



AgEcon SEARCH

RESEARCH IN AGRICULTURAL & APPLIED ECONOMICS

The World's Largest Open Access Agricultural & Applied Economics Digital Library

This document is discoverable and free to researchers across the globe due to the work of AgEcon Search.

Help ensure our sustainability.

Give to AgEcon Search

AgEcon Search

<http://ageconsearch.umn.edu>

aesearch@umn.edu

*Papers downloaded from **AgEcon Search** may be used for non-commercial purposes and personal study only. No other use, including posting to another Internet site, is permitted without permission from the copyright owner (not AgEcon Search), or as allowed under the provisions of Fair Use, U.S. Copyright Act, Title 17 U.S.C.*

No endorsement of AgEcon Search or its fundraising activities by the author(s) of the following work or their employer(s) is intended or implied.



Queen's Economics Department Working Paper No. 1299

Least Squares Model Averaging by Prediction Criterion

Tian Xie

Department of Economics
Queen's University
94 University Avenue
Kingston, Ontario, Canada
K7L 3N6

11-2012

Least Squares Model Averaging by Prediction Criterion (Job Market Paper)

Tian Xie *

November 8, 2012

Abstract

This paper proposes a new estimator for least squares model averaging. A model average estimator is a weighted average of common estimates obtained from a set of models. We propose computing weights by minimizing a model average prediction criterion (MAPC). We prove that the MAPC estimator is asymptotically optimal in the sense of achieving the lowest possible mean squared error. For statistical inference, we derive asymptotic tests for single hypotheses and joint hypotheses on the average coefficients for the “core” regressors. These regressors are of primary interest to us and are included in every approximation model. To improve the finite sample performance, we also consider bootstrap tests. In simulation experiments the MAPC estimator is shown to have significant efficiency gains over existing model selection and model averaging methods. We also show that the bootstrap tests have more reasonable rejection frequency than the asymptotic tests in small samples. As an empirical illustration, we apply the MAPC estimator to cross-country economic growth models.

JEL classification: C52, C53, O40

Keywords: Model Averaging, MAPC, Convex Optimization, Optimality, Statistical Inference

*Queen’s University, Department of Economics, 94 University Ave., Kingston, Ontario, Canada K7L 3N6, xietian@econ.queensu.ca. I am extremely grateful to James G. MacKinnon and Morten Ørregaard Nielsen for their valuable supervision. I am also very grateful to Bruce E. Hansen, Jeffrey S. Racine, Joris Pinkse, Donald W. K. Andrews, and Joon Y. Park for their comments and suggestions. Helpful comments were received from participants at the Canadian Economics Association Conference held in Calgary (2012) and the Canadian Econometrics Study Group Annual Meeting in Kingston (2012).

1 Introduction

Economists formulate approximation models to capture the effects or factors supported by the empirical data. However, different approximation models usually yield different empirical results, which give rise to model uncertainty. There are two popular approaches for dealing with model uncertainty: model selection and model averaging.

Model selection is a procedure through which the best model is selected from a set of approximation models. This procedure generally involves calculating a criterion function for all of the approximation models and ranking them accordingly. One of the most widely used criterion functions is the Akaike information criterion (AIC) proposed by Akaike (1973). There are multiple versions of the AIC, the simplest of which is composed of a log-likelihood maximum and a penalty term. A popular alternative to AIC is the Bayesian information criterion (BIC) developed by Schwarz (1978). BIC is constructed similarly to AIC, but with a stronger penalty for complexity. There are other methods based on various criteria. Examples of these methods include the Mallows Criterion (Mallows' C_p) by Mallows (1973), the prediction criterion by Amemiya (1980), and the focused information criterion (FIC) by Claeskens and Hjort (2003).

Model averaging is an alternative to model selection. Instead of selecting a single “winning” model, model averaging calculates the weighted average of a set of approximation models. Barnard (1963) first mentioned the concept of “model combination” in a paper studying airline passenger data. Leamer (1978) proposed the basic paradigm for Bayesian model averaging (BMA). Buckland, Burnham and Augustin (1997) suggested using exponential AIC estimates as model weights and proposed the model average AIC (MA-AIC). There is an increasing focus on BMA in current literature (Draper (1995), Raftery, Madigan and Hoeting (1997), Kass and Raftery (1995), etc.); for a literature review on this topic, see Hoeting, Madigan, Raftery and Volinsky (1999).

However, applying BMA can sometimes be difficult due to the “prior set-up” required: before using the BMA approach, researchers need to assign prior probability distributions to the parameters of each model and prior probabilities to each model. The correct and efficient assignment of these prior values can be controversial in the field of economics, although some recent applied works may provide guidance: see Sala-i-Martin, Doppelhofer and Miller (2004), Ley and Steel (2009), Liu and Maheu (2009) and Wright (2009).

Least squares model averaging is an alternative to BMA. Hansen (2007) proposed the Mallows model average (MMA) method based on the original Mallows criterion. An implementation of MMA in forecast combination was made in Hansen (2008). Hansen (2009) extended MMA to regressions with a possible structural break. Another extension was made to autoregression with a near unit root in Hansen (2010). Most of these works are based on homoskedastic error terms. Hansen and Racine (2011) proposed a jackknife model averaging (JMA) that considers heteroskedastic error settings. One limitation of MMA in Hansen (2007) is that the approximation models must be strictly nested in a way that depends on

the ordering of regressors. In response to Hansen (2007), Wan, Zhang and Zou (2010) proved that the optimality of MMA holds for continuous model weights with non-nested models.

There are alternatives to model selection and model averaging. Fan and Li (2001) studied penalized likelihood estimators. Knox, Stock and Watson (2004) proposed the empirical Bayes estimator. There is also a considerable amount of literature that concentrates on general-to-specific (GETS) modeling. The foundations of GETS modeling have been developed over last several decades; see Hendry (1976, 1980, 1983), Gilbert (1986), Pagan (1987), Hoover and Perez (1999) and Hendry and Krolzig (1999). Campos, Ericsson and Hendry (2005) provide an overview of, and selected bibliography regarding, GETS modeling.

In this paper, we propose a model average estimator with empirical weights computed through numerical minimization of a model average prediction criterion (MAPC). Our criterion can be seen as a model averaging version of the original prediction criterion proposed by Amemiya (1980). We prove that the MAPC estimator is asymptotically optimal in the sense of achieving the lowest possible mean squared error, which applies both to nested and to non-nested approximation models. We divide the regressors into two groups: the “core” regressors, which are of primary interest to us and are included in every approximation model; and the “potential” regressors, which are of marginal interest to us and are included in some but not all approximation models. For statistical inference, we derive asymptotic tests for single hypotheses and joint hypotheses on the core average coefficients. To improve the finite sample performance, we also consider bootstrap tests. In simulation experiments the MAPC estimator is shown to have significant efficiency gains over existing model selection and model averaging methods. We also show that the bootstrap tests have more reasonable rejection frequency than the asymptotic tests in small samples. As an empirical illustration, we apply the MAPC estimator to the cross-country economic growth models in Barro (1991).

This paper continues with an introduction of the framework of MAPC in Section 2. Section 3 proves the asymptotic optimality of the MAPC estimator and derives asymptotic tests and bootstrap tests for single hypotheses and joint hypotheses on the average coefficients. Section 4 presents simulation experiments. In Section 5, we provide an empirical application, in which the MAPC estimator is applied to the economic growth models in Barro (1991). Section 6 concludes the paper. Proofs are presented in Appendix A and a description of the data set used in Section 5 can be found in Appendix B.

2 Model Averaging Prediction Criterion

Let $(y_i, \mathbf{x}_i) : i = 1, \dots, n$ be a random sample, where y_i and $\mathbf{x}_i = [x_{i1}, x_{i2}, \dots]$ are real-valued. We let \mathbf{x}_i be countably infinite. The same design can be seen in Hansen (2007) and Wan et al. (2010). We assume the data generating process is

$$y_i = \mu_i + u_i, \tag{1}$$

where $\mu_i = \sum_{j=1}^{\infty} \beta_j x_{ij}$, $\mathbb{E}(u_i | \mathbf{x}_i) = 0$ and $\mathbb{E}(u_i^2 | \mathbf{x}_i) = \sigma^2$.

We consider a sequence of linear approximation models $m = 1, 2, \dots, M$. The concept of “approximation model” can be vague. In this paper, an approximation model m uses $k^{(m)}$ regressors belonging to \mathbf{x}_i such that

$$y_i = \sum_{j=1}^{k^{(m)}} \beta_j^{(m)} x_{ij}^{(m)} + u_i^{(m)} \quad \text{for } i = 1, \dots, n, \quad (2)$$

where $\beta_j^{(m)}$ is a coefficient in model m and $x_{ij}^{(m)}$ is a regressor in model m . Other forms of approximation models are beyond the scope of this paper. The approximation error for model m is defined as

$$b_i^{(m)} \equiv \mu_i - \sum_{j=1}^{k^{(m)}} \beta_j^{(m)} x_{ij}^{(m)} \quad \text{for } i = 1, \dots, n. \quad (3)$$

Therefore, as long as the approximation model is finite, it always contains non-zero approximation error.

Hansen (2007) assumed that the regressors \mathbf{x}_i were an ordered set and an approximation model m contained the first $k^{(m)}$ regressors from \mathbf{x}_i . As a result, models with fewer regressors would always nest within larger models, which made the $k^{(m)}$ such that

$$0 \leq k^{(1)} < k^{(2)} < \dots < k^{(m)} < \dots < k^{(M)}.$$

This nested model set-up was demonstrated to be unnecessary by Wan et al. (2010) and Hansen and Racine (2011). In this paper, we place no such restrictions on the orders of \mathbf{x}_i and the $k^{(m)}$. Approximation models in our paper can be either nested or non-nested, which makes our method more widely applicable.

The DGP (1) and approximation model (2) can be represented in the following matrix forms:

$$\mathbf{y} = \boldsymbol{\mu} + \mathbf{u}$$

and

$$\mathbf{y} = \mathbf{X}^{(m)} \boldsymbol{\beta}^{(m)} + \mathbf{u}^{(m)},$$

where \mathbf{y} is $n \times 1$, $\boldsymbol{\mu}$ is $n \times 1$, $\mathbf{X}^{(m)}$ is $n \times k^{(m)}$ with the ij^{th} element being $x_{ij}^{(m)}$, $\boldsymbol{\beta}^{(m)}$ is $k^{(m)} \times 1$ and $\mathbf{u}^{(m)}$ is the error term for model m . Let \mathbf{P} stand for a projection matrix. For an approximation model m , we have

$$\mathbf{P}^{(m)} = \mathbf{X}^{(m)} \left(\mathbf{X}^{(m)\top} \mathbf{X}^{(m)} \right)^{-1} \mathbf{X}^{(m)\top}. \quad (4)$$

Therefore, the least squares estimate of $\boldsymbol{\mu}$ from model m is $\hat{\boldsymbol{\mu}}^{(m)} = \mathbf{P}^{(m)} \mathbf{y}$.

Let $\mathbf{w} = [w^{(1)}, \dots, w^{(M)}]^\top$ be a weight vector in the unit simplex in \mathbb{R}^M and define

$$\mathbf{H}_M \equiv \left\{ \mathbf{w} \in [0, 1]^M : \sum_{m=1}^M w^{(m)} = 1 \right\},$$

where \mathbf{H}_M is a continuous set. Note that Hansen (2007) and Hansen and Racine (2011) assumed a discrete set \mathbf{H}_M^* for \mathbf{w} , in which

$$\mathbf{H}_M^*(N) \equiv \left\{ w^{(m)} \in \left[0, \frac{1}{N}, \frac{2}{N}, \dots, 1 \right] : \sum_{m=1}^M w^{(m)} = 1 \right\}$$

for some fixed integer N . We will discuss this further in Section 3.

We define the model average estimator of $\boldsymbol{\mu}$ as

$$\boldsymbol{\mu}(\mathbf{w}) \equiv \sum_{m=1}^M w^{(m)} \hat{\boldsymbol{\mu}}^{(m)} = \sum_{m=1}^M w^{(m)} \mathbf{P}^{(m)} \mathbf{y}. \quad (5)$$

For the sake of simplicity, we define a weighted average projection matrix $\mathbf{P}(\mathbf{w})$ as

$$\mathbf{P}(\mathbf{w}) \equiv \sum_{m=1}^M w^{(m)} \mathbf{P}^{(m)}.$$

Accordingly, (5) can be simplified to

$$\boldsymbol{\mu}(\mathbf{w}) = \mathbf{P}(\mathbf{w}) \mathbf{y}. \quad (6)$$

The effective number of parameters $k(\mathbf{w})$ in model averaging estimation is defined as

$$k(\mathbf{w}) \equiv \sum_{m=1}^M w^{(m)} k^{(m)}, \quad (7)$$

which is a weighted sum of the $k^{(m)}$. Note that $k(\mathbf{w})$ is not necessarily an integer.

We propose the model average prediction criterion (MAPC):

$$\text{MAPC}_n(\mathbf{w}) = (\mathbf{y} - \boldsymbol{\mu}(\mathbf{w}))^\top (\mathbf{y} - \boldsymbol{\mu}(\mathbf{w})) \left(\frac{n + k(\mathbf{w})}{n - k(\mathbf{w})} \right), \quad (8)$$

where $\boldsymbol{\mu}(\mathbf{w})$ and $k(\mathbf{w})$ are defined in (6) and (7). MAPC can be understood as the model averaging version of the prediction criterion by Amemiya (1980). Like most model selection criteria and model averaging criteria, MAPC follows the idea of parsimony and balances between the fitness and the size of a model. MAPC can be used to calculate the empirical

weight vector $\hat{\mathbf{w}}$, in which

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbf{H}_M} \text{MAPC}_n(\mathbf{w}).$$

According to Hansen (2007), the Mallows' model average (MMA) criterion is

$$\text{MMA}_n(\mathbf{w}) = (\mathbf{y} - \boldsymbol{\mu}(\mathbf{w}))^\top (\mathbf{y} - \boldsymbol{\mu}(\mathbf{w})) + 2\sigma^2 k(\mathbf{w}). \quad (9)$$

The empirical weights $\hat{\mathbf{w}}$ can be selected by minimizing (9) subject to $\mathbf{w} \in \mathbf{H}_M^*$. The MMA criterion is composed of an averaged sum of squared residuals and a penalty term for complexity. Note that the penalty term includes an unknown σ^2 that must be replaced by a sample estimate.

For convenience in calculations, we rewrite both criteria. First, denote $\hat{\mathbf{u}}^{(m)}$ as an $n \times 1$ estimated residual vector from model m . Let $\hat{\mathbf{U}}$ be an $n \times M$ matrix consisting of these residuals such that $\hat{\mathbf{U}} \equiv [\hat{\mathbf{u}}^{(1)}, \hat{\mathbf{u}}^{(2)}, \dots, \hat{\mathbf{u}}^{(M)}]$. Define an $M \times 1$ vector \mathbf{k} which contains the number of parameters from each model such that $\mathbf{k} \equiv [k^{(1)}, k^{(2)}, \dots, k^{(M)}]^\top$. Then, the MAPC in (8) can be written as

$$\text{MAPC}_n(\mathbf{w}) = \mathbf{w}^\top \hat{\mathbf{U}}^\top \hat{\mathbf{U}} \mathbf{w} \left(\frac{n + \mathbf{k}^\top \mathbf{w}}{n - \mathbf{k}^\top \mathbf{w}} \right). \quad (10)$$

Likewise, the MMA criterion in (9) becomes

$$\text{MMA}_n(\mathbf{w}) = \mathbf{w}^\top \hat{\mathbf{U}}^\top \hat{\mathbf{U}} \mathbf{w} + 2\sigma^2 \mathbf{k}^\top \mathbf{w}. \quad (11)$$

3 Asymptotic Properties

In this section, we first prove the asymptotic optimality of the MAPC estimator by showing that it achieves the lowest possible mean squared error. Then, we derive asymptotic tests for single hypotheses and joint hypotheses on the average coefficients. We also recommend bootstrap tests for improved inference in finite samples. Proofs for this section are presented in Appendix A.

3.1 Asymptotic Optimality

We start by listing some properties of $\mathbf{P}(\mathbf{w})$.

Lemma 1 Define $\mathbf{M}(\mathbf{w}) \equiv \mathbf{I} - \mathbf{P}(\mathbf{w})$. We have

- (i) $\text{Tr}(\mathbf{P}(\mathbf{w})) = \sum_{m=1}^M w^{(m)} k^{(m)} = k(\mathbf{w})$,
- (ii) $\lambda_{\max}(\mathbf{P}(\mathbf{w})) \leq 1$, where $\lambda_{\max}(\cdot)$ returns the largest eigenvalue of its argument,

$$(iii) \quad \|\mathbf{P}(\mathbf{w}^*)\mathbf{M}(\mathbf{w})\boldsymbol{\mu}\|^2 \leq \|\mathbf{M}(\mathbf{w})\boldsymbol{\mu}\|^2 \text{ for any } \mathbf{w}^*, \mathbf{w} \in \mathbf{H}_M,$$

$$(iv) \quad \text{Tr}[\mathbf{P}(\mathbf{w})\mathbf{P}(\mathbf{w}^*)\mathbf{P}(\mathbf{w})] \leq \text{Tr}[\mathbf{P}(\mathbf{w})\mathbf{P}(\mathbf{w})] \text{ for any } \mathbf{w}^*, \mathbf{w} \in \mathbf{H}_M.$$

The matrix $\mathbf{P}(\mathbf{w})$ is not a traditional projection matrix. For example, it is not idempotent. Therefore, we cannot further simplify Lemma 1 (iv).

Define the average mean squared error as $L_n(\mathbf{w}) \equiv (\boldsymbol{\mu}(\mathbf{w}) - \boldsymbol{\mu})^\top (\boldsymbol{\mu}(\mathbf{w}) - \boldsymbol{\mu})$ and the conditional average mean squared error as $R_n(\mathbf{w}) \equiv \mathbb{E}(L_n(\mathbf{w})|\mathbf{X})$. The same definitions can be found in Li (1987), Hansen (2007) and Wan et al. (2010). We investigate the asymptotic properties of $R_n(\mathbf{w})$ in the next lemma.

