# EMPIRICAL BAYES ESTIMATION IN A MULTIPLE LINEAR REGRESSION MODEL*

R. S. SINGH

## Summary

Estimation of the vector $\beta$ of the regression coefficients in a multiple linear regression $Y = X\beta + \varepsilon$ is considered when $\beta$ has a completely unknown and unspecified distribution and the error-vector $\varepsilon$ has a multivariate standard normal distribution. The optimal estimator for $\beta$, which minimizes the overall mean squared error, cannot be constructed for use in practice. Using $X$, $Y$ and the information contained in the observation-vectors obtained from $n$ independent past experiences of the problem, (empirical Bayes) estimators for $\beta$ are exhibited. These estimators are compared with the optimal estimator and are shown to be asymptotically optimal. Estimators asymptotically optimal with rates near $O(n^{-1})$ are constructed.

## 1. Introduction

Empirical Bayes (EB) approach to a statistical problem, introduced by Robbins [22], and later developed by Johns [7], Robbins [23], [24], Samuel [25] and Johns and Van Ryzin [8], [9], among others, has recently drawn considerable attention in the literature. Suppose there is a pair $(X, \theta)$ of random variables, where $X$ is observable and $\theta$, called the parameter, is unobservable. The r.v. $X$ given $\theta$ has a specified distribution $P_\theta$ on the observation space $\mathcal{X}$ and $\theta$ has a completely *unknown* and *unspecified* distribution $G$ on the parameter space $\Theta$. Based on an observation on $X$ (which could be a sufficient statistic for $\theta$), the problem is to decide about $\theta$ with respect to a nonnegative loss function. If the prior distribution $G$ were known, the statistician would use the *Bayes procedure* $\phi_G$ which achieves the *minimum Bayes risk*, say $R(G)$, with respect to $G$. (We assume such a $\phi_G$ exists). But since

$G$ is not known, the *optimal procedure* $\phi_G$ is not available for the use to the statistician. In the EB context we assume that the above problem, called the *component problem*, has occurred independently in the past, say $n$ times, so that there are $n+1$ independent pairs $(X_1, \theta_1)$, $\cdots, (X_n, \theta_n)$ and $(X, \theta)$. The object is to utilize the information contained in the past observations $(X_1, \cdots, X_n)$ and the present observation $X$ to produce a decision rule for the present parameter $\theta$ so that for large $n$ this rule is "nearly" as good as the unavailable optimal rule $\phi_G$ in the sense that the overall risk, say $R_n$, of this rule approximates the *Bayes envelope* $R(G)$ achieved by $\phi_G$. If $R_n \to R(G)$ as $n \to \infty$ then the rule is called asymptotically optimal (a.o.), (Robbins [22]). If for a $\delta > 0$ $R_n - R(G) = O(n^{-\delta})$ as $n \to \infty$ we will say that the rule is a.o. with rate $O(n^{-\delta})$.

Johns and Van Ryzin [8], [9] considered EB linear loss two-action hypothesis testing problems in one parameter discrete [8] and continuous [9] exponential families and exhibited a.o. test procedures with rates near $O(n^{-1})$ in both cases. Lin [13], [14] considered EB squared error loss estimation (SELE) in the two families considered by Johns and Van Ryzin and obtained EB estimators with rates near $O(n^{-1/3})$. Yu [34] and Singh [27], [29] exhibited a.o. EB estimators in the one parameter Lebesgue exponential family with rates near $O(n^{-1/3})$, $O(n^{-2/5})$ and $O(n^{-1})$ respectively. O'Bryan and Susarla [20] considered EB estimation with varying sample sizes in component problems in the univariate normal distribution and obtained estimators a.o. with rates near $O(n^{-1/3})$. Recently Monte-Carlo simulations performed by Clemmer and Krutchkoff [2], Maritz [16], Martz and Krutchkoff [18], Griffin and Krutchkoff [5], Bennet and Martz [1] and Maritz and Lwin [17], among others, have shown how for certain priors the EB procedures often perform better compared to the usual procedures whenever there is at least one past experience of the problem.

We, in this paper, consider EB approach to the SELE problem in a multiple linear regression model $Y = X\beta + \varepsilon$ with loss function $L(\beta, \hat{\beta}) = (\beta - \hat{\beta})'(\beta - \hat{\beta})$. Using the information contained in $(Y_1, \cdots, Y_n)$ from the past problems and $(Y, X)$ from the present problem we will exhibit for each $n > 0$ two classes of estimators $\hat{\phi}$ and $\tilde{\phi}$ for $\beta$, one for the case when nothing is known about the support of the prior distribution and the other for the case when it is known that the prior distribution has a compact support. We show that $\hat{\phi}$ are a.o. with rates $O(n^{-1+\eta})$ uniformly over the class of all priors satisfying certain moment conditions dependent on $\eta$, whereas $\tilde{\phi}$ are shown to be a.o. with rates $O(n^{-1+\eta})$ uniformly over the class of all priors with compact support. Thus we have given procedures for constructing EB estimators a.o. with

rates arbitrarily close to $O(n^{-1})$ in the above multiple linear regression model. As pointed out in Singh [29] EB procedures a.o. with rates $O(n^{-1})$ have not yet been exhibited in any EB problem other than that involving discrete distributions.

Martz and Krutchkoff [18] have considered EB estimation in the general linear model considered here. They have, however, not proved any consistency or asymptotic optimality of their estimators, but have shown, through a Monte-Carlo simulation, that for certain priors their EB estimators perform better than the usual least square estimator. Wind [33] has considered EB estimation of $\beta$ when the error vector $\boldsymbol{\varepsilon}$ is assumed to have $\mathbf{0}$ mean and covariance $\sigma^2 \boldsymbol{I}$ but is not assumed to take on a specific parametric form, e.g. Normal. For priors specifying their means and variances or specifying their second moments, he has exhibited restricted asymptotically optimal EB estimators of $\beta$, that is estimators whose Bayes risks converge to the risk of the restricted minimax estimator at each component stage.

Some other important related contributions worth mentioning here are due to Stein [31], James and Stein [6], Cogburn [3], Kantor [10], Wind [32], Efron and Morris [4] and Rao [21], among others, where problems related to estimation of regression coefficient $\beta$ or of a multivariate normal distribution mean are considered.

