

## MODEL AVERAGING BASED ON KULLBACK-LEIBLER DISTANCE

Xinyu Zhang<sup>1</sup>, Guohua Zou<sup>1,2</sup> and Raymond J. Carroll<sup>3</sup>

<sup>1</sup>Chinese Academy of Sciences, <sup>2</sup>Capital Normal University and <sup>3</sup>Texas A&M University

### Supplementary Material

This file contains proofs and additional simulation results as follows:

- Section S1: Proof of Theorem 1;
- Section S2: Proof of Theorem 2;
- Section S3: Proof of Theorem 3;
- Section S4: Proof of (2.4);
- Section S5: Proof of Theorem 4;
- Section S6: Calculation for  $\partial \hat{\eta}^T / \partial y$ ;
- Section S7: Discussion on Assumption (A.1) and its relationship with the normality of  $e$ ;
- Section S8: All Simulation Results of Example 1;
- Section S9: All Simulation Results of Example 2;
- Section S10: All Simulation Results of Example 3;
- Section S11: Simulation Results with a Uniformly Distributed Error Term;
- Section S12: Simulation Results with a Chi-squared Distributed Error Term;
- Section S13: Simulation Results with Coefficients Depending on the Sample Size.

### S1 Proof of Theorem 1

This proof is mainly an application of Stein's Lemma (Stein (1981)). Let  $y_i$ ,  $\mu_i$  and  $\hat{\mu}(w)_i$  be the  $i^{\text{th}}$  elements of  $y$ ,  $\mu$  and  $\hat{\mu}(w)$ , respectively. By the conditions stated in Theorem 1 and Stein's Lemma, we have

$$E_{f(y)} \left[ \frac{(y_i - \mu_i) \{y_i - \hat{\mu}(w)_i\}}{\hat{\sigma}^2} \right] = \sigma^2 E_{f(y)} \left( \frac{\partial [\{y_i - \hat{\mu}(w)_i\} / \hat{\sigma}^2]}{\partial y_i} \right) \quad (\text{S1.1})$$

and

$$E_{f(y)} \left\{ \frac{(y_i - \mu_i)^2}{\hat{\sigma}^2} \right\} = \sigma^2 E_{f(y)}(\hat{\sigma}^{-2}) + \sigma^4 E_{f(y)} \left\{ \frac{\partial^2(\hat{\sigma}^{-2})}{\partial y_i^2} \right\}. \quad (\text{S1.2})$$

Using (S1.1), it can be seen that

$$\begin{aligned} E_{f(y)} \left[ \frac{(y - \mu)^T \{y - \hat{\mu}(w)\}}{\hat{\sigma}^2} \right] &= \sigma^2 E_{f(y)} \left( \text{trace} \frac{\partial [\{y - \hat{\mu}(w)\} / \hat{\sigma}^2]}{\partial y^T} \right) \\ &= E_{f(y)}(n\sigma^2 \hat{\sigma}^{-2}) - E_{f(y)} \text{trace} \left\{ \sigma^2 \hat{\sigma}^{-2} \frac{\partial \hat{\mu}(w)}{\partial y^T} \right\} \\ &\quad - E_{f(y)} \left[ \sigma^2 \hat{\sigma}^{-4} \{y - \hat{\mu}(w)\}^T \frac{\partial \hat{\sigma}^2}{\partial y} \right]. \end{aligned} \quad (\text{S1.3})$$

Using (S1.2), we obtain

$$\begin{aligned} E_{f(y)} \left( \frac{\|y - \mu\|^2}{\hat{\sigma}^2} \right) &= n\sigma^2 E_{f(y)}(\hat{\sigma}^{-2}) + \sigma^4 E_{f(y)} \left[ \text{trace} \left\{ \frac{\partial^2(\hat{\sigma}^{-2})}{\partial y \partial y^T} \right\} \right] \\ &= E_{f(y)}(n\sigma^2 \hat{\sigma}^{-2}) + E_{f(y)} \left\{ 2\sigma^4 \hat{\sigma}^{-6} \text{trace} \left( \frac{\partial \hat{\sigma}^2}{\partial y} \frac{\partial \hat{\sigma}^2}{\partial y^T} \right) \right\} \\ &\quad - E_{f(y)} \left\{ \sigma^4 \hat{\sigma}^{-4} \text{trace} \left( \frac{\partial^2 \hat{\sigma}^2}{\partial y \partial y^T} \right) \right\}. \end{aligned} \quad (\text{S1.4})$$

In addition, it is straightforward to show that

$$\|\mu - \hat{\mu}(w)\|^2 = \|y - \hat{\mu}(w)\|^2 + \|y - \mu\|^2 - 2\{y - \hat{\mu}(w)\}^T (y - \mu). \quad (\text{S1.5})$$

Now, by combining (S1.3)-(S1.5), Theorem 1 is proved.

## S2 Proof of Theorem 2

Based on the proofs of Theorems 1' and 2 in Wan, Zhang and Zou (2010), and Assumptions (A.1), (A.2) and (A.4), to prove Theorem 2, we need only to verify that

$$\sup_{w \in \mathcal{W}} \left\{ \left| y^T P^T(w) \frac{\partial \hat{\sigma}^2}{\partial y} \Big| R_n^{-1}(w) \right\} = o_p(1). \quad (\text{S2.1})$$

By Assumptions (A.2), (A.3), and (A.5), we have

$$\begin{aligned}
 & \sup_{w \in \mathcal{W}} \left\{ \left| y^T P^T(w) \frac{\partial \hat{\sigma}^2}{\partial y} \right| R_n^{-1}(w) \right\} \leq \xi_n^{-1} \sup_{w \in \mathcal{W}} \left| y^T P^T(w) \hat{T} y \right| \\
 & = \xi_n^{-1} \sup_{w \in \mathcal{W}} \left| y^T \left\{ P^T(w) \hat{T} + \hat{T}^T P(w) \right\} y \right| / 2 \\
 & \leq \xi_n^{-1} \sup_{w \in \mathcal{W}} \left[ \lambda_{\max} \left\{ P^T(w) \hat{T} + \hat{T}^T P(w) \right\} \right] \|y\|^2 / 2 \\
 & \leq \xi_n^{-1} \sup_{w \in \mathcal{W}} \left[ \lambda_{\max} \{P(w)\} \lambda_{\max}(\hat{T}) \|y\|^2 \right] \\
 & \leq \xi_n^{-1} \sup_{w \in \mathcal{W}} \left\{ \max_{s \in \{1, \dots, S\}} \lambda_{\max}(P_{(s)}) \right\} \lambda_{\max}(\hat{T}) \|y\|^2 \\
 & = \xi_n^{-1} \max_{s \in \{1, \dots, S\}} \lambda_{\max}(P_{(s)}) \cdot n \lambda_{\max}(\hat{T}) \cdot n^{-1} \|y\|^2 = o_p(1). \tag{S2.2}
 \end{aligned}$$

This completes the proof.

### S3 Proof of Theorem 3

By the assumption that  $\mu$  is a linear function of  $X$ , we have

$$y^T (I_n - P)y = e^T (I_n - P)e,$$

where  $I_n - P$  is a symmetric idempotent matrix with rank  $n - m$ . So  $\sigma^2 / \{y^T (I_n - P)y\}$  is distributed as an inverse Chi-squared distribution with mean  $(n - m - 2)^{-1}$ , and thus

$$E_{f(y)} \left\{ \frac{\sigma^2}{\hat{\sigma}^2(y, k)} \right\} = \frac{k}{n - m - 2}. \tag{S3.1}$$

In addition,

$$\begin{aligned}
 \{y - \hat{\mu}(w)\}^T \frac{\partial \hat{\sigma}^2(y, k)}{\partial y} & = 2k^{-1} \{y - \hat{\mu}(w)\}^T (I_n - P)y \\
 & = 2k^{-1} y^T \sum_{s=1}^S w_s (I_n - P_{(s)}) (I_n - P)y \\
 & = 2k^{-1} y^T (I_n - P)y. \tag{S3.2}
 \end{aligned}$$

By applying (S3.1), (S3.2), and the facts that  $\partial \hat{\sigma}^2(y, k) / \partial y = 2k^{-1} (I_n - P)y$  and  $\partial^2 \hat{\sigma}^2(y, k) / (\partial y \partial y^T) = 2k^{-1} (I_n - P)$  to Theorem 1, we obtain Theorem 3.

## S4 Proof of (2.4)

Let  $\tilde{\mathcal{C}}^*(w) = \mathcal{C}^*(w) - \|e\|^2$  such that  $\hat{w} = \underset{w \in \mathcal{W}}{\operatorname{argmin}} \tilde{\mathcal{C}}^*(w)$ . Write  $\tilde{w} = \underset{w \in \mathcal{W}}{\operatorname{argmin}} E\{L_n(w)\}$ . It is straightforward to show that

$$\begin{aligned} L_n(\hat{w}) &= \tilde{\mathcal{C}}^*(\hat{w}) - c(\hat{w}) \\ &= \inf_{w \in \mathcal{W}} \tilde{\mathcal{C}}^*(w) - c(\hat{w}) \\ &= \inf_{w \in \mathcal{W}} \{L_n(w) + c(w)\} - c(\hat{w}) \\ &\leq L_n(\tilde{w}) + e'\{I_n - P(\tilde{w})\}\mu + \sigma^2 \operatorname{trace}\{P(\tilde{w})\} - e'P(\tilde{w})e - c(\hat{w}), \end{aligned}$$

which implies (2.4).

