

1 LEAST SQUARES MODEL AVERAGING 1

2
3 BY BRUCE E. HANSEN¹ 3

4 This paper considers the problem of selection of weights for averaging across least
5 squares estimates obtained from a set of models. Existing model average methods are
6 based on exponential Akaike information criterion (AIC) and Bayesian information
7 criterion (BIC) weights. In distinction, this paper proposes selecting the weights by
8 minimizing a Mallows criterion, the latter an estimate of the average squared error
9 from the model average fit. We show that our new Mallows model average (MMA)
10 estimator is asymptotically optimal in the sense of achieving the lowest possible squared
11 error in a class of discrete model average estimators. In a simulation experiment we
12 show that the MMA estimator compares favorably with those based on AIC and BIC
13 weights. The proof of the main result is an application of the work of Li (1987). 13

14 KEYWORDS: ???.

15 1. INTRODUCTION 15

16 THIS PAPER DEVELOPS a new model averaging estimator for least squares re- 16
17 gression. A model average estimator is a weighted average of estimates ob- 17
18 tained from different models. The goal in model averaging is to reduce estima- 18
19 tion variance while controlling omitted variable bias. We propose a Mallows 19
20 criterion for the selection of the model weights, an estimate of the squared er- 20
21 ror. The empirical weights are found by numerical minimization of this crite- 21
22 rion. We show that this method of weight selection is asymptotically optimal in 22
23 the sense that the fitted estimates asymptotically achieve the minimum squared 23
24 error in a class of discrete model average estimators. 24
25

26 Model selection has a long history in statistics and econometrics, and differ- 26
27 ent methods have been advocated based on distinct estimation criteria, includ- 27
28 ing Akaike information criterion (AIC; Akaike (1973)), Mallows' C_p (Mal- 28
29 lows (1973)), Bayesian information criterion (BIC; Schwarz (1978)), delete- 29
30 one cross-validation (Stone (1974)), generalized cross-validation (Craven and 30
31 Wahba (1979)), and the focused information criterion (Claeskens and Hjort 31
32 (2003)). For generalized method of moments and empirical likelihood estima- 32
33 tion, analogous criteria have been proposed by Andrews and Lu (2001), Hong, 33
34 Preston, and Shum (2003), and Hall, Inoue, Jana, and Shin (2005). 34

35 Model averaging is an alternative to model selection. There is a large 35
36 Bayesian literature and a growing frequentist literature. Seminal contributions 36
37 to Bayesian model averaging include those by Draper (1995) and Raftery, 37
38 Madigan, and Hoeting (1997); for literature reviews, see Hoeting, Madigan, 38
39 Raftery, and Volinsky (1999) and Raftery and Zheng (2003). Some applica- 39
40 tions in econometrics include works by Doppelhofer, Miller, and Sala-i-Martin 40
41 (2004), Brock and Durlauf (2001), Avramov (2002), Fernandez, Ley, and Steel 41
42

43 ¹Research supported by the National Science Foundation. I gratefully thank the co-editor, 43
44 three referees, and Benedickt Potscher for helpful comments. 44

1 (2001a, 2001b), Garratt, Lee, Pesaran, and Shin (2003), Brock, Durlauf, and 1
2 West (2003), and Wright (2003a, 2003b). In the frequentist literature, Buck- 2
3 land, Burnham, and Augustin (1997) and Burnham and Anderson (2002) sug- 3
4 gested exponential AIC weights. The risk properties of a similar class of esti- 4
5 mators was examined by Leung and Barron (2004). Yang (2001) and Yuan and 5
6 Yang (2005) proposed a mixing estimator. Hjort and Claeskens (2003) pro- 6
7 vided an asymptotic analysis of model average estimators in likelihood-based 7
8 models. 8

9 Shrinkage and parameter penalization are other alternatives to model selec- 9
10 tion and averaging. Some recent contributions include the lasso-type estima- 10
11 tors of Knight and Fu (2000), the penalized likelihood estimators of Fan and Li 11
12 (2001) and Fan and Peng (2004), and the empirical Bayes estimator of Knox, 12
13 Stock, and Watson (2004). 13

14 There is also a large literature that discusses the effects of model selection 14
15 on inference. Potscher (1991) showed that AIC selection results in distorted in- 15
16 ference. Kabaila (1995) examined the impact on confidence regions. Buhlmann 16
17 (1999) presented conditions under which post-model-selection (PMS) estima- 17
18 tors are adaptive. Leeb and Potscher (2003, 2005a, 2005b) examined the un- 18
19 conditional and conditional distribution of PMS estimators and argued that 19
20 they cannot be uniformly estimated. 20

21 The approach we take in this paper is similar to that of selecting the number 21
22 of terms in a series expansion. Andrews (1991a) and Newey (1997) studied the 22
23 convergence rates for series estimators and give conditions for asymptotic nor- 23
24 mality, but did not give rules for selection. Shibata (1980, 1981, 1983) demon- 24
25 strated the asymptotic optimality of AIC selection in the context of Gaussian 25
26 regressions. Shibata's analysis was extended to non-Gaussian autoregressions 26
27 by Lee and Karagrigoriou (2001). Li (1987) demonstrated the asymptotic op- 27
28 timality of model selection in homoskedastic linear regression using Mallows' 28
29 criterion, cross-validation, and generalized cross-validation. Andrews (1991b) 29
30 extended Li's results to the case of heteroskedastic errors. A thorough review 30
31 of the asymptotic properties of model selection criteria has been provided by 31
32 Shao (1997). The optimality criterion used in these papers was critiqued by 32
33 Kabaila (2002). 33

34 We propose a model average estimator with weights selected by minimizing 34
35 a Mallows criterion. Our main contribution is a demonstration that the Mal- 35
36 lows criterion is asymptotically equivalent to the squared error, and thus our 36
37 Mallows model average (MMA) estimator asymptotically achieves the lowest 37
38 possible squared error in the class of model average estimators. Our proof is 38
39 an application of Theorem 2.1 of Li (1987). 39

40 There are two important limitations of our results. First, we restrict atten- 40
41 tion to regressions with conditionally homoskedastic errors. Andrews (1991a, 41
42 1991b) showed that model selection by Mallows' criterion is not optimal un- 42
43 der heteroskedasticity. The optimality of MMA will similarly fail under het- 43
44 eroskedasticity. Second, our asymptotic theory restricts the model average 44

weights to a discrete set due to the difficulty of establishing uniformity over a weight vector whose dimension is unbounded. Developing weight selection methods that allow for heteroskedasticity and extending the proof technique to allow for continuous weights are important topics for future research.

Section 2 discusses the estimation framework and model average estimators. Section 3 calculates the average squared error of the model average estimator. Section 4 introduces the Mallows criterion for the model average estimator and its sampling properties. Section 5 presents simulation evidence in support of the new MMA estimator. Proofs of the results are presented in the Appendix. A Gauss program that calculates the MMA estimator is available on the author's webpage, www.ssc.wisc.edu/~bhansen.

