

## WALS Prediction

Magnus, J.R.; Wang, W.; Zhang, Xinyu

*Publication date:*  
2012

[Link to publication](#)

*Citation for published version (APA):*

Magnus, J. R., Wang, W., & Zhang, X. (2012). *WALS Prediction*. (CentER Discussion Paper; Vol. 2012-043). *Econometrics*.

### General rights

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the public portal for the purpose of private study or research
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the public portal

### Take down policy

If you believe that this document breaches copyright, please contact us providing details, and we will remove access to the work immediately and investigate your claim.

No. 2012-043

**WALS PREDICTION**

By

Jan R. Magnus, Wendun Wang, Xinyu Zhang

May 30, 2012

ISSN 0924-7815

# WALS prediction\*

May 29, 2012

Jan R. Magnus  
*Department of Econometrics & Operations Research,  
Tilburg University*

Wendun Wang  
*CentER, Tilburg University*

Xinyu Zhang  
*Academy of Mathematics & Systems Science,  
Chinese Academy of Sciences, Beijing*

---

\*E-mail addresses: magnus@uvt.nl (Jan Magnus), wangwendun@gmail.com (Wendun Wang), xinyu@amss.ac.cn (Xinyu Zhang).

**Corresponding author:**

Jan R. Magnus  
Department of Econometrics & OR  
Tilburg University  
PO Box 90153  
5000 LE Tilburg  
The Netherlands  
E-mail: magnus@uvt.nl

**Abstract:** Prediction under model uncertainty is an important and difficult issue. Traditional prediction methods (such as pretesting) are based on model selection followed by prediction in the selected model, but the reported prediction and the reported prediction variance ignore the uncertainty from the selection procedure. This paper proposes a weighted average least squares (WALS) prediction procedure that is not conditional on the selected model. Taking both model and error uncertainty into account, we also propose an appropriate estimate of the variance of the WALS predictor. Correlations among the random errors are explicitly allowed. Compared to other prediction averaging methods, the WALS predictor has important advantages both theoretically and computationally. Simulation studies show that the WALS predictor generally produces lower mean squared prediction errors than its competitors, and that the proposed estimator for the prediction variance performs particularly well when model uncertainty increases.

**AMS 2010 subject classification:** 62F99, 62J05, 62M20.

**Key words:** Model averaging; Model uncertainty; Bayesian analysis; Prediction.

# 1 Introduction

Suppose a ruler seeks advice on a specific parameter, say next year's inflation. He has twelve advisors, and each advisor provides an estimate. When all have left, the ruler has twelve estimates. In addition, he has an opinion about each advisor based on past experience and current performance. How does the ruler now obtain a single estimate? Let us consider two possibilities. The ruler may think: Whom do I trust most? Whose advice do I think most reliable? Then, he takes the advice of his most trusted advisor. This is the first method. Alternatively, he may consider all advisors useful, but not to the same degree. Some are more experienced and more clever than others, so they get a higher weight. Then, the ruler computes a weighted average of the twelve estimates. This is the second method.

While the second method appeals to common sense, econometric practice favors the first method. In econometric practice one typically first selects the 'best' model based on diagnostic tests (such as  $t$ -ratios,  $R^2$ , and various information criteria) and then computes estimates within this selected model. This is called 'pretesting'. There are many problems with this procedure (Magnus, 1999; Danilov and Magnus, 2004a,b), but the most important is that model selection and estimation are completely separated—just like the ruler only listening to his most trusted advisor—so that uncertainty in the model selection is ignored when reporting properties of the estimates.

In the second method, called 'model averaging', we average over estimates from different models, instead of estimating based on a single selected model. Model averaging not only appeals to common sense, but also has two major advantages. First, it avoids arbitrary thresholds (like 1.96), thus forcing continuity on a previously discontinuous estimator; second, it allows us to combine model selection and estimation into *one* procedure, thus moving from conditional to unconditional estimator characteristics.

Much of the model averaging literature has concentrated on estimation rather than on prediction. In this paper we concentrate on prediction (forecasting), which may in fact be a more appropriate application of model averaging, because the interpretation of coefficients changes with different models but the predictor always has the same interpretation. A substantial literature on the averaging of forecasts exists, going back to Bates and Granger (1969); see Granger (2003), Yang (2004), Elliott and Timmermann (2004), and Aiolfi and Timmermann (2006) for some recent contributions, and Hendry and Clements (2004) and Timmermann (2006) for recent reviews. Simulation and empirical studies indicate that predictors based on a set of models generally perform better than predictors obtained from a single model (Stock and Watson, 2004; Jackson and Karlsson, 2004; Bjørnland et

al., 2012).

Our paper has two main contributions. First, we introduce the prediction counterpart to the weighted average least squares (WALS) estimator proposed in Magnus et al. (2010) and study its properties in simulations. The WALS procedure avoids some of the problems encountered in standard Bayesian model averaging (BMA). In particular, the prior is based on a coherent notion of ignorance, thus avoiding normality of the prior and unbounded risk. Also, the computational burden increases linearly rather than exponentially with the number of regressors, and is therefore trivial compared to other model averaging estimators such as standard BMA, model-selection-based weights methods (Buckland et al., 1997; Hjort and Claeskens, 2003), exponential reweighting (Yang, 2004), or Mallows model averaging (Hansen, 2007, 2008). Our proposed method explicitly allows for correlation in the observations, including possible correlation between the errors in the realized sample and the predictive sample.

The second contribution of the paper is that we propose an estimate for the prediction variance taking model uncertainty into account, and evaluate the accuracy of this estimate. The typical researcher’s instinct is to favor a predictor with a small variance over one with a large variance. We argue that what we require is not a small but a ‘correct’ variance: in a situation with much noise a predictor with a small variance can cause much harm, while a truthfully reported large variance may lead to more prudent policy. In fact, one of the problems with the credibility of econometric predictions may be that our reported prediction variances are too small, and this is caused, at least in part, by the fact that model uncertainty is ignored. We shall see that WALS predictions may lead to higher variances, but that these variances are closer to the truth.

The paper is organized as follows. Sections 2–7 develop the theory. In Section 2 we set up the model and present the traditional predictor. The commonly employed conditional predictor is presented in Section 3, and the WALS predictor in Section 4. In Section 5 we discuss the computation of the WALS predictor based on the Laplace prior. An estimator for the variance of the WALS predictor is proposed in Section 6. Finally, in Section 7, we discuss the estimation of unknown parameters in the variance matrix of the random disturbances. Then, in Sections 8–11, we compare the WALS predictor with its most important competitors: unrestricted maximum likelihood, pretesting, ridge regression, and Mallows model averaging. Our comparison is conducted through a large number of Monte Carlo simulation experiments, controlling for sample size, parameter values, and variance specifications. The simulation results show that the WALS predictor typically has the lowest mean squared prediction error among the predictors considered, and that

the more uncertainty exists in the model, the better is the relative performance of WALS. Section 12 concludes.

## 2 The traditional predictor

Our framework is the linear regression model

$$y = X\beta + u, \quad (1)$$

where  $y$  is a vector of  $N$  observations on the dependent variable,  $X$  ( $N \times k$ ) is a matrix of regressors,  $u$  is a random vector of  $N$  unobservable disturbances, and  $\beta$  is a vector of  $k$  unknown parameters. We assume throughout that  $1 \leq k \leq N - 1$  and that  $X$  has full column-rank  $k$ . We are interested in some specific (possibly future) values of the regressors  $X_f$  ( $N_f \times k$ ), and we wish to predict the value  $y_f$  ( $N_f \times 1$ ) likely to be associated with  $X_f$ . We assume that  $y_f$  is generated by

$$y_f = X_f\beta + u_f, \quad (2)$$

and our task is to find a predictor  $\hat{y}_f$  of  $y_f$ .

In general the observations will be correlated, and we shall assume that

$$\begin{pmatrix} u \\ u_f \end{pmatrix} \sim N \left( \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} \Omega & C_f' \\ C_f & \Omega_f \end{pmatrix} \right), \quad (3)$$

where the variance of  $(u, u_f)$  is a positive definite  $(N + N_f) \times (N + N_f)$  matrix, whose component blocks  $\Omega$ ,  $C_f$ , and  $\Omega_f$  are functions of an  $m$ -dimensional unknown parameter vector  $\theta = (\theta_1, \dots, \theta_m)'$ . To simplify notation we treat the regressors as fixed (at least for the moment), but the theory applies also to random regressors if we condition appropriately.

The joint distribution of  $u$  and  $u_f$  in (3) implies that

$$E(u_f|u) = C_f\Omega^{-1}u, \quad \text{var}(u_f|u) = \Omega_f - C_f\Omega^{-1}C_f', \quad (4)$$

so that

$$E(y_f|y) = X_f\beta + C_f\Omega^{-1}(y - X\beta). \quad (5)$$

This leads to the traditional least squares predictor in the presence of a non-scalar variance matrix:

$$\hat{y}_f = X_f\hat{\beta} + C_f\Omega^{-1}(y - X\hat{\beta}), \quad (6)$$

where  $\hat{\beta} = (X'\Omega^{-1}X)^{-1}X'\Omega^{-1}y$  denotes the generalized least squares (GLS) estimator of  $\beta$ , and it is assumed (for the moment) that  $\theta$  is known; see

Whittle (1963, p. 53, Eq. (10)) for the general formula, and Johnston and DiNardo (1997, Sec. 6.8) and Ruud (2000, Sec. 19.7) for the special case where  $N_f = 1$  and the errors follow an AR(1) process. The predictor (6) is normally distributed with mean  $E(\hat{y}_f) = X_f\beta$  and variance

$$\text{var}(\hat{y}_f) = X_f(X'\Omega^{-1}X)^{-1}X'_f + C_f(\Omega^{-1} - \Omega^{-1}X(X'\Omega^{-1}X)^{-1}X'\Omega^{-1})C'_f \quad (7)$$

from which we see *inter alia* that the presence of the covariance  $C_f$  increases the variance of the predictor, and therefore that ignoring correlation leads to misleadingly precise predictions.

The prediction error  $\text{PE} := \hat{y}_f - y_f$  can be conveniently written as the sum of two independent random variables:

$$\text{PE} = (X_f - C_f\Omega^{-1}X)(\hat{\beta} - \beta) - (u_f - C_f\Omega^{-1}u), \quad (8)$$

and the traditional predictor  $\hat{y}_f$  is a good predictor in the sense that it is unbiased and that the prediction error has minimum variance

$$\begin{aligned} \text{var}(\text{PE}) &= (X_f - C_f\Omega^{-1}X)(X'\Omega^{-1}X)^{-1}(X_f - C_f\Omega^{-1}X)' \\ &\quad + \Omega_f - C_f\Omega^{-1}C'_f \end{aligned} \quad (9)$$

in the class of linear unbiased estimators.

### 3 The conditional predictor

The previous section assumes that the data-generation process (DGP) and the model coincide, which one might call the ‘traditional’ approach. In practice, the model is likely to be (much) smaller than the DGP. In this section we shall assume that the model is a special case of the DGP obtained by setting some of the  $\beta$ -parameters equal to zero. We do not know in advance which  $\beta$ -parameters should be set to zero and we use model selection diagnostics (such as  $t$ - and  $F$ -statistics) to arrive at a model that we like. Once we have obtained this model we derive the properties of the predictor *conditional* on the selected model and hence we ignore the noise generated by the model selection process. We call this the ‘conditional’ approach. This is not quite right of course, and we shall present a method which combines model selection and prediction in the next section.

We distinguish between *focus* regressors  $X_1$  (those we want in the model on theoretical or other grounds) and *auxiliary* regressors  $X_2$  (those we are less certain of), and write model (1) accordingly as

$$y = X_1\beta_1 + X_2\beta_2 + u, \quad (10)$$



so that  $X = (X_1 : X_2)$  and  $\beta = (\beta_1', \beta_2')'$ . Let  $k_1 \geq 0$  be the dimension of  $\beta_1$  and  $k_2 \geq 0$  the dimension of  $\beta_2$ , so that  $k = k_1 + k_2$ . Model selection takes place over the auxiliary regressors only. Since each of the  $k_2$  auxiliary regressors can either be included or not, we have  $2^{k_2}$  models to consider.

