

# Estimating a Mean Matrix: Boosting Efficiency by Multiple Affine Shrinkage

Rudolf Beran\*

University of California, Davis

beran@wald.ucdavis.edu

Revised December 2006

## Abstract

The unknown matrix  $M$  is the mean of the observed response matrix in a multivariate linear model with independent random errors. This paper constructs regularized estimators of  $M$  that dominate, in asymptotic risk, least squares fits to the model and to specified nested submodels. In the first construction, the response matrix is expressed as the sum of orthogonal components determined by the submodels; each component is replaced by an adaptive total least squares fit of possibly lower rank; and these fits are then summed. The second, lower risk, construction differs only in the second step: each orthogonal component is replaced by a modified Efron-Morris fit before summation. Singular value decompositions yield computable formulae for the estimators and their asymptotic and estimated risks. In the asymptotics, the row dimension of  $M$  tends to infinity while the column dimension remains fixed. Convergences are uniform when signal-to-noise ratio is bounded.

*Keywords:* total least squares, Efron-Morris fit, penalized least squares, regularization, rank reduction.

## 1 Introduction

Consider the multivariate linear model

$$Y = M + E, \quad M = XB. \tag{1.1}$$

---

\*This research was supported in part by National Science Foundation Grant DMS 0404547.

Here the observation matrix  $Y$  and the constant matrix  $M$  are both  $n \times q$  with  $n \geq q$ . The  $n \times d$  design matrix  $X$  is known and has rank  $p \leq \max\{d, n\}$  while the  $d \times q$  matrix  $B$  is unknown. The elements of the  $n \times q$  error matrix  $E$  are independent, identically distributed random variables, each having a  $N(0, \sigma^2)$  distribution with unknown positive variance  $\sigma^2$ . A basic problem is to estimate  $M$  and  $\sigma^2$  under this model. Section 5 will extend the results to a multivariate linear model with correlated errors.

The Frobenius norm of a matrix  $D$  is defined by  $|D|^2 = \text{tr}(DD') = \text{tr}(D'D)$ . Let  $\hat{M}$  denote any estimator of  $M$ . The quality of  $\hat{M}$  is assessed through the normalized quadratic loss

$$L(\hat{M}, M) = p^{-1}|\hat{M} - M|^2. \quad (1.2)$$

The risk of  $\hat{M}$  is then

$$R(\hat{M}, M, \sigma^2) = \text{EL}(\hat{M}, M), \quad (1.3)$$

the expectation being evaluated under model (1.1). Let  $X^+$  denote the pseudoinverse of  $X$ . The *least squares estimator*  $\hat{M}_{ls} = XX^+Y$  of  $M$  has risk  $q\sigma^2$ . It is known from Stein (1956) that  $\hat{M}_{ls}$  is inadmissible under quadratic loss when  $pq$  exceeds 2. This paper develops regularized estimators of  $M$  that trade off bias against variance so as to achieve lower risk, asymptotically in  $p$ , than  $\hat{M}_{ls}$  under model (1.1). The proposed estimators also dominate specified least squares submodel fits used in their construction.

We begin by considering a set of nested least squares submodel fits to the linear model. It is *not* assumed that any of the proper submodels is true. All risk calculations in this paper are done under model (1.1), except those for those done in Section 5 under a more general error model. Suppose that  $X_1, X_2, \dots, X_{s-1}$  are submodel design matrices whose range spaces are nested as follows:

$$\mathcal{R}(X_1) \subset \mathcal{R}(X_2) \subset \dots \subset \mathcal{R}(X_{s-1}) \subset \mathcal{R}(X). \quad (1.4)$$

Define

$$P_1 = X_1X_1^+, P_2 = X_2X_2^+ - X_1X_1^+, \dots, P_s = XX^+ - X_{s-1}X_{s-1}^+. \quad (1.5)$$

The  $\{P_k\}$  are orthogonal projections that are mutually orthogonal: each  $P_k$  is symmetric, idempotent and  $P_iP_j = 0$  whenever  $i \neq j$ . Because  $\sum_{k=1}^s P_k = XX^+$  and  $M = XB$  under model (1.1), it follows that

$$M = \sum_{k=1}^s P_k M, \quad \hat{M}_{ls} = \sum_{k=1}^s P_k Y = \sum_{k=1}^s P_k \hat{M}_{ls}. \quad (1.6)$$

*Rank reduction.* The first fundamental idea in this paper is to estimate each  $P_k M$  more efficiently by reducing, in a data-based manner, the rank of  $P_k Y = P_k \hat{M}_{ls}$ . The rank-reduced

estimators of the  $\{P_k M\}$  are then summed to obtain the *adaptive projection estimator* of  $M$ . More fully, suppose that the singular value decomposition of  $P_k Y$  is

$$P_k Y = \hat{U}_k \hat{L}_k \hat{V}_k' = \sum_{j=1}^q \hat{l}_{kj} \hat{u}_{kj} \hat{v}_{kj}', \quad (1.7)$$

where  $\hat{U}_k = [\hat{u}_{k1}, \dots, \hat{u}_{kq}]$  is  $n \times q$ ,  $\hat{V}_k = [\hat{v}_{k1}, \dots, \hat{v}_{kq}]$  is  $q \times q$ ,  $\hat{U}_k' \hat{U}_k = \hat{V}_k' \hat{V}_k = \hat{V}_k \hat{V}_k' = I_q$  and  $\hat{L}_k = \text{diag}(\hat{l}_{k1}, \dots, \hat{l}_{kq})$  with  $\hat{l}_{k1} \geq \dots \geq \hat{l}_{kq} \geq 0$ . Let  $\hat{\sigma}^2$  be a consistent estimator of  $\sigma^2$ . For each  $k$ , define

$$\hat{\tau}_k = p^{-1} \hat{\sigma}^2 \text{tr}(P_k), \quad \check{W}_k = p^{-1} (P_k Y)' P_k Y - \hat{\tau}_k I_q. \quad (1.8)$$

Note for later use that  $\sum_{k=1}^s \hat{\tau}_k = \hat{\sigma}^2$ . Equation (1.7) implies the spectral representation

$$\check{W}_k = \sum_{j=1}^q \check{w}_{kj} \hat{v}_{kj} \hat{v}_{kj}', \quad \check{w}_{kj} = p^{-1} \hat{l}_{kj}^2 - \hat{\tau}_k. \quad (1.9)$$

Let  $\check{w}_{kj+} = \max\{\check{w}_{kj}, 0\}$  and  $\check{w}_{kj-} = \min\{\check{w}_{kj}, 0\}$ . Define  $\hat{w}_{kj} = \check{w}_{kj+}$  and

$$\hat{W}_k = \check{W}_{k+} = \sum_{j=1}^q \hat{w}_{kj} \hat{v}_{kj} \hat{v}_{kj}', \quad \check{W}_{k-} = \sum_{j=1}^q \check{w}_{kj-} \hat{v}_{kj} \hat{v}_{kj}'. \quad (1.10)$$

Note that  $\hat{W}_k$  is positive semidefinite. The *adaptive projection estimator* of  $M$  is

$$\hat{M}_{pro} = \sum_{k=1}^s P_k Y \sum_{j: \hat{w}_{kj} > \hat{\tau}_k} \hat{v}_{kj} \hat{v}_{kj}' = \sum_{k=1}^s \sum_{j: \hat{w}_{kj} > \hat{\tau}_k} \hat{l}_{kj} \hat{u}_{kj} \hat{v}_{kj}', \quad (1.11)$$

the last expression using (1.7). Because  $\hat{M}_{pro} = \sum_{k=1}^s P_k \hat{M}_{pro} = X X^+ \hat{M}_{pro}$ , it satisfies the linear model constraint on  $M$  in model (1.1). The derivation of estimator (1.11) from the risk considerations in Sections 2 and 3 justifies its name.

Let  $\hat{h}_k = \#\{j: \hat{w}_{kj} > \hat{\tau}_k\}$ . By the Eckart-Young matrix approximation theorem, the  $k$ -th summand on the right side of (1.11),  $\sum_{j: \hat{w}_{kj} > \hat{\tau}_k} \hat{l}_{kj} \hat{u}_{kj} \hat{v}_{kj}'$ , gives, among matrices of rank not exceeding  $\hat{h}_k$ , the best approximation in Frobenius norm to  $P_k Y$ . In other words, the  $k$ -th summand is the total least squares approximation to  $P_k Y$  of rank not exceeding  $\hat{h}_k$ . For accounts of the total least squares problem and of its solution through the singular value decomposition, see Golub and Van Loan (1980, 1996), Van Huffel and Vandewalle (1991), and Van Huffel (2004). The latter references relate total least squares to the extensive statistical literature on errors-in-variables regression, including Gleser (1981) and Fuller (1987). Section 4 shows that  $\hat{M}_{pro}$  implements the asymptotically minimum risk strategy for estimating  $M$  through rank reduction of each  $P_k Y$  as  $p$  tends to infinity.