## S5 Proof of Theorem 4

Let  $\tilde{y} = \Omega^{-1/2}y$ ,  $\tilde{\mu} = \Omega^{-1/2}\mu$ ,  $\tilde{e} = \Omega^{-1/2}e$ , and  $\tilde{\Omega} = \Omega^{-1/2}\hat{\Omega}\Omega^{-1/2}$ . By simple calculations, we have

$$\begin{aligned} \mathcal{D}^*(w) &= \{\tilde{y} - \tilde{P}(w)\tilde{y}\}^T \{\tilde{y} - \tilde{P}(w)\tilde{y}\} + 2\operatorname{trace}\{\tilde{P}(w)\} \\ &\quad + \{\tilde{\mu} - \tilde{P}(w)\tilde{\mu}\}^T (\tilde{\Omega}^{-1} - I_n) \{\tilde{\mu} - \tilde{P}(w)\tilde{\mu}\} + \{\tilde{P}(w)\tilde{e}\}^T (\tilde{\Omega}^{-1} - I_n) \tilde{P}(w)\tilde{e} \\ &\quad - 2\{\tilde{\mu} - \tilde{P}(w)\tilde{\mu}\}^T (\tilde{\Omega}^{-1} - I_n) \tilde{P}(w)\tilde{e} + 2\tilde{e}^T (\tilde{\Omega}^{-1} - I_n) \{\tilde{\mu} - \tilde{P}(w)\tilde{\mu}\} \\ &\quad - 2\tilde{e}^T (\tilde{\Omega}^{-1} - I_n) \tilde{P}(w)\tilde{e} - 2\tilde{y}^T \tilde{P}^T(w) \tilde{\Omega}^{-1} \Omega^{-1/2} \hat{a} + \tilde{e}^T (\tilde{\Omega}^{-1} - I_n) \tilde{e}. \end{aligned} \quad (\text{S5.1})$$

It follows from Assumptions (A.2) and (B.2) that

$$\sup_{w \in \mathcal{W}} \lambda_{\max}\{\tilde{P}(w)\} = O(1). \quad (\text{S5.2})$$

From the proof of Theorem 1' in Wan, Zhang and Zou (2010), (S5.1), (S5.2), Assumptions (B.1)-(B.2), and the fact that  $\tilde{e}^T (\tilde{\Omega}^{-1} - I_n) \tilde{e}$  is unrelated to  $w$ , in order to prove (3.2), we need only to verify that

$$\sup_{w \in \mathcal{W}} [R_{\text{hetero},n}^{-1}(w) \{ \tilde{\mu} - \tilde{P}(w)\tilde{\mu} \}^T (\tilde{\Omega}^{-1} - I_n) \{ \tilde{\mu} - \tilde{P}(w)\tilde{\mu} \}] = o_p(1), \quad (\text{S5.3a})$$

$$\sup_{w \in \mathcal{W}} [R_{\text{hetero},n}^{-1}(w) \{ \tilde{P}(w)\tilde{e} \}^T (\tilde{\Omega}^{-1} - I_n) \tilde{P}(w)\tilde{e}] = o_p(1), \quad (\text{S5.3b})$$

$$\sup_{w \in \mathcal{W}} [R_{\text{hetero},n}^{-1}(w) \{ \tilde{\mu} - \tilde{P}(w)\tilde{\mu} \}^T (\tilde{\Omega}^{-1} - I_n) \tilde{P}(w)\tilde{e}] = o_p(1), \quad (\text{S5.3c})$$

$$\sup_{w \in \mathcal{W}} [R_{\text{hetero},n}^{-1}(w) \{ \tilde{e}^T (\tilde{\Omega}^{-1} - I_n) \{ \tilde{\mu} - \tilde{P}(w)\tilde{\mu} \} \}] = o_p(1), \quad (\text{S5.3d})$$

$$\sup_{w \in \mathcal{W}} \{ R_{\text{hetero},n}^{-1}(w) \{ \tilde{e}^T (\tilde{\Omega}^{-1} - I_n) \tilde{P}(w)\tilde{e} \} \} = o_p(1), \quad (\text{S5.3e})$$

$$\sup_{w \in \mathcal{W}} \{ R_{\text{hetero},n}^{-1}(w) \{ \tilde{y}^T \tilde{P}^T(w) \tilde{\Omega}^{-1} \Omega^{-1/2} \hat{a} \} \} = o_p(1). \quad (\text{S5.3f})$$

By (S5.2) and Assumptions (A.3) and (B.2), it is straightforward to show that

$$\begin{aligned} R_{\text{hetero},n}(w) &= \|\tilde{\mu} - \tilde{P}(w)\tilde{\mu}\|^2 + \text{trace}\{\tilde{P}(w)\tilde{P}^T(w)\} \\ &\geq \max\{\|\tilde{\mu} - \tilde{P}(w)\tilde{\mu}\|^2, \text{trace}\{\tilde{P}(w)\tilde{P}^T(w)\}\}, \end{aligned} \quad (\text{S5.4})$$

$$\begin{aligned} R_{\text{hetero},n}(w) &= \mu^T \{I_n - P(w)\}^T \Omega^{-1} \{I_n - P(w)\} \mu + \text{trace}\{P(w)\Omega P^T(w)\Omega^{-1}\} \\ &= O(n), \end{aligned} \quad (\text{S5.5})$$

and that

$$\sup_{w \in \mathcal{W}} \{R_{\text{hetero},n}^{-1}(w) \|\tilde{P}(w)\tilde{e}\|^2\} = o_p(1) + \sup_{w \in \mathcal{W}} [R_{\text{hetero},n}^{-1}(w) \text{trace}\{\tilde{P}(w)\tilde{P}^T(w)\}] \quad (\text{S5.6})$$

by the proof of Theorem 1' in Wan, Zhang and Zou (2010).

Using (S5.5), we see that  $n^{-1}\xi_{\text{hetero},n} = O(1)$ , which together with Assumption (B.3) implies that

$$\max_{i \in \{1, \dots, n\}} |\hat{\Omega}_{ii} - \Omega_{ii}| = o_p(1). \quad (\text{S5.7})$$

Using (S5.7) and Assumption (B.2), we have

$$\lambda_{\max}(\tilde{\Omega}^{-1} - I_n) = O_p\left(\max_{i \in \{1, \dots, n\}} |\hat{\Omega}_{ii} - \Omega_{ii}|\right) = o_p(1). \quad (\text{S5.8})$$

Thus, from (S5.4) and (S5.8), we obtain

$$\begin{aligned} &\sup_{w \in \mathcal{W}} [R_{\text{hetero},n}^{-1}(w) \{ \tilde{\mu} - \tilde{P}(w)\tilde{\mu} \}^T (\tilde{\Omega}^{-1} - I_n) \{ \tilde{\mu} - \tilde{P}(w)\tilde{\mu} \}] \\ &\leq \lambda_{\max}(\tilde{\Omega}^{-1} - I_n) \sup_{w \in \mathcal{W}} [R_{\text{hetero},n}^{-1}(w) \{ \tilde{\mu} - \tilde{P}(w)\tilde{\mu} \}^T \{ \tilde{\mu} - \tilde{P}(w)\tilde{\mu} \}] \\ &\leq \lambda_{\max}(\tilde{\Omega}^{-1} - I_n) = o_p(1), \end{aligned}$$

which is the result (S5.3a). Using (S5.6), (S5.8), and Assumption (B.4), it is seen that

$$\begin{aligned} &\sup_{w \in \mathcal{W}} [R_{\text{hetero},n}^{-1}(w) \{ \tilde{P}(w)\tilde{e} \}^T (\tilde{\Omega}^{-1} - I_n) \tilde{P}(w)\tilde{e}] \\ &\leq \lambda_{\max}(\tilde{\Omega}^{-1} - I_n) \sup_{w \in \mathcal{W}} \{ R_{\text{hetero},n}^{-1}(w) \|\tilde{P}(w)\tilde{e}\|^2 \} \\ &= \lambda_{\max}(\tilde{\Omega}^{-1} - I_n) \left( \sup_{w \in \mathcal{W}} [R_{\text{hetero},n}^{-1}(w) \text{trace}\{\tilde{P}(w)\tilde{P}^T(w)\}] + o_p(1) \right) \\ &= o_p(1), \end{aligned} \quad (\text{S5.9})$$

which is the result (S5.3b). Using (S5.4), (S5.6), (S5.8), and Assumption (B.4), we obtain

$$\begin{aligned} &\sup_{w \in \mathcal{W}} (R_{\text{hetero},n}^{-2}(w) \{ \tilde{\mu} - \tilde{P}(w)\tilde{\mu} \}^T (\tilde{\Omega}^{-1} - I_n) \tilde{P}(w)\tilde{e})^2 \\ &\leq \sup_{w \in \mathcal{W}} \{ R_{\text{hetero},n}^{-1}(w) \|(\tilde{\Omega}^{-1} - I_n) \tilde{P}(w)\tilde{e}\|^2 \} \\ &\leq \{ \lambda_{\max}(\tilde{\Omega}^{-1} - I_n) \}^2 \sup_{w \in \mathcal{W}} \{ R_{\text{hetero},n}^{-1}(w) \|\tilde{P}(w)\tilde{e}\|^2 \} = o_p(1), \end{aligned}$$

