

## EMPIRICAL BAYES ESTIMATES FOR A TWO-WAY CROSS-CLASSIFIED MODEL

BY LAWRENCE D. BROWN<sup>1</sup>, GOURAB MUKHERJEE<sup>2</sup>  
AND ASAF WEINSTEIN<sup>1</sup>

*University of Pennsylvania, University of Southern California and  
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 [*J. R. Stat. Soc., B* **34** (1972) 1–41]. 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, under suitable conditions. 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 suboptimally. 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.

**1. Introduction.** Multilevel cross-classified models are pervasive in statistics, with applications ranging from detecting sources of variability in medical research [Goldstein, Browne and Rasbash (2002)] to understanding micro–macro linkages in social studies [Mason, Wong and Entwisle (1983), Zaccarin and Rivellini (2002)]. These models offer a natural and flexible approach to specify meaningful 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 [Lindley (1962), Stein (1962)], the usefulness of such models in estimation problems involving multiple *nonrandom* effects has been well recognized.

---

Received May 2016; revised February 2017.

<sup>1</sup>Supported in part by grant NSF DMS-15-12084.

<sup>2</sup>Supported in part by the Zumberge individual award from the University of Southern California’s James H. Zumberge faculty research and innovation fund.

*MSC2010 subject classifications.* Primary 62C12; secondary 62C25, 62F10, 62J07.

*Key words and phrases.* Shrinkage estimation, empirical Bayes, two-way ANOVA, oracle optimality, Stein’s unbiased risk estimate (SURE), empirical BLUP.

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 *toward* a subspace [e.g., Sclove (1968), Oman (1982), Jiang, Nguyen and Rao (2011), Tan (2016)] or *within* a subspace [e.g., Lindley and Smith (1972), Rolph (1976), Kou and Yang (2017)]; 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 hyperparameters governing the prior. Any corresponding estimator that substitutes *data-dependent* values for the hyperparameters 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, so-called empirical BLUP (Best Linear Unbiased Predictors) estimators, using the usual maximum-likelihood approach in estimating the hyperparameters, are shown to perform suboptimally in the unbalanced case. Instead of using the maximum-likelihood criterion, we choose the values for the hyperparameters 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, Nickerson and Sen (1987), Donoho et al. (1995), Johnstone and Silverman (2004), Candès, Sing-Long and Trzasko (2013), to name a few]. Recently, Xie, Kou and Brown (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, Kou and Brown \(2012\)](#) 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 fact that the two-way case imposes structure on the mean vector, which is nontrivial to handle due to missingness or 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 computation 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 \(2017\)](#) 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 hyperparameter. 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: we begin with the all-cells-filled situation, and then present a more general model which allows missing observations. In Section 3, we show that our proposed estimation methodology is asymptotically optimal and is capable of recovering the directions and magnitude for optimal shrinkage. In Section 4, we report the results from extensive numerical experiments, and also 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 Bayes estimators.* Consider the following basic two-way cross-classified additive model with fixed effects:

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

$K_{ij}$  is the count in the  $(i, j)$ th cell;  $\sigma^2 > 0$  is assumed to be known; and  $\varepsilon_{ij}$  is an independent Gaussian noise term, which can be considered as the average of  $K_{ij}$  homoscedastic error terms. Throughout the paper, we write vectors in bold and matrices in capital letters. 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, but the vector of cell means  $\boldsymbol{\eta} = (\eta_{11}, \eta_{12}, \dots, \eta_{rc})^\top$  always is. Our goal is to estimate  $\boldsymbol{\eta}$  under 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),  $\alpha_i$  and  $\beta_j$  will be referred to as the  $i$ th row and the  $j$ th column 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 empty cells, is presented in Section 2.3. We would like to emphasize the focus in this section on the loss (2) rather than the weighted quadratic loss  $L_{r,c}^{\text{wgt}}(\boldsymbol{\eta}, \hat{\boldsymbol{\eta}}) = (rc)^{-1} \sum_{i=1}^r \sum_{j=1}^c K_{ij} (\hat{\eta}_{ij} - \eta_{ij})^2$ , 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.

The usual estimator of  $\boldsymbol{\eta}$  is the weighted least squares (WLS) estimator, also the maximum-likelihood estimator under (1). The WLS estimator is unbiased and minimax but can be substantially improved on by shrinkage estimators, particularly when  $r, c \rightarrow \infty$  [Draper and Van Nostrand (1979)]. As the starting point for the shrinkage estimators proposed in this paper, we consider a family of Bayes estimators with respect to the conjugate prior

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

where  $\sigma_A^2, \sigma_B^2$  are hyperparameters. 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.<sup>3</sup> Employing this prior is standard in Bayesian Analysis-of-Variance [e.g., Gelman et al. (2004), Chapter 15.6; Gelman (2005)]. In vector form, the two-level hierarchical model can be written as

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

<sup>3</sup>As in Lindley and Smith (1972) and Gelman (2005), we refer to within-batch effects in our model as exchangeable, although they are really i.i.d.

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  $Z_A = I_r \otimes 1_c$  and  $Z_B = 1_r \otimes I_c$ . 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 hyperparameters will be kept implicit. The marginal covariance of  $\mathbf{y}$  in (5) is  $\sigma^2 \Sigma$  where

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

The following is a standard result and is proved in Section S.2 of the supplementary materials [Brown, Mukherjee and Weinstein (2018)].

LEMMA 2.1. *For any fixed  $\mu \in \mathbb{R}$  and nonnegative  $\lambda_A, \lambda_B$ , 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  $\Sigma$  depends on the hyperparameters  $\lambda_A, \lambda_B$  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 obtained by fixing the values of the hyperparameters in (5),

$$(6) \quad \begin{aligned} \mathcal{S}(a, b) &= \{\hat{\boldsymbol{\eta}}(\mu, \lambda_A, \lambda_B) = \mathbf{y} - M \Sigma^{-1}(\mathbf{y} - \mathbf{1}\mu) : \mu \in [a, b], \\ &\lambda_A \geq 0, \lambda_B \geq 0\} \end{aligned}$$

for  $a, b \in \mathbb{R}$ . Note that above,  $\mu$  is restricted to lie in  $[a, b]$ . In practice, we will consider only estimators in  $\mathcal{S}[\tau] := \mathcal{S}(\hat{a}_\tau, \hat{b}_\tau)$  where  $\hat{a}_\tau(\mathbf{y}) = \text{quantile}\{y_{ij} : 1 \leq i \leq r, 1 \leq j \leq c; \tau/2\}$  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 hyperparameter  $\mu$  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 nonrobust; it also seems as 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 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. *Empirical Bayes estimators.* 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 typically estimated by generalized least squares, producing an empirical version of what is known as the Best Linear Unbiased Predictor (BLUP). There is extensive literature on the estimation of the variance components. For a textbook account, see, for example, Chapters 5 and 6 of [Searle, Casella and McCulloch \(1992\)](#); and Section 7.3.2 in [Fahrmeir et al. \(2013\)](#), who also discuss Bayesian inference in Section 7.4.2. Computational issues are discussed in [Harville \(1977\)](#). The main methods for estimation of variance components in linear mixed models are 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 \(2010\)](#)]. 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  $\hat{\eta}^{\text{ML}} = \hat{\boldsymbol{\eta}}(\hat{\mu}^{\text{ML}}, \hat{\lambda}_A^{\text{ML}}, \hat{\lambda}_B^{\text{ML}})$  and will be referred to as EBMLE (for Empirical Bayes Maximum-Likelihood).

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

$$(8) \quad \hat{\mu} = \hat{\mu}_1 \cdot I\{\hat{\mu}_1 \in [\hat{a}_\tau, \hat{b}_\tau]\} + \hat{a}_\tau \cdot I\{\hat{\mu}_1 < \hat{a}_\tau\} + \hat{b}_\tau \cdot I\{\hat{\mu}_1 > \hat{b}_\tau\},$$

where  $\hat{\mu}_1 = (\mathbf{1}^\top \hat{\Sigma}^{-1} \mathbf{y}) / (\mathbf{1}^\top \hat{\Sigma}^{-1} \mathbf{1})$ ; if  $\hat{\mu}_1 \in [\hat{a}_\tau, \hat{b}_\tau]$  and if  $\hat{\lambda}_a, \hat{\lambda}_b$  are both strictly positive, they satisfy

$$(9) \quad \begin{aligned} \text{tr}(\hat{\Sigma}^{-1} Z_A Z_A^\top) - \sigma^{-2} \mathbf{y}^\top (I - \hat{P})^\top \hat{\Sigma}^{-1} Z_A Z_A^\top \hat{\Sigma}^{-1} (I - \hat{P}) \mathbf{y} &= 0, \\ \text{tr}(\hat{\Sigma}^{-1} Z_B Z_B^\top) - \sigma^{-2} \mathbf{y}^\top (I - \hat{P})^\top \hat{\Sigma}^{-1} Z_B Z_B^\top \hat{\Sigma}^{-1} (I - \hat{P}) \mathbf{y} &= 0, \end{aligned}$$

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.2.1 of the supplement, which also contains the estimating equation for the case when  $\hat{\mu}_1 \notin [\hat{a}_\tau, \hat{b}_\tau]$ . If the solution to the estimating equations (9) includes a negative component, adjustments are needed in order to produce the maximum-likelihood estimates of the scale hyperparameters [see [McCulloch and Searle \(2001\)](#), Section 2.2b-iii for a discussion of the one-way case].

We propose an alternative method for estimating the shrinkage parameters. Following the approach of [Xie, Kou and Brown \(2012\)](#), for fixed  $\tau \in (0, 1]$  we choose

the shrinkage parameters by minimizing unbiased risk estimate (URE) over estimators  $\hat{\boldsymbol{\eta}}(\boldsymbol{\mu}, \lambda_A, \lambda_B)$  in  $\mathcal{S}(\tau)$ . By Lemma S.2.2 of the supplement, an unbiased estimate of the risk,  $R_{r,c}(\boldsymbol{\eta}, \hat{\boldsymbol{\eta}}(\boldsymbol{\mu}, \lambda_A, \lambda_B)) \triangleq (rc)^{-1} \mathbb{E} \|\hat{\boldsymbol{\eta}}(\boldsymbol{\mu}, \lambda_A, \lambda_B) - \boldsymbol{\eta}\|^2$  is given by  $\text{URE}(\boldsymbol{\mu}, \lambda_A, \lambda_B)$ , which equals

$$(10) \quad \frac{1}{rc} \{ \sigma^2 \text{tr}(M) - 2\sigma^2 \text{tr}(\Sigma^{-1} M^2) + (\mathbf{y} - \mathbf{1}\boldsymbol{\mu})^\top [\Sigma^{-1} M^2 \Sigma^{-1}] (\mathbf{y} - \mathbf{1}\boldsymbol{\mu}) \}.$$

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

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

The corresponding empirical Bayes estimator is  $\hat{\boldsymbol{\eta}}^{\text{URE}} = \hat{\boldsymbol{\eta}}(\hat{\boldsymbol{\mu}}^U, \hat{\lambda}_A^U, \hat{\lambda}_B^U)$ . 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. *URE based estimates of (11) satisfy the following estimating equations:*

$$(12) \quad \hat{\boldsymbol{\mu}} = \hat{\boldsymbol{\mu}}_1 \cdot I\{\hat{\boldsymbol{\mu}}_1 \in [\hat{a}_\tau, \hat{b}_\tau]\} + \hat{a}_\tau \cdot I\{\hat{\boldsymbol{\mu}}_1 < \hat{a}_\tau\} + \hat{b}_\tau \cdot I\{\hat{\boldsymbol{\mu}}_1 > \hat{b}_\tau\},$$

where  $\hat{\boldsymbol{\mu}}_1 = (\mathbf{1}^\top [\hat{\Sigma}^{-1} M^2 \hat{\Sigma}^{-1}] \mathbf{y}) / (\mathbf{1}^\top [\hat{\Sigma}^{-1} M^2 \hat{\Sigma}^{-1}] \mathbf{1})$ ; if  $\hat{\boldsymbol{\mu}}_1 \in [\hat{a}_\tau, \hat{b}_\tau]$  and if  $\hat{\lambda}_a, \hat{\lambda}_b$  are both strictly positive, they satisfy

$$(13) \quad \begin{aligned} & \text{tr}(\hat{\Sigma}^{-1} Z_A Z_A^\top \hat{\Sigma}^{-1} M^2) \\ & \quad - \sigma^{-2} \mathbf{y}^\top (I - \hat{P})^\top \hat{\Sigma}^{-1} Z_A Z_A^\top \hat{\Sigma}^{-1} M^2 \hat{\Sigma}^{-1} (I - \hat{P}) \mathbf{y} = 0, \\ & \text{tr}(\hat{\Sigma}^{-1} Z_B Z_B^\top \hat{\Sigma}^{-1} M^2) \\ & \quad - \sigma^{-2} \mathbf{y}^\top (I - \hat{P})^\top \hat{\Sigma}^{-1} Z_B Z_B^\top \hat{\Sigma}^{-1} M^2 \hat{\Sigma}^{-1} (I - \hat{P}) \mathbf{y} = 0, \end{aligned}$$

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.2.1 of the supplementary materials. Comparing the two systems of equations (9) and (13) without substituting the value of  $\boldsymbol{\mu}$ , 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 (see Section S.5 of the supplement for a discussion of the balanced case). In Section 3, 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 4 demonstrate that in certain situations EBMLE performs significantly worse.

REMARK 1. 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 in general. 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 nontrivial 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 URE is not necessarily convex in  $(\mu, \lambda_A, \lambda_B)$ . To bypass these issues, we minimize URE by conducting a grid search on  $(\lambda_A, \lambda_B)$ , and  $\mu$  is subsequently estimated by (12). For efficiency in handling large data sets, our implementation adopts some of the key computational elements from the `lme4` package [Section 5.4, Bates (2010)]; Details are provided in Section A.3 of the Appendix.

REMARK 2. 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 so-called 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 methodology we develop here incorporates the exchangeable structure in (3).

2.3. *A model with missing cells.* A more general model than (1) allows some cells to be empty. Hence, consider

$$(14) \quad y_{ij} = \eta_{ij} + \varepsilon_{ij}, \quad \eta_{ij} = \mu + \alpha_i + \beta_j, \varepsilon_{ij} \sim N(0, \sigma^2 K_{ij}^{-1}), (i, j) \in \mathcal{E},$$

where  $\mathcal{E} = \{(i, j) : K_{ij} \geq 1\} \subseteq \{1, \dots, r\} \times \{1, \dots, c\}$  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<sup>4</sup> 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 a 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.2 of the supplementary materials.

LEMMA 2.4. *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

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

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  $\mathcal{S}[\tau] := \mathcal{S}(\hat{a}_\tau, \hat{b}_\tau)$ , with  $\mathcal{S}(a, b)$  defined in (6). The EBMLE estimates the hyperparameters  $\mu, \lambda_A, \lambda_B$  based on the marginal likelihood  $y$  according to (3), where  $M, \Sigma^{-1}$  are as defined in the current

---

<sup>4</sup>A design is disconnected if there is a partition of the row effects into nonempty subsets such that any two members of distinct subsets never appear with the same column effect. A design is called connected if it is not disconnected.

section. As shown in Lemma S.2.3 of the supplement, an unbiased estimator of the point risk corresponding to (15),

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

is given by  $\text{URE}^Q(\mu, \lambda_A, \lambda_B)$  which is evaluated as

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

The URE estimates of the tuning parameters are

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

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

**3. Risk properties and asymptotic optimality of the URE based estimator.**

We now present the results that establish the optimality properties of our proposed URE-based estimator. We present the result for the general quadratic loss  $L_{r,c}^Q$  of the previous section with the matrix  $Q$  defined in (16). Replacing  $Q$  with  $I_{rc}$  will give the results for the all-cells-filled model (1), which are explained separately. In proving our theoretical results, we make the following assumptions:

A1. *On the parameter space:* We assume that the parameter  $\boldsymbol{\eta}_c$  in the complete model is estimable and satisfies the following second-order moment condition:

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

This assumption is very mild, and similar versions are widely used in the EB literature [see Assumption  $C'$  of Xie, Kou and Brown (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

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

where  $\nu_{r,c} = \max\{K_{ij} : (i, j) \in \mathcal{E}\} / \min\{K_{ij} : (i, j) \in \mathcal{E}\}$ . As shown in Lemma S.3.1 of the supplement,  $\lambda_1(Q)$  equals the largest eigenvalue of  $(Z_c^\top Z_c)(Z^\top Z)^\dagger$ . 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,  $\lambda_1((Z_c^\top Z_c)(Z^\top Z)^\dagger)$  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 S.3.1). Thus, in that case the assumption reduces to  $\lim_{r,c \rightarrow \infty} (rc)^{-1/8} (\log(rc))^2 v_{r,c} = 0$ . 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 slower rate than  $(rc)^{1/8} / (\log(rc))^2$ . This assumption on the design matrix is essential for our asymptotic optimality proofs. Section A.1 of the Appendix shows its role in our proofs and a detailed discussion about it is provided in the supplementary materials.

