

## AN EMPIRICAL BAYES APPROACH TO MULTIPLE LINEAR REGRESSION<sup>1</sup>

BY SERGE L. WIND

*A. T. & T., Management Sciences Division  
New York*

We consider estimation (subject to quadratic loss) of the vector of coefficients of a multiple linear regression model in which the error vector is assumed to have 0 mean and covariance matrix  $\sigma^2I$  but is not assumed to take on a specific parametric form, e.g., Normal. The vector of coefficients is taken to be randomly distributed according to some unknown prior. Restricted minimax solutions are exhibited relative to equivalence classes on the space of all prior probability distributions which group distributions with the same specified moments. In the context of the classic Empirical Bayes formulation, we determine restricted asymptotically optimal estimators—i.e., decision functions whose Bayes risks converge to the risk of the restricted minimax decision at each component stage.

**1. Introduction and summary.** Estimators are exhibited for the parameters in the multiple linear regression model, which, in addition to using data from the present set of observations, incorporate data obtained in previous (independent) experiments. These estimators have lower Bayes risk than the usual least squares estimators.

The coefficient parameters are assumed to be random variables, distributed according to an unknown prior distribution. The regression problem is assumed to occur repeatedly and independently, with the same prior throughout—the Empirical Bayes formulation of Robbins [10]. The classic Robbins Empirical Bayes approach has been applied to the general linear model with Normally distributed error variables, e.g., see Martz and Krutchkoff [8], where consistent estimators for the (single stage) Bayes rule are proposed as Empirical Bayes procedures. In this paper, no specific parametric form is assumed, and the problem is formulated in the generalized Empirical Bayes approach first considered by Robbins in Section 5 of [10] as an application to estimation of a binomial parameter, and then expanded upon by Cogburn [3]. In this formulation, a restricted minimax decision, defined on an equivalence class, or subset of the set of priors, which groups distributions with the same specified moments, is the decision function whose Bayes risk is less than or equal to the infimum (over all decisions) of the supremum (over an equivalence class) of the Bayes risk function. We determine restricted asymptotically optimal estimators, i.e., decision functions

---

Received March 17, 1971; revised June 6, 1972.

<sup>1</sup> This work is a component of the research submitted for a Ph. D. degree at Columbia University, Dept. of Mathematical Statistics, and was supported in part by National Science Foundation Grant NSF-GP-9640.

*AMS 1970 subject classifications.* Primary 62C10, 62J05; Secondary 62G05.

*Key words and phrases.* Empirical Bayes, regression, Stein–James estimators, restricted minimax.

whose Bayes risks approach the risk of the restricted minimax decision at each component stage.

Satisfying the criterion of restricted asymptotic optimality are components of "Stein-James estimators." These estimators are generalizations and modifications of the one proposed by James and Stein [5] for estimating the mean of a multivariate Normal distribution, and they have been considered as Empirical Bayes procedures by Kantor [6] and Efron and Morris [4] for Normal distributions, and by Cogburn [3] for the case where no parametric distribution is specified. For the latter instance, properties of Stein-James estimators are given by Cogburn [2] and Wind [14].

Definitions of terms and a formalization of this approach in a statistical framework, based on concepts defined in [3], are presented in Section 2.

The regression model is analyzed in Section 3. It is assumed that  $Y = X\beta + \varepsilon$ , where  $X$  is a matrix of known constants,  $\beta$  is a vector of unknown parameters, and  $\varepsilon$  is a random vector with mean 0 and covariance  $\sigma^2 I$ . For quadratic loss, restricted minimax solutions are exhibited for classes of prior distributions on  $\beta$  (Theorem 3.1) and for joint priors on  $(\beta', \sigma^2)$  (Theorem 3.2). For equivalence relations on the priors specifying (i) their means and variances or (ii) second moments, restricted asymptotically optimal decision functions are found in Theorem 3.3 if each independent component regression problem has the same diagonal  $X'X$  matrix and common unknown variance  $\sigma^2$ .

In Section 4, based on the work of Bhattacharya [1] and Sclove [11], a correspondence is developed between the problems of estimating regression coefficients and estimating location parameters. Restricted asymptotically optimal estimators of  $\beta$  relative to equivalences which group priors by their mean vectors and by all components of their covariance matrices are exhibited in Theorem 4.1.

**2. Definitions and statistical framework.** For the statistical decision problem with which we are concerned, the sample space, parameter space, decision space and their  $\sigma$ -fields are given respectively by  $(\mathcal{L}, \mathcal{N})$ ,  $(\Theta, \tau)$ , and  $(\mathcal{E}, D)$ . The family of distributions  $\{F_\theta(\cdot) : \theta \in \Theta\}$  satisfies  $F_\theta(B)$   $\tau$ -measurable for each  $B \in \mathcal{N}$ . The loss function  $L$  is  $\tau \times D$  measurable.  $R(\cdot, \delta)$  denotes the risk function of the randomized (behavioral) decision function  $\delta$ .  $\Delta$  is the space of all such decision procedures.  $L^*(\theta, \delta, x) = \int_{\mathcal{E}} L(\theta, c)\delta(x, dc)$ .  $\Pi$  denotes the set of all (prior) probability distributions on  $\tau$  and  $R^*(\pi, \delta) = \int_{\Theta} R(\theta, \delta)\pi(d\theta)$  is the Bayes risk function.