which implies the result (S5.3c). Using (S5.4), (S5.8), and Assumptions (B.2)-(B.3), it follows that

$$\begin{aligned} & \sup_{w \in \mathcal{W}} [R_{\text{hetero},n}^{-2}(w) \{\tilde{e}^T (\tilde{\Omega}^{-1} - I_n) \{\tilde{\mu} - \tilde{P}(w)\tilde{\mu}\}\}^2] \\ & \leq \{\lambda_{\max}(\tilde{\Omega}^{-1} - I_n)\}^2 \sup_{w \in \mathcal{W}} \{R_{\text{hetero},n}^{-1}(w)\} \|\tilde{e}\|^2 \\ & \leq \{\lambda_{\max}(\tilde{\Omega}^{-1} - I_n)\}^2 n \xi_{\text{hetero},n}^{-1} n^{-1} \|\tilde{e}\|^2 = o_p(1), \end{aligned}$$

which implies the result (S5.3d). Similarly, from (S5.4), (S5.6), (S5.8), and Assumptions (B.2)-(B.3), we have

$$\begin{aligned} & \sup_{w \in \mathcal{W}} [R_{\text{hetero},n}^{-2}(w) \{\tilde{e}^T (\tilde{\Omega}^{-1} - I_n) \tilde{P}(w) \tilde{e}\}^2] \\ & \leq \{\lambda_{\max}(\tilde{\Omega}^{-1} - I_n)\}^2 \xi_{\text{hetero},n}^{-1} \|\tilde{e}\|^2 \sup_{w \in \mathcal{W}} \{R_{\text{hetero},n}^{-1}(w)\} \|\tilde{P}(w) \tilde{e}\|^2 \\ & = o_p(1), \end{aligned}$$

which implies the result (S5.3e). Finally, using steps similar to (S2.2), it follows from (S5.2), (S5.7), and Assumptions (A.3), (B.2) and (B.5) that

$$\begin{aligned} & \sup_{w \in \mathcal{W}} \{R_{\text{hetero},n}^{-1}(w) |\tilde{y}^T \tilde{P}(w) \tilde{\Omega}^{-1} \Omega^{1/2} \hat{a}|\} = \sup_{w \in \mathcal{W}} \{R_{\text{hetero},n}^{-1}(w) |\tilde{y}^T \tilde{P}(w) \tilde{\Omega}^{-1} \Omega^{1/2} \hat{A} \tilde{y}|\} \\ & = \sup_{w \in \mathcal{W}} [R_{\text{hetero},n}^{-1}(w) |\tilde{y}^T \{\tilde{P}(w) \tilde{\Omega}^{-1} \Omega^{1/2} \hat{A} + \hat{A}^T \Omega^{1/2} \tilde{\Omega}^{-1} \tilde{P}(w)^T\} \tilde{y}|/2] \\ & \leq \lambda_{\max}\{\tilde{P}(w)\} \lambda_{\max}(\tilde{\Omega}^{-1}) \lambda_{\max}(\Omega^{1/2}) n^{-1} \|\tilde{y}\|^2 n \lambda_{\max}(\hat{A}) \xi_{\text{hetero},n}^{-1} \\ & = o_p(1), \end{aligned}$$

which is the result (S5.3f). This completes the proof.

## S6 Calculation for $\partial \hat{\eta}^T / \partial y$

Assuming that  $\mu = X\beta$ , the -2log-likelihood is

$$n \log 2\pi + \log |\Omega(\eta)| + (y - X\beta)^T \Omega^{-1}(\eta) (y - X\beta).$$

Write  $\hat{\eta} = (\hat{\eta}_1, \dots, \hat{\eta}_q)^T$ . Let  $\hat{\beta}$  be the ML estimator of  $\beta$ ,  $U_j(\hat{\eta}) = \partial \Omega^{-1}(\hat{\eta}) / \partial \hat{\eta}_j$ , and  $V_{jk}(\hat{\eta}) = \partial U_j(\hat{\eta}) / \partial \hat{\eta}_k$  for  $j, k = 1, \dots, q$ . From Magnus and Neudecker (1988), we have the following score equations:

$$\begin{cases} (y - X\hat{\beta})^T \Omega^{-1}(\hat{\eta}) X = 0, \\ (y - X\hat{\beta})^T U_j(\hat{\eta}) (y - X\hat{\beta}) - \text{trace}\{U_j(\hat{\eta}) \Omega(\hat{\eta})\} = 0, \quad j = 1, \dots, q. \end{cases} \quad (\text{S6.1})$$

Let  $K$  be a  $q \times n$  matrix with the  $j^{\text{th}}$  row  $K_j = (y - X\hat{\beta})^T U_j(\hat{\eta})$ ,  $Q$  be a  $q \times q$  matrix with the  $jk^{\text{th}}$  element  $Q_{jk} = (y - X\hat{\beta})^T V_{jk}(\hat{\eta}) (y - X\hat{\beta})$ , and  $B$  be a  $q \times q$  matrix with the  $jk^{\text{th}}$  element  $B_{jk} = \partial \text{trace}\{U_j(\hat{\eta}) \Omega(\hat{\eta})\} / \partial \hat{\eta}_k$ .

Taking derivatives with respect to  $y$  on both sides of (S6.1), we have

$$\begin{cases} \frac{\partial \hat{\eta}^T}{\partial y} K X + \Omega^{-1}(\hat{\eta}) X - \frac{\partial \hat{\beta}^T}{\partial y} X^T \Omega^{-1}(\hat{\eta}) X = 0, \\ \frac{\partial \hat{\eta}^T}{\partial y} Q^T + 2K^T - 2\frac{\partial \hat{\beta}^T}{\partial y} X^T K^T - \frac{\partial \hat{\eta}^T}{\partial y} B^T = 0, \end{cases} \quad (\text{S6.2})$$

by which, we have

$$\begin{aligned} \frac{\partial \hat{\eta}^T}{\partial y} &= [2\Omega^{-1}(\hat{\eta}) X \{X^T \Omega^{-1}(\hat{\eta}) X\}^{-1} X^T K^T - 2K^T] \\ &\quad \times [Q^T - B^T - 2K X \{X^T \Omega^{-1}(\hat{\eta}) X\}^{-1} X^T K^T]^{-1}, \end{aligned}$$

given the above inverses exist.

## S7 Discussion on Assumption (A.1) and its relationship with the normality of $e$

It is seen that

$$\begin{aligned} R_n(w) &= E\{L_n(w)\} = E\{\|\hat{\mu}(w) - \mu\|^2\} = E\{\|P(w)(\mu + e) - \mu\|^2\} \\ &= \|P(w)\mu - \mu\|^2 + \sigma^2 \text{trace}\{P(w)P^T(w)\}, \end{aligned}$$

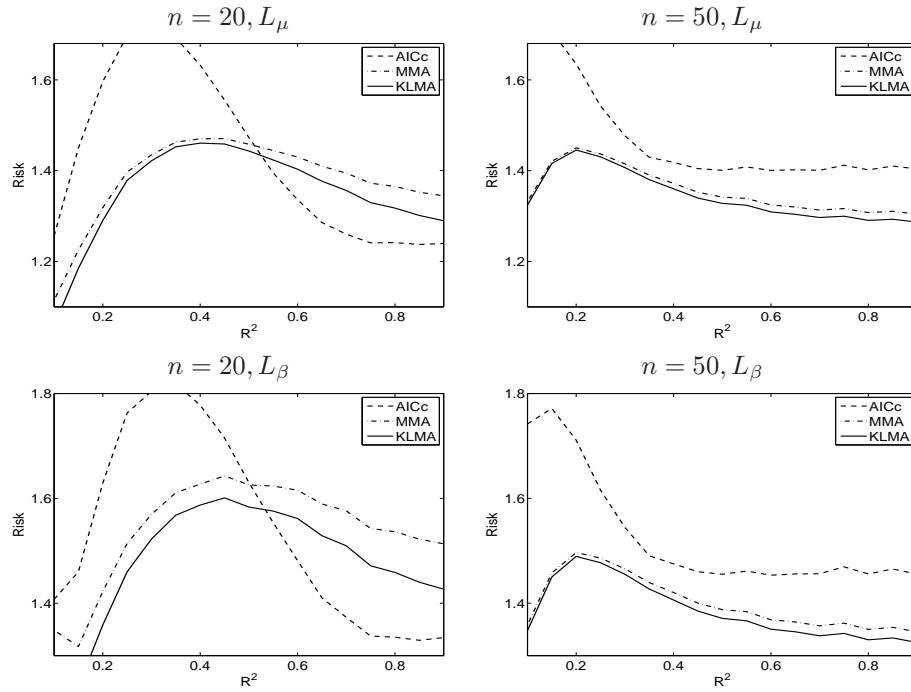
which depends on the mean and variance of  $e_i$ , but is unrelated to other properties of  $e_i$ . So given the mean zero and the variance  $\sigma^2$ , Assumption (A.1) depends on the distribution of  $e_i$  only through the integer  $G$ , which is determined by the first part of Assumption (A.1). The first part of Assumption (A.1) is just a moment condition, which is satisfied for any  $G < \infty$  when  $e_i$  follows normal, uniform or Chi-squared distribution. Thus, if Assumption (A.1) holds for normal  $e_i$ , then it also holds for uniform or Chi-squared  $e_i$ .