Suppose, for  $1 \leq k \leq s$ , that the singular value decomposition of  $P_k M$  is

$$P_k M = U_k L_k V_k' = \sum_{j=1}^q l_{kj} u_{kj} v_{kj}', \quad (1.12)$$

where  $U_k = [u_{k1}, \dots, u_{kq}]$  is  $n \times q$ ,  $V_k = [v_{k1}, \dots, v_{kq}]$  is  $q \times q$ ,  $U_k'U_k = V_k'V_k = V_kV_k' = I_q$  and  $L_k = \text{diag}(l_{k1}, \dots, l_{kq})$  with  $l_{k1} \geq \dots \geq l_{kq} \geq 0$ . For each  $k$ , define

$$\tau_k = p^{-1}\sigma^2 \text{tr}(P_k), \quad W_k = p^{-1}(P_kM)'P_kM. \quad (1.13)$$

Note for later use that  $\sum_{k=1}^s \tau_k = \sigma^2$ . Equation (1.12) implies the spectral representation

$$W_k = \sum_{j=1}^q w_{kj}v_{kj}v_{kj}', \quad w_{kj} = p^{-1}l_{kj}^2. \quad (1.14)$$

It is shown in Section 4 that the risk  $R(\hat{M}_{pro}, M, \sigma^2)$  converges asymptotically, as  $p$  tends to infinity, to

$$\sum_{k=1}^s \sum_{j=1}^q \min\{\tau_k, w_{kj}\} \leq q\sigma^2 = R(\hat{M}_{ls}, M, \sigma^2). \quad (1.15)$$

The convergence is uniform when  $p^{-1}|M|^2$  is bounded. It will be seen that  $\hat{M}_{pro}$  asymptotically dominates each of the submodel least squares estimators  $\{X_tX_t^+Y : 1 \leq t \leq s\}$  for  $M$ .

*Symmetric affine shrinkage.* It is evident from the right side of (1.11) that  $\hat{M}_{pro}$  applies a shrinkage factor that is either 0 or 1 to each summand in the singular value decomposition (1.7) of  $P_kY$ . Can a more sophisticated shrinkage strategy reduce asymptotic risk below that of  $\hat{M}_{pro}$ ? The answer is yes. The second fundamental idea in this paper is to estimate each  $P_kM$  more efficiently by applying data-based, symmetric affine shrinkage to the right side of each  $P_kY = P_k\hat{M}_{ls}$ . These affinely shrunk estimators of the  $\{P_kM\}$  are then summed to estimate  $M$ .

The *adaptive symmetric affine estimator* of  $M$  is

$$\hat{M}_{sym} = \sum_{k=1}^s P_kY \sum_{j=1}^q \hat{w}_{kj}(\hat{\tau}_k + \hat{w}_{kj})^{-1} \hat{v}_{kj} \hat{v}_{kj}' = \sum_{k=1}^s \sum_{j=1}^q \hat{w}_{kj}(\hat{\tau}_k + \hat{w}_{kj})^{-1} \hat{l}_{kj} \hat{u}_{kj} \hat{u}_{kj}', \quad (1.16)$$

the last expression using (1.7). Because  $\hat{M}_{sym} = \sum_{k=1}^s P_k\hat{M}_{sym} = XX^+\hat{M}_{sym}$ , it satisfies the linear model constraint on  $M$  in model (1.1). In view of (1.16),  $\hat{M}_{sym}$  applies a shrinkage factor that lies between 0 and 1 to each summand in the singular value decomposition (1.7) of  $P_kY$ . The derivation of estimator  $\hat{M}_{sym}$  from the risk considerations in Sections 2 and 3 justifies its name. Section 3 further shows that (1.16) is equivalent to estimating each  $P_kM$  by a modified Efron-Morris (1972) estimator based on  $P_kY$  and then summing over all  $k$ .

Section 4 shows that the risk  $R(\hat{M}_{sym}, M, \sigma^2)$  converges asymptotically, as  $p$  tends to infinity, to

$$\sum_{k=1}^s \sum_{j=1}^q \tau_k w_{kj} (\tau_k + w_{kj})^{-1}. \quad (1.17)$$

The convergence is uniform when  $p^{-1}|M|^2$  is bounded. It follows algebraically from (1.17) and (1.15) that the asymptotic risk of  $\hat{M}_{sym}$  lies between the asymptotic risk of  $\hat{M}_{pro}$  and one half that asymptotic risk.

## 2 Oracle Estimators

This section studies classes of candidate linear estimators for  $M$ , constructing within each class an estimator that minimizes the quadratic risk (1.3). These best candidate estimators are *oracle* estimators in that they depend on functions of the unknown parameters  $M$  and  $\sigma^2$ . The labeling of the oracle estimators foreshadows their linkage, in Section 3, with the adaptive projection and symmetric affine shrinkage estimators described in the Introduction.

### 2.1 Oracle linear estimators

For  $1 \leq k \leq s$ , let  $A_k$  be an arbitrary  $q \times q$  matrix and let  $A$  denote the concatenated matrix  $A = [A_1, A_2, \dots, A_s]$ . As candidate estimators for  $M$ , we first consider the linear estimators

$$\hat{M}(A) = \sum_{k=1}^s P_k Y A_k = \sum_{k=1}^s P_k \hat{M}_{ls} A_k. \quad (2.1)$$

Because  $\hat{M}(A) = \sum_{k=1}^s P_k \hat{M}(A) = X X^+ \hat{M}(A)$ , it satisfies the linear model constraint on  $M$  in (1.1). The loss (1.2) of candidate estimator  $\hat{M}(A)$  is

$$L(\hat{M}(A), M) = p^{-1} |\hat{M}(A) - M|^2 = p^{-1} \sum_{k=1}^s |P_k Y A_k - P_k M|^2. \quad (2.2)$$

The corresponding risk is

$$\begin{aligned} R(\hat{M}(A), M, \sigma^2) &= \text{EL}(\hat{M}(A), M) \\ &= \sum_{k=1}^s \text{tr}[\tau_k A_k A_k' + W_k (A_k - I_q)(A_k - I_q)']. \end{aligned} \quad (2.3)$$

Using derivative formulae compiled by Lütkepohl (1996),

$$\partial R(\hat{M}(A), M, \sigma^2) / \partial A_k = 2(\tau_k A_k + W_k A_k - W_k). \quad (2.4)$$

The derivative in (2.4) vanishes if and only if  $A_k$  equals

$$\tilde{A}_k = (\tau_k I_q + W_k)^{-1} W_k = I_q - \tau_k (\tau_k I_q + W_k)^{-1} = W_k (\tau_k I_q + W_k)^{-1}, \quad (2.5)$$

a symmetric matrix whose eigenvalues all lie in  $[0, 1]$ . Let  $\tilde{A} = [\tilde{A}_1, \dots, \tilde{A}_s]$ . Because the risk  $R(\hat{M}(A), M, \sigma^2)$  is convex in  $A$ , it is minimized by setting  $A = \tilde{A}$ .

## 2.2 Oracle affine shrinkage and projection estimators

Let  $\mathcal{A}_{sym}$  denote the set of all symmetric  $q \times q$  matrices whose eigenvalues lie in  $[0, 1]$ . The concatenated matrix  $A = [A_1, \dots, A_s]$  then lies in  $\mathcal{A}_{sym}^s$ . By the preceding paragraphs, it is reasonable to limit the search for low risk linear estimators to the *symmetric affine shrinkage* candidate estimators  $\{\hat{M}(A): A \in \mathcal{A}_{sym}^s\}$ . Let  $\mathcal{A}_{pro}$  be the subset of matrices in  $\mathcal{A}_{sym}$  that are orthogonal projections. These projections are the elements of  $\mathcal{A}_{sym}$  whose eigenvalues are either 0 or 1.

Let  $W = [W_1, W_2, \dots, W_s]$  and let  $\tau = (\tau_1, \dots, \tau_s)$ . For every  $A \in \mathcal{A}_{sym}^s$ , the risk (2.3) of  $\hat{M}(A)$  simplifies to

$$R(M(A), M, \sigma^2) = \sum_{k=1}^s \rho(A_k, \tau_k, W_k) = r(A, \tau, W) \quad (\text{say}), \quad (2.6)$$

where

$$\begin{aligned} \rho(A_k, \tau_k, W_k) &= \text{tr}[\tau_k A_k^2 + (I_q - A_k)^2 W_k] \\ &= \text{tr}[(A_k - \tilde{A}_k)^2 (\tau_k I_q + W_k)] + \tau_k \text{tr}(\tilde{A}_k). \end{aligned} \quad (2.7)$$

Of interest for subsequent developments are the following oracle estimators, obtained by minimizing risk over  $\mathcal{A}_{sym}$  and over  $\mathcal{A}_{pro}$ :

- The *oracle symmetric affine shrinkage estimator* of  $M$  is  $\tilde{M}_{sym} = \hat{M}(\tilde{A}_{sym})$ , where

$$\tilde{A}_{sym} = \underset{A \in \mathcal{A}_{sym}^s}{\text{argmin}} r(A, \tau, W) = [\tilde{A}_{sym,1}, \dots, \tilde{A}_{sym,s}]. \quad (2.8)$$

- The *oracle projection estimator* of  $M$  is  $\tilde{M}_{pro} = \hat{M}(\tilde{A}_{pro})$ , where

$$\tilde{A}_{pro} = \underset{A \in \mathcal{A}_{pro}^s}{\text{argmin}} r(A, \tau, W) = [\tilde{A}_{pro,1}, \dots, \tilde{A}_{pro,s}]. \quad (2.9)$$

The next theorem provides explicit formulae for the oracle estimators and their risks.