**DEFINITION 2.1.** The *Bayes envelope* function  $r^*$  defined on  $\Pi$  is  $r^*(\pi) = \inf_{\delta \in \Delta} R^*(\pi, \delta)$ .

**DEFINITION 2.2.**  $\delta_\pi$  is *Bayes* re  $\pi$  if  $R^*(\pi, \delta_\pi) = r^*(\pi)$ .

**DEFINITION 2.3.** Let  $T$  be contained in the class of all subsets of  $\Pi$ . A decision function  $\delta_T$  is *restricted minimax* relative to  $T$  if

$$R^*(\pi', \delta_T) \leq \inf_{\delta} \sup_{\pi \in T} R^*(\pi, \delta), \quad \forall \pi' \in T.$$

If a priori we know that  $\pi \in T$ , it is reasonable, in a minimax sense, to use a restricted minimax decision, if it exists. Thus, if we describe partial information about the prior distribution  $\pi$  by introducing an equivalence relation  $\sim$  on  $\Pi$ , we assume we know  $\pi$  up to  $\sim$ . We will denote the equivalence class generated by  $\pi$  by

$$\sim(\pi) = \{\pi' \in \Pi : \pi' \sim \pi\}.$$

DEFINITION 2.4. The *envelope risk function*  $r_{\sim}^*$ , relative to  $\sim$ , is

$$r_{\sim}^*(\pi) = \inf_{\delta} \sup_{\sim\pi} R^*(\pi', \delta),$$

and is the best (in a minimax sense) we can do given that we know  $\pi$  up to  $\sim$ .

Let us now describe the Empirical Bayes problem. Given  $(\Theta, \tau)$ ,  $(\mathcal{L}, \mathcal{K})$ , and  $(\mathcal{E}, D)$ , measurable spaces, and loss function  $L$  as above. Let  $(\theta_1, x_1), (\theta_2, x_2), \dots$  be a sequence of pairs of random variables, each pair independent of the others, with the  $\theta_i$  having a common a priori distribution  $\pi \in \Pi$ , and the distribution of  $x_i$  given  $\theta_i = \theta$  denoted by  $F_{\theta}(\cdot)$ . So, at the  $n$ th stage or component experiment, we can say that a real vector  $\theta^{(n)} = (\theta_1, \dots, \theta_n)$  is selected by choosing the components independently, each with the same distribution  $\pi$  and then  $x^{(n)} = (x_1, \dots, x_n)$  has probability distribution  $F_{\theta^{(n)}}(\cdot)$  with

$$F_{\theta^{(n)}}(x^{(n)}) = F_{\theta_1}(x_1) \cdot F_{\theta_2}(x_2) \cdot \dots \cdot F_{\theta_n}(x_n).$$

The parameter space is  $\Theta^{(n)}$  and the experiment space  $(\mathcal{L}^{(n)}, \mathcal{K}^{(n)})$ , where the superscript denotes the  $n$ -fold Cartesian product. Let decision function  $\delta_n$  be a transition probability mapping such that  $\delta_n(x^{(n)}, \cdot)$  is a probability measure on the space  $(\mathcal{E}, D)$  for each  $x^{(n)} \in \mathcal{L}^{(n)}$ , and  $\delta_n(\cdot, c)$  is  $\mathcal{K}^{(n)}$ -measurable in  $x^{(n)}$  for each  $c \in D$ ; i.e.,  $\delta_n : \mathcal{L}^{(n)} \rightarrow \mathcal{E}$ , whose form is functionally dependent on  $x^{(n)}$ . The overall expected loss of the  $n$ th stage estimator  $\delta_n$  is

$$R_n^*(\pi, \delta_n) = \int_{\Theta} \pi(d\theta_1) \cdot \dots \cdot \int_{\Theta} \pi(d\theta_n) \int_{\mathcal{L}^{(n)}} L_n^*(\theta_n, \delta_n, x^{(n)}) \cdot F_{\theta^{(n)}}(dx^{(n)}),$$

where

$$L_n^*(\theta_n, \delta_n, x^{(n)}) = \int_{\mathcal{E}} L(\theta_n, c) \delta_n(x^{(n)}, dc).$$

DEFINITION 2.5. Given a sequence of decision functions  $\{\delta_n\}$ ,  $n = 1, 2, \dots$  for the Empirical Bayes problem, if

$$R_n^*(\pi, \delta_n) \rightarrow \inf_{\delta_n} R_n^*(\pi, \delta_n)$$

for all  $\pi \in \Pi$ , then  $\{\delta_n\}$  is called *asymptotically optimal*. As noted in (3.6) of [3],

$$\inf_{\delta_n} R_n^*(\pi, \delta_n) = \inf_{\delta} R^*(\pi, \delta) \equiv r^*(\pi).$$

Suppose we select the equivalence relation that identifies those distributions of  $\Pi$  having the same specified set of moments, a characteristic of  $\pi$  which should be easily estimable from our  $x^{(n)}$  sample. We hope the Bayes risks of our estimators do as well asymptotically as if the equivalence class were known beforehand and at each stage we used the corresponding restricted minimax decision. Thus this criterion, described in [3], becomes

$$(2.1) \quad R_n^*(\pi, \delta_n) \rightarrow r_{\sim}^*(\pi) \quad \forall \pi.$$

DEFINITION 2.6. A sequence of estimators  $\{\delta_n\}$  satisfying (2.1) is called *restricted asymptotically optimal*.