*Asymptotic results.* Our decision theoretic optimality results depend on the following pointwise approximation of the true loss of estimators in  $\mathcal{S}[\tau]$  by the URE methodology. Consider the following design-dependent quantities:

$$(19) \quad d_{r,c} = m_{r,c}^7 v_{r,c}^3 \lambda_1^3(Q), \quad m_{r,c} = \log(rc).$$

Note that, as  $r, c \rightarrow \infty$ , by Assumption A2,  $d_{r,c} v_{r,c} \lambda_1(Q) = o(\sqrt{rc})$  which will be used afterward in the proofs for bounding the aforementioned pointwise loss approximation error. The following theorem shows that the unbiased risk estimator approximates the true loss at an asymptotic  $L_1$  error rate smaller than  $d_{r,c}^{-1}$  for every hyperparameter value in the set where  $\lambda_A, \lambda_B$  is nonnegative and the location hyperparameter  $\mu$  is restricted to the interval  $[-m_{r,c}, m_{r,c}]$ , which grows as  $r, c$  increase.

**THEOREM 3.1.** *Under Assumptions A1–A2, we have*

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

The set of all hyperparameters considered in  $\mathcal{S}[\tau]$  differs from hyperparameter set considered in the above theorem, 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 3.1 asymptotically covers all hyperparameters considered in  $\mathcal{S}[\tau]$  for any  $\tau \in (0, 1]$ . This explains intuitively why in choosing the hyperparameters by minimizing an unbiased risk estimate as in (17), we can expect the resulting estimate  $\hat{\boldsymbol{\eta}}(\hat{\mu}^{\text{UQ}}, \hat{\lambda}_A^{\text{UQ}}, \hat{\lambda}_B^{\text{UQ}})$  to have competitive performance.

*Decision theoretic optimality.* To compare the asymptotic performance of our proposed estimate, we define the oracle loss (OL) hyperparameter as

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

and the corresponding oracle rule

$$(20) \quad \tilde{\eta}_c^{\text{OL}} = Z_c Z_c^\dagger \hat{\eta}(\tilde{\mu}^{\text{OL}}, \tilde{\lambda}_A^{\text{OL}}, \tilde{\lambda}_B^{\text{OL}}).$$

Note that the oracle rule depends on the unknown cell means  $\eta_c$  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  $\eta_c^{\text{OL}}$ . 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 2.1 and 2.3, both Theorems 3.2 and 3.3 simultaneously cover the missing and all-cells-filled model.

**THEOREM 3.2.** *Under Assumptions A1–A2, for any  $\varepsilon > 0$  we have*

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

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 3.3.** *Under Assumptions A1–A2, the following holds:*

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

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

- (a)  $\lim_{r \rightarrow \infty, c \rightarrow \infty} P\{L_{r,c}^Q(\eta, \hat{\eta}^{\text{URE}}) \geq L_{r,c}^Q(\eta, \hat{\eta}(\hat{\mu}, \hat{\lambda}_A, \hat{\lambda}_B)) + \varepsilon\} = 0.$
- (b)  $\limsup_{r \rightarrow \infty, c \rightarrow \infty} R_{r,c}^Q(\eta, \hat{\eta}^{\text{URE}}) - R_{r,c}^Q(\eta, \hat{\eta}(\hat{\mu}, \hat{\lambda}_A, \hat{\lambda}_B)) \leq 0.$

Unlike the above two theorems, this corollary is based on the quadratic loss  $L^Q$ . 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 (3.1), 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  $\eta$  would have to be very close to the URE estimate).

*Proof Overview.* The pointwise loss estimation result of Theorem 3.1 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 3.1, display our moment-based convergence approach, where the concentration of relevant quantities about their respective mean is proved. To prove Theorem 3.1, we first show (Lemma 3.1) that for each estimator in  $\mathcal{S}[\tau]$  that shrinks toward the origin (i.e., with  $\mu$  set to 0) the URE methodology estimates the risk in  $L_2$  norm below the desired error rate. Thereafter, in Lemma 3.2 we prove that the loss of those estimators concentrate around their expected values (risk) when we have a large number of row and column effects.

LEMMA 3.1. *Under Assumptions A1–A2,*

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

LEMMA 3.2. *Under Assumptions A1–A2,*

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

If we restrict ourselves to only estimators in  $\mathcal{S}[\tau]$  that shrink toward the origin, then Theorem 3.1 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 of these results are provided in Section A.1 of the Appendix. The results for the weighted loss  $L_{r,c}^{\text{wgt}}$  (defined in Section 2) are discussed in Section S.3.1 of the supplement.

We would like to point out the qualitative differences between the results and the proof techniques employed here and those usually used for the one-way Normal mean estimation problem exemplified in Xie, Kou and Brown (2012). For the cell mean estimation problem in two-way layouts with unbalanced designs, there is no indexing on the parametric space under which the “row” effects and the “column” effects can be decoupled. Thus, in unbalanced designs (see Section S.5 of the supplement for discussion on balanced designs) the two-way cell mean estimation problem cannot be reduced to the one-way setting. The approach of Xie, Kou and Brown (2012), which would require showing uniform convergence of the difference between the URE and the loss over the entire set  $\mathcal{H} = \{\boldsymbol{\mu} \in [\hat{a}_\tau, \hat{b}_\tau]; \lambda_A, \lambda_B \geq 0\}$  of possible hyperparameter values, that is, showing  $L_1$  convergence of  $\sup_{(\boldsymbol{\mu}, \lambda_A, \lambda_B) \in \mathcal{H}} |\text{URE}_{r,c}^Q(\boldsymbol{\mu}, \lambda_A, \lambda_B) - L_{r,c}^Q(\boldsymbol{\eta}, \hat{\boldsymbol{\eta}}(\boldsymbol{\mu}, \lambda_A, \lambda_B))|$  to 0, cannot be trivially adapted to the two-way layout.

Instead, in Theorem 3.1 we show pointwise control on the expected absolute difference between the URE and the loss. The pointwise control was established through computations involving the variance of Gaussian quadratic forms, which greatly helps in tackling the difficulties encountered when passing to the two-way problem, where computations involving matrices are generally unavoidable. Specifically, we show that the expected absolute difference between the URE and the loss is asymptotically controlled at  $o(d_{r,c}^{-1})$  over the essential support of the hyperparameters. Thereafter, for establishing oracle optimality of the hyperparameters produced by our URE method, we leverage the differentiability of the loss as function of the hyperparameters. This allows us to concentrate on controlling the expected absolute difference uniformly over a discrete subset of  $\mathcal{H}$  whose cardinality diverges to infinity as  $r, c \rightarrow \infty$  but is of the order of  $O(d_{r,c})$  (see Appendix A.2 for details). The oracle optimality results follow by applying the pointwise control of  $o(d_{r,c}^{-1})$  on the expected absolute difference between the URE and the loss at every point on this discrete subset which has only  $O(d_{r,c})$  hyperparameter values; thus, ensuring at most  $o(d_{r,c}^{-1}) \cdot O(d_{r,c}) = o(1)$  cumulative error over the concerned discrete subset. Section A.2 contains the detailed proofs of the oracle optimality results of Theorems 3.2, 3.3 and that of Corollary 3.1.

**4. Empirical studies.** We carry out numerical experiments to compare the performance of the URE-based estimator to that of different estimators discussed in the previous sections. As the standard technique, we consider the weighted least squares estimator  $\hat{\eta}^{\text{LS}} = \hat{\mu}^{\text{LS}} \mathbf{1} + Z \hat{\theta}^{\text{LS}}$ , where  $(\hat{\mu}^{\text{LS}}, \hat{\theta}^{\text{LS}})$  is any pair that minimizes

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

The two-way shrinkage estimators reported are the maximum-likelihood empirical Bayes (EBML) estimator  $\hat{\eta}^{\text{ML}}$  and the URE-based estimator  $\hat{\eta}^{\text{URE}}$ , as well as versions of these two estimators which shrink toward the origin (i.e., with  $\mu$  fixed at 0); these are designated in Table 1 as “EBMLE (origin)” and “URE (origin).”

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 4*

	(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

We also consider an adaptation of the one-way estimator of [Xie, Kou and Brown \(2012\)](#) to the two-way layout, which estimates the scale hyperparameters based on two independent one-way shrinkage problems and shrinks toward a general data-driven location; details for this estimator, which we label “XKB” in [Table 1](#), appear in [Section S.5](#) of the supplement. As a benchmark, we consider the oracle rule  $\tilde{\boldsymbol{\eta}}^{\text{OL}} = \hat{\boldsymbol{\eta}}(\tilde{\boldsymbol{\mu}}^{\text{OL}}, \tilde{\lambda}_A^{\text{OL}}, \tilde{\lambda}_B^{\text{OL}})$ , where

$$(21) \quad (\tilde{\boldsymbol{\mu}}^{\text{OL}}, \tilde{\lambda}_A^{\text{OL}}, \tilde{\lambda}_B^{\text{OL}}) = \underset{\mu, \lambda_A \geq 0, \lambda_B \geq 0}{\arg \min} \|\mathbf{y} - M\boldsymbol{\Sigma}^{-1}(\mathbf{y} - \mu \cdot \mathbf{1}) - \boldsymbol{\eta}\|^2.$$

*Simulation Experiments.* We report results across 6 simulation scenarios. For each of these, 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  $K_{ij}$ . We then draw  $y_{ij} \sim N(\mu + \alpha_i + \beta_j, \sigma^2 K_{ij}^{-1})$  independently, fixing  $\mu = 0$  throughout and setting  $\sigma^2$  to some (known) constant value. This process is repeated for  $N = 100$  times 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 A.3](#) of the [Appendix](#), and the oracle “estimate” is computed employing a similar technique. The EBMLE estimate is computed using the R package `lme4` [[Bates \(2010\)](#)].

[Table 1](#) shows the mean square error (MSE) 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). [Figure 1](#) displays the MSE of the URE, EBMLE, LS, XKB and the Oracle loss (OL) estimators across the six experiments as the number of levels in the design varies. The figure shows how the estimation errors of the different estimators compare with the minimum achievable (oracle) error rates as the number of levels in the design increases. An additional simulation, examining the performance of the estimators under model misspecification, is included in the supplement.

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 the URE-based method is close to that of the oracle when the number of levels is large; when the number of levels is bigger 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 poorly, and is outperformed even by the “one-way” XKB estimator. In cases with dependency be-

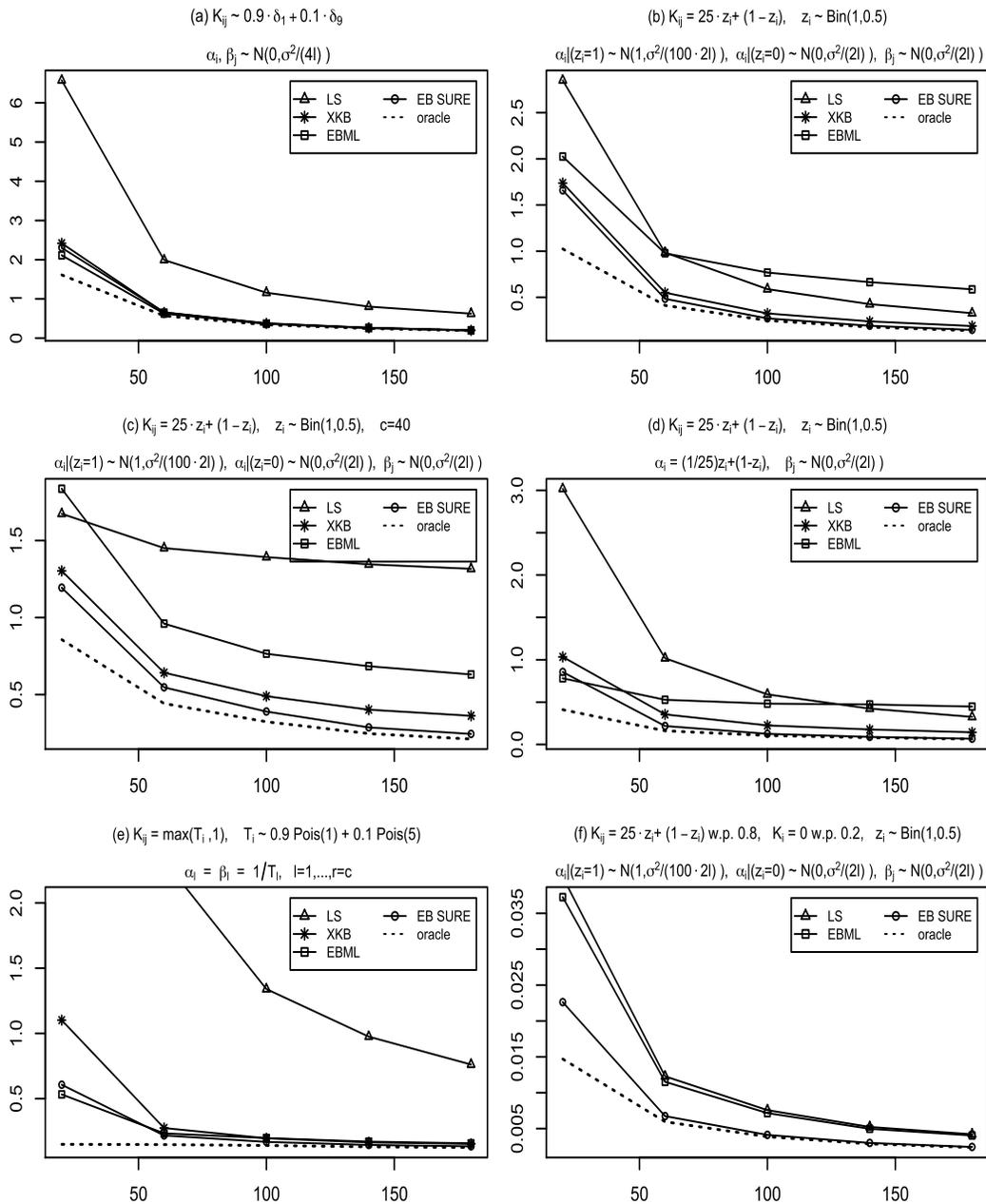


FIG. 1. Risk of the various estimators in the six simulation scenarios of Table 1. The ordinate shows the risk of the estimators while we vary  $L$  along the abscissa. In all experiments  $r = c = L$ , except (d) where  $L = r$  and  $c$  was fixed at 40. The sampling variance is  $\sigma^2 = 10$  in all experiments except (f), where  $\sigma^2 = 1$ .

tween the effects and the cell counts, even the LS estimator can be preferable to the EBMLE [experiments (b) and (d)]. A detailed description of each of the simulation examples is provided in the supplement to this article.