Lemma 2 *We have*

$$(i) \quad R_n(\mathbf{w}) \geq \|\mathbf{M}(\mathbf{w})\boldsymbol{\mu}\|^2,$$

$$(ii) \quad R_n(\mathbf{w}) \geq \sigma^2 \text{Tr}(\mathbf{P}(\mathbf{w})\mathbf{P}(\mathbf{w})).$$

Assumption 1 *For some fixed integer $1 \leq G < \infty$, we have $\mathbb{E}(|u_i|^{4G}|\mathbf{x}_i) \leq \kappa < \infty$.*

Assumption 1 is a bound condition on the conditional moments of the error term. It can be compared with the corresponding condition in Hansen (2007), in which

$$\mathbb{E}(|u_i|^{4(N+1)}|\mathbf{x}_i) \leq \kappa < \infty. \tag{12}$$

Note that (12) depends on model weights \mathbf{w} through Hansen's assumption that $w^{(m)}$ is restricted to the set $\{0, \frac{1}{N}, \frac{2}{N}, \dots, 1\}$ for some integer N .

Assumption 2 *As $n \rightarrow \infty$, $\xi_n^{-2G} M \sum_{m=1}^M (R_n(\mathbf{w}_m^0))^G \rightarrow 0$, where $\xi_n = \inf_{\mathbf{w} \in \mathbf{H}_M} R_n(\mathbf{w})$ and \mathbf{w}_m^0 is an $M \times 1$ vector of which the m^{th} element is one and the others are zeros.*

Assumption 2 is the convergence condition. It is the same as the convergence condition in Wan et al. (2010). A necessary condition for Assumption 2 to hold is $\xi_n \rightarrow \infty$, which indicates that there is no finite approximating model for which the bias is zero. Moreover, we assume that as $n \rightarrow \infty$, ξ_n^{2G} goes to infinity at a faster rate than $M \sum_{m=1}^M (R_n(\mathbf{w}_m^0))^G$. This is a relatively stronger assumption than the corresponding condition required by Hansen (2007), which only needs $\xi_n \rightarrow \infty$ as $n \rightarrow \infty$. Note that Hansen's (2007) assumptions are appropriate only for nested models and a discrete set for \mathbf{w} ; in contrast, our theorem is built on a more general set-up with non-nested models and continuous \mathbf{H}_M . In practice, Assumption 2 can be easily sustained by removing poor models prior to estimation. See Wan et al. (2010) for two explicit examples under which Assumption 2 holds.

Assumption 3 *As $n \rightarrow \infty$, $k^{(m)} \rightarrow \infty$ and $k^{(m)}/n \rightarrow 0$ for all m .*

Assumption 3 states that as n goes to infinity, $k^{(m)}$ goes to infinity at a slower rate for $m = 1, \dots, M$. Similar assumptions can be found in other papers, such as Shibata (1981) and Hansen (2007). Based on Assumption 3, we have

Lemma 3 *Let Assumption 3 hold. Then, $k(\mathbf{w})/n \rightarrow 0$ as $n \rightarrow \infty$.*

With the above lemmas and assumptions, we now show that the MAPC estimator is asymptotically optimal in the following theorem.

Theorem 1 *Let Assumptions 1, 2 and 3 hold. Then, as $n \rightarrow \infty$*

$$\frac{L_n(\hat{\mathbf{w}})}{L_n(\mathbf{w}_{opt})} \xrightarrow{p} 1,$$

where

$$\mathbf{w}_{opt} = \arg \inf_{\mathbf{w} \in \mathbf{H}_M} L_n(\mathbf{w}).$$

Theorem 1 states that by using the empirical weight vector $\hat{\mathbf{w}}$, the mean squared error is asymptotically equivalent to the lowest possible mean squared error. This implies that the MAPC estimator is asymptotically optimal in the class of model average estimators (5) where the weight vector belongs to the set \mathbf{H}_M .

Define the average estimate of σ^2 as $\hat{\sigma}^2(\mathbf{w})$:

$$\hat{\sigma}^2(\mathbf{w}) \equiv \frac{(\mathbf{y} - \hat{\boldsymbol{\mu}}(\mathbf{w}))^\top (\mathbf{y} - \hat{\boldsymbol{\mu}}(\mathbf{w}))}{n - k(\mathbf{w})} = \frac{\mathbf{w}^\top \hat{\mathbf{U}}^\top \hat{\mathbf{U}} \mathbf{w}}{n - \mathbf{k}^\top \mathbf{w}}.$$

By inserting $\hat{\sigma}^2(\mathbf{w})$ into (10) and rearranging the equation, we can rewrite MAPC as

$$\text{MAPC}_n(\mathbf{w}) = \mathbf{w}^\top \hat{\mathbf{U}}^\top \hat{\mathbf{U}} \mathbf{w} + 2\hat{\sigma}^2(\mathbf{w}) \mathbf{k}^\top \mathbf{w},$$

which is similar to MMA in (11) with σ^2 being replaced by an average estimate $\hat{\sigma}^2(\mathbf{w})$.

Theorem 2 *Let Assumption 3 hold. Then $\hat{\sigma}^2(\mathbf{w}) \xrightarrow{p} \sigma^2$ as $n \rightarrow \infty$.*

Note that $\hat{\sigma}^2(\mathbf{w})$ is a consistent estimator of σ^2 . Theorem 2 implies $\hat{\sigma}^2(\hat{\mathbf{w}}) \xrightarrow{p} \sigma^2$ as $n \rightarrow \infty$, where

$$\hat{\sigma}^2(\hat{\mathbf{w}}) = \frac{\hat{\mathbf{w}}^\top \hat{\mathbf{U}}^\top \hat{\mathbf{U}} \hat{\mathbf{w}}}{n - \mathbf{k}^\top \hat{\mathbf{w}}}. \quad (13)$$

The MMA criterion includes an infeasible σ^2 . Therefore, this must be computed with a sample estimate in practice. Hansen (2007, 2008) recommended using $\hat{\sigma}_L^2$ to replace the unknown σ^2 , where

$$\hat{\sigma}_L^2 = \frac{(\mathbf{y} - \hat{\boldsymbol{\mu}}^{(L)})^\top (\mathbf{y} - \hat{\boldsymbol{\mu}}^{(L)})}{n - k^{(L)}}$$

is the estimated σ^2 from a large approximation model L .¹ As a result, the MMA criterion in practice becomes

$$\text{MMA}_n(\mathbf{w}) = \mathbf{w}^\top \hat{\mathbf{U}}^\top \hat{\mathbf{U}} \mathbf{w} + 2\hat{\sigma}_L^2 \mathbf{k}^\top \mathbf{w}.$$

We can rewrite $\hat{\sigma}_L^2$ as an average estimator such that

$$\hat{\sigma}_L^2 = \frac{(\mathbf{y} - \hat{\boldsymbol{\mu}}^{(L)})^\top (\mathbf{y} - \hat{\boldsymbol{\mu}}^{(L)})}{n - k^{(L)}} = \frac{(\mathbf{y} - \hat{\boldsymbol{\mu}}(\mathbf{w}_L^0))^\top (\mathbf{y} - \hat{\boldsymbol{\mu}}(\mathbf{w}_L^0))}{n - \mathbf{k}^\top \mathbf{w}_L^0},$$

where \mathbf{w}_L^0 is a weight vector in which the L^{th} element is one and the others are zeros. In practice, using $\hat{\sigma}_L^2$ as an approximation to the infeasible σ^2 can be inefficient at times. As we show in the next section, simulation evidence indicates that the MMA estimator with $\hat{\sigma}_L^2$ yields a higher mean squared error than the MAPC estimator in many cases.

The MMA estimator is a two-step estimator since a sample estimate of σ^2 must be provided prior to estimation. In contrast, the MAPC estimator is a continuous updating estimator that requires only one step of calculation. Estimating $\hat{\mathbf{w}}$ from the MMA criterion with constraints is a classic quadratic programming problem, while estimating $\hat{\mathbf{w}}$ by the MAPC estimator is a convex optimization problem.²

3.2 Estimating the Variance-Covariance Matrix for $\hat{\boldsymbol{\beta}}(\hat{\mathbf{w}})$

The average coefficient $\hat{\boldsymbol{\beta}}(\hat{\mathbf{w}})$ is not easy to compute. The number of regressors, $k^{(m)}$, is usually not the same for different models; even if the $k^{(m)}$ are the same for certain models, the regressors must be different from one model to another. Either scenario complicates the computation.

Assume that there exists a model L within which all of the approximation models are nested.³ The average coefficient $\hat{\boldsymbol{\beta}}(\hat{\mathbf{w}})$ can be computed by

$$\hat{\boldsymbol{\beta}}(\hat{\mathbf{w}}) = \sum_{m=1}^M \hat{w}^{(m)} \left(\boldsymbol{\Gamma}^{(m)} \hat{\boldsymbol{\beta}}^{(m)} \right),$$

where $\boldsymbol{\Gamma}^{(m)}$ is $k^{(L)} \times k^{(m)}$ and plays the role of mapping the $k^{(m)} \times 1$ vector $\hat{\boldsymbol{\beta}}^{(m)}$ to $k^{(L)} \times 1$ by filling the extra parameters with 0. A convenient way to construct $\boldsymbol{\Gamma}^{(m)}$ is to use the following equation:

$$\boldsymbol{\Gamma}^{(m)} = \left(\mathbf{X}^{(L)\top} \mathbf{X}^{(L)} \right)^{-1} \mathbf{X}^{(L)\top} \mathbf{X}^{(m)}. \quad (14)$$

In this case, the rank of $\boldsymbol{\Gamma}^{(m)}$ is $k^{(m)}$ and each element in $\boldsymbol{\Gamma}^{(m)}$ is either 1 or 0.

¹In fact, Hansen (2007, 2008) used the largest approximation model in his simulation experiments.

²Note that convex optimization usually requires slightly more computation time than quadratic programming.

³If no such model exists, one can be easily created by including all regressors within it.

Although not shown for the sake of brevity, the variances and covariances in the rest of this subsection are in fact conditional on $\mathbf{X}^{(L)}$. Therefore, $\mathbf{X}^{(m)}$ can be assumed to be exogenous for all m . By straightforward algebra, the variance-covariance matrix of $\hat{\boldsymbol{\beta}}(\hat{\mathbf{w}})$ is

$$\begin{aligned}\text{Var}\left(\hat{\boldsymbol{\beta}}(\hat{\mathbf{w}})\right) &= \text{Var}\left(\sum_{m=1}^M \hat{w}^{(m)}\left(\boldsymbol{\Gamma}^{(m)}\hat{\boldsymbol{\beta}}^{(m)}\right)\right) \\ &= \sum_{m=1}^M \sum_{s=1}^M \hat{w}^{(m)}\hat{w}^{(s)}\text{Cov}\left(\boldsymbol{\Gamma}^{(m)}\hat{\boldsymbol{\beta}}^{(m)}, \boldsymbol{\Gamma}^{(s)}\hat{\boldsymbol{\beta}}^{(s)}\right).\end{aligned}\quad (15)$$

The right-hand-side of equation (15) is a linear combination of $\text{Cov}(\boldsymbol{\Gamma}^{(m)}\hat{\boldsymbol{\beta}}^{(m)}, \boldsymbol{\Gamma}^{(s)}\hat{\boldsymbol{\beta}}^{(s)})$. When s equals m , the covariance matrix becomes the variance-covariance matrix of $\boldsymbol{\Gamma}^{(m)}\hat{\boldsymbol{\beta}}^{(m)}$. Each $\boldsymbol{\Gamma}^{(m)}\hat{\boldsymbol{\beta}}^{(m)}$ includes possible model misspecification bias. We define the following $k^{(L)} \times 1$ vector

$$\mathbf{d}^{(m)} = \mathbb{E}\left(\boldsymbol{\Gamma}^{(m)}\hat{\boldsymbol{\beta}}^{(m)}\right) - \boldsymbol{\beta}(\mathbf{w}) \quad (16)$$

as the misspecification bias vector. We propose an estimator for $\text{Cov}(\boldsymbol{\Gamma}^{(m)}\hat{\boldsymbol{\beta}}^{(m)}, \boldsymbol{\Gamma}^{(s)}\hat{\boldsymbol{\beta}}^{(s)})$ in the following lemma.

Lemma 4 *For any approximation models m and s , in which m and s can represent the same model, we have*

$$\begin{aligned}&\widehat{\text{Cov}}\left(\boldsymbol{\Gamma}^{(m)}\hat{\boldsymbol{\beta}}^{(m)}, \boldsymbol{\Gamma}^{(s)}\hat{\boldsymbol{\beta}}^{(s)}\right) \\ &= \hat{\sigma}^2(\hat{\mathbf{w}})\boldsymbol{\Gamma}^{(m)}\left(\mathbf{X}^{(m)\top}\mathbf{X}^{(m)}\right)^{-1}\mathbf{X}^{(m)\top}\mathbf{X}^{(s)}\left(\mathbf{X}^{(s)\top}\mathbf{X}^{(s)}\right)^{-1}\boldsymbol{\Gamma}^{(s)\top} + \hat{\mathbf{d}}^{(m)}\left(\hat{\mathbf{d}}^{(s)}\right)^\top,\end{aligned}$$

where

$$\hat{\mathbf{d}}^{(m)} = \boldsymbol{\Gamma}^{(m)}\hat{\boldsymbol{\beta}}^{(m)} - \hat{\boldsymbol{\beta}}(\hat{\mathbf{w}}). \quad (17)$$

For a particular average coefficient, for example $\hat{\beta}_j(\hat{\mathbf{w}})$, the variance of $\hat{\beta}_j(\hat{\mathbf{w}})$ is

$$\text{Var}\left(\hat{\beta}_j(\hat{\mathbf{w}})\right) = \left[\text{Var}\left(\hat{\boldsymbol{\beta}}(\hat{\mathbf{w}})\right)\right]_{jj},$$

which is the j^{th} element on the diagonal of the variance-covariance matrix $\text{Var}\left(\hat{\boldsymbol{\beta}}(\hat{\mathbf{w}})\right)$. Buckland et al. (1997) proposed other estimators for $\text{Var}\left(\hat{\beta}_j(\hat{\mathbf{w}})\right)$. Some of the estimators are based on a restrictive assumption that there is perfect correlation between each $\boldsymbol{\beta}^{(m)}$. They also proposed computing $\text{Var}\left(\hat{\beta}_j(\hat{\mathbf{w}})\right)$ via a pairs bootstrap.

In practice, we want to include certain regressors in every approximation model because these regressors are of primary interest to us. Models without these regressors provide no useful information and therefore are not of interest to us. We name these regressors the core regressors. Let the $n \times k_c$ matrix \mathbf{X}_c represent the core regressors and let $\hat{\boldsymbol{\beta}}_c(\hat{\mathbf{w}})$ be the

corresponding averaged core coefficients. We define $\Gamma_c^{(m)}$ as a $k^{(m)} \times k_c$ matrix that plays the role of subtracting the $n \times k_c$ matrix \mathbf{X}_c out of any $\mathbf{X}^{(m)}$ such that $\mathbf{X}^{(m)}\Gamma_c^{(m)} = \mathbf{X}_c$. Similar to $\Gamma^{(m)}$ defined in (14), we can construct $\Gamma_c^{(m)}$ using $\Gamma_c^{(m)} = (\mathbf{X}^{(m)\top}\mathbf{X}^{(m)})^{-1}\mathbf{X}^{(m)\top}\mathbf{X}_c$. The variance-covariance matrix for $\hat{\beta}_c(\hat{\mathbf{w}})$ is then

$$\text{Var}\left(\hat{\beta}_c(\hat{\mathbf{w}})\right) = \Gamma_c^{(m)\top}\text{Var}\left(\hat{\beta}(\hat{\mathbf{w}})\right)\Gamma_c^{(m)}.$$

We name the regressors that are of marginal interest to us potential regressors. These potential regressors are not included in every approximation model. We let β_p be the corresponding coefficient. The coefficient β_p is a $k_p \times 1$ vector with $k_p \leq k^{(L)}$. We define the regressors that are not included in $\mathbf{X}^{(m)}$ as $\mathbf{X}^{(-m)}$ and the associated coefficient as β_{-m} . By definition, the set of β_{-m} belongs to the set of β_p for all m . Similar to the definition of $\Gamma^{(m)}$ in (14), we let $\Gamma^{(-m)} = (\mathbf{X}^{(L)\top}\mathbf{X}^{(L)})^{-1}\mathbf{X}^{(L)\top}\mathbf{X}^{(-m)}$ play the role of mapping the $(k^{(L)} - k^{(m)}) \times 1$ vector β_{-m} to $k^{(L)} \times 1$. How to distinguish potential regressors from core regressors is not the primary concern of this paper. We leave that for future research.

3.3 Asymptotic Inference and Bootstrap Based Inference

We start by listing more assumptions.

Assumption 4 *We have*

- (i) (\mathbf{x}_i, u_i) is an iid sequence,
- (ii) $\mathbb{E}(\mathbf{x}_i u_i) = \mathbf{0}$, $\mathbb{E}|x_{ij} u_i|^2 < \infty$, $\text{Var}(n^{-1/2}\mathbf{X}^\top \mathbf{u})$ is positive definite,
- (iii) As $n \rightarrow \infty$, $n^{-1}\mathbf{X}^{(m)\top}\mathbf{X}^{(m)} \xrightarrow{p} \mathbf{S}^{(m)}$ and $n^{-1}\mathbf{X}^{(m)\top}\mathbf{X}^{(s)} \xrightarrow{p} \mathbf{S}^{(m,s)} \forall m, s$, where both $\mathbf{S}^{(m)}$ and $\mathbf{S}^{(m,s)}$ are finite, deterministic matrices and $\mathbf{S}^{(m)}$ is also positive definite.

Assumption 5 $\beta_p = \mathbf{h}_p/\sqrt{n}$, where \mathbf{h}_p is a fixed vector.

Assumption 6 *The average coefficient $\beta(\mathbf{w})$ is a function of $\beta^{(m)}$ and β_{-m} that can be written as $\beta(\mathbf{w}) = \mathbf{f}(\beta^{(m)}, \beta_{-m})$, where $\beta(\mathbf{w}) = \mathbf{f}(\beta^{(m)}, \beta_{-m})$ is twice differentiable in a neighborhood of β_{-m} and $\mathbf{f}(\beta^{(m)}, \mathbf{0}) = \Gamma^m \beta^{(m)}$.*

Assumption 4 is a standard assumption about regressors and error terms. Assumption 5 follows the logic that the potential regressors are only of marginal interest to us as they are weakly correlated with \mathbf{y} . This idea of weak variables is similar to ideas put forth in the weak instruments literature. For example, the Assumption L_{Π} in Staiger and Stock (1997)

is almost identical to our Assumption 5. Assumption 6 is a standard assumption that allows us to use Taylor expansion on $\beta(\mathbf{w})$.

