

## ITERATIVE WEIGHTED LEAST SQUARES ESTIMATORS<sup>1</sup>

BY JIAHUA CHEN AND JUN SHAO

*University of Waterloo and University of Ottawa*

In a heteroscedastic linear model, we establish the asymptotic normality of the iterative weighted least squares estimators with weights constructed by using the within-group residuals obtained from the previous model fitting. An adaptive procedure is proposed which ensures that the iterative process stops after a finite number of iterations and produces an estimator asymptotically equivalent to the best estimator that can be obtained by using the iterative procedure. Theoretical and empirical results of the performance of the adaptive estimator are presented.

**1. Introduction.** One of the most useful models in statistical applications is the following general linear model:

$$(1.1) \quad y_{ij} = x'_{ij}\beta + e_{ij}, \quad j = 1, 2, \dots, n_i, \quad i = 1, 2, \dots, k,$$

where  $y_{ij}$  are responses,  $x_{ij}$  are values of a  $p$ -dimensional covariate,  $x'_{ij}$  is the transpose of  $x_{ij}$ ,  $\beta$  is a  $p$ -vector of unknown parameters and  $e_{ij}$  are random errors. Usually, the  $e_{ij}$  are assumed to be mutually independent,  $E(e_{ij}) = 0$  for all  $i, j$ , and the variances  $\text{var}(e_{ij})$  exist but are unknown and not necessarily equal.

Let  $\theta = g(\beta)$  be the parameter of interest, where  $g$  is a known function from  $R^p$  to  $R$ . If  $\text{var}(e_{ij}) = \sigma^2$  for all  $i, j$ , the customary estimator of  $\theta$  is the ordinary least squares estimator (OLSE) given by

$$\hat{\theta}_o = g(\hat{\beta}_o), \quad \hat{\beta}_o = \left( \sum_{i=1}^k \sum_{j=1}^{n_i} x_{ij}x'_{ij} \right)^{-1} \sum_{i=1}^k \sum_{j=1}^{n_i} x_{ij}y_{ij}.$$

In many situations the errors in the same group,  $e_{ij}$ ,  $j = 1, 2, \dots, n_i$ , have a common distribution. For example,  $x_{ij} = x_i$  for all  $j$  and  $y_{ij}$ ,  $j = 1, 2, \dots, n_i$ , are replications. But the errors from different groups,  $e_{ij}$  and  $e_{i'j'}$  with  $i \neq i'$ , may have different distributions. Hence  $\sigma_i^2 = \text{var}(e_{ij})$  are not necessarily the same. Under such a case, the OLSE may be improved by a weighted least squares estimator (WLSE)

$$(1.2) \quad \hat{\theta}_w = f(\hat{\beta}_w), \quad \hat{\beta}_w = \left( \sum_{i=1}^k \sum_{j=1}^{n_i} w_i x_{ij}x'_{ij} \right)^{-1} \sum_{i=1}^k \sum_{j=1}^{n_i} w_i x_{ij}y_{ij}$$

Received March 1991; revised March 1992.

<sup>1</sup>Research supported by the Natural Sciences and Engineering Research Council of Canada.

AMS 1991 subject classifications. Primary 62J05; secondary 60F05.

Key words and phrases. Asymptotic normality, combining groups, adaptive estimator, efficiency.

for some weights  $w_i > 0, i = 1, 2, \dots, k$ . If  $\sigma_i^2$  are known,

$$(1.3) \quad \tilde{\theta} = g(\tilde{\beta}), \quad \tilde{\beta} = \left( \sum_{i=1}^k \sum_{j=1}^{n_i} \sigma_i^{-2} x_{ij} x'_{ij} \right)^{-1} \sum_{i=1}^k \sum_{j=1}^{n_i} \sigma_i^{-2} x_{ij} y_{ij}$$

improves  $\hat{\theta}_o$  and has some optimality property. Since  $\sigma_i^2$  are unknown, we need to estimate  $\sigma_i^2$  by  $\hat{\sigma}_i^2$  and use  $w_i = \hat{\sigma}_i^{-2}$  in (1.2).

In some situations  $\sigma_i^2$  are related to the  $x_{ij}$  or  $x'_{ij}\beta, j = 1, 2, \dots, n_i$ , through a smooth but unknown function  $h$ , that is,  $\sigma_i^2 = h(x_{i1}, x_{i2}, \dots, x_{in_i})$ . In these situations  $\sigma_i^2$  can be estimated consistently and the resulting WLSE is asymptotically as efficient as  $\tilde{\theta}$  in (1.3) [Carroll (1982) and Davidian and Carroll (1987)]. However, there are situations where  $\sigma_i^2$  is unrelated to the  $x_{ij}$  and therefore this approach cannot be applied. For example, in a ‘‘common mean’’ problem [Fuller and Rao (1978)],

$$y_{ij} = \mu + e_{ij}.$$

The  $\sigma_i^2$  are unrelated to  $x_{ij}$  since  $x_{ij} = 1$  for all  $i, j$ . More generally, in a ‘‘common regression coefficients’’ problem [Box and Tiao (1973), pages 478–479], where for each  $i$ ,

$$y_{ij} = x'_{ij}\beta + e_{ij}, \quad j = 1, 2, \dots, n_i,$$

is a regression model and we wish to combine the data from  $k$  different communities (batches, days) to improve the accuracy of estimates, the heteroscedasticity in  $\sigma_i^2$  is usually caused by the variation among different communities (batches, days). In these situations estimators of  $\sigma_i^2$  such as the MINQUE and its modification were proposed and studied in the literature [Rao (1970), (1973)]. If  $k$  is fixed and  $\min_{i \leq k} n_i \rightarrow \infty$ , then these estimators are consistent and the resulting WLSE is as efficient as  $\tilde{\theta}$ . We confine our attention to the situation.

$$\sup_i n_i < \infty \quad \text{and} \quad k \text{ is large,}$$

which is particularly true if  $y_{ij}, j = 1, 2, \dots, n_i$ , are replications [Fuller and Rao (1978)]. Rao (1973) proposed a modified MINQUE of  $\sigma_i^2$ , the within-group average of residual squares

$$(1.4) \quad v_i(\hat{\beta}) = \frac{1}{n_i} \sum_{j=1}^{n_i} (y_{ij} - x'_{ij}\hat{\beta})^2,$$

where  $\hat{\beta}$  is an estimator of  $\beta$ . The WLSE constructed by using  $v_i$  is better than the WLSE constructed by using the MINQUE or the within-group sample variance [Rao (1973) and Carroll and Cline (1988)].

Because of the existence of a large number of nuisance parameters  $\sigma_i^2$ , estimators of  $\sigma_i^2$  are not consistent as  $k \rightarrow \infty$  and therefore the corresponding WLSE is not as efficient as the ‘‘estimator’’  $\tilde{\theta}$  in (1.3). Hence it is still possible to improve the WLSE. A natural approach suggested by many researchers is to

use the following iterative procedure:

- (1.5) Obtain an estimate of  $\beta \rightarrow$  obtain estimators of  $\sigma_i^2 \rightarrow$  feedback and repeat.

More precisely, the WLSE after the  $m$ th iteration is

$$(1.6) \quad \hat{\theta}_w^{(m)} = g(\hat{\beta}_w^{(m)}), \quad \hat{\beta}_w^{(m)} = \left( \sum_{i=1}^k \sum_{j=1}^{n_i} w_i^{(m)} x_{ij} x'_{ij} \right)^{-1} \sum_{i=1}^k \sum_{j=1}^{n_i} w_i^{(m)} x_{ij} y_{ij},$$

where

$$w_i^{(m)} = [v_i(\hat{\beta}_w^{(m-1)})]^{-1}, \quad i = 1, 2, \dots, k,$$

for  $m = 1, 2, \dots$ . The OLSE  $\hat{\beta}_o$  can be used as an initial estimator  $\hat{\beta}_w^{(0)}$ , in which case  $\hat{\beta}_w^{(1)}$  is the WLSE proposed in Fuller and Rao (1978).

This raises some interesting questions:

- (a) What is the asymptotic distribution of  $\hat{\theta}_w^{(m)}$  for each fixed  $m$ ?
- (b) Is there a finite integer  $m^*$  such that  $\hat{\theta}_w^{(m^*)}$  is the best among  $\hat{\theta}_w^{(m)}$ ,  $m = 1, 2, \dots$ ?
- (c) If such an  $m^*$  exists but is unknown, when should we stop the iterative process?

These problems are studied in the present paper. In the case where  $\sigma_i^2$  are related to the  $x_{ij}$  or  $x'_{ij}\beta$ , problems similar to (a)–(c) were studied in Carroll, Wu and Ruppert (1988), for the small sample case.

Our main findings are the following.