**3. Linear regression.** The multiple linear regression model consists of data

$$Z' = (Y, X),$$

where  $y_j = \beta'x_j + \varepsilon_j$ , ( $j = 1, 2, \dots, N$ ),  $x_j = (x_{j1}, \dots, x_{jp})'$ , with

$$(3.1) \quad E\varepsilon_j = 0 \quad \text{and} \quad E\varepsilon_j \varepsilon_k = \sigma^2 \delta_{jk},$$

where  $\delta_{jk}$  is the Kronecker delta. Or,

$$(3.2) \quad Y = X\beta + \varepsilon,$$

where

- $Y$  is an  $N \times 1$  random observed vector,
- $\varepsilon$  is an  $N \times 1$  random (error) vector with moments (3.1),
- $X$  is an  $N \times p$  matrix of fixed quantities, and
- $\beta$  is a  $p \times 1$  vector of unknown coefficients.

We will always assume that  $X$  is of full rank ( $\text{rank } X = p < N$ ). Let  $V = X'X$ . Then the LSE is  $\hat{\beta} = V^{-1}X'Y$ .

*Normality is not assumed* in (3.1); a partially non-parametric approach is taken by specifying only the first two moments of the error variable.

The problem considered is point estimation of  $\beta$ . A general loss function, under which methods of this paper apply, is

$$(3.3) \quad \sigma^{-2}(\beta - b)' \Delta (\beta - b),$$

with  $\Delta$  a known  $p \times p$  positive definite matrix. However, the loss function

$$(3.4) \quad \sigma^{-2}(\beta - b)' V (\beta - b)$$

will be used frequently, since it can be considered "natural" in the following sense. Suppose we have another set of data  $Z^*$ , independent of  $Z$ , but with the same distribution; i.e.,  $y_k^* = \beta'x_k^* + \varepsilon_k$ ,  $k = 1, \dots, N$  with moments (3.1). We are given  $X^*$  and wish to predict  $Y^*$ , by taking an estimator  $b$  of  $\beta$  which is a function only of  $Z$ . Then the mean squared error of prediction, conditional on  $b$ , appropriately transformed so that the loss of a perfect estimate is zero, can be shown, using the same proof as given by Stein [12] for the Normal linear regression model with random predictors, to be precisely (3.4).

Groups of prior distributions on the  $p$ -vector  $\beta$  are identified by specifying a set of moments of the prior. Assuming variance  $\sigma^2$  is known, the moments considered that determine an equivalence class are the  $p$  means and the  $p(p+1)/2$  elements of the covariance matrix. Each class must be identified by at least  $p$  moments—the second moments of the marginals—and can be identified by as many as  $(p^2 + 3p)/2$  moments, including the means and covariances.

Let  $\mu = (\mu_1, \dots, \mu_p)'$  and let  $\Sigma^*$  be a  $p \times p$  symmetric matrix with components  $\sigma_{ij}^*$ . For the  $i$ th regression coefficient, if only the second moment is specified in

a given equivalence class, set  $\sigma_{ii}^* = E\beta_i^2$  and  $\mu_i = 0$ ; if the variance and mean are specified, set  $\sigma_{ii}^* = \text{Var}(\beta_i)$  and  $\mu_i = E\beta_i$ . For  $j \neq i$ , let  $\sigma_{ij}^* = \text{Cov}(\beta_i, \beta_j)$  if the covariance and the means of  $\beta_i$  and  $\beta_j$  are specified in the class; let  $\sigma_{ij}^* = E\beta_i\beta_j$  if that moment is part of the equivalence and both the means of  $\beta_i$  and  $\beta_j$  are not part of the equivalence specification; let  $\sigma_{ij}^* = 0$  if a covariance is not specified in the class.

Restricted minimax decisions for the (single stage) regression model are given below in Theorems 3.1 (for priors on  $\beta$ ) and 3.2 (for joint priors on  $(\beta', \sigma^2)$ ). First, two useful lemmas are stated.

LEMMA 3.1. *The Bayes estimator  $\delta$  of  $\beta$  for generalized quadratic loss (3.3) and for specified parametric distributions on  $\varepsilon$  and  $\beta$  is  $\delta = E(\beta | \hat{\beta})$ , the mean of the posterior distribution of  $\beta$  given  $\hat{\beta}$ , and*

$$R^*(\pi, \delta) = \sigma^{-2} \sum_{i=1}^p \sum_{j=1}^p d_{ij} E[\text{Cov}(\beta_i, \beta_j | \hat{\beta})],$$

where  $d_{ij}$  are the components of  $\Delta$  and where the expectation is taken with respect to the marginal density of  $\hat{\beta}$ . (See, e.g., [7] pages 18–20.)