Let $\hat{\beta}_c^{(m)} = \Gamma_c^{(m)\top} \hat{\beta}^{(m)}$ be the core part of $\hat{\beta}^{(m)}$. We investigate the asymptotic distribution of $\hat{\beta}_c^{(m)}$ in the following lemma:

Lemma 5 *Let Assumptions 4, 5, and 6 hold. Then, as $n \rightarrow \infty$,*

(i) $\beta_{-m} = \mathbf{h}_{-m}/\sqrt{n} \forall m$, where \mathbf{h}_{-m} is a fixed vector.

(ii) $\sqrt{n} \left(\hat{\beta}_c^{(m)} - \beta_c \right) \xrightarrow{d} N \left(\delta_1^{(m)}, \sigma^2 \Gamma_c^{(m)\top} (\mathbf{S}^{(m)})^{-1} \Gamma_c^{(m)} \right)$, where

$$\delta_1^{(m)} \equiv \Gamma_c^{(m)\top} (\mathbf{S}^{(m)})^{-1} \mathbf{S}^{(m,-m)} \mathbf{h}_{-m}.$$

(iii) Define $\mathbf{F}_{\beta_{-m}} \equiv \partial \mathbf{f}(\beta^{(m)}, \beta_{-m})^\top / \partial \beta_{-m}$, then

$$\sqrt{n}(\beta_c - \beta_c(\mathbf{w})) \xrightarrow{p} \delta_2^{(m)},$$

where $\delta_2^{(m)} \equiv -\Gamma_c^{(L)\top} \left(\mathbf{F}_{\beta_{-m}} |_{\beta_{-m}=\mathbf{0}} \right) \mathbf{h}_{-m}$.

(iv) $\sqrt{n} \mathbf{d}^{(m)} \xrightarrow{p} \delta^{(m)}$, where $\delta^{(m)} \equiv \Gamma_c^{(m)\top} (\mathbf{S}^{(m)})^{-1} \mathbf{S}^{(m,-m)} \mathbf{h}_{-m} - \left(\mathbf{F}_{\beta_{-m}} |_{\beta_{-m}=\mathbf{0}} \right) \mathbf{h}_{-m}$.

(v) The asymptotic distribution of $\hat{\beta}_c^{(m)}$ is

$$\sqrt{n} \left(\hat{\beta}_c^{(m)} - \beta_c(\mathbf{w}) \right) \rightarrow_d \Lambda_c^{(m)} \sim N \left(\delta_c^{(m)}, \mathbf{V}_c^{(m)} \right), \quad (18)$$

where $\delta_c^{(m)} = \delta_1^{(m)} + \delta_2^{(m)}$ and $\mathbf{V}_c^{(m)} \equiv \sigma^2 \Gamma_c^{(m)\top} (\mathbf{S}^{(m)})^{-1} \Gamma_c^{(m)}$.

Based on Lemma 5, we derive the asymptotic distribution of the core average coefficient $\hat{\beta}_c(\hat{\mathbf{w}})$ conditional on $\hat{\mathbf{w}}$ in the following theorem:

Theorem 3 *Let Assumptions 4, 5, and 6 hold. Then, as $n \rightarrow \infty$,*

$$\sqrt{n} \left(\hat{\beta}_c(\hat{\mathbf{w}}) - \beta_c(\mathbf{w}) \right) \Big| \hat{\mathbf{w}} \rightarrow_d \Lambda_c = \sum_{m=1}^M \hat{w}^{(m)} \Lambda_c^{(m)} \sim N \left(\mathbf{0}, \Gamma_c^{(m)\top} \mathbf{V}(\hat{\mathbf{w}}) \Gamma_c^{(m)} \right),$$

where $\Lambda_c^{(m)}$ is defined in (18) and

$$\mathbf{V}(\hat{\mathbf{w}}) \equiv \sum_{m=1}^M \sum_{s=1}^M \hat{w}^{(m)} \hat{w}^{(s)} \left(\sigma^2 \Gamma_c^{(m)} (\mathbf{S}^{(m)})^{-1} \mathbf{S}^{(m,s)} (\mathbf{S}^{(s)})^{-1} \Gamma_c^{(s)\top} + \delta^{(m)} \delta^{(s)\top} \right).$$

The terms $\mathbf{S}^{(m)}$, $\mathbf{S}^{(m,s)}$, and $\delta^{(m)}$ are defined in Assumption 4 and Lemma 5.

There is joint convergence of $\sqrt{n}(\hat{\beta}_c^{(m)} - \beta_c(\mathbf{w}))$ and the stochastic weights $\hat{\mathbf{w}}$. The conditional asymptotic distribution $\mathbf{\Lambda}$ is normal, which implies that the unconditional asymptotic distribution of $\sqrt{n}(\hat{\beta}_c(\hat{\mathbf{w}}) - \beta_c(\mathbf{w}))$ is a mixed normal distribution. There is a large literature in time series that studies the mixed normal distribution and its inference (see Johansen (1995, pp.177–178) for a detailed explanation).

To test a single restriction, for example $\beta_j(\mathbf{w}) = \beta_{j0}$, we derive the t -statistic for the average core coefficient $\hat{\beta}_j(\hat{\mathbf{w}})$

$$t_{\hat{\beta}_j(\hat{\mathbf{w}})} \equiv \frac{\hat{\beta}_j(\hat{\mathbf{w}}) - \beta_{j0}}{\widehat{\text{Var}}(\hat{\beta}_j(\hat{\mathbf{w}}))^{1/2}},$$

where the estimated variance $\widehat{\text{Var}}(\hat{\beta}_j(\hat{\mathbf{w}}))$ is the j^{th} element on the diagonal of $\widehat{\text{Var}}(\hat{\beta}(\hat{\mathbf{w}}))$:

$$\begin{aligned} \widehat{\text{Var}}(\hat{\beta}(\hat{\mathbf{w}})) &= \sum_{m=1}^M \sum_{s=1}^M \hat{w}^{(m)} \hat{w}^{(s)} \left(\hat{\mathbf{d}}^{(m)} \left(\hat{\mathbf{d}}^{(s)} \right)^\top \right. \\ &\quad \left. + \hat{\sigma}^2(\hat{\mathbf{w}}) \mathbf{\Gamma}^{(m)} \left(\mathbf{X}^{(m)\top} \mathbf{X}^{(m)} \right)^{-1} \mathbf{X}^{(m)\top} \mathbf{X}^{(s)} \left(\mathbf{X}^{(s)\top} \mathbf{X}^{(s)} \right)^{-1} \mathbf{\Gamma}^{(s)\top} \right), \end{aligned}$$

where $\hat{\mathbf{d}}^{(m)}$ and $\hat{\mathbf{d}}^{(s)}$ are defined in (17) and $\hat{\sigma}^2(\hat{\mathbf{w}})$ is defined in (13). To test the joint null hypothesis that $\mathbf{R}\beta_c(\mathbf{w}) = \mathbf{r}$, where \mathbf{r} is $k_r \times 1$, we use the Wald statistic such that

$$W_{\mathbf{r}} \equiv \left(\mathbf{R}\hat{\beta}_c(\hat{\mathbf{w}}) - \mathbf{r} \right)^\top \left(\mathbf{R}\widehat{\text{Var}}(\hat{\beta}_c(\hat{\mathbf{w}}))\mathbf{R}^\top \right)^{-1} \left(\mathbf{R}\hat{\beta}_c(\hat{\mathbf{w}}) - \mathbf{r} \right).$$

We derive the asymptotic distribution of the t -statistic and the Wald statistic in the following theorem.

Theorem 4 *Let Assumptions 4, 5, and 6 hold. Then*

$$t_{\hat{\beta}_j(\hat{\mathbf{w}})} \longrightarrow_d \text{N}(0, 1) \quad \text{and} \quad W_{\mathbf{r}} \longrightarrow_d \chi^2(k_r).$$

We can construct the $1 - \alpha$ confidence interval for $\hat{\beta}_j(\hat{\mathbf{w}})$ in the classical way:

$$\left[\beta_{\text{lower}}, \beta_{\text{upper}} \right] = \left[\hat{\beta}_j(\hat{\mathbf{w}}) - \hat{s}_j z_{1-\alpha/2}, \hat{\beta}_j(\hat{\mathbf{w}}) + \hat{s}_j z_{1-\alpha/2} \right],$$

where \hat{s}_j is the standard error for $\hat{\beta}_j(\hat{\mathbf{w}})$ and $z_{1-\alpha/2}$ is the $1 - \alpha/2$ quantile of a standard normal distribution.

Since both test statistics are asymptotically pivotal, we can use a semiparametric bootstrap test to provide improved statistical inference in finite samples. We denote $\hat{\tau}$ as an estimated test statistic (t or Wald). We first compute the estimated residuals $\tilde{\mathbf{u}}$ by plugging in the estimated average coefficients $\tilde{\beta}(\tilde{\mathbf{w}})$ under the null hypothesis. We then resample $\tilde{\mathbf{u}}$

B times and obtain B bootstrap samples \mathbf{y}_l^* for $l = 1, \dots, B$. The value of B should satisfy the condition that $\alpha(1 + B)$ is an integer,⁴ where α is the desired level of significance. For each bootstrap sample \mathbf{y}_l^* , we compute a simulated test statistic $\hat{\tau}_l^*$ in exactly the same way that $\hat{\tau}$ was computed from the original data.

Following MacKinnon (2009), we can use the following equation to compute the bootstrap P value for a one-tail test that rejects in the upper tail, as is the case in the Wald test:

$$\hat{p}^*(\hat{\tau}) = \frac{1}{B} \sum_{l=1}^B I(\hat{\tau}_l^* > \hat{\tau}),$$

where $I(\cdot)$ is the indicator function, which takes the value 1 when its argument is true and takes the value 0 otherwise. If we assume that τ is symmetrically distributed around zero, as is the case in the t test, we can use the symmetric bootstrap:

$$\hat{p}^*(\hat{\tau}) = \frac{1}{B} \sum_{l=1}^B I(|\hat{\tau}_l^*| > |\hat{\tau}|).$$

If we are not willing to make the assumption that τ is symmetrically distributed around zero, we can instead use the equal-tail bootstrap:

$$\hat{p}^*(\hat{\tau}) = 2 \min \left(\frac{1}{B} \sum_{l=1}^B I(\hat{\tau}_l^* \leq \hat{\tau}), \frac{1}{B} \sum_{l=1}^B I(\hat{\tau}_l^* > \hat{\tau}) \right).$$

Bootstrap tests generally perform better than asymptotic tests, especially when working with a small sample size (see Section 4.2 as an example).

To construct symmetric bootstrap confidence intervals, we just need to invert the symmetric bootstrap test we described above. We can also use the bootstrap- t method to construct asymmetric bootstrap confidence intervals. We first sort the $\hat{\tau}_l^*$ in ascending order. We let the value of $c_{\alpha/2}^*$ be the value of number $\alpha(1 + B)/2$ bootstrap t -statistic in the sorted list. Similarly, the value of $c_{1-\alpha/2}^*$ is the value of number $(1 - \alpha/2)(1 + B)$ bootstrap t -statistic. Both $\alpha(1 + B)/2$ and $(1 - \alpha/2)(1 + B)$ should be integers. The bootstrap- t , or asymmetric equal-tail bootstrap confidence interval, is then

$$\left[\beta_{\text{lower}}, \beta_{\text{upper}} \right] = \left[\hat{\beta}_j(\hat{\mathbf{w}}) - \hat{s}_j c_{1-\alpha/2}^*, \hat{\beta}_j(\hat{\mathbf{w}}) + \hat{s}_j c_{\alpha/2}^* \right].$$

There are many ways to compute bootstrap confidence intervals. For a literature review, see DiCiccio and Efron (1996).

Computational cost can be a concern for model averaging estimation when the total number of approximation models is large. In this situation, we do not want to assign a huge

⁴If we are using a two-tailed test, it is helpful to make $\alpha(1+B)/2$ an integer for the purpose of constructing confidence intervals.

number for B . In fact, it is often possible to obtain reliable results without using a large value of B by using the iterative procedure proposed in Davidson and MacKinnon (2000).

4 Finite Sample Performance

This section contains two parts. In the first part, we investigate the finite sample performance of the MAPC estimator in a simulation experiment. In the second part, we compare the finite sample performance of different tests via rejection frequency under the same simulation design used in the first part. Contrary to the design in Hansen (2007), we propose a simulation experiment with non-nested models.

4.1 The MAPC Estimator Versus Other Estimators

The general unrestricted model is the simple regression model

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{u}, \quad (19)$$

where \mathbf{X} is $n \times k$, $\boldsymbol{\beta}$ is $k \times 1$ and \mathbf{u} is $n \times 1$. The number of regressors k increases as n increases but at a slower rate, in which $k = \text{round}(3+n^{1/5})$. We set the first column of \mathbf{X} , \mathbf{x}_1 , to be the intercept; the remaining \mathbf{x}_i are assumed to be correlated with each other and are generated by $N(\mathbf{0}, \boldsymbol{\Sigma})$, where $\boldsymbol{\Sigma}$ is a $(k-1) \times (k-1)$ symmetric matrix with all diagonal terms equal to 1 and all off-diagonal terms equal to ρ . The $k \times 1$ coefficients $\boldsymbol{\beta}$ are determined by $\boldsymbol{\beta} = [1/5, 1/5, 5/\sqrt{n}, 4/\sqrt{n}, \dots, (8-k)/\sqrt{n}]^\top$, where the first two regressors are core regressors and the remaining regressors are potential regressors. The error term \mathbf{u} is independent of the regressors \mathbf{X} and is distributed as $N(0, \sigma_u^2 \mathbf{I})$. The parameter σ_u controls the population $R^2 = \boldsymbol{\beta}_2^\top \boldsymbol{\Sigma} \boldsymbol{\beta}_2 / (\boldsymbol{\beta}_2^\top \boldsymbol{\Sigma} \boldsymbol{\beta}_2 + \sigma_u^2)$ so as to vary on a grid between 0.01 and 0.99, where $\boldsymbol{\beta}_2 = [\beta_2, \dots, \beta_k]^\top$. We consider six different sample sizes, in which $n = 25, 50, 100, 200, 400$ and 800. Other simulation results, which are not reported here, demonstrate that the findings are not sensitive to alternative distributions. We also find that the results are not sensitive to different values of ρ , which we will show later.

All submodels that are nested in the general unrestricted model (19) are treated as approximation models. Therefore, the approximation models are clearly non-nested in our experiment. The total number of approximation models, M , is equal to the total number of combinations made by all the potential regressors.

We study seven methods: (1) general-to-specific approach (GETS); (2) Akaike information criterion (AIC); (3) model averaging by AIC (MA-AIC); (4) model averaging by Bayesian information criterion (MA-BIC); (5) Mallows model averaging (MMA); (6) jackknife model averaging (JMA); and (7) model averaging by prediction criterion (MAPC).

The GETS approach aims to modify the general unrestricted model by removing irrelevant variables according to pre-determined criteria. In our experiment, we adopt a simple

GETS approach from Hendry and Nielsen (2007). We first estimate model (19). Then, regressors with the absolute value of the t -statistics smaller than $c_\alpha = 2$ are eliminated. If multiple t -statistics are smaller than 2, we eliminate the smallest one. The remaining regressors are retained and form a new model for the next-round test until no regressors can be eliminated.

The Akaike information criterion (AIC) for a model m is defined as

$$\text{AIC}^{(m)} = n \log(\hat{\sigma}_m^2) + 2k^{(m)}.$$

The model that achieves the lowest value among all of the estimated $\text{AIC}^{(m)}$ is selected. The model average AIC (MA-AIC) makes use of the estimated $\text{AIC}^{(m)}$ to compute the empirical weights, where

$$\hat{w}_{\text{AIC}}^{(m)} = \exp\left(-\frac{1}{2}\text{AIC}^{(m)}\right) / \sum_{m=1}^M \exp\left(-\frac{1}{2}\text{AIC}^{(m)}\right).$$

MA-BIC computes the empirical weights for its associated average estimator according to

$$\hat{w}_{\text{BIC}}^{(m)} = \exp\left(-\frac{1}{2}\text{BIC}^{(m)}\right) / \sum_{m=1}^M \exp\left(-\frac{1}{2}\text{BIC}^{(m)}\right),$$

where

$$\text{BIC}^{(m)} = n \log(\hat{\sigma}_m^2) + \log(n)k^{(m)}.$$

Jackknife model averaging (JMA) (Hansen and Racine 2011) is also known as leave-one-out cross-validation. As its name indicates, JMA requires the jackknife residuals for the average estimator. The jackknife residual vector for model m can be conveniently written as $\hat{\mathbf{u}}_{\mathbf{J}}^{(m)} = \mathbf{D}^{(m)}\hat{\mathbf{u}}^{(m)}$, where $\hat{\mathbf{u}}^{(m)}$ is the least squares residual vector and $\mathbf{D}^{(m)}$ is the $n \times n$ diagonal matrix with the i^{th} diagonal element equal to $(1 - h_i^{(m)})^{-1}$. The term $h_i^{(m)}$ is the i^{th} diagonal element of $\mathbf{P}^{(m)}$ defined in (4). Define an $n \times M$ matrix that collects all the jackknife residuals, in which $\hat{\mathbf{U}}_{\mathbf{J}} = [\hat{\mathbf{u}}_{\mathbf{J}}^{(1)}, \dots, \hat{\mathbf{u}}_{\mathbf{J}}^{(M)}]$. The least squares cross-validation criterion for JMA is simply

$$\text{CV}_n(\mathbf{w}) = \frac{1}{n} \mathbf{w}^\top \hat{\mathbf{U}}_{\mathbf{J}}^\top \hat{\mathbf{U}}_{\mathbf{J}} \mathbf{w} \quad \text{with} \quad \hat{\mathbf{w}} = \underset{\mathbf{w} \in \mathbf{H}_M^*}{\text{argmin}} \text{CV}_n(\mathbf{w}).$$

The MAPC estimator and the MMA estimator are presented in previous sections. The infeasible σ^2 in (11) is replaced by a sample estimate, $\hat{\sigma}_L^2$, from the largest approximation model (19).