1. It is shown in Section 2 that under some regularity conditions, the WLSE given by (1.2) with  $w_i = [v_i(\hat{\beta})]^{-1}$  given by (1.4) is asymptotically normal with mean  $\theta$ . As a consequence of this result, the WLSE  $\hat{\theta}_w^{(m)}$  for each fixed  $m$  is asymptotically normal with mean  $\theta$  and its asymptotic variance can be explicitly obtained.
2. It is shown in Section 3 that under some conditions, there exists a finite integer  $m^*$  such that in terms of its asymptotic efficiencies,  $\hat{\theta}_w^{(m^*)}$  is better than  $\hat{\theta}_w^{(m)}$  for any  $m \neq m^*$ . This  $m^*$  depends on  $\sigma_i^2$ ,  $i = 1, 2, \dots, k$ , and therefore is unknown.
3. An adaptive procedure is proposed in Section 3. Three main features of this adaptive procedure are: (i) It ensures that for given data, the iterative process (1.5) stops after a finite number of iterations; (ii) the resulting adaptive estimator of  $\theta$  has an optimality property, that is, it is better than any  $\hat{\beta}_w^{(m)}$  and is as good as the optimal  $\hat{\beta}_w^{(m^*)}$  if it exists and is unique; (iii) an estimate of the asymptotic variance of the adaptive estimator is obtained.
4. If  $n_i \leq 2$  for all  $i$ , it is shown in Section 4 that asymptotically the WLSE cannot improve the OLSE. We consider combining some groups so that each new group contains at least three observations. The adaptive proce-

ture can then be applied to the data with combined groups. The resulting adaptive estimator improves the OLSE in some cases.

5. A simulation study is presented in Section 5. The performance of the adaptive estimator is shown to be adequate in the simulation study.

**2. General asymptotic results.** Denote the  $t \times t$  identity matrix by  $I_t$ . Let

$$\begin{aligned} \mathbf{e}_i &= (e_{i1}, e_{i2}, \dots, e_{in_i})', \\ \mathbf{e} &= (\mathbf{e}'_1, \mathbf{e}'_2, \dots, \mathbf{e}'_k)', \\ u_i &= \mathbf{e}'_i \mathbf{e}_i / n_i, \\ \tilde{\mathbf{e}}_i &= \mathbf{e}_i / u_i, \\ \tilde{\mathbf{e}} &= (\tilde{\mathbf{e}}'_1, \tilde{\mathbf{e}}'_2, \dots, \tilde{\mathbf{e}}'_k)', \\ V_i &= \text{cov}(\mathbf{e}_i), \\ V &= \text{block diag}(V_1, V_2, \dots, V_k), \\ X_i &= (x_{i1}, x_{i2}, \dots, x_{in_i})', \\ X &= (X'_1, X'_2, \dots, X'_k)'. \end{aligned}$$

The following assumptions will be used.

**ASSUMPTION A.** There are positive constants,  $\sigma_o$ ,  $\sigma_\infty$ ,  $c_o$  and  $c_\infty$  and positive integers  $n_o$  and  $n_\infty$  such that for all  $i$  and  $j$ ,

$$\begin{aligned} \sigma_o I_{n_i} &\leq V_i \leq \sigma_\infty I_{n_i}, & c_o I_p &\leq k^{-1}(X'X), \\ n_o &\leq n_i \leq n_\infty, & x'_{ij} x_{ij} &\leq c_\infty. \end{aligned}$$

**ASSUMPTION B.**  $\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_k$  are independent.

**ASSUMPTION C.** There are positive constants  $d$  and  $\delta$  such that for all  $i, j$ ,

$$\begin{aligned} E \left( e_{ij} \left/ \sum_{j=1}^{n_i} e_{ij}^2 \right. \right) &= 0, \\ E |e_{ij}|^{2+\delta} &\leq d, & E \left( \sum_{j=1}^{n_i} e_{ij}^2 \right)^{-(1+\delta)} &\leq d. \end{aligned}$$

The first part of Assumption C reflects a certain degree of symmetry of the error distributions. Without this condition, the WLSE may be inconsistent [Carroll and Cline (1988), Theorem 3].

We first establish a general asymptotic result for a WLSE  $\hat{\beta}_w$ .

THEOREM 1. Suppose that Assumptions A, B and C hold. Let  $\hat{\beta}_w$  be given by (1.2) with  $w_i = [v_i(\hat{\beta})]^{-1}$ ,  $v_i(\hat{\beta})$  given by (1.4) with  $\hat{\beta}$  satisfying

$$(2.1) \quad \hat{\beta} - \beta = G_k^{-1}(A_k \tilde{\mathbf{e}} + B_k \mathbf{e}) + o_p(k^{-1/2}),$$

where  $G_k = \sum_{i=1}^k E u_i^{-1} X_i' X_i$ ,  $A_k$  and  $B_k$  satisfy  $A_k A_k' = O(k)$  and  $B_k B_k' = O(k)$ . Then

$$(2.2) \quad \hat{\beta}_w - \beta = G_k^{-1}(\tilde{A}_k \tilde{\mathbf{e}} + \tilde{B}_k \mathbf{e}) + o_p(k^{-1/2}),$$

where

$$\tilde{A}_k = X' + 2H_k G_k^{-1} A_k,$$

$$\tilde{B}_k = 2H_k G_k^{-1} B_k$$

and

$$H_k = E \left( \sum_{i=1}^k n_i^{-1} X_i' \tilde{\mathbf{e}}_i \tilde{\mathbf{e}}_i' X_i \right).$$

REMARK. It follows from (2.2) that

$$[G_k^{-1} \Sigma_k G_k^{-1}]^{-1/2} (\hat{\beta}_w - \beta) \rightarrow_d N(0, I_p),$$

where  $\Sigma_k = \text{var}(\tilde{A}_k \tilde{\mathbf{e}} + \tilde{B}_k \mathbf{e})$ . Result (2.2) was established in Shao (1989a) for the special case where  $\hat{\beta} = \hat{\beta}_o$  and  $e_{ij}$  are mutually independent.

PROOF OF THEOREM 1. Let  $v_i = v_i(\hat{\beta})$ . Then

$$(2.3) \quad \hat{\beta}_w - \beta = \left( \sum_{i=1}^k v_i^{-1} X_i' X_i \right)^{-1} \sum_{i=1}^k v_i^{-1} X_i' \mathbf{e}_i.$$

Note that

$$(2.4) \quad \begin{aligned} v_i^{-1} &= u_i^{-1} - u_i^{-2}(v_i - u_i) + u_i^{-2}v_i^{-1}(v_i - u_i)^2 \\ &= u_i^{-1} + 2n_i^{-1}u_i^{-2} \sum_{j=1}^{n_i} \phi_{ij} e_{ij} - n_i^{-1}u_i^{-2} \sum_{j=1}^{n_i} \phi_{ij}^2 + u_i^{-2}v_i^{-1}(v_i - u_i)^2, \end{aligned}$$

where  $\phi_{ij} = x'_{ij}(\hat{\beta} - \beta)$ . We now show that

$$(2.5) \quad \sum_{i=1}^k u_i^{-2}v_i^{-1}(v_i - u_i)^2 X_i' \mathbf{e}_i = o_p(k^{1/2})$$

and

$$(2.6) \quad \sum_{i=1}^k n_i^{-1}u_i^{-2} \sum_{j=1}^{n_i} \phi_{ij}^2 X_i' \mathbf{e}_i = o_p(k^{1/2}).$$

For  $(2(1 + \delta))^{-1} < \alpha_1 < \alpha_2 < 1/2$ , let

$$\mathbf{A}_k = \left\{ \min_i \max_j |e_{ij}| \geq k^{-\alpha_1}, \max_{i,j} |x'_{ij}(\hat{\beta} - \beta)| \leq k^{-\alpha_2} \right\}.$$

Let  $c$  be a generic constant. On the event  $\mathbf{A}_k$ ,

$$\begin{aligned} v_i &\geq \max_j |e_{ij} - x'_{ij}(\hat{\beta} - \beta)|^2/n_i \\ &\geq (1 - k^{\alpha_1 - \alpha_2})^2 \max_j e_{ij}^2/n_i \geq (1 - k^{\alpha_1 - \alpha_2})^2 \sum_{j=1}^{n_i} e_{ij}^2/n_i^2 \geq cu_i \end{aligned}$$

and therefore

$$u_i^{-2}v_i^{-1}(v_i - u_i)^2 \leq c(v_i - u_i)^2 u_i^{-3}.$$

Under Assumption A,

$$\begin{aligned} (v_i - u_i)^2 &= n_i^{-2} \left[ \sum_{j=1}^{n_i} \phi_{ij}^2 - 2 \sum_{j=1}^{n_i} \phi_{ij} e_{ij} \right]^2 \\ &\leq \left[ \sum_{j=1}^{n_i} \phi_{ij}^2 + 2 \left( \sum_{j=1}^{n_i} e_{ij}^2 \right)^{1/2} \left( \sum_{j=1}^{n_i} \phi_{ij}^2 \right)^{1/2} \right]^2 \\ &\leq c \left[ \|\hat{\beta} - \beta\|^4 + u_i \|\hat{\beta} - \beta\|^2 \right], \end{aligned}$$

where  $\|\hat{\beta} - \beta\|^2 = (\hat{\beta} - \beta)'(\hat{\beta} - \beta)$ . Consequently, on the event  $\mathbf{A}_k$ ,

$$\begin{aligned} \left| \sum_{i=1}^k u_i^{-2}v_i^{-1}(v_i - u_i)^2 X'_i \mathbf{e}_i \right| &\leq c \left( \|\hat{\beta} - \beta\|^4 \sum_{i=1}^k u_i^{-5/2} + \|\hat{\beta} - \beta\|^2 \sum_{i=1}^k u_i^{-3/2} \right) \\ &= o_p(k^{1/2}), \end{aligned}$$

since  $\|\hat{\beta} - \beta\| = O_p(k^{-1/2})$  and  $\max_i u_i^{-1} \leq c \max_i \min_j e_{ij}^{-2} \leq ck^{2\alpha_1} = o_p(k)$ .