**Theorem 2.1.** *The following expressions hold:*

$$\begin{aligned} \tilde{A}_{sym,k} &= \tilde{A}_k = \sum_{j=1}^q w_{kj} (\tau_k + w_{kj})^{-1} v_{kj} v'_{kj} \\ \tilde{M}_{sym} &= \sum_{k=1}^s P_k Y \tilde{A}_{sym,k} \\ R(\tilde{M}_{sym}, M, \sigma^2) &= \sum_{k=1}^s \tau_k \text{tr}(\tilde{A}_{sym,k}) = \sum_{k=1}^s \sum_{j=1}^q \tau_k w_{kj} (\tau_k + w_{kj})^{-1}. \end{aligned} \quad (2.10)$$

Moreover,

$$\begin{aligned}
\tilde{A}_{pro,k} &= \sum_{j:w_{kj}>\tau_k} v_{kj}v'_{kj} \\
\tilde{M}_{pro} &= \sum_{k=1}^s P_k Y \tilde{A}_{pro,k} \\
R(\tilde{M}_{pro}, M, \sigma^2) &= \sum_{k=1}^s \sum_{j=1}^q \min\{\tau_k, w_{kj}\}.
\end{aligned} \tag{2.11}$$

*Proof.* The three equations in (2.10) follow easily from (2.5), (2.6), (2.7), and the spectral representation (1.14).

Let  $\mathcal{A}_{pro}(t)$  denote the set of matrices in  $\mathcal{A}_{pro}$  whose rank is  $t$ . To find  $\tilde{A}_{pro}$ , we first minimize each risk component  $\rho(A_k, \tau_k, W_k)$  over all  $A_k \in \mathcal{A}_{pro}(t)$ , then minimize further over  $t$ . For each  $k$ , using (2.7), let

$$\tilde{A}_{pro,k}(t) = \operatorname{argmin}_{B \in \mathcal{A}_{pro}(t)} \rho(B, \tau_k, W_k) = \operatorname{argmin}_{B \in \mathcal{A}_{pro}(t)} |W_k^{1/2} - BW_k^{1/2}|^2. \tag{2.12}$$

Evidently,  $\operatorname{rank}(BW_k^{1/2}) \leq t$  under the constraint  $\operatorname{rank}(B) = t$ . By the Eckart-Young matrix approximation theorem and the spectral representation (1.14), the minimum norm approximation to  $W_k^{1/2}$  of rank not exceeding  $t$  is  $\sum_{j=1}^t w_{kj}^{1/2} v_{kj} v'_{kj}$ . It follows that  $\tilde{A}_{pro,k}(t) = \sum_{j=1}^t v_{kj} v'_{kj}$  achieves the minimum in (2.12). From this, the second line of (2.7), and (1.14),

$$\rho(\tilde{A}_{pro,k}(t), \tau_k, W_k) = \sum_{j=1}^q \{[I(j \leq t) - w_{kj}(\tau_k + w_{kj})^{-1}]^2 (\tau_k + w_{kj})\} + \tau_k \operatorname{tr}(\tilde{A}_k). \tag{2.13}$$

This, in turn, is minimized by selecting  $t$  so that the indicator  $I(j \leq t)$  equals 1 when  $w_{kj}(\tau_k + w_{kj})^{-1}$  exceeds  $1/2$  and equals 0 otherwise. The expression for  $\tilde{A}_{pro,k}$  in (2.11) follows.

Using that expression, (2.1), (2.6), and the spectral representation (1.14) yields the other two equations in (2.11).  $\square$

## 2.3 Remarks on candidate and oracle estimators

*Remark 1.* It follows from the risk formulae in Theorem 2.1 that the risks of the oracle affine projection and symmetric affine shrinkage estimators satisfy the inequalities

$$2^{-1}R(\tilde{M}_{pro}, M, \sigma^2) \leq R(\tilde{M}_{sym}, M, \sigma^2) \leq R(\tilde{M}_{pro}, M, \sigma^2) \leq R(\hat{M}_{ls}, M, \sigma^2) = q\sigma^2. \tag{2.14}$$

The oracle estimators  $\tilde{M}_{pro}$  and  $\tilde{M}_{sym}$  are not realizable because they depend on functions of the unknown parameters  $M$  and  $\sigma^2$ . However, Section 3 will construct data-based approximations to both oracle estimators that achieve the oracle risks asymptotically as  $p$  tends to infinity.

*Remark 2.* The class of candidate estimators  $\{\hat{M}(A): A \in \mathcal{A}_{sym}^s\}$  considered in Section 2.2 substantially enlarges the underlying set of least squares submodel fits specified in the Introduction. Indeed, if  $A_k = I_q$  for  $1 \leq k \leq t < s$  and  $A_k$  vanishes otherwise, then  $\hat{M}(A) = \sum_{k=1}^t P_k Y = X_t X_t^+ Y$ , the least squares estimator for  $M$  under the constraint  $M \in \mathcal{R}(X_t)$ . More generally, when some eigenvalues of  $A_k$  are strictly less than 1, then  $P_k Y A_k = P_k \hat{M}_{l_s} A_k$  shrinks toward 0 the corresponding components of the projection  $P_k \hat{M}_{l_s}$ , after the latter quantity is expanded in terms of the eigenvectors of  $A_k$ . Consequently, if  $A_k = I_q$  for  $1 \leq k \leq t < s$  and  $A_k \neq I_q$  otherwise but still lies in  $\mathcal{A}_{sym}$ , then  $\hat{M}(A)$  shrinks  $\hat{M}_{l_s}$  toward  $\sum_{k=1}^t P_k Y = X_t X_t^+ Y$ , the least squares estimator of  $M$  under the constraint  $M \in \mathcal{R}(X_t)$ .

*Remark 3.* The candidate estimators  $\{\hat{M}(A): A \in \mathcal{A}_s\}$  are regularized estimators of  $M$  because they can be derived as generalized penalized least squares estimators or limits thereof. Indeed, let the  $\{C_k: 1 \leq k \leq s\}$  be specified penalty matrices of row dimension  $q$  and let  $C = [C_1, C_2, \dots, C_s]$ . The *affinely penalized least squares candidate estimator*  $\hat{M}_{pls}$  minimizes

$$\begin{aligned} T(M) &= |Y - M|^2 + \sum_{k=1}^s |P_k M C_k|^2 \\ &= \text{tr}(Y'Y) - 2 \text{tr}(M'Y) + \text{tr}(M'M) + \sum_{k=1}^s \text{tr}(M C_k C_k' M' P_k) \end{aligned} \quad (2.15)$$

over all  $M$  that satisfy the linear model constraint  $M = XB$  in (1.1). This generalization of ordinary penalized least squares has  $s$  penalty terms in which the matrices  $\{C_k: 1 \leq k \leq s\}$  replace scalar penalty weights.

We will show that

$$\hat{M}_{pls} = \sum_{k=1}^s P_k Y (I_q + C_k C_k')^{-1}. \quad (2.16)$$

To verify this, observe that  $M = XB = XX^+M$  for some  $d \times q$  matrix  $B$  if and only if  $M = XX^+G = \sum_{k=1}^s P_k G$  for some  $n \times q$  matrix  $G$ . By substituting the last expression for  $M$  into the right side of (2.15) and then using derivative formulae in Lütkepohl (1996), we obtain

$$\partial T(M)/\partial G = 2\left(-\sum_{k=1}^s P_k Y + \sum_{k=1}^s P_k G + \sum_{k=1}^s P_k G C_k C_k'\right). \quad (2.17)$$

The derivative in (2.17) vanishes if and only if  $\sum_{k=1}^s P_k [G(I_q + C_k C_k') - Y] = 0$  or, equivalently, if and only if  $P_k G(I_q + C_k C_k') = P_k Y$  for every  $k$ . Hence, the minimizing value  $\hat{G}$  of  $G$  satisfies  $P_k \hat{G} = P_k Y (I_q + C_k C_k')^{-1}$  for every  $k$ . It follows that  $\hat{M}_{pls} = \sum_{k=1}^s P_k \hat{G}$  is given by (2.16).

As  $C_k$  varies,  $(I_q + C_k C_k')^{-1}$  generates all matrices in  $\mathcal{A}_{sym}^s$  whose eigenvalues lie in  $(0, 1]$ . Equations (2.16) and (2.1) indicate that each of the candidate estimators  $\{\hat{M}(A): A \in \mathcal{A}_{sym}^s\}$  can be expressed either as a penalized least squares estimator or as a limit of such.

### 3 Adaptive Estimators

This section devises adaptive estimators that are realizable data-based approximations to the oracle estimators derived in Section 2. The oracle construction is modified by replacing the unknown parameters  $\tau$  and  $W$  in the risk function  $r(A, \tau, W)$  with estimators. The resulting adaptive estimators coincide with the estimators  $\hat{M}_{pro}$  and  $\hat{M}_{sym}$  discussed in the Introduction. It will be seen in Section 4 that the risk of each adaptive estimator converges to that of its oracle counterpart as  $p$  tends to infinity.