LEMMA 3.2. *If  $\delta$  is Bayes with respect to some  $\pi^* \sim \pi$  and  $R^*(\pi', \delta)$  is constant for  $\sim (\pi)$ , then  $\delta$  is restricted minimax relative to  $\sim (\pi)$ .*

THEOREM 3.1. *With loss (3.4), for the regression model (3.2) with  $\sigma^2$  known, relative to a given equivalence class of the priors on  $\beta$ , the restricted minimax estimator is*

$$(3.5) \quad b = (\sigma^2(\Sigma^*)^{-1} + V)^{-1}(\sigma^2(\Sigma^*)^{-1}\mu + V\hat{\beta}),$$

and  $r^*(\pi) = \sum_i \sum_j v_{ij} m_{ij}$ , where  $m_{ij}$  are the components of  $M = (\sigma^2(\Sigma^*)^{-1} + V)^{-1}$ .

PROOF. Consider the parametric subfamilies of Normal distributions of  $\varepsilon$  and  $\beta$ . The posterior distribution of  $\beta$  given LSE  $\hat{\beta}$  is Normal with mean (3.5), which is the Bayes estimator, and covariance  $\sigma^2 M$  ([9] page 337). Since its risk is constant for the equivalence class,  $b$  is restricted minimax for the Normal subfamily and thus for all distributions with first two moments specified.

If  $\sigma^2$  is a random variable, an equivalence class of the joint priors on  $(\beta', \sigma^2)$  must be identified by at least  $p + 1$  moments (the second moments on  $\beta$  and  $E\sigma^2$ ) and by at most  $(p^2 + 3p + 2)/2$ .

THEOREM 3.2. *For the regression model (3.2), subject to loss  $(b - \beta)'V(b - \beta)$ , and with  $\rho = E\sigma^2$  as part of the specification of the equivalence class of the priors on  $(\beta', \sigma^2)$ , the restricted minimax estimator is*

$$(3.6) \quad b^* = (\rho(\Sigma^*)^{-1} + V)^{-1}(\rho(\Sigma^*)^{-1}\mu + V\hat{\beta}),$$

and

$$r^*(\pi) = \rho \sum_i \sum_j v_{ij} l_{ij}, \quad \text{where } L = (\rho(\Sigma^*)^{-1} + V)^{-1}.$$

PROOF. If  $\varepsilon$  is Normally distributed and the prior on  $(\beta', \sigma^{-2})$  is a Normal-gamma ([9] pages 343–345) with parameters  $(\mu, t, \rho(\Sigma^*)^{-1}, \nu)$  with  $\nu, t$  such that  $\rho = \nu t / (\nu - 2)$  and  $\text{rank}(\Sigma^*)^{-1} = p_1$ , the marginal posterior distribution of  $\beta$  given

$\hat{\beta}$  is a multivariate Student distribution ([9] pages 256–259) with mean  $b^*$  (3.6), which is constant risk Bayes, and covariance  $L\nu^*t^*/(\nu^* - 2)$ , where

$$\begin{aligned} \nu^* &= \nu + p_1 + N - p_2, \\ p_2 &= \text{rank}(L^{-1}), \\ t^* &= (\rho(\nu - 2) + \rho\mu'(\Sigma^*)^{-1}\mu + Y'Y - b^{*\prime}L^{-1}b^*)/\nu^*. \end{aligned}$$

But  $E(\nu^*t^*/(\nu^* - 2)) = \nu t/(\nu - 2) = \rho$ , with the expectation taken with respect to  $\hat{\beta}$ .

Consider (3.2) with orthogonal design—i.e.,  $V = X'X$  is diagonal—and those equivalence classes which identify prior distributions on  $\beta$  having the same specified finite means and variances or the same second moments of  $\beta$  (but with no specification of covariance terms). A typical equivalence relation in this subgroup considered is  $(\lambda_1, \dots, \lambda_p)$ , where  $\lambda_i$  is equal either to  $(\mu_i, \Psi_i^2)$  or  $\gamma_i$  for  $i = 1, \dots, p$ , and where  $\mu_i = E\beta_i$ ,  $\Psi_i^2 = \text{Var}(\beta_i)$ , and  $\gamma_i = E\beta_i^2$ . The corollary below follows as a special case of Theorem 3.1 with diagonal  $\Sigma^*$  and  $V$ .

**COROLLARY 3.1.** *Given the regression model (3.2) with orthogonal design,  $b = (b_1, \dots, b_p)'$  is restricted minimax relative to  $(\lambda_1, \dots, \lambda_p)$ , subject to loss (3.4), where*

$$(3.7) \quad \begin{aligned} b_i &= \frac{\Psi_i^2 \hat{\beta}_i + \sigma^2 \nu^{ii} \mu_i}{\sigma^2 \nu^{ii} + \Psi_i^2} & \text{if } \lambda_i &= (\mu_i, \Psi_i^2) \\ &= \frac{\gamma_i \hat{\beta}_i}{\sigma^2 \nu^{ii} + \gamma_i} & \text{if } \lambda_i &= \gamma_i \end{aligned}$$

for  $i = 1, \dots, p$ , and  $R^*(\pi, b) = \sum \phi_i (\phi_i + \sigma^2 \nu^{ii})^{-1}$ , where

$$\begin{aligned} \phi_i &= \Psi_i^2 & \text{if } \lambda_i &= (\mu_i, \Psi_i^2) \\ &= \gamma_i & \text{if } \lambda_i &= \gamma_i. \end{aligned}$$

