}}_A, \widehat{\boldsymbol{\lambda}}_B))\} \geq \mathbb{E}\{L^Q(\boldsymbol{\eta}, \boldsymbol{\eta}^{\text{OL}})\}$ . Unlike in the above two theorems, here we only have optimality over the loss  $L^Q$  defined over the observed cells with  $Q$  in (16). As explained in Section 3, the loss  $L^Q$  for the observed cells is the same as the (normalized) sum-of-squares loss over all (observed and missing)  $rc$  cell means for an estimator of the form  $Z_c Z^\dagger \widehat{\boldsymbol{\eta}}^S(\widehat{\boldsymbol{\mu}}, \widehat{\boldsymbol{\lambda}}_A, \widehat{\boldsymbol{\lambda}}_B)$ , where  $\widehat{\boldsymbol{\mu}}, \widehat{\boldsymbol{\lambda}}_A, \widehat{\boldsymbol{\lambda}}_B$  are any estimates of the hyperparameters.

We end this section by proving the following interesting property of the  $Q$  matrix.

LEMMA A.6. For  $Q$  defined in (15) we have  $\lambda_1(Q) = \lambda_1((Z_c^T Z_c)(Z^T Z)^\dagger)$ . Also,  $\lambda_1(Q) \geq 1$  and  $\lambda_1(Q) = 1$  if  $Z = Z_c$ .

**Proof of Lemma A.6.** By definition (15) we have

$$\lambda_1(Q) = \lambda_1((Z_c Z^\dagger)^T Z_c Z^\dagger) = \lambda_1(Z_c Z^\dagger (Z_c Z^\dagger)^T) = \lambda_1(Z_c (Z^T Z)^\dagger Z_c^T)$$

where the last equality follows as  $Z^\dagger = (Z^T Z)^\dagger Z^T$  and so  $(Z^\dagger)^T Z^\dagger = (Z^T Z)^\dagger$ . Thus we have  $\lambda_1(Q) = \lambda_1((Z_c^T Z_c)(Z^T Z)^\dagger)$ .

If  $Z = Z_c$ , then  $\lambda_1(Q) = \lambda_1((Z_c^T Z_c)(Z_c^T Z_c)^\dagger) = 1$  by definition of Moore-Penrose inverse. We will prove by contradiction that  $\lambda_1(Q) \geq 1$  for any  $Q$  under which  $\boldsymbol{\eta}$  is estimable. If possible assume  $\lambda_1(Q) < 1$  which would imply  $(Z_c^T Z_c)^{1/2} (Z^T Z)^\dagger (Z_c^T Z_c)^{1/2} \prec I$ . Again, as  $\boldsymbol{\eta}$  is estimable,  $\text{rank}(Z_c^T Z_c) = \text{rank}(Z^T Z) = r + c - 1$ . The last two inferences combined suggest that  $\lambda_j(Z^T Z) > \lambda_j(Z_c^T Z_c)$  for some  $j \in \{1, \dots, r + c - 1\}$ . By the Cauchy interlacing theorem, this is a contradiction as  $Z$  was produced by deleting rows of  $Z_c$ , and so  $Z^T Z$  is a compression of  $Z_c^T Z_c$ .

## APPENDIX B: SECTION ?? DETAILS: URE COMPUTATIONS

By definition,  $\Sigma = Z\Lambda\Lambda^T Z^T + M$ . We apply the matrix inverse identity to get

$$(44) \quad \Sigma^{-1} = M^{-1} - M^{-1}Z\Lambda(\Lambda^T Z^T M^{-1}Z\Lambda + I_q)^{-1}\Lambda^T Z^T M^{-1}.$$

Hence, we have

$$\begin{aligned} M\Sigma^{-1} &= I_{rc} - Z\Lambda(\Lambda^T Z^T M^{-1}Z\Lambda + I_q)^{-1}\Lambda^T Z^T M^{-1} \\ M\Sigma^{-1}M &= M - Z\Lambda(\Lambda^T Z^T M^{-1}Z\Lambda + I_q)^{-1}\Lambda^T Z^T. \end{aligned}$$

Using the above, we get

$$\text{tr}(\Sigma^{-1} \mathbf{A} \mathbf{A}^T) = \text{tr}(M\Sigma^{-1}M) = \text{tr}(M) - \text{tr}(Z\Lambda(\Lambda^T Z^T M^{-1}Z\Lambda + I_q)^{-1}\Lambda^T Z^T).$$