2. MODEL AVERAGE ESTIMATOR

Let $(y_i, x_i) : i = 1, \dots, n$ be a random sample, where y_i is real-valued while $x_i = (x_{1i}, x_{2i}, \dots)$ is countably infinite. The model is the homoskedastic linear regression

$$(1) \quad y_i = \mu_i + e_i,$$

$$(2) \quad \mu_i = \sum_{j=1}^{\infty} \theta_j x_{ji},$$

$$(3) \quad E(e_i | x_i) = 0,$$

$$(4) \quad E(e_i^2 | x_i) = \sigma^2.$$

We assume $E\mu_i^2 < \infty$ and that (2) converges in mean square. The linearity of (2) is not essential to the idea model averaging, but it greatly simplifies the algebraic calculations. Because the elements of x_i may be terms in a series expansion, (2) includes nonparametric regression.

Consider a sequence of approximating models $m = 1, 2, \dots$, where the m th model uses the first k_m elements of x_i , where $0 \leq k_1 < k_2 < \dots$. The m th approximating model is

$$(5) \quad y_i = \sum_{j=1}^{k_m} \theta_j x_{ji} + b_{mi} + e_i,$$

where the approximation error is $b_{mi} = \sum_{j=k_m+1}^{\infty} \theta_j x_{ji}$. In matrix notation, $Y = X_m \Theta_m + b_m + e$, where $Y = (y_1, \dots, y_n)'$, X_m is the $n \times k_m$ matrix with ij th element x_{ji} , $\Theta_m = (\theta_1, \dots, \theta_{k_m})'$, $b_m = (b_{m1}, \dots, b_{mn})'$, and $e = (e_1, \dots, e_n)'$.

Lurking behind (5) is an explicit ordering of the regressors x_{ji} . In some cases (such as a series expansion) this may not be troubling, but in other cases a natural ordering of the regressors may not be obvious. In practice, it may be

feasible to order the regressors by groups, and this may be a common application of model averaging.

Let $M = M_n \leq n$ be an integer for which $X'_{k_M} X_{k_M}$ is invertible. For all $m \leq M$, the least squares estimate of Θ_m is $\hat{\Theta}_m = (X'_m X_m)^{-1} X'_m Y$. Let $W = (w_1, \dots, w_M)'$ be a weight vector in the unit simplex in \mathbb{R}^M :

$$(6) \quad \mathcal{H}_n = \left\{ W \in [0, 1]^M : \sum_{m=1}^M w_m = 1 \right\}.$$

A model average estimator of Θ_M is

$$(7) \quad \hat{\Theta} = \sum_{m=1}^M w_m \begin{pmatrix} \hat{\Theta}_m \\ 0 \end{pmatrix}.$$

A model average estimator bears some resemblance to a shrinkage estimator. This can be seen most plainly when the regressors are orthogonal. In this case, the j th element of the model average estimator $\hat{\Theta}$ is the j th element of the unconstrained estimator $\hat{\Theta}_M$ multiplied by $\sum_{m=j}^M w_m$. Thus the coefficient estimates shrink toward zero, with the degree of shrinkage increasing with j . However, in the standard case where the regressors are not orthogonal, such a simple representation is not possible.

In the m th approximating model (5), let $\mu_m = X_m \Theta_m$ so that $\mu = \mu_m + b_m$, where $\mu = (\mu_1, \dots, \mu_n)'$. The estimate of μ in the m th approximating model is $\hat{\mu}_m = X_m \hat{\Theta}_m = P_m Y$, where $P_m = X_m (X'_m X_m)^{-1} X'_m$. The model average estimate of μ is $\hat{\mu}(W) = X_M \hat{\Theta} = P(W) Y$, where $P(W) = \sum_{m=1}^M w_m P_m$ is the implied "hat" matrix.

Because the matrix $P(W)$ plays an important role in the algebraic structure of the model average estimator, we discuss here some of its properties. Note that $P(W)$ is symmetric but generally not idempotent. Let $\lambda_{\max}(A)$ denote the largest eigenvalue of A and define

$$(8) \quad \Gamma_M = \begin{bmatrix} k_1 & k_1 & k_1 & \cdots & k_1 \\ k_1 & k_2 & k_2 & \cdots & k_2 \\ k_1 & k_2 & k_3 & \cdots & k_3 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ k_1 & k_2 & k_3 & \cdots & k_M \end{bmatrix}.$$

LEMMA 1: *We have:*

- (i) $\text{tr}(P(W)) = \sum_{m=1}^M w_m k_m \equiv k(W)$;
- (ii) $\text{tr}(P(W)P(W)) = \sum_{m=1}^M \sum_{l=1}^M w_m w_l \min(k_l, k_m) = W' \Gamma_M W$;
- (iii) $\lambda_{\max}(P(W)) \leq 1$.

3. SQUARED ERROR

Define the average squared error $L_n(W) = (\hat{\mu}(W) - \mu)'(\hat{\mu}(W) - \mu)$ and conditional squared error $R_n(W) = E(L_n(W)|X)$, where $X = \{x_1, \dots, x_n\}$.

LEMMA 2: *We have*

$$(9) \quad R_n(W) = W'(A_n + \sigma^2 \Gamma_M)W,$$

where Γ_M is defined in (8),

$$(10) \quad A_n = \begin{bmatrix} a_1 & a_2 & a_3 & \cdots & a_M \\ a_2 & a_2 & a_3 & \cdots & a_M \\ a_3 & a_3 & a_3 & \cdots & a_M \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ a_M & a_M & a_M & \cdots & a_M \end{bmatrix},$$

and $a_m = b'_m(I - P_m)b_m$. Furthermore, $a_l \geq a_m$ for $l \leq m$, and $A_n + \sigma^2 \Gamma_M > 0$ if $a_1 > 0$.

Lemma 2 shows that the conditional squared error $R_n(W)$ is a quadratic function in the weight vector W , an ellipsoid in \mathbb{R}^M centered at the zero vector. It is interesting to observe that the optimal weight vector W , which minimizes $R_n(W)$, necessarily puts non-zero weight on at least two models, except in the special case that $a_1 = a_M$. To see this, suppose that $M = 2$, in which case $R_n(W)$ is uniquely minimized by $w_1 = (1 + (a_1 - a_2)/\sigma^2(k_2 - k_1))^{-1}$, which is in $(0, 1)$ unless $a_1 = a_2$.

4. THE MALLOW'S CRITERION

The Mallows criterion for the model average estimator is

$$(11) \quad C_n(W) = (Y - X_M \hat{\Theta})'(Y - X_M \hat{\Theta}) + 2\sigma^2 k(W),$$

where $k(W)$ defined in Lemma 1 is the effective number of parameters. Definition (11) depends on the unknown σ^2 . We discuss below the replacement of σ^2 with an estimate.