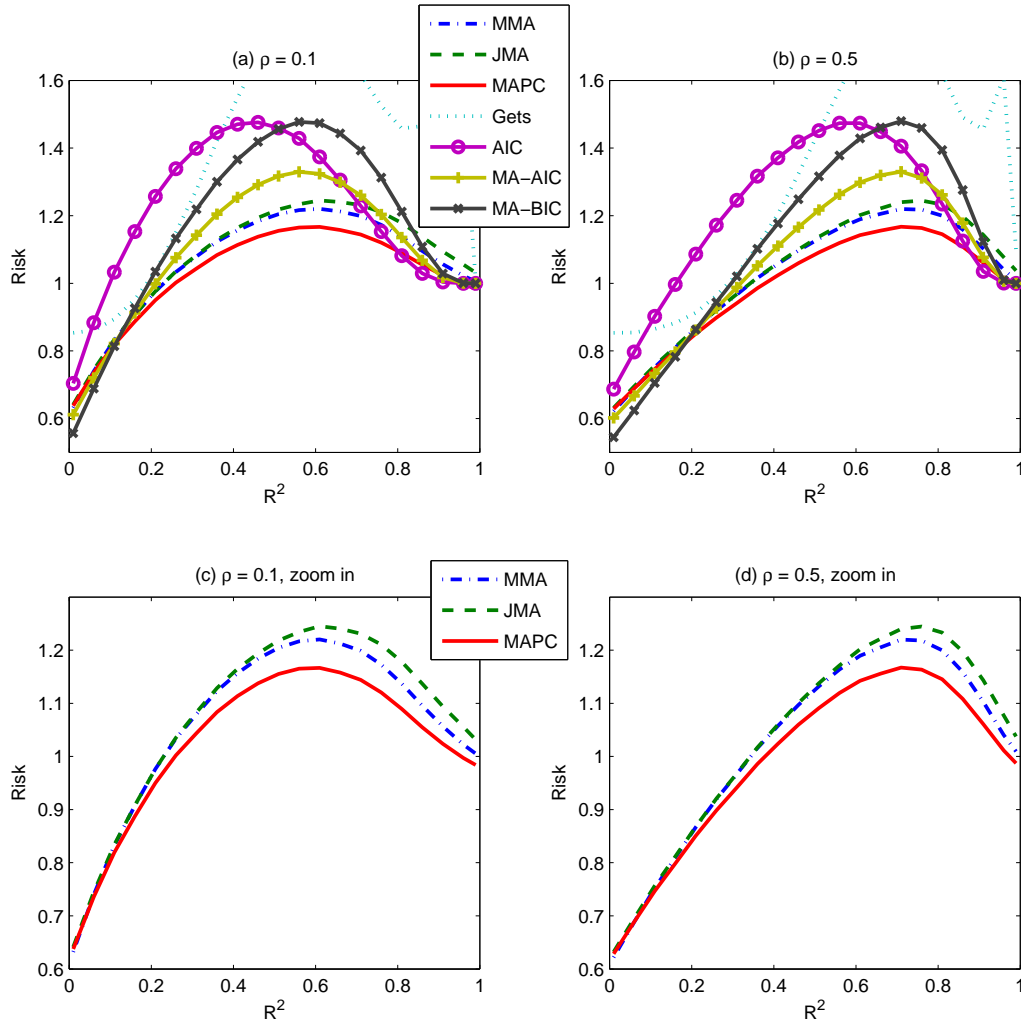
We define the risk of an estimator as the mean squared error such that

$$\text{Risk} \equiv \frac{1}{n} (\boldsymbol{\mu}(\hat{\mathbf{w}}) - \boldsymbol{\mu}_0)^\top (\boldsymbol{\mu}(\hat{\mathbf{w}}) - \boldsymbol{\mu}_0)$$

where $\boldsymbol{\mu}(\hat{\boldsymbol{w}})$ is the averaged $\boldsymbol{\mu}$ by the estimated $\hat{\boldsymbol{w}}$ and $\boldsymbol{\mu}_0$ is the true value of $\boldsymbol{\mu}$ (feasible in simulation). We compute risks for all seven estimators and average across 100,000 simulation draws. We normalize risks by dividing by the risk of the infeasible optimal least squares estimator (the OLS estimator of model (19)). We present the risk calculations for $n = 25$ with two different values of ρ (0.1 and 0.5) in Figures 1(a) and 1(b). The R^2 is presented on the x -axis and the risk is displayed on the y -axis. The dash-dotted line, dashed line, solid line, star, cross, circle, and x-mark correspond to MMA, JMA, MAPC, GETS, AIC, MA-AIC, and MA-BIC, respectively.

Figures 1(a) and 1(b) show similar results. As $R^2 \rightarrow 1$, risks of all seven methods tend to converge to 1, which suggests that all corresponding estimators are converging to the infeasible optimal least squares estimator. In many cases, the GETS method shows a much

Figure 1: Simulation Results for $n = 25$



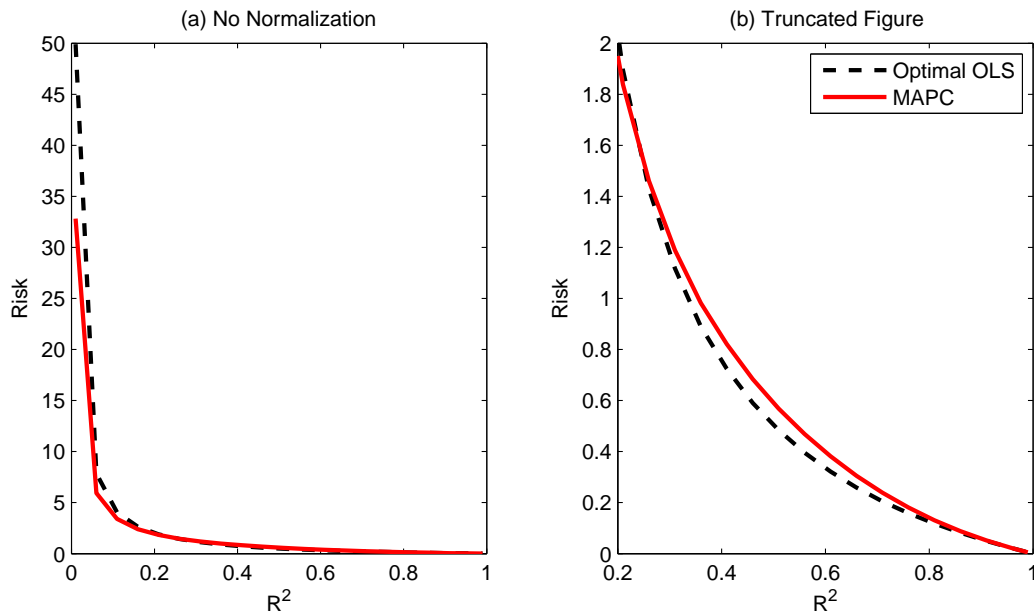
higher risk relative to other methods. We notice that when R^2 is close to 1, AIC yields the lowest risk. This means that the true model (19) is frequently selected by the AIC method, when there is sufficient information. However, AIC has a relatively poor performance for other values of R^2 . The MA-AIC estimator achieves lower risk than the MA-BIC estimator in most cases, except when R^2 is close to 0. In extreme cases, when R^2 is close to 0 or 1, MA-AIC and MA-BIC yield lower risks than MMA, JMA, and MAPC. However, they are outperformed by MMA, JMA, and MAPC in most other cases that are more reasonable in practice.

It is instructive to compare the performance of MMA, JMA, with MAPC. In most cases, these three methods have better performance when compared to other methods. We provide Figures 1(c) and 1(d) to show only these three methods. For all values of R^2 , the MAPC estimator achieves lower risks than the MMA estimator and the JMA estimator. This suggests that, in finite samples, MAPC is more efficient than MMA and JMA.

We notice that all the risk curves are hump-shaped. As we mentioned, all the risk curves are normalized by dividing by the risk of the infeasible optimal OLS estimator. The hump shape is caused by this normalization. To demonstrate the reason for this, we use the MAPC estimator as an example and present its estimated risk without normalization (denoted by solid lines) in Figure 2, along with the estimated risk of the optimal OLS estimator (denoted by dashed lines).

In Figure 2, part (a) represents the estimated risks for the MAPC estimator along with the OLS estimator. When R^2 is low, both estimators yield high risks. A low R^2 means

Figure 2: The Estimated Risks without Normalization



that the model does not fit the data very well. As a consequence, the estimated risks for the OLS estimator are high. The MAPC estimator, however, considers other approximation models as opposed to limiting itself to the unrestricted model (19). Therefore, as we see from the figure, the MAPC estimator yields smaller risk when R^2 is low. The corresponding normalized risk for MAPC is smaller than 1. As R^2 increases, the risk decreases for both estimators and the gap between the two curves shrinks. This implies that when R^2 is low, the normalized risk for the MAPC estimator increases as R^2 increases from 0.

Part (b) of Figure 2 is the truncated version of part (a) with y -axis being 0 to 2 and x -axis being 0.2 to 1. As R^2 increases, the estimated risk for the OLS estimator decreases and eventually becomes smaller than the MAPC risks. Therefore, the corresponding normalized risk for MAPC is greater than 1. As we can see from the figure, the gap between the two curves first expands and then shrinks. The two methods yield almost identical risks when $R^2 = 0.99$. This implies that the normalized risk curve is hump-shaped and that the normalized risk for the MAPC estimator is very close to 1 when $R^2 = 0.99$.

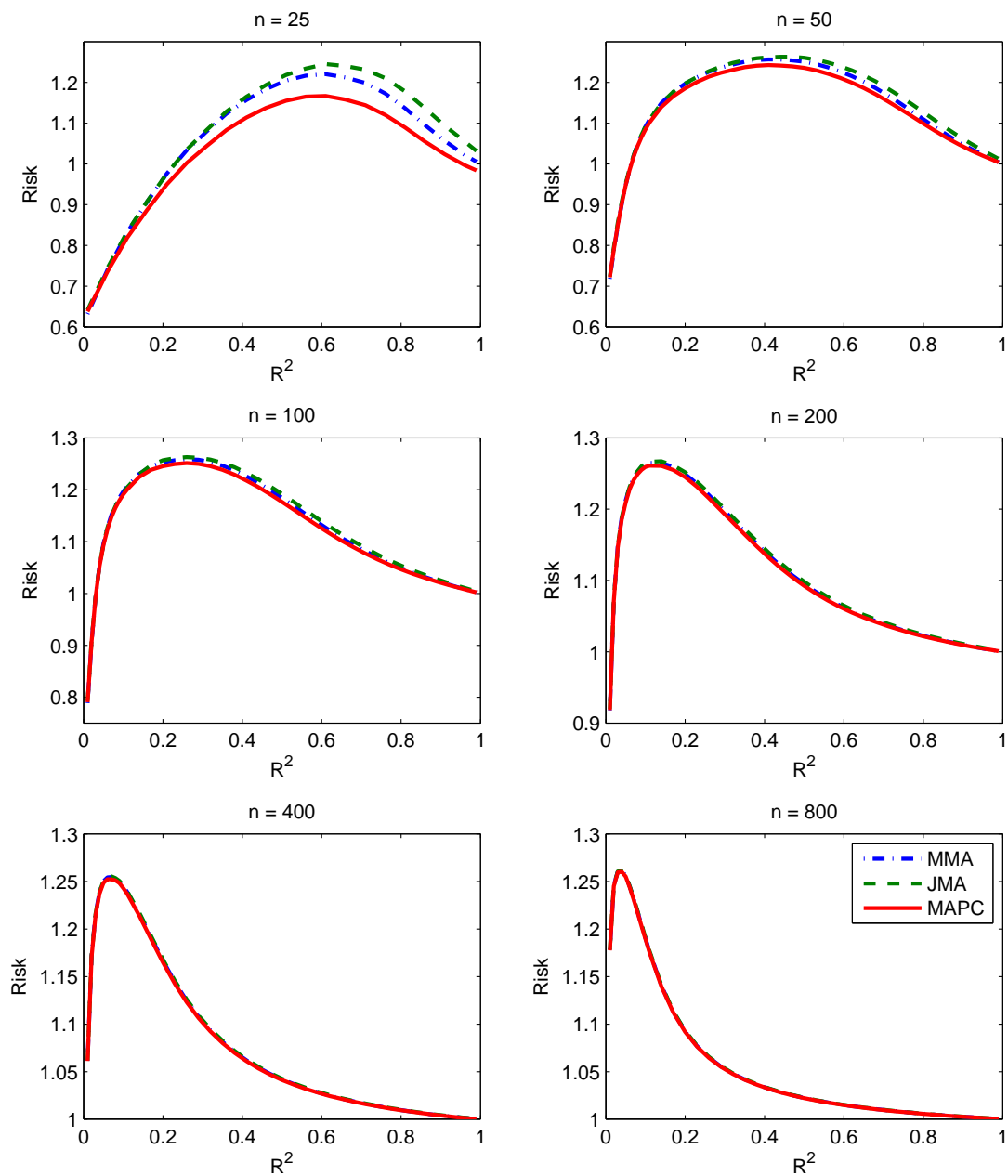
Risk calculations for various sample sizes ($n = 25, 50, 100, 200, 400$, and 800) are displayed in Figure 3. In this case, we set $\rho = 0.1$. To keep the figure uncluttered, only MMA, JMA and MAPC are displayed and represented by the dash-dotted line, dashed, and solid lines, respectively. In each panel, the MAPC estimator has a better performance than the MMA and JMA estimator by yielding lower risks. As n increases, the index on the y -axis shrinks and MMA and JMA merge with MAPC. This implies that, as $n \rightarrow \infty$, these three estimators converge to the infeasible optimal least squares estimator. We also notice that as n increases, the hump shifts from the right to the left. Again, this is caused by the normalization. As n becomes larger, the performance of the OLS estimator also gets better. Therefore, the values of the normalized risks increases when R^2 is low.

4.2 Evaluate Tests via Rejection Frequency

Buckland et al. (1997) proposed using a pairs bootstrap method to estimate the variance of an average coefficient. To implement this method, we first resample the data matrix and obtain B bootstrap resamples. Then, we apply model averaging to each resample and obtain B bootstrap estimates. Finally, we estimate the sample variance of the B bootstrap estimates as the estimated variance of the average coefficient. We can easily derive asymptotic tests (t and Wald) using the pairs bootstrap sample variance, henceforth pairs tests.

In Section 3.3, we derived the asymptotic tests (t and Wald) and the bootstrap tests for the average core coefficient. In this section, we compare the finite sample performance of these tests with pairs tests via the rejection frequency. We use the same simulation design we proposed in Section 4.1. We concentrate on the second coefficient, β_2 and the null hypothesis is that $\beta_2 = 1/5$, which is the true value of β_2 . We generate 10,000 simulation draws for 15 sample sizes: $n = 25, 30, 35, \dots, 95$. For each simulation draw, we compute the P values from the asymptotic tests, the semiparametric bootstrap tests and pairs tests for each sample size. In order to reduce the costs of doing the Monte Carlo experiments, the total number of

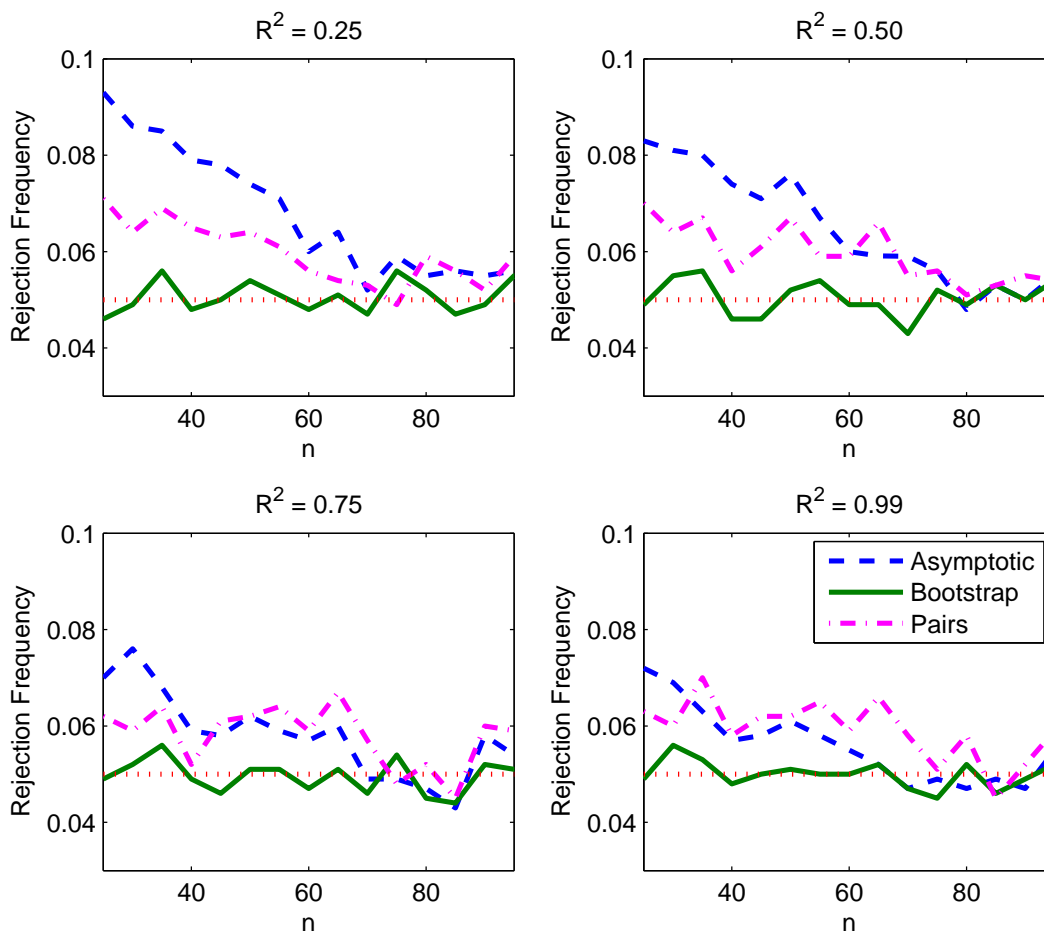
Figure 3: Simulation Results for Various Sample Sizes



bootstraps for each bootstrap test and pairs test is $B = 199$. We then record the frequencies with which tests based on the three p values reject at the 0.05 level. We consider four different R^2 : 0.25, 0.5, 0.75, and 0.99. The results of this simulation are presented in the following figure.