## 2. The model and preliminaries

We consider the following multiple linear regression model

$$(2.1) \qquad\qquad \boldsymbol{Y} = \boldsymbol{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$$

where $\boldsymbol{Y}$ is an $l \times 1$ vector of random observations, $\boldsymbol{X}$ is an $l \times p$ matrix of known constants so that $\boldsymbol{X}'\boldsymbol{X}$ is invertible, $\boldsymbol{\beta}$ is the $p \times 1$ vector of unknown regression coefficients having an *unknown* and *unspecified* prior distribution $G$ on $(E^p, \mathscr{B}^p)$, $\mathscr{B}^p$ being the Borel-field of subsets of $E^p$, the $p$-dimentional Euclidean space, and $\boldsymbol{\varepsilon}$ is an $l \times 1$ vector of unobservable random variables. We will assume that the conditional distribution of $\boldsymbol{\varepsilon}$ given $\boldsymbol{\beta}$ is $N_p(\mathbf{0}, \sigma^2 \boldsymbol{I})$, the $p$-variate normal distribution with mean vector $\mathbf{0}$ and variance-covariance matrix $\sigma^2 \boldsymbol{I}$, where $\sigma^2$ is assumed to be known. The object is to estimate $\boldsymbol{\beta}$ with respect to the quadratic loss function

$$(2.2) \qquad\qquad L(\boldsymbol{\beta}, \boldsymbol{\phi}) = \|\boldsymbol{\beta} - \boldsymbol{\phi}\|_2^2$$

where $\boldsymbol{\phi}$ is an estimate of $\boldsymbol{\beta}$, and for a $1 \times p$ vector $\boldsymbol{t}' = (t_1, \cdots, t_p)$ and $m > 0$, $\|\boldsymbol{t}\|_m = \left( \sum_{i=1}^{p} |t_i|^m \right)^{1/m}$.

For an estimator $\boldsymbol{\phi}$ of $\boldsymbol{\beta}$, let $R(\boldsymbol{\phi}, G)$ be the (Bayes) risk of $\boldsymbol{\phi}$ with

respect to the prior distribution $G$, i.e.,

(2.3)                              $R(\phi, G) = \mathrm{E} \|\phi - \beta\|_2^2$

where E stands for the expectation operator with respect to all the random variables involved in the expression. The *Bayes envelope* with respect to $G$ is given by

$$R(G) = \inf_{\phi} R(\phi, G)$$

where the inf is taken over the set of all estimators $\phi$ for which (2.3) exist. The estimator which achieves the Bayes envelope $R(G)$ is the *Bayes estimator*, also called *optimal estimator* (o.e.) $\phi_G$ given by

(2.4)                              $\phi_G(Y) = \mathrm{E}(\beta | Y)$ .

Thus $R(\phi_G, G) = R(G)$. Notice that $R(G)$ can be exactly achieved only if the prior distribution $G$ is known and $\beta$ is estimated by the o.e. $\phi_G$. Unfortunately $G$ is completely unknown, and hence $\phi_G$ is not available to us for use. This motivates us to use the EB approach to exhibit estimators whose risks are close to $R(G)$ achived by $\phi_G$.

## 3. Empirical Bayes approach

Suppose we have incurred $n$ independent experiences of the above estimation problem, called the *component problem*, in the past. That is we have independent pairs

$$\{Y_1, \beta_1\}, \cdots, \{Y_n, \beta_n\}$$

from the past experiments, and $(Y, \beta)$ from the present experiment, with

$$Y_i = X\beta_i + \varepsilon_i , \qquad i = 1, \cdots, n .$$

The vectors $Y_i$, $\beta_i$, $\varepsilon_i$ behave like $Y$, $\beta$, $\varepsilon$ described above; $\beta_1, \cdots, \beta_n$ and $\beta$, are i.i.d. according to the same unknown prior distribution $G$; and given $\beta_1, \cdots, \beta_n$ and $\beta$, the vectors $\varepsilon_1, \cdots, \varepsilon_n$ are i.i.d. according to $N_p(0, \sigma^2 I)$ and $Y_i \sim N_p(X\beta_i, \sigma^2 I)$. We will exhibit (EB) estimators for $\beta$ on the basis of $X$, the past observations $Y_1, \cdots, Y_n$ and the present observation $Y$; and show that for large $n$ these estimators can be considered as good as the unavailable o.e. $\phi_G$.

In Section 4 we will show how the problem of approximation of the optimal estimator can be reduced to the problem of estimation of the marginal density of $\hat{\beta}$, the maximum likelihood estimate of $\beta$, and its first partial derivatives. In Section 5 we exhibit for each integer $r > 1$ uniformly mean squared consistent estimators of the marginal den-

sity of $\hat{\beta}$ and its partial derivatives. Using these estimates we will introduce two EB estimators $\hat{\phi}$ and $\tilde{\phi}$ of $\beta$, one for the case when nothing is known about the support of the prior distribution $G$ of $\beta$ and the other for the case when the support of $G$ is a known compact subset of $E^p$. In Section 6 we will show that the excess risk of $\tilde{\phi}$ over the minimum risk $R(G)$ is of order $n^{-(r-1)\gamma/(2r+p)}$ for every $0<\gamma<2$, and that of $\hat{\phi}$ is of order $n^{-(r\gamma-2)/(2r+p)}$ for every $\gamma$ in $(0,2)$ for which $\int \|\beta\|_1^{pr+/(2-r)} dG(\beta)$ is finite. Thus the rates of convergence can be sharpened to any desired degree by choosing appropriately a larger value of $r$.

## 4.  Reduction in the problem

Recall from Section 1 that an EB estimator based on the observations from the past $n$ experiences of the problem and the present observation is a.o. if its risk approaches to the minimum Bayes risk $R(G)$ as $n$ gets large. The following known lemma therefore reduces the problem of searching an a.o. estimator for $\beta$ to the one of searching a quadratic mean squared consistent estimator of the o.e. $\phi_G$. The univariate version of the lemma is proved (under the present condition) in Singh [29].

LEMMA 4.1.  *Let the prior distribution $G$ be such that $R(G)<\infty$. Let $\phi$ be an arbitrary $p$-vector statistic.  Then*

$$R(\phi, G)-R(G)=\mathrm{E}\,\|\phi-\phi_G\|_2^2$$

*where the expectation is taken with respect to all the random variables involved in the expression.*

PROOF.  Let $\mathrm{E}_*$ stand for the expectation operator conditional on $Y$ and all other random variables involved in $\phi$. Then

$$\mathrm{E}_*\,[(\phi-\beta)'(\phi-\beta)]=(\phi-\phi_G)'(\phi-\phi_G)+\mathrm{E}_*\,\{(\phi_G-\beta)'(\phi_G-\beta)$$
$$+2(\phi-\phi_G)'(\phi_G-\beta)\}\ .$$

Since $R(G)=\mathrm{E}\,(\phi_G-\beta)'(\phi_G-\beta)<\infty$, the second term above is $\mathrm{E}_*\,(\phi_G-\beta)'$ $\cdot(\phi_G-\beta)+2(\phi-\phi_G)'\,\mathrm{E}_*\,(\phi_G-\beta)=\mathrm{E}_*\,\|\phi_G-\beta\|_2^2$ since $\mathrm{E}_*\,(\beta)=\phi_G$. Thus

$$\mathrm{E}_*\,\|\phi-\beta\|_2^2=\|\phi-\phi_G\|_2^2+\mathrm{E}_*\,\|\phi_G-\beta\|_2^2\ .$$

Now taking expectations on both sides we get the conclusion of the lemma.