Note that

$$(2.7) \quad P\{\mathbf{A}_k^c\} \leq P\left\{ \max_{i,j} |x'_{ij}(\hat{\beta} - \beta)| \geq k^{-\alpha_2} \right\} + \sum_{i=1}^k P\left\{ \max_j |e_{ij}| \leq k^{-\alpha_1} \right\}.$$

Since  $\|\hat{\beta} - \beta\| = O_p(k^{-1/2})$ , the first term on the right side of (2.7) goes to 0. The second term on the right side of (2.7) also goes to 0, since

$$\begin{aligned} P\left\{ \max_j |e_{ij}| \leq k^{-\alpha_1} \right\} &\leq P\{u_i \leq k^{-2\alpha_1}\} \leq P\{u_i^{-1} \geq k^{2\alpha_1}\} \\ &\leq E u_i^{-(1+\delta)} k^{-2\alpha_1(1+\delta)} \end{aligned}$$

and  $2\alpha_1(1 + \delta) > 1$ . This proves  $P(\mathbf{A}_k) \rightarrow 1$ , and hence (2.5) holds. The proof of (2.6) is similar.

From (2.4)–(2.6),

$$(2.8) \quad \sum_{i=1}^k v_i^{-1} X'_i \mathbf{e}_i = \sum_{i=1}^k u_i^{-1} X'_i e_i + 2 \left[ \sum_{i=1}^k n_i^{-1} X'_i \tilde{\mathbf{e}}_i \tilde{\mathbf{e}}'_i X_i \right] (\hat{\beta} - \beta) + o_p(k^{1/2}).$$

Following the argument in the proof of (2.5), we obtain that

$$k^{-1} \left( \sum_{i=1}^k v_i^{-1} X_i' X_i - G_k \right) \rightarrow_p 0,$$

which, together with Assumption A, implies that

$$k \left[ \left( \sum_{i=1}^k v_i^{-1} X_i' X_i \right)^{-1} - G_k^{-1} \right] \rightarrow_p 0.$$

Also, it is easy to see that

$$k^{-1} \left( \sum_{i=1}^k n_i^{-1} X_i' \tilde{\mathbf{e}}_i \tilde{\mathbf{e}}_i' X_i - H_k \right) \rightarrow_p 0.$$

Hence the result follows from (2.3) and (2.8).  $\square$

Let  $\{A_k^{(0)}\}$  and  $\{B_k^{(0)}\}$  be the two sequences of  $p \times (\sum_{i=1}^k n_i)$  matrices satisfying  $(A_k^{(0)})(A_k^{(0)})' = O(k)$  and  $(B_k^{(0)})(B_k^{(0)})' = O(k)$ . For any positive integer  $m$ , let

$$A_k^{(m)} = \sum_{t=0}^{m-1} (2H_k G_k^{-1})^t X' + (2H_k G_k^{-1})^m A_k^{(0)}$$

and

$$B_k^{(m)} = (2H_k G_k^{-1})^m B_k^{(0)},$$

where  $G_k$  and  $H_k$  are given in Theorem 1. Applying the  $\delta$ -method and Theorem 1 repeatedly with  $A_k = A_k^{(t-1)}$ ,  $B_k = B_k^{(t-1)}$ ,  $\hat{\beta} = \hat{\beta}_w^{(t-1)}$  and  $\hat{\beta}_w = \hat{\beta}_w^{(t)}$ ,  $t = 1, 2, \dots, m$ , we obtain the following result for the WLSE  $\hat{\theta}_w^{(m)}$  given in (1.6).

**THEOREM 2.** *Assume the conditions in Theorem 1 and that (2.1) holds with  $\hat{\beta} = \hat{\beta}_w^{(0)}$ . Assume also that  $g$  is differentiable at  $\beta$ . Then for any fixed  $m$ ,*

$$(2.9) \quad \left\{ \nabla g(\beta) G_k^{-1} \Sigma_k^{(m)} G_k^{-1} [\nabla g(\beta)]' \right\}^{-1/2} (\hat{\theta}_w^{(m)} - \theta) \rightarrow_d N(0, 1),$$

where  $\Sigma_k^{(m)} = \text{Var}(A_k^{(m)} \tilde{\mathbf{e}} + B_k^{(m)} \mathbf{e})$ .

We now discuss some important special cases.

1. When  $\hat{\beta}_w^{(0)} = \hat{\beta}_o$  (the OLSE), (2.1) is satisfied since

$$A_k^{(0)} = 0 \quad \text{and} \quad B_k^{(0)} = G_k (X'X)^{-1} X'.$$

Then for  $m \geq 1$ ,

$$A_k^{(m)} = \sum_{t=0}^{m-1} (2H_k G_k^{-1})^t X',$$

$$B_k^{(m)} = (2H_k G_k^{-1})^m G_k (X'X)^{-1} X'.$$

2. If  $e_{ij}/\sigma_i$ , for  $j = 1, 2, \dots, n_i; i = 1, 2, \dots, k$ , are independent and identically distributed (i.i.d.),  $n_i = n_o$  for all  $i$  and if for  $j \neq j'$ ,

$$(2.10) \quad E(e_{ij}e_{ij'})/u_i^2 = 0,$$

then  $H_k = n_o^{-1}G_k$ ,  $G_k = \rho(n_o)X'V^{-1}X$ , where  $\rho(n_o) = \sigma_i^2 E u_i^{-1}$ ,

$$\begin{aligned} A_k^{(m)} &= \sum_{t=0}^{m-1} \left(\frac{2}{n_o}\right)^t X' + \left(\frac{2}{n_o}\right)^m A_k^{(0)} \\ &= \left(1 - \frac{2}{n_o}\right)^{-1} \left[1 - \left(\frac{2}{n_o}\right)^m\right] X' + \left(\frac{2}{n_o}\right)^m A_k^{(0)} \end{aligned}$$

and

$$B_k^{(m)} = \left(\frac{2}{n_o}\right)^m B_k^{(0)}.$$

3. If  $e_{ij}/\sigma_i$  are i.i.d. and for  $j \neq j'$ ,

$$(2.11) \quad E e_{ij} e_{ij'} / u_i = 0,$$

then

$$\Sigma_k^{(m)} = A_k^{(m)} E(\tilde{\mathbf{e}}\tilde{\mathbf{e}}') (A_k^{(m)})' + A_k^{(m)} (B_k^{(m)})' + B_k^{(m)} (A_k^{(m)})' + B_k^{(m)} V (B_k^{(m)})'.$$

Note that (2.10) and (2.11) hold if  $e_{ij}/\sigma_i$  are i.i.d. with a symmetric distribution. Combining (1)–(3), we have the following result. The proof is straightforward and omitted.

**THEOREM 3.** *Suppose that the conditions in Theorem 2 and (2.10) and (2.11) hold. Assume the  $e_{ij}/\sigma_i$  are i.i.d.,  $n_i = n_o$  for all  $i$  and  $\hat{\beta}_w^{(0)} = \hat{\beta}_o$ . Then for any  $m$ ,*

$$(\tau_k^{(m)}/k)^{-1/2} (\hat{\theta}_w^{(m)} - \theta) \rightarrow_d N(0, 1),$$

with

$$\begin{aligned} \tau_k^{(m)} &= k \left[ \frac{\varphi(m)}{\rho(n_o)} \nabla g(\beta) (X'V^{-1}X)^{-1} \nabla g(\beta)' \right. \\ &\quad \left. + \psi(m) \nabla g(\beta) (X'X)^{-1} (X'VX) (X'X)^{-1} \nabla g(\beta)' \right], \end{aligned}$$

where

$$\varphi(m) = \left(1 - \frac{2}{n_o}\right)^{-2} \left(1 - \left(\frac{2}{n_o}\right)^m\right)^2 + 2 \left(1 - \frac{2}{n_o}\right)^{-1} \left(1 - \left(\frac{2}{n_o}\right)^m\right) \left(\frac{2}{n_o}\right)^m$$



and

$$\psi(m) = \left(\frac{2}{n_o}\right)^{2m}.$$

**3. The efficiency of the WLSE and an adaptive estimator.** We now study the asymptotic efficiency of  $\hat{\theta}_w^{(m)}$ . For simplicity, we confine our attention to the situation where  $n_i = n_o$  for all  $i$  and  $e_{ij}/\sigma_i$  are i.i.d. Let

$$b_k = k \left[ \nabla g(\beta) (X'X)^{-1} (X'VX) (X'X)^{-1} \nabla g(\beta)' \right]$$

and

$$\tilde{b}_k = k \left[ \nabla g(\beta) (X'V^{-1}X)^{-1} \nabla g(\beta)' \right].$$