In the following, we discuss the assumption

$$S \xi_n^{-2G} \sum_{s=1}^S R_n^G(w_s^0) = o(1). \quad (\text{S7.1})$$

First,  $\xi_n \rightarrow \infty$  is a necessary condition of the assumption (S7.1). As Hansen and Racine (2012) remarked, this condition requires that all finite dimensional models be approximations. Let  $\eta_n = \max_{s \in \{1, \dots, S\}} R_n(w_s^0)$ . A sufficient condition of the assumption (S7.1) is  $S^2(\xi_n^{-2}\eta_n)^G = o(1)$ . For a fixed rate of  $\xi_n \rightarrow \infty$ , the slower the rates of  $S \rightarrow \infty$  and  $\eta_n \rightarrow \infty$ , the faster the rate of  $S^2(\xi_n^{-2}\eta_n)^G \rightarrow 0$ . Practically, the rates of  $S \rightarrow \infty$  and  $\eta_n \rightarrow \infty$  can be reduced by removing the very poor models at the outset prior to model averaging.

In fact, the assumption (S7.1) is introduced by Wan, Zhang and Zou (2010), where its rationality was discussed in detail. Since that, this condition has been used in model averaging studies such as Liu and Okui (2013).

**S8 All Simulation Results of Example 1**Figure S.1: Results for Example 1: risk comparisons under  $L_\mu$  and  $L_\beta$  as a function of  $R^2$ .



### S9 All Simulation Results of Example 2

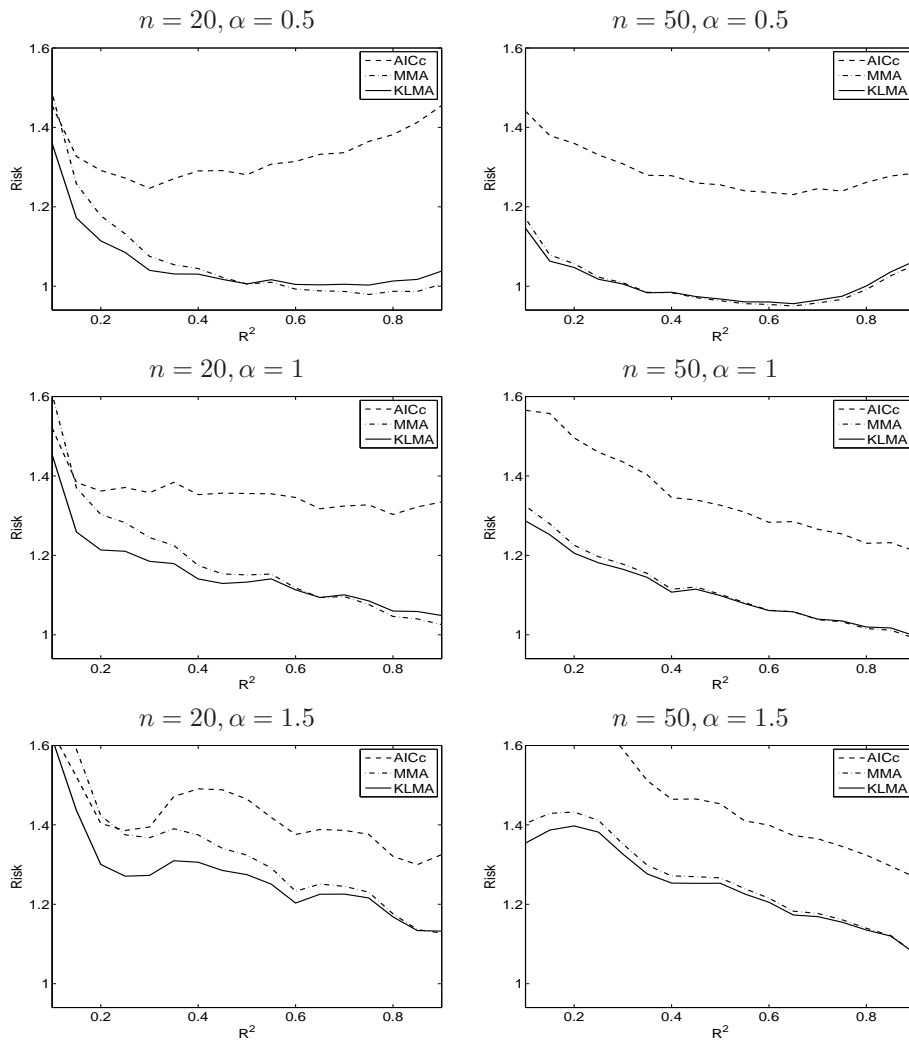
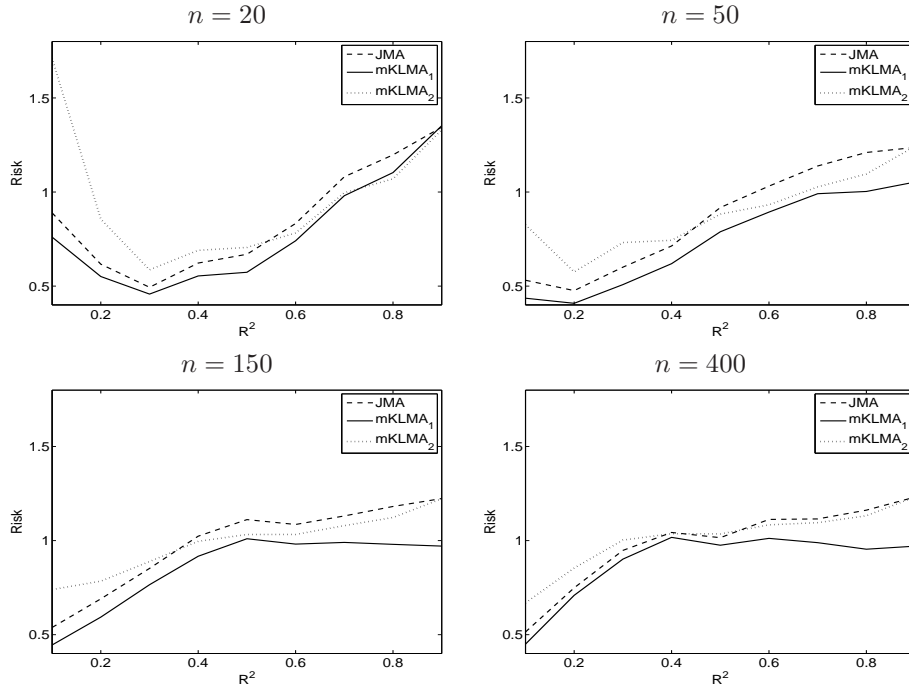


Figure S.2: Results for Example 2: risk comparisons under  $L_\mu$  as a function of  $R^2$ .

**S10 All Simulation Results of Example 3**Figure S.3: Results for Example 3: risk comparisons under  $L_\mu$  as a function of  $R^2$ .

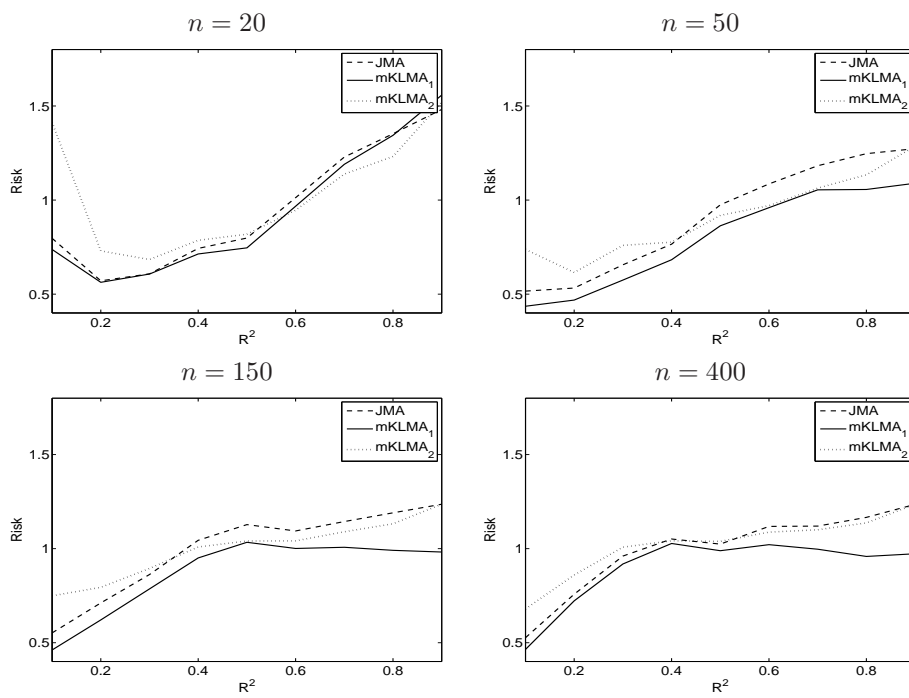


Figure S.4: Results for Example 3: risk comparisons under  $L_\beta$  as a function of  $R^2$ .