**PROPERTY 3.1.** For estimating  $\beta_i$ , with marginal prior  $\pi_i$ , and letting  $b_{i1}$  denote the restricted minimax solution (3.7) for  $\lambda_i = (\mu_i, \Psi_i^2)$ , and  $b_{i2}$  denote (3.7) for  $\lambda_i = \gamma_i$ , then

$$R^*(\pi_i, b_{i1}) = \frac{\Psi_i^2}{\Psi_i^2 + \sigma^2 \nu^{ii}} \leq \frac{\gamma_i}{\gamma_i + \sigma^2 \nu^{ii}} = R^*(\pi_i, b_{i2}),$$

with equality only if  $\mu_i = 0$ . Both risks are less than one, the risk of LSE  $\hat{\beta}_i$ .

Thus far, the treatment has been Bayesian. However, suppose the parameters of the prior are unknown, but we have available data from previous regression experiments—the “ $q$ -stage Empirical Bayes regression problem”: Given a sequence of independent regression problems with data  $Z_1, Z_2, \dots, Z_q$  and, respectively, regression coefficients  $\beta_{(1)}, \beta_{(2)}, \dots, \beta_{(q)}$ , where  $Z_i = \left\{ \begin{pmatrix} y_{1i} \\ \vdots \\ y_{Ni} \end{pmatrix} \begin{pmatrix} x_{1i} \\ \vdots \\ x_{pi} \end{pmatrix} \right\}$  is the data of the  $i$ th stage, where  $x_k = (x_{k1}, \dots, x_{kp})'$  and  $\beta_{(i)} = (\beta_{1i}, \dots, \beta_{pi})'$  and

$$(3.8) \quad y_{kj} = \sum_{i=1}^p \beta_{ij} x_{ki} + \epsilon_{kj} = \beta'_{(j)} x_k + \epsilon_{kj},$$

where

- $y_{kj}$  is the  $k$ th observation ( $k = 1, \dots, N$ ) of the  $j$ th component stage ( $j = 1, \dots, q$ );
- $\beta_{ij}$  is the  $i$ th parameter ( $i = 1, \dots, p$ ) of the  $j$ th stage;
- $x_{ki}$  is the independent variable corresponding to the  $k$ th observation and  $i$ th regression coefficient;
- $\varepsilon_{kj}$  has moments (3.1) and  $\varepsilon_{ki}, \varepsilon_{ej}$  are independent for  $j \neq i$ , for all  $k$  and  $c$ .

Each experiment has the same independent variables, the same number of observations  $N$ , and common variance  $\sigma^2$ .  $(\beta_{(1)}, Z_1), (\beta_{(2)}, Z_2), \dots$  is a sequence of pairs of sets of random variables, each pair independent of the others.  $\beta_{(1)}, \beta_{(2)}, \dots$  is a sequence of independent realizations of random variables, each distributed according to the same (unknown) prior distribution.

Set

$$(3.9) \quad \hat{\sigma}^2 = \frac{1}{q(N-p) + 2} \sum_{j=1}^q \sum_{i=1}^N (y_{ij} - \sum_{k=1}^p \hat{\beta}_{kj} x_{ik})^2,$$

$$(3.10) \quad B_{iq}(\hat{\beta}_{iq}) = \left(1 - \frac{(q-1)\hat{\sigma}^2 v^{ii}}{\sum_j (\hat{\beta}_{ij} - \bar{\beta}_i)^2}\right)^+ (\hat{\beta}_{iq} - \bar{\beta}_i) + \bar{\beta}_i,$$

$$(3.11) \quad D_{iq}(\hat{\beta}_{iq}) = \left(1 - \frac{q\hat{\sigma}^2 v^{ii}}{\sum_j \hat{\beta}_{ij}^2}\right)^+ \hat{\beta}_{iq},$$

with

$$(3.12) \quad \bar{\beta}_i = q^{-1} \sum_{j=1}^q \hat{\beta}_{ij},$$

where  $\hat{\beta}_{ij}$  is the  $j$ th stage LSE of  $\beta_{ij}$ , and  $Z^+ \equiv \max(0, Z)$ .

**THEOREM 3.3.** *Given the  $q$ -stage Empirical Bayes regression problem with orthogonal design, with the error variables  $\varepsilon$  drawn from a family of uniformly square integrable distributions with bounded fourth moments, the sequence of restricted asymptotically optimal decision functions relative to  $(\lambda_1, \dots, \lambda_p)$  for estimating  $\beta_{1q}, \dots, \beta_{pq}$  is  $\{\hat{b}_{1q}, \dots, \hat{b}_{pq}\}$ ,  $q = 1, \dots$  where  $\lambda_i = (\mu_i, \Psi_i^2)$  or  $\gamma_i$  and*

$$(3.13) \quad \begin{aligned} \hat{b}_{iq} &= B_{iq} && \text{if } \lambda_i = (\mu_i, \Psi_i^2) \text{ and } q \geq 3 \\ &= D_{iq} && \text{if } \lambda_i = \gamma_i \text{ and } q \geq 3 \\ &= \hat{\beta}_{iq} && \text{if } q = 1, 2 \end{aligned}$$

with  $B_{iq}$  defined in (3.10) and  $D_{iq}$  in (3.11).