The Mallows criterion may be used to select the weight vector W . Define

$$(12) \quad \hat{W} = \arg \min_{W \in \mathcal{H}_n} C_n(W),$$

the empirical Mallows selected weight vector. Because there is no closed-form solution to (12), the weight vector must be found numerically. For this calculation, it is convenient to write (11) in the following form. Let \hat{e}_m be the $n \times 1$ residual vector from the m th model, let $\bar{e} = (\hat{e}_1, \dots, \hat{e}_M)$ be the $n \times M$ matrix

1 collection of these residuals, and let $K = (k_1, \dots, k_M)'$ be the $M \times 1$ vector of
2 the number of parameters in the M models. Then (11) equals

$$3 \quad (13) \quad C_n(W) = W' \bar{e}' \bar{e} W + 2\sigma^2 K' W, \quad 4$$

5 which is linear-quadratic in W . The solution (12) minimizes (13) subject to
6 the nonnegativity and summation constraints (6). This is a classic quadratic
7 programming problem for which numerical algorithms are readily available.
8 (For example, in the Gauss programming language, the procedure QPROG is
9 appropriate.) The solution may be a unit vector or an interior value. If M is
10 moderately large, a typical solution \hat{W} can put zero weight on many of the in-
11 dividual models. The Mallows model average estimator is (7) using the weight
12 vector \hat{W} .

13 We present two justifications for the Mallows criterion. Our first is the classic
14 observation that $C_n(W)$ is an unbiased estimate of the expected squared error
15 plus a constant.

16 LEMMA 3: *We have*

$$17 \quad (14) \quad EC_n(W) = EL_n(W) + n\sigma^2. \quad 18$$

19 Our second justification is that if the weights are restricted to a discrete set,
20 the empirical Mallows weight vector asymptotically minimizes the squared er-
21 ror. Specifically, for some integer N , let the weights w_m be restricted to the
22 set $\{0, \frac{1}{N}, \frac{2}{N}, \dots, 1\}$ and let $\mathcal{H}_n(N)$ be the subset of \mathcal{H}_n restricted to this set of
23 weights. Let

$$24 \quad \hat{W}_N = \arg \min_{W \in \mathcal{H}_n(N)} C_n(W) \quad 25$$

26 be the Mallows weight vector, the choice obtained by minimizing the Mallows
27 criterion over the discrete weight set $\mathcal{H}_n(N)$.

28 The following result is an application of Theorem 2.1 of Li (1987), who
29 showed the asymptotic optimality of Mallows' criterion for model selection.

30 THEOREM 1: *As $n \rightarrow \infty$, if*

$$31 \quad (15) \quad \xi_n = \inf_{W \in \mathcal{H}_n} R_n(W) \rightarrow \infty \quad 32$$

33 *almost surely and for some fixed integer $N < \infty$,*

$$34 \quad (16) \quad E(|e_i|^{4(N+1)} | x_i) \leq \kappa < \infty, \quad 35$$

36 *then*

$$37 \quad (17) \quad \frac{L_n(\hat{W}_N)}{\inf_{W \in \mathcal{H}_n(N)} L_n(W)} \rightarrow_p 1. \quad 38$$

1 Note that the theorem places no restriction on M , the largest model included 1
2 in the model average (other than the requirement that $X'_{k_M} X_{k_M}$ is invertible). 2
3 Thus M may be fixed as $n \rightarrow \infty$ or $M = M_n$ may diverge to infinity. 3

4 Theorem 1 shows that the squared error obtained using the Mallows weight 4
5 vector \hat{W}_N is asymptotically equivalent to the infeasible optimal weight vector. 5
6 This means that the MMA estimator is asymptotically optimal in the class of 6
7 model average estimators (7) where the weight vector W is restricted to the set 7
8 $\mathcal{H}_n(N)$. 8

9 The restriction of \mathcal{H}_n to $\mathcal{H}_n(N)$ can be made less binding by picking N large, 9
10 which can be done as long as the conditional moment bound (16) holds. This 10
11 restriction is imposed because the proof of (17) requires that $C_n(W)$ is 11
12 asymptotically equivalent to $L_n(W)$ uniformly over W . The trouble is that the 12
13 dimension of the set \mathcal{H}_n is unbounded when $M_n \rightarrow \infty$ as $n \rightarrow \infty$, rendering 13
14 conventional proof methods inapplicable. 14

15 Theorem 1 requires condition (15), which specifies that there is no finite 15
16 approximating model m for which the bias is zero. This assumption is conven- 16
17 tional for nonparametric regression. For example, if $\gamma_m \sim m^{-\alpha}$, then we have 17
18 the explicit rate $\xi_n \sim n^{1/(1+2\alpha)}$. If (15) fails, then MMA will not satisfy the opti- 18
19 mality (17). 19

20 In practice, σ^2 is unknown, so (11) needs to be computed with a sample esti- 20
21 mate. One choice is $\hat{\sigma}_K^2 = (n - K)^{-1} (Y - X_K \hat{\theta}_K)' (Y - X_K \hat{\theta}_K)$, where $k_K = K$ 21
22 corresponds to a “large” approximating model. Other estimators for σ^2 have 22
23 been proposed in the nonparametric regression literature. Lemma 3 continues 23
24 to hold if $\hat{\sigma}_K^2$ is unbiased for σ^2 , which holds if $b_K = 0$, so the K th approximat- 24
25 ing model has no bias. Theorem 1 holds as stated as long as $\hat{\sigma}_K^2$ is consistent 25
26 for σ^2 , which is valid as shown next. 26
27

28 THEOREM 2: *If $K \rightarrow \infty$ and $K/n \rightarrow 0$ as $n \rightarrow \infty$, then $\hat{\sigma}_K^2 \rightarrow_p \sigma^2$ as $n \rightarrow \infty$.* 28
29

30 5. FINITE SAMPLE INVESTIGATION 30
31

32 We now investigate the finite sample mean squared error of the our model 32
33 average estimator in a simple simulation experiment. The setting is the infinite- 33
34 order regression $y_i = \sum_{j=1}^{\infty} \theta_j x_{ji} + e_i$. We set $x_{1i} = 1$ to be the intercept; the 34
35 remaining x_{ji} are independent and identically distributed $N(0, 1)$. The error 35
36 e_i is $N(0, 1)$ and independent of x_i . (Other experiments, not reported, 36
37 showed that the results are not sensitive to alternative distributions for the 37
38 regressors and regression error.) The parameters are determined by the rule 38
39 $\theta_j = c\sqrt{2\alpha}j^{-\alpha-1/2}$. The population $R^2 = c^2/(1 + c^2)$ is controlled by the 39
40 parameter c . 40

41 The sample size is varied between $n = 50, 150, 400$, and $1,000$. The parame- 41
42 ter α is varied between $0.5, 1.0$, and 1.5 . The larger α implies that the coeffi- 42
43 cients θ_j decline more quickly with j . The number of models M is determined 43
44 by the rule $M = 3n^{1/3}$ (so $M = 11, 16, 22$, and 30 for the four sample sizes). 44