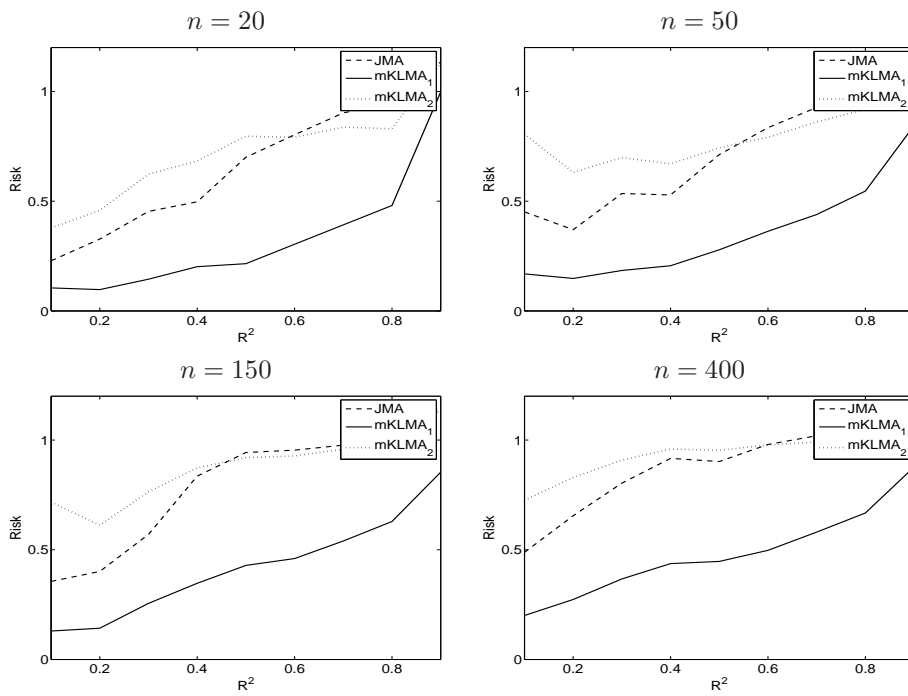


Figure S.5: Results for Example 3: risk comparisons under  $L_{\text{hetero}, \mu}$  as a function of  $R^2$ .

## S11 Simulation Results with a Uniformly Distributed Error Term

In this section, we repeat Examples 1-3 but using a uniformly distributed error term.

| Example 1                                  | Example 2                           | Example 3   |
|--|-------------------------------------|---|
| $\sigma \text{Uniform}(-3^{1/2}, 3^{1/2})$ | $\text{Uniform}(-3^{1/2}, 3^{1/2})$ | $\exp(\eta X_{2i})^{1/2} \text{Uniform}(-3^{1/2}, 3^{1/2})$ |

The following figures show new results.

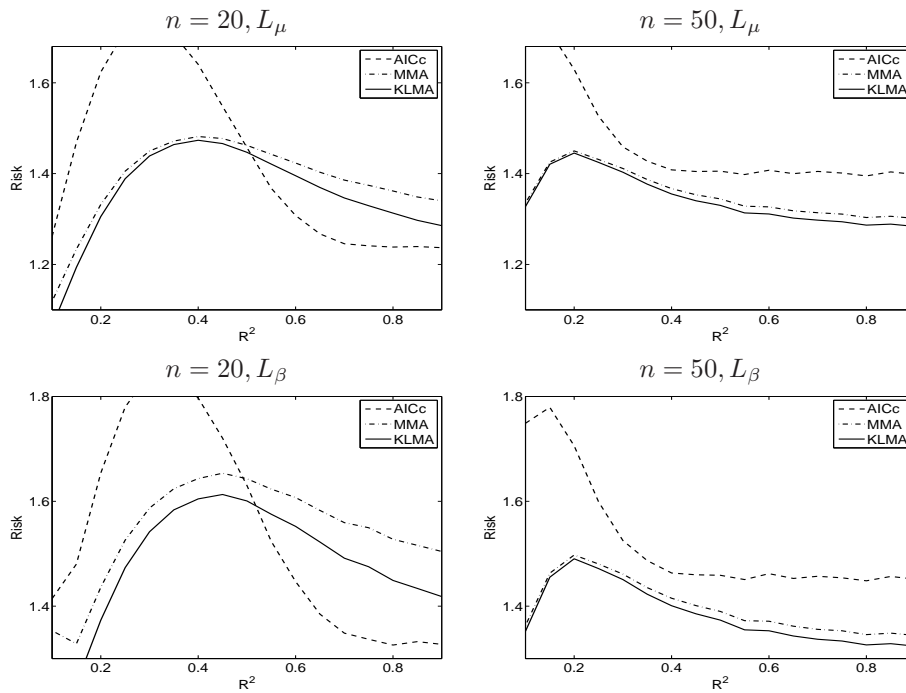


Figure S.6: Results for Example 1 with a uniformly distributed error term: risk comparisons under  $L_\mu$  and  $L_\beta$  as a function of  $R^2$ .

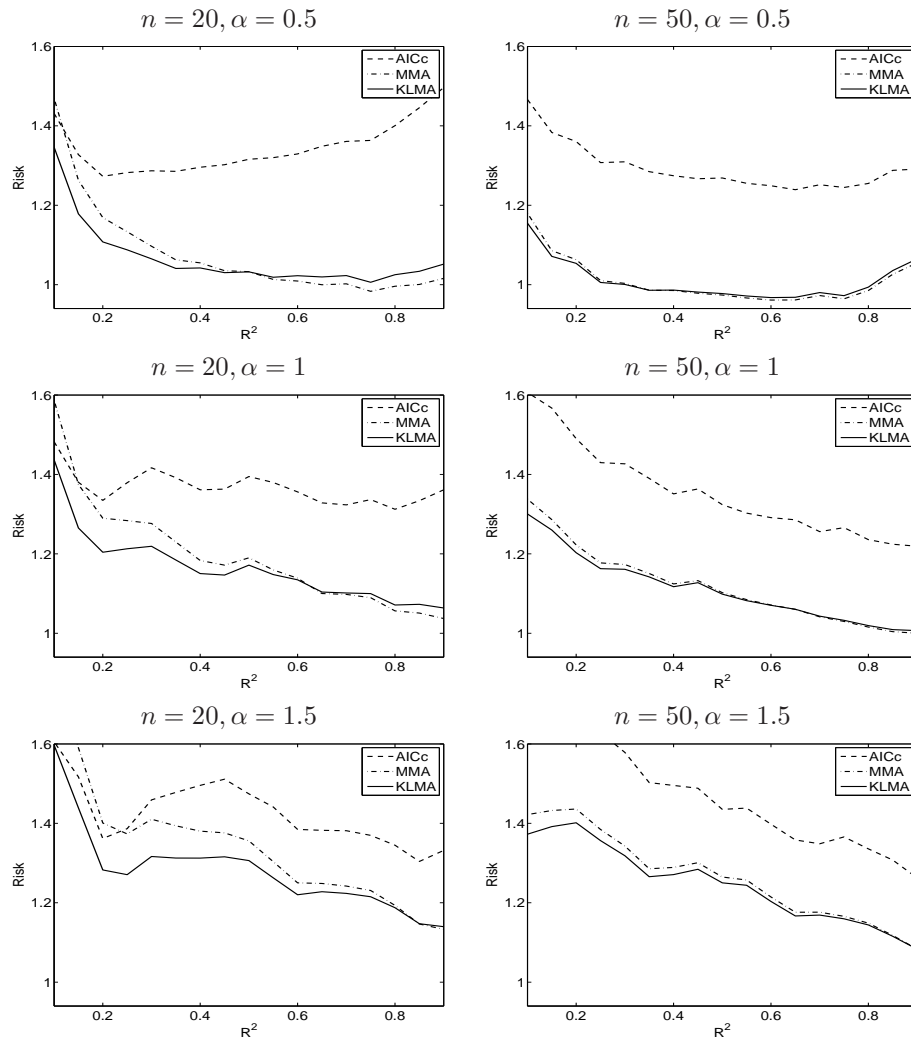


Figure S.7: Results for Example 2 with a uniformly distributed error term: risk comparisons under  $L_\mu$  as a function of  $R^2$ .

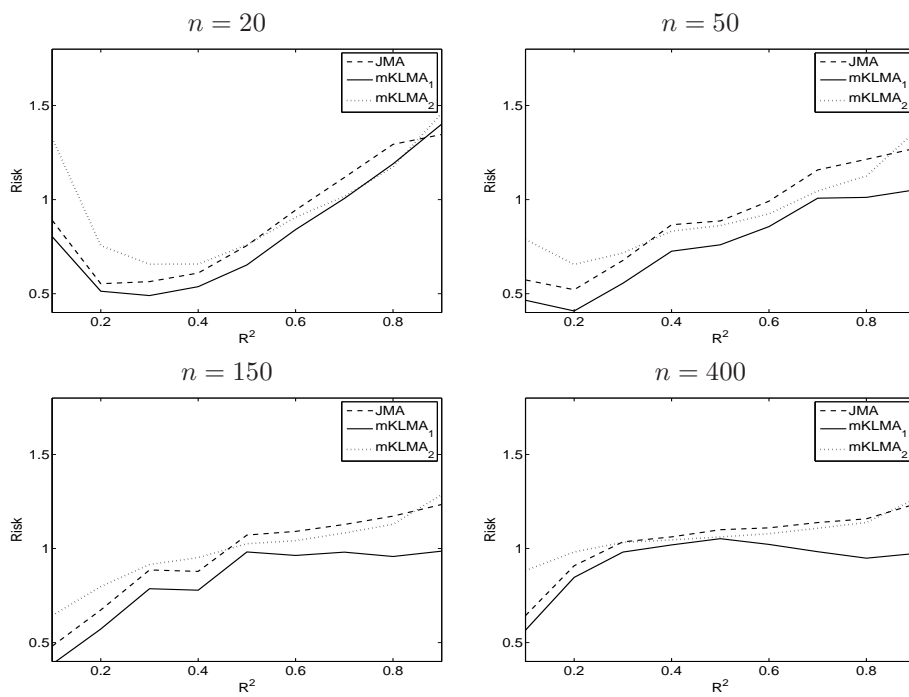


Figure S.8: Results for Example 3 with a uniformly distributed error term: risk comparisons under  $L_\mu$  as a function of  $R^2$ .