TABLE 2  
*Estimated fixed effect ( $\mu$ ) and components  $\lambda_{\text{county}}$  and  $\lambda_{\text{hour}}$ , which determine shrinkage*

	$\mu$	County	Hour
EBMLE	1.10	0.57	0.05
URE	0.78	0.07	0.80

*Real Data Analysis.* In addition to the simulation study, we analyzed real data collected on Nitrate levels measured in water sources across the US. The data was obtained from the Water Quality Portal cooperative (<http://waterqualitydata.us/>). We estimated the (log transformed) average Nitrate level based on the location of the water resource and time when the measurement was taken. Our data set includes a total of 858 observations categorized into 108 counties and 9 different hour-slots (8–16). The data is highly unbalanced: 57% of the cells are empty, and the cell counts among the nonempty cells vary between 1 to 12. Both the EBMLE and the URE estimators shrink the LS estimates in this example, but the shrinkage factors are quite different. There is a marked difference between the methods in the estimates of the two variance components, as Table 2 shows. Note specifically that EBMLE applies much more shrinkage to “hour”; this is in accord with the fact that the p-value for “hour” in a two-way ANOVA (using Type-II sums-of-squares) is not significant (p-value = 0.25). In terms of the *fitted* values, EBMLE applies more shrinkage, although the differences are not very big. 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 two-way additive model (1) is correctly specified. In the analysis with validation, EBMLE had a smaller estimated total squared error (0.42 vs. 0.5, presented as fraction of TSE of LS), but in the data-informed simulation, the URE achieved smaller squared error (0.72 vs. 0.81). As both EBMLE and the URE estimators are designed for the additive model, the results from the latter analysis might be considered a better basis for comparison between the methods. The two experiments are described in details in the supplement.

**5. 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 suboptimal 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 via efficient pointwise risk estimation which reduces to uniformly controlling the variance of the associated quadratic forms. Our proof techniques can be extended to  $k$ -way additive models, although computational difficulty of our proposed method might become a problem when  $k$  is even moderately large. It would be interesting to investigate whether our proof techniques can be extended to tackle URE-based shrinkage estimation in more complex models (non-Gaussian or misspecified) by using the general quadratic form concentration results established in [Dicker and Erdogdu \(2017\)](#).

## APPENDIX

**A.1. Pointwise loss estimation by URE: Proof details.** Here, we present the proofs of [Theorem 3.1](#), [Lemma 3.1](#) and [Lemma 3.2](#). Henceforth, we assume  $\sigma = 1$ . It is done mainly for the ease of presentation, and the proofs can easily be modified for any known value of  $\sigma$ . We fix the following notation. Denote by  $\sigma_k(A)$  the  $k$ th largest singular value of a matrix  $A$ . Denote by  $\lambda_k(B)$  the  $k$ th largest eigenvalue of a symmetric matrix  $B$ . Also,  $G \doteq M\Sigma^{-1}$  and  $H \doteq G^\top QG = \Sigma^{-1}MQM\Sigma^{-1}$  where  $M$ ,  $\Sigma^{-1}$  and  $Q$  are defined in [Section 2.3](#). Let  $W = M^{1/2}\Sigma^{-1}M^{1/2}$ . As  $0 < M \preceq \Sigma$  we have  $W \preceq I$ , and also  $W^2 \preceq I$ . We will use the following result of [Searle, Casella and McCulloch \(1992\)](#) ([Theorem S4](#), page 467).

LEMMA A.1 (Central moments of Gaussian Quadratic Forms).

If  $\mathbf{y} \sim N(\boldsymbol{\eta}, V)$ ,

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

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

The proof of Theorem 3.1 for the case when the general effect  $\mu = 0$  follows directly from the results of Lemma 3.1 and Lemma 3.2 as  $\mathbb{E}\{\text{URE}^Q_{r,c}(0, \lambda_A, \lambda_B) - L^Q_{r,c}(\boldsymbol{\eta}, \hat{\boldsymbol{\eta}}(0, \lambda_A, \lambda_B))\}^2$  is bounded above by

$$2\mathbb{E}\{\text{URE}^Q_{r,c}(0, \lambda_A, \lambda_B) - R^Q_{r,c}(\boldsymbol{\eta}, \hat{\boldsymbol{\eta}}(0, \lambda_A, \lambda_B))\}^2 \\ + 2\mathbb{E}\{L^Q_{r,c}(\boldsymbol{\eta}, \hat{\boldsymbol{\eta}}(0, \lambda_A, \lambda_B)) - R^Q_{r,c}(\boldsymbol{\eta}, \hat{\boldsymbol{\eta}}(0, \lambda_A, \lambda_B))\}^2.$$

In fact, this proves Theorem 3.1 in the stronger  $L_2$  norm (with its correspondingly adjusted rate). We now concentrate on proving the lemmas and will thereafter prove the theorem for the general ( $\mu \neq 0$ ) case.

**PROOF OF LEMMA 3.1.** 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

$$(22) \quad \mathbb{E}[\text{URE}^Q(0, \lambda_A, \lambda_B) - R^Q_{r,c}(\boldsymbol{\eta}; \hat{\boldsymbol{\eta}}(0, \lambda_A, \lambda_B))]^2 = \text{Var}[\text{URE}^Q(0, \lambda_A, \lambda_B)].$$

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

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

and so the RHS of (22) 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

$$(23) \quad \text{Var}[\text{URE}^Q(0, \lambda_A, \lambda_B)] = (rc)^{-2} \{2\text{tr}(HMHM) + 4\boldsymbol{\eta}^\top HMH\boldsymbol{\eta}\}.$$

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.6 of the supplement),  $HMH$  is also p.s.d. Thus,  $\boldsymbol{\eta}^\top HMH\boldsymbol{\eta} \leq \lambda_1(HMH)\|\boldsymbol{\eta}\|^2$ . Now, using the bound  $\lambda_1(HMH) \leq \nu_{r,c}\lambda_1^2(Q)$  on the largest eigenvalue of  $HMH$  (for detailed derivation, see Lemma S.3.2 in the supplements) we arrive at the following upper bound:

$$(rc)^{-2} d_{r,c}^2 \sup_{\lambda_A, \lambda_B \geq 0} \boldsymbol{\eta}^\top HMH\boldsymbol{\eta} \leq (rc)^{-2} d_{r,c}^2 \nu_{r,c} \lambda_1^2(Q) \|\boldsymbol{\eta}\|^2$$

which, under Assumptions A1 and A2, converges to 0 as  $r, c \rightarrow \infty$ . From Lemma S.3.2, we have  $\text{tr}(HMHM) \leq rc \cdot \lambda_1^2(Q)$ . Hence, as  $r, c \rightarrow \infty$ , by Assumption A2 the first term in (23) 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.  $\square$

**PROOF OF LEMMA 3.2.** As the risk is the expectation of the loss, to prove the lemma we need to show

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

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

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

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 3.1.

For the first term  $L_1$ , we have

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

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

$$\begin{aligned} 4^{-1} \text{Var}(L_4) &= \text{Var}(\boldsymbol{\eta}^\top QG\mathbf{y}) = \boldsymbol{\eta}^\top QGMG^\top Q\boldsymbol{\eta} \leq \lambda_1(QGMG^\top Q) \|\boldsymbol{\eta}\|^2 \\ &\leq v_{r,c} \lambda_1^2(Q) \|\boldsymbol{\eta}\|^2, \end{aligned}$$

where the last inequality follows from Lemma S.3.2. By Assumptions A1–A2, the above is  $o(r^2c^2d_{r,c}^{-2})$ .

The third term requires detailed analysis. It breaks into two components:

$$\begin{aligned} (24) \quad \text{Var}[(rc)^{-1}L_3] &= 4(rc)^{-2} \text{Var}(\mathbf{y}^\top QG\mathbf{y}) \\ &= 8(rc)^{-2} \text{tr}(\widetilde{G}M\widetilde{G}M) + 16(rc)^{-2} \boldsymbol{\eta}^\top \widetilde{G}M\widetilde{G}\boldsymbol{\eta}, \end{aligned}$$

where  $\widetilde{G} = QG + G^\top Q$  is a symmetric matrix. By Lemma S.3.2, the first term on the RHS of (24) is bounded above by  $32(rc)^{-1} \lambda_1^2(Q)$  and hence, by Assumption A2 is uniformly bounded by  $o(d_{r,c}^{-2})$  as  $r, c \rightarrow \infty$ . We now concentrate on the second term of the RHS of (24). Note that  $(rc)^{-2} \boldsymbol{\eta}^\top \widetilde{G}M\widetilde{G}\boldsymbol{\eta} \leq (rc)^{-2} \boldsymbol{\eta}^\top \boldsymbol{\eta} \sigma_1(\widetilde{G}M\widetilde{G})$ . If we can uniformly bound the largest eigenvalue of  $\widetilde{G}M\widetilde{G}$  as  $o(rcd_{r,c}^{-2})$  then by Assumption A1 the second term of the RHS of (24) is  $o(d_{r,c}^{-2})$  as  $r, c \rightarrow \infty$ . As  $\widetilde{G} = QG + G^\top Q = QM\Sigma^{-1} + \Sigma^{-1}MQ$ , we have

$$\widetilde{G}M\widetilde{G} = H_1 + H_1^\top + H_2 + H_3,$$

where

$$\begin{aligned} H_1 &= QM\Sigma^{-1}MQM\Sigma^{-1}, \\ H_2 &= QM\Sigma^{-1}M\Sigma^{-1}MQ, \quad H_3 = \Sigma^{-1}MQM\Sigma^{-1}. \end{aligned}$$

To uniformly bound the eigenvalues of  $\tilde{G}M\tilde{G}$  as desired before we just show that for each of  $i = 1, \dots, 3$ ,  $(rc)^{-1}d_{r,c}^{-2}\sigma_1(H_i) \rightarrow 0$  as  $r, c \rightarrow \infty$ .

Note that,  $H_2, H_3$  is symmetric but  $H_1$  is not. By Lemma S.3.2, we have  $\max\{\sigma_1(H_1), \lambda_1(H_3)\} \leq v_{r,c}\lambda_1^2(Q)$ , from which the desired asymptotic control on the variance follows. For  $H_2$  using  $W^2 \leq I$ , 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. This completes the proof of the lemma.  $\square$

**PROOF OF THEOREM 3.1.** As shown in the beginning of this section, by the above two lemmas the theorem easily follows for the case when  $\mu = 0$ . We now prove the theorem for the general case. First, note that for arbitrary fixed  $\mu \in \mathbb{R}$ , the loss  $(\hat{\boldsymbol{\eta}}(\mu, \lambda_A, \lambda_B) - \boldsymbol{\eta})^\top Q(\hat{\boldsymbol{\eta}}(\mu, \lambda_A, \lambda_B) - \boldsymbol{\eta})$  decomposes into the following components:

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

Comparing it with the definition of URE, we have

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

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 show:

$$(25) \quad \lim_{\substack{r \rightarrow \infty \\ c \rightarrow \infty}} \sup_{\substack{|\mu| \leq m_{r,c} \\ \lambda_A, \lambda_B \geq 0}} d_{r,c} \cdot (rc)^{-1} \cdot \mathbb{E}|\mu \mathbf{1}^\top G^\top Q(\mathbf{y} - \boldsymbol{\eta})| = 0.$$

Noting that for any fixed  $\eta$  the random variable  $F = \mathbf{1}^\top G^\top Q(\mathbf{y} - \boldsymbol{\eta})$  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 we show

$$(26) \quad \lim_{r \rightarrow \infty, c \rightarrow \infty} m_{r,c} \cdot d_{r,c} \cdot (rc)^{-1} \cdot \sup_{\lambda_A, \lambda_B \geq 0} \{\text{Var}(|\mathbf{1}^\top G^\top Q(\mathbf{y} - \boldsymbol{\eta})|)\}^{1/2} = 0.$$

The above is true based on our assumptions as we can upper bound the variance of  $|F|$  by

$$\text{Var}(F) \leq \mathbf{1}^\top G^\top Q M Q G \mathbf{1} \leq rc \lambda_1(G^\top Q M Q G) \leq rc v_{r,c} \lambda_1^2(Q),$$

where the last inequality follows from Lemma S.3.2. This based on Assumption A2 establishes (26) and completes the proof of the theorem.  $\square$

**A.2. Proof of the oracle optimality results.** The proofs of Theorems 3.2, 3.3 and Corollary 3.1 are presented here. For that purpose, we first construct a discrete subset of the set of hyperparameters and define analogous versions of the URE and oracle estimators over it.

*Discretization.* In (18) and (20), the URE and oracle estimators involve minimizing the hyperparameters  $(\mu, \lambda_A, \lambda_B)$  simultaneously over  $\hat{T}_{r,c} = [\hat{a}_\tau, \hat{b}_\tau] \times [0, \infty] \times [0, \infty]$  where the range of the location hyperparameter  $\mu$  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 afterward. 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

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

and the corresponding oracle rule by  $\tilde{\boldsymbol{\eta}}^{\text{OD}} = Z_c Z^\dagger \hat{\boldsymbol{\eta}}(\tilde{\mu}^{\text{OD}}, \tilde{\lambda}_A^{\text{OD}}, \tilde{\lambda}_B^{\text{OD}})$ . 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 hyperparameters 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:

$$(27) \quad (\hat{\mu}^{\text{UD}}, \hat{\lambda}_A^{\text{UD}}, \hat{\lambda}_B^{\text{UD}}) = \underset{(\mu, \lambda_A, \lambda_B) \in \{\mu_1, \mu_2\} \times \{\lambda_1, \lambda_2\} \times \{\lambda_3, \lambda_4\}}{\text{arg min}} \text{URE}^Q(\mu, \lambda_A, \lambda_B).$$

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

*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 hyperparameter  $\mu$  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

$$(28) \quad \delta_{r,c}^{[1]} = \{m_{r,c}^{4/3} \cdot \nu_{r,c} \cdot \lambda_1(Q)\}^{-1}.$$

For constructing the grid  $\Theta_{r,c}^{[2]}$  on the scale hyperparameter, we consider the following transformation  $\tilde{\lambda}_A = (1 + \lambda_A)^{-1/2}$ . Note that  $\tilde{\lambda}_A \in [0, 1]$  as  $\lambda_A$  varies over  $[0, \infty]$ . Let  $\tilde{\lambda}_A[k] = (k - 1)\delta_{r,c}^{[2]}$ ,  $n_2 = \lceil \{\delta_{r,c}^{[2]}\}^{-1} \rceil$ . Consider the equispaced grid on  $\tilde{\lambda}_A$  between 0 and 1 at a spacing of  $\delta_{r,c}^{[2]}$ :  $\{0 = \tilde{\lambda}_A[1] < \tilde{\lambda}_A[2] < \dots < \tilde{\lambda}_A[n_2] \leq 1\}$ . The grid on  $\tilde{\lambda}_A$  is then retransformed to produce the grid  $\Theta_{r,c}^{[2]}$  on the scale hyperparameter  $\lambda_A$  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  $\tilde{\lambda}_B$  scale. The spaces were chosen as

$$(29) \quad \delta_{r,c}^{[2]} = \delta_{r,c}^{[3]} = \{m_{r,c}^{7/3} \cdot \nu_{r,c} \cdot \lambda_1(Q)\}^{-1}.$$

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 our decision theoretic results.

LEMMA A.2. *Under Assumptions A1–A2 for any fixed  $\varepsilon > 0$  as  $r, c \rightarrow \infty$ ,*

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

LEMMA A.3. *Under Assumptions A1–A2 for any fixed  $\varepsilon > 0$  as  $r, c \rightarrow \infty$ ,*

- A.  $P\{\text{URE}^Q(\hat{\mu}^{\text{UD}}, \hat{\lambda}_A^{\text{UD}}, \hat{\lambda}_B^{\text{UD}}) - \text{URE}^Q(\hat{\mu}^{\text{UQ}}, \hat{\lambda}_A^{\text{UQ}}, \hat{\lambda}_B^{\text{UQ}}) > \varepsilon\} \rightarrow 0$ ,
- B.  $\mathbb{E}[\text{URE}^Q(\hat{\mu}^{\text{UD}}, \hat{\lambda}_A^{\text{UD}}, \hat{\lambda}_B^{\text{UD}}) - \text{URE}^Q(\hat{\mu}^{\text{UQ}}, \hat{\lambda}_A^{\text{UQ}}, \hat{\lambda}_B^{\text{UQ}})] \rightarrow 0$ .