1 The coefficient c was selected to control the population R^2 to vary on a grid 1
2 between 0.1 and 0.9. 2

3 We consider five estimators: (1) AIC model selection (AIC), (2) Mallows' 3
4 model selection (Mallows), (3) smoothed AIC (S-AIC), (4) smoothed BIC 4
5 (S-BIC), and (5) Mallows' model averaging (MMA). The AIC criterion for 5
6 model Θ_m is $AIC_m = n \ln \hat{\sigma}_m^2 + 2m$. The AIC model selection estimator is $\hat{\Theta}_{\hat{m}}$, 6
7 where \hat{m} minimizes AIC_m . S-AIC was introduced by Buckland, Burnham, and 7
8 Augustin (1997) and embraced by Burnham and Anderson (2002) and Hjort 8
9 and Claeskens (2003). It is the least squares model average estimator (7) 9
10 with the weights $w_m = \exp(-\frac{1}{2}AIC_m) / \sum_{j=1}^M \exp(-\frac{1}{2}AIC_m)$. S-BIC is a simplified 10
11 form of Bayesian model averaging. It is the least squares model average es- 11
12 timator (7) with the weights $w_m = \exp(-\frac{1}{2}BIC_m) / \sum_{j=1}^M \exp(-\frac{1}{2}BIC_m)$, where 12
13 $BIC_m = n \ln \hat{\sigma}_m^2 + \ln(n)m$. 13
14 14

15 To evaluate the estimators, we compute the risk (expected squared error). 15
16 We do this by computing averages across 100,000 simulation draws. For each 16
17 parameterization, we normalize the risk by dividing by the risk of the infeasible 17
18 optimal least squares estimator (the risk of the best-fitting model m). 18

19 The risk calculations are displayed in Figures 1–3 for $\alpha = 0.5, 1.0,$ and $1.5,$ 19
20 respectively. In each figure, the four panels display sample sizes. In each panel, 20
21 risk (expected squared error) is displayed on the y axis and the population R^2 is 21
22 displayed on the x axis. The two dotted lines correspond to AIC and Mallows 22
23 selection. The dashed, dash-dotted, and solid lines correspond to S-AIC, S- 23
24 BIC, and MMA, respectively. 24

25 In each panel, the AIC and Mallows selection methods have quite similar 25
26 risk. The smoothed AIC estimator achieves a lower risk than AIC model selec- 26
27 tion, which is consistent with the findings in the earlier literature. The S-AIC 27
28 and MMA estimators are nearly equivalent for the case $\alpha = 1.5$ and large n ; 28
29 otherwise, MMA achieves a lower risk than S-AIC. In many cases, its normal- 29
30 ized risk is less than 1, meaning that it is lower than that of infeasible optimal 30
31 model selection. 31

32 It is also instructive to contrast the performance of the MMA and S-BIC 32
33 estimators. The MMA estimator achieves lower risk in most cases, but S-BIC 33
34 has lower risk when n and R^2 are small, and its relative performance improves 34
35 when α is large. In particular, S-BIC has much lower risk when $\alpha = 1.5$ and 35
36 $n = 50$. Their relative performance depends strongly on sample size, with the 36
37 S-BIC estimator showing increasing relative risk and the MMA showing de- 37
38 creasing relative risk, as n increases. In many cases, however, the risk of the 38
39 S-BIC estimator is quite poor relative to the other methods. 39
40 40

41 *Dept. of Economics, University of Wisconsin, 1180 Observatory Drive, Madison,* 41
42 *WI 53706, U.S.A.; bhansen@ssc.wisc.edu.* 42

43 43
44 *Manuscript received January, 2006; final revision received August, 2006.* 44

LEAST SQUARES MODEL AVERAGING

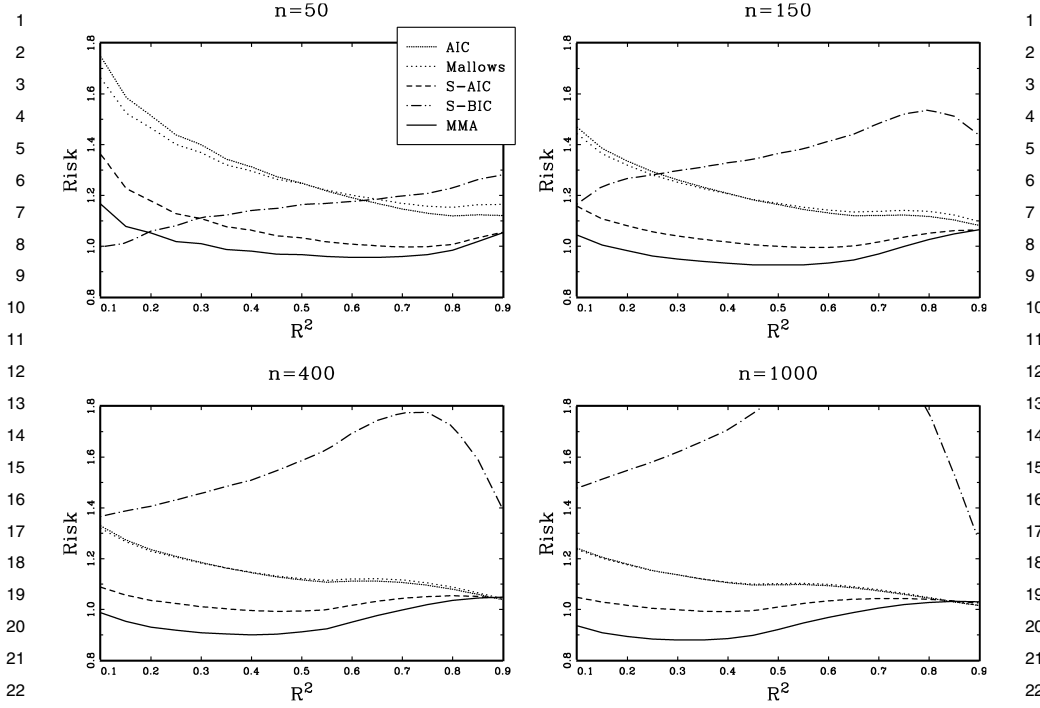


FIGURE 1.— $\alpha = 0.5$.

APPENDIX

PROOF OF LEMMA 1: Parts (i) and (ii) follow from the facts that $\text{tr}(P_m) = k_m$, $\text{tr}(P_m P_l) = \text{tr}(P_{\min(k_l, k_m)}) = \min(k_l, k_m)$, and simple algebra. Part (iii) uses the fact that P_m is idempotent so that

$$\lambda_{\max}(P(W)) = \max_{\eta} \frac{\eta' P(W) \eta}{\eta' \eta} \leq \sum_{m=1}^M w_m \max_{\eta} \frac{\eta' P_m \eta}{\eta' \eta} = 1. \quad Q.E.D.$$

PROOF OF LEMMA 2: Note that $\mu - \hat{\mu}(W) = (I - P(W))\mu - P(W)e$ and thus

$$(18) \quad L_n(W) = \mu'(I - P(W))(I - P(W))\mu - 2e'P(W)B_n W + e'P(W)P(W)e.$$