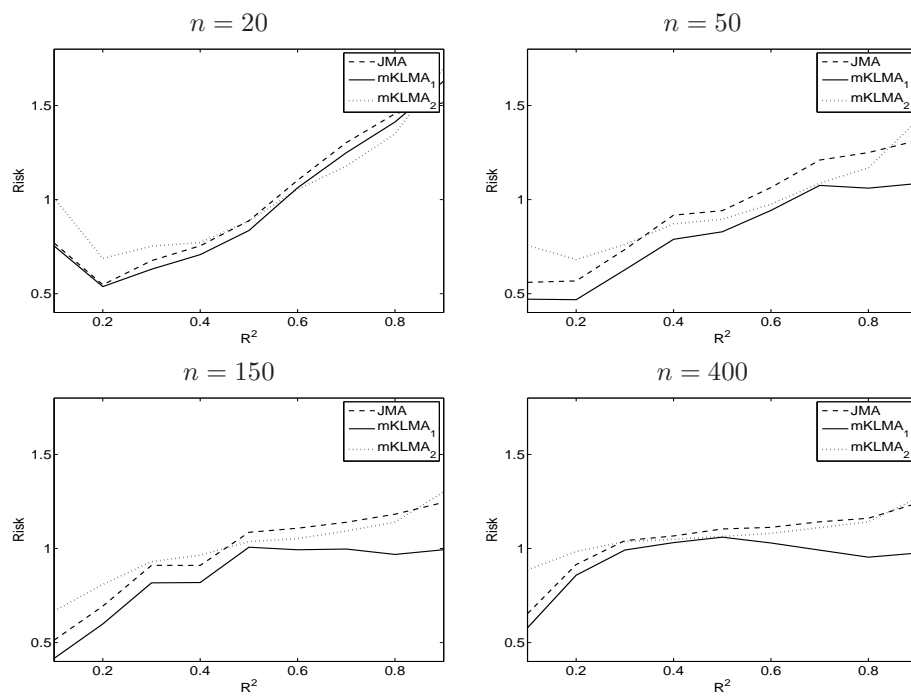


Figure S.9: Results for Example 3 with a uniformly distributed error term: risk comparisons under  $L_\beta$  as a function of  $R^2$ .



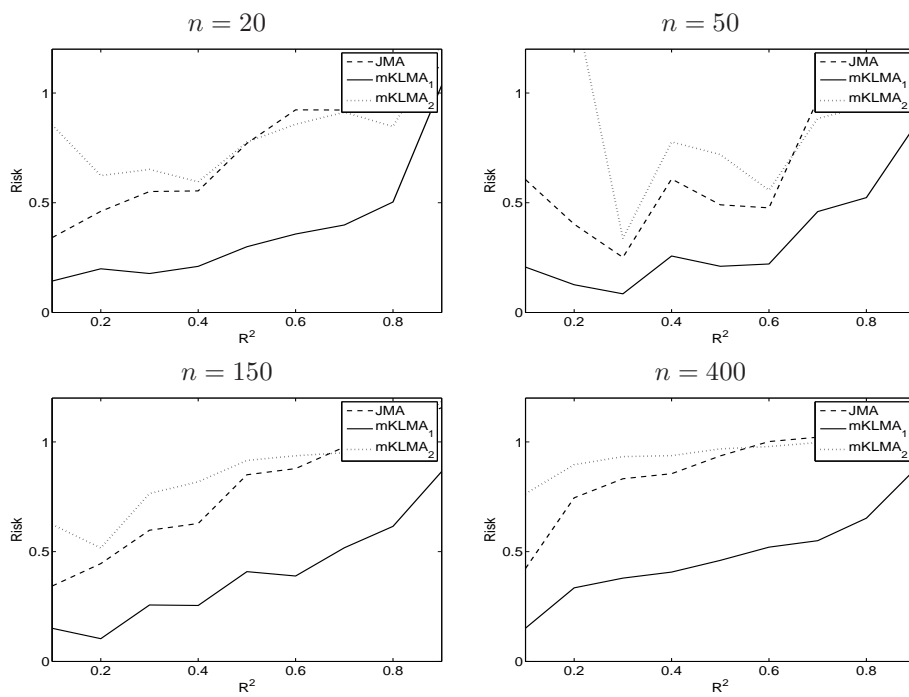


Figure S.10: Results for Example 3 with a uniformly distributed error term: risk comparisons under  $L_{\text{hetero}, \mu}$  as a function of  $R^2$ .

## S12 Simulation Results with a Chi-squared Distributed Error Term

In this section, we repeat Examples 1-3 but using a Chi-squared distributed error term:

| Example 1                           | Example 2                    | Example 3  |
|-------------------------------------|------------------------------|--|
| $\sigma 8^{-1/2} \{\chi^2(4) - 4\}$ | $8^{-1/2} \{\chi^2(4) - 4\}$ | $\exp(\eta X_{2i})^{1/2} 8^{-1/2} \{\chi^2(4) - 4\}$ |

The following figures show new results.

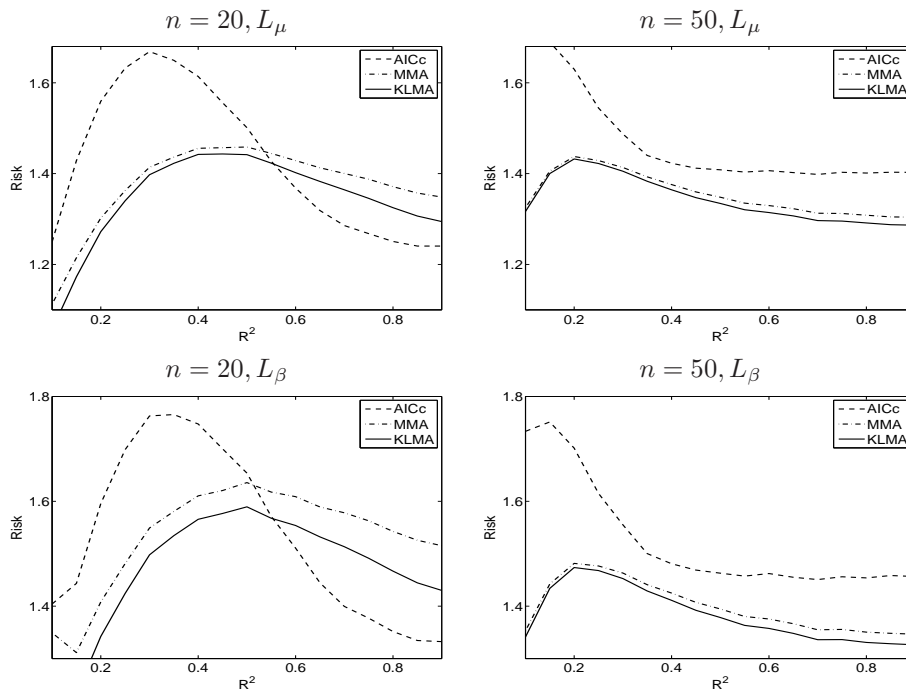


Figure S.11: Results for Example 1 with a Chi-squared distributed error term: risk comparisons under  $L_\mu$  and  $L_\beta$  as a function of  $R^2$ .