Note that  $b_k/k$  and  $\tilde{b}_k/k$  are the asymptotic variances of  $\hat{\theta}_o$  and  $\tilde{\theta}$ , respectively. From Theorem 3, the asymptotic variance of  $\hat{\theta}_w^{(m)}$  for a fixed  $m$  is

$$(3.1) \quad \tau_k^{(m)}/k = \left[ \frac{\varphi(m)}{\rho(n_o)} \tilde{b}_k + \psi(m) b_k \right] / k,$$

where  $\varphi(m)$  and  $\psi(m)$  are given in Theorem 3:

LEMMA 1. When  $n_o \geq 3$ ,

- (i)  $\varphi(m)$  is strictly increasing in  $m$ .
- (ii) The function

$$\Delta(m) = \frac{\varphi(m+1) - \varphi(m)}{\psi(m) - \psi(m+1)}$$

is strictly increasing in  $m$  and  $\Delta(m) \rightarrow \infty$  as  $m \rightarrow \infty$ .

PROOF. (i) From the expression of  $\varphi(m)$ , we need only consider the function

$$\begin{aligned} & \left[ \left( 1 - \left( \frac{2}{n_o} \right)^m \right)^2 + 2 \left( 1 - \frac{2}{n_o} \right) \left( 1 - \left( \frac{2}{n_o} \right)^m \right) \left( \frac{2}{n_o} \right)^m \right] \\ & = 1 - \frac{4}{n_o} x_m + \left( \frac{4}{n_o} - 1 \right) x_m^2, \end{aligned}$$

where  $x_m = (2/n_o)^m$ . Since

$$\frac{d}{dx} \left[ 1 - \frac{4}{n_o} x + \left( \frac{4}{n_o} - 1 \right) x^2 \right] = 2 \left( \frac{4}{n_o} - 1 \right) x - \frac{4}{n_o} < 0,$$

when  $x \leq 2/3$ , the function  $1 - (4/n_o)x + ((4/n_o) - 1)x^2$  is strictly decreasing in  $x$ . This implies  $\varphi(m)$  is strictly increasing in  $m$ .

(ii) Similarly, with  $x_m = (2/n_o)^m$ ,

$$\Delta(m) = \left(1 - \frac{2}{n_o}\right)^{-2} \left[ \frac{4(x_m - x_{m+1})}{n_o(x_m^2 - x_{m+1}^2)} + \left(1 - \frac{4}{n_o}\right) \right],$$

which is strictly increasing in  $m$ . Finally,

$$\lim_{m \rightarrow \infty} \frac{x_m - x_{m+1}}{(x_m^2 - x_{m+1}^2)} = \lim_{m \rightarrow \infty} \frac{1}{x_m + x_{m+1}} = \infty. \quad \square$$

Using this lemma, we can compare the asymptotic efficiencies of  $\hat{\theta}_w^{(m)}$ ,  $m = 0, 1, 2, \dots$ . Throughout this section we assume  $n_o \geq 3$ . The case of  $n_o = 2$  is treated in Section 4.

**THEOREM 4.** *Suppose that  $b = \lim_{k \rightarrow \infty} b_k$  and  $\tilde{b} = \lim_{k \rightarrow \infty} \tilde{b}_k$  exist. Then there exists an integer  $m^*$  such that*

$$\lim_{k \rightarrow \infty} \tau_k^{(m^*)} / \tau_k^{(m)} \leq 1$$

for all  $m$  with equality holding at most for  $m = m^*$  or  $m^* + 1$ .

**PROOF.** Note that

$$\tau_k^{(m+1)} - \tau_k^{(m)} = \frac{\psi(m) - \psi(m+1)}{\rho(n_o)} \tilde{b}_k \left[ \Delta(m) - \frac{\rho(n_o)b_k}{\tilde{b}_k} \right].$$

From Assumption A,  $0 < b < \infty$  and  $0 < \tilde{b} < \infty$ . Since  $\psi(m) - \psi(m+1) > 0$ , whether or not  $\lim_{k \rightarrow \infty} \tau_k^{(m+1)} / \tau_k^{(m)} \leq 1$  depends on the sign of

$$t(m) = \Delta(m) - \frac{\rho(n_o)b}{\tilde{b}}.$$

From Lemma 1,  $t(m)$  can change sign only once as  $m$  increases. Thus, as  $m$  increases, the quantity

$$\lim_{k \rightarrow \infty} (\tau_k^{(m+1)} / \tau_k^{(m)} - 1)$$

will change sign only once. Since  $\Delta(m) \rightarrow \infty$  as  $m \rightarrow \infty$  (Lemma 1), there always exists a unique  $m^*$  such that

$$\lim_{k \rightarrow \infty} \frac{\tau_k^{(m^*)}}{\tau_k^{(m)}} \begin{cases} < 1 & \text{if } m < m^* \text{ or } m > m^* + 1 \\ \leq 1 & \text{if } m = m^* + 1. \end{cases}$$

This proves the theorem.  $\square$

From Theorem 4, in terms of asymptotic variance,  $\hat{\theta}_w^{(m^*)}$  is the optimal estimator of  $\theta$  among  $\hat{\theta}_w^{(m)}$ ,  $m = 1, 2, \dots$ , and  $m^*$  is the optimal number of iterations. Note that  $m^*$  can be as small as 0.  $\hat{\theta}_w^{(m^*+1)}$  is also optimal if  $\lim_{k \rightarrow \infty} (\tau_k^{(m^*+1)} / \tau_k^{(m^*)}) = 1$ , which does not occur frequently in practice.

In general,  $b_k/\tilde{b}_k$  measures the degree of heteroscedasticity of the model (1.1), that is,  $b_k/\tilde{b}_k$  is large when  $\sigma_i$  are very different and  $b_k/\tilde{b}_k$  is close to 1 if  $\sigma_i$  are nearly the same ( $b_k/\tilde{b}_k \geq 1$  by Jensen's inequality). Since  $\Delta(m)$  is increasing as  $m$  increases, our result shows that the more different the  $\sigma_i$  are, the more iterations we need.

The optimal  $m^*$  is generally unknown and therefore we need to estimate it using the data. Let

$$a_k = \tilde{b}_k/\rho(n_o).$$

Then, by (3.1),

$$\tau_k^{(m)} = \varphi(m)a_k + \psi(m)b_k,$$

where  $\varphi(m)$  and  $\psi(m)$  are known and  $a_k$  and  $b_k$  are independent of  $m$ . Suppose that  $\hat{a}_k$  and  $\hat{b}_k$  are consistent estimators of  $a_k$  and  $b_k$ , respectively. That is,

$$(3.2) \quad \hat{a}_k - a_k \rightarrow_p 0 \quad \text{and} \quad \hat{b}_k - b_k \rightarrow_p 0.$$

Then  $\hat{\tau}_k^{(m)} = \varphi(m)\hat{a}_k + \psi(m)\hat{b}_k$  is consistent for  $\tau_k^{(m)}$  for each fixed  $m$ .

Let  $\hat{r}_k = \hat{b}_k/\hat{a}_k$  and  $r_k = b_k/a_k$ . Then, by (3.2),

$$(3.3) \quad \hat{r}_k - r_k \rightarrow_p 0.$$

Define

$$\hat{m}_k = \min\{m: \Delta(m) \geq \hat{r}_k, m = 0, 1, 2, \dots\}.$$

**THEOREM 5.** Assume that (3.3) holds and  $\lim_{k \rightarrow \infty} r_k = r$ . Let  $m^*$  be as given in Theorem 4. Then

$$(3.4) \quad \lim_{k \rightarrow \infty} P(\hat{m}_k = m^* \text{ or } \hat{m}_k = m^* + 1) = 1.$$

Furthermore, if  $m^*$  is the unique optimum ( $\lim_{k \rightarrow \infty} \tau_k^{(m^*+1)}/\tau_k^{(m^*)} > 1$ ), then

$$(3.5) \quad \lim_{k \rightarrow \infty} P(\hat{m}_k = m^*) = 1.$$

**PROOF.** From Lemma 1 and the proof of Theorem 4,

$$\Delta(m^* - 1) < r \quad \text{and} \quad \Delta(m^* + 1) > r.$$

Then

$$P(\hat{m}_k \leq m^* - 1) \leq P(\Delta(m^* - 1) \geq \hat{r}_k) \rightarrow 0,$$

since  $P(r - \varepsilon \geq \hat{r}_k) \rightarrow 0$  for  $\varepsilon = r - \Delta(m^* - 1)$ .