In Figure 4, asymptotic t tests, bootstrap tests, and pairs tests are represented by dashed lines, solid lines, and dash-dotted lines, respectively. Although not shown for the sake of brevity, we also compare the three tests using the Wald statistic. Simulation results are very similar to what is shown here. The results suggest that both asymptotic test and pairs test tend to overreject severely when n is small, although their performance improves quickly as n increases. Note that the pairs test becomes more computer intensive when the bootstrap method is used. In contrast, the bootstrap tests overreject only very slightly on average for all values of R^2 . The rejection frequencies are always very close to 0.05. Overall, the bootstrap tests outperform the asymptotic tests and pairs tests, especially when the sample

Figure 4: Rejection Frequency for Three Tests



size is small. One drawback of the bootstrap tests is that they require more computation time than the asymptotic tests.

Simulation results also suggest that asymptotic tests and pairs tests have similar performance. The performance of both tests are sensitive to different values of R^2 . When $R^2 = 0.25$, rejection frequencies yielded by the pairs test are closer to 0.05 compared with the asymptotic test. As R^2 increases, the performance of both tests improves. When $R^2 = 0.75$, both tests yield similar rejection frequencies for most values of n . In the extreme case of $R^2 = 0.99$, we see that the asymptotic test has better performance than the pairs test by yielding rejection frequencies closer to 0.05 for most values of n .

5 Empirical Application

In this section, we apply our MAPC estimator to the economic growth models in Barro (1991), which proposed a group of non-nested approximation models to analyze the relation between economic growth rate and a number of variables in a cross section of countries. The models considered by Barro (1991) are a small subset of the models we consider. The main results in Barro (1991) are summarized in Table 1.

We create our data set based on the Barro and Lee (1994) updated growth data for 98 countries. Variables that are used in our estimation are GR6085, GDP60, SEC60, PRIM60,

Table 1: A Summary of Barro’s (1991) Arguments

- The growth rate of real per capital GDP (GR6085) is positively related to the *initial human capital*^a
- GR6085 is negatively related to the initial (1960) level of real per capita GDP (GDP60)
- GR6085 is negatively related to the share of government consumption in GDP (g^c/y)
- GR6085 is positively related to the share of total investment on GDP (i/y)
- GR6085 is negatively related to the measures of *political instability*^b
- GR6085 is negatively related to the proxy for *market distortions*^c

^a Proxied by 1960 primary (PRIM60) and secondary (SEC60) school-enrollment rates.

^b Proxied by the number of revolutions and coups per year (REV) and the number per million population of political assassinations per year (AS).

^c Proxied by the magnitude of PPPI60 deviation (PPI60DEV), where PPPI60 is the 1960 PPP ratio based on the investment deflator.

g^c/y , REV, AS, PPI60DEV, GDP60SQ, RPRI, RSEC, AFRICA, LAT.AMER, i/y , and FERTNET. Means and standard deviations are reported in Appendix B, along with detailed definitions for all of the variables. There are eight variables that are of primary interest to us (core variables, including a constant term) and the other seven variables are only of marginal interest to us (potential variables).

We first revisit Barro's (1991) approach using our updated data set. The average growth rate between 1960–1985 (GR6085) is the common dependent variable. There are eight models that are summarized from Barro (1991). Table 2 shows regression results of these models, which are labeled from [1] to [8]. The smallest model is model [1], which only contains the

Table 2: Regressions for Per Capita Growth

	[1]	[2]	[3]	[4]	[5]	[6]	[7]	[8]
Const.	0.0287 *(0.0076)	0.0304 (0.0076)	0.0368 (0.0103)	0.0355 (0.0103)	0.0332 (0.0069)	0.0218 (0.0071)	0.0449 (0.0119)	0.0438 (0.0111)
GDP60	-0.0473 (0.0091)	-0.0826 (0.0284)	-0.0495 (0.0093)	-0.0490 (0.0093)	-0.0435 (0.0080)	-0.0489 (0.0083)	-0.0554 (0.0085)	-0.0501 (0.0080)
SEC60	0.0331 (0.0108)	0.0352 (0.0109)	0.0321 (0.0108)	0.0306 (0.0109)	0.0107 (0.0103)	0.0213 (0.0102)	0.0127 (0.0106)	0.0013 (0.0102)
PRIM60	0.0202 (0.0071)	0.0237 (0.0076)	0.0196 (0.0072)	0.0187 (0.0072)	0.0249 (0.0069)	0.0071 (0.0072)	0.0037 (0.0071)	0.0110 (0.0075)
g^c/y	-0.1071 (0.0286)	-0.1102 (0.0286)	-0.1060 (0.0286)	-0.1047 (0.0286)	-0.0800 (0.0260)	-0.1104 (0.0261)	-0.1071 (0.0255)	-0.0853 (0.0244)
REV	-0.0185 (0.0067)	-0.0183 (0.0067)	-0.0197 (0.0068)	-0.0202 (0.0068)	-0.0144 (0.0060)	-0.0129 (0.0062)	-0.0123 (0.0061)	-0.0111 (0.0057)
AS	-0.0441 (0.0184)	-0.0430 (0.0183)	-0.0440 (0.0183)	-0.0428 (0.0183)	-0.0242 (0.0170)	-0.0401 (0.0168)	-0.0375 (0.0164)	-0.0254 (0.0160)
PPI60DEV	-0.0074 (0.0036)	-0.0080 (0.0037)	-0.0071 (0.0036)	-0.0073 (0.0037)	-0.0082 (0.0033)	-0.0045 (0.0034)	-0.0051 (0.0033)	-0.0061 (0.0032)
GDP60SQ	-	0.0377 (0.0287)	-	-	-	-	-	-
RPRI	-	-	-0.1929 (0.1646)	-0.2637 (0.1770)	-	-	-	-
RSEC	-	-	-	0.2452 (0.2267)	-	-	-	-
AFRICA	-	-	-	-	-0.0141 (0.0037)	-	-	-0.0114 (0.0035)
LAT.AMER	-	-	-	-	-0.0171 (0.0035)	-	-	-0.0126 (0.0035)
i/y	-	-	-	-	-	0.0955 (0.0218)	0.0866 (0.0216)	0.0640 (0.0208)
FERTNET	-	-	-	-	-	-	-0.0035 (0.0015)	-0.0025 (0.0014)
R_c^2	0.4679	0.4780	0.4760	0.4829	0.6015	0.5624	0.5888	0.6578

* Values in parentheses are standard errors.

eight core variables, while the largest is model [8], which contains twelve variables (eight cores and four potentials). We also report the standard error associated with each estimated coefficient and the centered R^2 for each model.

One problem with Barro’s (1991) approach is the huge fluctuation of estimates (coefficient and P value) for the core variable across different models using different potential variables. Take SEC60 for an example. The estimated coefficient in regression [2] is 0.0352 with a standard error of 0.0109, which means that SEC60 is highly significant. However, the estimated coefficient of the same variable drops to 0.0013 in regression [8] with a standard error of 0.0102, which implies that SEC60 is highly insignificant. By using different models (potential variables), we may end up with contradictory results.

To solve this model uncertainty problem, we conduct model average estimation using the MAPC estimator. Following Barro (1991), we include the constant term, GDP60, SEC60, PRIM60, g^c/y , REV, AS, and PPI60DEV in every approximation model as the core variables and consider every possible combination of the potential variables. Since there are seven such variables, the total number of combinations is $2^7 = 128$. Therefore, we include 128 approximation models in our estimation.

In Table 3, we report the estimated average coefficients for the core variables and their standard errors. We also estimate the P value for each coefficient equal to zero using the t -test we derived in Section 3. We also calculate the bootstrap P value for each average coefficient. We set $B = 9999$ for all variables. We also show the top eight models with associated weights in Table 4.

The result of the MAPC estimation provides significant evidence to support Barro’s (1991) argument from the perspective of model average estimation. All estimates yield the same signs as Barro’s, which implies similar relations between the growth rate and potential variables. The estimated coefficient for GDP60 is -0.0713, which indicates a negative relationship between growth rate and the initial per capita product. We notice that the magnitude of this estimate is larger than those estimated by Barro’s (1991) models (except model [2]). Both SEC60 and PRIM60 are proxies for the initial human capital. Positive signs on these proxies indicate a positive relation between the growth rate and the initial human capital. The estimate on g^c/y is negative, which indicates a negative relation between the

Table 3: Results for the MAPC Estimation

	Estimates	s.e.	Asymptotic	Bootstrap
Const.	0.0438	0.0089	0.0000	0.0002
GDP60	-0.0713	0.0148	0.0000	0.0100
SEC60	0.0070	0.0097	0.4715	0.5445
PRIM60	0.0136	0.0069	0.0483	0.0911
g^c/y	-0.0902	0.0240	0.0002	0.0020
REV	-0.0122	0.0056	0.0289	0.0361
AS	-0.0281	0.0156	0.0716	0.0861
PPI60DEV	-0.0064	0.0031	0.0376	0.0501

Table 4: Top Eight Models in MAPC Estimation

	Assigned Weights							
	0.0814	0.0795	0.0792	0.0626	0.0547	0.0522	0.0518	0.0504
Const.	+	+	+	+	+	+	+	+
GDP60	+	+	+	+	+	+	+	+
SEC60	+	+	+	+	+	+	+	+
PRIM60	+	+	+	+	+	+	+	+
g^c/y	+	+	+	+	+	+	+	+
REV	+	+	+	+	+	+	+	+
AS	+	+	+	+	+	+	+	+
PPI60DEV	+	+	+	+	+	+	+	+
GDP60SQ	-	+	-	+	+	+	+	-
RPRI	-	-	-	-	-	-	+	-
RSEC	-	-	-	-	-	-	-	-
AFRICA	+	+	+	+	+	-	+	+
LAT.AMER	+	+	+	+	-	-	+	+
i/y	+	+	+	+	+	+	+	-
FERTNET	-	+	+	-	+	+	+	+

The “+” indicates that the corresponding parameters are included in the model and the “-” indicates the opposite.

growth rate and the share of government consumption in GDP. Political instability, which is proxied by REV and AS, reflects negative effects on the growth rate.

We also compare the relative out-of-sample predictive efficiency of the MAPC estimator with other estimators. For the original data sample $n = 98$, we shuffle the sample into a training set of n_1 and an evaluation set of size $n_2 = n - n_1$. We apply the seven methods (GETS, AIC, MA-AIC, MA-BIC, MMA, JMA, and MAPC) in Section 4 to the training set. We evaluate the selected models and computed estimates via mean squared prediction error (MSPE):

$$\text{MSPE} = \frac{1}{n_2} (\mathbf{y}_2 - \mathbf{X}_2 \hat{\boldsymbol{\beta}}_1)^\top (\mathbf{y}_2 - \mathbf{X}_2 \hat{\boldsymbol{\beta}}_1),$$

where $(\mathbf{y}_2, \mathbf{X}_2)$ is the evaluation set, n_2 is the number of observations of the evaluation set, and $\hat{\boldsymbol{\beta}}_1$ is the estimated coefficients by a particular method based on the training set. We normalize the MSPEs by the MSPE by the MAPC estimator. We repeat the entire procedure 10001 times and report the median MSPE. We vary n_1 and consider $n_1 = 20, 30, \dots, 80$. Table 5 reports the relative out-of-sample predictive efficiency. Entries larger than one indicate inferior performance relative to the MAPC estimator.

Table 5 suggests that the MAPC estimator delivers models that have better out-of-sample predictive efficiency than those by the six existing methods. As n_1 increases, the relative performance for all six methods improves. We notice that the MA-AIC method works better than MMA. As mentioned in Hansen and Racine (2012), the MMA method is somewhat sensitive to the preliminary estimate of σ^2 needed for its computation, and relying on a “large” approximating model may not be sufficient to deliver optimal results. In contrast,

Table 5: Relative Out-of-sample Predictive Efficiency

n_1	GETS	AIC	MA-AIC	MA-BIC	MMA	JMA
20	1.1808	1.0168	1.0052	1.0290	1.0111	1.0100
30	1.1801	1.0154	1.0049	1.0284	1.0101	1.0098
40	1.1680	1.0148	1.0046	1.0283	1.0100	1.0096
50	1.1665	1.0139	1.0046	1.0281	1.0099	1.0095
60	1.1659	1.0137	1.0032	1.0281	1.0097	1.0094
70	1.1642	1.0131	1.0024	1.0274	1.0096	1.0094
80	1.1630	1.0109	1.0023	1.0271	1.0093	1.0070

our MAPC estimator does not require a preliminary estimate, as it is a continuous-updating estimator.

6 Conclusion

Our MAPC estimator computes the weighted average across all approximation models; in doing so, it reduces the chance of a poor model being selected. The eight models selected from Barro (1991) are a small subset of the 128 models in the MAPC estimation. The idea of considering every possible approximation model helps avoid the risk of omitting important variables or retaining irrelevant variables. Both simulation and application suggest that as a continuous-updating estimator, our MAPC estimator yields better efficiency than the two-steps estimator MMA, especially when sample sizes are small. The asymptotic test statistics we derived in Section 3 work well in large sample sizes and bootstrap tests are highly recommended for small sample sizes.

A Proof

Proof of Lemma 1. To prove (i), we make use of $\text{Tr}(\mathbf{P}^{(m)}) = k^{(m)}$, then

$$\text{Tr}(\mathbf{P}(\mathbf{w})) = \text{Tr}\left(\sum_{m=1}^M w^{(m)} \mathbf{P}^{(m)}\right) = \sum_{m=1}^M w^{(m)} \text{Tr}(\mathbf{P}^{(m)}) = \sum_{m=1}^M w^{(m)} k^{(m)} = k(\mathbf{w}).$$

To prove (ii), we first note that projection matrix $\mathbf{P}^{(m)}$ is idempotent, then for an $n \times 1$ vector $\boldsymbol{\eta}$

$$\max_{\boldsymbol{\eta}} \frac{\boldsymbol{\eta}^\top \mathbf{P}^{(m)} \boldsymbol{\eta}}{\boldsymbol{\eta}^\top \boldsymbol{\eta}} = 1 \quad \text{for all } m.$$

By the definition of eigenvalue, we have

$$\lambda_{\max}(\mathbf{P}(\mathbf{w})) = \max_{\boldsymbol{\eta}} \frac{\boldsymbol{\eta}^\top \mathbf{P}(\mathbf{w}) \boldsymbol{\eta}}{\boldsymbol{\eta}^\top \boldsymbol{\eta}} \leq \sum_{m=1}^M w^{(m)} \left(\max_{\boldsymbol{\eta}} \frac{\boldsymbol{\eta}^\top \mathbf{P}^{(m)} \boldsymbol{\eta}}{\boldsymbol{\eta}^\top \boldsymbol{\eta}} \right) = 1.$$

Parts (iii) and (iv) can be obtained using result (ii) such that

$$\|\mathbf{P}(\mathbf{w}^*) \mathbf{M}(\mathbf{w}) \boldsymbol{\mu}\|^2 \leq \lambda_{\max}^2(\mathbf{P}(\mathbf{w}^*)) \|\mathbf{M}(\mathbf{w}) \boldsymbol{\mu}\|^2 \leq \|\mathbf{M}(\mathbf{w}) \boldsymbol{\mu}\|^2$$

and

$$\text{Tr}\left[(\mathbf{P}(\mathbf{w}^*) \mathbf{P}(\mathbf{w}^*)) \mathbf{P}(\mathbf{w}) \mathbf{P}(\mathbf{w})\right] \leq \lambda_{\max}^2(\mathbf{P}(\mathbf{w}^*)) \text{Tr}[\mathbf{P}(\mathbf{w}) \mathbf{P}(\mathbf{w})] \leq \text{Tr}[\mathbf{P}(\mathbf{w}) \mathbf{P}(\mathbf{w})].$$

■

Proof of Lemma 2. $R_n(\mathbf{w})$ can be written as

$$\begin{aligned} R_n(\mathbf{w}) &= \mathbb{E} \left[(\mathbf{P}(\mathbf{w}) \boldsymbol{\mu} + \mathbf{P}(\mathbf{w}) \mathbf{u} - \boldsymbol{\mu})^\top (\mathbf{P}(\mathbf{w}) \boldsymbol{\mu} + \mathbf{P}(\mathbf{w}) \mathbf{u} - \boldsymbol{\mu}) \mid \mathbf{X} \right] \\ &= \boldsymbol{\mu}^\top \mathbf{M}(\mathbf{w}) \mathbf{M}(\mathbf{w}) \boldsymbol{\mu} + \sigma^2 \text{Tr}(\mathbf{P}(\mathbf{w}) \mathbf{P}(\mathbf{w})) \\ &= \|\mathbf{M}(\mathbf{w}) \boldsymbol{\mu}\|^2 + \sigma^2 \text{Tr}(\mathbf{P}(\mathbf{w}) \mathbf{P}(\mathbf{w})). \end{aligned} \tag{20}$$

Both terms in (20) are non-negative, which implies Lemma 2. ■

Proof of Lemma 3. The proof is straightforward. If $k^{(m)}/n \rightarrow 0$ for all m , then the linear combination of $k^{(m)}$, $k(\mathbf{w}) = \sum_{m=1}^M w^{(m)} (k^{(m)}/n) \rightarrow 0$. ■

Proof of Theorem 1. Our proof follows the techniques derived in Li (1987) and Wan et

al. (2010). Rewrite $\text{MAPC}_n(\mathbf{w})$ to contain $L_n(\mathbf{w})$,

$$\begin{aligned} \text{MAPC}_n(\mathbf{w}) &= (\boldsymbol{\mu} - \boldsymbol{\mu}(\mathbf{w}) + \mathbf{u})^\top (\boldsymbol{\mu} - \boldsymbol{\mu}(\mathbf{w}) + \mathbf{u}) \left(\frac{n + k(\mathbf{w})}{n - k(\mathbf{w})} \right) \\ &= L_n(\mathbf{w}) + 2\boldsymbol{\mu}^\top \mathbf{M}(\mathbf{w})\mathbf{u} + 2L_n(\mathbf{w}) \left(\frac{k(\mathbf{w})}{n - k(\mathbf{w})} \right) \\ &\quad + 4\boldsymbol{\mu}^\top \mathbf{M}(\mathbf{w})\mathbf{u} \left(\frac{k(\mathbf{w})}{n - k(\mathbf{w})} \right) + \mathbf{u}^\top \mathbf{u} \left(\frac{n + k(\mathbf{w})}{n - k(\mathbf{w})} \right). \end{aligned} \quad (21)$$

As $n \rightarrow \infty$, $(n + k(\mathbf{w})) / (n - k(\mathbf{w})) \rightarrow 1$ by Lemma 3, which implies that the last term in (21) is independent of \mathbf{w} . Therefore, $\hat{\mathbf{w}}$ minimizes

$$L_n(\mathbf{w}) + 2\boldsymbol{\mu}^\top \mathbf{M}(\mathbf{w})\mathbf{u} + 2L_n(\mathbf{w}) \left(\frac{k(\mathbf{w})}{n - k(\mathbf{w})} \right) + 4\boldsymbol{\mu}^\top \mathbf{M}(\mathbf{w})\mathbf{u} \left(\frac{k(\mathbf{w})}{n - k(\mathbf{w})} \right). \quad (22)$$

If we can show that as $n \rightarrow \infty$, all terms in (22) except $L_n(\mathbf{w})$ are negligible compared with $L_n(\mathbf{w})$, for any $\mathbf{w} \in \mathbf{H}_M$, then the asymptotic optimality of $\hat{\mathbf{w}}$ is established.