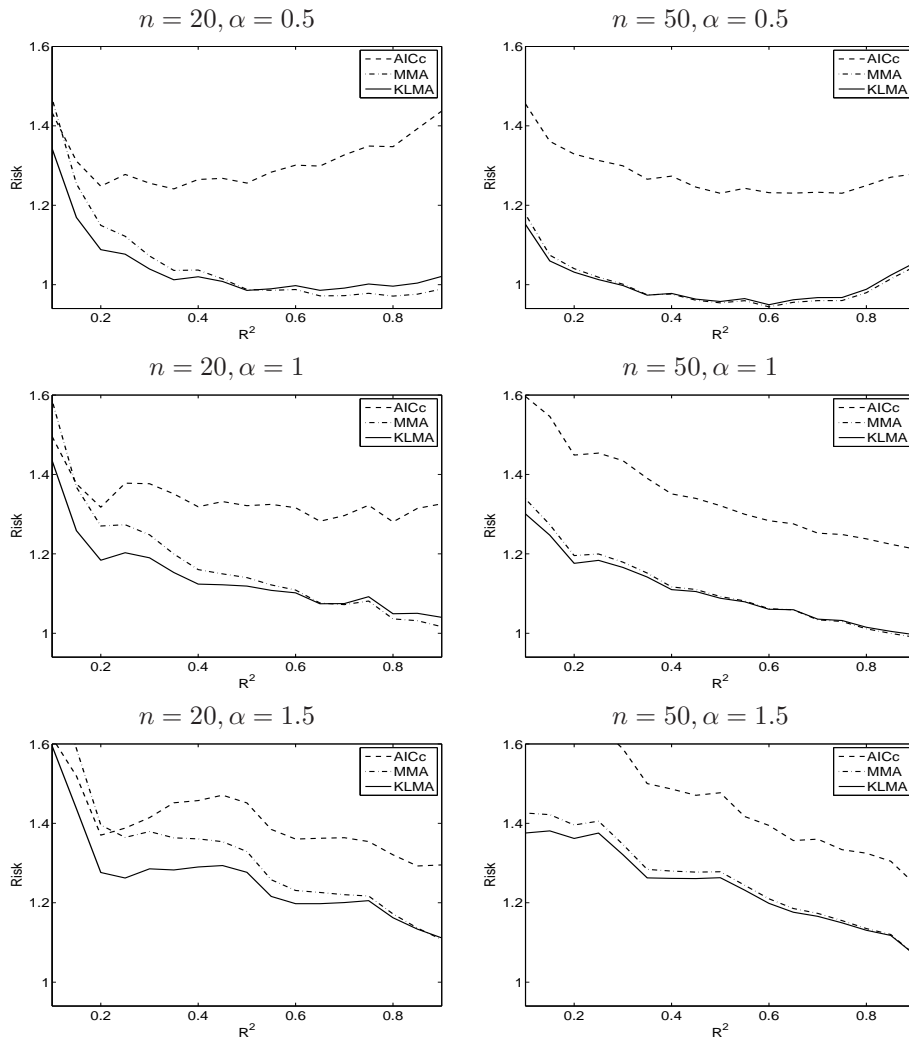


Figure S.12: Results for Example 2 with a Chi-squared distributed error term: risk comparisons under  $L_\mu$  as a function of  $R^2$ .

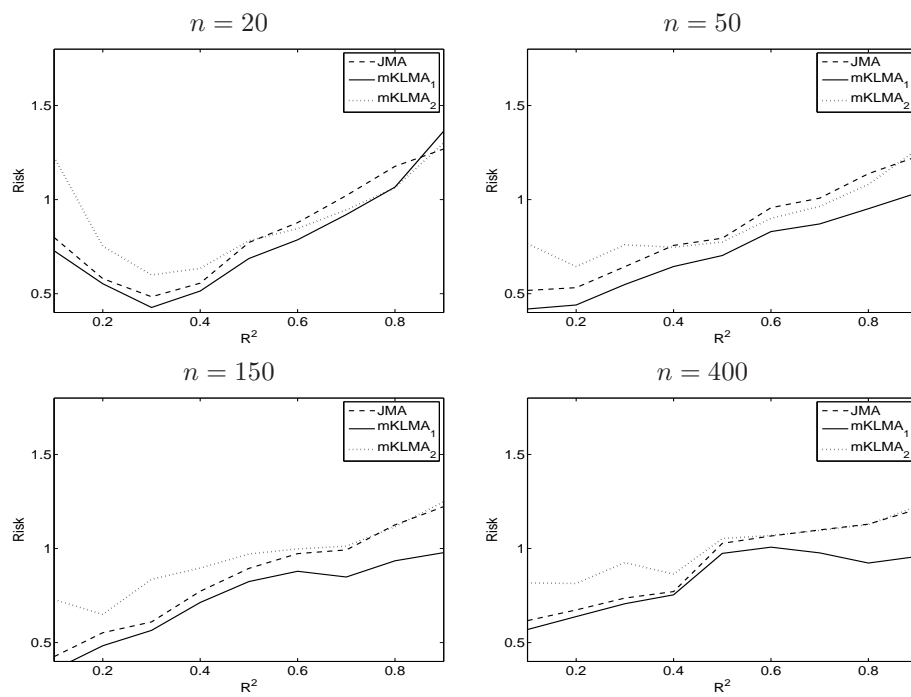


Figure S.13: Results for Example 3 with a Chi-squared distributed error term: risk comparisons under  $L_\mu$  as a function of  $R^2$ .

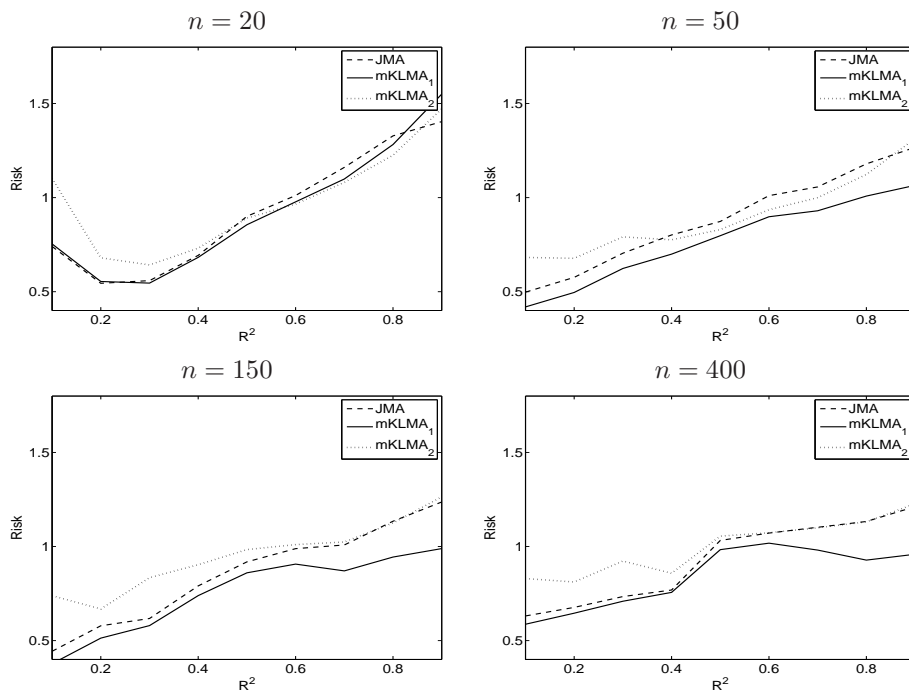


Figure S.14: Results for Example 3 with a Chi-squared distributed error term: risk comparisons under  $L_\beta$  as a function of  $R^2$ .

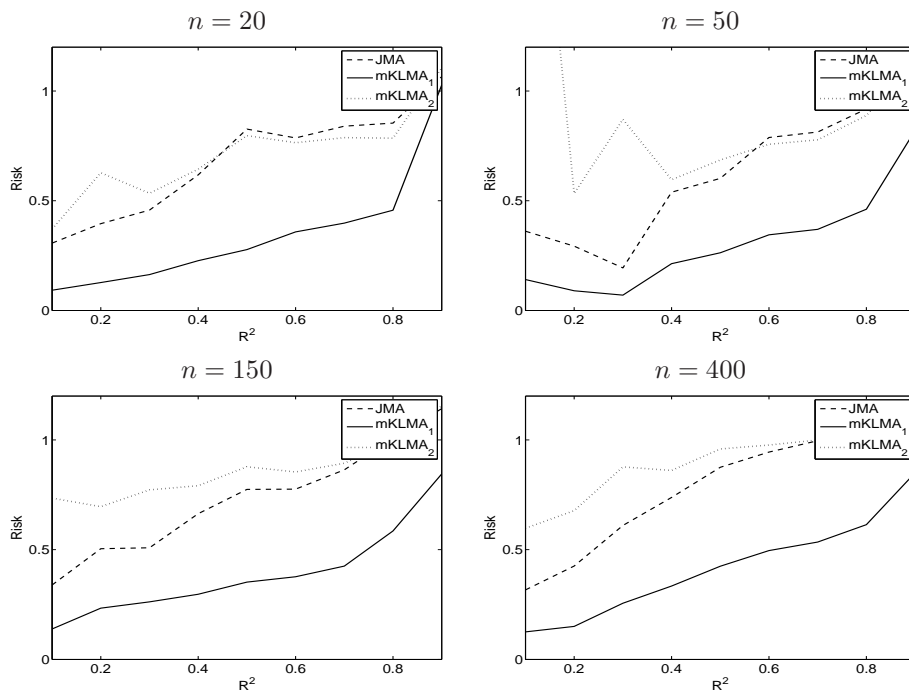


Figure S.15: Results for Example 3 with a Chi-squared distributed error term: risk comparisons under  $L_{\text{hetero}, \mu}$  as a function of  $R^2$ .

### S13 Simulation Results with Coefficients Depending on the Sample Size

In this section, we repeat Examples 1 and 3, but let the coefficients depend on the sample size. Specifically, we set  $\beta = (1, 2, 3, 2/\sqrt{n}, 2/\sqrt{n}, 2/\sqrt{n}, 2/\sqrt{n})^T$ . The following figures show new results.