Lemma A.2 above shows that the difference in the loss between the true estimators and their discretized versions is asymptotically controlled at any prefixed level. It uses the following two lemmas. Lemma A.3 above shows that the URE values for the estimator is also asymptotically close for the discretized version. The proofs of all these lemmas (A.2, A.3, A.4 and A.5) are provided in the supplementary materials.

LEMMA A.4. *Under Assumption A1 in model (14), for any fixed  $\tau \in (0, 1]$  and  $m_{r,c}$  of (19), the event  $A_{r,c}(Y) = \{[\hat{a}_\tau(Y), \hat{b}_\tau(Y)] \subseteq [-m_{r,c}, m_{r,c}]\}$  satisfies  $P\{A_{r,c}\} \rightarrow 1$  as  $r, c \rightarrow \infty$ .*

LEMMA A.5. *Under Assumptions A1–A2, for any fixed  $\tau \in (0, 1]$  and  $m_{r,c}$  of (19), the event  $A_{r,c}(\mathbf{Y}) = \{[\hat{a}_\tau(\mathbf{Y}), \hat{b}_\tau(\mathbf{Y})] \subseteq [-m_{r,c}, m_{r,c}]\}$  satisfies:*

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

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

PROOF OF THEOREM 3.2. We know that

$$P\{L(\boldsymbol{\eta}_c, \hat{\boldsymbol{\eta}}_c^{\text{URE}}) \geq L(\boldsymbol{\eta}_c, \tilde{\boldsymbol{\eta}}_c^{\text{OL}}) + \varepsilon\} \leq P\{L(\boldsymbol{\eta}_c, \hat{\boldsymbol{\eta}}_c^{\text{URE}}) \geq L(\boldsymbol{\eta}_c, \tilde{\boldsymbol{\eta}}_c^{\text{OD}}) + \varepsilon/2\} \\ + P\{L(\boldsymbol{\eta}_c, \tilde{\boldsymbol{\eta}}_c^{\text{OD}}) \geq L(\boldsymbol{\eta}_c, \tilde{\boldsymbol{\eta}}_c^{\text{OL}}) + \varepsilon/2\}.$$

By Lemma A.2, the second term converges to 0. The first term is less than

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

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

$$P\{\text{URE}^{\text{Q}}(\hat{\mu}^{\text{UD}}, \hat{\lambda}_A^{\text{UD}}, \hat{\lambda}_B^{\text{UD}}) \leq \text{URE}^{\text{Q}}(\tilde{\mu}^{\text{OD}}, \tilde{\lambda}_A^{\text{OD}}, \tilde{\lambda}_B^{\text{OD}}) + \varepsilon/8\} \rightarrow 0 \\ (30) \quad \text{as } r, c \rightarrow \infty.$$

Thus, showing  $P\{L(\boldsymbol{\eta}_c, \hat{\boldsymbol{\eta}}_c^{\text{UD}}) \geq L(\boldsymbol{\eta}, \tilde{\boldsymbol{\eta}}_c^{\text{OD}}) + \varepsilon/4\} \rightarrow 0$  as  $r, c \rightarrow \infty$  can be reduced to showing  $\lim_{r,c \rightarrow \infty} P\{A(\mathbf{y}; \boldsymbol{\eta}_c) \geq B(\mathbf{y}; \boldsymbol{\eta}_c) + \varepsilon/8\} = 0$  where

$$A(\mathbf{y}; \boldsymbol{\eta}_c) = L(\boldsymbol{\eta}_c, \hat{\boldsymbol{\eta}}_c^{\text{UD}}) - \text{URE}^{\text{Q}}(\hat{\mu}^{\text{UD}}, \hat{\lambda}_A^{\text{UD}}, \hat{\lambda}_B^{\text{UD}}),$$

and

$$B(\mathbf{y}; \boldsymbol{\eta}_c) = L(\boldsymbol{\eta}_c, \tilde{\boldsymbol{\eta}}_c^{\text{OD}}) - \text{URE}^{\text{Q}}(\tilde{\mu}^{\text{OD}}, \tilde{\lambda}_A^{\text{OD}}, \tilde{\lambda}_B^{\text{OD}}).$$

As  $L(\boldsymbol{\eta}_c, \hat{\boldsymbol{\eta}}_c^{\text{UD}}) = L^{\text{Q}}(\boldsymbol{\eta}, \hat{\boldsymbol{\eta}}^{\text{UD}})$  and  $L(\boldsymbol{\eta}_c, \tilde{\boldsymbol{\eta}}_c^{\text{OD}}) = L^{\text{Q}}(\boldsymbol{\eta}, \hat{\boldsymbol{\eta}}(\tilde{\mu}^{\text{OD}}, \tilde{\lambda}_A^{\text{OD}}, \tilde{\lambda}_B^{\text{OD}}))$ , by using Markov's inequality we get

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

Now, by triangle inequality the RHS above is upper bounded by

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

As  $|\Theta_{r,c}| = O(d_{r,c})$  by Theorem 3.1, the above expression converges to zero when  $r, c \rightarrow \infty$ . This completes the proof of the theorem.  $\square$

**PROOF OF THEOREM 3.3.** Decompose the loss as the sum of the following components:  $L(\eta_c, \hat{\eta}_c^{\text{URE}}) - L(\eta_c, \hat{\eta}_c^{\text{UD}})$ ,  $L(\eta_c, \tilde{\eta}_c^{\text{OD}}) - L(\eta_c, \tilde{\eta}_c^{\text{OL}})$  and  $L(\eta_c, \hat{\eta}_c^{\text{UD}}) - L(\eta_c, \tilde{\eta}_c^{\text{OD}})$ . By Lemma A.2, expectations of the absolute values of the first two terms converge to 0 as  $r, c \rightarrow \infty$ . Write the third term as

$$\begin{aligned} & \{L(\eta_c, \hat{\eta}_c^{\text{UD}}) - \text{URE}^{\text{Q}}(\hat{\mu}^{\text{UD}}, \hat{\lambda}_A^{\text{UD}}, \hat{\lambda}_B^{\text{UD}})\} \\ & - \{L(\eta_c, \tilde{\eta}_c^{\text{OD}}) - \text{URE}^{\text{Q}}(\tilde{\mu}^{\text{OD}}, \tilde{\lambda}_A^{\text{OD}}, \tilde{\lambda}_B^{\text{OD}})\} \\ & + \{\text{URE}^{\text{Q}}(\hat{\mu}^{\text{UD}}, \hat{\lambda}_A^{\text{UD}}, \hat{\lambda}_B^{\text{UD}}) - \text{URE}^{\text{Q}}(\tilde{\mu}^{\text{OD}}, \tilde{\lambda}_A^{\text{OD}}, \tilde{\lambda}_B^{\text{OD}})\}. \end{aligned}$$

Now, by definition  $\text{URE}^{\text{Q}}(\hat{\mu}^{\text{UD}}, \hat{\lambda}_A^{\text{UD}}, \hat{\lambda}_B^{\text{UD}}) \leq \text{URE}^{\text{Q}}(\tilde{\mu}^{\text{OD}}, \tilde{\lambda}_A^{\text{OD}}, \tilde{\lambda}_B^{\text{OD}})$  which, combined with Lemma A.3, suggests that the last term above has asymptotically non-positive expectation. Thus, there exist constants  $\kappa_{r,c}$  with  $\kappa_{r,c} \rightarrow 0$  as  $r, c \rightarrow \infty$  such that  $\mathbb{E}\{L(\eta_c, \hat{\eta}_c^{\text{UD}}) - L(\eta_c, \tilde{\eta}_c^{\text{OD}})\}$  is

$$\begin{aligned} & \leq 2\mathbb{E}\left\{\sup_{(\mu, \lambda_A, \lambda_B) \in \Theta_{r,c}} |L^{\text{Q}}(\eta, \hat{\eta}(\mu, \lambda_A, \lambda_B)) - \text{URE}^{\text{Q}}(\mu, \lambda_A, \lambda_B)|\right\} + \kappa_{r,c} \\ & \leq 2\mathbb{E}\left\{\sum_{(\mu, \lambda_A, \lambda_B) \in \Theta_{r,c}} |L^{\text{Q}}(\eta, \hat{\eta}(\mu, \lambda_A, \lambda_B)) - \text{URE}^{\text{Q}}(\mu, \lambda_A, \lambda_B)|\right\} + \kappa_{r,c} \\ & \leq 2|\Theta_{r,c}| \sup_{|\mu| \in m_{r,c}; \lambda_A, \lambda_B \geq 0} \mathbb{E}\{|L^{\text{Q}}(\eta, \hat{\eta}(\mu, \lambda_A, \lambda_B)) - \text{URE}^{\text{Q}}(\mu, \lambda_A, \lambda_B)|\} + \kappa_{r,c}. \end{aligned}$$

As  $|\Theta_{r,c}| = O(d_{r,c})$ , the above expression tends to zero when  $r, c \rightarrow \infty$  by Theorem 3.1. This completes the proof of Theorem 3.3.  $\square$

**PROOF OF COROLLARY 3.1.** (a) and (b) are direct consequences, respectively, of Theorems 3.2 and 3.3, since  $L^{\text{Q}}(\eta, \hat{\eta}(\hat{\mu}, \hat{\lambda}_A, \hat{\lambda}_B)) \geq L^{\text{Q}}(\eta, \eta^{\text{OL}})$ , and hence, also  $\mathbb{E}\{L^{\text{Q}}(\eta, \hat{\eta}(\hat{\mu}, \hat{\lambda}_A, \hat{\lambda}_B))\} \geq \mathbb{E}\{L^{\text{Q}}(\eta, \eta^{\text{OL}})\}$ . Unlike in the above two theorems, here we only have optimality over the loss  $L^{\text{Q}}$  defined over the observed cells with  $Q$  in (16). As explained in Section 2.3, the loss  $L^{\text{Q}}$  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  $Z_c Z^{\dagger} \hat{\eta}(\hat{\mu}, \hat{\lambda}_A, \hat{\lambda}_B)$ , where  $\hat{\mu}, \hat{\lambda}_A, \hat{\lambda}_B$  are any estimates of the hyper-parameters.  $\square$

**A.3. Computations for implementing the URE method.** By definition,  $\Sigma = Z\Lambda\Lambda^{\text{T}}Z^{\text{T}} + M$ . By the matrix inverse identity, we have  $\Sigma^{-1} = M^{-1} - M^{-1}Z\Lambda(\Lambda^{\text{T}}Z^{\text{T}}M^{-1}Z\Lambda + I_q)^{-1}\Lambda^{\text{T}}Z^{\text{T}}M^{-1}$ . Hence, we get

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

Using the above, we get

$$\text{tr}(\Sigma^{-1}M^2) = \text{tr}(M\Sigma^{-1}M) = \text{tr}(M) - \text{tr}(Z\Lambda(\Lambda^\top Z^\top M^{-1}Z\Lambda + I_q)^{-1}\Lambda^\top Z^\top).$$

Therefore, the URE expression in (10) multiplied by  $rc$  can be written as

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

In computing the above expression:

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

2.  $(\Lambda^\top Z^\top M^{-1}Z\Lambda + I_q)^{-1}$  is computed efficiently employing a sparse Cholesky factorization of  $\Lambda^\top Z^\top M^{-1}Z\Lambda + I_q$  similar 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

$$(31) \quad M\Sigma^{-1}\mathbf{x} = \mathbf{x} - Z\Lambda(\Lambda^\top Z^\top M^{-1}Z\Lambda + I_q)^{-1}\Lambda^\top Z^\top(M^{-1}\mathbf{x}),$$

where (31) 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  $(\Lambda^\top Z^\top)(M^{-1}\mathbf{x})$ , and so on.

**Acknowledgement.** We thank Tony Cai, Samuel Kou and Art Owen for helpful discussions. We would also like to thank the Associate Editor and three referees for constructive suggestions to shorten and improve the paper.

## SUPPLEMENTARY MATERIAL

**Supplement to “Empirical Bayes estimates for a two-way cross-classified model”** (DOI: [10.1214/17-AOS1599SUPP](https://doi.org/10.1214/17-AOS1599SUPP); .pdf). The supplement [Brown, Mukherjee and Weinstein (2018)] contains detailed proofs of the lemmas that were used in the [Appendix](#) for proving the results in Section 3; and derivations and further discussions on the results of Sections 2 and 4.

## REFERENCES

- BATES, D. M. (2010). `lme4`: Mixed-effects modeling with R. <http://lme4.r-forge.r-project.org/book>.  
 BROWN, L. D., MUKHERJEE, G. and WEINSTEIN, A. (2018). Supplement to “Empirical Bayes estimates for a two-way cross-classified model.” DOI:[10.1214/17-AOS1599SUPP](https://doi.org/10.1214/17-AOS1599SUPP).  
 CANDÈS, E. J., SING-LONG, C. A. and TRZASKO, J. D. (2013). Unbiased risk estimates for singular value thresholding and spectral estimators. *IEEE Trans. Signal Process.* **61** 4643–4657. [MR3105401](#)  
 DICKER, L. H. (2013). Optimal equivariant prediction for high-dimensional linear models with arbitrary predictor covariance. *Electron. J. Stat.* **7** 1806–1834. [MR3084672](#)

- DICKER, L. H. and ERDOGDU, M. A. (2017). Flexible results for quadratic forms with applications to variance components estimation. *Ann. Statist.* **45** 386–414. [MR3611496](#)
- DONOHO, D. L., JOHNSTONE, I. M., KERKYACHARIAN, G. and PICARD, D. (1995). Wavelet shrinkage: Asymptopia? *J. R. Stat. Soc., B* **57** 301–369. [MR1323344](#)
- DRAPER, N. R. and VAN NOSTRAND, R. C. (1979). Ridge regression and James–Stein estimation: Review and comments. *Technometrics* **21** 451–466. [MR0555086](#)
- EFRON, B. and MORRIS, C. (1972). Empirical Bayes on vector observations: An extension of Stein’s method. *Biometrika* **59** 335–347. [MR0334386](#)
- EFRON, B. and MORRIS, C. (1973). Stein’s estimation rule and its competitors—an empirical Bayes approach. *J. Amer. Statist. Assoc.* **68** 117–130. [MR0388597](#)
- FAHRMEIR, L., KNEIB, T., LANG, S. and MARX, B. (2013). *Regression: Models, Methods and Applications*. Springer, Heidelberg. [MR3075546](#)
- GELMAN, A. (2005). Analysis of variance—why it is more important than ever. *Ann. Statist.* **33** 1–53. [MR2157795](#)
- GELMAN, A., CARLIN, J. B., STERN, H. S. and RUBIN, D. B. (2004). *Bayesian Data Analysis*, 2nd ed. Chapman & Hall/CRC, Boca Raton, FL. [MR2027492](#)
- GHOSH, M., NICKERSON, D. M. and SEN, P. K. (1987). Sequential shrinkage estimation. *Ann. Statist.* **15** 817–829. [MR0888442](#)
- GOLDSTEIN, H., BROWNE, W. and RASBASH, J. (2002). Multilevel modelling of medical data. *Stat. Med.* **21** 3291–3315.
- HARVILLE, D. A. (1977). Maximum likelihood approaches to variance component estimation and to related problems. *J. Amer. Statist. Assoc.* **72** 320–340. [MR0451550](#)
- HENDERSON, C. (1984). ANOVA, MIVQUE, REML, and ML algorithms for estimation of variances and covariances. In *Statistics: An Appraisal: Proceedings 50th Anniversary Conference* (H. A. David and H. T. David, eds.) 257–280. The Iowa State University Press, Ames, IA.
- JAMES, W. and STEIN, C. (1961). Estimation with quadratic loss. In *Proc. 4th Berkeley Sympos. Math. Statist. and Prob., Vol. I* 361–379. Univ. California Press, Berkeley, CA. [MR0133191](#)
- JIANG, J., NGUYEN, T. and RAO, J. S. (2011). Best predictive small area estimation. *J. Amer. Statist. Assoc.* **106** 732–745. [MR2847987](#)
- 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. *Ann. Statist.* **32** 1594–1649. [MR2089135](#)
- KOU, S. C. and YANG, J. J. (2017). Optimal shrinkage estimation in heteroscedastic hierarchical linear models. In *Big and Complex Data Analysis. Contrib. Stat.* 249–284. Springer, Cham. [MR3674720](#)
- LI, K.-C. (1986). Asymptotic optimality of  $C_L$  and generalized cross-validation in ridge regression with application to spline smoothing. *Ann. Statist.* **14** 1101–1112. [MR0856808](#)
- LINDLEY, D. (1962). Discussion of the paper by Stein. *J. R. Stat. Soc., B* **24** 265–296.
- LINDLEY, D. V. and SMITH, A. F. M. (1972). Bayes estimates for the linear model. *J. R. Stat. Soc., B* **34** 1–41. [MR0415861](#)
- MASON, W. M., WONG, G. Y. and ENTWISLE, B. (1983). Contextual analysis through the multilevel linear model. *Sociol. Method.* **1984** 72–103.
- MCCULLOCH, C. E. and SEARLE, S. R. (2001). *Generalized, Linear, and Mixed Models*. Wiley, New York. [MR1884506](#)
- OMAN, S. D. (1982). Shrinking towards subspaces in multiple linear regression. *Technometrics* **24** 307–311. [MR0687188](#)
- RASBASH, J. and GOLDSTEIN, H. (1994). Efficient analysis of mixed hierarchical and cross-classified random structures using a multilevel model. *J. Educ. Behav. Stat.* **19** 337–350.
- ROLPH, J. E. (1976). Choosing shrinkage estimators for regression problems. *Comm. Statist. Theory Methods* **A5** 789–802. [MR0436481](#)