Similarly,

$$P(\hat{m}_k \geq m^* + 2) \leq P(\Delta(m^* + 1) < \hat{r}_k) \rightarrow 0.$$

This proves (3.4).

Now assume  $m^*$  is the unique optimum. Then  $\Delta(m^*) > r$  and

$$P(\hat{m}_k = m^* + 1) \leq P(\Delta(m^*) < \hat{r}_k) \rightarrow 0,$$

which with (3.4) implies (3.5).  $\square$

Result (3.5) indicates that if we stop the iterative procedure (1.5) after  $\hat{m}_k$  iterations, we always stop at the right time when  $k$  is large enough. This leads to the following adaptive estimator of  $\theta$ :

$$\hat{\theta}_\alpha = \hat{\theta}_w^{(\hat{m}_k)}.$$

This adaptive estimator has the desired optimality property, provided that the estimators  $\hat{a}_k$  and  $\hat{b}_k$  satisfying (3.2) can be found. Furthermore, an estimator of the accuracy of  $\hat{\theta}_\alpha$  is

$$\hat{\tau}_k/k,$$

with

$$\hat{\tau}_k = \varphi(\hat{m}_k)\hat{a}_k + \psi(\hat{m}_k)\hat{b}_k.$$

**THEOREM 6.** *Assume that (3.2) holds.*

(i) *As  $k \rightarrow \infty$ ,*

$$P(\hat{\theta}_\alpha = \hat{\theta}_w^{(m^*)} \text{ or } \hat{\theta}_\alpha = \hat{\theta}_w^{(m^*+1)}) \rightarrow 1.$$

*If in addition,  $m^*$  is the unique optimum, then*

$$P(\hat{\theta}_\alpha = \hat{\theta}_w^{(m^*)}) \rightarrow 1.$$

(ii) *Let  $\tau^* = \lim_{k \rightarrow \infty} \tau_k^{(m^*)}$ . Then*

$$\hat{\tau}_k - \tau^* \rightarrow_p 0.$$

**PROOF.** The results follow directly from (3.4) and (3.5).  $\square$

In practice when there is a computational limitation, one may force the iterative process (1.5) to stop after at most  $M$  steps, where  $M$  is an integer. That is,

$$\begin{aligned} \tilde{m}_k &= \min\{m: \Delta(m) \geq \hat{r}_k, m = 0, 1, 2, \dots, M\} \\ &= M \text{ if } \Delta(M) \leq \hat{r}_k \end{aligned}$$

and

$$(3.6) \quad \tilde{\theta}_\alpha = \hat{\theta}_w^{(\tilde{m}_k)}.$$

From the previous proofs,  $\tilde{\theta}_\alpha$  is the best among  $\hat{\theta}_w^{(m)}$ ,  $m = 0, 1, 2, \dots, M$ .

Estimators  $\hat{a}_k$  and  $\hat{b}_k$  satisfying (3.2) are needed to carry out this adaptive procedure. We suggest the use of the following consistent estimators:

$$(3.7) \quad \hat{a}_k = k \sum_{i=1}^k \sum_{j=1}^{n_i} (1 - h_{ij})(\hat{\theta}_w^{(i,j)} - \hat{\theta}_w^{(1)})^2,$$

$$(3.8) \quad \hat{b}_k = k \sum_{i=1}^k \sum_{j=1}^{n_i} (1 - c_{ij})(\hat{\theta}_o^{(i,j)} - \hat{\theta}_o)^2,$$

where

$$h_{ij} = n_i^{-1} v_i^{-1} x'_{ij} (\sum_{i=1}^k v_i^{-1} X'_i X_i)^{-1} x_{ij}, \quad v_i = v_i(\hat{\beta}_o),$$

$$c_{ij} = n_i^{-1} x'_{ij} (X'X)^{-1} x_{ij},$$

$$\hat{\theta}_w^{(i,j)} = g(\hat{\beta}_w^{(i,j)}),$$

$$\hat{\beta}_w^{(i,j)} = \hat{\beta}_w^{(1)} - (1 - h_{ij})^{-1} v_i^{-1} \left( \sum_{i=1}^k v_i^{-1} X'_i X_i \right)^{-1} x_{ij} (y_{ij} - x'_{ij} \hat{\beta}_w^{(1)})$$

and

$$\hat{\theta}_o^{(i,j)} = g(\hat{\beta}_o^{(i,j)}), \quad \hat{\beta}_o^{(i,j)} = \hat{\beta}_o - (1 - c_{ij})^{-1} (X'X)^{-1} x_{ij} (y_{ij} - x'_{ij} \hat{\beta}_o).$$

Note that  $\hat{b}_k$  is a consistent (weighted) jackknife estimator of  $b_k$ , the asymptotic variance of  $k^{1/2}(\hat{\theta}_o - \theta)$  [Shao (1989b)]. Also,  $\hat{a}_k$  is asymptotically equivalent to

$$k \left[ \nabla g(\beta) \left( \sum_{i=1}^k u_i^{-1} X'_i X_i \right)^{-1} \nabla g(\beta)' \right]$$

[Shao (1989b)] and hence  $\hat{a}_k$  is consistent for  $\bar{b}_k/\rho(n_o) = a_k$  because of  $Eu_i^{-1} = \sigma_i^{-2}\rho(n_o)$ .

In some cases we need to consider simultaneous estimation of  $\theta = g(\beta)$ , where  $g$  is a  $q$ -dimensional vector-valued function  $\beta$ ,  $1 \leq q \leq p$ . For example, we may need to set a joint confidence region for several functions of  $\beta$ . When  $g$  is a vector-valued function, we may need different numbers of iterations for different components of  $g$ , which is undesirable unless we can reduce the problem to several one-dimensional problems (e.g., applying the Bonferroni method, we may obtain a joint confidence region for  $\theta$  by taking the product of confidence intervals for components of  $\theta$ ). In the following we discuss some extensions of the one-dimensional results in Theorems 4–6 to the case of vector-valued  $g$ .

We shall use the same notation  $\hat{\theta}_w^{(m)} = g(\hat{\beta}_w^{(m)})$ ,  $\nabla g$ ,  $a_k$ ,  $b_k$  and  $\tau_k^{(m)}$ , but now  $g$  and  $\hat{\theta}_w^{(m)}$  are  $q$ -vectors,  $\nabla g$  is a  $q \times p$  matrix and  $a_k$ ,  $b_k$  and  $\tau_k^{(m)}$  are  $q \times q$  matrices. Using the delta method, one can show that the result in Theorem 3 still holds [with  $N(0, 1)$  replaced by  $N(0, I_q)$ ]. Since  $\tau_k^{(m)}$  is a matrix, we need a scalar characteristic of  $\tau_k^{(m)}$  to measure the efficiency of  $\hat{\theta}_w^{(m)}$ . For example, we may consider a quadratic loss function used frequently in simultaneous estimation:

$$L(\theta, d) = (d - \theta)' Q(\theta)(d - \theta),$$

where  $Q(\theta)$  is a  $q \times q$  positive definite matrix function continuous at the true parameter  $\theta$ . The asymptotic risk function of  $\hat{\theta}_w^{(m)}$  is

$$\begin{aligned} E_A L(\theta, \hat{\theta}_w^{(m)}) &= E_A (\hat{\theta}_w^{(m)} - \theta)' Q(\theta) (\hat{\theta}_w^{(m)} - \theta) = \text{trace}[Q(\theta) \tau_k^{(m)} / k] \\ &= \{\varphi(m) \text{trace}[Q(\theta) a_k] + \psi(m) \text{trace}[Q(\theta) b_k]\} / k, \end{aligned}$$

where  $E_A$  is the expectation under the asymptotic distribution of  $\hat{\theta}_w^{(m)} - \theta$ . The scalars  $\text{trace}[Q(\theta)a_k]$  and  $\text{trace}[Q(\theta)b_k]$  can be consistently estimated by  $\text{trace}[Q(\hat{\theta}_o)\hat{a}_k]$  and  $\text{trace}[Q(\hat{\theta}_o)\hat{b}_k]$ , where  $\hat{a}_k$  and  $\hat{b}_k$  are still given by (3.7) and (3.8) (with the square replaced by the vector product). Then the results in Theorems 4–6 can be directly extended to this case. For example, we estimate the optimal number of iterations by

$$\hat{n}_k = \min \left\{ m : \Delta(m) \geq \frac{\text{trace}[Q(\hat{\theta}_o)\hat{b}_k]}{\text{trace}[Q(\hat{\theta}_o)\hat{a}_k]}, m = 0, 1, 2, \dots \right\}.$$

Another widely used measure of efficiency is  $|\tau_k^{(m)}|$ , the determinant of  $\tau_k^{(m)}$ , which is also called the generalized asymptotic variance of  $\hat{\theta}_w^{(m)}$ . We prove some results similar to those in Theorems 4–6.