**PROOF.** Follows from a general result given as Corollary 2.3 in [14], where Stein–James estimators of this form are considered. (See also [13].)

**PROPERTY 3.2.** The form of the Stein–James estimators (3.10) and (3.11) is suggested by substituting in (3.7) for the now unknown parameters  $\mu_i, \Psi_i^2, \gamma_i$ ,

and  $\sigma^2$ , estimators  $\hat{\mu}_i = \bar{\beta}_i$  (3.12),  $\hat{\Psi}_i^2$ ,  $\hat{\gamma}_i$ , and  $\hat{\sigma}^2$  (3.9) respectively, with

$$\hat{\Psi}_i^2 = [\sum_{j=1}^q (\hat{\beta}_{ij} - \bar{\beta}_i)^2 - (q - 1)\hat{\sigma}^2 v^{ii}]^+ / (q - 1),$$

$$\hat{\gamma}_i = [\sum_{j=1}^q \hat{\beta}_{ij}^2 - q\hat{\sigma}^2 v^{ii}]^+ / q.$$

**THEOREM 3.4.** *Under the conditions of Theorem 3.3, if the variance of  $\varepsilon$  varies with each component problem and a joint prior on  $(\beta', \sigma^2)$  is assumed, then relative to  $(\beta - b)'V(\beta - b)$ , (3.13) is restricted asymptotically optimal with  $\hat{\sigma}^2$  (3.9) replaced by*

$$(3.14) \quad \bar{\sigma}^2 = \frac{1}{q(N - p + 2)} \sum_{i=1}^N \sum_j^q (y_{ij} - \sum_{k=1}^p \hat{\beta}_{kj} x_{ik})^2.$$

**PROOF.** Follows from Corollary 2.4 of [14].

**PROPERTY 3.3.** The results of the above two Theorems hold if in (3.8) the restrictions that the number of observations  $N$  and the independent observations  $X$  remain the same for each experiment are relaxed subject to the condition that  $X'X$  remains constant for each component problem. Then, in place of (3.9), we estimate  $\sigma^2$  by

$$\hat{\sigma}^2 = \frac{1}{\sum_{j=1}^q (N_j - p) + 2} \sum_{j=1}^q \sum_{k=1}^{N_j} (y_{kj} - \sum_{i=1}^p \hat{\beta}_{ij} x_{kij})^2,$$

where  $x_{kij}$  denotes the  $k$ th observation ( $k = 1, \dots, N_j$ ) of the  $i$ th parameter ( $i = 1, \dots, p$ ) at the  $j$ th stage ( $j = 1, \dots, q$ ), and  $N_j$  is the number of observations at stage  $j$ .

**PROPERTY 3.4.** The form of the estimator  $B_{iq}$  (3.10) suggests that we are making a preliminary test of the hypothesis  $H_1: \beta_{ij} = \bar{\beta}_i$  for all  $j = 1, \dots, q$ . We accept  $H_1$  if  $\sum_j (\hat{\beta}_{ij} - \bar{\beta}_i)^2 \leq (q - 1)\hat{\sigma}^2 v^{ii}$ , and then estimate  $\beta_{ij}$  by  $\bar{\beta}_i$ . If we accept  $H_1$ , we are really pooling independent consistent estimators  $\hat{\beta}_{i1}, \dots, \hat{\beta}_{iq}$  and using  $\bar{\beta}_i$ .

**4. Restricted asymptotically optimal solutions.** In this section, a correspondence between estimating regression coefficients in the  $q$ -stage Empirical Bayes formulation and estimating mean vectors also in a  $q$ -stage Empirical Bayes setup is established. This correspondence is employed in Theorem 4.1 to present restricted asymptotically optimal decision functions for the regression problem.

Given, as specified in (3.8), a sequence of  $q$  independent regression problems with data  $Z_1, \dots, Z_q$ , we are interested in estimating the  $p \times q$  matrix  $\beta$ . For convenience in handling the data, (3.8) is rewritten in the following form: Observe a sample of  $N$  observations  $Y_1, \dots, Y_N$  on the dependent variables, where the  $q$ -dimensional vector  $Y_k = (y_{k1}, \dots, y_{kq})'$  is distributed with mean  $\beta'x_k$  and nonsingular covariance matrix  $\Sigma = \sigma^2 I$ . Set the  $q \times N$  matrix  $Y = (Y_1, \dots, Y_N)$  and the  $p \times N$  matrix  $X' = (x_1, \dots, x_N)$ . Then, conditional on  $\beta$ ,  $EY = \beta'X'$ . The loss function is

$$(4.1) \quad L_1(\beta, \mathbf{b}) = \text{tr } \sigma^{-2}(\mathbf{b} - \beta)\Delta(\mathbf{b} - \beta),$$

the  $q$ -variate analog to (3.3), with  $\text{tr}$  denoting the trace.



The least squares estimate of  $\beta$  is  $\hat{\beta} = V^{-1}U$ , with  $U = \sum_{k=1}^N x_k Y_k'$ ,  $V = X'X = \sum_{k=1}^N x_k x_k'$ . Estimate  $\sigma^2$  by (3.9), or identically, set

$$(4.2) \quad \hat{\sigma}^2 = (q(N - p) + 2)^{-1} \text{tr } S,$$

with

$$S = T - U'V^{-1}U = T - \hat{\beta}'V\hat{\beta}, \quad T = \sum_{k=1}^N Y_k Y_k'.$$

The transformations on the regression problem given below were suggested by the mappings in [1] and [11] which assumed a Normal process, but the imposition here of a prior distribution necessitated some modifications.