Since $\sigma^2$ is known, without loss of generality we can (and do) take $\sigma^2=1$. The least square estimator of $\beta$ is

(4.1) $$\hat{\beta} = \Sigma X'Y \qquad \text{where } \Sigma = (X'X)^{-1} .$$

The estimator (4.1) is the maximum likelihood estimator (MLE) as well as the minimum variance unbiased estimator (MVUE) of $\beta$.

To reduce the problem further we will now show that the o.e. in the EB context $E(\beta | Y_1, \cdots, Y_n, Y)$ is the same as the o.e. (2.4) in the component problem. Then we will show how the o.e. (2.4) can be expressed as a function of the marginal p.d.f. of $\hat{\beta}$ and its first order partial derivatives.

Notice that the conditional distribution of $\hat{\beta}$ given $\beta$ is $N_p(\beta, \Sigma)$. Hence $\hat{\beta}$ is sufficient for $\beta$ and

(4.2) $$\phi_G = E(\beta | Y) = E(\beta | \hat{\beta}) = \phi_G(\hat{\beta}) .$$

If $f(\hat{\beta} | \beta)$ denotes the conditional p.d.f. of $\hat{\beta}$ given $\beta$, then

(4.3) $$f(\hat{\beta} | \beta) = \frac{|\Sigma|^{1/2}}{(2\pi)^{p/2}} \exp \left\{ -\frac{1}{2} (\hat{\beta} - \beta)' \Sigma^{-1} (\hat{\beta} - \beta) \right\}$$

where $|\Sigma| = \det \Sigma$. The marginal p.d.f. of $\hat{\beta}$ is given by

(4.4) $$f(\hat{\beta}) = \int f(\hat{\beta} | \beta) dG(\beta) .$$

For $j = 1, \cdots, n$, let $\hat{\beta}_j$ be the MLE of $\beta_j$ in the $j$th experiment that we had in the past. Then

$$\hat{\beta}_j = \Sigma X'Y_j , \qquad j = 1, \cdots, n$$

and $\hat{\beta}_1, \cdots, \hat{\beta}_n$ and $\hat{\beta}$ are i.i.d. according to the marginal p.d.f. given by (4.4). Thus the *optimal estimator in the* EB *context* given by $E(\beta | Y_1, \cdots, Y_n, Y) = E(\beta | \hat{\beta}_1, \cdots, \hat{\beta}_n, \hat{\beta})$ is simply $E(\beta | \hat{\beta}) = \phi_G$. Thus the estimators which are as good as the o.e. $\phi_G$ in the component problem are also as good as the o.e. in the EB context.

It follows from Theorem 2.9 of Lehmann [12] that the $m$th order partial derivative of $f(\hat{\beta})$ with respect to the $i$th component $\hat{\beta}_i$ of $\hat{\beta}$ is

(4.5) $$\frac{\partial^m f(\hat{\beta})}{\partial \hat{\beta}_i^m} = \int \frac{\partial^m f(\hat{\beta} | \beta)}{\partial \hat{\beta}_i^m} dG(\beta)$$

$$\text{for all } m = 1, 2, \cdots \text{ and } i = 1, 2, \cdots, p .$$

Now notice that

$$\frac{\partial f(\hat{\beta} | \beta)}{\partial \hat{\beta}_i} = \begin{bmatrix} \partial f(\hat{\beta} | \beta) / \partial \hat{\beta}_1 \\ \vdots \\ \partial f(\hat{\beta} | \beta) / \partial \hat{\beta}_p \end{bmatrix} = \Sigma^{-1} (\beta - \hat{\beta}) f(\hat{\beta} | \beta) .$$

Therefore

$$(4.6) \qquad \beta f(\hat{\beta}|\beta) = \hat{\beta} f(\hat{\beta}|\beta) + \Sigma(\partial f(\hat{\beta}|\beta)/\partial\hat{\beta}) \ .$$

Hence from (4.2), (4.5) and (4.6) it follows that $\phi_G$ can be written as

$$(4.7) \qquad \phi_G = \mathrm{E}(\beta|\hat{\beta}) = (f(\hat{\beta}))^{-1} \int \beta f(\hat{\beta}|\beta) dG(\beta) = \hat{\beta} + \Sigma q(\beta)$$

where $q' = (q_1, \cdots, q_p)$ and

$$(4.8) \qquad q_s = f_s/f \quad \text{and} \quad f_s(\hat{\beta}) = \partial f(\hat{\beta})/\partial\hat{\beta}_s$$

with $\hat{\beta}_s$ denoting the $s$th component of $\hat{\beta}$.

Thus the problem of approximating the minimum expected loss estimator $\phi_G$ reduces to the problem of estimating the marginal p.d.f. $f$ and its first order partial derivatives $f_1, \cdots, f_p$. The representation (4.7) of $\phi_G$ is also noted in Martz and Krutchkoff [18] and Singh [28].

## 5. Proposed empirical Bayes estimators and an important result

In this section first we will exhibit estimators of the marginal p.d.f. $f$ and its first partial derivatives $f_s$. We will then obtain a bound for the mean square errors of these estimates. Finally, on the basis of these estimates we will exhibit two classes of EB estimates for $\beta$, one for the case when nothing is known about the support of the prior distribution and the other for the case when the prior distribution has compact support in $E^p$, the $p$-dimensional Euclidean space.

### 5.1. *Estimation of $f$ and $f_s$*

For an integer $r > 1$ and for $i = 0, 1$, let $\mathcal{K}_i^r$ be the class of all Borel-measurable bounded functions vanishining outside $(0, 1)$ such that for $K_0$ in $\mathcal{K}_0^r$,

$$\int y^j K_0(y) dy = \begin{cases} 1 & \text{if } j=0 \\ 0 & \text{if } j=1, \cdots, r-1 \end{cases}$$

and for $K_1$ in $\mathcal{K}_1^r$,

$$\int y^j K_1(y) dy = \begin{cases} 1 & \text{if } j=1 \\ 0 & \text{if } j=0, 2, 3, \cdots, r-1 \ . \end{cases}$$