**THEOREM 7.** *Suppose that  $a_k \rightarrow a > 0$  (positive definite) and  $b_k \rightarrow b > 0$ .*

(i) *There exists an integer  $m^*$  such that  $\lim_{k \rightarrow \infty} (|\tau_k^{(m^*)}|/|\tau_k^{(m)}|) \leq 1$  for all  $m$  with equality holding for at most finitely many  $m$ 's.*

(ii) *Let  $\hat{\tau}_k^{(m)} = \varphi(m)\hat{a}_k + \psi(m)\hat{b}_k$  and  $\hat{n}_k$  be a minimizer of  $|\hat{\tau}_k^{(m)}|$ , that is,  $|\hat{\tau}_k^{(\hat{n}_k)}| = \min_m |\hat{\tau}_k^{(m)}|$ . Then*

$$(3.9) \quad \lim_{k \rightarrow \infty} P(\hat{n}_k \in \mathbf{M}) = 1,$$

where  $\mathbf{M} = \{m : \lim_{k \rightarrow \infty} (|\tau_k^{(m^*)}|/|\tau_k^{(m)}|) = 1\}$ .

**PROOF.** (i) Since  $a > 0$  and  $b > 0$ , there exists a nonsingular matrix  $\Gamma$  such that

$$\Gamma'[\varphi(m)a + \psi(m)b]\Gamma = \varphi(m)I_q + \psi(m)\Lambda,$$

where  $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_q)$  with  $0 < \lambda_1 \leq \dots \leq \lambda_q$ . Then

$$\lim_{k \rightarrow \infty} |\tau_k^{(m)}| = |\Gamma|^2 |\varphi(m)I_q + \psi(m)\Lambda| = |\Gamma|^2 \prod_{t=1}^q [\varphi(m) + \psi(m)\lambda_t].$$

Note that

$$\begin{aligned} & [\varphi(m+1) + \psi(m+1)\lambda_t] - [\varphi(m) + \psi(m)\lambda_t] \\ &= [\psi(m) - \psi(m+1)][\Delta(m) - \lambda_t] \end{aligned}$$

and by Lemma 1,  $\Delta(m) - \lambda_t$  can change sign only once as  $m$  increases. Thus

$$\frac{\varphi(m+1) + \psi(m+1)\lambda_t}{\varphi(m) + \psi(m)\lambda_t} \begin{cases} < 1 & \text{if } m < m(t) \\ > 1 & \text{if } m > m(t) + 1 \end{cases}$$

for an integer  $m(t)$ ,  $t = 1, \dots, q$ . Since  $\lambda_1 \leq \dots \leq \lambda_q$ ,  $m(1) \leq \dots \leq m(q)$ . Hence

$$\lim_{k \rightarrow \infty} \frac{|\tau_k^{(m+1)}|}{|\tau_k^{(m)}|} \begin{cases} < 1 & \text{if } m < m(1) \\ > 1 & \text{if } m > m(q) + 1. \end{cases}$$

This proves that  $\{\lim_{k \rightarrow \infty} |\tau_k^{(m)}|, m = 0, 1, 2, \dots\}$  has a minimum at some  $m^*$  satisfying  $m(1) \leq m^* \leq m(q) + 1$ . Hence the result in (i) holds.

(ii) Since  $\varphi(m)$  and  $\psi(m)$  are bounded,

$$\hat{\tau}_k^{(m)} \rightarrow_p \varphi(m)a + \psi(m)b \quad \text{uniformly in } m.$$

This implies (3.9).  $\square$

Thus, in terms of the generalized asymptotic variance,  $\hat{n}_k$  in (3.9) can be used to stop the iterative procedure.

**4. Asymptotics for the case of  $n_i \leq 2$ .** One of the conditions required for the asymptotic results in the previous sections is that for any  $i$ ,

$$(4.1) \quad E \left( \sum_{j=1}^{n_i} e_{ij}^2 \right)^{-(1+\delta)} \leq d$$

for some constants  $\delta$  and  $d$ . This condition is satisfied for most error distributions if  $n_i \geq 3$  [see Proposition 4.1 in Shao (1989a)]. Often, when  $n_i \leq 2$ ,  $Eu_i^{-1} = \infty$  and hence (4.1) does not hold. Furthermore, Lemma 1 in Section 3 requires  $n_o \geq 3$ . Hence the results in the previous sections cannot be applied when  $n_i \leq 2$ . The problem of having  $Eu_i^{-1} = \infty$  can be avoided by considering

$$\tilde{u}_i = \frac{1}{n_i} \sum_{j=1}^{n_i} e_{ij}^2 + k^{-1}$$

and

$$\tilde{v}_i = \tilde{v}_i(\hat{\beta}) = \frac{1}{n_i} \sum_{j=1}^{n_i} (y_{ij} - x'_{ij}\hat{\beta})^2 + k^{-1},$$

since  $E\tilde{u}_i^{-1}$  exists for any  $n_i$ . However, the WLSE with weights  $w_i = \tilde{v}_i^{-1}$  is not asymptotically better than the OLSE. More precisely, we have the following result.

**THEOREM 8.** *Suppose that Assumption A holds and  $n_i = n_o = 2$  for all  $i$ . Suppose also that the  $e_{ij}/\sigma_i$  are i.i.d. with a symmetric density which is positive and continuous at 0. Let  $\hat{\beta}_w$  be given by (1.2) with  $w_i = \tilde{v}_i^{-1}$  and let  $\hat{\beta} = \hat{\beta}_o$  be the OLSE. Then*

$$(4.2) \quad \hat{\beta}_w - \beta = \hat{\beta}_o - \beta + o_p(k^{-1/2}).$$

**PROOF.** Under the symmetry condition, when  $i_1 \neq i_2$ ,  $\mathbf{e}'_{i_1} X_{i_1} X'_{i_2} \mathbf{e}_{i_2} / \tilde{u}_{i_1} \tilde{u}_{i_2}$  has zero expectation. Let  $f(x)$  be the density function of  $e_{ij}/\sigma_i$ . Since  $f$  is positive and continuous at 0,

$$E\tilde{u}_i^{-1} \sim \int \frac{cf(x)f(y)}{k^{-1} + x^2 + y^2} dx dy \sim \int_0^1 \int_0^1 \frac{ck}{1 + kx^2 + ky^2} dx dy = O(\log k).$$

Therefore,

$$E \left( \sum_{i=1}^k \tilde{u}_i^{-1} X_i' \mathbf{e}_i \right)^2 = \sum_{i=1}^k E \mathbf{e}_i' X_i X_i' \mathbf{e}_i / \tilde{u}_i^2 = O(kE(e_{11}^2 / \tilde{u}_1^2)) = O(k \log k),$$

where the second equality follows from the boundedness of  $X_i$  and  $\sigma_i$ . Consequently,

$$(4.3) \quad \sum_{i=1}^k \tilde{u}_i^{-1} X_i' \mathbf{e}_i = O_p((k \log k)^{1/2}) = o_p(k^{1/2} \log k)$$

and

$$(4.4) \quad \sum_{i=1}^k \tilde{v}_i^{-1} X_i' \mathbf{e}_i = \sum_{i=1}^k (\tilde{v}_i^{-1} - \tilde{u}_i^{-1}) X_i' \mathbf{e}_i + o_p(k^{1/2} \log k).$$

Note that

$$(4.5) \quad \begin{aligned} \sum_{i=1}^k (\tilde{v}_i^{-1} - \tilde{u}_i^{-1}) X_i' \mathbf{e}_i &= \sum_{i=1}^k (\tilde{u}_i \tilde{v}_i)^{-1} X_i' \mathbf{e}_i \mathbf{e}_i' X_i (\hat{\beta} - \beta) \\ &\quad - \frac{1}{2} \sum_{i=1}^k (\tilde{v}_i \tilde{u}_i)^{-1} X_i' \mathbf{e}_i \| X_i (\hat{\beta} - \beta) \|^2. \end{aligned}$$

We show that the second term on the right side of (4.5) is negligible. Let  $\alpha$  be the set of the  $i$ 's satisfying

$$(4.6) \quad \max\{|e_{i1}|, |e_{i2}|\} \geq k^{-1/2}(\log k)^{1/4}.$$

Since  $\max_{i,j} [x'_{ij}(\hat{\beta} - \beta)]^2 = O_p(k^{-1})$ , we may focus on the event of  $\max_{i,j} [x'_{ij}(\hat{\beta} - \beta)]^2 \leq k^{-1}(\log k)^{1/2}/4$ . Then for each  $i \in \alpha$ ,  $\max_{i,j} [x'_{ij}(\hat{\beta} - \beta)]^2 \leq \max_j e_{ij}^2/4$  and

$$\begin{aligned} \tilde{v}_i &\geq k^{-1} + \max_j (y_{ij} - x'_{ij} \hat{\beta})^2 / 2 \\ &\geq k^{-1} + \left\{ \max_j e_{ij}^2 / 2 - \max_{i,j} [x'_{ij}(\hat{\beta} - \beta)]^2 \right\} / 2 \\ &\geq k^{-1} + \max_j e_{ij}^2 / 8 \geq \tilde{u}_i / 16. \end{aligned}$$