- SCLOVE, S. L. (1968). Improved estimators for coefficients in linear regression. *J. Amer. Statist. Assoc.* **63** 596–606. [MR0237057](#)
- SEARLE, S. R., CASELLA, G. and MCCULLOCH, C. E. (1992). *Variance Components*. Wiley, New York. [MR1190470](#)
- STEIN, C. M. (1962). Confidence sets for the mean of a multivariate normal distribution. *J. R. Stat. Soc., B* **24** 265–296. [MR0148184](#)
- TAN, Z. (2016). Steinized empirical Bayes estimation for heteroscedastic data. *Statist. Sinica* **26** 1219–1248. [MR3559950](#)
- XIE, X., KOU, S. C. and BROWN, L. D. (2012). SURE estimates for a heteroscedastic hierarchical model. *J. Amer. Statist. Assoc.* **107** 1465–1479. [MR3036408](#)
- ZACCARIN, S. and RIVELLINI, G. (2002). Multilevel analysis in social research: An application of a cross-classified model. *Stat. Methods Appl.* **11** 95–108.

L. D. BROWN  
DEPARTMENT OF STATISTICS  
UNIVERSITY OF PENNSYLVANIA  
400 JON M. HUNTSMAN HALL  
3730 WALNUT STREET  
PHILADELPHIA, PENNSYLVANIA 19104  
USA  
E-MAIL: [lbrown@wharton.upenn.edu](mailto:lbrown@wharton.upenn.edu)

G. MUKHERJEE  
DEPARTMENT OF DATA SCIENCES AND OPERATIONS  
MARSHALL SCHOOL OF BUSINESS  
UNIVERSITY OF SOUTHERN CALIFORNIA  
LOS ANGELES, CALIFORNIA 90089-0809  
USA  
E-MAIL: [gourab@usc.edu](mailto:gourab@usc.edu)

A. WEINSTEIN  
DEPARTMENT OF STATISTICS  
SEQUOIA HALL, 390 SERRA MALL  
STANFORD UNIVERSITY  
STANFORD, CALIFORNIA 94305-4065  
USA  
E-MAIL: [asafw@stanford.edu](mailto:asafw@stanford.edu)

SUPPLEMENT TO “EMPIRICAL BAYES ESTIMATES FOR A  
TWO-WAY CROSS-CLASSIFIED MODEL”

BY LAWRENCE D. BROWN

*University of Pennsylvania*

BY GOURAB MUKHERJEE

*University of Southern California*

AND

BY ASAF WEINSTEIN

*Stanford University*

The supplement contains derivations and further discussions on the results whose proofs were not provided in the main paper. Details for the empirical studies and additional insights on the asymptotic theory are also presented here.

**S.1. Introduction.** Additional proofs, supporting results and discussions related to Section 2, 3 and 4 of the main paper are presented in the following three sections, which are organized according to their corresponding main paper section numbers. Thereafter, Section S.5 provides discussions on URE estimation in two-way layouts with balanced designs and Section S.6 contains a list of basic results used in this paper.

**S.2. Detailed remarks and proofs for results of Section 2.**

LEMMA S.2.1. *The Bayes estimate of  $\boldsymbol{\eta}$  in the hierarchical Gaussian model (3) is*

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

**Proof.** The distribution of  $(\mathbf{y}, \boldsymbol{\theta})^\top$  is Gaussian because it is a linear transformation of  $(\boldsymbol{\theta}, \boldsymbol{\epsilon})^\top$ . Now, since  $\text{cov}(\boldsymbol{\theta}, \mathbf{y}) = \sigma^2 \Lambda \Lambda^\top Z^\top$ ,  $\mathbb{E}[\boldsymbol{\theta}|\mathbf{y}] = \text{cov}(\boldsymbol{\theta}, \mathbf{y})[\text{cov}(\mathbf{y})]^{-1}(\mathbf{y} - \mu \mathbf{1}) =$

$\Lambda\Lambda^\top Z^\top \Sigma^{-1}(\mathbf{y} - \mu\mathbf{1})$ . Hence,

$$\begin{aligned}\mathbb{E}[\boldsymbol{\eta}|\mathbf{y}] &= \mathbb{E}[\mathbf{1} \cdot \mu + Z\boldsymbol{\theta}|\mathbf{y}] = \mu\mathbf{1} + Z\mathbb{E}[\boldsymbol{\theta}|\mathbf{y}] \\ &= \mathbf{1} \cdot \mu + (\Sigma - M)\Sigma^{-1}(\mathbf{y} - \mathbf{1} \cdot \mu) = \mathbf{y} - M\Sigma^{-1}(\mathbf{y} - \mathbf{1} \cdot \mu).\end{aligned}$$

LEMMA S.2.2. *An unbiased estimate of the risk  $R_{r,c}(\boldsymbol{\eta}, \hat{\boldsymbol{\eta}}(\mu, \lambda_A, \lambda_B))$  is*

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

**Proof.** This is immediate from the formula in Berger (1985, p. 362) after noticing that  $\mathbf{y}|\boldsymbol{\eta} \sim N_{rc}(\boldsymbol{\eta}, \sigma^2 M)$  and writing  $\hat{\boldsymbol{\eta}}(\mu, \lambda_A, \lambda_B) = \mathbf{y} - \sigma^2 M(\sigma^2 \Sigma)^{-1}(\mathbf{y} - \mu \cdot \mathbf{1})$ .

S.2.1. *Estimating Equations for Model (1).*

**Estimating Equations for the ML method.** ML estimates are computed based on the likelihood of  $\mathbf{y}$  in the hierarchical model (3). Our derivation is similar to the analysis conducted in Chapter 6.3, 6.4, 6.8 and 6.12 of Searle and McCulloch (2001). Since  $\mathbf{y} \sim N_{rc}(\mathbf{1}\mu, \sigma^2 \Sigma)$ , its density is given by:

$$(S.2.1) \quad f(\mathbf{y}) = \frac{1}{(2\pi\sigma^2)^{rc/2} |\Sigma|^{1/2}} \exp \left\{ -\frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{1}\mu)^t \Sigma^{-1} (\mathbf{y} - \mathbf{1}\mu) \right\}$$

and the corresponding log-likelihood is

$$(S.2.2) \quad l(\mu, \boldsymbol{\theta}) = -(rc)/2 \cdot \log(2\pi\sigma^2) - \frac{1}{2} \log |\Sigma| - \frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{1}\mu)^t \Sigma^{-1} (\mathbf{y} - \mathbf{1}\mu)$$

Using chain rule, we have

$$(S.2.3) \quad \frac{\partial l}{\partial \mu} \stackrel{(7)}{=} -\frac{1}{\sigma^2} (\mathbf{y} - \mathbf{1}\mu)^t \Sigma^{-1} \frac{\partial \{\mathbf{y} - \mathbf{1}\mu\}}{\partial \mu} = \frac{1}{\sigma^2} (\mathbf{y} - \mathbf{1}\mu)^t \Sigma^{-1} \mathbf{1}$$

Also,

$$\begin{aligned}(S.2.4) \quad \frac{\partial l}{\partial \lambda_A^2} &\stackrel{(R8)}{=} -\frac{1}{2} \text{tr} \left( \Sigma^{-1} \frac{\partial \Sigma}{\partial \lambda_A^2} \right) - \frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{1}\mu)^t \left[ \frac{\partial \Sigma^{-1}}{\partial \lambda_A^2} \right] (\mathbf{y} - \mathbf{1}\mu) \\ &= -\frac{1}{2} \left\{ \text{tr} \left( \Sigma^{-1} \frac{\partial \Sigma}{\partial \lambda_A^2} \right) + \frac{1}{\sigma^2} (\mathbf{y} - \mathbf{1}\mu)^t \left[ \frac{\partial \Sigma^{-1}}{\partial \lambda_A^2} \right] (\mathbf{y} - \mathbf{1}\mu) \right\} \\ &\stackrel{(R9)}{=} -\frac{1}{2} \left\{ \text{tr} \left( \Sigma^{-1} \frac{\partial \Sigma}{\partial \lambda_A^2} \right) - \frac{1}{\sigma^2} (\mathbf{y} - \mathbf{1}\mu)^t \Sigma^{-1} \left[ \frac{\partial \Sigma}{\partial \lambda_A^2} \right] \Sigma^{-1} (\mathbf{y} - \mathbf{1}\mu) \right\} \\ &= -\frac{1}{2} \left\{ \text{tr} \left( \Sigma^{-1} Z_A Z_A^\top \right) - \frac{1}{\sigma^2} (\mathbf{y} - \mu\mathbf{1})^t \Sigma^{-1} Z_A Z_A^\top \Sigma^{-1} (\mathbf{y} - \mu\mathbf{1}) \right\}\end{aligned}$$

where in the last equality we use the fact that

$$(S.2.5) \quad \Sigma = \lambda_A^2 Z_A Z_A^\top + \lambda_B^2 Z_B Z_B^\top + \sigma^2 M.$$

On equating to zero, we get from (S.2.3) that the optimal estimate of the location parameter is given by

$$(S.2.6) \quad \hat{\mu}_1 = \frac{\mathbf{1}^\top \Sigma^{-1} \mathbf{y}}{\mathbf{1}^\top \Sigma^{-1} \mathbf{1}},$$

the GLS estimate of  $\mu$ . If  $\hat{\mu}_1 \notin [\hat{a}_\tau, \hat{b}_\tau]$ ,  $\hat{\mu}$  takes the nearest boundary value in the set. From (S.2.4), we get the estimating equations for  $\lambda_A$ ,

$$(S.2.7) \quad \text{tr} \left( \Sigma^{-1} Z_A Z_A^\top \right) - \frac{1}{\sigma^2} (\mathbf{y} - \mathbf{1}\mu)^t \Sigma^{-1} Z_A Z_A^\top \Sigma^{-1} (\mathbf{y} - \mathbf{1}\mu) = 0.$$

By symmetry, taking the partial derivative w.r.t.  $\lambda_B^2$  gives

$$(S.2.8) \quad \text{tr} \left( \Sigma^{-1} Z_B Z_B^\top \right) - \frac{1}{\sigma^2} (\mathbf{y} - \mathbf{1}\mu)^t \Sigma^{-1} Z_B Z_B^\top \Sigma^{-1} (\mathbf{y} - \mathbf{1}\mu) = 0.$$

If  $\hat{\mu}_1 \in [\hat{a}_\tau, \hat{b}_\tau]$ , plugging (S.2.6) into (S.2.7) and (S.2.8) gives the estimating equations for  $\lambda_A^2$  and  $\lambda_B^2$  as

$$(S.2.9) \quad \text{tr} \left( \hat{\Sigma}^{-1} Z_A Z_A^\top \right) - \frac{1}{\sigma^2} \mathbf{y}^\top (I - P)^t \hat{\Sigma}^{-1} Z_A Z_A^\top \hat{\Sigma}^{-1} (I - P) \mathbf{y} = 0$$

$$(S.2.10) \quad \text{tr} \left( \hat{\Sigma}^{-1} Z_B Z_B^\top \right) - \frac{1}{\sigma^2} \mathbf{y}^\top (I - P)^t \hat{\Sigma}^{-1} Z_B Z_B^\top \hat{\Sigma}^{-1} (I - P) \mathbf{y} = 0$$

where  $P$  is the Generalized Least Square projection matrix:

$$(S.2.11) \quad P = \mathbf{1} (\mathbf{1}^\top \hat{\Sigma}^{-1} \mathbf{1})^{-1} \mathbf{1}^\top \hat{\Sigma}^{-1}.$$

**Estimating Equations for the URE method.** For URE estimates, note that in (10), in comparison to (S.2.2),  $\Sigma^{-1} M^2 V^{-1}$  replaces  $\Sigma^{-1}$ . Hence the partial derivative w.r.t.  $\mu$  vanishes for

$$(S.2.12) \quad \hat{\mu}_1 = \frac{\mathbf{1}^\top [\Sigma^{-1} M^2 V^{-1}] \mathbf{y}}{\mathbf{1}^\top [\Sigma^{-1} M^2 V^{-1}] \mathbf{1}}.$$

Again, if  $\hat{\mu}_1 \notin [\hat{a}_\tau, \hat{b}_\tau]$  it takes the nearest boundary value of the set. Furthermore,

$$\begin{aligned} \frac{\partial}{\partial \lambda_A^2} \text{URE} &\stackrel{(R10)}{=} -2\sigma^2 \text{tr} \left( \frac{\partial \Sigma^{-1}}{\partial \lambda_A^2} M^2 \right) + (\mathbf{y} - \mathbf{1}\mu)^t \left\{ \frac{\partial \Sigma^{-1}}{\partial \lambda_A^2} M^2 \Sigma^{-1} + \Sigma^{-1} M^2 \frac{\partial \Sigma^{-1}}{\partial \lambda_A^2} \right\} (\mathbf{y} - \mathbf{1}\mu) \\ &= -2\sigma^2 \text{tr} \left( \frac{\partial \Sigma^{-1}}{\partial \lambda_A^2} M^2 \right) + 2(\mathbf{y} - \mathbf{1}\mu)^t \left[ \frac{\partial \Sigma^{-1}}{\partial \lambda_A^2} M^2 \Sigma^{-1} \right] (\mathbf{y} - \mathbf{1}\mu) \\ &\stackrel{(R9)}{=} 2\sigma^2 \text{tr} \left( \Sigma^{-1} \frac{\partial \Sigma}{\partial \lambda_A^2} \Sigma^{-1} M^2 \right) - 2(\mathbf{y} - \mathbf{1}\mu)^t \left[ \Sigma^{-1} \frac{\partial \Sigma}{\partial \lambda_A^2} \Sigma^{-1} M^2 \Sigma^{-1} \right] (\mathbf{y} - \mathbf{1}\mu) \\ (S.2.13) \quad &= 2\sigma^2 \text{tr} (\Sigma^{-1} Z_A Z_A^\top \Sigma^{-1} M^2) - 2(\mathbf{y} - \mathbf{1}\mu)^t [\Sigma^{-1} Z_A Z_A^\top \Sigma^{-1} M^2 \Sigma^{-1}] (\mathbf{y} - \mathbf{1}\mu). \end{aligned}$$