Theorem 1 is valid, if, as $n \rightarrow \infty$

$$\sup_{\mathbf{w} \in \mathbf{H}_M} \left| \frac{\boldsymbol{\mu}^\top \mathbf{M}(\mathbf{w})\mathbf{u}}{R_n(\mathbf{w})} \right| \rightarrow_p 0, \quad (23)$$

$$\sup_{\mathbf{w} \in \mathbf{H}_M} \left| \frac{L_n(\mathbf{w})}{R_n(\mathbf{w})} - 1 \right| \rightarrow_p 0, \quad (24)$$

$$\sup_{\mathbf{w} \in \mathbf{H}_M} \left| \frac{\boldsymbol{\mu}^\top \mathbf{M}(\mathbf{w})\mathbf{u}}{R_n(\mathbf{w})} \left(\frac{k(\mathbf{w})}{n - k(\mathbf{w})} \right) \right| \rightarrow_p 0, \quad (25)$$

$$\sup_{\mathbf{w} \in \mathbf{H}_M} \left| \frac{L_n(\mathbf{w})}{R_n(\mathbf{w})} \left(\frac{k(\mathbf{w})}{n - k(\mathbf{w})} \right) \right| \rightarrow_p 0. \quad (26)$$

We shall prove (23) first. Given any $\delta > 0$, by triangular inequality, Bonferroni's inequality we have

$$\begin{aligned} \mathbb{P} \left\{ \sup_{\mathbf{w} \in \mathbf{H}_M} \left| \frac{\boldsymbol{\mu}^\top \mathbf{M}(\mathbf{w})\mathbf{u}}{R_n(\mathbf{w})} \right| > \delta \right\} &\leq \mathbb{P} \left\{ \sup_{\mathbf{w} \in \mathbf{H}_M} \sum_{m=1}^M w^{(m)} \left| \boldsymbol{\mu}^\top \mathbf{M}^{(m)}\mathbf{u} \right| > \delta \xi_n \right\} \\ &= \mathbb{P} \left\{ \max_m \left| \boldsymbol{\mu}^\top \mathbf{M}^{(m)}\mathbf{u} \right| > \delta \xi_n \right\} \\ &\leq \sum_{m=1}^M \mathbb{P} \left\{ \left| \boldsymbol{\mu}^\top \mathbf{M}(\mathbf{w}_m^0)\mathbf{u} \right| > \delta \xi_n \right\}, \end{aligned}$$

which, by Chebyshev's inequality, is no greater than

$$\sum_{m=1}^M \mathbb{E} \left\{ \frac{(\boldsymbol{\mu}^\top \mathbf{M}(\mathbf{w}_m^0) \mathbf{u})^{2G}}{\delta^{2G} \xi_n^{2G}} \right\}.$$

By Theorem 2 of Whittle (1960) and our Lemma 2(i), we observe that

$$\begin{aligned} \sum_{m=1}^M \mathbb{E} \left\{ \frac{(\boldsymbol{\mu}^\top \mathbf{M}(\mathbf{w}_m^0) \mathbf{u})^{2G}}{\delta^{2G} \xi_n^{2G}} \right\} &\leq C_1 \delta^{-2G} \xi_n^{-2G} \sum_{m=1}^M \|\mathbf{M}(\mathbf{w}_m^0) \boldsymbol{\mu}\|^{2G} \\ &\leq C_1 \delta^{-2G} \xi_n^{-2G} \sum_{m=1}^M (R_n(\mathbf{w}_m^0))^G \end{aligned}$$

for some constant C_1 . The last term above goes to zero by Assumption 2. Thus, (23) is proved.

To prove (24), we first see

$$\begin{aligned} &\sup_{\mathbf{w} \in \mathbf{H}_M} \left| \frac{L_n(\mathbf{w})}{R_n(\mathbf{w})} - 1 \right| \rightarrow_p 0 \\ \Leftrightarrow &\sup_{\mathbf{w} \in \mathbf{H}_M} \left| \frac{\mathbf{u}^\top \mathbf{P}(\mathbf{w}) \mathbf{P}(\mathbf{w}) \mathbf{u} - \sigma^2 \text{Tr}(\mathbf{P}(\mathbf{w}) \mathbf{P}(\mathbf{w})) - 2\boldsymbol{\mu}^\top \mathbf{M}(\mathbf{w}) \mathbf{P}(\mathbf{w}) \mathbf{u}}{R_n(\mathbf{w})} \right| \rightarrow_p 0. \end{aligned}$$

Then, it suffices to prove

$$\sup_{\mathbf{w} \in \mathbf{H}_M} \left| \frac{\boldsymbol{\mu}^\top \mathbf{M}(\mathbf{w}) \mathbf{P}(\mathbf{w}) \mathbf{u}}{R_n(\mathbf{w})} \right| \rightarrow_p 0 \quad (27)$$

and

$$\sup_{\mathbf{w} \in \mathbf{H}_M} \left| \frac{\mathbf{u}^\top \mathbf{P}(\mathbf{w}) \mathbf{P}(\mathbf{w}) \mathbf{u} - \sigma^2 \text{Tr}(\mathbf{P}(\mathbf{w}) \mathbf{P}(\mathbf{w}))}{R_n(\mathbf{w})} \right| \rightarrow_p 0. \quad (28)$$

By Chebyshev's inequality and Theorem 2 of Whittle (1960), given any $\delta > 0$,

$$\begin{aligned} \mathbb{P} \left\{ \sup_{\mathbf{w} \in \mathbf{H}_M} \left| \frac{\boldsymbol{\mu}^\top \mathbf{M}(\mathbf{w}) \mathbf{P}(\mathbf{w}) \mathbf{u}}{R_n(\mathbf{w})} \right| > \delta \right\} &\leq \sum_{m=1}^M \sum_{l=1}^M \mathbb{E} \left\{ \frac{(\boldsymbol{\mu}^\top \mathbf{M}(\mathbf{w}_m^0) \mathbf{P}(\mathbf{w}_l^0) \mathbf{u})^{2G}}{\delta^{2G} \xi_n^{2G}} \right\} \\ &\leq C_2 \delta^{-2G} \xi_n^{-2G} \sum_{m=1}^M \sum_{l=1}^M (\boldsymbol{\mu}^\top \mathbf{M}(\mathbf{w}_m^0) \mathbf{P}(\mathbf{w}_l^0) \mathbf{M}(\mathbf{w}_m^0) \boldsymbol{\mu})^G \\ &= C_2 \delta^{-2G} \xi_n^{-2G} \sum_{m=1}^M \sum_{l=1}^M \|\mathbf{P}(\mathbf{w}_l^0) \mathbf{M}(\mathbf{w}_m^0) \boldsymbol{\mu}\|^{2G}, \end{aligned}$$

where C_2 is a constant. By Lemma 1 (iii), Lemma 2 (i) and Assumption 2, we obtain

$$\begin{aligned} C_2 \delta^{-2G} \xi_n^{-2G} \sum_{m=1}^M \sum_{l=1}^M \|\mathbf{P}(\mathbf{w}_l^0) \mathbf{M}(\mathbf{w}_m^0) \boldsymbol{\mu}\|^{2G} &\leq C_2 \delta^{-2G} \xi_n^{-2G} \sum_{m=1}^M \sum_{l=1}^M \|\mathbf{M}(\mathbf{w}_m^0) \boldsymbol{\mu}\|^{2G} \\ &\leq C_2 \delta^{-2G} \xi_n^{-2G} M \sum_{m=1}^M (R_n(\mathbf{w}_m^0))^G \rightarrow 0. \end{aligned}$$

Likewise, by Chebyshev's inequality, Theorem 2 of Whittle (1960), Lemma 1 (iii), Lemma 2 (ii) and Assumption 2, given any $\delta > 0$, we observe that, for a constant C_3 ,

$$\begin{aligned} &\mathbb{P} \left\{ \sup_{\mathbf{w} \in \mathbf{H}_M} \left| \frac{\mathbf{u}^\top \mathbf{P}(\mathbf{w}) \mathbf{P}(\mathbf{w}) \mathbf{u} - \sigma^2 \text{Tr}(\mathbf{P}(\mathbf{w}) \mathbf{P}(\mathbf{w}))}{R_n(\mathbf{w})} \right| > \delta \right\} \\ &\leq \sum_{m=1}^M \sum_{l=1}^M \mathbb{E} \left\{ \frac{[\mathbf{u}^\top \mathbf{P}(\mathbf{w}_l^0) \mathbf{P}(\mathbf{w}_m^0) \mathbf{u} - \sigma^2 \text{Tr}(\mathbf{P}(\mathbf{w}_l^0) \mathbf{P}(\mathbf{w}_m^0))]^{2G}}{\delta^{2G} \xi_n^{2G}} \right\} \\ &\leq C_3 \delta^{-2G} \xi_n^{-2G} \sum_{m=1}^M \sum_{l=1}^M [\text{Tr}(\mathbf{P}(\mathbf{w}_m^0) \mathbf{P}(\mathbf{w}_l^0) \mathbf{P}(\mathbf{w}_m^0))]^G \\ &\leq C_3 \delta^{-2G} \xi_n^{-2G} \sum_{m=1}^M \sum_{l=1}^M [\text{Tr}(\mathbf{P}(\mathbf{w}_m^0) \mathbf{P}(\mathbf{w}_m^0))]^G \\ &\leq \frac{C_3}{\sigma^2} \delta^{-2G} \xi_n^{-2G} M \sum_{m=1}^M [R_n(\mathbf{w}_m^0)]^G \rightarrow 0. \end{aligned}$$

Proving (25) and (26) becomes straightforward once we validate (23) and (24). We obtain, as $n \rightarrow \infty$, that

$$\sup_{\mathbf{w} \in \mathbf{H}_M} \left| \frac{\boldsymbol{\mu}^\top \mathbf{M}(\mathbf{w}) \mathbf{u}}{R_n(\mathbf{w})} \cdot \frac{k(\mathbf{w})}{n - k(\mathbf{w})} \right| \leq \sup_{\mathbf{w} \in \mathbf{H}_M} \left| \frac{\boldsymbol{\mu}^\top \mathbf{M}(\mathbf{w}) \mathbf{u}}{R_n(\mathbf{w})} \right| \cdot \sup_{\mathbf{w} \in \mathbf{H}_M} \left| \frac{k(\mathbf{w})}{n - k(\mathbf{w})} \right| \rightarrow_p 0$$

and

$$\sup_{\mathbf{w} \in \mathbf{H}_M} \left| \frac{L_n(\mathbf{w})}{R_n(\mathbf{w})} \cdot \frac{k(\mathbf{w})}{n - k(\mathbf{w})} \right| \leq \sup_{\mathbf{w} \in \mathbf{H}_M} \left| \frac{L_n(\mathbf{w})}{R_n(\mathbf{w})} \right| \cdot \sup_{\mathbf{w} \in \mathbf{H}_M} \left| \frac{k(\mathbf{w})}{n - k(\mathbf{w})} \right| \rightarrow_p 0.$$

This completes the proof of Theorem 1. ■

Proof of Theorem 2. Since $\hat{\mathbf{u}}^{(m)} = \mathbf{y} - \hat{\boldsymbol{\mu}}^{(m)} = \mathbf{y} - \mathbf{P}^{(m)} \mathbf{y} = \mathbf{M}^{(m)} \mathbf{u} + \mathbf{M}^{(m)} \boldsymbol{\mu}$ for model m , the average estimate of \mathbf{u} becomes

$$\mathbf{u}(\mathbf{w}) = \sum_{m=1}^M w^{(m)} (\mathbf{y} - \hat{\boldsymbol{\mu}}^{(m)}) = \sum_{m=1}^M w^{(m)} \mathbf{M}^{(m)} (\boldsymbol{\mu} + \mathbf{u}) = \mathbf{M}(\mathbf{w}) (\boldsymbol{\mu} + \mathbf{u}).$$

Then, the average estimate for σ^2 becomes

$$\begin{aligned}\sigma^2(\mathbf{w}) &= \frac{\mathbf{u}(\mathbf{w})^\top \mathbf{u}(\mathbf{w})}{n - k(\mathbf{w})} = \frac{(\boldsymbol{\mu} + \mathbf{u})^\top \mathbf{M}(\mathbf{w}) \mathbf{M}(\mathbf{w}) (\boldsymbol{\mu} + \mathbf{u})}{n - k(\mathbf{w})} \\ &= \frac{\mathbf{u}^\top \mathbf{M}(\mathbf{w}) \mathbf{M}(\mathbf{w}) \mathbf{u}}{n - k(\mathbf{w})} + \frac{\boldsymbol{\mu}^\top \mathbf{M}(\mathbf{w}) \mathbf{M}(\mathbf{w}) \boldsymbol{\mu}}{n - k(\mathbf{w})} + \frac{2\boldsymbol{\mu}^\top \mathbf{M}(\mathbf{w}) \mathbf{M}(\mathbf{w}) \mathbf{u}}{n - k(\mathbf{w})}.\end{aligned}\quad (29)$$

Therefore, Theorem 2 is valid if the following hold: as $n \rightarrow \infty$,

$$\sup_{\mathbf{w} \in \mathbf{H}_M} \left| \frac{\mathbf{u}^\top \mathbf{M}(\mathbf{w}) \mathbf{M}(\mathbf{w}) \mathbf{u}}{n - k(\mathbf{w})} \right| \rightarrow_p \sigma^2, \quad (30)$$

$$\sup_{\mathbf{w} \in \mathbf{H}_M} \left| \frac{\boldsymbol{\mu}^\top \mathbf{M}(\mathbf{w}) \mathbf{M}(\mathbf{w}) \boldsymbol{\mu}}{n - k(\mathbf{w})} \right| \rightarrow_p 0, \quad (31)$$

$$\sup_{\mathbf{w} \in \mathbf{H}_M} \left| \frac{\boldsymbol{\mu}^\top \mathbf{M}(\mathbf{w}) \mathbf{M}(\mathbf{w}) \mathbf{u}}{n - k(\mathbf{w})} \right| \rightarrow_p 0. \quad (32)$$

Equation (30) is equivalent to

$$\sup_{\mathbf{w} \in \mathbf{H}_M} \left| \frac{\mathbf{u}^\top \mathbf{M}(\mathbf{w}) \mathbf{u}}{n - k(\mathbf{w})} - \frac{\mathbf{u}^\top \mathbf{P}(\mathbf{w}) \mathbf{M}(\mathbf{w}) \mathbf{u}}{n - k(\mathbf{w})} \right| \rightarrow_p \sigma^2. \quad (33)$$

To prove (30), it suffices to show, as $n \rightarrow \infty$, that

$$\sup_{\mathbf{w} \in \mathbf{H}_M} \left| \frac{\mathbf{u}^\top \mathbf{M}(\mathbf{w}) \mathbf{u}}{n - k(\mathbf{w})} \right| \rightarrow_p \sigma^2 \quad (34)$$

and

$$\sup_{\mathbf{w} \in \mathbf{H}_M} \left| \frac{\mathbf{u}^\top \mathbf{P}(\mathbf{w}) \mathbf{M}(\mathbf{w}) \mathbf{u}}{n - k(\mathbf{w})} \right| \rightarrow_p 0. \quad (35)$$

First, because $\mathbb{E}(\mathbf{u}^\top \mathbf{M}(\mathbf{w}) \mathbf{u}) = \sigma^2 (n - k(\mathbf{w}))$, by Theorem 2 of Whittle (1960),

$$\begin{aligned}\mathbb{E} \left| \mathbf{u}^\top \mathbf{M}(\mathbf{w}) \mathbf{u} - \sigma^2 (n - k(\mathbf{w})) \right|^2 &\leq C_4 \text{Tr}(\mathbf{M}(\mathbf{w}) \mathbf{M}(\mathbf{w})) \\ &\leq C_4 \text{Tr}(\mathbf{M}(\mathbf{w})) \\ &= C_4 (n - k(\mathbf{w})),\end{aligned}$$

where C_4 is some constant. Thus, for any $\delta > 0$, by Markov's inequality and Lemma 3,

$$\begin{aligned}\mathbb{P} \left\{ \sup_{\mathbf{w} \in \mathbf{H}_M} \left| \frac{\mathbf{u}^\top \mathbf{M}(\mathbf{w}) \mathbf{u}}{n - k(\mathbf{w})} - \sigma^2 \right| > \delta \right\} &\leq \frac{\mathbb{E} \left| \mathbf{u}^\top \mathbf{M}(\mathbf{w}) \mathbf{u} - \sigma^2 (n - k(\mathbf{w})) \right|^2}{\delta^2 (n - k(\mathbf{w}))^2} \\ &\leq \frac{C_4}{\delta^2 (n - k(\mathbf{w}))} \rightarrow 0.\end{aligned}$$

Second,

$$\begin{aligned}
\mathbb{P} \left\{ \sup_{\mathbf{w} \in \mathbf{H}_M} \left| \frac{\mathbf{u}^\top \mathbf{P}(\mathbf{w}) \mathbf{M}(\mathbf{w}) \mathbf{u}}{n - k(\mathbf{w})} \right| > \delta \right\} &\leq \mathbb{P} \left\{ \sup_{\mathbf{w} \in \mathbf{H}_M} \left| \frac{\mathbf{u}^\top \mathbf{P}(\mathbf{w}) \mathbf{u}}{n - k(\mathbf{w})} \right| > \delta \right\} \\
&= \mathbb{P} \left\{ \sup_{\mathbf{w} \in \mathbf{H}_M} \sum_{m=1}^M w^{(m)} \left| \frac{\mathbf{u}^\top \mathbf{P}^{(m)} \mathbf{u}}{n - k(\mathbf{w})} \right| > \delta \right\} \\
&\leq \mathbb{P} \left\{ \sup_{\mathbf{w} \in \mathbf{H}_M} \max_m \left| \frac{\mathbf{u}^\top \mathbf{P}^{(m)} \mathbf{u}}{n - k(\mathbf{w})} \right| > \delta \right\} \rightarrow 0
\end{aligned}$$

since $(n - k(\mathbf{w}))^{-1} (\mathbf{u}^\top \mathbf{P}^{(m)} \mathbf{u}) \rightarrow_p 0$ as $n \rightarrow \infty$ for all m . (30) is obtained.