Let $0 < h = h(n)$ be a function of $n$ such that $h \to 0$ as $n \to \infty$. For $j = 1, \cdots, n$, let $\hat{\beta}_{j1}, \cdots, \hat{\beta}_{jp}$ denote the components of the MLE $\hat{\beta}_j$ of $\beta_j$. At $x' = (x_1, \cdots, x_p)$, we estimate $f$ by

$$(5.1) \qquad \hat{f}(\boldsymbol{x}) = (nh^p)^{-1} \sum_{j=1}^{n} \left\{ \underset{(i=1)}{\overset{p}{\times}} K_0\left(\frac{\hat{\beta}_{ji} - x_i}{h}\right) \right\}$$

and $f_s(\boldsymbol{x}) = \partial f(\boldsymbol{x})/\partial x_s$ by

$$(5.2) \qquad \hat{f}_s(\boldsymbol{x}) = (nh^{p+1})^{-1} \sum_{j=1}^{n} \left\{ K_1\left(\frac{\hat{\beta}_{js} - x_s}{h}\right) \underset{\substack{i=1 \\ i \neq s}}{\overset{p}{\times}} \left( K_0\left(\frac{\hat{\beta}_{ji} - x_i}{h}\right) \right) \right\}.$$

Estimates (5.1) and (5.2) are special cases of those given in Singh ([26], p. 37 and [30]) where the author has given nonparametric method of estimation of a mixed partial derivatives of a multivariate density. For $0 < \gamma \leq 2$, the following theorem proves the $\gamma$th mean consistency of estimators $\hat{f}$ and $\hat{f}_s$ of $f$ and $f_s$ respectively for each $s = 1, \cdots, p$. The numbers $c_1, c_2, \cdots$ below are absolute constants.

THEOREM 5.1. *For every* $0 < \gamma \leq 2$,

$$(5.3) \qquad \mathrm{E}\,|\hat{f}(\boldsymbol{x}) - f(\boldsymbol{x})|^{\gamma} \leq c_1 f^{\gamma}(\boldsymbol{x})\{h^{r\gamma}B^{\gamma}(\boldsymbol{x}) + (nh^p)^{-\gamma/2}C^{\gamma/2}(\boldsymbol{x})\}$$

*and*

$$(5.4) \qquad \mathrm{E}\,|\hat{f}_s(\boldsymbol{x}) - f_s(\boldsymbol{x})|^{\gamma} \leq c_2 f^{\gamma}(\boldsymbol{x})\{h^{(r-1)\gamma}B^{\gamma}(\boldsymbol{x}) + (nh^{p+2})^{-\gamma/2}C^{\gamma/2}(\boldsymbol{x})\}$$

*where, with*

$$(5.5) \qquad \lambda_0 = \text{max. root } \Sigma^{-1}.$$

$$(5.6) \qquad t_1(\hat{\boldsymbol{\beta}}, \boldsymbol{\beta}) = \exp(h\lambda_0 \|\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}\|_1) \sum_{j=0}^{r} (\|\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}\|_1^j + hp^j)$$

*and*

$$(5.7) \qquad t_2(\hat{\boldsymbol{\beta}}, \boldsymbol{\beta}) = \exp(h\lambda_0 \|\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}\|_1)(f(\hat{\boldsymbol{\beta}}))^{-1},$$

*B and C are given by*

$$B(\boldsymbol{x}) = \mathrm{E}\,[t_1(\hat{\boldsymbol{\beta}}, \boldsymbol{\beta}) \,|\, \hat{\boldsymbol{\beta}} = \boldsymbol{x}]$$

*and*

$$C(\boldsymbol{x}) = \mathrm{E}\,[t_2(\hat{\boldsymbol{\beta}}, \boldsymbol{\beta}) \,|\, \hat{\boldsymbol{\beta}} = \boldsymbol{x}].$$

*Remark* 5.1. Notice that $f^2(\boldsymbol{x})B^2(\boldsymbol{x})$ and $f^2(\boldsymbol{x})C(\boldsymbol{x})$ are bounded in $\boldsymbol{x}$ by a constant independent of $n$ and $G$. Hence if $h$ is taken proportional to $n^{-1/(p+2r)}$, then Theorem 5.1 proves that $\sup_{\boldsymbol{x}} \mathrm{E}\,(\hat{f}(\boldsymbol{x}) - f(\boldsymbol{x}))^2 \leq$ .const. $n^{-2r/(p+2r)}$ and $\sup_{\boldsymbol{x}} \mathrm{E}\,(\hat{f}_s(\boldsymbol{x}) - f_s(\boldsymbol{x}))^2 \leq$ const. $n^{-2(r-1)/(p+2r)}$ uniformly for all $n \geq 1$ and for all prior distributions.

PROOF OF THE THEOREM. We will prove (5.3). Proof of (5.4) follows analogously. Liapunov's inequality followed by $c_r$-inequality (Loève

[15], p. 155) gives

$$(5.8) \qquad \mathrm{E}\,|\hat{f}(\boldsymbol{x}) - f(\boldsymbol{x})|^r \leq \{(\mathrm{E}\,\hat{f}(\boldsymbol{x}) - f(\boldsymbol{x}))^2 + \mathrm{var}\,(\hat{f}(\boldsymbol{x}))\}^{r/2}$$
$$\leq |\mathrm{E}\,\hat{f}(\boldsymbol{x}) - f(\boldsymbol{x})|^r + (\mathrm{var}\,(\hat{f}(\boldsymbol{x})))^{r/2}\ .$$

We now obtain bounds for each term on the right hand side of (5.8). Since $\hat{\boldsymbol{\beta}}_1, \cdots, \hat{\boldsymbol{\beta}}_n$ are i.i.d. with marginal p.d.f. $f$, it follows that

$$\mathrm{E}\,\hat{f}(\boldsymbol{x}) = \int \cdots \int \Big( \underset{i=1}{\overset{p}{\times}} K_0(z_i) \Big) f(\boldsymbol{x} + h\boldsymbol{z}) dz_1 \cdots dz_p\ .$$

Substituting $f(\boldsymbol{x} + h\boldsymbol{z})$ by its $r$th order Taylor expansion about $\boldsymbol{x}$ with Lagrange-form of the remainder at the $r$th term and then making use of the orthogonality properties of $K_0$ and the fact that $K_0$ vanishes outside $(0, 1)$, we get

$$(5.9) \quad |\mathrm{E}\,\hat{f}(\boldsymbol{x}) - f(\boldsymbol{x})| \leq \mathrm{const.}\ h^r \sup_{\|\boldsymbol{z}\|_1 \leq p} \Big| \Big\{ \sum_{\substack{i_1=0 \\ i_1+\cdots+i_p=r}}^{r} \cdots \sum_{i_p=0}^{r} \frac{\partial^r f(\boldsymbol{t})}{\partial t_1^{i_1} \cdots \partial t_p^{i_p}} \Big| \boldsymbol{t} =$$

$$(\boldsymbol{x} + h\boldsymbol{z}) \Big\} \Big|\ .$$

Let $\boldsymbol{C}_i$ be the $i$th column of $\Sigma^{-1}$. Using Schwarz-inequality at the second step below we get for some constants $a_0, a_1, \cdots$