Hence, on equating (S.2.4) to zero we obtain

$$(S.2.14) \quad \text{tr}(\Sigma^{-1}Z_A Z_A^T \Sigma^{-1}M^2) - \frac{1}{\sigma^2}(\mathbf{y} - \mathbf{1}\mu)^t[\Sigma^{-1}Z_A Z_A^T \Sigma^{-1}M^2 \Sigma^{-1}](\mathbf{y} - \mathbf{1}\mu) = 0.$$

By symmetry, equating the partial derivative w.r.t.  $\lambda_B^2$  to zero gives

$$(S.2.15) \quad \text{tr}(\Sigma^{-1}Z_B Z_B^T \Sigma^{-1}M^2) - \frac{1}{\sigma^2}(\mathbf{y} - \mathbf{1}\mu)^t[\Sigma^{-1}Z_B Z_B^T \Sigma^{-1}M^2 \Sigma^{-1}](\mathbf{y} - \mathbf{1}\mu) = 0.$$

If  $\hat{\mu}_1 \in [\hat{a}_\tau, \hat{b}_\tau]$ , plugging (S.2.12) into (S.2.14) and (S.2.15) gives the estimating equations for  $\lambda_A^2, \lambda_B^2$  as

$$(S.2.16) \quad \text{tr}\left(\hat{\Sigma}^{-1}Z_A Z_A^T \hat{\Sigma}^{-1}M^2\right) - \frac{1}{\sigma^2}\mathbf{y}^T(I - P)^t \Sigma^{-1}Z_A Z_A^T \hat{\Sigma}^{-1}M^2 \hat{\Sigma}^{-1}(I - P)\mathbf{y} = 0$$

$$(S.2.17) \quad \text{tr}\left(\hat{\Sigma}^{-1}Z_B Z_B^T \hat{\Sigma}^{-1}M^2\right) - \frac{1}{\sigma^2}\mathbf{y}^T(I - P)^t \Sigma^{-1}Z_B Z_B^T \hat{\Sigma}^{-1}M^2 \hat{\Sigma}^{-1}(I - P)\mathbf{y} = 0$$

where  $P$  is given by:

$$(S.2.18) \quad P = \mathbf{1}(\mathbf{1}^T \hat{\Sigma}^{-1}M^2 \hat{\Sigma}^{-1} \mathbf{1})^{-1} \mathbf{1}^T \hat{\Sigma}^{-1}M^2 \hat{\Sigma}^{-1}.$$

**Proof of Lemma 2.4.** This is an immediate consequence of Theorem 5 in Searle (1966), because  $\boldsymbol{\eta}_c = Z_c \boldsymbol{\theta}$  is estimable if and only if  $v^T \boldsymbol{\theta}$  is estimable for each row  $v$  of  $Z_c$ .

LEMMA S.2.3. *An unbiased estimator of the generalized risk  $R_{r,c}^Q(\boldsymbol{\eta}, \hat{\boldsymbol{\eta}}(\mu, \lambda_A, \lambda_B))$  of estimators of the form  $\hat{\boldsymbol{\eta}}(\mu, \lambda_A, \lambda_B)$  is given by*

$$\begin{aligned} \text{URE}^Q(\mu, \lambda_A, \lambda_B) &= \sigma^2 \text{tr}(QM) - 2\sigma^2 \text{tr}(\Sigma^{-1}MQM) \\ &\quad + (\mathbf{y} - \mu \mathbf{1})^t [\Sigma^{-1}MQM \Sigma^{-1}](\mathbf{y} - \mu \mathbf{1}). \end{aligned}$$

**Proof.** Similar to Lemma S.2.2.

### S.3. Section 3: Proof details and further insights.

In this section, we provide additional details for Section 3 and its associated appendix A.1. We first prove the following interesting property of the  $Q$  matrix defined in (15).

LEMMA S.3.1. *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 S.3.1.** 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  $Z_c^T Z_c$ .

We next prove the seven inequalities that we had used in our proofs in Appendix A.1. Then, we provide proof of Lemma A.4. Thereafter, detailed proofs for lemmas A.2, A.3, A.5 used in the appendix for proving the results of Section 3 are provided. We end this section by providing insights for the weighted loss case and with some remarks on the assumptions used in our theory.

LEMMA S.3.2. *With matrices  $G, H$  defined in Appendix A.1 and  $M, \Sigma^{-1}$  and  $Q$  as defined in Section 2.3 and  $\tilde{Q} = M^{\frac{1}{2}} Q M^{\frac{1}{2}}$  we have the following inequalities:*

$$(S.3.1) \quad \lambda_1(HMH) \leq \lambda_1(M^{-1})\lambda_1^2(\tilde{Q}) .$$

$$(S.3.2) \quad \text{tr}(HMHM) \leq rc\lambda_1^2(\tilde{Q}) .$$

$$(S.3.3) \quad \lambda_1(QGMG^T Q) \leq \lambda_1(M^{-1})\lambda_1^2(\tilde{Q}) .$$

$$(S.3.4) \quad \sigma_1(H_1) \leq \lambda_1(M^{-1})\lambda_1^2(\tilde{Q}) .$$

$$(S.3.5) \quad \lambda_1(H_3) \leq \lambda_1(M^{-1})\lambda_1^2(\tilde{Q}) .$$

$$(S.3.6) \quad \text{tr}(\tilde{G}M\tilde{G}M) \leq 4rc\lambda_1^2(\tilde{Q}) \text{ where } \tilde{G} = QG + G^T Q .$$

$$(S.3.7) \quad \lambda_1(G^T Q M Q G) \leq \lambda_1(M^{-1})\lambda_1^2(\tilde{Q}) .$$

In Appendix A.1, the above bounds are used with  $\lambda_1(\tilde{Q}) \leq \lambda_1(Q)$  and  $\nu_{r,c} = \lambda_1(M^{-1})$ .

**Proof of Lemma S.3.2.** For the first inequality note that

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

The last inequality above uses  $W^2 \preceq I$ . Again, by R6 of Section S.6, the RHS above equals

$\lambda_1(M^{\frac{1}{2}}QM\Sigma^{-1}\Sigma^{-1}MQM^{\frac{1}{2}})$ . Thus, we have

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

where, the inequality above follows by using  $WM^{-1}W \preceq \lambda_1(M^{-1})I$ .

The second inequality stated in the lemma follows as

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

where the last equation above follows by using R6 again.

The third inequality uses

$$\begin{aligned}\lambda_1(QGMG^TQ) &= \lambda_1(QM\Sigma^{-1}M\Sigma^{-1}MQ) = \lambda_1(QM^{\frac{1}{2}}W^2M^{\frac{1}{2}}Q) \\ &\leq \lambda_1(QM^{\frac{1}{2}}M^{\frac{1}{2}}Q) \leq \lambda_1(M^{-1})\lambda_1^2(M^{\frac{1}{2}}QM^{\frac{1}{2}}).\end{aligned}$$

For the proof of the forth inequality 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 above uses  $W \preceq I$ . For the fifth inequality:

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

For the sixth inequality, denote  $\dot{G} = QG$ . Thus,  $\tilde{G} = \dot{G} + \dot{G}^\top$ , and so,

$$(S.3.8) \quad \text{tr}(\tilde{G}M\tilde{G}M) = \text{tr}(\dot{G}M\dot{G}M) + \text{tr}(\dot{G}^\top M\dot{G}^\top M) + 2\text{tr}(\dot{G}M\dot{G}^\top M).$$

Substituting the expression of  $\dot{G}$  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

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

for any  $\Sigma^{-1}$ . For the second term in (S.3.8) note that  $\text{tr}(\dot{G}^\top M\dot{G}^\top M) = \text{tr}(M\dot{G}M\dot{G}) = \text{tr}(\dot{G}M\dot{G}M)$ . 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

$$\text{tr}(\dot{G}M\dot{G}^\top M) \leq \text{tr}(QMQM) = \text{tr}(M^{\frac{1}{2}}QM^{\frac{1}{2}})^2 = rc \cdot \lambda_1^2(M^{\frac{1}{2}}QM^{\frac{1}{2}})$$

and the result is proved.

Finally, the last inequality stated in this lemma follows as

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

**Proof of Lemma A.4.** Without loss of generality we can assume that there are no missing cells. As  $\hat{a}_\tau, \hat{b}_\tau$  is the  $\tau/2$  th and  $(1 - \tau/2)$  th quantile of  $\mathbf{y}$ :

$$\begin{aligned} \max(|\hat{a}_\tau|, |\hat{b}_\tau|) &\leq \mathbf{quantile}(|y_{ij}| : (i, j) \in \mathcal{E}; 1 - \tau/2) \\ &= \mathbf{quantile}(|\eta_{ij}| + |\epsilon_{ij}| : (i, j) \in \mathcal{E}; 1 - \tau/2) \end{aligned}$$

where  $\epsilon_{ij}$  are i.i.d. standard normal variables. The RHS is bounded above by:

$$\max \{ |\eta_{ij}| + |\epsilon_{ij}| : (i, j) \in \mathcal{E} \text{ and } |\eta_{ij}| \leq q_\tau(|\eta|), |\epsilon_{ij}| \leq q_\tau(|\epsilon|) \} \leq q_\tau(|\eta|) + q_\tau(|\epsilon|),$$

where  $q_\tau(|\eta|) = \mathbf{quantile}(|\eta_{ij}| : (i, j) \in \mathcal{E}; 1 - \tau/2)$  and  $q_\tau(|\epsilon|) = \mathbf{quantile}(|\epsilon_{ij}| : (i, j) \in \mathcal{E}, 1 - \tau/2)$ . Thus,

$$\max(|\hat{a}_\tau|, |\hat{b}_\tau|) \leq q_\tau(|\eta|) + q_\tau(|\epsilon|).$$

Again,

$$q_\tau(|\eta|) \leq \max\{1, \mathbf{quantile}(\eta_{ij}^2 : (i, j) \in \mathcal{E}; 1 - \tau/2)\} \leq \max\left\{1, \frac{1}{\tau/2 \cdot RC} \sum_{i,j} \eta_{ij}^2\right\} < \infty$$

which follows from Assumption A1. The second inequality above is due to the fact that the highest possible value of the  $1 - \tau/2$  quantile of a series of positive numbers with a constraint on their sum is attained when all the values above that quantile are all same.

Also, as sample quantiles are asymptotically normally distributed we have:

$$(rc)^{1/2} \cdot (q_\tau(|\eta|) - x_0) \sim N(0, 8^{-1}\tau(1 - \tau/2)\phi^{-2}(x_0)) \text{ where } x_0 = \Phi^{-1}(1 - \tau/4).$$

Thus, we have  $P(\max(|\hat{a}_\tau|, |\hat{b}_\tau|) \leq \log(rc)) \rightarrow 1$  as  $r, c \rightarrow \infty$ . This, completes the proof of the lemma.

**Proof of Lemma A.2** With a slight abuse of notation, we use  $L^Q(\boldsymbol{\eta}, \hat{\boldsymbol{\eta}}(\mu, \tilde{\lambda}_A, \tilde{\lambda}_B))$  to denote the loss  $L^Q(\boldsymbol{\eta}, \hat{\boldsymbol{\eta}}(\mu, \lambda_A, \lambda_B)) = L^Q(\boldsymbol{\eta}, \hat{\boldsymbol{\eta}}(\mu, 1/\tilde{\lambda}_A^2 - 1, 1/\tilde{\lambda}_B^2 - 1))$ . As  $L^Q(\boldsymbol{\eta}, \hat{\boldsymbol{\eta}}(\mu, \tilde{\lambda}_A, \tilde{\lambda}_B))$  is everywhere differentiable, for any triplet  $(\mu, \tilde{\lambda}_A, \tilde{\lambda}_B)$  and any point  $(\mu[i], \tilde{\lambda}_A[j], \tilde{\lambda}_B[k])$  on the grid  $\Theta_{r,c}$  we have:

$$\begin{aligned} &|L^Q(\boldsymbol{\eta}, \hat{\boldsymbol{\eta}}(\mu, \tilde{\lambda}_A, \tilde{\lambda}_B)) - L^Q(\boldsymbol{\eta}, \hat{\boldsymbol{\eta}}(\mu[i], \tilde{\lambda}_A[j], \tilde{\lambda}_B[k]))| \\ &\leq D_{r,c}^{[1]} \cdot |\mu - \mu[i]| + D_{r,c}^{[2]} \cdot |\tilde{\lambda}_A - \tilde{\lambda}_A[j]| + D_{r,c}^{[3]} \cdot |\tilde{\lambda}_B - \tilde{\lambda}_B[k]|, \end{aligned}$$

where,

$$(S.3.9) \quad D_{r,c}^{[1]}(\boldsymbol{\eta}, \mathbf{y}) = \sup_{|\mu| \leq m_{r,c}; \tilde{\lambda}_A, \tilde{\lambda}_B \in [0,1]} \left| \frac{\partial}{\partial \mu} L^Q(\boldsymbol{\eta}, \hat{\boldsymbol{\eta}}(\mu, \tilde{\lambda}_A, \tilde{\lambda}_B)) \right|,$$

$$(S.3.10) \quad D_{r,c}^{[2]}(\boldsymbol{\eta}, \mathbf{y}) = \sup_{|\mu| \leq m_{r,c}; \tilde{\lambda}_A, \tilde{\lambda}_B \in [0,1]} \left| \frac{\partial}{\partial \tilde{\lambda}_A} L^Q(\boldsymbol{\eta}, \hat{\boldsymbol{\eta}}(\mu, \tilde{\lambda}_A, \tilde{\lambda}_B)) \right| \text{ and,}$$

$$(S.3.11) \quad D_{r,c}^{[3]}(\boldsymbol{\eta}, \mathbf{y}) = \sup_{|\mu| \leq m_{r,c}; \tilde{\lambda}_A, \tilde{\lambda}_B \in [0,1]} \left| \frac{\partial}{\partial \tilde{\lambda}_B} L^Q(\boldsymbol{\eta}, \hat{\boldsymbol{\eta}}(\mu, \tilde{\lambda}_A, \tilde{\lambda}_B)) \right|.$$

Thus, based on the construction of the grid  $\Theta_{r,c}$  we have for any triplet  $(\mu, \lambda_A, \lambda_B) \in [-m_{r,c}, m_{r,c}] \otimes [0, \infty] \otimes [0, \infty]$ :

$$(S.3.12) \quad \inf_{(\mu[i], \lambda_A[j], \lambda_B[k]) \in \Theta_{r,c}} |L^Q(\boldsymbol{\eta}, \hat{\boldsymbol{\eta}}(\mu, \lambda_A, \lambda_B)) - L^Q(\boldsymbol{\eta}, \hat{\boldsymbol{\eta}}(\mu[i], \tilde{\lambda}_A[j], \tilde{\lambda}_B[k]))|$$

$$(S.3.13) \quad \leq D_{r,c}^{[1]}(\boldsymbol{\eta}, \mathbf{y}) \cdot \delta_{r,c}^{[1]} + D_{r,c}^{[2]}(\boldsymbol{\eta}, \mathbf{y}) \cdot \delta_{r,c}^{[2]} + D_{r,c}^{[3]}(\boldsymbol{\eta}, \mathbf{y}) \cdot \delta_{r,c}^{[3]} = D_{r,c}(\boldsymbol{\eta}, \mathbf{y}) \text{ (say).}$$

Thus, on the set  $A_{r,c}(\mathbf{Y}) = \{[\hat{a}_\tau, \hat{b}_\tau] \subseteq [-m_{r,c}, m_{r,c}]\}$  we have:

$$|L^Q(\boldsymbol{\eta}, \hat{\boldsymbol{\eta}}(\mu^{\text{OD}}, \lambda_A^{\text{OD}}, \lambda_B^{\text{OD}})) - L^Q(\boldsymbol{\eta}, \hat{\boldsymbol{\eta}}(\mu^{\text{OL}}, \lambda_A^{\text{OL}}, \lambda_B^{\text{OL}}))| \leq D_{r,c}(\boldsymbol{\eta}, \mathbf{y}).$$

By the construction of  $\Theta_{r,c}$  as shown afterwards in Lemma S.3.3 we have:  $\mathbb{E}[D_{r,c}(\boldsymbol{\eta}, \mathbf{y}) I\{A_{r,c}(\mathbf{Y})\}] \rightarrow 0$  as  $r, c \rightarrow \infty$  under assumptions A1-A2. It implies by Markov's inequality that  $P(D_{r,c}(\boldsymbol{\eta}, \mathbf{y}) > \epsilon \text{ and } A_{r,c}(\mathbf{Y})) \rightarrow 0$  as  $r, c \rightarrow \infty$ . These coupled with Lemmas A.4 and A.5 provide us the results A and B of the lemma.