#### 3.1 Estimated risk and adaptation

Let  $\hat{\sigma}^2$  be an  $L_1$ -consistent estimator of  $\sigma^2$  in a sense to be made precise in Theorem 4.1. The strategy is to estimate the unknown risk function  $r(A, \tau, W)$  by plugging in the estimators  $\hat{\tau}$  and  $\hat{W}$  defined by (1.8) and (1.10).

Let

$$\hat{A}_k = (\hat{\tau}_k I_q + \hat{W}_k)^{-1} \hat{W}_k = I_q - \hat{\tau}_k (\hat{\tau}_k I_q + \hat{W}_k)^{-1} = \hat{W}_k (\hat{\tau}_k I_q + \hat{W}_k)^{-1}, \quad (3.1)$$

a matrix that lies in  $\mathcal{A}_{sym}$ . Let  $\hat{A} = [\hat{A}_1, \dots, \hat{A}_s]$ . For every  $A \in \mathcal{A}_{sym}^s$ , define the *estimated risk* of candidate estimator  $\hat{M}(A)$  by analogy with (2.6) and (2.7):

$$r(A, \hat{\tau}, \hat{W}) = \sum_{k=1}^s \rho(A_k, \hat{\tau}_k, \hat{W}_k), \quad (3.2)$$

where

$$\begin{aligned} \rho(A_k, \hat{\tau}_k, \hat{W}_k) &= \text{tr}[\hat{\tau}_k A_k^2 + (I_q - A_k)^2 \hat{W}_k] \\ &= \text{tr}[(A_k - \hat{A}_k)^2 (\hat{\tau}_k I_q + \hat{W}_k)] + \hat{\tau}_k \text{tr}(\hat{A}_k). \end{aligned} \quad (3.3)$$

Corresponding to the oracle estimators discussed in Section 2.2 are the following adaptive estimators, obtained by minimizing estimated risk over  $\mathcal{A}_{sym}$  and over  $\mathcal{A}_{pro}$ :

- The *adaptive symmetric affine estimator* of  $M$  is  $\hat{M}_{sym} = \hat{M}(\hat{A}_{sym})$ , where

$$\hat{A}_{sym} = \underset{A \in \mathcal{A}_{sym}^s}{\text{argmin}} r(A, \hat{\tau}, \hat{W}) = [\hat{A}_{sym,1}, \dots, \hat{A}_{sym,s}]. \quad (3.4)$$

- The *adaptive projection estimator* of  $M$  is  $\hat{M}_{pro} = \hat{M}(\hat{A}_{pro})$ , where

$$\hat{A}_{pro} = \underset{A \in \mathcal{A}_{pro}^s}{\text{argmin}} r(A, \hat{\tau}, \hat{W}) = [\hat{A}_{pro,1}, \dots, \hat{A}_{pro,s}]. \quad (3.5)$$

The estimated risks of these two estimators are, respectively,

$$\hat{R}(\hat{M}_{sym}) = r(\hat{A}_{sym}, \hat{\tau}, \hat{W}), \quad \hat{R}(\hat{M}_{pro}) = r(\hat{A}_{pro}, \hat{\tau}, \hat{W}). \quad (3.6)$$

The next theorem shows that these adaptive estimators coincide with the adaptive projection and adaptive symmetric affine estimators discussed in the Introduction and provides formulae for their estimated risks.

**Theorem 3.1.** *The following expressions hold:*

$$\begin{aligned}\hat{A}_{sym,k} &= \hat{A}_k = \sum_{j=1}^q \hat{w}_{kj} (\hat{\tau}_k + \hat{w}_{kj})^{-1} \hat{v}_{kj} \hat{v}'_{kj} \\ \hat{M}_{sym} &= \sum_{k=1}^s P_k Y \sum_{j=1}^q \hat{A}_{sym,k} = \sum_{k=1}^s \sum_{j=1}^q \hat{w}_{kj} (\hat{\tau}_k + \hat{w}_{kj})^{-1} \hat{l}_{kj} \hat{u}_{kj} \hat{v}'_{kj} \\ \hat{R}(\hat{M}_{sym}) &= \sum_{k=1}^s \hat{\tau}_k \text{tr}(\hat{A}_{sym,k}) = \sum_{k=1}^s \sum_{j=1}^q \hat{\tau}_k \hat{w}_{kj} (\hat{\tau}_k + \hat{w}_{kj})^{-1},\end{aligned}\tag{3.7}$$

with  $\hat{A}_k$  defined in (3.1). Moreover,

$$\begin{aligned}\hat{A}_{pro,k} &= \sum_{j: \hat{w}_{kj} > \hat{\tau}_k} \hat{v}_{kj} \hat{v}'_{kj} \\ \hat{M}_{pro} &= \sum_{k=1}^s P_k Y \hat{A}_{pro,k} = \sum_{k=1}^s \sum_{j: \hat{w}_{kj} > \hat{\tau}_k} \hat{l}_j \hat{u}_{kj} \hat{v}'_{kj} \\ \hat{R}(\hat{M}_{pro}) &= \sum_{k=1}^s \sum_{j=1}^q \min\{\hat{\tau}_k, \hat{w}_{kj}\}.\end{aligned}\tag{3.8}$$

*Proof.* Because  $\hat{W}_k$  was constructed to lie in  $\mathcal{A}_{sym}$ , just like  $W_k$ , the argument that proved Theorem 2.1 also works for Theorem 3.1. The rightmost expressions for  $\hat{M}_{sym}$  and  $\hat{M}_{pro}$  use the singular values decomposition of  $P_k Y$ , given in (1.7).  $\square$

It follows from the estimated risk formulae in Theorem 3.1 that

$$2^{-1} \hat{R}(\hat{M}_{pro}) \leq \hat{R}(\hat{M}_{sym}) \leq \hat{R}(\hat{M}_{pro}) \leq \hat{R}(\hat{M}_{ls}) = q\hat{\sigma}^2.\tag{3.9}$$

Section 4 shows that each estimated risk in (3.9) converges asymptotically, as  $p$  tends to infinity, to the corresponding actual risk, which asymptotically equals the oracle risk. In this sense, the estimated risk inequalities (3.9) approximate the estimated risk inequalities (2.14).

## 3.2 Examples

The examples of model (1.1) in this section have two purposes:

- To relate, by specialization, the adaptive estimator  $\hat{M}_{sym}$  and  $\hat{M}_{pro}$  to earlier work on shrinkage estimators;

- To illustrate ways in which  $\hat{M}_{sym}$  and  $\hat{M}_{pro}$  go beyond the earlier procedures. The case study in Section 4.3 develops this point on data.

*Example 1.* (Multiple scalar shrinkage). Consider the scalar response case,  $q = 1$ . We write  $Y = y$ ,  $M = m$ , both column vectors, and  $A_k = a_k$ , a scalar for each  $k$ . Then,  $\hat{W}_k$  is the non-negative scalar  $\hat{w}_k = [p^{-1}(P_k y)'(P_k y) - \tau_k]_+$ . The adaptive affine shrinkage estimator of  $m$  is

$$\hat{m}_{sym} = \sum_{k=1}^s \hat{w}_k (\hat{\tau}_k + \hat{w}_k)^{-1} P_k y \quad (3.10)$$

and it has estimated risk

$$\hat{R}(\hat{m}_{sym}) = \sum_{k=1}^s \hat{\tau}_k \hat{w}_k (\hat{\tau}_k + \hat{w}_k)^{-1}. \quad (3.11)$$

It may be verified that, as  $p$  tends to infinity,  $\hat{m}_{sym}$  converges to a multiple scalar shrinkage estimator constructed by Stein (1966).

In this scalar case, the adaptive projection estimator of  $m$  is the vector

$$\hat{m}_{pro} = \sum_{k: \hat{w}_k > \hat{\tau}_k} P_k y, \quad (3.12)$$

with estimated risk

$$\hat{R}(\hat{m}_{pro}) = \sum_{k=1}^s \min\{\hat{\tau}_k, \hat{w}_k\}. \quad (3.13)$$

The projection estimator  $\hat{m}_{pro}$  is a multiple submodel selection estimator that simplifies  $\hat{m}_{sym}$  as follows: the shrinkage factor  $\hat{w}_k (\hat{\tau}_k + \hat{w}_k)^{-1}$  in (3.10) is replaced by 1 whenever it exceeds 1/2 and is replaced by 0 otherwise. The asymptotics in Section 4 establish that the risk estimators (3.11) and (3.13) both converge to the respective true risks under model (1.1).