In addition to the regressors that are always in the model ( $X_1$ ) and those that are sometimes in the model ( $X_2$ ), there are also regressors that are never in the model (say  $X_3$ ), even though they are in the DGP. This is because the modeler is ignorant about these regressors or has no access to the necessary data. We disregard this situation for the moment, but return to it in Section 8.

We assume (at first) that  $\theta$  and hence  $\Omega$  is known. It is convenient to semi-orthogonalize the regression model as follows. Let

$$M_1^* := \Omega^{-1} - \Omega^{-1}X_1(X_1'\Omega^{-1}X_1)^{-1}X_1'\Omega^{-1}, \quad (11)$$

where we notice that the matrix  $\Omega^{1/2}M_1^*\Omega^{1/2}$  is idempotent. Let  $P$  be an orthogonal matrix and  $\Lambda$  a diagonal matrix with positive diagonal elements such that  $P'X_2'M_1^*X_2P = \Lambda$ . Next define the transformed auxiliary regressors and the transformed auxiliary parameters as

$$X_2^* := X_2P\Lambda^{-1/2}, \quad \beta_2^* := \Lambda^{1/2}P'\beta_2. \quad (12)$$

Then  $X_2^*\beta_2^* = X_2\beta_2$ , so that we can write (10) equivalently as

$$y = X_1\beta_1 + X_2^*\beta_2^* + u. \quad (13)$$

The result of this transformation is that the new design matrix  $(X_1 : X_2^*)$  is 'semi-orthogonal' in the sense that  $X_2^{*'}M_1^*X_2^* = I_{k_2}$  and this has important advantages that will become clear shortly.

### 3.1 Estimation in model $\mathcal{M}_i$

Our strategy will be to estimate  $(\beta_1, \beta_2^*)$  rather than  $(\beta_1, \beta_2)$ . Each of the  $k_2$  components of  $\beta_2^*$  can either be included or not included in the model and this gives rise to  $2^{k_2}$  models. A specific model is identified through a  $k_2 \times (k_2 - k_{2i})$  selection matrix  $S_i$  of full column-rank, where  $0 \leq k_{2i} \leq k_2$ , so that  $S_i' = (I_{k_2 - k_{2i}} : 0)$  or a column-permutation thereof. Our first interest is in the GLS estimator of  $(\beta_1, \beta_2^*)$  in the  $i$ -th model, that is, in the GLS estimator of  $(\beta_1, \beta_2^*)$  under the restriction  $S_i'\beta_2^* = 0$ .

Let  $\mathcal{M}_i$  represent model (13) under the restriction  $S_i'\beta_2^* = 0$ , and let  $\hat{\beta}_{1(i)}$  and  $\hat{\beta}_{2(i)}^*$  denote the GLS estimators of  $\beta_1$  and  $\beta_2^*$  under  $\mathcal{M}_i$ . Extending

Danilov and Magnus (2004a, Lemmas A1 and A2), the GLS estimators of  $\beta_1$  and  $\beta_2^*$  under  $\mathcal{M}_i$  may be written as (see also Magnus et al., 2011):

$$\hat{\beta}_{1(i)} = (X_1' \Omega^{-1} X_1)^{-1} X_1' \Omega^{-1} y - Q^* W_i b_2^*, \quad \hat{\beta}_{2(i)}^* = W_i b_2^*, \quad (14)$$

respectively, where

$$b_2^* := X_2^{*'} M_1^* y, \quad Q^* := (X_1' \Omega^{-1} X_1)^{-1} X_1' \Omega^{-1} X_2^*, \quad W_i := I_{k_2} - S_i S_i'. \quad (15)$$

Note that  $b_2^*$  is simply the GLS estimator of  $\beta_2^*$  in the unrestricted model, and that  $W_i$  is a diagonal  $k_2 \times k_2$  matrix with  $k_{2i}$  ones and  $(k_2 - k_{2i})$  zeros on the diagonal. The  $j$ -th diagonal element of  $W_i$  equals zero if  $\beta_{2j}^*$  (the  $j$ -th component of  $\beta_2^*$ ) is restricted to zero, and equals one otherwise. If  $k_{2i} = k_2$  then  $W_i = I_{k_2}$ . The diagonality of  $W_i$  is a direct consequence of the semi-orthogonal transformation.

The distributions of  $\hat{\beta}_{1(i)}$  and  $\hat{\beta}_{2(i)}^*$  are then

$$\hat{\beta}_{1(i)} \sim N_{k_1} (\beta_1 + Q^* S_i S_i' \beta_2^*, (X_1' \Omega^{-1} X_1)^{-1} + Q^* W_i Q^{*'}), \quad (16)$$

$$\hat{\beta}_{2(i)}^* \sim N_{k_2} (W_i \beta_2^*, W_i), \quad (17)$$

and  $\text{cov}(\hat{\beta}_{1(i)}, \hat{\beta}_{2(i)}^*) = -Q^* W_i$ . The residual vector  $e_i := y - X_1 \hat{\beta}_{1(i)} - X_2^* \hat{\beta}_{2(i)}^*$  is given by  $e_i = \Omega D_i^* y$ , where  $D_i^* := M_1^* - M_1^* X_2^* W_i X_2^{*'} M_1^*$  and  $\Omega^{1/2} D_i^* \Omega^{1/2}$  is a symmetric idempotent matrix of rank  $n - k_1 - k_{2i}$ . It follows that:

- all models that include the  $j$ -th column of  $X_2^*$  as a regressor have the same estimator of  $\beta_{2j}^*$ , namely  $b_{2j}^*$ , irrespective of which other columns of  $X_2^*$  are included;
- the estimators  $b_{21}^*, b_{22}^*, \dots, b_{2k_2}^*$  are independent; and
- the residuals of the  $i$ -th model  $\mathcal{M}_i$  depend on  $y$  only through  $M_1^* y$ .

### 3.2 Prediction in model $\mathcal{M}_i$

Next we wish to predict  $N_f$  (possibly future) values  $y_f$ , based on values of the regressors  $X_{1f}$  ( $N_f \times k_1$ ) and  $X_{2f}$  ( $N_f \times k_2$ ). Corresponding to  $X_2^*$  we define  $X_{2f}^* := X_{2f} P \Lambda^{-1/2}$ , so that

$$\begin{pmatrix} y \\ y_f \end{pmatrix} = \begin{pmatrix} X_1 & X_2^* \\ X_{1f} & X_{2f}^* \end{pmatrix} \begin{pmatrix} \beta_1 \\ \beta_2^* \end{pmatrix} + \begin{pmatrix} u \\ u_f \end{pmatrix}, \quad (18)$$

where the errors  $(u, u_f)$  are distributed as in (3). From (5) we obtain

$$E(y_f | y) = X_{1f} \beta_1 + X_{2f}^* \beta_2^* + C_f \Omega^{-1} (y - X_1 \beta_1 - X_2^* \beta_2^*), \quad (19)$$

leading to the predictor in model  $\mathcal{M}_i$ , using (14),

$$\begin{aligned}\hat{y}_f^{(i)} &= X_{1f}\hat{\beta}_{1(i)} + X_{2f}^*\hat{\beta}_{2(i)} + C_f\Omega^{-1}(y - X_1\hat{\beta}_{1(i)} - X_2^*\hat{\beta}_{2(i)}) \\ &= X_{1f}(X_1'\Omega^{-1}X_1)^{-1}X_1'\Omega^{-1}y + C_fM_1^*y + Z_fW_ib_2^*,\end{aligned}\quad (20)$$

where

$$Z_f := (X_{2f}^* - X_{1f}Q^*) - C_f\Omega^{-1}(X_2^* - X_1Q^*). \quad (21)$$

**Theorem 1:** The predictor  $\hat{y}_f^{(i)}$  follows a normal distribution with

$$E(\hat{y}_f^{(i)} - X_{1f}\beta_1 - X_{2f}\beta_2) = -Z_f(I - W_i)\beta_2^*$$

and

$$\begin{aligned}\text{var}(\hat{y}_f^{(i)}) &= X_{1f}(X_1'\Omega^{-1}X_1)^{-1}X_{1f}' + C_fM_1^*C_f' + Z_fW_iZ_f' \\ &\quad + C_fM_1^*X_2^*W_iZ_f' + Z_fW_iX_2^{*'}M_1^*C_f'.\end{aligned}$$

**Proof:** Using the facts that  $y = X_1\beta_1 + X_2^*\beta_2^* + u$  and  $b_2^* = \beta_2^* + X_2^{*'}M_1^*u$ , we write the predictor as

$$\begin{aligned}\hat{y}_f^{(i)} &= X_{1f}\beta_1 + X_{2f}^*\beta_2^* - Z_f(I - W_i)\beta_2^* \\ &\quad + X_{1f}(X_1'\Omega^{-1}X_1)^{-1}X_1'\Omega^{-1}u + (C_f' + X_2^*W_iZ_f')'M_1^*u,\end{aligned}$$

where we notice that the last two terms are independent. The results follow.  $\parallel$

The prediction error  $\text{PE}^{(i)} := \hat{y}_f^{(i)} - y_f$  can now be written as

$$\text{PE}^{(i)} = Z_{1f}(X_1'\Omega^{-1}X_1)^{-1}X_1'\Omega^{-1}u + Z_f(W_ib_2^* - \beta_2^*) - v_f, \quad (22)$$

where

$$Z_{1f} := X_{1f} - C_f\Omega^{-1}X_1, \quad v_f := u_f - C_f\Omega^{-1}u. \quad (23)$$

Since  $v_f$  and  $u$  are uncorrelated, and  $X_1'\Omega^{-1}u$  and  $b_2^*$  are also uncorrelated, we find that  $\text{PE}^{(i)}$  is the sum of three *independent* random variables.

**Theorem 2:** The prediction error  $\text{PE}^{(i)}$  follows a normal distribution with

$$E(\text{PE}^{(i)}) = -Z_f(I - W_i)\beta_2^*$$

and

$$\text{var}(\text{PE}^{(i)}) = Z_{1f}(X_1'\Omega^{-1}X_1)^{-1}Z_{1f}' + Z_fW_iZ_f' + \Omega_f - C_f\Omega^{-1}C_f',$$

and hence the mean squared prediction error  $\text{MSPE}^{(i)} := \text{MSE}(\text{PE}^{(i)})$  is

$$\text{MSPE}^{(i)} = Z_{1f}'(X_1'\Omega^{-1}X_1)^{-1}Z_{1f}' + Z_f'\Delta_i Z_f' + \Omega_f - C_f'\Omega^{-1}C_f',$$

where

$$\Delta_i := W_i + (I - W_i)\beta_2^*\beta_2^{*'}(I - W_i).$$

**Proof:** The results follow directly from (22).  $\parallel$

The best model is therefore the one where the matrix  $\Delta_i$  is as ‘small’ as possible. Since  $W_i$  is a diagonal matrix with only zeros and ones on the diagonal,  $\Delta_i$  is ‘small’ if the selected model  $\mathcal{M}_i$  includes precisely those regressors  $x_{2j}^*$  of  $X_2^*$  whose corresponding parameter  $\beta_{2j}^*$  is larger than one in absolute value. Since the  $\beta_{2j}^*$  are ‘theoretical’  $t$ -ratios, this result corresponds exactly to econometric intuition.

This econometric intuition is based on the following fact. Consider the partitioned regression model (10) with  $\text{var}(u) = \sigma^2 I_n$ . Let  $e_r$  denote the residual vector when  $y$  is regressed on  $X_1$  only, and let  $e_u$  denote the residual vector when  $y$  is regressed on  $X = (X_1 : X_2)$ . Then, under the null hypothesis that  $\beta_2 = 0$ , the test statistic

$$F = \frac{(e_r'e_r - e_u'e_u)/k_2}{e_u'e_u/(N - k)} \quad (24)$$

is distributed as  $F(k_2, N - k)$ . This is a standard result (Johnston and DiNardo, 1997, p. 97). Now define the least squares estimator for  $\sigma^2$  in the full model as  $s_u^2 = e_u'e_u/(N - k)$ , and the adjusted  $R^2$  as  $\bar{R}_u^2 = 1 - s_u^2/\sigma_y^2$ , where  $\sigma_y^2 := \sum_{n=1}^N (y_n - \bar{y})^2/(N - 1)$ . In the restricted model (where  $\beta_2 = 0$ ) define  $s_r^2$  and  $\bar{R}_r^2$  accordingly. Then,

$$\frac{1 - \bar{R}_r^2}{1 - \bar{R}_u^2} = \frac{s_r^2}{s_u^2} = 1 + \frac{k_2}{N - k_1}(F - 1), \quad (25)$$

and hence

$$s_u^2 < s_r^2 \iff \bar{R}_u^2 > \bar{R}_r^2 \iff F > 1. \quad (26)$$

As a special case ( $k_2 = 1$ ), we find that adding one regressor will decrease  $s^2$  and increase  $\bar{R}^2$  if and only if the  $t$ -statistic of the corresponding parameter is larger than one in absolute value.