Again, note that by definition (27), on the set  $A_{r,c}(\mathbf{Y})$  we have:

$$\begin{aligned} & |L^Q(\boldsymbol{\eta}, \hat{\boldsymbol{\eta}}(\mu^{\text{UD}}, \lambda_A^{\text{UD}}, \lambda_B^{\text{UD}})) - L^Q(\boldsymbol{\eta}, \hat{\boldsymbol{\eta}}(\mu^{\text{URE}}, \lambda_A^{\text{URE}}, \lambda_B^{\text{URE}}))| \\ & \leq D_{r,c}^{[1]}(\boldsymbol{\eta}, \mathbf{y}) \cdot \delta_{r,c}^{[1]} + D_{r,c}^{[2]}(\boldsymbol{\eta}, \mathbf{y}) \cdot \delta_{r,c}^{[2]} + D_{r,c}^{[3]}(\boldsymbol{\eta}, \mathbf{y}) \cdot \delta_{r,c}^{[3]} = D_{r,c}(\boldsymbol{\eta}, \mathbf{y}). \end{aligned}$$

and so the results C and D of the lemma follows using Lemma A.4 and result B of Lemma A.5.

LEMMA S.3.3. *With  $D_{r,c}(\boldsymbol{\eta}, \mathbf{y})$  defined in (S.3.9)-(S.3.12), for any  $\boldsymbol{\eta}$  obeying assumption A1 and under assumption A2 on the design we have:*

$$\mathbb{E}[D_{r,c}(\boldsymbol{\eta}, \mathbf{y})] \rightarrow 0 \text{ as } r, c \rightarrow \infty.$$

**Proof of Lemma S.3.3.** First, note that the quadratic loss is

$$L^Q(\boldsymbol{\eta}, \hat{\boldsymbol{\eta}}(\mu, \lambda_A, \lambda_B)) = (rc)^{-1}(\boldsymbol{\eta} - \mathbf{y} + G\mathbf{y} - \mu G\mathbf{1})^T Q(\boldsymbol{\eta} - \mathbf{y} + G\mathbf{y} - \mu G\mathbf{1}),$$

where  $G = M\Sigma^{-1}$  and  $\Sigma = (\lambda_A Z_A Z_A^\top + \lambda_B Z_B Z_B^\top + M)$  involves the scale parameters. Differentiating the loss with respect to  $\mu$  we have:

$$\begin{aligned} \frac{\partial}{\partial \mu} L^Q(\boldsymbol{\eta}, \hat{\boldsymbol{\eta}}(\mu, \lambda_A, \lambda_B)) &= (rc)^{-1} \frac{\partial}{\partial \mu} \{ \mu^2 \mathbf{1}^\top G^\top Q G \mathbf{1} - 2\mu \mathbf{1}^\top G^\top Q(\boldsymbol{\eta} - \mathbf{y} + G\mathbf{y}) \} \\ &= (rc)^{-1} (2\mu \mathbf{1}^\top G^\top Q G \mathbf{1} - 2\mathbf{1}^\top G^\top Q(\boldsymbol{\eta} - \mathbf{y} + G\mathbf{y})). \end{aligned}$$

Note that,

$$(rc)^{-1} |\mu| \mathbf{1}^\top G^\top Q G \mathbf{1} \leq m_{r,c} \lambda_1(H) \text{ where } H = G^\top Q G$$

and by calculations in Section A.1 of the appendix it follows that  $\lambda_1(H) \leq \nu_{r,c} \lambda_1(Q)$  for any  $\lambda_A, \lambda_B \geq 0$ .

Also,  $\mathbf{1}^\top G^\top Q(\boldsymbol{\eta} - \mathbf{y} + G\mathbf{y}) \sim N(\mathbf{1}^\top H \boldsymbol{\eta}, \mathbf{1}^\top G^\top Q(I - G^\top)M(I - G)G\mathbf{1})$  and by moment calculations similar to Section A.1 we have:

$$(rc)^{-1} \mathbb{E}\{|\mathbf{1}^\top G^\top Q(\boldsymbol{\eta} - \mathbf{y} + G\mathbf{y})|\} \leq O(\nu_{r,c} \lambda_1(Q)) \text{ for any } \lambda_A, \lambda_B \geq 0.$$

Therefore,  $D_{r,c}^{[1]}(\boldsymbol{\eta}, \mathbf{y}) \leq O(m_{r,c} \nu_{r,c} \lambda_1(Q))$  and so,  $\mathbb{E}\{D_{r,c}^{[1]}(\boldsymbol{\eta}, \mathbf{y}) \delta_{r,c}^{[1]}\} \rightarrow 0$  as  $r, c \rightarrow \infty$ .

Now, we concentrate on the scale hyper-parameters. Differentiating the loss with respect to  $\lambda_A$  we have:

$$\begin{aligned} \frac{\partial}{\partial \lambda_A} L^Q(\boldsymbol{\eta}, \hat{\boldsymbol{\eta}}(\mu, \lambda_A, \lambda_B)) &= (\mathbf{y} - \mu \mathbf{1})^\top \frac{\partial(G^\top Q G)}{\partial \lambda_A} (\mathbf{y} - \mu \mathbf{1}) + 2(\mathbf{y} - \mu \mathbf{1})^\top \frac{\partial G^\top}{\partial \lambda_A} Q(\boldsymbol{\eta} - \mathbf{y}), \\ \text{where, } \frac{\partial}{\partial \lambda_A} (G^\top Q G) &= \frac{\partial G^\top}{\partial \lambda_A} Q G + G^\top Q \frac{\partial G}{\partial \lambda_A} \quad \text{and} \quad \frac{\partial G}{\partial \lambda_A} = M \Sigma^{-1} Z_A Z_A^\top \Sigma^{-1}. \end{aligned}$$

Again, note that for the transformed scale hyper-parameter  $\tilde{\lambda}_A$ :

$$\begin{aligned} \frac{\partial}{\partial \tilde{\lambda}_A} L^Q(\boldsymbol{\eta}, \hat{\boldsymbol{\eta}}(\mu, \tilde{\lambda}_A, \tilde{\lambda}_B)) &= \frac{\partial}{\partial \lambda_A} L^Q(\boldsymbol{\eta}, \hat{\boldsymbol{\eta}}(\mu, \lambda_A, \lambda_B)) \times \frac{\partial \lambda_A}{\partial \tilde{\lambda}_A} \\ &= -2(1 + \lambda_A)^{3/2} \frac{\partial}{\partial \lambda_A} L^Q(\boldsymbol{\eta}, \hat{\boldsymbol{\eta}}(\mu, \lambda_A, \lambda_B)). \end{aligned}$$

Note that the change of scale to  $\tilde{\lambda}$  was chosen cleverly such that not only the range of  $\tilde{\lambda}$  is bounded but also the subsequent change in scale does not lead to the derivative to blow up as  $\lambda_A$  varies over 0 to  $\infty$ . As such:

$$\begin{aligned} -\frac{1}{2} M^{-1} \frac{\partial G}{\partial \tilde{\lambda}_A} &= (1 + \lambda_A)^{3/2} \Sigma^{-1} Z_A Z_A^\top \Sigma^{-1} \\ &\preceq \left( \lambda_A (1 + \lambda_A)^{-3/4} Z_A Z_A^\top + \lambda_B (1 + \lambda_A)^{-3/4} Z_B Z_B^\top + (1 + \lambda_A)^{-3/4} M \right)^{-1}. \end{aligned}$$

As  $\lambda_A \rightarrow \infty$ ,  $(1 + \lambda_A)^{-3/4} M$  becomes negligible but  $\lambda_A (1 + \lambda_A)^{-3/4} Z_A Z_A^\top$  contributes massively and using moment calculations similar to Section A.1 of the Appendix, it can

be shown that:  $\mathbb{E}\{D_{r,c}^{[2]}(\boldsymbol{\eta}, \mathbf{y})\} \leq O(m_{r,c}^2 \nu_{r,c} \lambda_1(Q))$ . Similar, calculations hold for the other scale hyper-parameter. Combining the bounds on the three hyper-parameters, we get:  $\mathbb{E}[D_{r,c}(\boldsymbol{\eta}, \mathbf{y})] \rightarrow 0$  as  $r, c \rightarrow \infty$ .

**Proof of Lemma A.3** The proof is very similar to that of Lemma A.2 and is avoided here to prevent repetition.

**Proof of Lemma A.5** To prove the  $L_1$  convergence results of the lemma, we apply Cauchy-Schwarz inequality and convert our problem to showing convergence of the products of the respective expected values. As such,

$$\begin{aligned} \mathbb{E}\{|L_{r,c}(\boldsymbol{\eta}_c, \tilde{\boldsymbol{\eta}}_c^{\text{OD}}) - L_{r,c}(\boldsymbol{\eta}_c, \tilde{\boldsymbol{\eta}}_c^{\text{OL}})| \cdot I\{A_{r,c}(\mathbf{Y})\}\} &\leq 2\mathbb{E}\{|L_{r,c}(\boldsymbol{\eta}_c, \tilde{\boldsymbol{\eta}}_c^{\text{OD}})| I\{A_{r,c}(\mathbf{Y})\}\} \\ &\leq 2\{\mathbb{E}\{L_{r,c}^Q(\boldsymbol{\eta}, \tilde{\boldsymbol{\eta}}^{\text{OD}})\}^2 P(A_{r,c}(\mathbf{Y}))\}^{1/2}. \end{aligned}$$

Based on the calculations made in the proof of Lemma A.4, it follows that  $P(A_{r,c}(\mathbf{Y})) = O((rc)^{-1})$ . Using moment bounding techniques used in Section A.1, under assumptions A1 and A2, it can be shown that  $(rc)^{-1}\mathbb{E}\{L_{r,c}^Q(\boldsymbol{\eta}, \tilde{\boldsymbol{\eta}}^{\text{OD}})\}^2$ ,  $(rc)^{-1}\mathbb{E}\{L_{r,c}^Q(\boldsymbol{\eta}, \hat{\boldsymbol{\eta}}^{\text{URE}})\}^2$  and  $(rc)^{-1}\mathbb{E}\{L_{r,c}^Q(\boldsymbol{\eta}, \hat{\boldsymbol{\eta}}^{\text{UD}})\}^2$  all converges to 0 as  $r, c \rightarrow \infty$  which will complete the proof of the lemma.

S.3.1. *Brief Outline of the results for the Weighted loss case.* We now briefly discuss estimation under weighted loss  $L_{r,c}^{\text{wgt}}(\boldsymbol{\eta}, \hat{\boldsymbol{\eta}})$  defined in Section 2. For simplicity, we describe the case where there are no unobserved cells. Under this weighted loss, applying the following linear transformation

$$\tilde{\mathbf{y}} = M^{-1/2}\mathbf{y}, \quad \tilde{\boldsymbol{\eta}} = M^{-1/2}\boldsymbol{\eta}, \quad \tilde{Z} = M^{-1/2}Z, \quad \tilde{\boldsymbol{\mu}}\mathbf{1} = M^{-1/2}\boldsymbol{\mu}\mathbf{1}$$

the problem reduces to estimating  $\tilde{\boldsymbol{\eta}}$  from  $\tilde{\mathbf{y}} \sim N(\tilde{\boldsymbol{\eta}}, \sigma^2 I)$  under the usual sum-of-squares loss. As the problem can be converted into a homoskedastic case, estimation here is easier than the cases discussed before. Assuming the hierarchical Gaussian prior structure like before, the complete Bayes model is given by:

$$\tilde{\boldsymbol{\eta}} \sim N_{rc}(1\tilde{\boldsymbol{\mu}}, \sigma^2 M^{-1/2} Z \Lambda \Lambda^\top Z^\top M^{-1/2})$$

and the corresponding Bayes estimate of  $\tilde{\boldsymbol{\eta}}$  is

$$\hat{\boldsymbol{\eta}} = \tilde{\mathbf{y}} - \tilde{V}^{-1}(\tilde{\mathbf{y}} - 1\tilde{\boldsymbol{\mu}}), \quad \text{where } \tilde{V} = M^{-1/2} Z \Lambda \Lambda^\top Z^\top M^{-1/2} + I$$

which unlike the shrinkage matrix in (5) is symmetric. The oracle optimality proof can be worked out following in verbatim the proofs with the  $L^Q$  loss. However, in this case due to the presence of symmetric shrinkage matrix, the estimation problem reduces to the easier situation when  $\nu_{r,c} = 1$ .

S.3.2. *Discussions on the relevance of the Assumptions made.* Here, we discuss the genesis of Assumption A2 in our asymptotic optimality proofs. Our Assumption A1 is not very restrictive and so discussions on it is avoided here. On the other hand, assumption A2 put an asymptotic control on the imbalance in our design matrix as  $r, c \rightarrow \infty$ . It is peculiar to the two-way nature of the problem and is not usually seen in the huge literature around shrinkage estimation of the normal mean in the one-way problem.

Assumption A2 is needed in several parts of our proof. Let us concentrate on Lemma 3.1 which shows that our risk estimation strategy indeed approximates the true risk uniformly well for estimators with the location hyper-parameter  $\mu$  set at 0. By equation (23), the approximation error was exacted evaluated to be:

$$\mathbb{E}\left\{\text{URE}_{r,c}^Q(0, \lambda_A, \lambda_B) - R_{r,c}^Q(\boldsymbol{\eta}, \widehat{\boldsymbol{\eta}}(0, \lambda_A, \lambda_B))\right\}^2 = (rc)^{-2}\{2\text{tr}(HMHM) + 4\boldsymbol{\eta}^tHMH\boldsymbol{\eta}\}.$$

We need to show that the RHS is  $o(d_{r,c}^2)$  uniformly over any choices of the scale hyper-parameters and for all  $\boldsymbol{\eta}$  satisfying Assumption A1. Recall,  $d_{r,c}^2$  rate of control of the square error was needed due to the discretization process. We concentrate on the component  $\boldsymbol{\eta}^tHMH\boldsymbol{\eta}$ . Based on the equality condition on the von-Neumann trace inequality we can say that

$$\boldsymbol{\eta}^tHMH\boldsymbol{\eta} = \text{tr}(\{HMH\}\{\boldsymbol{\eta}\boldsymbol{\eta}^T\}) = \lambda_1(HMH)\boldsymbol{\eta}^T\boldsymbol{\eta}$$

when the eigen vector corresponding to the largest eigen value of  $HMH$  matches  $\boldsymbol{\eta}/\boldsymbol{\eta}^T\boldsymbol{\eta}$ . This, can indeed happen as for uniform convergence we not only have to consider all possible values  $\boldsymbol{\eta}$  but also all possible values of the  $H$  matrix as  $\lambda_A, \lambda_B$  changes. To simplify further let us assume  $Q = I$ . We now provide heuristic reasons why  $\lambda_1(HMH)$  can be close to the upper bound  $\lambda_1^{-1}(M)$  that we use for it in our proofs. As shown before:

$$\lambda_1(HMH) = \sigma_1^2(M^{-1/2}WM).$$

Now,  $M$  is a diagonal matrix with  $0 \prec M \preceq I$  and  $0 \preceq W \preceq I$ .  $W$  depends on  $\lambda_A, \lambda_B$  as they vary over  $[0, \infty]^2$ . We relax the range and consider  $M$  and  $W$  to be any possible p.d. diagonal matrix and n.n.d. matrix respectively. It is difficult to gauge the degree of tightness of this relaxation as  $M$  and  $W$  are related, but we can expect them to be close as  $\lambda_A$  and  $\lambda_B$  span over the entire first quadrant. Simplifying the scenario further assume a  $2 \times 2$  situation where

$$M = \begin{bmatrix} 1 & 0 \\ 0 & b \end{bmatrix} \quad \text{and} \quad W = \begin{bmatrix} w_{11} & w_{12} \\ w_{12} & w_{22} \end{bmatrix}.$$

where  $b \in (0, 1]$  and  $w_{11}, w_{12}$  and  $w_{22}$  are chosen such that  $0 \preceq W \preceq I$ . Thus,  $(M^{-1/2}WM)(M^{-1/2}WM)^t$  is given by:

$$\begin{bmatrix} c_{11} = w_{11}^2 + b^4 a_{12}^2 & c_{12} = b^{-1} w_{11} w_{21} + b^3 w_{22} w_{12} \\ c_{12} & c_{22} = b^{-2} w_{12}^2 + b^2 w_{22}^2 \end{bmatrix}$$

and its eigenvalues are given by:

$$2^{-1} \{ (c_{11} + c_{22}) \pm \sqrt{(c_{11} + c_{22})^2 + 4c_{12}^2} \}.$$

We would like to evaluate the maximum of the eigenvalue as  $b$  decreases. We consider finding the eigenvalue asymptotically as  $b \rightarrow 0$ . Under the asymptotic regime  $b \rightarrow 0$ , we have:

$$c_{11} \sim w_{11}^2; \quad c_{12} \sim b^{-1} w_{11} w_{21}, \quad \text{and} \quad c_{22} \sim b^{-2} w_{12}^2.$$

Thus, for any fixed positive value of  $w_{11}, w_{12}$  the highest eigenvalue is of the order of  $b^{-1} = \lambda_1(M^{-1})$  as  $b$  approaches zero.