There exists a nonsingular  $p \times p$  matrix  $C$  such that

$$(4.3) \quad C'C = V \quad \text{and} \quad C'DC = \Delta.$$

Let  $W^{(1)} = C'^{-1}X'$ , and the rows of  $W^{(1)}$  are orthogonal. Following Sclove [11], select the  $(N - p) \times N$  matrix  $W^{(2)}$  such that  $W' = (W^{(1)'}, W^{(2)'})$  is an  $N \times N$  orthogonal matrix. Set  $Z^* = YW'$ . Then  $Z^* = (Z^{(1)}, Z^{(2)}) = Y(W^{(1)'}, W^{(2)'})$ , where  $Z^{(1)}$  is of dimension  $q \times p$  and  $Z^{(2)}$  is  $q \times (N - p)$ . Let  $Z_i = (Z_{1i}, \dots, Z_{qi})'$  be a  $q$ -vector for  $i = 1, \dots, N$ . Then  $Z^* = (Z_1, \dots, Z_p, Z_{p+1}, \dots, Z_N)$ . It is easy to verify that  $\hat{\beta} = C^{-1}Z^{(1)'}$ ,  $EZ^{(1)} = (E\beta)'C'$ ,

$$S = Z^{(2)}Z^{(2)'} = \sum_{k=p+1}^N Z_k Z_k',$$

and  $EZ^{(2)} = 0$ , with  $E\beta$  the mean of  $\beta$ . Set the  $q \times p$  matrix  $\theta = (\theta_1, \dots, \theta_p) = \beta'C' = E(Z^{(1)} | \beta)$ . Then, for  $i = 1, \dots, p$ ,  $Z_i = \hat{\beta}'(C')_i$  has mean  $\theta_i$  and covariance  $\Sigma = \sigma^2 I$ , conditional on  $\beta$ , where  $(C')_i$  denotes the  $i$ th column of  $C'$ . For  $i = p + 1, \dots, N$ ,  $Z_i$  has mean 0 and covariance  $\Sigma$ . Conditional on  $\beta$ ,  $Z_i, Z_j$  are mutually orthogonal and  $Z^{(1)}$  and  $Z^{(2)}$  lie in orthogonal spaces.

Let  $P$  be a  $p \times p$  orthogonal matrix which rotates the axis of the distribution such that for each row of the transformed variable

$$(4.4) \quad K = Z^{(1)}P',$$

the elements are pairwise uncorrelated. Recall  $\Sigma^* = \text{Cov}(\beta)$ .

Thus  $P(\sigma^2 I + C\Sigma^*C')P'$  is the diagonal covariance matrix for the data of each component experiment. ( $P$  is the matrix such that the elements of  $P(Z^{(1)'})_i$ , for each  $i = 1, \dots, q$ , are uncorrelated unconditionally where  $(Z^{(1)'})_i$  is the  $i$ th column of  $Z^{(1)'}$ .) Let  $K = (K_1, \dots, K_p)$ , each  $K_i$  a  $q$ -vector.  $K_i, K_j$  ( $j \neq i$ ) are unconditionally uncorrelated and the components of each  $K_i$  are mutually independent.

Under transformations (4.3) and (4.4),  $K = \hat{\beta}'C^* = YW^{(1)'}P'$  with  $C^* = PC$  and  $C^*C^* = V$ . Let  $\hat{\theta}$  be an estimator of  $\theta = \beta'C'$  and let  $\hat{\eta}$  be an estimator of  $\eta = \beta'C^*$ ; define the corresponding estimator  $\mathbf{b}$  of  $\beta$  as  $\mathbf{b} = C^{*-1}\hat{\eta}'$ .

PROPERTY 4.1. For estimating  $\eta$ , loss (4.1) is

$$(4.5) \quad L_2(\eta, \hat{\eta}) = \sigma^{-2} \sum_{i=1}^q d_i(\hat{\eta}_i - \eta_i)'(\hat{\eta}_i - \eta_i) = L_1(\beta, \mathbf{b}).$$

We can now see how the problem of estimating  $\beta$  with loss (4.1) is related to

the problem of estimating mean  $\eta$  subject to (4.5) with data  $K$ , with uncorrelated column vectors  $K_1, \dots, K_p$ , with components of each  $K_i$  independent, and data  $S$ , which lies in a space orthogonal to  $K$ . For an estimator  $\varphi(K_1, \dots, K_p, S)$  of  $\eta$ , define the corresponding estimator for  $\beta$ ,

$$(4.6) \quad \Psi(\hat{\beta}, S) = C^{*-1}\varphi'(\hat{\beta}'(C^{*'})_1, \dots, \hat{\beta}'(C^{*'})_p, S).$$

With respect to losses (4.1) and (4.5), let

$$R_1(\beta, \Psi) = EL_1(\beta, \Psi) \quad \text{and} \quad R_2(\eta, \varphi) = EL_2(\eta, \varphi).$$

PROPERTY 4.2.  $R_2(\eta, \varphi) = R_1(\beta, \Psi)$ .