## 4 The WALS predictor

The problem, of course, is that we don’t know which model to choose. Given estimates  $\hat{\beta}_{2j}^*$  of the  $k_2$  components  $\beta_{2j}^*$  of  $\beta_2^*$ , we could include the regressor  $x_{2j}^*$  if  $|\hat{\beta}_{2j}^*| > 1$ , and exclude it otherwise. This would lead to a *pretest*

estimator with well-established poor properties. These poor properties stem primarily from the fact that the pretest estimator is ‘kinked’; it has a discontinuity at one. This is not only mathematically undesirable but also intuitively: If  $\hat{\beta}_{2j}^* = 0.99$  we exclude  $x_{2j}^*$ ; if  $\hat{\beta}_{2j}^* = 1.01$  we include it. It would seem better to include  $x_{2j}^*$  ‘continuously’ in such a way that the higher is  $|\hat{\beta}_{2j}^*|$ , the more of  $x_{2j}^*$  is included in our model. This is precisely the idea behind model averaging. The additional benefit of model averaging is that we develop the theory taking into account both model uncertainty and parameter uncertainty. In other words, we think of model selection and parameter estimation as *one* combined procedure, so that the reported standard errors reflect both types of uncertainty.

Thus motivated, we define the WALS predictor of  $y_f$  as

$$\hat{y}_f = \sum_{i=1}^{2^{k_2}} \lambda_i \hat{y}_f^{(i)}, \quad (27)$$

where the sum is taken over all  $2^{k_2}$  different models obtained by setting a subset of the  $\beta_2^*$ ’s equal to zero, and the  $\lambda_i$ ’s are weight-functions satisfying certain minimal regularity conditions, namely

$$\lambda_i \geq 0, \quad \sum_{i=1}^{2^{k_2}} \lambda_i = 1, \quad \lambda_i = \lambda_i(M_1^* y). \quad (28)$$

This leads to the following definition.

**Definition 1 (WALS predictor):** The WALS predictor of  $y_f$  is given by

$$\hat{y}_f := X_{1f}(X_1' \Omega^{-1} X_1)^{-1} X_1' \Omega^{-1} y + C_f M_1^* y + Z_f \hat{\beta}_2^*,$$

where  $\hat{\beta}_2^* := W b_2^*$  and  $W := \sum_i \lambda_i W_i$ .

Note that, while the  $W_i$ ’s are non-random diagonal matrices, the matrix  $W$  is random (but still diagonal) because it depends on the random  $\lambda_i$ ’s. The prediction error  $\text{PE} := \hat{y}_f - y_f$  now takes the form

$$\text{PE} = Z_{1f}(X_1' \Omega^{-1} X_1)^{-1} X_1' \Omega^{-1} u + Z_f(\hat{\beta}_2^* - \beta_2^*) - v_f, \quad (29)$$

and we present its moments in the following ‘equivalence’ theorem.

**Theorem 3 (Equivalence theorem):** If the weights  $\lambda_i$  satisfy condition (28), then the WALS prediction error PE has the following expectation, variance and mean squared error:

$$\text{E}(\text{PE}) = Z_f \text{E}(\hat{\beta}_2^* - \beta_2^*),$$

$$\text{var}(\text{PE}) = Z_{1f}(X_1'\Omega^{-1}X_1)^{-1}Z_{1f}' + Z_f \text{var}(\hat{\beta}_2^*)Z_f' + \Omega_f - C_f\Omega^{-1}C_f',$$

and hence

$$\text{MSE}(\text{PE}) = Z_{1f}(X_1'\Omega^{-1}X_1)^{-1}Z_{1f}' + Z_f \text{MSE}(\hat{\beta}_2^*)Z_f' + \Omega_f - C_f\Omega^{-1}C_f'.$$

**Proof:** The key ingredient is that  $\text{cov}(M_1^*u, X_1'\Omega^{-1}u)$  and  $\text{cov}(u, v_f)$  are both zero. In addition, the  $\lambda_i$  (and hence  $W$ ) depend only on  $M_1^*y$  so that  $\hat{\beta}_2^* = Wb_2^*$  also depends only on  $M_1^*y$ . Hence, the three random variables  $X_1'\Omega^{-1}u$ ,  $\hat{\beta}_2^*$ , and  $v_f$  are all independent of each other. The results follow.  $\parallel$

The equivalence theorem tells us that the WALS predictor  $\hat{y}_f$  will be a ‘good’ predictor of  $y_f$  in the mean squared error sense if and only if  $\hat{\beta}_2^*$  is a ‘good’ estimator of  $\beta_2^*$ . That is, if we can find  $\lambda_i$ ’s such that  $\hat{\beta}_2^*$  is an ‘optimal’ estimator of  $\beta_2^*$ , then *the same*  $\lambda_i$ ’s will provide an ‘optimal’ predictor of  $y_f$ .

Next we obtain expressions for the bias and variance of the predictor itself, under the assumption that the diagonal elements of  $W$  depend only on  $b_2^* = X_2^*M_1^*y$  rather than only on  $M_1^*y$ .

**Theorem 4:** If the diagonal elements  $w_j$  of  $W$  depend only on  $b_2^*$ , then the WALS predictor  $\hat{y}_f$  has the following bias and variance:

$$\text{E}(\hat{y}_f - X_{1f}\beta_1 - X_{2f}\beta_2) = Z_f \text{E}(\hat{\beta}_2^* - \beta_2^*)$$

and

$$\begin{aligned} \text{var}(\hat{y}_f) &= X_{1f}(X_1'\Omega^{-1}X_1)^{-1}X_{1f}' + C_fM_1^*C_f' + Z_f \text{var}(\hat{\beta}_2^*)Z_f' \\ &\quad + C_fM_1^*X_2^* \text{cov}(b_2^*, \hat{\beta}_2^*)Z_f' + Z_f \text{cov}(\hat{\beta}_2^*, b_2^*)X_2^{*'}M_1^*C_f'. \end{aligned}$$

Under the stronger assumption that  $w_j$  depends only on  $b_{2j}^*$ , the  $k_2 \times k_2$  matrices  $\text{var}(\hat{\beta}_2^*)$  and  $\text{cov}(b_2^*, \hat{\beta}_2^*)$  are both diagonal.

**Proof:** The bias follows directly from Theorem 3. Noting that

$$\text{cov}(X_1'\Omega^{-1}y, M_1^*y) = X_1'M_1^* = 0, \quad \text{cov}(X_1'\Omega^{-1}y, \hat{\beta}_2^*) = 0,$$

Definition 1 implies that

$$\begin{aligned} \text{var}(\hat{y}_f) &= X_{1f}(X_1'\Omega^{-1}X_1)^{-1}X_{1f}' + C_fM_1^*C_f' + Z_f \text{var}(\hat{\beta}_2^*)Z_f' \\ &\quad + C_f \text{cov}(M_1^*y, \hat{\beta}_2^*)Z_f' + Z_f \text{cov}(\hat{\beta}_2^*, M_1^*y)C_f'. \end{aligned}$$

Since  $\Omega^{1/2}M_1^*\Omega^{1/2}$  is idempotent, we can write

$$\Omega^{1/2}M_1^*\Omega^{1/2} = AA', \quad A'A = I_{n-k_1}.$$

Define  $y^* := A'\Omega^{-1/2}y$  and  $B_1 := A'\Omega^{-1/2}X_2^*$ , so that  $y^* \sim N(B_1\beta_2^*, I_{n-k_1})$ . Since  $B_1'B_1 = I_{k_2}$  there exists an  $(n-k_1) \times (n-k)$  matrix  $B_2$ , such that  $B := (B_1 : B_2)$  is orthogonal. This allows us to write

$$M_1^*y = \Omega^{-1/2}A(B_1B_1' + B_2B_2')y^*, \quad \hat{\beta}_2^* = WB_1'y^*,$$

so that

$$\begin{aligned} \text{cov}(M_1^*y, \hat{\beta}_2^*) &= \text{cov}(\Omega^{-1/2}AB_1B_1'y^*, WB_1'y^*) + \text{cov}(\Omega^{-1/2}AB_2B_2'y^*, WB_1'y^*) \\ &= M_1^*X_2^* \text{cov}(b_2^*, \hat{\beta}_2^*) + \Omega^{-1/2}AB_2 \text{cov}(B_2'y^*, WB_1'y^*) \\ &= M_1^*X_2^* \text{cov}(b_2^*, \hat{\beta}_2^*), \end{aligned}$$

because  $B_1'y^*$  and  $B_2'y^*$  are independent, and the diagonal elements  $w_j$  of  $W$  depend only on  $X_2^{*'}M_1^*y = B_1'y^*$ .

Finally, if  $w_j$  depends only on  $b_{2j}^*$ , then

$$\text{cov}(b_{2i}^*, w_j b_{2j}^*) = 0, \quad \text{cov}(w_i b_{2i}^*, w_j b_{2j}^*) = 0 \quad (i \neq j),$$

because  $b_{2i}^*$  and  $b_{2j}^*$  are independent. In that case both  $\text{cov}(b_2^*, \hat{\beta}_2^*)$  and  $\text{cov}(\hat{\beta}_2^*, \hat{\beta}_2^*)$  are diagonal. This completes the proof.  $\parallel$

## 5 Computation of the WALS predictor based on prior ignorance

The WALS predictor proposed in Definition 1 can not be computed unless we know  $W = \sum_i \lambda_i W_i$ . Because of the semi-orthogonal transformation, we do know that  $W$  is diagonal, say  $W = \text{diag}(w_1, \dots, w_{k_2})$ . There are  $2^{k_2}$   $\lambda_i$ 's, but there are only  $k_2$   $w_j$ 's. These are functions of the  $\lambda_i$ 's, but we can not identify the  $\lambda_i$ 's from the  $w_j$ 's. This does not matter because we are not interested in the  $\lambda_i$ 's as we are not interested in selecting the 'best' model. We are only interested in the 'best' predictor.

The  $k_2$  components  $b_{2j}^*$  of  $b_2^*$  are independent with  $\text{var}(b_{2j}^*) = 1$ . Therefore, if we choose  $w_j$  to be a function of  $b_{2j}^*$  only, then the components  $\hat{\beta}_{2j}^* = w_j b_{2j}^*$  of  $\hat{\beta}_2^*$  will also be independent, and our  $k_2$ -dimensional problem reduces to  $k_2$  one-dimensional problems. The one-dimensional problem is simply how to estimate  $\beta_{2j}^*$  using only the information that  $b_{2j}^* \sim N(\beta_{2j}^*, 1)$ .