*Example 2.* (Efron-Morris estimator and total least squares). Let  $n = p$ ,  $X = I_p$ , and  $s = 1$ . For simplicity, suppose that  $\sigma^2$  is known to be 1, so that  $\hat{\sigma}^2 = 1$ . Then  $P_1 = XX^+ = I_p$ ,  $\tau_1 = 1$ , and  $\hat{W}_1$  adjusts  $\check{W}_1 = p^{-1}Y'Y - I_q$  to be positive semidefinite. In this case,

$$\hat{M}_{sym} = Y \hat{W}_1 (I_q + \hat{W}_1)^{-1}. \quad (3.14)$$

is a modification of the estimator  $Y \check{W}_1 (I_q + \check{W}_1)^{-1} = Y [I_q - p(Y'Y)^{-1}]$ . For  $p$  much larger than  $q$ , this second expression nearly coincides with the Efron-Morris estimator

$$\hat{M}_{EM} = Y [I_q - (p - q - 1)(Y'Y)^{-1}], \quad (3.15)$$

developed from a different perspective by Efron and Morris (1972). Under asymptotics where  $p$  tends to infinity while response dimension  $q$  stays fixed,  $\hat{M}_{sym}$  and  $\hat{M}_{EM}$  have the same asymptotic risk.

Over certain neighborhoods of  $M = 0$  that are defined in terms of the eigenvalues of  $W_1$ , the estimator  $\hat{M}_{sym}$  is asymptotically minimax as  $p$  tends to infinity. The least squares estimator  $\hat{M}_{ls}$  is not. This result in Beran (1999) applies Pinsker's (1980) theorem to a canonical transformation of the present example.

In the present example,  $\hat{M}_{pro}$ , given by (3.8) with  $s = 1$ , is the adaptive total least squares estimator for  $M$  of rank  $\hat{h}_1 = \#\{j: \hat{w}_{1j} > \hat{\tau}_1\}$ . The total least squares problem and its solution through the singular value decomposition were described by Golub and Van Loan (1980, 1996), by Van Huffel and Vandewalle (1991), and by Van Huffel (2004). The asymptotics in Section 4 provide a risk rationale for choosing the rank of the total least squares fit to be  $\hat{h}_1$ . As in Example 1,  $\hat{M}_{pro}$  is a simplification of the modified Efron-Morris estimator. The risk penalty incurred by  $\hat{M}_{pro}$  relative to  $\hat{M}_{sym}$  can be estimated asymptotically as  $p$  increases. For another treatment of this example, see Beran (2007).

Example 2 reveals, more generally, that  $\hat{M}_{sym}$  is the sum of  $s$  modified Efron-Morris estimators, one fitted to each  $P_k Y$ . Similarly,  $\hat{M}_{pro}$  is the sum of  $s$  adaptive total least squares estimators, one fitted to each  $P_k Y$ .

*Example 3.* (Two-way MANOVA model). Consider the two-way layout with  $q$ -variate responses where factors 1 and 2 have, respectively,  $p_1$  and  $p_2$  levels. For each factor level combination  $(i, j)$ , the observed responses are  $1 \times q$  row vectors  $\{y_{ijh}: 1 \leq h \leq n_{ij}\}$ . The model asserts that

$$y_{ijh} = \beta_{ij} + e_{ijh}, \quad 1 \leq h \leq n_{ij}, \quad (3.16)$$

where  $\beta_{ij}$  is a constant  $1 \times q$  vector and  $e_{ijh}$  is a random  $1 \times q$  error vector.

Let  $n = \sum_{i=1}^{p_1} \sum_{j=1}^{p_2} n_{ij}$  and  $p = p_1 p_2$ . Form the  $n \times q$  response matrix  $Y$  by stacking the row vectors  $\{y_{ijh}\}$  in the following order: replication label  $h$  varies most quickly, the level  $i$  of factor 1 varies next most quickly, and the level  $j$  of factor 2 varies most slowly. Form the error matrix  $E$  analogously. Form the  $p \times q$  constant matrix  $B$  by stacking the row vectors  $\{\beta_{ij}\}$  in the following order: the level  $i$  of factor 1 varies most quickly, and the level  $j$  of factor 2 varies most slowly. Let  $X$  be the  $n \times p$  data incidence matrix, a matrix with orthogonal columns whose elements are either 0 or 1, that links each component of  $B$  to the corresponding responses in  $Y$ . Then  $Y$  in this example satisfies model (1.1)

For  $i = 1, 2$ , consider the  $p_i \times 1$  vector  $u_i = p_i^{-1/2}(1, \dots, 1)'$  and the orthogonal projections  $J_i = u_i u_i'$  and  $H_i = I_{p_i} - J_i$ . Let

$$Q_0 = J_2 \otimes J_1, \quad Q_1 = J_2 \otimes H_1, \quad Q_2 = H_2 \otimes J_1, \quad Q_{12} = H_2 \otimes H_1. \quad (3.17)$$

The two-way MANOVA decomposition of the cell means  $B$  into overall means, main effects, and interactions is the identity  $B = Q_0 B + Q_1 B + Q_2 B + Q_{12} B$ . Customary nested submodel design matrices for two-way MANOVA are

$$X_1 = X Q_0, \quad X_2 = X(Q_0 + Q_1) \text{ or } X(Q_0 + Q_2), \quad X_3 = X(Q_0 + Q_1 + Q_2). \quad (3.18)$$

In this example,  $\hat{M}_{sym}$  and  $\hat{M}_{pro}$  provide superior estimation of  $M = XB$ .

*Example 4.* (Multivariate regression). Let  $x = (x_1, x_2, \dots, x_n)'$  be a non-random covariate vector associated with the rows of the response matrix  $Y$ . In model (1.1), let  $X$  be the matrix  $[x^0, x^1, \dots, x^{p-1}]$ , the powers being computed componentwise. For  $1 \leq k \leq s-1$ , let  $X_k$  be the submatrix consisting of the first  $c_k$  columns of  $X$ , where  $1 \leq c_1 \leq \dots \leq c_{s-1} \leq p-1$ . By construction,  $\hat{M}_{sym}$  asymptotically dominates in risk the least squares polynomial fits of degrees  $c_1 - 1, c_2 - 1, \dots, p - 1$ . The polynomial basis in this example can be replaced by other bases, such as discrete spline bases.

## 4 Asymptotic Theory

This section develops conditions under which the adaptive estimators behave asymptotically in loss and risk like their oracle counterparts and under which the estimated risks of the adaptive estimators converge to their actual risks. The usefulness of  $\hat{M}_{sym}$  and  $\hat{M}_{pro}$  for data analysis is explored in the case study of Section 4.3. The assumption that the rows of the error matrix  $E$  are independent  $N(0, \sigma^2 I_q)$  random vectors is important to the proofs below because it ensures the distributional property that follows (4.4). For the special case when  $s = 1$  and  $P_1 = I_n$ , analogs of Theorems 4.1 and 4.2 have been proved by Beran (2007), using an empirical process approach that does not require a multivariate normality assumption. However, the approach of that paper does not suffice to handle the case  $s > 1$ .

### 4.1 Convergence of loss, risk, and estimated risk

The first theorem establishes that the loss and the estimated risk of candidate estimator  $\hat{M}(A)$  both converge, uniformly over  $A \in \mathcal{A}_{sym}^s$ , to the estimator's risk as  $p$  tends to infinity.

**Theorem 4.1.** *Suppose that model (1.1) holds and that  $\hat{\sigma}^2$  satisfies, for every finite  $c > 0$ ,*

$$\lim_{p \rightarrow \infty} \sup_{p^{-1}|M|^2 \leq c} \mathbb{E}|\hat{\sigma}^2 - \sigma^2| = 0. \quad (4.1)$$

*Then, for every finite  $c > 0$  and fixed integers  $q$  and  $s$ ,*

$$\lim_{p \rightarrow \infty} \sup_{p^{-1}|M|^2 \leq c} \mathbb{E} \sup_{A \in \mathcal{A}_{sym}^s} |L(\hat{M}(A), M) - r(A, \tau, W)| = 0 \quad (4.2)$$

*and*

$$\lim_{p \rightarrow \infty} \sup_{p^{-1}|M|^2 \leq c} \mathbb{E} \sup_{A \in \mathcal{A}_{sym}^s} |r(A, \hat{\tau}, \hat{W}) - r(A, \tau, W)| = 0. \quad (4.3)$$

*Proof.* Each  $n \times n$  projection matrix  $P_k$  has spectral decomposition  $P_k = U_k U_k'$ , where  $U_k$  is  $n \times t_k$  with  $t_k = \text{tr}(P_k)$  and  $U_k' U_k = I_{t_k}$ . Let

$$Y_k = U_k' Y, \quad M_k = U_k' M, \quad E_k = U_k' E. \quad (4.4)$$

Under model (1.1), the elements of the  $t_k \times q$  matrix  $E_k$  are independent identically distributed  $N(0, \sigma^2)$  random variables. Moreover,  $W_k = p^{-1} M_k' M_k$  and  $\tau_k = p^{-1} \sigma^2 t_k$ .