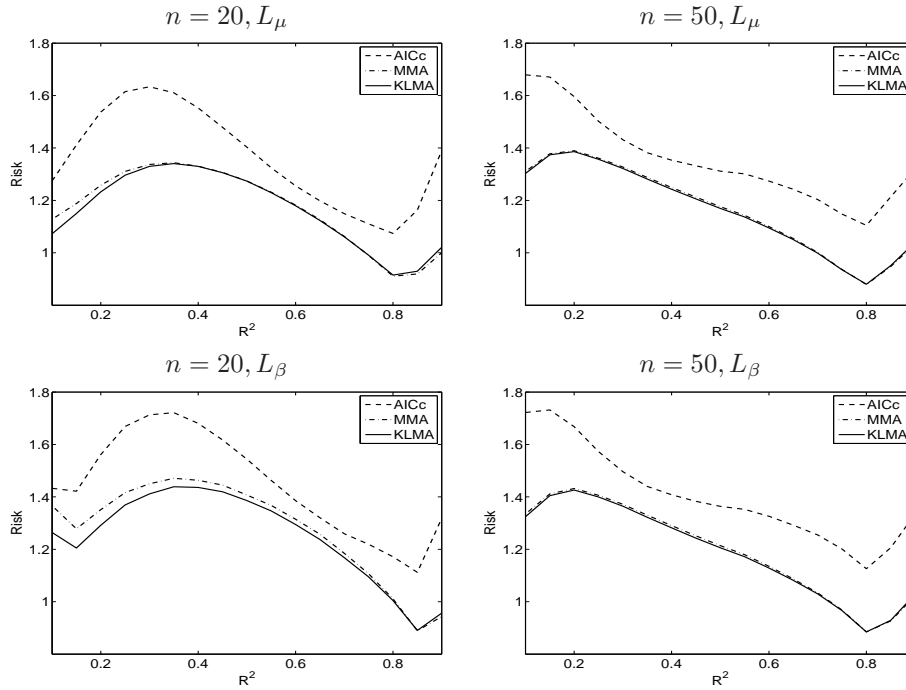


Figure S.16: Results for Example 1 with the coefficients depending on the sample size: risk comparisons under  $L_\mu$  and  $L_\beta$  as a function of  $R^2$ .

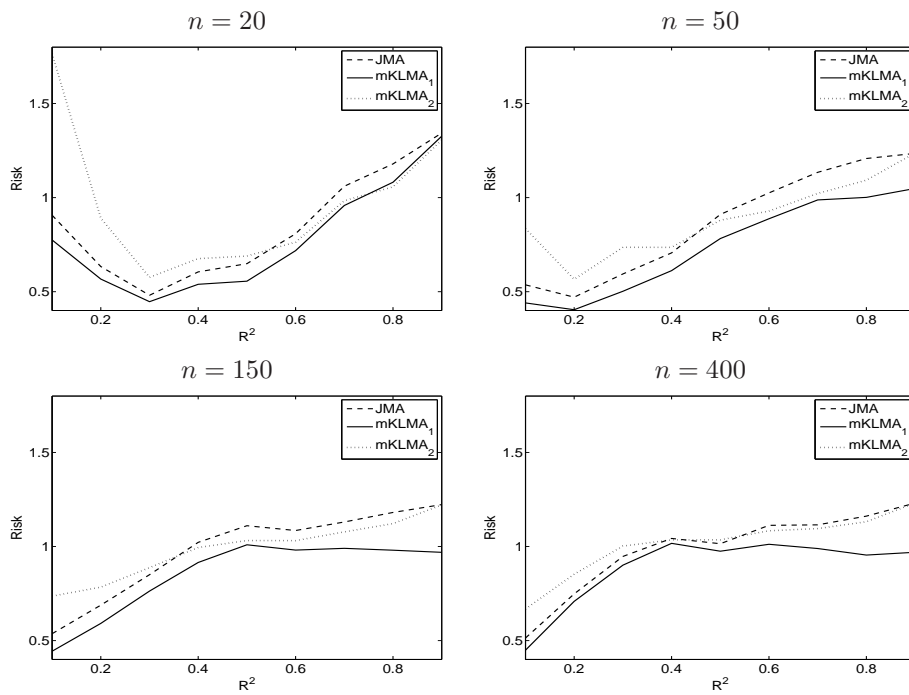


Figure S.17: Results for Example 3 with the coefficients depending on the sample size: risk comparisons under  $L_\mu$  as a function of  $R^2$ .



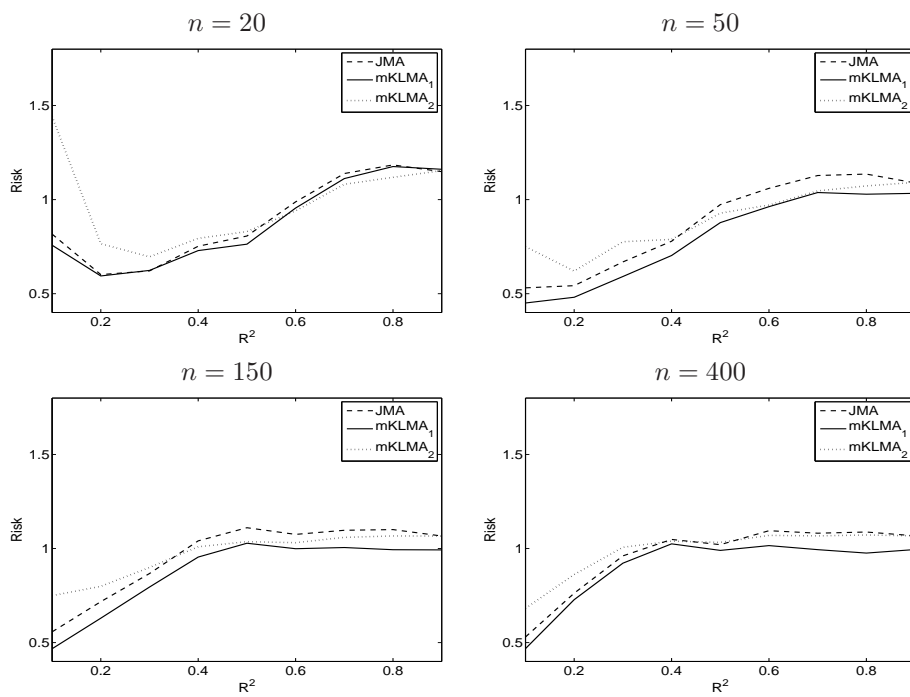


Figure S.18: Results for Example 3 with the coefficients depending on the sample size: risk comparisons under  $L_\beta$  as a function of  $R^2$ .

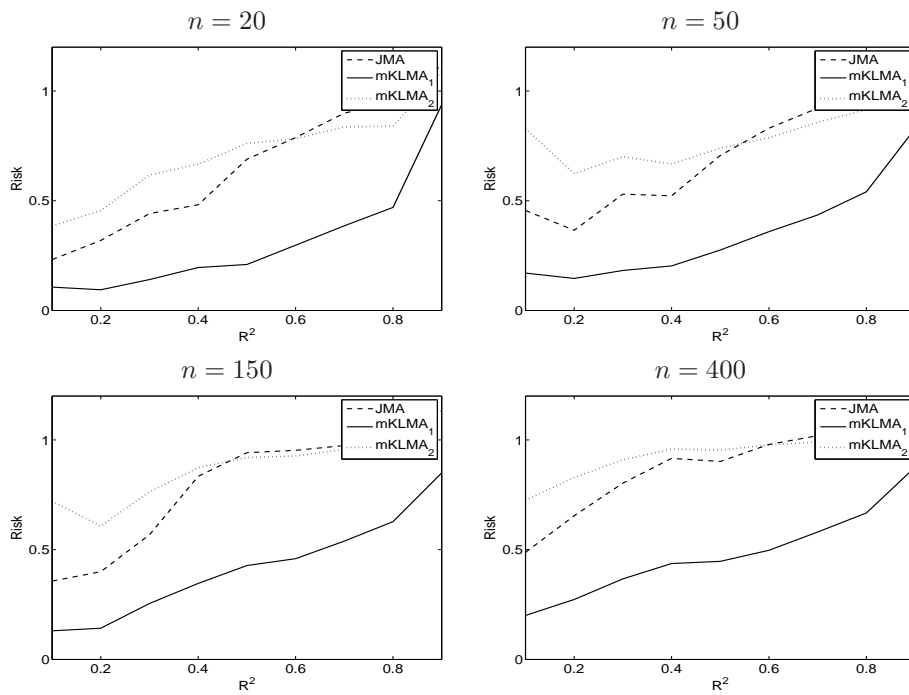


Figure S.19: Results for Example 3 with the coefficients depending on the sample size: risk comparisons under  $L_{\text{hetero}, \mu}$  as a function of  $R^2$ .

**References**

- Hansen, B. E. and Racine, J. (2012). Jackknife model averaging. *Journal of Econometrics* **167**, 38–46.
- Liu, Q. and Okui, R. (2013). Heteroskedasticity-robust  $C_p$  model averaging. *The Econometrics Journal* **16**, 463–472.
- Magnus, J. R. and Neudecker, H. (1988). *Matrix Differential Calculus with Applications in Statistics and Econometrics*. John Wiley & Sons, New York.
- Stein, C. M. (1981). Estimation of the mean of a multivariate normal distribution. *Annals of Statistics* **9**, 1135–1151.
- Wan, A. T. K., Zhang, X., and Zou, G. (2010). Least squares model averaging by Mallows criterion. *Journal of Econometrics* **156**, 277–283.