This seemingly trivial question was addressed in Magnus (2002), who proposed the 'Laplace' estimator. This estimator is obtained by combining the normal likelihood with the Laplace prior,

$$b_{2j}^* | \beta_{2j}^* \sim N(\beta_{2j}^*, 1), \quad \pi(\beta_{2j}^*) = (c/2) \exp(-c|\beta_{2j}^*|), \quad (30)$$

where  $c$  is a positive constant. The Laplace estimator is now defined as the resulting posterior expectation  $\hat{\beta}_{2j}^* := \mathbb{E}(\beta_{2j}^* | b_{2j}^*)$ . It is admissible, has bounded risk, has good properties around  $|\beta_{2j}^*| = 1$ , and is near-optimal in terms of minimax regret. It is also easily computable. The mean and variance of  $\beta_{2j}^* | b_{2j}^*$  are given in Theorem 1 of Magnus et al. (2010). The mean is

$$\hat{\beta}_{2j}^* = \mathbb{E}(\beta_{2j}^* | b_{2j}^*) = b_{2j}^* - c \cdot h(b_{2j}^*) \quad (31)$$

with

$$h(x) := \frac{e^{-cx}\Phi(x-c) - e^{cx}\Phi(-x-c)}{e^{-cx}\Phi(x-c) + e^{cx}\Phi(-x-c)}, \quad (32)$$

and the variance  $v_j := \text{var}(\beta_{2j}^* | b_{2j}^*)$  is

$$v_j = v(b_{2j}^*) = 1 + c^2(1 - h^2(b_{2j}^*)) - \frac{c(1 + h(b_{2j}^*))\phi(b_{2j}^* - c)}{\Phi(b_{2j}^* - c)}, \quad (33)$$

where  $\phi$  and  $\Phi$  denote the density function and the cumulative distribution function of the standard-normal distribution, respectively.

The weights  $w_j$  are defined implicitly by  $\hat{\beta}_{2j}^* = w_j b_{2j}^*$  and are thus given by

$$w_j = w(b_{2j}^*) = 1 - \frac{c \cdot h(b_{2j}^*)}{b_{2j}^*}. \quad (34)$$

Each  $w_j$  satisfies  $w(-b_{2j}^*) = w(b_{2j}^*)$  and increases monotonically between  $w(0)$  and  $w(\infty) = 1$ . Hence,  $\hat{\beta}_{2j}^*$  is a shrinkage estimator, and we have

$$w(0)|b_{2j}^*| < |\hat{\beta}_{2j}^*| < |b_{2j}^*|. \quad (35)$$

In particular, when  $c = \log 2$ , we find that  $w(0) = 0.5896$  which defines the maximum allowable shrinkage.

The hyperparameter  $c$  is chosen as  $c = \log 2$ , because this implies

$$\Pr(\beta_{2j}^* > 0) = \Pr(\beta_{2j}^* < 0), \quad \Pr(|\beta_{2j}^*| > 1) = \Pr(|\beta_{2j}^*| < 1). \quad (36)$$

What this means is that we assume a priori ignorance about whether  $\beta_{2j}^*$  is positive or negative, and also about whether  $|\beta_{2j}^*|$  is larger or smaller than one. These seem natural properties for a prior in our context, because we don't know a priori whether the  $\beta_2^*$  coefficients are positive or negative, and we don't know either whether adding a specific column of  $X_2^*$  to the model will increase or decrease the mean squared error of the predictors. Such a prior thus captures prior ignorance in a natural way. Given the choice of the weights  $w_j$  and hence of the estimator  $\hat{\beta}_2^*$ , the WALs predictor  $\hat{y}_f$  can be computed.



## 6 Moments of the WALS predictor

The moments of the WALS predictor are given in Theorem 4, but the expressions provided there depend on unknown quantities. Under the assumption that the weights  $w_j$  are specified as in (34), and hence depend on  $b_{2j}^*$  only, we estimate these unknown quantities as follows.

**Theorem 5:** If the diagonal elements  $w_j$  of  $W$  depend only on  $b_{2j}^*$  as specified in (34), then the expected bias of the WALS predictor  $\hat{y}_f$ , based on prior densities  $\pi(\beta_{2j}^*)$ , is zero:

$$\mathbb{E}(\mathbb{E}(\hat{y}_f - X_{1f}\beta_1 - X_{2f}\beta_2 | \beta_2^*)) = 0.$$

**Proof:** According to Theorem 4, the prediction bias, conditional on  $\beta_2^*$ , is

$$\mathbb{E}(\hat{y}_f - X_{1f}\beta_1 - X_{2f}\beta_2^* | \beta_2^*) = Z_f \mathbb{E}(\hat{\beta}_2^* - \beta_2^* | \beta_2^*).$$

Further,

$$\begin{aligned} \mathbb{E}(\hat{\beta}_{2j}^* - \beta_{2j}^*) &= \mathbb{E}\left(\mathbb{E}(\hat{\beta}_{2j}^* - \beta_{2j}^* | \beta_{2j}^*)\right) \\ &= \mathbb{E}\left(\mathbb{E}(b_{2j}^* - \beta_{2j}^* | \beta_{2j}^*)\right) - c \cdot \mathbb{E}\left(\mathbb{E}(h(b_{2j}^*) | \beta_{2j}^*)\right) = 0, \end{aligned}$$

because  $\mathbb{E}(h(b_{2j}^*) | \beta_{2j}^*)$  is antisymmetric in  $\beta_{2j}^*$  and  $\pi(\beta_{2j}^*)$  is symmetric in  $\beta_{2j}^*$ . Hence the expected bias of  $\hat{y}_f$  vanishes.  $\parallel$

The variance of  $\hat{y}_f$  is given in Theorem 4. Under the assumption that the weights  $w_j$  depend only on  $b_{2j}^*$ , the matrices  $\text{var}(\hat{\beta}_2^*)$  and  $\text{cov}(b_2^*, \hat{\beta}_2^*)$  are both diagonal. Hence it suffices to discuss the estimation of  $\text{var}(\hat{\beta}_{2j}^*)$  and  $\text{cov}(b_{2j}^*, \hat{\beta}_{2j}^*)$ . The variance in the posterior distribution of  $\beta_{2j}^* | b_{2j}^*$  is given by  $v_j$  in (33), and hence provides the obvious estimate of  $\text{var}(\hat{\beta}_{2j}^*)$ . It is less obvious how to find an appropriate estimate of  $\text{cov}(b_{2j}^*, \hat{\beta}_{2j}^*)$ . We propose

$$w_j = \widehat{\text{cov}}(b_{2j}^*, \hat{\beta}_{2j}^*) = \widehat{\text{cov}}(b_{2j}^*, w(b_{2j}^*)b_{2j}^*). \quad (37)$$

Since  $\text{var}(b_{2j}^*) = 1$ , this would be a perfect estimate if  $w_j$  were a constant. Now,  $w_j$  depends on  $b_{2j}^*$  and is therefore not a constant. Still, its variation is very small compared to the variation in  $b_{2j}^*$ . The correlation associated with the covariance is

$$\widehat{\text{corr}}(b_{2j}^*, \hat{\beta}_{2j}^*) = \frac{\widehat{\text{cov}}(b_{2j}^*, \hat{\beta}_{2j}^*)}{\sqrt{\widehat{\text{var}}(b_{2j}^*)\widehat{\text{var}}(\hat{\beta}_{2j}^*)}} = \frac{w(b_{2j}^*)}{\sqrt{v(b_{2j}^*)}}, \quad (38)$$

since we estimate  $\text{var}(\hat{\beta}_{2j}^*)$  by  $v_j = v(b_{2j}^*)$ . The estimated correlation is therefore always positive (in fact, larger than 0.7452) and smaller than one, such that when  $b_{2j}^*$  approaches  $\pm\infty$  the correlation approaches one.

We conclude that a suitable estimator for the variance of the WALS predictor is given by

$$\begin{aligned} \widehat{\text{var}}(\hat{y}_f) &= X_{1f}(X_1'\Omega^{-1}X_1)^{-1}X_{1f}' + C_fM_1^*C_f' + Z_fVZ_f' \\ &\quad + C_fM_1^*X_2^*WZ_f' + Z_fWX_2^*M_1^*C_f', \end{aligned} \quad (39)$$

where  $V$  and  $W$  are diagonal  $k_2 \times k_2$  matrices whose  $j$ -th diagonal elements  $v_j$  and  $w_j$  are given in (33) and (34), respectively. Having thus obtained estimators for all unknown quantities, the prediction variance can be computed.

## 7 Unknown variance matrix

We have thus far assumed  $\Omega$  is known, whereas in practice  $\Omega$  is of course unknown. One may estimate  $\theta$  based on the unrestricted model by minimizing

$$\varphi(\theta) := \log |\Omega| + y'(\Omega^{-1} - \Omega^{-1}X(X'\Omega^{-1}X)^{-1}X'\Omega^{-1})y \quad (40)$$

with respect to  $\theta$ . This leads to the maximum likelihood estimator  $\hat{\theta}$  of  $\theta$ , through which we also obtain an estimator  $\hat{\Omega} = \Omega(\hat{\theta})$  of  $\Omega$ . Note that the gradient of  $\varphi$  is the  $m \times 1$  vector whose  $i$ -th component is given by

$$\frac{\partial \varphi(\theta)}{\partial \theta_i} = \text{tr} \left( \Omega^{-1} \frac{\partial \Omega}{\partial \theta_i} \right) - (M^*y)' \frac{\partial \Omega}{\partial \theta_i} (M^*y), \quad (41)$$

where

$$M^* = M_1^*(\Omega - X_2^*X_2^{*'})M_1^*. \quad (42)$$

Therefore,  $\hat{\theta}$  depends on  $y$  only through  $M_1^*y$  and the same holds for  $\hat{\Omega}$ .

## 8 Simulation setup

Sections 2–7 contain the theoretical framework. Our next task is to evaluate the performance of the WALS predictor in a number of common situations and in comparison with other often-used predictors. In the current section we describe the setup of our simulation experiment. The simulation results are presented in Section 9. Many extensions of the benchmark setup were considered and some of these are summarized in Sections 10 and 11.

## 8.1 Five methods

In the simulations we compare the performance of the WALS predictor to four commonly-used methods: unrestricted maximum likelihood (ML), pretesting (PT), ridge regression (Ridge), and Mallows model averaging (MMA). We briefly describe each method below.

Unrestricted maximum likelihood simply estimates the unrestricted model (with *all* auxiliary regressors). There is no model selection here, and hence no noise associated with the model selection procedure. On the other hand, the noise associated with the estimation procedure will be large because of the large number of parameters.

Pretest estimation is a long-standing practice in applied econometrics, perhaps because pretest estimators are ‘logical outcomes of the increased diagnostic testing of assumptions advocated in many econometric circles’ (Poirier, 1995, p. 522). Pretest estimators and predictors do not follow textbook OLS or GLS properties, because the reported predictor is biased and its variance is only correct *conditional* on the selected model. One would expect the unconditional (‘true’) variance to be larger, because of the model selection noise. Giles and Giles (1993) provide a comprehensive review of the pretest literature. In pretest prediction one first selects the model based on diagnostic testing, and then predicts under the selected model. The choice of critical values of the pretest has received much attention (Toyoda and Wallace, 1976; Ohtani and Toyoda, 1980; Wan and Zou, 2003). Here we use the *stepwise fit* routine in Matlab, one of the most popular pretest methods. This routine begins with a forward selection procedure based on an initial model, then employs backward selection to remove variables. The steps are repeated until no additions or deletions of variables are indicated. We treat the model that includes only the focus regressors as the initial model and let the routine select the auxiliary regressors according to statistical significance. We choose the significance level for adding a variable to be 0.05 and for removing a variable to be 0.10.