Let  $\hat{F}_k = p^{-1} E_k' M_k$  and  $\hat{G}_k = p^{-1} E_k' E_k - \tau_k I_q$ . Using the respective definitions (1.8) and (1.13) of  $\check{W}_k$  and  $W_k$  gives

$$\begin{aligned} \check{W}_k - W_k &= p^{-1} Y_k' Y_k - \hat{\tau}_k I_q - p^{-1} M_k' M_k \\ &= \hat{F}_k + \hat{F}_k' + \hat{G}_k + (\tau_k - \hat{\tau}_k) I_q. \end{aligned} \quad (4.5)$$

Because the elements of  $E_k$  are independent  $N(0, \sigma^2)$  variables and  $|M_k|^2 \leq |M|^2$ ,

$$\sup_{p^{-1}|M|^2 < c} \text{E}|\hat{F}_k|^2 = O(p^{-1}), \quad \sup_{p^{-1}|M|^2 < c} \text{E}|\hat{G}_k|^2 = O(p^{-1}). \quad (4.6)$$

From (4.1) and the definitions (1.8) and (1.13) of  $\hat{\tau}_k$  and  $\tau_k$ ,

$$\lim_{p \rightarrow \infty} \sup_{p^{-1}|M|^2 \leq c} \text{E}|\hat{\tau}_k - \tau_k| = 0. \quad (4.7)$$

Recall the definitions in (1.10) of  $\hat{W}_k = \check{W}_{k+}$  and  $\check{W}_{k-}$ . Note that  $\text{tr}(W_k \check{W}_{k-}) = \sum_{j=1}^q \check{w}_{kj} - \hat{v}'_{kj} W_k \hat{v}_{kj} \leq 0$  because the matrix  $W_k$  is positive semidefinite. Consequently,

$$\begin{aligned} |\check{W}_k - W_k|^2 &= |\check{W}_{k+} - W_k|^2 + |\check{W}_{k-}|^2 - 2 \text{tr}(W_k \check{W}_{k-}) \\ &\geq |\check{W}_{k+} - W_k|^2 + |\check{W}_{k-}|^2 \\ &\geq |\check{W}_{k+} - W_k|^2 = |\hat{W}_k - W_k|^2. \end{aligned} \quad (4.8)$$

From (4.5) to (4.8),

$$\lim_{p \rightarrow \infty} \sup_{p^{-1}|M|^2 \leq c} \text{E}|\hat{W}_k - W_k| \leq \lim_{p \rightarrow \infty} \sup_{p^{-1}|M|^2 \leq c} \text{E}|\check{W}_k - W_k| = 0. \quad (4.9)$$

Let  $\{B_i : 1 \leq i \leq h\}$  be any  $q \times q$  matrices. Then,  $|\text{tr}[\prod_{i=1}^h B_i]| \leq \prod_{i=1}^h |B_i|$  (cf. Lütkepohl (1996), p. 111). If  $A_k \in \mathcal{A}_{sym}$ , its spectral decomposition implies that  $|A_k^2| \leq q^{1/2}$ ,  $|(I_q - A_k)^2| \leq q^{1/2}$ , and  $|A_k(I_q - A_k)| \leq 4^{-1} q^{1/2}$ . It follows from (2.6) that

$$\begin{aligned} |r(A, \tau, \hat{W}) - r(A, \tau, W)| &= \left| \sum_{k=1}^s \text{tr}[(I_q - A_k)^2 (\hat{W}_k - W_k)] \right| \\ &\leq \sum_{k=1}^s |(I_q - A_k)^2| |\hat{W}_k - W_k| \leq q^{1/2} \sum_{k=1}^s |\hat{W}_k - W_k|. \end{aligned} \quad (4.10)$$

On the other hand,

$$|r(A, \hat{\tau}, \hat{W}) - r(A, \tau, \hat{W})| = \left| \sum_{k=1}^s (\hat{\tau}_k - \tau_k) \text{tr}(A_k^2) \right| \leq q|\hat{\sigma}^2 - \sigma^2|. \quad (4.11)$$

Inequalities (4.9), (4.10), (4.11) and theorem condition (4.1) together imply limit (4.3).

Let  $J = |L(\hat{M}(A), M) - r(A, \tau, W)|$ . From (2.2) and (2.6),

$$\begin{aligned} J &= \left| \sum_{k=1}^s \{p^{-1}|Y_k A_k - M_k|^2 - \text{tr}[\tau_k A_k^2 + (I - A_k)^2 p^{-1} M_k' M_k]\} \right| \\ &= \left| \sum_{k=1}^s \text{tr}[A_k^2 \hat{G}_k + 2(A_k^2 - A_k) \hat{F}_k'] \right| \leq q^{1/2} \sum_{k=1}^s [|\hat{G}_k| + 2^{-1} |\hat{F}_k|]. \end{aligned} \quad (4.12)$$

This and (4.6) imply limit (4.2).  $\square$

From Theorem 4.1 follows

**Theorem 4.2.** *Suppose that the assumptions of Theorem 4.1 hold. Let  $T$  denote any one of  $L(\hat{M}_{sym}, M)$  or  $L(\tilde{M}_{sym}, M)$  or  $\hat{R}(\hat{M}_{sym})$ . Then, for every finite  $c > 0$  and for every fixed integers  $q$  and  $s$ ,*

$$\lim_{p \rightarrow \infty} \sup_{p^{-1}|M|^2 \leq c} \mathbb{E}|T - R(\tilde{M}_{sym}, M, \sigma^2)| = 0. \quad (4.13)$$

Consequently,

$$\lim_{p \rightarrow \infty} \sup_{p^{-1}|M|^2 \leq c} \mathbb{E}|L(\hat{M}_{sym}, M) - L(\tilde{M}_{sym}, M)| = 0 \quad (4.14)$$

and

$$\lim_{p \rightarrow \infty} \sup_{p^{-1}|M|^2 \leq c} |R(\hat{M}_{sym}, M, \sigma^2) - R(\tilde{M}_{sym}, M, \sigma^2)| = 0. \quad (4.15)$$

Replacing the subscript “sym” by the subscript “pro” in these assertions is valid.

*Proof.* Recall that the adaptive affine symmetric shrinkage estimator  $\hat{M}_{sym} = \hat{M}(\hat{A})$  with  $\hat{A} = [\hat{A}_1, \hat{A}_2, \dots, \hat{A}_k]$  while the oracle affine shrinkage estimator  $\tilde{M}_{sym} = \tilde{M}(\tilde{A})$  with  $\tilde{A} = [\tilde{A}_1, \tilde{A}_2, \dots, \tilde{A}_k]$ . That (4.13) holds for  $T = L(\tilde{M}_{sym}, M)$  is immediate from (4.2).

Limit (4.3) implies that, for  $A = \hat{A}$  and for  $A = \tilde{A}$ ,

$$\lim_{p \rightarrow \infty} \sup_{p^{-1}|M|^2 \leq c} \mathbb{E}|r(A, \hat{\tau}, \hat{W}) - r(A, \tau, W)| = 0 \quad (4.16)$$

and, through the respective minimizing properties of  $\hat{A}$  and  $\tilde{A}$ ,

$$\lim_{p \rightarrow \infty} \sup_{p^{-1}|M|^2 \leq c} \mathbb{E}|r(\hat{A}, \hat{\tau}, \hat{W}) - r(\tilde{A}, \tau, W)| = 0. \quad (4.17)$$

Because  $r(\tilde{A}, \tau, W) = R(\tilde{M}_{sym}, M, \sigma^2)$  and  $r(\hat{A}, \hat{\tau}, \hat{W}) = \hat{R}(\hat{M}_{sym})$ , (4.17) establishes that limit (4.13) holds for  $T = \hat{R}(\hat{M}_{sym})$ .

Moreover, (4.16) and (4.17) entail

$$\lim_{p \rightarrow \infty} \sup_{p^{-1}|M|^2 \leq c} \mathbb{E}|r(\hat{A}, \tau, W) - r(\tilde{A}, \tau, W)| = 0. \quad (4.18)$$

On the other hand, limit (4.2) gives

$$\lim_{p \rightarrow \infty} \sup_{p^{-1}|M|^2 \leq c} \mathbb{E}|L(\hat{M}(\hat{A}), M) - r(\hat{A}, \tau, W)| = 0. \quad (4.19)$$

Combining (4.19) with (4.18) yields (4.13) for  $T = L(\hat{M}_{sym}, M)$ . Limit (4.14) is immediate from (4.13) and implies (4.15).

Finally, because  $\mathcal{A}_{pro} \subset \mathcal{A}_{sym}$ , the argument just given applies as well to the adaptive and oracle projection estimators.  $\square$

By (4.14) and (4.15), the loss or risk of the adaptive estimator  $\hat{M}_{sym}$  or  $\hat{M}_{pro}$  converges asymptotically to the loss or risk of the corresponding oracle estimator. Moreover, by (4.13), the plug-in risk estimator  $\hat{R}(\hat{M}_{sym})$  or  $\hat{R}(\hat{M}_{pro})$  converges asymptotically to the actual loss or risk of the corresponding adaptive estimator. Thus, it is meaningful to compare the estimated risks of  $\hat{M}_{sym}$  and  $\hat{M}_{pro}$ , given in Theorem 3.1, with one another and with the estimated risk  $q\hat{\sigma}^2$  of the least squares estimator  $\hat{M}_{ls} = XX^+Y$ .

## 4.2 Variance estimation

This section considers three estimators for the variance  $\sigma^2$ , giving conditions for each under which it is  $L_1$ -consistent. The first estimator is useful when  $n$  substantially exceeds  $p$ . The other two estimators express alternative approaches available when  $n = p$  and  $X = I_p$ .