$$\frac{\partial^{i_k} \hat{f}(\boldsymbol{\beta} | \boldsymbol{\beta})}{\partial \hat{\beta}_i^{i_k}} = \Big| \sum_{j=0}^{i_k} a_j (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})' (\boldsymbol{C}_i)^j f(\hat{\boldsymbol{\beta}} | \boldsymbol{\beta}) \Big|$$

$$\leq (\lambda_0 \vee 1)^{i_k} \max_{0 \leq j \leq i_k} |a_j| \Big( \sum_{j=0}^{i_k} \|\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}\|_1^j \Big) f(\hat{\boldsymbol{\beta}} | \boldsymbol{\beta})$$

where $\lambda_0$ is given by (5.5). This inequality and (4.5) applied to (5.9) give

$$(5.10) \quad |\mathrm{E}\,\hat{f}(\boldsymbol{x}) - f(\boldsymbol{x})| \leq c_3 h^r \int \Big( \sum_{j=0}^{r} \sup_{\|\boldsymbol{z}\|_1 \leq p} \|\boldsymbol{x} + h\boldsymbol{z} - \boldsymbol{\beta}\|_1^j f(\boldsymbol{x} + h\boldsymbol{z} | \boldsymbol{\beta}) \Big) dG(\boldsymbol{\beta})\ .$$

Also, since $\sup\limits_{\|\boldsymbol{z}\|_1 \leq p} |\boldsymbol{z}' \Sigma^{-1} (\boldsymbol{x} - \boldsymbol{\beta})| \leq \lambda_0 \|\boldsymbol{x} - \boldsymbol{\beta}\|_1$, we have

$$(5.11) \qquad \sup_{\|\boldsymbol{z}\|_1 \leq p} \|\boldsymbol{x} + h\boldsymbol{z} - \boldsymbol{\beta}\|_1^j f(\boldsymbol{x} + h\boldsymbol{z} | \boldsymbol{\beta})$$
$$\leq 2^{j-1} (\|\boldsymbol{x} - \boldsymbol{\beta}\|_1^j + h p^j) f(\boldsymbol{x} | \boldsymbol{\beta}) \exp (h \lambda_0 \|\boldsymbol{x} - \boldsymbol{\beta}\|_1)\ .$$

This last inequality applied to (5.10) gives

$$(5.12) \qquad \mathrm{E}\,|\hat{f}(\boldsymbol{x}) - f(\boldsymbol{x})| \leq \mathrm{const.}\ f(\boldsymbol{x}) B(\boldsymbol{x})$$

where $B(\boldsymbol{x})$ is as given in the theorem. Now consider $\mathrm{var}\,(\hat{f}(\boldsymbol{x}))$. Since $\hat{\boldsymbol{\beta}}_1, \cdots, \hat{\boldsymbol{\beta}}_n$ are i.i.d. with marginal p.d.f. $f$; and for a r.v. $T$, $\mathrm{var}\,(t) \leq \mathrm{E}\,(T^2)$,

$$(5.3) \qquad \operatorname{var}(\hat{f}(\boldsymbol{x})) \leq (nh^p)^{-1} \int_0^1 \cdots \int_0^1 \left( \underset{i=1}{\overset{p}{\times}} K_0^2(z_i) \right) f(\boldsymbol{x}+hz) dz_1 \cdots dz_p$$

$$\leq \operatorname{const.} (nh^p)^{-1} \int \sup_{\|z\|_1 \leq p} f(\boldsymbol{x}+hz \,|\, \beta) dG(\beta)$$

$$\leq \operatorname{const.} (nh^p)^{-1} f(\boldsymbol{x}) C(\boldsymbol{x})$$

where $C(\boldsymbol{x})$ is as given in the theorem, and the last inequality follows by arguments used for (5.11). Proof of (5.3) is now complete from (5.8), (5.12) and (5.13).

### 5.2. *Proposed* EB *estimators for* $\beta$

Lemma 4.1, followed by Equation (4.7), shows that the problem of obtaining an a.o. estimator of $\beta$ is essentially the problem of obtaining a quadratic-mean consistent estimator of $\boldsymbol{q}(\hat{\beta})$ where $\boldsymbol{q}=(q_1, \cdots, q_p)$ and $q_s = f_s/f$. Having noted that $\hat{f}$ and $\hat{f}_s$, given by (5.1) and (5.2) are mean square consistent estimators of $f$ and $f_s$ respectively, we are now able to propose our EB estimators of $\beta$ for the two cases indicated below.

*Notation.* For $b>0$ and for a number $a$, $[a]_b$ will stand for $-b$, $a$ or $b$ according as $a < -b$, $|a| \leq b$ or $a > b$. For any vector $\boldsymbol{t}' = (t_1, \cdots, t_p)$, $[\boldsymbol{t}]_c$ means $([t_1]_c, \cdots, [t_p]_c)$.

*Case* 1. If nothing is known about the support of the prior distribution $G$, then our proposed EB estimator for $\beta$ is

$$(5.14) \qquad \hat{\boldsymbol{\phi}}(\hat{\beta}) = \hat{\beta} + \Sigma \hat{\boldsymbol{q}}(\hat{\beta})$$

where $\hat{\boldsymbol{q}}' = (\hat{q}_1, \cdots, \hat{q}_p)$, and for $s = 1, \cdots, p$,

$$(5.15) \qquad \hat{q}_s(\hat{\beta}) = [\hat{f}_s(\hat{\beta})/\hat{f}(\hat{\beta})]_{h^{-1}} .$$

*Case* 2. If it is known that the support of $G$ is in a compact subset $[-A, A]^p$ of $E^p$ and $A$ is known, then our proposed EB estimator for $\beta$ is

$$(5.16) \qquad \tilde{\boldsymbol{\phi}}(\hat{\beta}) = [\hat{\beta} + \Sigma \tilde{\boldsymbol{q}}(\hat{\beta})]_A$$

where $\tilde{\boldsymbol{q}}' = (\tilde{q}_1, \cdots, \tilde{q}_p)$ and for $s = 1, \cdots, p$

$$(5.17) \qquad \tilde{q}_s(\hat{\beta}) = \hat{f}_s(\hat{\beta})/\hat{f}(\hat{\beta}) .$$

### 6. Asymptotic optimality of the EB estimators $\hat{\phi}$ and $\tilde{\phi}$ and speed of convergence

In this section we will show that the EB estimators $\hat{\phi}$ and $\tilde{\phi}$ of

$\beta$ are asymptotically as good as the unavailable optimal estimator $\phi_G$ whatever be the $G$ as long as it satisfies a certain moment condition. Theorems in this section demonstrate that EB estimators a.o. with rates arbitrarily close to $O(n^{-1})$ can be constructed. Using a variation of (2.3), we will denote the risks of $\hat{\phi}$ and $\tilde{\phi}$ by $R_n(\hat{\phi}, G)$ and $R_n(\tilde{\phi}, G)$ respectively.