To prove (31), we see that for any approximation model m ,

$$\mathbb{E} \left[\frac{\left(\mathbf{b}^{(m)} \right)^\top \mathbf{M}^{(m)} \mathbf{b}^{(m)}}{n - k(\mathbf{w})} \right] \leq \mathbb{E} \left[\frac{\left(\mathbf{b}^{(m)} \right)^\top \mathbf{b}^{(m)}}{n - k(\mathbf{w})} \right] \leq \frac{n}{n - k(\mathbf{w})} \mathbb{E} \left(b_i^{(m)} \right)^2 \rightarrow 0$$

since $k^{(m)} \rightarrow \infty$ as $n \rightarrow \infty$ and the square integrability of $\mu_i^{(m)}$ implies $\mathbb{E} \left(b_i^{(m)} \right)^2 \rightarrow 0$ as $k^{(m)} \rightarrow \infty$. This implies

$$\frac{\left(\mathbf{b}^{(m)} \right)^\top \mathbf{M}^{(m)} \mathbf{b}^{(m)}}{n - k(\mathbf{w})} \rightarrow_p 0 \tag{36}$$

for all m . Therefore, by (36), we observe, for any $\delta > 0$, that

$$\begin{aligned}
\mathbb{P} \left\{ \sup_{\mathbf{w} \in \mathbf{H}_M} \left| \frac{\boldsymbol{\mu}^\top \mathbf{M}(\mathbf{w}) \mathbf{M}(\mathbf{w}) \boldsymbol{\mu}}{n - k(\mathbf{w})} \right| > \delta \right\} &\leq \mathbb{P} \left\{ \sup_{\mathbf{w} \in \mathbf{H}_M} \left| \frac{\boldsymbol{\mu}^\top \mathbf{M}(\mathbf{w}) \boldsymbol{\mu}}{n - k(\mathbf{w})} \right| > \delta \right\} \\
&= \mathbb{P} \left\{ \sup_{\mathbf{w} \in \mathbf{H}_M} \sum_{m=1}^M w^{(m)} \left| \frac{\boldsymbol{\mu}^\top \mathbf{M}^{(m)} \boldsymbol{\mu}}{n - k(\mathbf{w})} \right| > \delta \right\} \\
&\leq \mathbb{P} \left\{ \sup_{\mathbf{w} \in \mathbf{H}_M} \max_m \left| \frac{\boldsymbol{\mu}^\top \mathbf{M}^{(m)} \boldsymbol{\mu}}{n - k(\mathbf{w})} \right| > \delta \right\} \\
&= \mathbb{P} \left\{ \sup_{\mathbf{w} \in \mathbf{H}_M} \max_m \left| \frac{\left(\mathbf{b}^{(m)} \right)^\top \mathbf{M}^{(m)} \left(\mathbf{b}^{(m)} \right)}{n - k(\mathbf{w})} \right| > \delta \right\} \rightarrow 0.
\end{aligned}$$

Finally, since $(n - k(\mathbf{w}))^{-1} \left(\boldsymbol{\mu}^\top \mathbf{M}^{(m)} \mathbf{u} \right) \rightarrow_p 0$ for all model m , we obtain

$$\begin{aligned} \mathbb{P} \left\{ \sup_{\mathbf{w} \in \mathbf{H}_M} \left| \frac{\boldsymbol{\mu}^\top \mathbf{M}(\mathbf{w}) \mathbf{M}(\mathbf{w}) \mathbf{u}}{n - k(\mathbf{w})} \right| > \delta \right\} &\leq \mathbb{P} \left\{ \sup_{\mathbf{w} \in \mathbf{H}_M} \left| \frac{\boldsymbol{\mu}^\top \mathbf{M}(\mathbf{w}) \mathbf{u}}{n - k(\mathbf{w})} \right| > \delta \right\} \\ &\leq \mathbb{P} \left\{ \sup_{\mathbf{w} \in \mathbf{H}_M} \max_m \left| \frac{\boldsymbol{\mu}^\top \mathbf{M}^{(m)} \mathbf{u}}{n - k(\mathbf{w})} \right| > \delta \right\} \rightarrow_p 0. \end{aligned}$$

Therefore, we conclude that $\sigma^2(\mathbf{w}) \rightarrow_p \sigma^2$. ■

Proof of Lemma 4. Since the covariance matrix is conditional on $\mathbf{X}^{(L)}$, we can assume that $\mathbf{X}^{(m)}$ is exogenous for all m . The covariance matrix can be written as

$$\begin{aligned} &\text{Cov} \left(\Gamma^{(m)} \hat{\boldsymbol{\beta}}^{(m)}, \Gamma^{(s)} \hat{\boldsymbol{\beta}}^{(s)} \right) \\ &= \mathbb{E} \left[\left(\Gamma^{(m)} \hat{\boldsymbol{\beta}}^{(m)} - \boldsymbol{\beta}(\mathbf{w}) \right) \left(\Gamma^{(s)} \hat{\boldsymbol{\beta}}^{(s)} - \boldsymbol{\beta}(\mathbf{w}) \right)^\top \right] \\ &= \mathbb{E} \left(\left(\Gamma^{(m)} (\hat{\boldsymbol{\beta}}^{(m)} - \boldsymbol{\beta}^{(m)}) + \mathbf{d}_1^{(m)} \right) \left(\Gamma^{(s)} (\hat{\boldsymbol{\beta}}^{(s)} - \boldsymbol{\beta}^{(s)}) + \mathbf{d}_1^{(s)} \right)^\top \right), \end{aligned} \quad (37)$$

where $\mathbf{d}_1^{(m)} = \Gamma^{(m)} \boldsymbol{\beta}^{(m)} - \boldsymbol{\beta}(\mathbf{w})$ is non-random. First, we have

$$\begin{aligned} &\Gamma^{(m)} \left(\hat{\boldsymbol{\beta}}^{(m)} - \boldsymbol{\beta}^{(m)} \right) \\ &= \Gamma^{(m)} \left(\left(\mathbf{X}^{(m)\top} \mathbf{X}^{(m)} \right)^{-1} \mathbf{X}^{(m)\top} \left(\mathbf{X}^{(m)} \boldsymbol{\beta}^{(m)} + \mathbf{X}^{(-m)} \boldsymbol{\beta}^{(-m)} + \mathbf{u} \right) - \boldsymbol{\beta}^{(m)} \right) \\ &= \Gamma^{(m)} \left(\mathbf{X}^{(m)\top} \mathbf{X}^{(m)} \right)^{-1} \mathbf{X}^{(m)\top} \mathbf{X}^{(-m)} \boldsymbol{\beta}^{(-m)} + \Gamma^{(m)} \left(\mathbf{X}^{(m)\top} \mathbf{X}^{(m)} \right)^{-1} \Gamma^{(m)\top} \mathbf{X}^\top \mathbf{u}. \end{aligned}$$

Define

$$\begin{aligned} \mathbf{A}^{(m)} &\equiv \Gamma^{(m)} \left(\mathbf{X}^{(m)\top} \mathbf{X}^{(m)} \right)^{-1} \mathbf{X}^{(m)\top} \mathbf{X}^{(-m)} \boldsymbol{\beta}^{(-m)}, \\ \mathbf{B}^{(m)} &\equiv \Gamma^{(m)} \left(\mathbf{X}^{(m)\top} \mathbf{X}^{(m)} \right)^{-1} \Gamma^{(m)\top} \mathbf{X}^\top \mathbf{u}. \end{aligned}$$

We expand the expectation in (37) and obtain

$$\begin{aligned} &\mathbb{E} \left(\mathbf{A}^{(m)} \mathbf{A}^{(s)\top} + \mathbf{A}^{(m)} \mathbf{B}^{(s)\top} + \mathbf{A}^{(m)} \mathbf{d}_1^{(s)\top} + \mathbf{B}^{(m)} \mathbf{A}^{(s)\top} + \mathbf{B}^{(m)} \mathbf{B}^{(s)\top} \right. \\ &\quad \left. + \mathbf{B}^{(m)} \mathbf{d}_1^{(s)\top} + \mathbf{d}_1^{(m)} \mathbf{A}^{(s)\top} + \mathbf{d}_1^{(m)} \mathbf{B}^{(s)\top} + \mathbf{d}_1^{(m)} \mathbf{d}_1^{(s)\top} \right). \end{aligned} \quad (38)$$

Note that $\mathbb{E}(\mathbf{B}^{(m)}) = \mathbf{0}$ since $\mathbb{E}(\mathbf{u}|\mathbf{X}) = 0$. This implies that $\mathbb{E}(\mathbf{A}^{(m)} \mathbf{B}^{(s)\top}) = \mathbf{0}$, $\mathbb{E}(\mathbf{B}^{(m)} \mathbf{A}^{(s)\top}) =$

$\mathbf{0}$, $\mathbb{E}(\mathbf{B}^{(m)} \mathbf{d}_1^{(s)\top}) = \mathbf{0}$, and $\mathbb{E}(\mathbf{d}_1^{(m)} \mathbf{B}^{(s)\top}) = \mathbf{0}$. Also,

$$\mathbb{E} \left(\mathbf{B}^{(m)} \mathbf{B}^{(s)\top} \right) = \sigma^2 \mathbf{\Gamma}^{(m)} \left(\mathbf{X}^{(m)\top} \mathbf{X}^{(m)} \right)^{-1} \mathbf{X}^{(m)\top} \mathbf{X}^{(s)} \left(\mathbf{X}^{(s)\top} \mathbf{X}^{(s)} \right)^{-1} \mathbf{\Gamma}^{(s)\top}.$$

Finally, (38) is equivalent to

$$\begin{aligned} & \sigma^2 \mathbf{\Gamma}^{(m)} \left(\mathbf{X}^{(m)\top} \mathbf{X}^{(m)} \right)^{-1} \mathbf{X}^{(m)\top} \mathbf{X}^{(s)} \left(\mathbf{X}^{(s)\top} \mathbf{X}^{(s)} \right)^{-1} \mathbf{\Gamma}^{(s)\top} \\ & + \left(\mathbf{A}^{(m)} + \mathbf{d}_1^{(m)} \right) \left(\mathbf{A}^{(s)} + \mathbf{d}_1^{(s)} \right)^\top. \end{aligned}$$

Note that $\mathbf{A}^{(m)} + \mathbf{d}_1^{(m)} = \mathbb{E} \left(\mathbf{\Gamma}^{(m)} \hat{\boldsymbol{\beta}}^{(m)} \right) - \boldsymbol{\beta}(\mathbf{w}) = \mathbf{d}^{(m)}$. In practice, we replace the infeasible σ^2 and $\mathbf{d}^{(m)}$ with their sample estimates. This completes the proof. \blacksquare

Proof of Lemma 5. Part (i) follows the fact that the set of $\boldsymbol{\beta}_{-m}$ belongs to the set of $\boldsymbol{\beta}_p$ for all m . Therefore, $\sqrt{n} \boldsymbol{\beta}_{-m} \rightarrow \mathbf{h}_{-m}$ by Assumption 5.

To prove part (ii), we first expand the estimator $\hat{\boldsymbol{\beta}}^{(m)}$ around $\boldsymbol{\beta}^{(m)}$ and obtain

$$\begin{aligned} & \sqrt{n} \left(\hat{\boldsymbol{\beta}}^{(m)} - \boldsymbol{\beta}^{(m)} \right) \\ & = \sqrt{n} \left(\mathbf{X}^{(m)\top} \mathbf{X}^{(m)} \right)^{-1} \mathbf{X}^{(m)\top} \mathbf{X}^{(-m)} \boldsymbol{\beta}^{(-m)} + \sqrt{n} \left(\mathbf{X}^{(m)\top} \mathbf{X}^{(m)} \right)^{-1} \mathbf{\Gamma}^{(m)\top} \mathbf{X}^\top \mathbf{u}. \end{aligned}$$

By Assumption 5 and Lemma 5(i),

$$\sqrt{n} \left(\mathbf{X}^{(m)\top} \mathbf{X}^{(m)} \right)^{-1} \mathbf{X}^{(m)\top} \mathbf{X}^{(-m)} \boldsymbol{\beta}^{(-m)} \xrightarrow{p} \left(\mathbf{S}^{(m)} \right)^{-1} \mathbf{S}^{(m,-m)} \mathbf{h}_{-m} \equiv \boldsymbol{\delta}_1^{(m)}.$$

By Assumption 4 and central limit theorem,

$$\sqrt{n} \left(\mathbf{X}^{(m)\top} \mathbf{X}^{(m)} \right)^{-1} \mathbf{\Gamma}^{(m)\top} \mathbf{X}^\top \mathbf{u} \xrightarrow{d} \mathbf{N} \left(\mathbf{0}, \sigma^2 \left(\mathbf{S}^{(m)} \right)^{-1} \right).$$

This implies $\sqrt{n} (\hat{\boldsymbol{\beta}}_c^{(m)} - \boldsymbol{\beta}_c) = \sqrt{n} \mathbf{\Gamma}_c^{(m)\top} (\hat{\boldsymbol{\beta}}^{(m)} - \boldsymbol{\beta}^{(m)}) \xrightarrow{d} \mathbf{N}(\boldsymbol{\delta}_1^{(m)}, \sigma^2 \mathbf{\Gamma}_c^{(m)\top} (\mathbf{S}^{(m)})^{-1} \mathbf{\Gamma}_c^{(m)})$.

To prove part (iii), we first have $\boldsymbol{\beta}_{-m} = O(n^{-1/2})$ by Assumption 5. Under Assumption 6, we apply Taylor expansion to $\boldsymbol{\beta}(\mathbf{w})$ and obtain

$$\boldsymbol{\beta}(\mathbf{w}) = \mathbf{f}(\boldsymbol{\beta}^{(m)}, \boldsymbol{\beta}_{-m}) = \mathbf{f}(\boldsymbol{\beta}^{(m)}, \mathbf{0}) + \left(\mathbf{F}_{\boldsymbol{\beta}_{-m}} \big|_{\boldsymbol{\beta}_{-m}=\mathbf{0}} \right) \boldsymbol{\beta}_{-m} + \frac{1}{2} \boldsymbol{\beta}_{-m}^\top \bar{\mathbf{H}} \boldsymbol{\beta}_{-m}, \quad (39)$$

where $\bar{\mathbf{H}}$ is the Hessian matrix $\mathbf{H} \equiv \partial^2 \mathbf{f}(\cdot) / \partial \boldsymbol{\beta}_{-m} \partial \boldsymbol{\beta}_{-m}^\top$ evaluated at an intermediate point between $\mathbf{0}$ and $\boldsymbol{\beta}_{-m}$. Since $\boldsymbol{\beta}_{-m} = O(n^{-1/2})$, the last term in (39) is $O(n^{-1})$. With

$f(\boldsymbol{\beta}^{(m)}, \mathbf{0}) = \boldsymbol{\Gamma}^{(m)}\boldsymbol{\beta}^{(m)}$, we have

$$\begin{aligned}\sqrt{n} \left(\boldsymbol{\Gamma}^{(m)}\boldsymbol{\beta}^{(m)} - \boldsymbol{\beta}(\mathbf{w}) \right) &= \sqrt{n} \left(- \left(\mathbf{F}_{\boldsymbol{\beta}_{-m}} \big|_{\boldsymbol{\beta}_{-m}=\mathbf{0}} \right) \boldsymbol{\beta}_{-m} - O(n^{-1}) \right) \\ &\rightarrow - \left(\mathbf{F}_{\boldsymbol{\beta}_{-m}} \big|_{\boldsymbol{\beta}_{-m}=\mathbf{0}} \right) \mathbf{h}_{-m}.\end{aligned}$$

Therefore, $\sqrt{n} (\boldsymbol{\beta}_c - \boldsymbol{\beta}_c(\mathbf{w})) = \sqrt{n} \boldsymbol{\Gamma}_c^{(m)\top} \left(\boldsymbol{\Gamma}^{(m)}\boldsymbol{\beta}^{(m)} - \boldsymbol{\beta}(\mathbf{w}) \right) \rightarrow \boldsymbol{\delta}_2$.