Therefore, (10) can be written as

$$\widehat{\text{URE}} = -\sigma^2 \text{tr}(\mathbf{A} \mathbf{A}^T) + 2\sigma^2 \text{tr}\{(\Lambda^T Z^T M^{-1}Z\Lambda + I_q)^{-1}(\Lambda^T Z^T Z\Lambda)\} + \|M\Sigma^{-1}(\mathbf{y} - \mathbf{1}\mu)\|^2.$$

In computing (46):

1. The middle term is computed as the sum of the *elementwise* product of  $(\Lambda^T Z^T M^{-1}Z\Lambda + I_q)^{-1}$  and  $\Lambda^T Z^T Z\Lambda$ , using the property  $\text{tr}(A^T B) = \sum_{i,j} A_{ij} B_{ij}$

2.  $(\Lambda^T Z^T M^{-1} Z \Lambda + I_q)^{-1}$  is computed efficiently employing a sparse Cholesky factorization of  $\Lambda^T Z^T M^{-1} Z \Lambda + I_q$  similarly to the implementation in the lme4 package in R.
3. The quantity  $\min_{\mu} \|M\Sigma^{-1}(\mathbf{y} - \mathbf{1}\mu)\|^2$  is computed by regressing  $M\Sigma^{-1}\mathbf{y}$  on  $M\Sigma^{-1}\mathbf{1}_{rc}$  using the `lm` function in R. In doing that, the vector  $M\Sigma^{-1}\mathbf{x}$  (for  $\mathbf{x} = \mathbf{y}$  and  $\mathbf{x} = \mathbf{1}_{rc}$ ) is computed as:

$$(47) \quad M\Sigma^{-1}\mathbf{x} = \mathbf{x} - Z\Lambda(\Lambda^T Z^T M^{-1} Z \Lambda + I_q)^{-1} \Lambda^T Z^T (M^{-1}\mathbf{x})$$

where (47) is implemented proceeding “from right to left” to always compute a product of a matrix and a *vector*, instead of two matrices: First find  $M^{-1}\mathbf{x}$ , then find  $(\Lambda^T Z^T)(M^{-1}\mathbf{x})$ , and so on.

## APPENDIX C: SUPPLEMENTARY MATERIALS

Further derivations and discussions of the results of Sections 2, 3 and 4 are provided in the supplementary materials (Brown *et al.*, 2016).

## REFERENCES

- Bates, D., Maechler, M., Bolker, B., and Walker, S. (2014). *lme4: Linear mixed-effects models using Eigen and S4*. R package version 1.1-7.
- Bates, D. M. (2010). lme4: Mixed-effects modeling with r. <http://lme4.r-forge.r-project.org/book>.
- Brown, L., Mukherjee, G., and A, W. (2016). Supplement to “empirical bayes estimates for a 2-way cross-classified additive model” .
- Candes, E., Sing-Long, C. A., and Trzasko, J. D. (2013). Unbiased risk estimates for singular value thresholding and spectral estimators. *Signal Processing, IEEE Transactions on* **61**, 19, 4643–4657.
- Dey, A. (1986). *Theory of block designs*. J. Wiley.
- Dicker, L. H. (2013). Optimal equivariant prediction for high-dimensional linear models with arbitrary predictor covariance. *Electronic Journal of Statistics* **7**, 1806–1834.
- Donoho, D. L., Johnstone, I. M., Kerkyacharian, G., and Picard, D. (1995). Wavelet shrinkage: asymptopia? *Journal of the Royal Statistical Society. Series B (Methodological)* 301–369.
- Draper, N. R. and Van Nostrand, R. C. (1979). Ridge regression and james-stein estimation: review and comments. *Technometrics* **21**, 4, 451–466.
- Efron, B. and Morris, C. (1972). Empirical bayes on vector observations – an extension of stein’s method. *Biometrika* **59**, 2, 335–347.
- Efron, B. and Morris, C. (1973). Stein’s estimation rule and its competitors: an empirical bayes approach. *Journal of the American Statistical Association* **68**, 341, 117–130.
- Ghosh, M., Nickerson, D. M., and Sen, P. K. (1987). Sequential shrinkage estimation. *The Annals of Statistics* 817–829.
- Goldstein, H., Browne, W., and Rasbash, J. (2002). Multilevel modelling of medical data. *Statistics in medicine* **21**, 21, 3291–3315.
- Henderson, C. (1984). Anova, mivque, reml, and ml algorithms for estimation of variances and covariances. In *Statistics: An Appraisal: Proceedings 50th Anniversary Conference (David HA, David HT, eds)*, The Iowa State University Press, Ames, IA, 257–280.