*Least squares variance estimator.* Let  $\hat{E} = Y - \hat{M}_{ls} = (I_n - XX^+)Y$  be the matrix of residuals after the least squares fit. Write  $\hat{E} = [\hat{e}_1, \dots, \hat{e}_n]'$ , where  $\hat{e}_i'$  denotes the  $i$ -th row of  $\hat{E}$ . The classical *least squares estimator* of  $\sigma^2$  is

$$\hat{\sigma}_{ls}^2 = [q(n-p)]^{-1} |\hat{E}|^2 = [q(n-p)]^{-1} \sum_{i=1}^n |\hat{e}_i|^2. \quad (4.20)$$

Because the residual matrix  $\hat{E} = (I_n - XX^+)E$ , it does not depend on  $M$ . The estimator  $\hat{\sigma}_{ls}^2$  is unbiased for  $\sigma^2$  and has the  $L_1$ -consistency property

$$\lim_{n-p \rightarrow \infty} \sup_{p^{-1}|M|^2 \leq c} \mathbb{E}|\hat{\sigma}_{ls}^2 - \sigma^2| = 0. \quad (4.21)$$

*First difference variance estimator.* Express  $Y = [y_1, \dots, y_n]'$  and  $M = [m_1, \dots, m_n]'$  in terms of their rows. When  $n = p$ ,  $X = I_p$ , and row-to-row variation in  $M$  is expected to be slow, it is reasonable to consider the *first difference variance estimator*

$$\hat{\sigma}_{dif}^2 = [2q(p-1)]^{-1} \sum_{i=2}^p |y_i - y_{i-1}|^2. \quad (4.22)$$

The bias of this estimator is  $b(M) = \sum_{i=2}^p |m_i - m_{i-1}|^2$ . Let  $\{\epsilon_p\}$  be any sequence such that  $\lim_{p \rightarrow \infty} \epsilon_p = 0$ . Then  $\hat{\sigma}_{dif}^2$  has the  $L_1$ -consistency property

$$\lim_{p \rightarrow \infty} \sup_{p^{-1}|M|^2 \leq c, b(M) \leq \epsilon_p} \mathbb{E}|\hat{\sigma}_{dif}^2 - \sigma^2| = 0. \quad (4.23)$$

Variance estimation based on higher-order differences of the  $\{y_i\}$  is analogous to that just outlined.

*Smallest singular values variance estimator.* When  $n = p$ ,  $X = I_p$ , and  $M$  is expected to be nearly singular, it is reasonable to consider the *smallest singular values variance estimator*

$$\hat{\sigma}_{ssv}^2 = p^{-1} \sum_{k=1}^s \hat{l}_{kq}^2. \quad (4.24)$$

Recalling (1.7), (1.12), and the definition (1.13) of  $\tau_k$ , observe that

$$\begin{aligned} |p^{-1}\hat{l}_{kq}^2 - p^{-1}l_{kq}^2 - \tau_k| &= \left| \min_{|v|=1} p^{-1}v'(P_k Y)'P_k Y v - \min_{|v|=1} [p^{-1}v'(P_k M)'P_k M v - \tau_k v'v] \right| \\ &\leq \sup_{|v|=1} |v'[p^{-1}(P_k Y)'P_k Y - p^{-1}(P_k M)'P_k M - \tau_k I_q]v| \\ &\leq |p^{-1}(P_k Y)'P_k Y - p^{-1}(P_k M)'P_k M - \tau_k I_q|. \end{aligned} \quad (4.25)$$

By a calculation akin to (4.5) and (4.6), it follows from (4.24) and (4.25) that

$$\lim_{p \rightarrow \infty} \sup_{p^{-1}|M|^2 \leq c} \mathbb{E}|\hat{\sigma}_{ssv}^2 - p^{-1} \sum_{k=1}^s l_{kq}^2 - \sigma^2| = 0. \quad (4.26)$$

Let  $\{\epsilon_p\}$  be any sequence such that  $\lim_{p \rightarrow \infty} \epsilon_p = 0$ . Then  $\hat{\sigma}_{ssv}^2$  is  $L_1$ -consistent in the sense that

$$\lim_{p \rightarrow \infty} \sup_{p^{-1}|M|^2 \leq c, p^{-1} \sum_{k=1}^s l_{kq}^2 \leq \epsilon_p} \mathbb{E}|\hat{\sigma}_{ssv}^2 - \sigma^2| = 0. \quad (4.27)$$

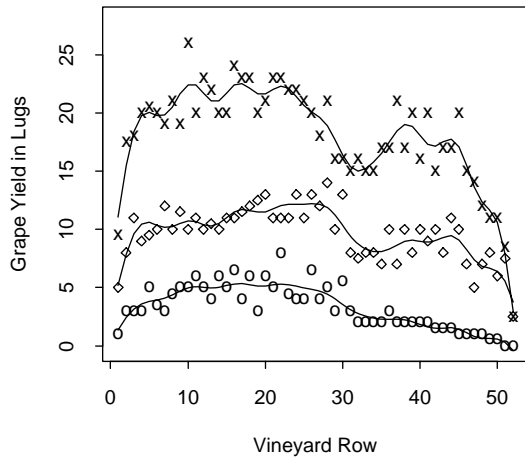
Variance estimators based on the smallest singular values of only some of the  $\{P_k Y\}$ , can be defined and analyzed similarly.

Theorems 4.1 and 4.2 are readily adjusted to use the consistency properties in (4.21) or (4.23) or (4.27) in place of (4.1).

### 4.3 Viticultural case study

The data matrix  $Y$  in this case study is  $52 \times 3$ . Row  $i$  of  $Y$  reports the grape yields harvested in three different years from row  $i$  of a vineyard with 52 rows. The data is taken from Chatterjee, Handcock, and Simonoff (1995). The grape yields, measured in lugs of grapes harvested from each row, are plotted in Figure 1, using a different plotting character for each of the three years. Both year-to-year and row-to-row changes in viticulture affect

Adaptive Symmetric Affine Shrinkage



Adaptive Projection

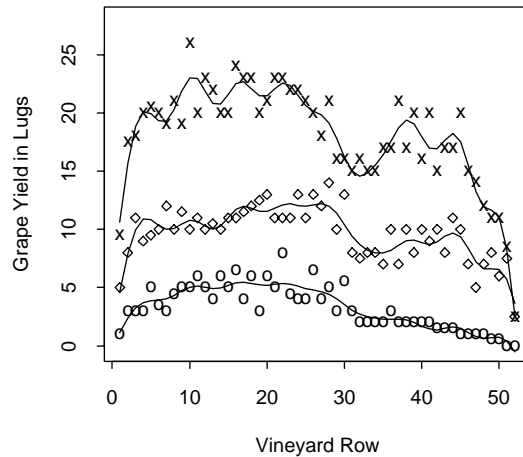


Figure 1: Linearly interpolated adaptive symmetric affine shrinkage and adaptive projection fits to the trivariate grape yields. The  $\{P_k\}$  are defined in terms of a third-difference penalty basis.

the observed yields. The analysis seeks to bring out patterns in the row harvest yields that persist across years.

Let  $D$  denote a matrix with  $n$  columns such that, for any column vector  $x \in R^n$ , the penalty function  $\pi(x) = |Dx|^2 = x'D'Dx$  measures the roughness of  $x$ . For instance,  $D$  may be the  $t$ -th difference operator, defined below, when the components of  $x$  are equally spaced. If adjacent coordinates of  $x$  vary slowly in terms of the difference operator, then the penalty  $\pi(x)$  is small. Construct an orthogonal *penalty basis* for  $R^n$  as follows:

- Find the unit vector  $\gamma \in R^n$  that minimizes the penalty  $\pi(\gamma)$ . This smoothest vector is  $\gamma_n$ , where the  $\{\gamma_j: 1 \leq j \leq n\}$  are the eigenvectors of  $D'D$ , ordered so that the associated eigenvalues go from largest to smallest. In case of tied eigenvalues, the corresponding eigenvectors are selected and ordered by imposing an additional rule.
- Find the unit vector  $\gamma \in R^n$  that, subject to the constraint  $\gamma \perp \gamma_n$ , minimizes the penalty  $\pi(\gamma)$ . This second smoothest vector is  $\gamma_{n-1}$ .
- Continue sequential constrained minimization in order to obtain the orthogonal penalty basis matrix  $U = [u_1, u_2, \dots, u_n]$ , where  $u_i = \gamma_{n-i+1}$ .

The columns of  $U$  are the eigenvectors of  $D'D$ , ordered from the smallest to the largest eigenvalue of  $D'D$ .

Consider the  $(n - 1) \times n$  matrix  $\Delta(n) = \{\delta_{i,j}\}$  in which  $\delta_{i,i} = 1$ ,  $\delta_{i,i+1} = -1$  for every  $i$  and all other entries are zero. Define the  $t$ -th difference operator  $D_t(n)$  through

$$D_1(n) = \Delta(n), \quad D_t(n) = \Delta(n - t + 1)D_{t-1}(n) \quad \text{for } 2 \leq t \leq n - 1. \quad (4.28)$$

In the present case study, where  $n = 52$  and  $q = 3$ , we set  $D$  to be the third difference operator  $D_3(52)$ . The penalty function  $\pi$  then penalizes departures from locally quadratic behavior and the columns of  $U$  form a discrete spline basis.