THEOREM 6.1. *Let h be proportional to* $n^{-1/(p+2r)}$. *If for some* $0 < \gamma < 2$, *the prior distribution G satisfies*

$$(6.0) \qquad \int \|\beta\|_1^{pr+/(2-r)} dG(\beta) < \infty ,$$

*then*

$$(6.1) \qquad R_n(\hat{\phi}, G) - R(G) \leq c_4 n^{-(r\gamma-2)/(2r+p)} \qquad for\ every\ n \geq 1$$

*where* $c_4$ *is independent of* $n$ *and* $\gamma$.

PROOF. Fix a $\gamma$ in $(0, 2)$. Let $\lambda_0^*$ denote the max. root of $\Sigma$. Then Lemma 4.1 followed by (4.7) and (5.14) gives

$$(6.2) \qquad \begin{aligned} R_n(\hat{\phi}, G) - R(G) &= \mathrm{E}\,\|\hat{\phi} - \phi_G\|_2^2 \\ &= \mathrm{E}\,\|\Sigma(\hat{q}(\beta) - \hat{q(\beta)})\|_2^2 \\ &\leq \lambda_0^{*2} \sum_{s=1}^p \mathrm{E}\,[\hat{q}_s(\hat{\beta}) - q_s(\hat{\beta})]^2 . \end{aligned}$$

Temporarily, let $t \neq 0$ and $t'$ be numbers and $T$ and $T'$ be random variables. Then we prove that for every $a > 0$,

$$(6.3) \qquad \mathrm{E}\left|\frac{t'}{t} - \left[\frac{T'}{T}\right]_a\right|^2 \leq 8|t'/t|^2 I(|t'/t| > a) + 16a^{2-r}|t|^{-r}$$
$$\times \{\mathrm{E}\,|T' - t'|^r + 2a^r\,\mathrm{E}\,|T - t|^r\} .$$

(The inequality is true for $\gamma = 2$ too). To prove (6.3), notice that for numbers $b_1$ and $b_2$,

$$|b_1 - [b_2]_a| \leq (|b_1 - b_2| \wedge 2a) I(|b_1| \leq a) + (a + |b_1|) I(|b_1| > a) ,$$

and by a use of $c_r$-inequality (Loève [15], p. 155) we get

$$|b_1 - [b_2]_a|^2 \leq 2(2a)^{2-r}(|b_1 - b_2| \wedge 2a)^r I(|b_1| \leq a) + 8|b_1|^2 I(|b_1| > a) .$$

This inequality followed by the lemma in the Appendix of Singh [28] gives (6.3).

Thus (6.3) in conjunction with the definitions $q_s = f_s/f$ from (4.8) and $\hat{q}_s = [\hat{f}_s/\hat{f}]_{h^{-1}}$ from (5.14) and Theorem 5.1 gives

(6.4)                          $E[\hat{q}_s(\hat{\beta}) - q_s(\hat{\beta})]^2 \leqq c_5(D_1 + D_2 + D_3)$

where, with $B$ and $C$ as given in Theorem 5.1,

$$D_1 = E\{|q_s(\hat{\beta})|^2 I(|q_s(\hat{\beta})| > h^{-1})\}$$

$$D_2 = h^{rr-2} E\{B^r(\hat{\beta})\}$$

and                          $D_3 = n^{-r/2} h^{-2-(rp/2)} E\{C^{r/2}(\hat{\beta})\}$ .

First consider $D_1$. Hölder-inequality followed by Markov-inequality gives for a $\delta > 1$,

(6.5)          $D_1 \leqq (E|q_s(\hat{\beta})|^{2\delta})^{\delta^{-1}} (E I(|q_s(\hat{\beta})| > h^{-1}))^{1-\delta^{-1}}$

               $\leqq h^{rr-2}(E|q_s(\hat{\beta})|^{2\delta})^{\delta^{-1}} (E|q_s(\hat{\beta})|^{(rr-2)/(1-\delta^{-1})})^{1-\delta^{-1}}$ .

Notice that if $C_s$ is the $s$th column of $\Sigma^{-1}$, then

$$\left| \frac{\partial f(\hat{\beta}|\beta)}{\partial \hat{\beta}_s} \right| = |(\hat{\beta} - \beta)' C_s| f(\hat{\beta}|\beta) \leqq \|C_s\|_1 \|\hat{\beta} - \beta\|_1 f(\hat{\beta}|\beta)$$ .

Thus by (4.5), $|q_s(\hat{\beta})| = |f_s(\hat{\beta})/f(\hat{\beta})| \leqq \|C_s\|_1 E(\|\hat{\beta} - \beta\|_1|\hat{\beta})$. Thus by Hölder-inequality for every $\eta \geqq 1$

$$E|q_s(\hat{\beta})|^\eta \leqq \|C_s\|_1^\eta E\|\hat{\beta} - \beta\|_1^\eta < \infty$$ .

Similarly it can be shown that for every $0 < \eta < 1$, $E|q_s(\hat{\beta})|^\eta < \infty$. Hence from (6.5) we have

(6.6)                          $h^{-(rr-2)} D_1 < \infty$      for all $n$ .

Now to see that $h^{-(rr-2)} D_2$ is finite, notice that for an integer $j \geqq 0$ and any $\eta \geqq 1$, Hölder inequality yields

(6.7)          $E\{E(\|\hat{\beta} - \beta\|_1^j \exp(c_6\|\hat{\beta} - \beta\|_1)|\hat{\beta})\}^\eta$

               $\leqq E[\|\hat{\beta} - \beta\|_1^{j\eta} \exp(c_6\eta\|\hat{\beta} - \beta\|_1)] < \infty$ .

Similarly, it can be shown that the left hand side of (6.7) for $0 < \eta < 1$ too is finite. Thus from the definition of $B(\hat{\beta})$ given in Theorem 5.1, we conclude that

(6.8)                          $h^{-(rr-2)} D_2 < \infty$      whatever be $n$ .

Finally consider $D_3$. Notice that from (5.7)

$$E(C^{r/2}(\hat{\beta})) = E[(f(\hat{\beta}))^{-r/2}(E\{\exp(h\lambda_0\|\hat{\beta} - \beta\|_1)|\hat{\beta}\})^{r/2}]$$ .

Take an $\eta$ such that $\gamma/2 < \eta < 1$. Then by Hölder inequality

(6.9)                          $E(C^{r/2}(\hat{\beta})) \leqq c_7[E(f(\hat{\beta}))^{-\eta}]^{r/2\eta}$ .