Ridge regression (Hoerl and Kennard, 1970) is a common shrinkage technique, originally designed to address multicollinearity. Since the focus parameters are always in the model, we only penalize the auxiliary parameters. The ridge estimator is then obtained by minimizing

$$\phi(\beta_1, \beta_2) = (y - X_1\beta_1 - X_2\beta_2)'(y - X_1\beta_1 - X_2\beta_2) + \kappa\beta_2'\beta_2. \quad (43)$$

Letting

$$E_1 = \begin{pmatrix} I_{k_1} & 0_{k_1 \times k_2} \\ 0_{k_2 \times k_1} & 0_{k_2 \times k_2} \end{pmatrix}, \quad E_2 = \begin{pmatrix} 0_{k_1 \times k_1} & 0_{k_1 \times k_2} \\ 0_{k_2 \times k_1} & I_{k_2} \end{pmatrix}, \quad (44)$$

the solution can be written as

$$\hat{\beta}(\kappa) = (X'X + \kappa E_2)^{-1} X'y, \quad (45)$$

where  $\kappa$  is the tuning parameter. Alternatively we obtain the ridge estimator in a Bayesian framework as the mean in the posterior distribution of  $\beta|(X'X)^{-1}X'y$  by combining the data density  $(X'X)^{-1}X'y|\beta \sim N(\beta, \sigma^2(X'X)^{-1})$  with the partially informative prior  $\beta/\sigma \sim N(0, (1/\epsilon)E_1 + (1/\kappa)E_2)$  and letting  $\epsilon \rightarrow 0$ . Following Golub et al. (1979), we choose the tuning parameter  $\kappa$  by minimizing the generalized cross validation criterion

$$\text{GCV}(\kappa) = \frac{(y - \Xi(\kappa)y)'(y - \Xi(\kappa)y)}{(N - \text{tr} \Xi(\kappa))^2}, \quad \Xi(\kappa) = X(X'X + \kappa E_2)^{-1} X'. \quad (46)$$

Finally, Mallows model averaging, proposed by Hansen (2007), averages over estimators using weights obtained by minimizing the Mallows criterion

$$C(\lambda) = (y - P(\lambda)y)'(y - P(\lambda)y) + 2\sigma^2 \text{tr} P(\lambda), \quad (47)$$

where  $\lambda = (\lambda_1, \dots, \lambda_{2k_2})$ ,  $P(\lambda) = \sum_i \lambda_i X^{(i)}(X^{(i)'}X^{(i)})^{-1}X^{(i)'}$ , and  $X^{(i)}$  is the regressor matrix in model  $\mathcal{M}_i$ . Note that we do not assume an explicit ordering of the regressors, as Hansen does. An explicit ordering has the computational advantage that it reduces the number of weights from  $2^{k_2}$  to  $k_2$ , but it is typically not practical in applications. When the submodels are strictly nested, Hansen (2007) proved that the MMA estimator is asymptotically optimal in a given class of model averaging estimators. Wan et al. (2010) extended the optimality to non-nested models.

All predictors explicitly account for possible correlation in the random disturbances. In particular, the WALS predictor is obtained using Definition 1, and the predictors of the other four predictors are all computed from

$$\hat{y}_f = X_f \hat{\beta} + C_f \Omega^{-1}(y - X \hat{\beta}), \quad (48)$$

where  $\hat{\beta}$  depends on the chosen method. For ML (unrestricted model, no model selection), the predictor is linear in  $y$  and the associated variance is easily computed. For PT and Ridge, the predictor is not linear in  $y$ , but the reported variance is calculated as if the predictor were linear in  $y$ , following common practice. The variance for WALS is estimated from (39) while the variance for MMA can not be computed.

## 8.2 Data-generation process

We generate the data in three steps. First, we design the regressor matrix  $X = (X_1 : X_2 : X_3)$ , where  $X_1$  and  $X_2$  contain the focus and auxiliary

variables, while  $X_3$  contains the regressors that are omitted by the researcher (from *every* model) either because of ignorance or because of data limitations. The DGP and the largest (unrestricted) model are therefore not necessarily the same in the simulations. This is important because it brings us one step closer to econometric practice. In the benchmark DGP we consider six regressors with  $k_1 = 2$ ,  $k_2 = 3$ , and  $k_3 = 1$ , such that

$$X_1 = (x_1, x_2), \quad X_2 = (x_3, x_4, x_5) \quad X_3 = (x_6), \quad (49)$$

where  $x_1$  is the intercept. Since  $k_2 = 3$  we have  $2^3 = 8$  possible models. In the benchmark, all regressors, except the intercept, are generated by independent standard-normal distributions, and they are treated as fixed, so that each replication uses the same realization of the regressors once they have been generated. In Section 11 we shall consider extensions where we have a large number of regressors and the regressors are correlated or non-normally distributed.

Table 1: Theoretical  $t$ -ratios for benchmark model

$T$	Auxiliary			Omitted
	$t_3$	$t_4$	$t_5$	$t_6$
$T_1$	1.2	0.9	1.1	0.0
$T_2$	1.2	1.7	0.7	0.9
$T_3$	1.2	0.9	1.0	2.5
$T_4$	2.0	2.5	2.7	0.0
$T_5$	0.4	0.2	0.5	0.0

Next, we simulate the parameters  $\beta_j$  ( $j = 1, \dots, 6$ ) corresponding to regressors  $x_1, \dots, x_6$ . For the auxiliary and omitted regressors  $x_3, \dots, x_6$  we set these parameters indirectly by controlling the ‘theoretical’  $t$ -ratios, as follows. If we estimate the focus variables and just one auxiliary variable  $x_j$ , we obtain an estimated coefficient  $\hat{\beta}_j$  with variance  $\text{var}(\hat{\beta}_j) = (x_j' M_1^* x_j)^{-1}$ . This implies a  $t$ -ratio  $\hat{t}_j = \hat{\beta}_j \sqrt{x_j' M_1^* x_j}$ . The ‘theoretical’  $t$ -ratio is now defined as

$$t_j = \beta_j \sqrt{x_j' M_1^* x_j} \quad (j = 3, \dots, 6). \quad (50)$$

The values of the  $t_j$  are important (especially whether  $|t_j| > 1$  or  $|t_j| < 1$ ), because they determine whether adding an auxiliary regressor to the model will increase or decrease the root mean squared prediction error (the square root of the mean squared prediction error); see Theorem 2. We consider five combinations, see Table 1. Given  $x_j$  and  $t_j$ , we then obtain the parameters  $\beta_j$  ( $j = 3, \dots, 6$ ). Three of the five cases ( $T_1, T_4, T_5$ ) have no omitted variables.

In  $T_1$  the  $t$ -ratios of the auxiliary variables are close to 1, in  $T_4$  the  $t$ -ratios are large, and in  $T_5$  they are small. The other two cases ( $T_2, T_3$ ) have an omitted variable. The value of  $t_6$  is either close to one ( $T_2$ ) or large ( $T_3$ ).

Regarding the focus parameters we let  $\beta_1 = \beta_2 = \nu \sqrt{\sum_{j=3}^6 \beta_j^2}$  for three values of  $\nu$ : 1, 2, and 3. Since the prediction performance is hardly affected by this choice, we shall report for  $\nu = 2$  only.

Finally, we generate the error terms, based on (3), from a normal distribution with mean zero and variance  $\Omega_{all}$ . We consider three specifications of  $\Omega_{all}$ : homoskedasticity, heteroskedasticity, and autocorrelation. More precisely,

- homoskedasticity:  $\Omega_{all} = \sigma^2 I_{n+n_f}$  with  $\sigma^2 \in \{0.25, 1.00\}$ ;
- heteroskedasticity:  $\Omega_{all} = \text{diag}[\exp(\tau x_2)]$  with  $\tau \in \{0.2, 0.7\}$ ;
- autocorrelation: AR(1) with  $\sigma^2 = 1.0$  and  $\rho \in \{0.3, 0.8\}$ .

### 8.3 Comparison of prediction methods

We evaluate the five methods by comparing the predictors and the estimated variances of the predictors. To compare the predictors produced by the five methods, we consider the deviation between the predictor  $\hat{y}_f$  and the true value  $y_f$ . A direct comparison is, however, misleading because there is a component common to all procedures. Hence we compute a modified version of the root mean squared prediction error,

$$\sqrt{\frac{1}{R} \sum_{r=1}^R \left( \hat{y}_f^{(r)} - y_f^{(r)} + (u_f - C_f \Omega^{-1} u) \right)' \left( \hat{y}_f^{(r)} - y_f^{(r)} + (u_f - C_f \Omega^{-1} u) \right)} \quad (51)$$

where  $\hat{y}_f^{(r)}$  and  $y_f^{(r)}$  are the predictor and the true value in the  $r$ -th replication. We follow Hansen (2008) and subtract  $u_f - C_f \Omega^{-1} u$  from the prediction error, because it is common across prediction methods and independent of  $u$ , hence independent of  $\hat{\beta} - \beta$ .

To compare the prediction variances is more subtle. We could just compare the magnitudes of

$$\frac{1}{R} \sum_{r=1}^R \text{var}(\hat{y}_f^{(r)}), \quad (52)$$

which would tell us whether one method reports more precise predictions than another. This is of interest, but more important than whether the reported prediction variance is small is whether the prediction variance is

*correct*. It is easy to find predictors with small variances, but this does not make them good predictors.

Thus we wish to determine how close the estimated variance is to the ‘true’ variance, and this is measured by the RMSE of the prediction variance,

$$\sqrt{\frac{1}{R} \sum_{r=1}^R \left( \text{var}(\hat{y}_f^{(r)}) - V_T \right)^2}, \quad (53)$$

where  $V_T$  denotes the ‘true’ variance, that is, the actual variance of the predictor. Since different methods give different predictors, the ‘true’ variance of the predictor varies across methods. We estimate  $V_T$  by obtaining  $R_v = 100$  predictors from the replications, and then computing the sample variance of these predictors,

$$V_T := \frac{1}{R_v - 1} \sum_{r=1}^{R_v} \left( \hat{y}_f^{(r)} - \frac{1}{R_v} \sum_{r=1}^{R_v} \hat{y}_f^{(r)} \right)^2. \quad (54)$$

We consider training samples of size  $N = 100$  and  $N = 300$ , and a prediction sample of size  $N_f = 10$ . The simulation results are obtained by computing averages across  $R = 3000$  draws.

## 9 Simulation results: The benchmark

We compare the predictors by considering two sample sizes ( $N = 100$ ,  $N = 300$ ), five sets of parameter values ( $T_1, \dots, T_5$ ), six specifications of  $\Omega_{all}$ , and five methods. Each method is presented relative to ML, that is, we present the RMSE of each method divided by the RMSE of ML. An entry greater than one thus indicates an inferior performance relative to the ML method.

The RMSEs of the predictors are presented in Table 2. WALS comes out best in 39 out of 60 cases (65%), followed by Ridge (27%) and ML (8%). The pretest and MMA predictors never dominate. The dominance of WALS occurs for each of the specifications of  $\Omega_{all}$ , though slightly less in the autocorrelation case than in the homo- and heteroskedastic cases.

In  $T_1$  WALS dominates in all 12 cases, and in  $T_2$  in 11/12 cases. This shows that WALS performs well when the  $t$ -ratios of the auxiliary variables are close to one, even when the model possibly omits one variable with a  $t$ -ratio close to one. If the omitted variable has a stronger impact on the dependent variable, as in  $T_3$ , WALS still works best in 9/12 cases followed by Ridge (3/12). This suggests that omitting important regressors may affect