Also

$$\begin{aligned} E \tilde{u}_i^{-3/2} &\sim \int \frac{cf(x)f(y)}{(k^{-1} + x^2 + y^2)^{3/2}} dx dy \\ &\sim \int_0^1 \int_0^1 \frac{c\sqrt{k}}{(1 + kx^2 + ky^2)^{3/2}} d(\sqrt{k}x) d(\sqrt{k}y) = O(k^{1/2}). \end{aligned}$$

Hence

$$(4.7) \quad \| X_i(\hat{\beta} - \beta) \|^2 \times \| X_i' \mathbf{e}_i (\tilde{u}_i \tilde{v}_i)^{-1} \| \leq O_p(k^{-1}) \tilde{u}_i^{-3/2} = O_p(k^{-1/2})$$



uniformly for  $i \in \alpha$ . Let  $\alpha^c$  be the complement of  $\alpha$ . Since  $[x'_{ij}(\hat{\beta} - \beta)]^2 \leq 2(y_{ij} - x'_{ij}\hat{\beta})^2 + 2e_{ij}^2$ ,

$$(4.8) \quad \begin{aligned} \left\| X_i(\hat{\beta} - \beta) \right\|^2 \times \left\| X'_i \mathbf{e}_i (\tilde{u}_i \tilde{v}_i)^{-1} \right\| &= O_p(\|\mathbf{e}_i\|/\tilde{u}_i + \|\mathbf{e}_i\|/\tilde{v}_i) \\ &\leq O_p(k\|\mathbf{e}_i\|) = O_p(k^{1/2}(\log k)^{1/4}) \end{aligned}$$

for  $i \in \alpha^c$ . Let  $N$  be the number of elements in  $\alpha^c$ . Then  $N$  is a random variable with binomial distribution having parameters  $(k, p_k)$ . The density of  $e_{ij}/\sigma_i$  being continuous at 0 implies that

$$p_k = P\{|e_{i1}| < k^{-1/2}(\log k)^{1/4}, |e_{i2}| < k^{-1/2}(\log k)^{1/4}\} = O(k^{-1}(\log k)^{1/2}).$$

Therefore,  $N = O_p((\log k)^{1/2})$  and

$$(4.9) \quad \begin{aligned} &\sum_{i=1}^k (\tilde{u}_i \tilde{v}_i)^{-1} X'_i \mathbf{e}_i \left\| X_i(\hat{\beta} - \beta) \right\|^2 \\ &= \sum_{i \in \alpha} O_p(k^{-1/2}) + \sum_{i \in \alpha^c} O_p(k^{1/2}(\log k)^{1/4}) \\ &= O_p(k^{1/2}) + O_p(k^{1/2}(\log k)^{3/4}) = o_p(k^{1/2}(\log k)). \end{aligned}$$

From (4.3), (4.4), (4.5) and (4.9), we get

$$(4.10) \quad \sum_{i=1}^k \tilde{v}_i^{-1} X'_i \mathbf{e}_i = \sum_{i=1}^k (\tilde{u}_i \tilde{v}_i)^{-1} X'_i \mathbf{e}_i \mathbf{e}'_i X_i(\hat{\beta} - \beta) + o_p(k^{1/2} \log k).$$

Similarly, for  $i \in \alpha$ ,

$$\begin{aligned} \left\| \sum_{i \in \alpha} (\tilde{u}_i^{-1} - \tilde{v}_i^{-1}) X'_i X_i \right\| &\leq \sum_{i \in \alpha} \left[ \left\| X_i(\beta - \hat{\beta}) \right\|^2 + 2|\mathbf{e}'_i X_i(\hat{\beta} - \beta)| X'_i X_i \right] / \tilde{u}_i^2 \\ &\leq c \|\hat{\beta} - \beta\|^2 \sum_{i \in \alpha} \tilde{u}_i^{-2} + c \|\hat{\beta} - \beta\| \sum_{i \in \alpha} \tilde{u}_i^{-3/2} = O_p(k), \end{aligned}$$

since  $E\tilde{u}_i^{-2} = O(k)$  and  $E\tilde{u}_i^{-3/2} = O(k^{1/2})$ . We hence obtain

$$\begin{aligned} \left\| \sum_{i=1}^k (\tilde{u}_i^{-1} - \tilde{v}_i^{-1}) X'_i X_i \right\| &= O_p(k) + \left\| \sum_{i \in \alpha^c} (\tilde{u}_i^{-1} - \tilde{v}_i^{-1}) X'_i X_i \right\| \\ &= O_p(k) + O_p(k(\log k)^{1/2}) = o_p(k \log k), \end{aligned}$$

because each of  $\tilde{u}_i$  or  $\tilde{v}_i$  are larger than  $k^{-1}$ . From  $E\tilde{u}_i^{-2} = O(k)$ , we get

$$\begin{aligned} E \left\| \sum_{i=1}^k (\tilde{u}_i^{-1} - E\tilde{u}_i^{-1}) X'_i X_i \right\|^2 &\leq c \sum_{i=1}^k E(\tilde{u}_i^{-1} - E\tilde{u}_i^{-1})^2 = O(k^2), \\ \left\| \sum_{i=1}^k (\tilde{u}_i^{-1} - E\tilde{u}_i^{-1}) X'_i X_i \right\| &= O_p(k) = o_p(k \log k). \end{aligned}$$

Clearly,

$$\left\| \sum_{i=1}^k E\tilde{u}_i^{-1}X_i'X_i \right\| \geq ck \log k.$$

Therefore,

$$\begin{aligned} \hat{\beta}_w - \beta &= \left( \sum_{i=1}^k \tilde{v}_i^{-1}X_i'X_i \right)^{-1} \sum_{i=1}^k \tilde{v}_i^{-1}X_i'e_i \\ &= \left( \sum_{i=1}^k E\tilde{u}_i^{-1}X_i'X_i \right)^{-1} \sum_{i=1}^k (\tilde{u}_i\tilde{v}_i)^{-1}X_i'e_i e_i'X_i(\hat{\beta}_o - \beta) + o_p(k^{-1/2}). \end{aligned}$$

Using the same technique as that in establishing (4.10), we can show that  $(\tilde{u}_i\tilde{v}_i)^{-1}e_i e_i'$  in the preceding equation can be replaced by  $E\tilde{u}_i^{-1}I_2$ . This proves (4.2).  $\square$

Result (4.2) indicates that when  $n_i = 2$  for all  $i$ , the asymptotic variance of  $\hat{\beta}_w^{(m)}$  is the same as that of  $\hat{\beta}_o$  if  $\hat{\beta}_o$  is used as the initial estimator. Therefore, asymptotically  $\hat{\beta}_w^{(m)}$  does not improve  $\hat{\beta}_o$  for any  $m \geq 1$ . Despite this asymptotic result,  $\hat{\beta}_w^{(m)}$  and  $\hat{\beta}_o$  may have different performances for fixed  $k$  (see the simulation result in Section 5).

Note that Theorems 1 and 2 do not require that the within-group errors have the same variance. Hence we may combine some groups with  $n_i \leq 2$  so that all the new groups contain at least three observations and then apply the adaptive procedure.

In some cases there is auxiliary information about variance pattern. For example, we may consider some physical background: The variances may vary with time, community, batch, and so on. A general rule for combining groups is to combine groups with similar variances since combining groups with the same variance generally increases the efficiency of the WLSE.

Since the adaptive procedure picks the OLSE if the OLSE is better than the WLSE, the adaptive estimator with combined groups is asymptotically as good as the OLSE and improves the OLSE when the heteroscedasticity is severe. In Section 5 we show by simulation that with a correct combining method, the adaptive estimator with combined groups significantly improves the OLSE and the WLSE with noncombined groups.

Another method for improving the OLSE when  $n_i$  are very small is the empirical Bayes method. Assuming that the  $\sigma_i$  is random, Hooper (1990) proposed the use of empirical Bayes estimators (EBE). He showed (by asymptotic theory and by simulation) that the EBE improves the Fuller and Rao WLSE  $\hat{\theta}_w^{(1)}$ , provided that we can correctly specify the type of distribution of the random  $\sigma_i$ . Similarly, the method of combining groups requires some auxiliary information about the  $\sigma_i$ . In the case where  $n_i \leq 2$ , it is actually hard

to improve the OLSE without having any auxiliary information about the variances.

**5. Simulation results.** We study by simulation the finite sample bias and root mean square error (rmse) of the OLSE, one-step WLSE [WLSE(1)], two-step WLSE [WLSE(2)] and the adaptive estimator (AWLSE). To illustrate, we force the iterative process to stop after at most two steps so that the adaptive estimator given by (3.6) with  $M = 2$  is used. We consider a common mean problem

$$(5.1) \quad y_{ij} = \mu + e_{ij}, \quad j = 1, 2, \dots, n_o, i = 1, 2, \dots, k,$$

where  $k = 40$ ,  $n_o = 3$  or  $4$  and  $e_{ij}/\sigma_i$  are i.i.d. random variables with common  $N(0, 1)$  distribution. Without loss of generality, we take  $\mu = 0$  in our simulation. Four models with different variance patterns are considered:

MODEL 1 (Homoscedastic model).  $\sigma_i = 1$  for all  $i$ .