Lemma 1 and assumption (4) imply that

$$E(e'P(W)P(W)e|X) = \sigma^2 \text{tr}(P(W)P(W)) = \sigma^2 W' \Gamma_M W.$$

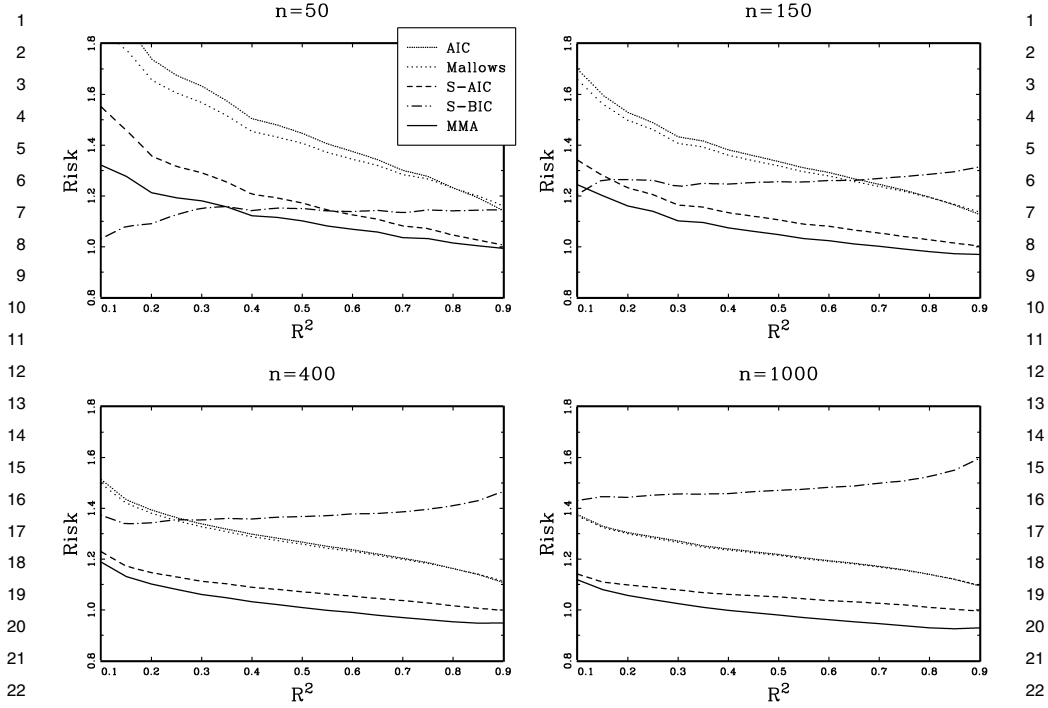


FIGURE 2.— $\alpha = 1.0$.

Taking conditional expectations of (18), we obtain

$$E(L_n(W)|X) = \mu'(I - P(W))(I - P(W))\mu + W'\sigma^2\Gamma_M W.$$

Define $b_m^* = (I - P_m)\mu = (I - P_m)b_m$ and $B_n = [b_1^*, \dots, b_M^*]$. Then

$$(19) \quad (I - P(W))\mu = \sum_{m=1}^M w_m b_m^* = B_n W.$$

Note that for $l \leq m$, $P_l P_m = P_l$ and $(I - P_m)b_l = (I - P_m)b_m$. Then

$$b_l^{*'} b_m^* = b_l'(I - P_l)(I - P_m)b_m = b_l'(I - P_m)b_m = b_m'(I - P_m)b_m = a_m$$

and thus $B_n' B_n = A_n$. It follows that $\mu'(I - P(W))(I - P(W))\mu = W' B_n' B_n W = W' A_n W$ and we obtain (9). Furthermore, for $l \leq m$ note that

$$\begin{aligned} b_m'(I - P_m)b_m &= b_l'(I - P_m)b_l = b_l'(I - P_l)b_l - b_l' P_m (I - P_l) P_m b_l \\ &\leq b_l'(I - P_l)b_l \end{aligned}$$

and thus $a_m \geq a_l$ as claimed.

LEAST SQUARES MODEL AVERAGING

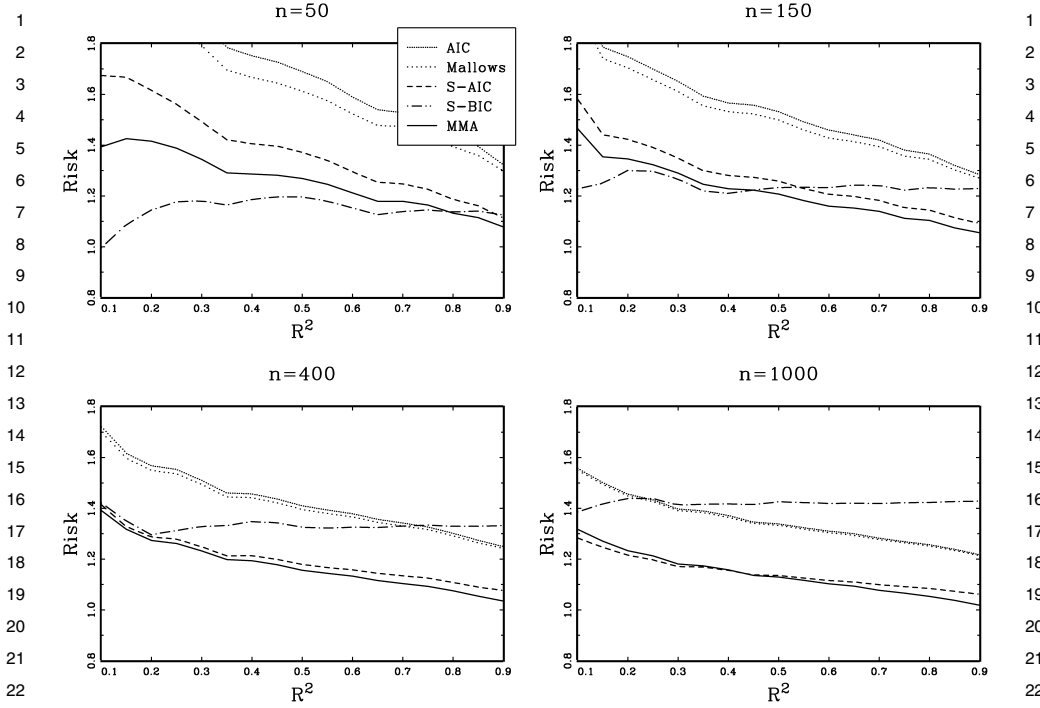


FIGURE 3.— $\alpha = 1.5$.

We now show that $A_n + \sigma^2 \Gamma_M > 0$, which holds if, for all $\alpha \neq 0$, $\alpha'(A_n + \sigma^2 \Gamma_M)\alpha > 0$. If $\alpha = \iota_1$ is the first unit vector, then $\alpha'(A_n + \sigma^2 \Gamma)\alpha = a_1^2 + \sigma^2 k_1^2 > 0$. Otherwise, if $\alpha \neq \iota_1$, note that $\alpha' A_n \alpha = \alpha' B_n' B_n \alpha \geq 0$, and by the definition of Γ and some algebraic manipulations,

$$\begin{aligned} \alpha' \Gamma_M \alpha &= k_1 \left(\sum_{m=1}^M \alpha_m \right)^2 + (k_2 - k_1) \left(\sum_{m=2}^M \alpha_m \right)^2 + \dots + (k_M - k_{M-1}) \alpha_M^2 \\ &> 0. \end{aligned}$$

Thus $\alpha'(A_n + \sigma^2 \Gamma_M)\alpha > 0$ as required. *Q.E.D.*

PROOF OF LEMMA 3: By straightforward algebra,

$$(20) \quad C_n(W) - L_n(W) = e'e + 2e'(I - P(W))\mu - 2(e'P(W)e - \sigma^2 k(W)).$$