Table 2: RMSE of predictor relative to ML, benchmark model

$N$	$T$	WALS	PT	Ridge	MMA	WALS	PT	Ridge	MMA
<i>Homoskedasticity</i>									
$\sigma^2 = 0.25$					$\sigma^2 = 1.00$				
100	$T_1$	<b>0.8644</b>	1.0570	0.9109	0.9416	<b>0.8981</b>	1.1589	1.0641	1.0172
	$T_2$	<b>0.8819</b>	1.0311	0.8979	0.9406	<b>0.9287</b>	1.1373	1.0350	1.0178
	$T_3$	<b>0.9525</b>	1.2025	1.0866	1.0769	<b>0.9366</b>	1.1001	0.9576	1.0031
	$T_4$	1.0024	1.2756	1.0296	1.1579	<b>0.9569</b>	1.2758	1.0521	1.1267
	$T_5$	0.8190	0.8600	<b>0.7556</b>	0.8019	0.7939	0.8280	<b>0.6906</b>	0.7680
300	$T_1$	<b>0.8990</b>	1.1152	0.9899	0.9940	<b>0.9295</b>	1.0421	0.9398	0.9727
	$T_2$	<b>0.9598</b>	1.0819	1.0237	1.0156	<b>0.9271</b>	1.0758	0.9648	0.9934
	$T_3$	<b>0.9523</b>	1.0341	0.9612	0.9841	<b>0.9642</b>	1.0480	0.9929	0.9977
	$T_4$	<b>0.9516</b>	1.1490	0.9977	1.0853	1.0513	1.2362	1.0481	1.1492
	$T_5$	0.8004	0.8571	<b>0.7195</b>	0.7878	0.8404	0.8777	<b>0.7823</b>	0.8238
<i>Heteroskedasticity</i>									
$\tau = 0.2$					$\tau = 0.7$				
100	$T_1$	<b>0.9038</b>	1.0766	0.9401	0.9708	<b>0.8926</b>	1.0964	0.9786	0.9815
	$T_2$	<b>0.8846</b>	1.1140	0.9936	0.9847	<b>0.8743</b>	1.0639	0.8699	0.9418
	$T_3$	<b>0.9736</b>	1.0682	1.0071	1.0169	<b>0.8752</b>	0.9326	0.8681	0.8849
	$T_4$	<b>0.9628</b>	1.2428	0.9650	1.1144	<b>0.9558</b>	1.2430	0.9679	1.1158
	$T_5$	0.8160	0.8633	<b>0.7494</b>	0.7976	0.8148	0.8633	<b>0.7486</b>	0.7948
300	$T_1$	<b>0.8928</b>	1.1104	0.9847	0.9917	<b>0.8834</b>	1.0718	0.8988	0.9581
	$T_2$	<b>0.9012</b>	1.0801	0.9810	0.9744	<b>0.8675</b>	1.0014	0.8694	0.9134
	$T_3$	0.8884	0.9702	<b>0.8807</b>	0.9069	<b>0.9421</b>	1.0325	0.9530	0.9754
	$T_4$	1.0596	1.2178	1.0485	1.1456	<b>0.9011</b>	1.2004	0.9741	1.1056
	$T_5$	0.8043	0.8431	<b>0.7209</b>	0.7806	0.9132	0.9360	<b>0.8847</b>	0.9058
<i>Autocorrelation</i>									
$\rho = 0.3$					$\rho = 0.8$				
100	$T_1$	<b>0.8875</b>	1.0543	0.9384	0.9586	<b>0.9809</b>	1.0182	0.9903	0.9958
	$T_2$	<b>0.8980</b>	1.0885	0.9001	0.9643	0.9659	1.0058	<b>0.9608</b>	0.9774
	$T_3$	0.8823	0.9723	<b>0.8398</b>	0.9110	<b>0.9990</b>	1.0232	1.0126	1.0111
	$T_4$	1.0089	1.1629	1.0127	1.0917	<b>0.9979</b>	1.0511	1.0109	1.0278
	$T_5$	0.8634	0.8971	<b>0.8077</b>	0.8510	0.9632	0.9715	<b>0.9525</b>	0.9600
300	$T_1$	<b>0.9188</b>	1.0868	0.9904	0.9934	<b>0.9760</b>	1.0241	0.9891	0.9947
	$T_2$	<b>0.9399</b>	1.0771	0.9975	0.9926	<b>0.9574</b>	0.9988	0.9603	0.9717
	$T_3$	0.9031	0.9819	<b>0.8991</b>	0.9243	<b>0.9856</b>	1.0164	0.9920	0.9975
	$T_4$	1.0460	1.1804	1.0423	1.1207	<b>0.9988</b>	1.0623	1.0060	1.0373
	$T_5$	0.8453	0.8735	<b>0.7829</b>	0.8265	0.9788	0.9828	<b>0.9711</b>	0.9764



the prediction ability of WALS, which is a point worthy of further investigation; see Section 11.

When the  $t$ -ratios of the auxiliary variables are much larger than one, as in  $T_4$ , then WALS is still the best, but this is the only case where ML also performs well. This makes sense, because model uncertainty plays a smaller role here. We note that increasing the parameter values in  $\Omega_{all}$  ( $\sigma^2$ ,  $\tau$ , and  $\rho$ , respectively) improves the relative performance of WALS over ML in  $T_4$ . Larger parameter values in  $\Omega_{all}$  imply more noise in the model, and the superiority of WALS seems stronger then. The possibility that the degree of model uncertainty affects the relative performance of different methods also warrants further investigation; see Section 10.

In the opposite case where the  $t$ -ratios of the auxiliary variables are much smaller than one, as in  $T_5$ , WALS is not the best. Here the Ridge predictor always dominates, and ML is always the worst. Again, there is little model uncertainty. The unrestricted model (ML) is not appropriate, but shrinkage towards the restricted estimator (with only the focus regressors) makes sense, and this is what Ridge does. Next we compare the performance of the prediction variance. We first consider the magnitude of the estimated variance itself, then we ask how close the estimated variance is to the ‘true’ variance. The MMA method is not included in this comparison because there is no procedure known to us to compute this variance. In the boxplots of Figure 1, the central mark is the median, the edges of each box indicate the 25-th and 75-th percentiles, the whiskers extend to the most extreme data points not considered outliers, and outliers are plotted individually.

We consider six representative cases. Judging by the median of the estimated variance, ML has the largest variance, followed by WALS, while the variance of the Ridge and PT predictors are both smaller than WALS. This is in accordance with intuition, because ML includes all regressors, while pretesting and ridge are based on the selected model or the selected parameter, while ignoring variation caused by the selection procedure. The WALS predictor has a relatively large variance (but still smaller than ML), because it does take the uncertainty in the selection procedure into account.

We note that the estimated variances for WALS and ML are more concentrated on their median values than those of Ridge and PT, and that the distributions of the latter two methods are also characterized by a strong asymmetry. The difference between the four variance estimates is relatively small when there is little model uncertainty ( $T_4$ ), and more pronounced when model uncertainty is large ( $T_1$ ). As discussed above, a variance estimate is a good estimate, not when it is small, but when it provides the correct information about the precision of the predictor. If this precision happens to be low, then we need to provide a high value for the variance estimate. Table 3

Figure 1: Estimated variance in the benchmark model ( $N = 100$ )

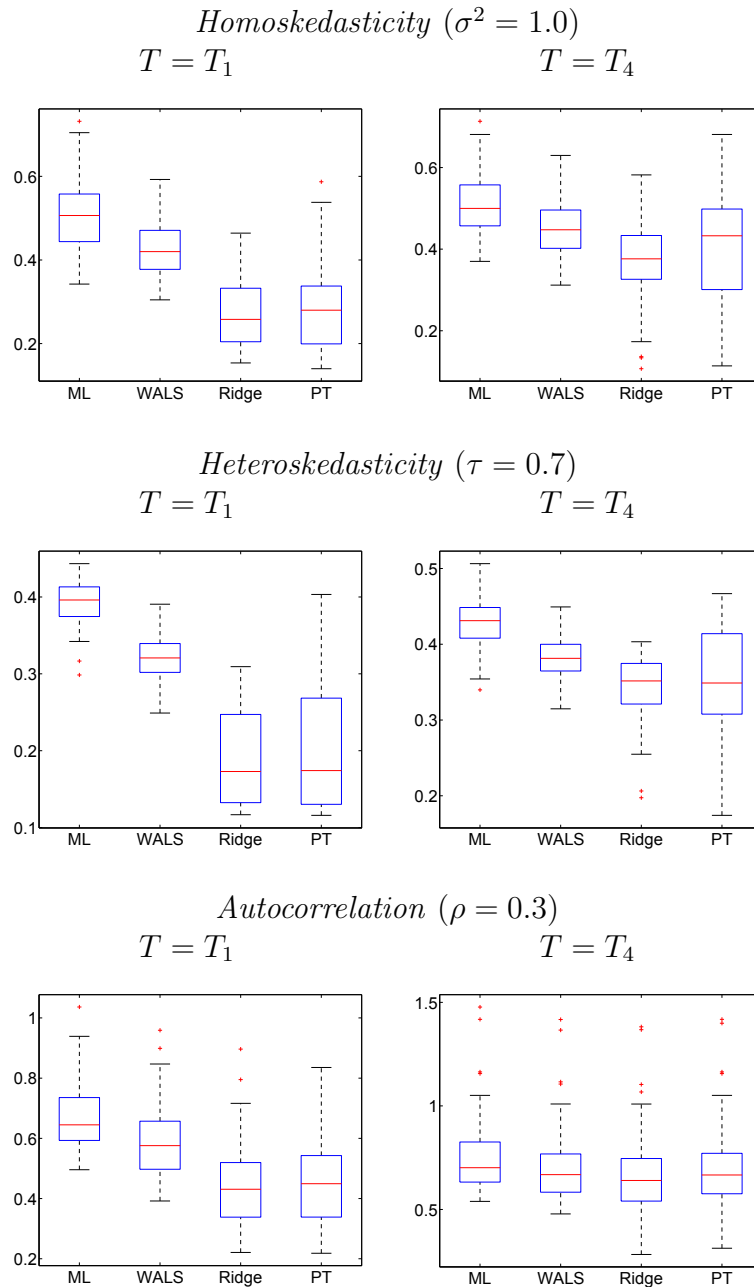


Table 3: RMSE of prediction variance relative to ML, benchmark model

$N$	$T$	WALS	PT	Ridge	WALS	PT	Ridge
<i>Homoskedasticity</i>							
		$\sigma^2 = 0.25$			$\sigma^2 = 1.00$		
100	$T_1$	<b>0.8155</b>	10.417	2.1262	<b>0.8050</b>	9.7509	2.1211
	$T_2$	<b>0.9011</b>	12.307	2.8462	<b>0.8606</b>	11.722	2.7199
	$T_3$	1.0760	8.5721	1.9561	1.0421	9.7699	2.0451
	$T_4$	<b>0.7755</b>	16.040	2.1837	<b>0.7972</b>	16.747	2.2023
	$T_5$	0.7765	3.0180	<b>0.6894</b>	0.7916	3.0839	<b>0.7061</b>
300	$T_1$	1.2139	16.351	3.3701	1.2947	15.494	3.0947
	$T_2$	1.1475	16.967	3.7704	1.3978	18.440	4.0967
	$T_3$	1.3510	15.950	3.3156	1.4509	15.176	3.3333
	$T_4$	1.0001	21.714	2.5714	<b>0.8938</b>	20.070	2.5663
	$T_5$	1.4811	5.3371	1.0851	1.2495	4.9051	1.0459
<i>Heteroskedasticity</i>							
		$\tau = 0.2$			$\tau = 0.7$		
100	$T_1$	1.1174	17.938	4.6738	<b>0.8849</b>	12.547	3.3035
	$T_2$	1.0030	23.113	5.3124	<b>0.9509</b>	18.382	4.2944
	$T_3$	1.1227	25.641	5.3652	<b>0.9648</b>	18.941	4.8476
	$T_4$	1.0083	17.748	3.3204	<b>0.9628</b>	15.244	2.8077
	$T_5$	1.2374	5.1309	1.6143	<b>0.9888</b>	4.5029	1.5134
300	$T_1$	1.3897	18.963	3.8448	1.1613	16.107	3.3285
	$T_2$	1.3279	20.426	4.4850	1.2645	17.335	4.3203
	$T_3$	1.4954	19.552	3.9925	1.3081	16.980	3.2964
	$T_4$	1.2030	27.566	3.4499	1.0624	22.479	2.7482
	$T_5$	1.1378	5.2663	1.1331	1.2631	5.2922	1.2149
<i>Autocorrelation</i>							
		$\rho = 0.3$			$\rho = 0.8$		
100	$T_1$	1.0605	2.5141	1.2681	1.0000	1.0156	1.0064
	$T_2$	1.0347	2.9666	1.3984	1.0011	1.0141	1.0072
	$T_3$	1.0381	2.8648	1.3693	<b>0.9984</b>	1.0213	1.0093
	$T_4$	1.0190	3.6270	1.2758	1.0007	1.0162	1.0047
	$T_5$	1.0633	1.4584	1.0620	1.0000	1.0088	1.0049
300	$T_1$	1.0060	1.6103	1.0716	<b>0.9983</b>	1.0050	1.0003
	$T_2$	1.0191	1.6745	1.1390	<b>0.9997</b>	1.0065	1.0016
	$T_3$	1.0080	1.6066	1.0980	<b>0.9991</b>	1.0087	1.0016
	$T_4$	1.0068	1.7774	1.0753	1.0006	1.0031	1.0009
	$T_5$	1.0258	1.1822	1.0238	1.0000	1.0035	1.0019

gives the RMSE of the estimated prediction variance, as given in (53), again relative to ML. On the left-side of the table (where the parameters  $\sigma^2$ ,  $\tau$ , and  $\rho$  are relatively small), the RMSE ratios (relative to ML) are, on average, 1.10 for WALS, 2.43 for Ridge, and 10.98 for PT. On the right-side (where the parameter values are larger, hence more uncertainty), the RMSE ratios are 1.05 for WALS, 2.19 for Ridge, and 9.44 for PT. The main conclusion from the table is therefore that ML and WALS provide the best estimates of the prediction variance, while Ridge and especially PT generally report a variance which is misleadingly small. While WALS provides a much better estimate of the forecast than ML, the variance of the forecast is slightly more accurately estimated in ML than in WALS.