It follows that (4.6) preserves dominance relations among estimators, and the theorem below can be verified.

THEOREM 4.1. *Given the  $q$ -stage Empirical Bayes regression problem (3.8) with the error variables  $\varepsilon$  drawn from a family of uniformly square integrable distributions with bounded fourth moments, the restricted asymptotically optimal estimator at stage  $q$ , relative to a given equivalence class of priors of  $\beta$  specified by its first two moments and loss (3.4), is*

$$C^{*-1}(\delta_{q1}(K_{q1}), \dots, \delta_{qp}(K_{qp}))',$$

with  $K = \hat{\beta}'C^{*}$ , and for  $i = 1, \dots, p$ ,

$$\begin{aligned} \delta_{qi}(K_{qi}) &= \delta_i(K_{qi}) && \text{if } q \geq 3 \text{ and } \rho_i \neq 0 \\ &= \delta_i^*(K_{qi}) && \text{if } q \geq 3 \text{ and } \rho_i = 0 \\ &= \bar{K}_i && \text{if } q < 3, \end{aligned}$$

$$\delta_i(K_{qi}) = \left(1 - \frac{(q-1)\hat{\sigma}^2}{\sum (K_{ji} - \bar{K}_i)^2}\right)^+ (K_{qi} - \bar{K}_i) + \bar{K}_i,$$

$$\delta_i^*(K_{qi}) = \left(1 - \frac{q\hat{\sigma}^2}{\sum_j K_{ji}^2}\right)^+ K_{qi},$$

$$\bar{K}_i = q^{-1} \sum_{j=1}^q K_{ji},$$

$$\rho_i = \sum_{d=1}^q c_{id}^* \mu_d,$$

with  $\mu_d$  as specified prior to Theorem 3.1 and  $\hat{\sigma}^2$  given by (4.2).

PROPERTY 4.3. If the component problems do not have common variance  $\sigma^2$ , let  $\Sigma$  be a diagonal matrix with elements  $\sigma_1^2, \dots, \sigma_q^2$ . For a joint prior on  $(\beta', \sigma^2)$ , replace  $\hat{\sigma}^2$  in Theorem 4.1 with (3.14), and the results obtain subject to loss  $(b - \beta)' \Delta (b - \beta)$ .

**Acknowledgment.** I wish to express my appreciation to Professors Herbert Robbins and A. J. Baranchik for their encouragement and guidance of my Ph. D. dissertation research, of which this work is a part.

REFERENCES

[1] BHATTACHARYA, P. K. (1966). Estimating the mean of a multivariate normal population with general quadratic loss function. *Ann. Math. Statist.* 37 1819-1824.

- [2] COGBURN, R. (1965). On the estimation of a multivariate location parameter with squared error loss, in *Bernoulli (1723), Bayes (1763) and Laplace (1813) Anniversary Volume*, ed. J. Neyman and L. LeCam. Springer-Verlag, Berlin, 24–29.
- [3] COGBURN, R. (1967). Stringent solutions to statistical decision problems. *Ann. Math. Statist.* **38** 447–463.
- [4] EFRON, B. and MORRIS, C. (1972). Limiting the risk of Bayes and empirical Bayes estimators—Part II: The empirical Bayes case. *J. Amer. Statist. Assoc.* **67** 130–139.
- [5] JAMES, W. and STEIN, C. (1961). Estimation with quadratic loss. *Proc. Fourth Berkeley Symp. Math. Statist. Prob.* **1**. Univ. of California Press, 361–379.
- [6] KANTOR, M. (1967). Estimating the mean of a multivariate normal distribution with applications to time series and empirical Bayes estimation. Ph. D. dissertation, Columbia Univ.
- [7] MARTZ, H. (1967). Empirical Bayes estimation in multiple linear regression. Ph. D. dissertation, Virginia Polytechnic Institute.
- [8] MARTZ, H. and KRUTCHKOFF, R. (1969). Empirical Bayes estimators in a multiple linear regression model. *Biometrika* **56** 367–374.
- [9] RAIFFA, H. and SCHLAIFER, R. (1961). *Applied Statistical Decision Theory*. The M.I.T. Press.
- [10] ROBBINS, H. (1964). The empirical Bayes approach to statistical decision problems. *Ann. Math. Statist.* **35** 1–20.
- [11] SCLOVE, S. (1967). Improved estimation of regression parameters. Technical Report No. 125, Dept. of Statistics, Stanford Univ.
- [12] STEIN, C. (1960). Multiple regression, in *Contributions to Probability and Statistics—Essays in Honor of Harold Hotelling*. Stanford Univ. Press, 424–443.
- [13] WIND, S. (1970). An empirical Bayes approach to the multiple linear regression problem. Ph. D. dissertation, Columbia Univ.
- [14] WIND, S. (1972). Stein–James estimators of a multivariate location parameter. *Ann. Math. Statist.* **43** 340–343.

AMERICAN TELEPHONE AND TELEGRAPH COMPANY  
MANAGEMENT SCIENCES DIVISION  
130 JOHN STREET  
ROOM 1761  
NEW YORK, NEW YORK 10038