To specify the multivariate linear model and the submodels, set  $p = 20$  and  $s = 12$ . Define  $X$  to be the first 20 columns of  $U$ . Because these columns are orthonormal,  $X^+ = X'$ . Under this model, the least squares estimate (4.8) of  $\sigma^2$  is  $\hat{\sigma}_{ls}^2 = 1.839$ , with 96 degrees of freedom. For  $1 \leq k \leq 11$ , define  $X_k$  to be the first  $k$  columns of  $X$ . Then  $P_k = u_k u_k'$  for  $1 \leq k \leq 11$  and  $P_{12} = XX' - X_{11}X_{11}' = \sum_{k=12}^{20} u_k u_k'$ .

Figure 1 plots the adaptive estimates of mean grape yields for each year, adding linear interpolation between adjacent estimated row means to display their trend. The adaptive symmetric affine shrinkage and adaptive projection estimates are similar visually. Both estimates reveal shared patterns in the fitted harvests for each of the three years. Large dips in estimated mean grape yield occur in the outermost rows and near row 33; and smaller fluctuations occur elsewhere.

The estimated risk  $\hat{R}(\hat{M}_{sym})$  of the adaptive symmetric affine shrinkage estimator is 1.504 while the larger estimated risk  $\hat{R}(\hat{M}_{pro})$  of the adaptive projection estimator is 1.833. Both estimated risks are less than one third of the estimated risk  $\hat{R}(\hat{M}_{ls}) = 3\hat{\sigma}_{ls}^2 = 5.517$  of the least squares estimator. The ordering of the estimated risks agrees with the inequalities (3.9). Moreover, Theorem 4.2 provides grounds for regarding the estimated risks as useful guides to the relative performance of the three estimators.

In this example, the shrinkage factors  $\{\hat{w}_{kj}(\hat{\tau}_k + \hat{w}_{kj})^{-1}: 1 \leq k \leq 12, 1 \leq j \leq 3\}$  that define  $\hat{M}_{sym}$  through (1.16) are as follows. For  $j = 1$ , their values, as  $k$  ranges from 1 to 12 are: 1.00, 1.00, 1.00, .94, .99, .96, .96, .94, .40, .76, .90, .69. For  $j > 1$ , their values are all very near 0. Thus,  $\hat{M}_{sym}$  essentially replaces each  $P_k Y$  in the least squares fit  $\hat{M}_{ls} = \sum_{k=1}^{12} P_k Y$  by its total least squares approximation of rank one and then shrinks these approximations by varying amounts. The shrinkage is greatest when  $k = 9$  and remains notable when  $k = 10$  or 12.

The 0-1 shrinkage factors  $\{I[\hat{w}_{kj} > \hat{\tau}_k]: 1 \leq k \leq 12, 1 \leq j \leq 3\}$  that define  $\hat{M}_{pro}$  through (1.11) display a closely related pattern. For  $j = 1$ , their values as  $k$  ranges from 1 to 12 are 1 when  $k \neq 9$  and 0 when  $k = 9$ . For  $j > 1$ , their values are all 0. Thus,  $\hat{M}_{pro}$  discards  $P_9 Y$  from the least squares fit  $\sum_{k=1}^{12} P_k Y$  and replaces very other  $P_k Y$  by its total least squares approximation of rank one.

Setting  $p = 20$  in specifying  $X$  was a design choice, made outside the theory of this paper, that mediates between sufficient richness in  $\mathcal{R}(X)$  to represent the unknown, possibly

irregular, dependence of mean grape yield on row number and sufficient degrees of freedom to estimate  $\sigma^2$  adequately.

## 5 Correlated Errors

The classical multivariate linear model asserts that the rows of the error matrix  $E$  in (1.1) are independent, identically distributed random vectors, each having a  $N(0, \Sigma)$  distribution. Estimation in this model with more general covariance structure can be mapped into the special case treated in the previous sections. The least squares estimator of  $M$  is still  $XX^+Y$ . Suppose that the covariance matrix  $\Sigma$  is positive definite and known. Consider the *candidate affine shrinkage estimators*

$$\hat{M}(A, \Sigma) = \sum_{k=1}^s P_k Y \Sigma^{-1/2} A_k \Sigma^{1/2}, \quad (5.1)$$

where  $A$  lies in  $\mathcal{A}_{sym}^s$  or in  $\mathcal{A}_{pro}^s$ . The quadratic risk of  $\hat{M}(A, \Sigma)$  is defined to be

$$R(\hat{M}(A, \Sigma), M, \Sigma) = p^{-1} \mathbf{E}[\text{tr}\{(\hat{M}(A, \Sigma) - M)\Sigma^{-1}(\hat{M}(A, \Sigma) - M)'\}]. \quad (5.2)$$

Let  $Z = Y\Sigma^{-1/2}$ ,  $N = M\Sigma^{-1/2}$ , and  $\hat{N}(A, \Sigma) = \hat{M}(A, \Sigma)\Sigma^{-1/2}$ . Then

$$\hat{N}(A, \Sigma) = \sum_{k=1}^s P_k Z A_k \quad (5.3)$$

and

$$R(\hat{M}(A, \Sigma), M, \Sigma) = \mathbf{E}|\hat{N}(A, \Sigma) - N|^2 = p^{-1} \sum_{k=1}^s \mathbf{E}|P_k Z A_k - P_k N|^2. \quad (5.4)$$

The rows of  $Z$  are independent, identically distributed random vectors, each having a  $N(0, I_q)$  distribution. Thus, when  $\Sigma$  is known, the foregoing analysis reduces the problem of oracle or adaptive estimation of  $M$ , under the linear model and loss function of this section, into the problem already treated in this paper when  $\sigma^2$  equals 1.

If the covariance matrix  $\Sigma$  is unknown and  $n$  is substantially greater than  $p$ , we may replace  $\Sigma$ , in the definition of candidate estimator  $\hat{M}(A, \Sigma)$  and in its estimated risk, by

$$\hat{\Sigma}_{ls} = (n - p)^{-1} Y'(I_n - XX^+)Y. \quad (5.5)$$

If both  $p$  and  $n - p$  tend to infinity, arguments from Section 4 show that the loss and estimated risk of  $\hat{M}(A, \hat{\Sigma}_{ls})$  converge together in probability, uniformly over all  $A \in \mathcal{A}_{sym}^s$ . Hence, the loss of the adaptive estimator  $\hat{M}(\hat{A}_{sym}, \hat{\Sigma}_{ls})$ , where  $\hat{A}_{sym}$  minimizes the estimated risk of  $\hat{M}(A, \hat{\Sigma}_{ls})$  over all  $A \in \mathcal{A}_{sym}^s$ , converges in probability to the minimum loss achievable over the candidate estimator class. Counterparts to Theorems 4.1 and 4.2 hold when the estimator of  $\Sigma$  has the  $L_1$ -consistency properties used in Beran (1999).

## References

- Beran, R. (1999). Superefficient estimation of multivariate trend. *Mathematical Methods of Statistics*. **8** 166–180.
- Beran, R. (2007). Adaptive estimators of a mean matrix: total least squares versus total shrinkage. *Econometric Theory*. In Press.
- Chatterjee, S., Handcock, M. S., and Simonoff, J. S. (1995) *A Casebook for a First Course in Statistics and Data Analysis*. Wiley, New York.
- Efron, B. and Morris, C. (1972). Empirical Bayes on vector observations: an extension of Stein’s method. *Biometrika* **59** 335–347.
- Fuller, W. A. (1987). *Error Measurement Models*. Wiley, New York.
- Gleser, L. J. (1981). Estimation in a multivariate “errors in variables” regression model: large sample results. *Annals of Statistics*. **9** 24–44.
- Golub, G. H. and Van Loan, C. F. (1980). An analysis of the total least squares problem. *SIAM Journal on Numerical Analysis*. **17** 883–893.
- Golub, G. H. and Van Loan, C. F. (1996). *Matrix Computations* (third edition). Johns Hopkins University Press, Baltimore.
- Lütkepohl, H. (1996). *Handbook of Matrices*. Wiley, Chicester.
- Pinsker, M. S. (1980). Optimal filtration of square-integrable signals in Gaussian noise. *Problems of Information Transmission* **16** 120–133.
- Stein, C. (1956). Inadmissibility of the usual estimator for the mean of a multivariate normal distribution. In *Proceedings of the Third Berkeley Symposium on Mathematical Statistics and Probability* (J. Neyman, ed.) 197–206. University of California Press, Berkeley.
- Stein, C. (1966). An approach to the recovery of inter-block information in balanced incomplete block designs. In *Festschrift for Jerzy Neyman* (F. N. David, ed.) 351–364. Wiley, New York.
- Van Huffel, S. (2004). Total least squares and errors-in-variables modeling: bridging the gap between statistics, computational mathematics and engineering. In *Compstat 2004: Proceedings in Computational Statistics* (J. Antoch, ed.) 539–555. Physica-Verlag, Heidelberg.

Van Huffel, S. and Vandewalle, J. (1991). *The Total Least Squares Problem: Computational Aspects and Analysis*. SIAM Publications, Philadelphia.