- James, W. and Stein, C. (1961). Estimation with quadratic loss. In *Proceedings of the fourth Berkeley symposium on mathematical statistics and probability*, vol. 1, 361–379.
- Jiang, J., Nguyen, T., and Rao, J. S. (2011). Best predictive small area estimation. *Journal of the American Statistical Association* **106**, 494, 732–745.
- Johnstone, I. M. (2011). Gaussian estimation: Sequence and wavelet models. *Unpublished manuscript*.
- Johnstone, I. M. and Silverman, B. W. (2004). Needles and straw in haystacks: Empirical bayes estimates of possibly sparse sequences. *Annals of Statistics* 1594–1649.
- Kou, S. and Yang, J. J. (2015). Optimal shrinkage estimation in heteroscedastic hierarchical linear models. *arXiv preprint arXiv:1503.06262*.
- Li, K.-C. (1986). Asymptotic optimality of  $cl$  and generalized cross-validation in ridge regression with application to spline smoothing. *The Annals of Statistics* 1101–1112.
- Lindley, D. (1962). Discussion of the paper by stein. *J. Roy. Statist. Soc. Ser. B* **24**, 265–296.
- Lindley, D. V. and Smith, A. F. (1972). Bayes estimates for the linear model. *Journal of the Royal Statistical Society. Series B (Methodological)* 1–41.
- Mason, W. M., Wong, G. Y., and Entwisle, B. (1983). Contextual analysis through the multilevel linear model. *Sociological methodology* **1984**, 72–103.
- Oman, S. D. (1982). Shrinking towards subspaces in multiple linear regression. *Technometrics* **24**, 4, 307–311.
- Rasbash, J. and Goldstein, H. (1994). Efficient analysis of mixed hierarchical and cross-classified random structures using a multilevel model. *Journal of Educational and Behavioral statistics* **19**, 4, 337–350.
- Rolph, J. E. (1976). Choosing shrinkage estimators for regression problems. *Communications in Statistics-Theory and Methods* **5**, 9, 789–802.
- Slove, S. L. (1968). Improved estimators for coefficients in linear regression. *Journal of the American Statistical Association* **63**, 322, 596–606.
- Slove, S. L., Morris, C., and Radhakrishnan, R. (1972). Non-optimality of preliminary-test estimators for the mean of a multivariate normal distribution. *The Annals of Mathematical Statistics* 1481–1490.
- Searle, S. R., Casella, G., and McCulloch, C. E. (2009). *Variance components*, vol. 391. John Wiley & Sons.
- Searle, S. R. and McCulloch, C. E. (2001). *Generalized, linear and mixed models*. Wiley.
- Stein, C. M. (1962). Confidence sets for the mean of a multivariate normal distribution. *Journal of the Royal Statistical Society. Series B (Methodological)* 265–296.
- Tan, Z. (2014). Steinized empirical bayes estimation for heteroscedastic data. *Statistica Sinica*, to appear.
- Xie, X., Kou, S., and Brown, L. D. (2012). Sure estimates for a heteroscedastic hierarchical model. *Journal of the American Statistical Association* **107**, 500, 1465–1479.
- Zaccarin, S. and Rivellini, G. (2002). Multilevel analysis in social research: an application of a cross-classified model. *Statistical Methods and Applications* **11**, 1, 95–108.

ADDRESS OF THE FIRST AUTHOR:  
 DEPARTMENT OF STATISTICS,  
 UNIVERSITY OF PENNSYLVANIA,  
 400 JON M. HUNTSMAN HALL,  
 3730 WALNUT STREET,  
 PHILADELPHIA, PA 19104,  
 E-MAIL: [lbrown@wharton.upenn.edu](mailto:lbrown@wharton.upenn.edu)

ADDRESS OF THE SECOND AUTHOR:  
 DEPARTMENT OF DATA SCIENCES AND OPERATIONS,  
 MARSHALL SCHOOL OF BUSINESS,  
 UNIVERSITY OF SOUTHERN CALIFORNIA,  
 LOS ANGELES, CA 90089-0809,  
 E-MAIL: [gourab@usc.edu](mailto:gourab@usc.edu)



ADDRESS OF THE THIRD AUTHOR:  
DEPARTMENT OF STATISTICS,  
SEQUOIA HALL, 390 SERRA MALL,  
STANFORD UNIVERSITY,  
STANFORD, CA 94305-4065,  
E-MAIL: [asafw.at.stanford@gmail.com](mailto:asafw.at.stanford@gmail.com)