Part (iv) follows Lemma 5 (ii) and (iii). We have

$$\begin{aligned}\sqrt{n} \mathbf{d}^{(m)} &= \sqrt{n} \left(\mathbb{E}(\boldsymbol{\Gamma}^{(m)}\hat{\boldsymbol{\beta}}^{(m)}) - \boldsymbol{\beta}(\mathbf{w}) \right) \\ &\xrightarrow{p} \boldsymbol{\Gamma}^{(m)} \left(\mathbf{S}^{(m)} \right)^{-1} \mathbf{S}^{(m,-m)} \mathbf{h}_{-m} - \left(\mathbf{F}_{\boldsymbol{\beta}_{-m}} \big|_{\boldsymbol{\beta}_{-m}=\mathbf{0}} \right) \mathbf{h}_{-m} = \boldsymbol{\delta}^{(m)}.\end{aligned}$$

Part (v) is the main result of Lemma 5.

$$\begin{aligned}\sqrt{n} \left(\hat{\boldsymbol{\beta}}_c^{(m)} - \boldsymbol{\beta}_c(\mathbf{w}) \right) &= \sqrt{n} \left(\hat{\boldsymbol{\beta}}_c^{(m)} - \boldsymbol{\beta}_c \right) + \sqrt{n} (\boldsymbol{\beta}_c - \boldsymbol{\beta}_c(\mathbf{w})) \\ &\xrightarrow{d} \mathbf{N} \left(\boldsymbol{\delta}_c^{(m)}, \sigma^2 \boldsymbol{\Gamma}_c^{(m)\top} \left(\mathbf{S}^{(m)} \right)^{-1} \boldsymbol{\Gamma}_c^{(m)} \right) \sim \boldsymbol{\Lambda}_c^{(m)}.\end{aligned} \quad \blacksquare$$

Proof of Theorem 3. It is straightforward to show that the asymptotic distribution of $\boldsymbol{\Lambda}_c$ is a normal distribution, since

$$\sqrt{n} \left(\hat{\boldsymbol{\beta}}_c(\hat{\mathbf{w}}) - \boldsymbol{\beta}_c(\mathbf{w}) \right) \big|_{\hat{\mathbf{w}}} = \sum_{m=1}^M \hat{w}^{(m)} \sqrt{n} \left(\hat{\boldsymbol{\beta}}_c^{(m)} - \boldsymbol{\beta}_c(\mathbf{w}) \right) \big|_{\hat{\mathbf{w}}} \rightarrow_d \sum_{m=1}^M \hat{w}^{(m)} \boldsymbol{\Lambda}_c^{(m)}$$

is a linear combination of normal distribution $\boldsymbol{\Lambda}_c^{(m)}$ for all m . Given $\mathbb{E} \left(\hat{\boldsymbol{\beta}}_c(\hat{\mathbf{w}}) - \boldsymbol{\beta}_c(\mathbf{w}) \right) = \mathbf{0}$, the mean of $\boldsymbol{\Lambda}_c$ is $\mathbf{0}$. For the variance of $\boldsymbol{\Lambda}_c$, we have $\mathbf{V}_c(\hat{\mathbf{w}}) = \boldsymbol{\Gamma}_c^{(m)\top} \mathbf{V}(\hat{\mathbf{w}}) \boldsymbol{\Gamma}_c^{(m)}$, where

$$\begin{aligned}\mathbf{V}(\hat{\mathbf{w}}) &= \text{plim}_{n \rightarrow \infty} \text{Var} \left(\sqrt{n} \left(\hat{\boldsymbol{\beta}}(\hat{\mathbf{w}}) - \boldsymbol{\beta}(\mathbf{w}) \right) \big|_{\hat{\mathbf{w}}} \right) \\ &= \text{plim}_{n \rightarrow \infty} n \sum_{m=1}^M \sum_{s=1}^M \hat{w}^{(m)} \hat{w}^{(s)} \left(\mathbf{d}^{(m)} \mathbf{d}^{(s)\top} \right) \\ &\quad + \sigma^2 \boldsymbol{\Gamma}^{(m)} \left(\mathbf{X}^{(m)\top} \mathbf{X}^{(m)} \right)^{-1} \mathbf{X}^{(m)\top} \mathbf{X}^{(s)} \left(\mathbf{X}^{(s)\top} \mathbf{X}^{(s)} \right)^{-1} \boldsymbol{\Gamma}^{(s)\top} \big|_{\hat{\mathbf{w}}} \\ &= \text{plim}_{n \rightarrow \infty} \sum_{m=1}^M \sum_{s=1}^M \hat{w}^{(m)} \hat{w}^{(s)} \left(\sqrt{n} \mathbf{d}^{(m)} \sqrt{n} \mathbf{d}^{(s)\top} \right) \\ &\quad + \sigma^2 \boldsymbol{\Gamma}^{(m)} \left(n^{-1} \mathbf{X}^{(m)\top} \mathbf{X}^{(m)} \right)^{-1} n^{-1} \mathbf{X}^{(m)\top} \mathbf{X}^{(s)} \left(n^{-1} \mathbf{X}^{(s)\top} \mathbf{X}^{(s)} \right)^{-1} \boldsymbol{\Gamma}^{(s)\top}.\end{aligned}$$

We have proved in Lemma 5 that $\sqrt{n}\mathbf{d}^{(m)} \xrightarrow{p} \boldsymbol{\delta}^{(m)}$. Therefore,

$$\mathbf{V}(\hat{\mathbf{w}}) = \sum_{m=1}^M \sum_{s=1}^M \hat{w}^{(m)} \hat{w}^{(s)} \left(\sigma^2 \boldsymbol{\Gamma}^{(m)} (\mathbf{S}^{(m)})^{-1} \mathbf{S}^{(m,s)} (\mathbf{S}^{(s)})^{-1} \boldsymbol{\Gamma}^{(s)\top} + \boldsymbol{\delta}^{(m)} \boldsymbol{\delta}^{(s)\top} \right),$$

where $\mathbf{S}^{(m)}$, $\mathbf{S}^{(m,s)}$, and $\boldsymbol{\delta}^{(m)}$ are defined in Lemma 5. ■

Proof of Theorem 4. Since $\sqrt{n}(\hat{\boldsymbol{\beta}}_c(\hat{\mathbf{w}}) - \boldsymbol{\beta}_c(\mathbf{w})) | \hat{\mathbf{w}} \rightarrow_d \mathbf{N}(\mathbf{0}, \mathbf{V}_c(\hat{\mathbf{w}}))$ and $\mathbf{V}_c(\hat{\mathbf{w}})$ is a linear function of $\hat{\mathbf{w}}$, it is equivalent to have $\sqrt{n}(\hat{\boldsymbol{\beta}}_c(\hat{\mathbf{w}}) - \boldsymbol{\beta}_c(\mathbf{w})) | \mathbf{V}_c(\hat{\mathbf{w}}) \rightarrow_d \mathbf{N}(\mathbf{0}, \mathbf{V}_c(\hat{\mathbf{w}}))$.

We can rewrite the t -statistic as

$$t_{\beta_j(\mathbf{w})} \equiv \frac{\hat{\beta}_j(\hat{\mathbf{w}}) - \beta_{j0}}{\widehat{\text{Var}}(\hat{\beta}_j(\hat{\mathbf{w}}))^{1/2}} = \frac{\sqrt{n} \left(\hat{\beta}_j(\hat{\mathbf{w}}) - \beta_{j0} \right)}{\left(n \widehat{\text{Var}}(\hat{\beta}_j(\hat{\mathbf{w}})) \right)^{1/2}}.$$

Following Theorem 3, for fixed $\mathbf{V}_c(\hat{\mathbf{w}})$, we have

$$\sqrt{n} \left(\hat{\beta}_j(\hat{\mathbf{w}}) - \beta_{j0} \right) = \sqrt{n} \left(\hat{\beta}_j(\hat{\mathbf{w}}) - \beta_{j0} \right) | \mathbf{V}_c(\hat{\mathbf{w}}) \rightarrow_d \Lambda_j$$

with mean 0 and variance $V_j(\hat{\mathbf{w}})$, which is the j^{th} element on the diagonal of $\mathbf{V}_c(\hat{\mathbf{w}})$, and

$$n \widehat{\text{Var}}(\hat{\beta}_j(\hat{\mathbf{w}})) = \widehat{\text{Var}} \left(\sqrt{n} (\hat{\beta}_j(\hat{\mathbf{w}}) - \beta_{j0}) | \mathbf{V}_c(\hat{\mathbf{w}}) \right) \rightarrow_p V_j(\hat{\mathbf{w}}).$$

Therefore, for fixed $\mathbf{V}_c(\hat{\mathbf{w}})$, the conditional asymptotic distribution of $t_{\beta_j(\mathbf{w})}$ is clearly $\mathbf{N}(0, 1)$. Since this conditional distribution does not depend on $\mathbf{V}_c(\hat{\mathbf{w}})$, it also holds marginally.

Similarly, for the Wald statistic, we have

$$W_{\mathbf{r}} = \left(\mathbf{R} \hat{\boldsymbol{\beta}}_c(\hat{\mathbf{w}}) - \mathbf{r} \right)^\top \left(\mathbf{R} \widehat{\text{Var}}(\hat{\boldsymbol{\beta}}_c(\hat{\mathbf{w}})) \mathbf{R}^\top \right)^{-1} \left(\mathbf{R} \hat{\boldsymbol{\beta}}_c(\hat{\mathbf{w}}) - \mathbf{r} \right) \rightarrow_d \chi^2(k_r)$$

for fixed $\mathbf{V}_c(\hat{\mathbf{w}})$. Since this conditional distribution does not depend on $\mathbf{V}_c(\hat{\mathbf{w}})$, it also holds marginally. ■

B Description of Data Set

	Variable	Mean	Std.dev
1.	GR6085	0.0219	0.0198
2.	GDP60	1.9146	1.8195
3.	SEC60	0.2314	0.2117
4.	PRIM60	0.7562	0.2710
5.	g^c/y	0.1080	0.0555
6.	REV	0.1806	0.2355
7.	AS	0.0279	0.0764
8.	PPPI60	0.8373	0.4716
9.	PPI60DEV	0.2884	0.3778
10.	GDP60SQ	6.9427	11.4995
11.	RPRI	35.7462	9.1282
12.	RSEC	19.3646	6.4440
13.	AFRICA	0.2653	0.4438
14.	LAT.AMER	0.2347	0.4260
15.	i/y	0.1912	0.0796
16.	FERT	4.8704	1.7520
17.	MORT01	0.0810	0.0503
18.	FERTNET	4.3946	1.4383

Definitions:

1. GR6085: Growth rate of real per capita GDP from 1960 to 1985.
2. GDP60: 1960 value of real per capita GDP.
3. SEC60: 1960 secondary-school enrollment rate.
4. PRIM60: 1960 primary-school enrollment rate.
5. g^c/y : Average from from 1970 to 1985 of the ratio of real government consumption (exclusive of defense and education) to real GDP.
6. REV: Number of revolutions and coups per year (1960-1985).
7. AS: Number of assassinations per million population per year (1960-1985).
8. PPPI60: 1960 PPP value for the investment deflator (U.S. = 1.0).
9. PPI60DEV: Magnitude of the deviation of PPI60 from the sample mean.
10. GDP60SQ: Square of GDP60.
11. RPRI: Student-teacher ratio in primary schools in 1960.

12. RSEC: Student-teacher ratio in secondary schools in 1960.
13. AFRICA: Dummy variable for sub-Saharan Africa.
14. LAT.AMER: Dummy variable for LATIN America.
15. i/y : Average from 1960-1985 of the ratio of real domestic investment to real GDP.
16. FERT: Total fertility rate, average of 1965 and 1985.
17. MORT01: Mortality rate for age 0 through 1, average of 1965 and 1985.
18. FERTNET: $FERT \times (1 - MORT01)$.

References

- AKAIKE, H. (1973): “Information theory and an extension of the maximum likelihood principle,” *Second International Symposium on Information Theory*, pp. 267–281.
- AMEMIYA, T. (1980): “Selection of Regressors,” *International Economic Review*, 21(2), 331–354.
- BARNARD, G. A. (1963): “New methods of quality control,” *Journal of Royal Statistical Society Series: A*, 126, 255.
- BARRO, R. J. (1991): “Economic Growth in a Cross Section of Countries,” *The Quarterly Journal of Economics*, 106(2), 407–443.
- BARRO, R. J., AND J.-W. LEE (1994): “Sources of economic growth,” *Carnegie-Rochester Conference Series on Public Policy*, 40(1), 1–46.
- BUCKLAND, S. T., K. P. BURNHAM, AND N. H. AUGUSTIN (1997): “Model Selection: An Integral Part of Inference,” *Biometrics*, 53(2), pp. 603–618.
- BURNHAM, K. P., AND D. R. ANDERSON (2002): *Model Selection and Multimodel Inference: A Practical Information-Theoretic Approach*, chap. 4, pp. 153–169. Springer.
- CLAESKENS, G., AND N. L. HJORT (2003): “The focused information criterion,” *Journal of the American Statistical Association*, 98, 900–916.
- CLYDE, M., AND E. I. GEORGE (2004): “Model uncertainty,” *Statistical Science*, 19, 81–94.
- DAVIDSON, R., AND J. MACKINNON (2000): “Bootstrap tests: how many bootstraps?,” *Econometric Reviews*, 19(1), 55–68.
- DI CICCIO, T. J., AND B. EFRON (1996): “Bootstrap Confidence Intervals,” *Statistical Science*, 11(3), pp. 189–212.
- DRAPER, D. (1995): “Assessment and propagation of model uncertainty,” *Journal of Royal Statistical Society Series: B*, 57, 45–97.
- FAN, J., AND R. LI (2001): “Variable Selection via Nonconcave Penalized Likelihood and Its Oracle Properties,” *Journal of the American Statistical Association*, 96, 1348–1360.
- HANSEN, B. E. (2007): “Least Squares Model Averaging,” *Econometrica*, 75(4), 1175–1189.
- (2008): “Least-squares forecast averaging,” *Journal of Econometrics*, 146, 342–350.
- (2009): “Averaging estimators for regressions with a possible structural break,” *Econometric Theory*, 25, 1498–1514.

- (2010): “Averaging estimators for autoregressions with a near unit root,” *Journal of Econometrics*, 158, 142–155.
- HANSEN, B. E., AND J. S. RACINE (2012): “Jackknife model averaging,” *Journal of Econometrics*, 167(1), 38 – 46.
- HENDRY, D. F. (1976a): “The structure of simultaneous equations estimators,” *Journal of Econometrics*, 4(1), 51–88.
- (1976b): “The Structure of Simultaneous Equations Estimators,” *Journal of Econometrics*, 4(1), 51–88.
- (1980): “Econometrics – Alchemy or Science?,” *Economica*, 47(188), 387–406.
- (1983): “Econometric Modeling: The ‘Consumption Function’ in Retrospect,” *Scottish Journal of Political Economy*, 30(3), 193–220.
- HENDRY, D. F., AND H.-M. KROLZIG (1999): “Improving on ‘Data Mining Reconsidered’ by K.D. Hoover and S.J. Perez,” *Econometrics Journal*, 2, 41–58.
- HENDRY, D. F., AND B. NIELSEN (2007): *Econometric Modeling: A Likelihood Approach*, chap. 19, pp. 286–301. Princeton University Press.
- HJORT, N. L., AND G. CLAESKENS (2003): “Frequentist Model Average Estimators,” *Journal of the American Statistical Association*, 98(464), 879–899.
- HOETING, J. A., D. MADIAN, A. E. RAFTERY, AND C. T. VOLINSKY (1999): “Bayesian Model Averaging: A Tutorial,” *Statistical Science*, 14(4), 382–417.
- HOOVER, K. D., AND S. J. PEREZ (1999): “Data mining reconsidered: encompassing and the general-to-specific approach to specification search,” *Econometrics Journal*, 2(2), 167–191.
- JOHANSEN, S. (1995): “13. Asymptotic Properties of the Estimators,” *Likelihood-Based Inference in Cointegrated Vector Autoregressive Models*, pp. 177–178.
- KASS, R. E., AND A. E. RAFTERY (1995): “Bayes Factors,” *Journal of the American Statistical Association*, 90(430), 773–794.
- KNOX, T., J. H. STOCK, AND M. W. WATSON (2004): “Empirical Bayes Regression With Many Regressors,” *working paper*.
- KROLZIG, H.-M., AND D. F. HENDRY (2001): “Computer automation of general-to-specific model selection procedures,” *Journal of Economic Dynamics & Control*, 25, 831–866.
- LEAMER, E. (1978): *Specification Searches*. Wiley, New York.

- LEY, E., AND M. F. STEEL (2009): “On the effect of prior assumptions in Bayesian model averaging with applications to growth regression,” *Journal of Applied Econometrics*, 24(4), 651–674.
- LI, K.-C. (1987): “Asymptotic Optimality for C_p , C_L , Cross-validation and Generalized Cross-validation: Discrete Index Set,” *The Annals of Statistics*, 15(3), 958–975.
- LIU, C., AND J. M. MAHEU (2009): “Forecasting realized volatility: a Bayesian model-averaging approach,” *Journal of Applied Econometrics*, 24(5), 709–733.
- MACKINNON, J. G. (2009): “Bootstrap Hypothesis Testing,” in *Handbook of Computational Econometrics*, ed. by D. A. Belsley, and J. Kontoghiorghes, chap. 6, pp. 183–213. Wiley.
- MALLOWS, C. L. (1973): “Some Comments on C_p ,” *Technometrics*, 15(4), 661–675.
- PAGAN, A. R. (1987): “Three Econometric Methodologies: A Critical Appraisal,” *Journal of Economic Surveys*, 1(1), 3–24.
- RAFTERY, A. E., D. MADIGAN, AND C. T. VOLINSKY (1997): “Bayesian model averaging for linear regression models,” *Journal of American statistical association*, 92, 179–191.
- SALA-I-MARTIN, X., G. DOPPELHOFER, AND R. MILLER (2004): “Determinants of Long-Term Growth: A Bayesian Averaging of Classical Estimates (BACE) Approach,” *American Economic Review*, 94, 813–835.
- SCHWARZ, G. (1978): “Estimating the Dimension of a Model,” *The Annals of Statistics*, 6, 461–464.
- SHIBATA, R. (1981): “An Optimal Selection of Regression Variables,” *Biometrika*, 68(1), 45–54.
- STAIGER, D., AND J. H. STOCK (1997): “Instrumental Variables Regression with Weak Instruments,” *Econometrica*, 65(3), 557–586.
- WAN, A. T., X. ZHANG, AND G. ZOU (2010): “Least Squares model averaging by Mallows criterion,” *Journal of Econometrics*, 156, 277–283.
- WELCH, B. L. (1938): “The Significance of the Difference Between Two Means when the Population Variances are Unequal,” *Biometrika*, 29(3/4), pp. 350–362.
- WHITTLE, P. R. (1960): “Bounds for the Moments of Linear and Quadratic Forms in Independent Variables,” *Theory of Probability and Its Applications*, 5, 302–305.
- WRIGHT, J. H. (2009): “Forecasting US inflation by Bayesian model averaging,” *Journal of Forecasting*, 28, 131–144.