#### S.4. Details for the empirical studies of Section 4.

##### S.4.1. Details for simulation studies.

(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 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  $\sigma^2 = 100$ . 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  $\sigma^2 = 100$ . 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  $\sigma^2 = 100$ . 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.

Next we describe a simulation study that was mentioned in the main article but not included in Figure 1. Here we evaluate the performance of the different estimators under misspecification of the model. The setup is analogous to that of simulation example (a), but an interaction term is added to the main effects. Thus, now  $y_{ij} \sim N(\mu + \alpha_i + \beta_j + \gamma_{ij}, \sigma^2 K_{ij}^{-1})$  where  $\alpha_i, \beta_j \sim N(0, \sigma^2/(4l))$ ,  $\gamma_{ij} \sim N(0, \sigma^2/4)$  and  $K_{ij}$  are drawn independently; while the estimators (and the oracle) are based on the additive model (no interactions). Figure 1 below displays the results. When the number of rows (and columns) is sufficiently large, the risk of each of the estimators is almost constant, and all but the LS estimator achieve approximately the oracle risk. We tried to fit also the “full” LS estimator, which includes interaction terms, and its risk was significantly higher in our simulation than that of the additive LS estimator (and the other estimators compared here).

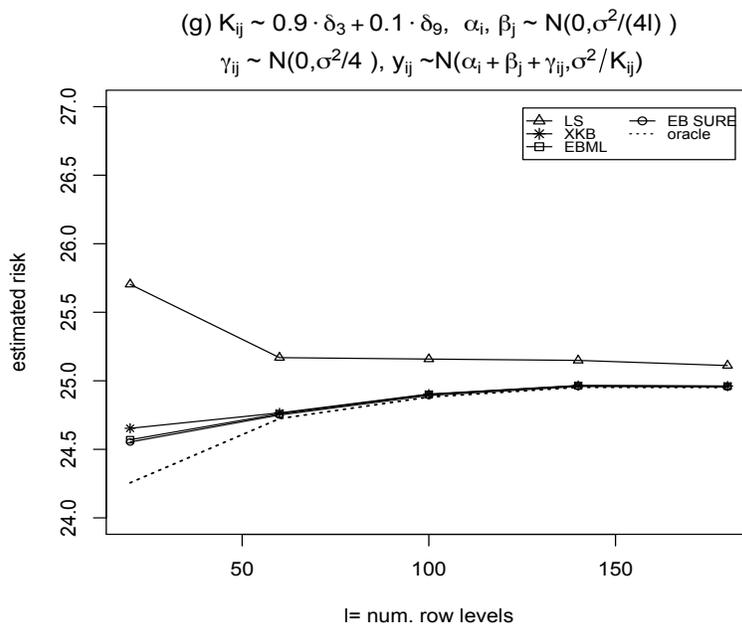


FIG 1. Risk of the various estimators in a simulation under misspecification.

S.4.2. Details for 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

$$(S.4.1) \quad y_{ijk} = \eta_{ij} + \epsilon_{ijk}, \quad \eta_{ij} = \mu + \alpha_i + \beta_j \quad k = 1, \dots, K_{ij}$$

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 (S.4.1).

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 1 shows the estimates of the relative variance components  $\lambda_A$  and  $\lambda_B$ , corresponding to hour and county, respectively, as well as the estimates of the fixed term  $\mu$ , for each of the shrinkage estimators. There is a

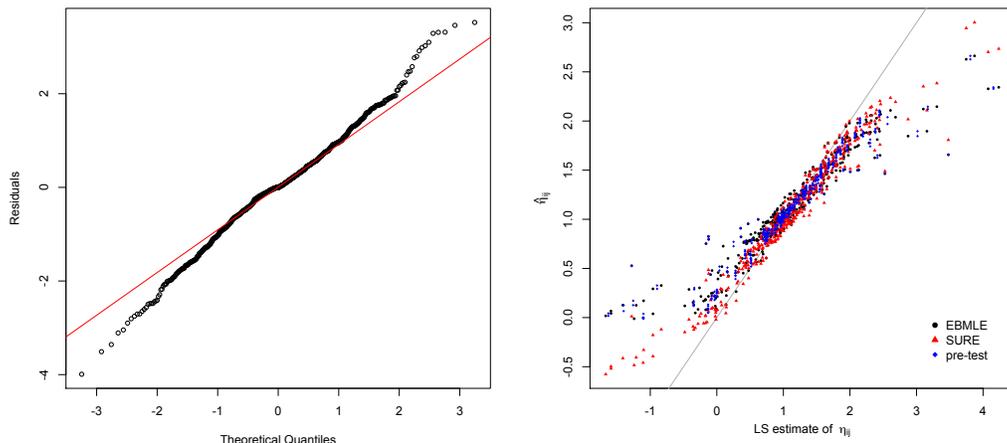


FIG 2. *Left: Normal Q-Q plot for the residuals of the LS fit to the two-way model for water data. Right: Plot of Shrinkage estimates vs. LS estimates of the cell means. The horizontal coordinate is the LS estimate and the vertical coordinate is an alternative estimate: EBMLE, URE or LS-county. EBMLE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

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 1 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 1

*Estimated fixed effect ( $\mu$ ) and components  $\lambda_{county}$  and  $\lambda_{hour}$ , which determine shrinkage. This is the same as Table 2 in the main article.*

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 (S.4.1) 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 (S.4.1),  $\hat{\boldsymbol{\eta}}_c^{\text{LS}(2)}$  is an unbiased estimator of  $\boldsymbol{\eta}_c$  and

$$(S.4.2) \quad \text{SSPE}[\hat{\boldsymbol{\eta}}_c^{(1)}] = \|\hat{\boldsymbol{\eta}}_c^{(1)} - \hat{\boldsymbol{\eta}}_c^{\text{LS}(2)}\|^2$$

is the Sum of Squared Prediction Error of  $\hat{\boldsymbol{\eta}}_c^{(1)}$ . Instead of averaging (S.4.2) 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 (S.4.1) 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 2 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 2

*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 (S.4.1), 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 (S.4.1) 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 2 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.

**S.5. URE in the balanced design.** 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

$$(S.5.1) \quad m = \mu + \alpha. + \beta., \quad a_i = \alpha_i - \alpha., \quad b_j = \beta_j - \beta.$$

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

$\widehat{\boldsymbol{\eta}}(y_{..}, \lambda_A, \lambda_B)$ , obtained by substituting the mean of  $\mathbf{y}$  for  $\boldsymbol{\mu}$  in (5), is

$$(S.5.2) \quad \widehat{\boldsymbol{\eta}}_{ij}(y_{..}, \lambda_A, \lambda_B) = \widehat{m}^{\text{LS}} + c_\alpha(\lambda_A) \widehat{a}_i^{\text{LS}} + c_\beta(\lambda_B) \widehat{b}_j^{\text{LS}},$$

where  $\widehat{m}^{\text{LS}} = y_{..}$ ,  $\widehat{\mathbf{a}}_i^{\text{LS}} = y_{i.} - y_{..}$ ,  $\widehat{\mathbf{b}}_j^{\text{LS}} = y_{.j} - y_{..}$  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 decomposes as  $R(\boldsymbol{\eta}, \widehat{\boldsymbol{\eta}}(y_{..}, \lambda_A, \lambda_B))$  equals

$$(S.5.3) \quad \frac{1}{rc} \mathbb{E} \left\{ \sum_{i=1}^r \sum_{j=1}^c [(\widehat{m}^{\text{LS}} - m) + (c_\alpha \widehat{a}_i^{\text{LS}} - a_i) + (c_\beta \widehat{b}_j^{\text{LS}} - b_j)]^2 \right\}$$

$$(S.5.4) \quad = \frac{1}{rc} \mathbb{E} \left\{ rc(\widehat{m}^{\text{LS}} - m)^2 + c \sum_{i=1}^r (c_\alpha \widehat{a}_i^{\text{LS}} - a_i)^2 + r \sum_{j=1}^c (c_\beta \widehat{b}_j^{\text{LS}} - b_j)^2 \right\}$$

$$(S.5.5) \quad = \mathbb{E} \left\{ (\widehat{m}^{\text{LS}} - m)^2 \right\} + \frac{1}{r} \mathbb{E} \left\{ \sum_{i=1}^r (c_\alpha \widehat{a}_i^{\text{LS}} - a_i)^2 \right\} + \frac{1}{c} \mathbb{E} \left\{ \sum_{j=1}^c (c_\beta \widehat{b}_j^{\text{LS}} - b_j)^2 \right\}$$

where equality (S.5.4) is due to orthogonality of the vectors corresponding to the three sums-of-squares. Note that that independence of  $\widehat{m}^{\text{LS}}, \widehat{\mathbf{a}}_i^{\text{LS}}, \widehat{\mathbf{b}}_j^{\text{LS}}$ , which holds in the balanced case, is not needed in (S.5.3)-(S.5.5). Specifically, (S.5.4) holds also for unbalanced design because of the side conditions satisfied by  $a, b$  and  $\widehat{a}_i^{\text{LS}}, \widehat{b}_j^{\text{LS}}$ ; and (S.5.6) holds, with some known covariance matrices, in general for the generalized least squares estimators. Hence the calculation goes through for unbalanced data as well.

Consequently, one obtains URE by writing URE for each of summands above. Hence, minimizing URE jointly over  $(c_\alpha, c_\beta)$  therefore consists of minimizing separately the “row” term over  $c_\alpha$  and the “column” term over  $c_\beta$ . Since

$$(S.5.6) \quad \widehat{m}^{\text{LS}} \sim N(m, \sigma^2 \lambda_m^2), \quad \widehat{\mathbf{a}}^{\text{LS}} \sim N_r(a, \sigma^2 \Lambda_a), \quad \widehat{\mathbf{b}}^{\text{LS}} \sim N_c(b, \sigma^2 \Lambda_b),$$

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  $\widehat{\mathbf{a}}^{\text{LS}}$  for the rows and  $\widehat{\mathbf{b}}^{\text{LS}}$  for the columns. These have closed-form solutions and are easily seen to coincide with the URE solutions.

In the unbalanced case, (S.5.2) no longer holds, and so the Bayes estimates for  $\mathbf{a}$  and  $\mathbf{b}$  are each functions of both  $\widehat{\mathbf{a}}^{\text{LS}}$  and  $\widehat{\mathbf{b}}^{\text{LS}}$ . We can nevertheless use shrinkage estimators of the form (S.5.2) 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

$$(S.5.7) \quad \widehat{\eta}_{ij}^{\text{XKB}} = \widehat{m}^{\text{LS}} + \widehat{c}_\alpha^{\text{XKB}} \widehat{a}_i^{\text{LS}} + \widehat{c}_\beta^{\text{XKB}} \widehat{b}_j^{\text{LS}}, \quad 1 \leq i \leq r, \quad 1 \leq j \leq c,$$

$$(S.5.8) \quad \text{where, } \widehat{c}_\alpha^{\text{XKB}} = \arg \min_{c_\alpha \in [0,1]} \text{URE} \left\{ \sum_{i=1}^r (c_\alpha \widehat{a}_i^{\text{LS}} - a_i)^2 \right\},$$

$$(S.5.9) \quad \widehat{c}_\beta^{\text{XKB}} = \arg \min_{c_\beta \in [0,1]} \text{URE} \left\{ \sum_{j=1}^c (c_\beta \widehat{b}_j^{\text{LS}} - b_j)^2 \right\}.$$

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.

**S.6. A list of some basic results used in our proofs.** The following basic matrix algebra results are used in our proofs:

- R1. For p.s.d. matrices  $A, B$ , if  $0 \prec B \preceq A$ , then  $A^{-1} \preceq B^{-1}$  and  $\lambda_k(B) \leq \lambda_k(A)$  for any  $k$ .
- R2. For p.s.d matrices  $A, B$ ,  $BAB$  is also p.s.d.
- R3. For p.s.d matrices  $A, B$ ,  $\lambda_k(AB) \leq \lambda_k(A) \cdot \lambda_k(B)$  for any  $k$ .
- R4. For any matrices  $C$  and  $D$ ,  $\sigma_1(CD) \leq \sigma_1(C) \cdot \sigma_1(D)$ .
- R5. (Von Neumann Trace inequality) If  $C$  and  $D$  are  $n \times n$  Hermitian matrices then:

$$\sum_{i=1}^n \lambda_i(A) \lambda_{n-i+1}(B) \leq \text{tr}(AB) \leq \sum_{i=1}^n \lambda_i(A) \lambda_i(B).$$

Equality holds on the right when  $B = \sum_{i=1}^n \lambda_i(B) U_i U_i^*$ , and equality holds on the left when  $B = \sum_{i=1}^n \lambda_{n-i+1}(B) U_i U_i^*$  where  $U_i$  is the right eigenvector of  $A$  for the eigenvalue  $\lambda_i(A), i = 1, \dots, n$ .

- R6. For any matrix  $C$ ,  $\sigma_1(C^\top C) = \sigma_1(CC^\top)$

The following facts about derivatives involving matrix expressions are used in our paper.

For matrices  $U, B$  and  $V$  where  $B$  is independent of  $x$  we have:

- R7.  $\frac{\partial}{\partial x} \{x^\top Bx\} = x^\top (B + B^\top)$

$$\text{R8. } \frac{\partial}{\partial x} \log |A| = \text{tr}(A^{-1} \frac{\partial A}{\partial x})$$

$$\text{R9. } \frac{\partial}{\partial x} A^{-1} = -A^{-1} \frac{\partial A}{\partial x} A^{-1}$$

$$\text{R10. } \frac{\partial}{\partial x} \{UBV\} = \frac{\partial U}{\partial x} BV + UB \frac{\partial V}{\partial x}$$

## References.

- Berger, J. O. (1985). *Statistical decision theory and Bayesian analysis*. Springer.
- Sclove, 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. (1966). Estimable functions and testable hypotheses in linear models. Tech. Rep. BU-213-M, Cornell University, Biometrics Unit.
- 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.
- 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.