ML performs particularly well when  $N$  is large (because of the asymptotic behavior of ML estimates and predictions) and when the variance parameters are small. The relative performance of WALS prediction variance estimates is improved by increasing the variance of the error terms. This suggests that more model uncertainty makes WALS prediction more attractive. In the benchmark setup, where we have assumed deterministic regressors and coefficients, there is not much model uncertainty. If we raise the model uncertainty, for example by introducing random regressors or random coefficients or by increasing the variance of the errors, then one would expect the WALS estimates, which incorporate the model uncertainty, to be more accurate than ML. We shall analyze this idea further in the next section.

## 10 Simulation results: More uncertainty

In this section we extend the benchmark setup by introducing additional randomness in the model. This is achieved by allowing for random regressors or random coefficients or by increasing the variance of errors.

### 10.1 Random regressors

We first consider the model with random but exogenous regressors. This is a common extension in simulation designs, and particularly useful in applications where one wishes to model dynamic economic behavior. The only difference with the benchmark is that we generate a new set of  $X$ 's from  $N(0, \sigma_x^2)$  in every replication, so that each realization of the  $y$ -series involves a new realization of the  $X$ -series. (The introduction of  $\sigma_x^2$  is unimportant, because the RMSE is invariant to its value.) The generation of  $X$  is independent of the errors.

Allowing the regressors to be random increases the RMSE of the forecast

in each method. The relative performance of the five predictors is similar to the benchmark case. In particular, the WALS predictor has the lowest RMSE in  $T_1$ ,  $T_2$ , and  $T_3$ , about 5% lower than the RMSE of the ML predictor. In case  $T_5$ , the ridge predictor has the lowest RMSE under all error structures, around 10% lower than the ML predictor. In contrast to the benchmark results, allowing random regressors improves the relative performance of WALS over ML in  $T_4$ , because more randomness decreases the importance of the auxiliary variables. The main difference between the random regressor model and the benchmark model is in the prediction variance, and we report its RMSE in Table 4. WALS now produces the most accurate prediction variance in all cases, including  $T_4$  and  $T_5$ . This remarkable performance of WALS is due to the fact that randomness in the regressors raises model uncertainty, which in turn increases the variation of the predictor, that is, the true variance. The prediction variance of WALS explicitly incorporates such model uncertainty, in contrast to pretesting, ridge regression, and ML.

## 10.2 Random coefficients

Next we consider the situation where the coefficients of the explanatory variables are subject to random variation, that is,

$$y_t = \sum_{j=1}^6 x_{tj}(\beta_j + v_{tj}) + u_t \quad (t = 1, 2, \dots, N), \quad (55)$$

where the  $v_{tj}$ 's are independent unobserved random disturbances, distributed as  $N(0, \sigma_v^2)$ . Such models date back to Rubin (1950), Hildreth and Houck (1968), Swamy (1970), Froehlich (1973), and others, who discussed parameter estimation and provided empirical applications. Prediction in random coefficient models is studied, *inter alia*, in Bondeson (1990) and Beran (1995). We can rewrite (55) as

$$y_t = \sum_{j=1}^6 x_{tj}\beta_j + \zeta_t, \quad \zeta_t = \sum_{j=1}^6 x_{tj}v_{tj} + u_t \quad (56)$$

where  $\zeta_t$  is normally distributed with mean zero and variance  $\sigma_\zeta^2 = \sigma_u^2 + \sigma_v^2 \sum_j x_{tj}^2$ . This shows that introducing variation in the coefficients increases the variance of the errors. We assume that the researcher is ignorant of the random coefficients and misspecifies them as fixed. Hence the model is the benchmark model, but the DGP has changed. How do the predictors respond to this situation?

Table 4: RMSE of prediction variance relative to ML, random regressor model

$N$	$T$	WALS	PT	Ridge	WALS	PT	Ridge
<i>Homoskedasticity</i>							
		$\sigma^2 = 0.25$			$\sigma^2 = 1.00$		
100	$T_1$	<b>0.7499</b>	1.0219	0.8056	<b>0.7467</b>	1.0119	0.8007
	$T_2$	<b>0.7958</b>	1.0096	0.8365	<b>0.7952</b>	1.0110	0.8362
	$T_3$	<b>0.8866</b>	1.0101	0.9149	<b>0.8899</b>	1.0161	0.9220
	$T_4$	<b>0.8487</b>	0.9951	0.8990	<b>0.8497</b>	0.9929	0.8987
	$T_5$	<b>0.5091</b>	0.9267	0.5336	<b>0.5064</b>	0.9021	0.5229
300	$T_1$	<b>0.7486</b>	1.0219	0.8064	<b>0.7435</b>	1.0165	0.8011
	$T_2$	<b>0.7978</b>	1.0107	0.8407	<b>0.7990</b>	1.0101	0.8399
	$T_3$	<b>0.8930</b>	1.0170	0.9222	<b>0.8889</b>	1.0115	0.9179
	$T_4$	<b>0.8473</b>	0.9953	0.8989	<b>0.8474</b>	0.9938	0.8988
	$T_5$	<b>0.5123</b>	0.9513	0.5505	<b>0.5147</b>	0.9417	0.5481
<i>Heteroskedasticity</i>							
		$\tau = 0.2$			$\tau = 0.7$		
100	$T_1$	<b>0.7448</b>	1.0296	0.8095	<b>0.7421</b>	1.0342	0.8128
	$T_2$	<b>0.7914</b>	1.0123	0.8381	<b>0.7950</b>	1.0129	0.8448
	$T_3$	<b>0.8862</b>	1.0080	0.9149	<b>0.8835</b>	1.0067	0.9137
	$T_4$	<b>0.8461</b>	0.9954	0.9010	<b>0.8495</b>	0.9925	0.9044
	$T_5$	<b>0.5125</b>	0.9515	0.5525	<b>0.5621</b>	0.9278	0.5855
300	$T_1$	<b>0.7444</b>	1.0192	0.8054	<b>0.7444</b>	1.0192	0.8054
	$T_2$	<b>0.7987</b>	1.0162	0.8450	<b>0.7987</b>	1.0162	0.8450
	$T_3$	<b>0.8906</b>	1.0139	0.9191	<b>0.8906</b>	1.0139	0.9191
	$T_4$	<b>0.8490</b>	0.9942	0.9014	<b>0.8490</b>	0.9942	0.9014
	$T_5$	<b>0.5146</b>	0.9614	0.5548	<b>0.5146</b>	0.9614	0.5548
<i>Autocorrelation</i>							
		$\rho = 0.3$			$\rho = 0.8$		
100	$T_1$	<b>0.7469</b>	1.0318	0.8150	<b>0.9809</b>	1.0059	0.9890
	$T_2$	<b>0.7997</b>	1.0241	0.8494	<b>0.9683</b>	1.0049	0.9785
	$T_3$	<b>0.8866</b>	1.0175	0.9177	<b>0.9646</b>	1.0065	0.9775
	$T_4$	<b>0.8479</b>	0.9986	0.9022	<b>0.9099</b>	1.0004	0.9433
	$T_5$	<b>0.6863</b>	0.9902	0.7205	<b>0.9975</b>	1.0082	1.0033
300	$T_1$	<b>0.7629</b>	1.0151	0.8199	<b>0.9924</b>	1.0003	0.9950
	$T_2$	<b>0.8054</b>	1.0122	0.8483	<b>0.9894</b>	1.0058	0.9940
	$T_3$	<b>0.8913</b>	1.0082	0.9201	<b>0.9848</b>	0.9996	0.9898
	$T_4$	<b>0.8478</b>	0.9934	0.9005	<b>0.9465</b>	1.0023	0.9666
	$T_5$	<b>0.8013</b>	0.9948	0.8239	0.9985	0.9990	<b>0.9977</b>

Regarding the accuracy of the predictors, we find similar results as in the random regressor model. The WALS predictor has the lowest RMSE in cases  $T_1$ – $T_4$ , while the ridge predictor is the best under  $T_5$ . This demonstrates good performance of the WALS predictor when the  $t$ -ratios of the auxiliary variables are close to one, even when the coefficients are misspecified. The accuracy of the estimated prediction variance is shown in Figure 2 as a function of  $\sigma_v^2$ . Increasing  $\sigma_v^2$  raises the model uncertainty as well as the degree of misspecification, thus lowering the accuracy of all predictions. The variance estimates obtained from pretesting have a much larger RMSE than those from other methods, and they are also more volatile. Ridge regression generally produces somewhat better variance estimates. Most accurate are ML and WALS, and their variance accuracy is close when  $\sigma_v^2$  is small. When  $\sigma_v^2 = 0$  (the benchmark), ML is more accurate than WALS, but as  $\sigma_v^2$  increases, the RMSE of WALS increases slower than the RMSE of ML, and when  $\sigma_v^2 > 0.03$  the accuracy of WALS variance estimates is higher than ML. These results confirm that WALS behaves well in the presence of a large degree of model uncertainty. Viewed differently, WALS is more robust than pretesting, ridge, and ML.

### 10.3 Increase in the variance of errors

Finally, we consider an increase in the variance of the errors by changing a parameter in  $\Omega_{all}$ . We only consider the homoskedastic and the heteroskedastic cases. Under homoskedasticity we can increase the error variance by increasing  $\sigma^2$ ; under heteroskedasticity case by increasing  $\tau$ . Figure 3 shows how the RMSE of the prediction variance changes as the parameters  $\sigma^2$  and  $\tau$  increase. In both cases, WALS and ML outperform Ridge and, in particular, PT. When the error variance is small, the prediction variances produced by WALS and ML show similar accuracy. But as the error variance increases, the WALS prediction variance is more accurate than ML.

Note that increasing the error variance affects the RMSE of the prediction variance in different ways: it increases the RMSE in the homoskedastic case but reduces the RMSE in the heteroskedastic case. This is because in the design of the heteroskedastic variance,  $\Omega_{all} = \exp(\tau x_2)$  is a function of  $x_2$ . Increasing  $\tau$  leads to a smaller estimated coefficient  $\hat{\beta}_2$  since the estimation process cannot distinguish between increasing the error variance from increasing the variation in  $x_2$ .

In summary, more model uncertainty leads to a better performance of WALS relative to the other methods.