Lemma 1 and assumption (4) imply that

$$(21) \quad E(e'P(W)e|X) = \sigma^2 \text{tr}(P(W)) = \sigma^2 k(W).$$

1 Taking expectations of (20), Equation (14) follows directly. Q.E.D. 1

2
3 PROOF OF THEOREM 1: Theorem 2.1 of Li (1987) established (17) for a 3
4 broad class of linear estimators. It is sufficient to verify that his equations (A.1), 4
5 (A.2), and (A.3) hold almost surely, conditional on X . Indeed, (A.1) is implied 5
6 by part (iii) of Lemma 1, and (A.2) holds by (16). It remains to show (A.3), 6
7 which in our notation is 7

$$9 \quad (22) \quad \sum_{W \in \mathcal{H}_n(N)} R_n(W)^{-(N+1)} \rightarrow 0 \quad 9$$

10
11 almost surely as $n \rightarrow \infty$. 11

12 For integers $1 \leq j_1 \leq j_2 \leq \dots \leq j_N$, let W_{j_1, j_2, \dots, j_N} be the weight vector that sets 12
13 $w_{j_l} = 1/N$ for $l = 1, \dots, N$, and the remainder zero. We can write 13

$$14 \quad \mathcal{H}_n(N) = \{W_{j_1, j_2, \dots, j_N} : 1 \leq j_1 \leq j_2 \leq \dots \leq j_N \leq M\}. \quad 14$$

15
16 The restriction of the weights to the form $1/N$ is without loss of generality, 16
17 because the weak ordering of the integers j_k allows ties. We then have 17

$$18 \quad (23) \quad \sum_{W \in \mathcal{H}_n(N)} R_n(W)^{-(N+1)} \leq \sum_{j_N=1}^{\infty} \sum_{j_{N-1}=1}^{j_N} \dots \sum_{j_1=1}^{j_2} R_n(W_{j_1, j_2, \dots, j_N})^{-(N+1)}. \quad 18$$

19
20 Now break the sum into two groups based on whether $k_{j_N} < \xi_n$ or $k_{j_N} \geq \xi_n$. For 20
21 the first group (which has less than ξ_n^N elements), use the bound $R_n(W) \geq \xi_n$ 21
22 from (15) and for the second group, use the simple bound 22

$$23 \quad R_n(W_{j_1, j_2, \dots, j_N}) \geq \sigma^2 W'_{j_1, j_2, \dots, j_N} \Gamma_M W_{j_1, j_2, \dots, j_N} \geq \frac{\sigma^2}{N^2} k_{j_N} \geq \frac{\sigma^2}{N^2} j_N, \quad 23$$

24
25 where the first inequality is implied by (9) and the second uses the definitions 25
26 of Γ_M and W_{j_1, j_2, \dots, j_N} . 26

27 Using these bounds, 27

$$28 \quad \begin{aligned} & \sum_{j_N=1}^{\infty} \sum_{j_{N-1}=1}^{j_N} \dots \sum_{j_1=1}^{j_2} R_n(W_{j_1, j_2, \dots, j_N})^{-(N+1)} \\ & \leq \xi_n^{-1} + \sum_{j_N=\xi_n}^{\infty} \sum_{j_{N-1}=1}^{j_N} \dots \sum_{j_1=1}^{j_2} \left(\frac{\sigma^2}{N^2} j_N \right)^{-(N+1)} \\ & \leq \xi_n^{-1} + \left(\frac{\sigma^2}{N^2} \right)^{-(N+1)} \sum_{j_N=\xi_n}^{\infty} j_N^{-2} \end{aligned} \quad 28$$

29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44

$$\begin{aligned} &\leq \xi_n^{-1} + \left(\frac{\sigma^2}{N^2}\right)^{-(N+1)} \xi_n^{-1} \\ &\rightarrow 0 \end{aligned}$$

almost surely as $n \rightarrow \infty$. Together with (23), this establishes (22) as desired. *Q.E.D.*

PROOF OF THEOREM 2: Since $\hat{e}_K = Y - X_K \hat{\Theta}_K = (I - P_K)e + (I - P_K)b_K$, we see that

$$(24) \quad \hat{\sigma}_K^2 = \frac{1}{n-K} e'(I - P_K)e + \frac{1}{n-K} b'_K(I - P_K)b_K + 2\frac{1}{n-K} e'(I - P_K)b_K.$$

We examine the terms on the right side of (24). First, because $Ee'(I - P_K)e = \sigma^2(n - K)$, by Theorem 2 of Whittle (1960),

$$\begin{aligned} E|e'(I - P_K)e - \sigma^2(n - K)|^2 &\leq C_2 \kappa^{1/(N+\delta)} \text{tr}((I - P_K)(I - P_K)) \\ &= C_2 \kappa^{1/(N+\delta)} (n - K). \end{aligned}$$

Thus for any $\delta > 0$, by Markov's inequality,

$$\begin{aligned} P\left(\left|\frac{1}{n-K} e'(I - P_K)e - \sigma^2\right| > \delta\right) &\leq \frac{E|e'(I - P_K)e - \sigma^2(n - K)|^2}{\delta^2(n - K)^2} \\ &\leq \frac{C_2 \kappa^{1/(N+\delta)}}{\delta^2(n - K)} \rightarrow 0 \end{aligned}$$

so $(n - K)^{-1} e'(I - P_K)e \rightarrow_p \sigma^2$. Second,

$$\frac{1}{n-K} E(b'_K(I - P_K)b_K) \leq \frac{n}{n-K} E b_{K_i}^2 \rightarrow 0$$

since $K \rightarrow \infty$ as $n \rightarrow \infty$ and the square integrability of μ_i implies $E b_{K_i}^2 \rightarrow 0$ as $K \rightarrow \infty$. This implies $(n - K)^{-1} b'_K(I - P_K)b_K \rightarrow_p 0$. Similarly, the third term on the right side of (24) is $o_p(1)$ and we conclude that $\hat{\sigma}_K^2 \rightarrow_p \sigma^2$. *Q.E.D.*

REFERENCES

AKAIKE, H. (1973): "Information Theory and an Extension of the Maximum Likelihood Principle," in *Second International Symposium on Information Theory*, ed. by B. Petroc and F. Csake. Budapest: Akademiai Kiado, ??-??. MR0483125[1]

ANDREWS, D. W. K. (1991a): "Asymptotic Normality of Series Estimators for Nonparametric and Semiparametric Regression Models," *Econometrica*, 59, 307-345. MR1097531[2]

_____ (1991b): "Asymptotic Optimality of Generalized C_L , Cross-Validation, and Generalized Cross-Validation in Regression with Heteroskedastic Errors," *Journal of Econometrics*, 47, 359-377. MR1097743[2]