Since $E(f(\hat{\beta}))^{-\eta} = \int f^{1-\eta}(\hat{\beta})d\hat{\beta}$; $E[(f(\hat{\beta}))^{-\eta}I(\|\hat{\beta}\|_1 \leq p)]$ is finite. According to our hypothesis on $G$, there exists an $\varepsilon > 0$ such that $E\|\beta\|_1^{p\eta(2-\eta)^{-1}+\varepsilon} < \infty$. Take $\xi = p - 1 + \varepsilon(1-\eta)\eta^{-1}$. Then by Hölder inequality

$$E(f(\hat{\beta}))^{-\eta}I(\|\hat{\beta}\| \geq p) \leq I_1 \cdot I_2$$

where

$$I_1 = \left\{ \int \|\hat{\beta}\|_1^{-1-\xi}I(\|\hat{\beta}\|_1 > p)d\hat{\beta} \right\}^\eta$$

and

$$I_2 = \{E\|\hat{\beta}\|_1^{(1+\xi)\eta/(1-\eta)}\}^{1-\eta}.$$

Now by a use of the multivariate polar coordinate transformation (see Kendall and Stuart [11], pp. 246–247) and the inequalities $\|\hat{\beta}\|_2^2 \leq \|\hat{\beta}\|_1^2 \leq p\|\hat{\beta}\|_2^2$ it can be shown that $I_1 < \infty$ since $\xi + 1 > p$. Further

$$I_2^{1/(1-\eta)} \leq \text{const.} (1 + E(\|\beta\|_1^{(1+\xi)\eta/(1-\eta)}) < \infty$$

by our hypothesis, since $(1+\xi)\eta/(1-\eta) = p\eta(1-\eta)^{-1} + \varepsilon$. Hence we conclude that

(6.10) $\qquad n^{\tau/2}h^{2+(\tau p/2)}D_3 < \infty \qquad$ for all $n \geq 1$.

Now the proof of the theorem is complete from (6.2), (6.4), (6.6), (6.8) and (6.10).

*Remark* 6.1. Notice that the moment assumption (6.0) on the prior distribution $G$ is used in the proof of the theorem only to prove

(6.11) $\qquad E(f(\hat{\beta}))^{-\eta} < \infty \qquad$ for an $\eta$ in $(0, 1)$.

Professor Dennis Gilliland has provided the following counter-example showing that (6.11) is not true for all $G$ whatever be $\eta$ in $(1/2, 1)$. Consider only the univariate case. It suffices to show that with $\varepsilon = 1 - \eta$, in $(0, 1/2]$ there exists a $G$ for which the integral,

$$I = \int \left[ \int \exp\left( -\frac{1}{2}(x-\beta)^2 dG(\beta) \right) \right]^\cdot dx$$

is not finite. Consider a prior $G$ which puts mass $c/k^2$ on $k = 1, 2, \cdots$ where $c = \left( \sum_1^\infty 1/k^2 \right)^{-1}$. Then

$$I = c^\varepsilon \int \left[ \sum_{k=1}^\infty \frac{\exp((x-k)^2/2)}{k^2} \right]^\cdot dx.$$

But $\left[ \sum_{k=1}^\infty \{\exp(-(x-k)^2/2)/k^2\} \right]^\cdot \geq \exp(-\varepsilon(x-k)^2/2)/k^{2\varepsilon}$ for all $x$ and $k$.

Hence applying this inequality to the interval $k-1/2 \leq x \leq k+1/2$, we see that

$$I \geq c^{\varepsilon} \lim_{K \to \infty} \sum_{k=1}^{K} \left\{ \int_{k-1/2}^{k+1/2} \exp\left(-\varepsilon(x-k)^2/2\right) dx/k^{2\varepsilon} \right\} .$$

But since $\int_{k-1/2}^{k+1/2} \exp\left(-\varepsilon(x-k)^2/2\right) dx \geq \exp(-\varepsilon/8)$ for all $k$,

$$I \geq c^{\varepsilon} \exp(-\varepsilon/8) \lim_{K \to \infty} \sum_{k=1}^{K} 1/k^{2\varepsilon} = \infty \qquad \text{for } \varepsilon \leq 1/2 .$$

Now we will prove the asymptotic optimality of our second EB estimator $\tilde{\phi}$ of $\beta$.

THEOREM 6.2. *Let $h$ be proportional to $n^{-1/(2r+p)}$. If the support of the prior distribution $G$ is in some known compact subset $[-A, A]^p$ of $E^p$ and $\tilde{\phi}$ is given by (5.16), then for every $0 < \gamma < 2$*

$$(6.12) \qquad R_n(\tilde{\phi}, G) - R(G) \leq c_8 n^{-(r-1)\gamma/(2r+p)}$$

*where $c_8$ is independent of $n$ and $\gamma$.*

PROOF. From (4.7) and (5.16), the $i$th components of $\hat{\phi}_G(\hat{\beta})$ and $\tilde{\phi}(\hat{\beta})$ can be expressed as ratios $t_i(\hat{\beta})/f(\hat{\beta})$ and $[T_i(\hat{\beta})/\hat{f}(\hat{\beta})]_A$ for $i = 1, \cdots, p$. Since the support of $G$ is in $[-A, A]^p$, $t_i/f$ is in $[-A, A]$ for each $i = 1, \cdots, p$. Therefore, by Lemma 4.1,

$$(6.13) \qquad R_n(\tilde{\phi}, G) - R(G) = \mathrm{E} \|\tilde{\phi}(\hat{\beta}) - \hat{\phi}_G(\hat{\beta})\|_2^2$$
$$\leq \sum_{i=1}^{p} \mathrm{E}\left( \left| \frac{t_i(\hat{\beta})}{f(\hat{\beta})} - \frac{T_i(\hat{\beta})}{\hat{f}(\hat{\beta})} \right| \wedge 2A \right)^2 .$$

Notice that $t_i$ is a linear combination of $f, f_1, \cdots, f_p$ and that $T_i$ of $\hat{f}, \hat{f}_1, \cdots, \hat{f}_p$; the coefficient of $f(\hat{\beta})$ in $t_i(\hat{\beta})$ and that of $\hat{f}(\hat{\beta})$ in $T_i(\hat{\beta})$ being the $i$th component of $\hat{\beta}$. Therefore, (6.13) followed by the lemma in the Appendix of Singh [28] and Theorem 5.1 here gives for $0 < \gamma < 2$

$$(6.14) \qquad R_n(\hat{\phi}, G) - R(G) \leq c_9 n^{-(r-1)\gamma/(p+2r)} (\mathrm{E}\, B^r(\hat{\beta}) + \mathrm{E}\, C^{r/2}(\hat{\beta})) .$$

The proof of the theorem is complete since in the proof of Theorem 6.1 it is shown that $\mathrm{E}\, B^r(\hat{\beta})$ and $\mathrm{E}\, C^{r/2}(\hat{\beta})$ are finite constants for $0 < \gamma < 2$.