MODEL 2.

$i$	1-10	11-20	21-30	31-40
$\sigma_i$	1	1.5	2.25	3.375

MODEL 3.

$i$	1-10	11-20	21-30	31-40
$\sigma_i$	0.5	1	2	4

MODEL 4.  $\{\sigma_i - 0.1, i = 1, \dots, 40\}$  is a random sample from Beta(0.5, 0.5).

Models 1-4 are arranged in increasing order of heteroscedasticity. To examine the performances of various estimators in the case where  $n_o = 2$ , we also consider the following model:

MODEL 5. Model (5.1) but  $n_o = 2$ ,  $k = 80$ .  $\{\sigma_i - 0.1, i = 1, \dots, 40\}$  and  $\{\sigma_i - 0.1, i = 41, \dots, 80\}$  are two ordered random samples from Beta(0.5, 0.5).

Under Model 5, in addition to the OLSE and WLSE, we also consider the estimators based on the data with combined groups and denote the one-step, two-step and adapted estimators by CWLSE(1), CWLSE(2) and CAWLSE, respectively. Two combining methods are used:

1. Randomly combining every two groups;
2. Combining the  $i$ th group with the  $(40 + i)$ th group,  $i = 1, 2, \dots, 40$ .

Note that for the second method, we actually assume that we know that the variances for the  $i$ th group and the  $(40 + i)$ th group are close. This allows us

to see what is the best we can do by combining groups. On the other hand, the first method (random combining) is the most inefficient combining method and hence we can see the worst.

The results in Table 1 are based on 10,000 repetitions. In addition to the biases and rmse's of various estimators, Table 1 also shows the number of times (frequency) that AWLSE equals OLSE, WLSE(1) and WLSE(2) in 10,000 repetitions.

The following is a summary of the simulation results:

1. The performance (in terms of rmse and bias) of the adaptive estimator AWLSE is generally good and is always close to the best estimator among the OLSE, one-step WLSE and two-step WLSE. In all cases, the biases are negligible, due to the fact that the OLSE and WLSE are unbiased when the errors have symmetric distributions.
2. In terms of the rmse, the improvement of the adaptive estimator over the OLSE can be as large as 48% when the heteroscedasticity is severe.
3. Except for Model 2 with  $n_o = 3$ , the adaptive procedure picks the winner with a high probability (in many cases the probability is over 90%).
4. When  $n_o = 2$  (Model 5), the WLSE(1) improves the OLSE in terms of rmse (about 20%). This does not conflict with Theorem 8 which is an asymptotic ( $k \rightarrow \infty$ ) result. Also, one of the conditions in Theorem 8 is that the ratio  $\gamma_k = (\max_i \sigma_i^2)/(\min_i \sigma_i^2)$  is bounded. In Model 5, however, we find  $\gamma_k = 118$  and  $\gamma_k/k = 1.5$  which is not small. Generally speaking, the WLSE(1) may still improve the OLSE when  $n_o = 2$  provided  $\gamma_k/k$  is not relatively small.
5. In Model 5, when an inefficient combining method (randomly combining) is used, the performance of the CWLSE is almost the same as that of the WLSE. On the other hand, when the groups are combined with some prior knowledge (the second method), the CWLSE and CAWLSE have very good performances: They improve the OLSE in terms of rmse (about 45%) and also improve the WLSE.

TABLE 1  
*Simulation results  $b = \text{bias} \times 10^4$ ,  $r = \text{root mean square error} \times 10^2$ ,  
 $f = \text{frequency distribution of AWLSE}$*

<b>Model 1</b>				
$(n_o, k) = (3, 40)$	OLSE	WLSE(1)	WLSE(2)	AWLSE
$b$	-4.6118	-3.8018	-3.8918	-4.4189
$r$	9.1800	10.1963	11.587	9.3039
$f$	9028	660	312	
$(n_o, k) = (4, 40)$	OLSE	WLSE(1)	WLSE(2)	AWLSE
$b$	-7.7597	-7.8325	-7.5796	-7.1013
$r$	7.8718	8.7922	9.7372	7.9627
$f$	9222	703	75	

TABLE 1  
(Continued)

<b>Model 2</b>				
$(n_o, k) = (3, 40)$	OLSE	WLSE(1)	WLSE(2)	AWLSE
$b$	2.6376	-2.0820	-4.1923	-4.7963
$r$	20.2841	17.8672	18.4310	19.0425
$f$	3193	4326	2481	
$(n_o, k) = (4, 40)$	OLSE	WLSE(1)	WLSE(2)	AWLSE
$b$	-11.1978	-10.6913	-9.9395	-13.6378
$r$	17.5123	14.6337	15.1137	14.9461
$f$	549	6954	2506	
<b>Model 3</b>				
$(n_o, k) = (3, 40)$	OLSE	WLSE(1)	WLSE(2)	AWLSE
$b$	8.9063	-0.0209	-3.2834	-1.8743
$r$	21.0282	14.0737	12.4129	12.5488
$f$	1	200	9799	
$(n_o, k) = (4, 40)$	OLSE	WLSE(1)	WLSE(2)	AWLSE
$b$	-8.1214	-7.5647	-6.5526	-6.4433
$r$	18.2416	10.6878	9.4770	9.4825
$f$	0	5	9995	
<b>Model 4</b>				
$(n_o, k) = (3, 40)$	OLSE	WLSE(1)	WLSE(2)	AWLSE
$b$	1.5179	-0.4538	-1.5582	-1.1074
$r$	6.6878	4.1942	3.5880	3.6581
$f$	5	112	9883	
$(n_o, k) = (4, 40)$	OLSE	WLSE(1)	WLSE(2)	AWLSE
$b$	-7.8392	-2.7950	-0.6531	-0.6434
$r$	5.7416	3.1513	2.7081	2.7100
$f$	0	3	9997	
<b>Model 5</b>				
Noncombining				
$(n_o, k) = (2, 80)$	OLSE	WLSE(1)	WLSE(2)	
$b$	-2.7945	-3.4378	-3.6127	
$r$	5.6685	4.5148	4.0397	
Combining method (1)				
$(n_o, k) = (2, 80)$	OLSE	CWLSE(1)	CWLSE(2)	CAWLSE
$b$	-2.7945	-1.6845	-0.4240	-1.4609
$r$	5.6685	4.2820	4.2597	4.3298
$f$	163	4623	5214	
Combining method (2)				
$(n_o, k) = (2, 80)$	OLSE	CWLSE(1)	CWLSE(2)	CAWLSE
$b$	-2.7945	-2.7571	-2.4675	-2.5271
$r$	5.6685	3.4309	3.1166	3.1293
$f$	0	74	9926	

**Acknowledgments.** The authors would like to thank an Associate Editor and two referees for their helpful comments and suggestions.

### REFERENCES

- BOX, G. E. P. and TIAO, G. C. (1973). *Bayesian Inference in Statistical Analysis*. Addison-Wesley, Reading, MA.
- CARROLL, R. J. (1982). Adapting for heteroscedasticity in linear models. *Ann. Statist.* **10** 1224–1233.
- CARROLL, R. J. and CLINE, D. B. H. (1988). An asymptotic theory for weighted least-squares with weights estimated by replication. *Biometrika* **75** 35–43.
- CARROLL, R. J., WU, C. F. J. and RUPPERT, D. (1988). The effect of estimating weights in weighted least squares. *J. Amer. Statist. Assoc.* **83** 1045–1054.
- DAVIDIAN, M. and CARROLL, R. J. (1987). Variance function estimation. *J. Amer. Statist. Assoc.* **82** 1079–1091.
- FULLER, W. A. and RAO, J. N. K. (1978). Estimation for a linear regression model with unknown diagonal covariance matrix. *Ann. Statist.* **6** 1149–1158.
- HOOPER, P. M. (1990). Iterative weighted least squares estimation in heteroscedastic linear regression models. Preprint.
- RAO, C. R. (1970). Estimation of heteroscedastic variances in linear models. *J. Amer. Statist. Assoc.* **65** 161–172.
- RAO, J. N. K. (1973). On the estimation of heteroscedastic variances. *Biometrics* **29** 11–24.
- SHAO, J. (1989a). Asymptotic distribution of the weighted least squares estimator. *Ann. Inst. Statist. Math.* **41** 365–382.
- SHAO, J. (1989b). Jackknifing weighted least squares estimators. *J. Roy. Statist. Soc. Ser. B* **51** 139–156.

DEPARTMENT OF STATISTICS  
AND ACTUARIAL SCIENCE  
UNIVERSITY OF WATERLOO  
WATERLOO, ONTARIO N2L 3G1  
CANADA

DEPARTMENT OF MATHEMATICS  
UNIVERSITY OF OTTAWA  
OTTAWA, ONTARIO K1N 6N5  
CANADA