<LS_link>

<LS_link>

- 1 ANDREWS, D. W. K., AND B. LU (2001): “Consistent Model and Moment Selection Procedures for
2 GMM Estimation with Application to Dynamic Panel Data Models,” *Journal of Econometrics*,
3 101, 123–164. [MR1805875](#)[1]
4 AVRAMOV, D. (2002): “Stock Return Predictability and Model Uncertainty,” *Journal of Finance*,
5 64, 423–458. [1]
6 BROCK, W., AND S. DURLAUF (2001): “Growth Empirics and Reality,” *World Bank Economic
7 Review*, 15, 229–272. [1]
8 BROCK, W., S. DURLAUF, AND K. D. WEST (2003): “Policy Analysis in Uncertain Economic En-
9 vironments,” *Brookings Papers on Economic Activity*, 1, 235–322. [2]
10 BUCKLAND, S. T., K. P. BURNHAM, AND N. H. AUGUSTIN (1997): “Model Selection: An Integral
11 Part of Inference,” *Biometrics*, 53, 603–618. [2,8]
12 BUHLMANN, P. (1999): “Efficient and Adaptive Post-Model-Selection Estimators,” *Journal of Sta-
13 tistical Planning and Inference*, 79, 1–9. [MR1704215](#)[2]
14 BURNHAM, K. P., AND D. R. ANDERSON (2002): *Model Selection and Multimodel Inference:
15 A Practical Information—Theoretic Approach*. Berlin: Springer-Verlag. [MR1919620](#)[2,8]
16 CLAESKENS, G., AND N. L. HJORT (2003): “The Focused Information Criterion,” *Journal of the
17 American Statistical Association*, 98, 900–916. [MR2041482](#)[1]
18 CRAVEN, P., AND G. WAHBA (1979): “Smoothing Noisy Data with Spline Functions: Estimating
19 the Correct Degree of Smoothing by the Method of Generalized Cross-Validation,” *Numerische
20 Mathematik*, 31, 377–403. [MR0516581](#)[1]
21 DOPPELHOFER, G., R. MILLER, AND X. SALA-I-MARTIN (2004): “Determinants of Long-Term
22 Growth: A Bayesian Averaging of Classical Estimates (BACE) Approach,” *American Eco-
23 nomic Review*, 94, 813–835. [1]
24 DRAPER, D. (1995): “Assessment and Propagation of Model Uncertainty,” *Journal of the Royal
25 Statistical Society, Ser. B*, 57, 45–70. [MR1325378](#)[1]
26 FAN, J., AND R. LI (2001): “Variable Selection via Nonconcave Penalized Likelihood and Its Or-
27 acle Properties,” *Journal of the American Statistical Association*, 96, 1348–1360. [MR1946581](#)[2]
28 FAN, J., AND H. PENG (2004): “Nonconcave Penalized Likelihood with a Diverging Number of
29 Parameters,” *The Annals of Statistics*, 32, 928–961. [MR2065194](#)[2]
30 FERNANDEZ, C., E. LEY, AND M. F. J. STEEL (2001a): “Benchmark Priors for Bayesian Model
31 Averaging,” *Journal of Econometrics*, 100, 381–427. [MR1820410](#)[2]
32 ——— (2001b): “Model Uncertainty in Cross-Country Growth Regressions,” *Journal of Applied
33 Econometrics*, 16, 563–576. [2]
34 GARRATT, A., K. LEE, M. H. PESARAN, AND Y. SHIN (2003): “Forecasting Uncertainties in
35 Macroeconomic Modelling: An Application to the UK Economy,” *Journal of the American
36 Statistical Association*, 98, 829–838. [MR2055491](#)[2]
37 HALL, A. R., A. INOUE, K. JANA, AND C. SHIN (2005): “Information in Generalized Method
38 of Moments Estimation and Entropy-Based Moment Selection,” *Journal of Econometrics*, ??,
39 ??–??. [1]
40 HJORT, N. L., AND G. CLAESKENS (2003): “Frequentist Model Average Estimators,” *Journal of
41 the American Statistical Association*, 98, 879–899. [MR2041481](#)[2,8]
42 HOETING, J. A., D. MADIGAN, A. E. RAFTERY, AND C. T. VOLINSKY (1999): “Bayesian Model
43 Averaging: A Tutorial,” *Statistical Science*, 14, 382–417. [MR1765176](#)[1]
44 HONG, H., B. PRESTON, AND M. SHUM (2003): “Generalized Empirical Likelihood-Based
Model Selection Criteria for Moment Condition Models,” *Econometric Theory*, 19, 923–943.
[MR2015971](#)[1]
KABAILA, P. (1995): “The Effect of Model Selection on Confidence Regions and Prediction Re-
gions,” *Econometric Theory*, 11, 537–549. [MR1349934](#)[2]
——— (2002): “On Variable Selection in Linear Regression,” *Econometric Theory*, 18, 913–925.
[MR1918328](#)[2]
KNIGHT, K., AND W. FU (2000): “Asymptotics for Lasso-Type Estimators,” *The Annals of Statistics*,
28, 1356–1378. [MR1805787](#)[2]