*Remark* 6.2. It may be recalled that the constants $C_4$ in (6.1) and $C_8$ in (6.12) are independent of $n$ and $r$. Further, the rate result in (6.1) in uniform over the class of all priors $G$ satisfying (6.0), and that in (6.12) is uniform over the class of all priors whose supports are a

subset of $[-A, A]^p$.

For each integer $r > 1$ we have exhibited two sets of EB estimators $\hat{\phi}$ and $\tilde{\phi}$ for the general multiple linear regression model $Y = X\beta + \varepsilon$. Theorem 6.1 gives sufficient condition on the distribution $G$ of $\beta$ under which $\hat{\phi}$ is a.o. with rates $O(n^{-(rr-2)/(2r+p)})$ for every choice of $\gamma$ in $(0, 2)$. Theorem 6.2 shows that if $G$ has a compact support then estimators $\tilde{\phi}$ a.o. with improved rates $O(n^{-(r-1)\gamma/(2r+p)})$ for any $\gamma$ in $(0, 2)$ can be constructed. This paper, thus, suggests how to construct EB estimators a.o. with rates arbitrarily close to $O(n^{-1})$ in the regression model $Y = X\beta + \varepsilon$.

## Acknowledgments

UNIVERSITY OF GUELPH

## REFERENCES

[1] Bennett, Kemble G. and Martz, H. F. (1972). A continuous empirical Bayes smoothing technique, *Biometrika*, **59**, 361-368.

[2] Clemmer, B. A. and Krutchkoff, Richard G. (1968). The use of empirical Bayes estimators in a linear regression model, *Biometrika*, **55**, 525-534.

[3] Cogburn, R. (1965). On the estimation of a multivariate location parameter with squared error loss, *Bernoulli (1723), Bayes (1763) and Laplace (1813) Anniversary Volume* (eds. J. Neyman and L. LeCam), Springer-Verlag, Berlin, 24-29.

[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. Ass.*, **67**, 130-139.

[5] Griffin, Barry S. and Krutchkoff, Richard G. (1971). Optimal linear estimators: an empirical Bayes version with application to the binomial distribution, *Biometrika*, **58**, 195-201.

[6] James, W. and Stein, C. (1961). Estimation with quadratic loss, *Proc. Fourth Berkeley Symp. Math. Statist. Prob.*, **1**, Univ. California Press, 361-379.

[7] Johns, M. V., Jr. (1957). Nonparametric empirical Bayes procedures, *Ann. Math. Statist.*, **28**, 649-669.

[8] Johns, M. V., Jr. and Van Ryzin, J. R. (1971). Convergence rates for empirical Bayes two-action problems I: Discrete case, *Ann. Math. Statist.*, **42**, 1521-1539.

[9] Johns, M. V., Jr. and Van Ryzin, J. R. (1972). Convergence rates for empirical Bayes two-action problems II: Continuous case, *Ann. Math. Statist.*, **43**, 934-947.

[10] 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.

[11] Kendall, Maurice and Stuart, Alan. (1969). *The Advanced Theory of Statistics, 3rd ed. Vol. I*, Hafner Pub. Co., New York.

[12] Lehmann, E. L. (1959). *Testing of Statistical Hypothesis*, Wiley, New York.

[13] Lin, P. E. (1972). Rates of convergence in empirical Bayes estimation problems: Discrete case, *Ann. Inst. Statist. Math.*, **24**, 319-325.

[14] Lin, P. E. (1975). Rates of convergence in empirical Bayes estimation problems: Continuous case, *Ann. Statist.*, **3**, 155-164.

[15] Loève, Michel (1963). *Probability Theory*, *3rd ed.*, Van Nostrand, Princeton.

[16] Maritz, J. S. (1969). Empirical Bayes estimation for continuous distributions, *Biometrika*, **56**, 349-359.

[17] Maritz, J. S. and Lwin, T. (1975). Construction of simple empirical Bayes estimators, *J. R. Statist. Soc.*, B, **75**, 421-425.

[18] Martz, H. and Krutchkoff, R. (1969). Empirical Bayes estimators in a multiple linear regression model, *Biometrika*, **56**, 367-374.

[19] Neyman, J. (1962). Two breakthroughs in the theory of statistical decision making, *Rev. Int. Statist. Inst.*, **30**, 11-27.

[20] O'Bryan, Thomas E. and Susarla, V. (1976). Rates in the empirical Bayes estimation problem with nonidentical components: Case of normal distributions, *Ann. Inst. Statist. Math.*, **28**, 389-397.

[21] Rao, C. Radhakrishna (1975). Simultaneous estimation of parameters in different linear models and applications to biometric problems, *Biometrika*, **31**, 545-554.

[22] Robbins, Herbert (1955). An empirical Bayes approach to statistics, *Proc. 3rd Berkeley Symp. Math. Statist. Prob.*, **1**, 157-163, University California Press.

[23] Robbins, Herbert (1963). The empirical Bayes approach to the testing of statistical hypothesis, *Rev. Int. Statist. Inst.*, **31**, 195-208.

[24] Robbins, Herbert (1964). The empirical Bayes approach to statistical decision problems, *Ann. Math. Statist.*, **35**, 1-20.

[25] Samuel, E. (1963). An empirical Bayes approach to the testing of certain parametric hypotheses, *Ann. Math. Statist.*, **34**, 1370-1385.

[26] Singh, R. S. (1974). Estimation of derivatives of average of $\mu$-densities and sequence-compound estimation in exponential families, RM-318, Dept. Statist. Prob., Michigan State University.

[27] Singh, R. S. (1976). Empirical Bayes estimation with convergence rates in non-continuous Lebesgue-exponential families, *Ann. Statist.*, **4**, 431-439.

[28] Singh, R. S. (1977). Applications of estimators of a density and its derivatives to certain statistical problems, *J. R. Statist. Soc.*, B, **39**, 357-363.

[29] Singh, R. S. (1979). Empirical Bayes estimation in Lebesgue-exponential families with rates near the best possible rate, *Ann. Statist.*, **7**, 890-902.

[30] Singh, R. S. (1981). Speed of convergence in nonparametric estimation of a multivariate $\mu$-density and its mixed partial derivatives, *J. Statist. Plann. Inf.*, **5**, 287-298.

[31] Stein, C. (1960). Multiple regression, *Contributions to Probability and Statistics—Essays in Honor of Harold Hotelling*, Stanford University Press, 424-443.

[32] Wind, S. (1972). Stein-James estimators of a multivariate location parameter, *Ann. Math. Statist.*, **43**, 340-343.

[33] Wind, S. (1973). An empirical Bayes approach to multiple linear regression, *Ann. Statist.*, **1**, 93-103.

[34] Yu, Benito (1971). Rates of convergence in empirical Bayes two-action and estimation problems and in sequence-compound estimation problems, RM-279, Dept. Statist. Prob., Michigan State University.