- 1 KNOX, T., J. H. STOCK, AND M. W. WATSON (2004): "Empirical Bayes Regression with Many
2 Regressors," Working Paper, ??? . [2]
- 3 LEE, S., AND A. KARAGRIGORIOU (2001): "An Asymptotically Optimal Selection of the Order of
4 a Linear Process," *Sankhyā*, Ser. A, 63, 93–106. [MR1898551\[2\]](#) <LS_link>
- 5 LEEB, H., AND B. M. PÖTSCHER (2003): "The Finite-Sample Distribution of Post-Model-
6 Selection Estimators and Uniform versus Non-Uniform Approximations," *Econometric Theory*,
7 19, 100–142. [MR1965844\[2\]](#) <LS_link>
- 8 ——— (2005a): "Model Selection and Inference: Facts and Fiction," *Econometric Theory*, 21,
9 21–59. [MR2153856\[2\]](#) <LS_link>
- 10 ——— (2005b): "Can One Estimate the Conditional Distribution of Post-Model-Selection Es-
11 timators?" *The Annals of Statistics*, ??, ??–??. [2] <LS_link>
- 12 LEUNG, G., AND A. R. BARRON (2004): "Information Theory and Mixing Least-Squares Regres-
13 sions," Working Paper, ??? . [2]
- 14 LI, K.-C. (1987): "Asymptotic Optimality for C_p , C_L , Cross-Validation and Generalized Cross-
15 Validation: Discrete Index Set," *The Annals of Statistics*, 15, 958–975. [MR0902239\[2,6,12\]](#) <LS_link>
- 16 MALLOWS, C. L. (1973): "Some Comments on C_p ," *Technometrics*, 15, 661–675. [1] <LS_link>
- 17 NEWBY, W. K. (1997): "Convergence Rates and Asymptotic Normality for Series Estimators,"
18 *Journal of Econometrics*, 79, 147–168. [MR1457700\[2\]](#) <LS_link>
- 19 POTSCHER, B. M. (1991): "Effects of Model Selection on Inference," *Econometric Theory*, 7,
20 163–185. [MR1128410\[2\]](#) <LS_link>
- 21 RAFTERY, A. E., D. MADIGAN, AND J. A. HOETING (1997): "Bayesian Model Averaging for Re-
22 gression Models," *Journal of the American Statistical Association*, 92, 179–191. [MR1436107\[1\]](#) <LS_link>
- 23 RAFTERY, A. E., AND Y. ZHENG (2003): "Long-Run Performance of Bayesian Model Averaging,"
24 Working Paper, University of Washington. [1]
- 25 SCHWARZ, G. (1978): "Estimating the Dimension of a Model," *The Annals of Statistics*, 6,
26 461–464. [MR0468014\[1\]](#) <LS_link>
- 27 SHAO, J. (1997): "An Asymptotic Theory for Linear Model Selection," *Statistica Sinica*, 7,
28 221–264. [MR1466682\[2\]](#) <LS_link>
- 29 SHIBATA, R. (1980): "Asymptotically Efficient Selection of the Order of the Model for Estimating
30 Parameters of a Linear Process," *The Annals of Statistics*, 8, 147–164. [MR0557560\[2\]](#) <LS_link>
- 31 ——— (1981): "An Optimal Selection of Regression Variables," *Biometrika*, 68, 45–54.
32 [MR0614940\[2\]](#) <LS_link>
- 33 ——— (1983): "Asymptotic Mean Efficiency of a Selection of Regression Variables," *Annals of*
34 *the Institute of Statistical Mathematics*, 35, 415–423. [MR0739383\[2\]](#) <LS_link>
- 35 STONE, M. (1974): "Cross-Validatory Choice and Assessment of Statistical Predictions," *Journal*
36 *of the Royal Statistical Society*, Ser. B, 36, 276–278. [MR0356377\[1\]](#) <LS_link>
- 37 WHITTLE, P. (1960): "Bounds for the Moments of Linear and Quadratic Forms in Independent
38 Variables," *Theory of Probability and Its Applications*, 5, 302–305. [MR0133849\[13\]](#) <LS_link>
- 39 WRIGHT, J. H. (2003a): "Bayesian Model Averaging and Exchange Rate Forecasts," *International*
40 *Finance Discussion Papers*, 779, Board of Governors, Federal Reserve Board. [2] <LS_link>
- 41 ——— (2003b): "Forecasting US Inflation by Bayesian Model Averaging," *International Finance*
42 *Discussion Papers*, 780, Board of Governors, Federal Reserve Board. [2] <LS_link>
- 43 YANG, Y. (2001): "Adaptive Regression by Mixing," *Journal of the American Statistical Association*,
44 96, 574–586. [MR1946426\[2\]](#) <LS_link>
- 45 YUAN, Z., AND Y. YANG (2005): "Combining Linear Regression Models: When and How?" *Jour-
46 nal of the American Statistical Association*, 100, 1202–1214. [MR2236435\[2\]](#) <LS_link>

- 1 KABAILA, P. (1995). The effect of model selection on confidence regions and prediction regions. 1
2 *Econometric Theory* **11** 537–549. MR1349934 (97a:62147) 2
- 3 KABAILA, P. (2002). On variable selection in linear regression. *Econometric Theory* **18** 913–925. 3
4 MR1918328 (2003d:62156) 4
- 5 KNIGHT, K. AND FU, W. (2000). Asymptotics for lasso-type estimators. *Ann. Statist.* **28** 1356–1378. 5
6 MR1805787 (2002a:62099) 6
- 7 Not Found! 7
- 8 LEE, S. AND KARAGRIGORIOU, A. (2001). An asymptotically optimal selection of the order of a 8
9 linear process. *Sankhyā Ser. A* **63** 93–106. MR1898551 (2003c:62150) 9
- 10 LEEB, H. AND PÖTSCHER, B. M. (2003). The finite-sample distribution of post-model-selection 10
11 estimators and uniform versus nonuniform approximations. *Econometric Theory* **19** 100–142. 11
12 MR1965844 (2004a:62045) 12
- 13 LEEB, H. AND PÖTSCHER, B. M. (2005). Model selection and inference: facts and fiction. *Econo-* 13
14 *metric Theory* **21** 21–59. MR2153856 14
- 15 Not Found! 15
- 16 Not Found! 16
- 17 LI, K.-C. (1987). Asymptotic optimality for C_p , C_L , cross-validation and generalized cross- 17
18 validation: discrete index set. *Ann. Statist.* **15** 958–975. MR902239 (89c:62112) 18
- 19 Not Found! 19
- 20 NEWBY, W. K. (1997). Convergence rates and asymptotic normality for series estimators. *J.* 20
21 *Econometrics* **79** 147–168. MR1457700 (98g:62033) 21
- 22 PÖTSCHER, B. M. (1991). Effects of model selection on inference. *Econometric Theory* **7** 163–185. 22
23 MR1128410 (92h:62048) 23
- 24 RAFTERY, A. E., MADIGAN, D., AND HOETING, J. A. (1997). Bayesian model averaging for linear 24
25 regression models. *J. Amer. Statist. Assoc.* **92** 179–191. MR1436107 25
- 26 Not Found! 26
- 27 SCHWARZ, G. (1978). Estimating the dimension of a model. *Ann. Statist.* **6** 461–464. MR0468014 27
28 (57 #7855) 28
- 29 SHAO, J. (1997). An asymptotic theory for linear model selection. *Statist. Sinica* **7** 221–264. With 29
30 comments and a rejoinder by the author. MR1466682 (99m:62104) 30
- 31 SHIBATA, R. (1980). Asymptotically efficient selection of the order of the model for estimating 31
32 parameters of a linear process. *Ann. Statist.* **8** 147–164. MR557560 (81c:62099) 32
- 33 SHIBATA, R. (1981). An optimal selection of regression variables. *Biometrika* **68** 45–54. 33
34 MR614940 (84a:62103a) 34
- 35 SHIBATA, R. (1983). Asymptotic mean efficiency of a selection of regression variables. *Ann. Inst.* 35
36 *Statist. Math.* **35** 415–423. MR739383 (86j:62158) 36
- 37 STONE, M. (1974). Cross-validators choice and assessment of statistical predictions. *J. Roy. Statist.* 37
38 *Soc. Ser. B* **36** 111–147. With discussion by G. A. Barnard, A. C. Atkinson, L. K. Chan, A. P. 38
39 Dawid, F. Downton, J. Dickey, A. G. Baker, O. Barndorff-Nielsen, D. R. Cox, S. Giesser, D. 39
40 Hinkley, R. R. Hocking, and A. S. Young, and with a reply by the authors. MR0356377 (50 40
41 #8847) 41
- 42 WHITTLE, P. (1960). Bounds for the moments of linear and quadratic forms in independent 42
43 variables. *Teor. Veroyatnost. i Primenen.* **5** 331–335. MR0133849 (24 #A3673) 43
- 44 Not Found! 44
- 45 Not Found! 45
- 46 YANG, Y. (2001). Adaptive regression by mixing. *J. Amer. Statist. Assoc.* **96** 574–588. MR1946426 46
47 (2003k:62127) 47
- 48 YUAN, Z. AND YANG, Y. (2005). Combining linear regression models: when and how? *J. Amer.* 48
49 *Statist. Assoc.* **100** 1202–1214. MR2236435 49
- 50 50
- 51 51
- 52 52
- 53 53
- 54 54