Figure 2: RMSE of prediction variance in random coefficient model ( $N = 100$ )

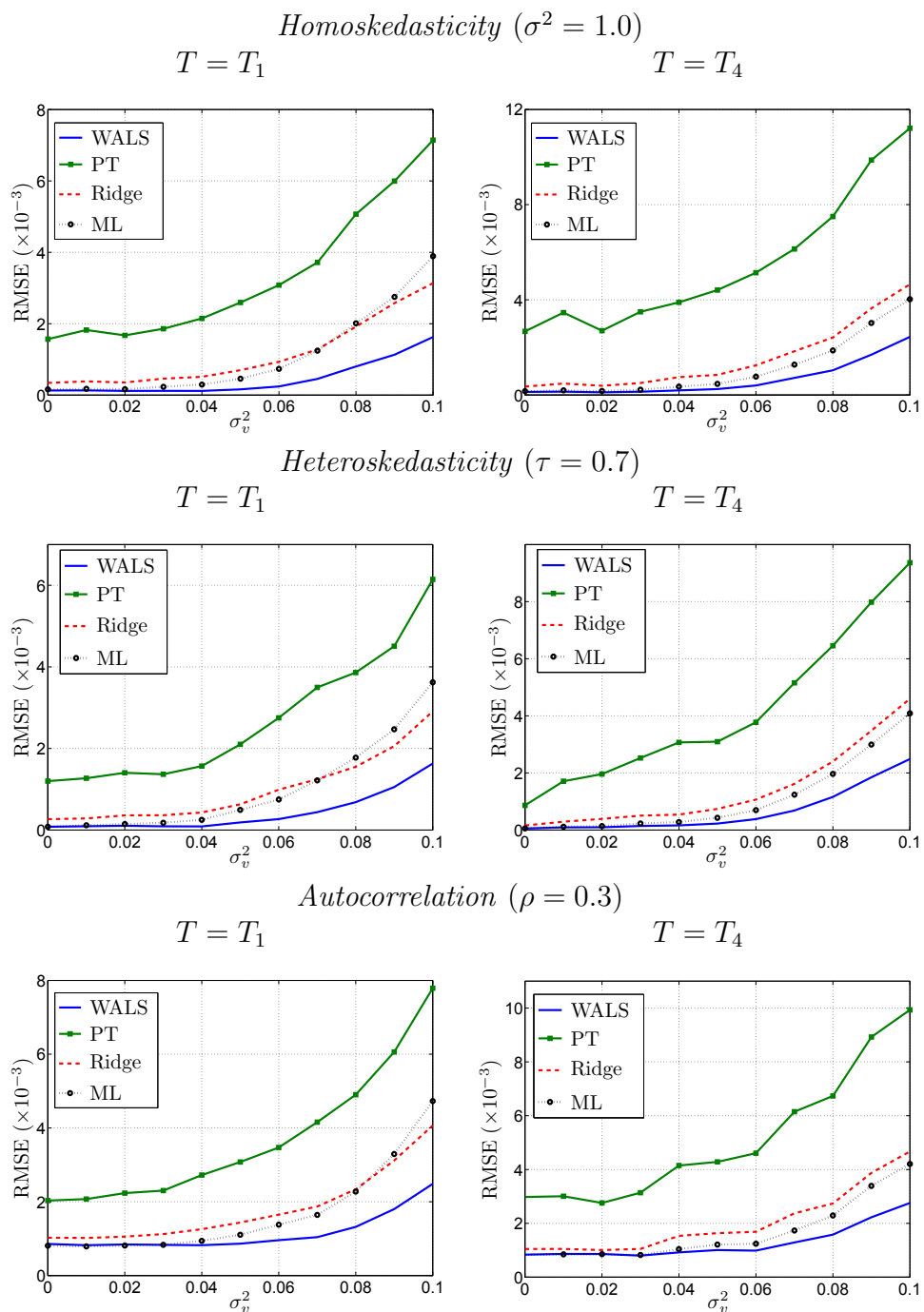
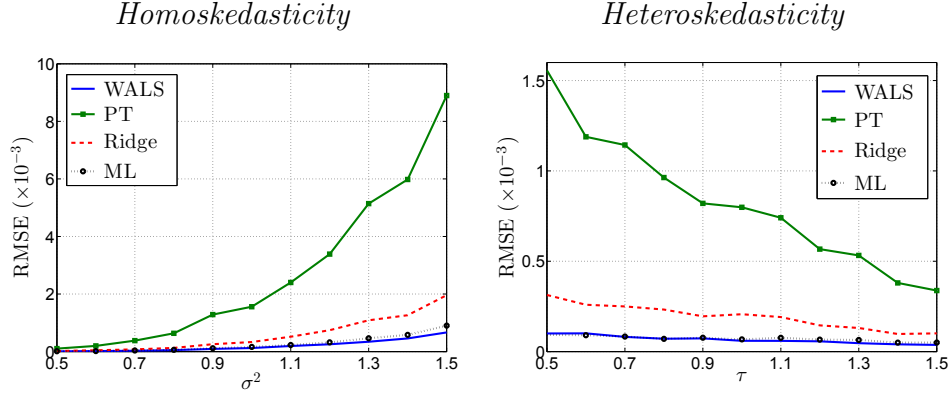




Figure 3: RMSE of prediction variance: homoskedastic versus heteroskedastic ( $N = 100, T = T_1$ )



## 11 Simulation results: More regressors

In Sections 9 and 10 we assumed two focus regressors, three auxiliary regressors, and one omitted regressor. In practical applications the number of regressors is likely to be larger. In this section we extend the benchmark framework by assuming  $k_2 = 12$  auxiliary regressors and  $k_3 = 3$  omitted regressors, while keeping the same number  $k_1 = 2$  of focus regressors. The large number of auxiliary regressors will increase the model uncertainty, because we now have  $2^{12} = 4096$  different models to consider compared to  $2^3 = 8$  in the benchmark. When introducing new variables we have to specify the ‘theoretical’  $t$ -ratios which are used to compute the values of the  $\beta$ -parameters. We consider four combinations, as in Table 5.

Table 5: Theoretical  $t$ -ratios for extended model

$T$	Auxiliary	Omitted
	$t_3-t_{14}$	$t_{15}-t_{17}$
$T_{L1}$	1.2, 0.9, 1.0, 1.3, 1.2, 1.5, 1.6, 1.2, 1.1, 0.8, 1.5, 1.4	0.0, 0.0, 0.0
$T_{L2}$	1.2, 0.9, 1.0, 1.3, 1.2, 1.5, 1.6, 1.2, 1.1, 0.8, 1.5, 1.4	2.4, 2.8, 2.0
$T_{L3}$	1.2, 0.9, 1.0, 2.3, 2.2, 2.5, 2.6, 2.1, 2.0, 0.5, 2.5, 1.4	0.0, 0.0, 0.0
$T_{L4}$	1.2, 0.9, 1.0, 0.7, 1.2, 0.5, 0.6, 2.2, 0.3, 0.8, 0.5, 1.2	0.0, 0.0, 0.0

In  $T_{L1}$  all auxiliary variables have  $t$ -ratios close to one and there are no omitted variables. In  $T_{L2}$  we have the same  $t$ -ratios for the auxiliary variables but now there are also omitted variables. In  $T_{L3}$  many of the auxiliary variables have ‘large’  $t$ -ratios, while in  $T_{L4}$  many of the  $t$ -ratios are ‘small’.

Only  $T_{L2}$  has omitted variables and they are all important. We combine this larger data set with the benchmark setup, random regressor DGP, and random coefficient DGP, again under each of the three error structures. We compare WALS, Ridge, and PT with ML. We do not compute MMA because the computational burden is too high when  $k_2$  is large. Some representative simulation results are presented in Table 6. Regarding the predictor, we see that WALS and Ridge perform best, better than ML and much better than PT. When we have fixed coefficients (and fixed or random regressors), WALS performs well in  $T_{L1}$  and  $T_{L3}$ , and Ridge dominates in  $T_{L2}$  and  $T_{L4}$ . With random coefficients WALS dominates more heavily. These findings, especially the comparison between  $T_{L3}$  (large  $t$ -ratios) and  $T_{L4}$  (small  $t$ -ratios), are in line with the small-data results. We see, as before, that WALS performs especially well under uncertainty and misspecification. Regarding the prediction variance, WALS performs best, followed by ML and Ridge, and much better than PT. Here too, the simulation results show that more uncertainty works in favor of WALS.

We briefly consider two other extensions, both analyzed in the context of the small data set: dependence among the regressors and non-normality. Dependence is introduced through a multivariate normal distribution with correlation 0.5, while the non-normal regressors are obtained from a Student distribution with five degrees of freedom. We experiment (separately) with these two extensions in the benchmark model and also in models with more uncertainty. The simulation results are largely similar to the case with normal and uncorrelated regressors. In particular, the WALS predictor is the most accurate when  $t$ -ratios are close to one, and the WALS prediction variance is particularly reliable when there is additional uncertainty.

## 12 Conclusion

This paper has introduced a new method of prediction averaging using weighted average least squares (WALS). We have argued that pretesting—the currently dominant prediction method—is dangerous, because it ignores the noise associated with model selection. Indeed, our simulation results demonstrate that pretesting performs very badly. Model averaging is an attractive method in that it allows us to combine model selection and prediction into one procedure. Within the model averaging methods we proposed the WALS predictor and also an estimate for its variance. Our predictor explicitly allows for correlation in the errors.

We have compared the WALS predictor with four competing predictors (unrestricted ML, pretesting, ridge regression, Mallows model averaging) in

Table 6: RMSE relative to ML, many auxiliary regressors ( $N = 100$ )

$T$	<i>Homoskedasticity</i> ( $\sigma^2 = 1.0$ )			<i>Heteroskedasticity</i> ( $\tau = 0.7$ )			<i>Autocorrelation</i> ( $\rho = 0.3$ )		
	WALS	PT	Ridge	WALS	PT	Ridge	WALS	PT	Ridge
<i>Benchmark model: fixed X, fixed <math>\beta</math></i>									
<i>Predictor</i>									
$T_{L1}$	<b>0.8611</b>	1.2088	0.8982	<b>0.8570</b>	1.1664	0.9411	<b>0.9284</b>	1.0807	0.9561
$T_{L2}$	0.9289	1.0789	<b>0.9234</b>	0.8755	1.0167	<b>0.8328</b>	0.9398	1.0637	<b>0.9229</b>
$T_{L3}$	<b>0.9625</b>	1.3057	0.9883	<b>0.9047</b>	1.2300	0.9294	<b>0.9052</b>	1.1423	0.9077
$T_{L4}$	0.8285	1.0579	<b>0.7997</b>	0.8035	0.9935	<b>0.7745</b>	0.8993	1.0058	<b>0.8760</b>
<i>Prediction variance</i>									
$T_{L1}$	<b>0.3440</b>	14.835	0.8088	<b>0.7854</b>	22.311	1.6711	1.3175	13.641	1.4356
$T_{L2}$	1.1154	17.104	<b>0.9023</b>	1.1196	25.152	1.4803	1.2224	15.158	1.5576
$T_{L3}$	<b>0.4239</b>	20.146	0.9557	<b>0.6404</b>	27.417	1.2445	1.1870	16.359	1.2598
$T_{L4}$	<b>0.3452</b>	8.6851	0.6884	<b>0.6729</b>	12.559	1.0391	1.2984	7.9238	1.3222
<i>Random regressor model: random X, fixed <math>\beta</math></i>									
<i>Predictor</i>									
$T_{L1}$	<b>0.9787</b>	1.0249	0.9791	0.9828	1.0136	<b>0.9823</b>	<b>0.9810</b>	1.0188	0.9821
$T_{L2}$	0.9769	1.0182	<b>0.9761</b>	0.9811	1.0106	<b>0.9801</b>	0.9778	1.0132	<b>0.9753</b>
$T_{L3}$	<b>0.9854</b>	1.0365	0.9874	<b>0.9891</b>	1.0249	0.9908	<b>0.9868</b>	1.0327	0.9876
$T_{L4}$	0.9745	1.0023	<b>0.9715</b>	0.9821	1.0000	<b>0.9800</b>	0.9770	1.0022	<b>0.9754</b>
<i>Prediction variance</i>									
$T_{L1}$	0.7713	1.0240	<b>0.7595</b>	<b>0.7634</b>	1.0329	0.7830	<b>0.7650</b>	1.0538	0.7779
$T_{L2}$	0.8569	1.0177	<b>0.8489</b>	0.8491	1.0024	0.8433	0.8523	1.0167	<b>0.8438</b>
$T_{L3}$	<b>0.8324</b>	1.0074	0.8443	<b>0.8277</b>	1.0003	0.8568	<b>0.8253</b>	1.0184	0.8498
$T_{L4}$	0.7365	1.0358	<b>0.7195</b>	<b>0.7299</b>	1.0636	0.7516	<b>0.7303</b>	1.0899	0.7540
<i>Random coefficient model: fixed X, random <math>\beta</math></i>									
<i>Predictor</i>									
$T_{L1}$	<b>0.9896</b>	1.0216	0.9960	0.9913	1.0050	<b>0.9907</b>	<b>0.9844</b>	1.0179	0.9876
$T_{L2}$	<b>0.9926</b>	1.0400	1.0166	<b>0.9747</b>	0.9978	0.9836	0.9683	1.0027	<b>0.9616</b>
$T_{L3}$	<b>0.9748</b>	1.0309	0.9775	<b>0.9941</b>	1.0203	0.9966	<b>0.9966</b>	1.0333	0.9988
$T_{L4}$	<b>0.9763</b>	1.0157	0.9893	<b>0.9846</b>	1.0067	0.9866	0.9839	1.0055	<b>0.9833</b>
<i>Prediction variance</i>									
$T_{L1}$	<b>0.1869</b>	6.9005	0.5812	<b>0.1577</b>	7.6836	0.7459	<b>0.4234</b>	8.3586	0.9135
$T_{L2}$	<b>0.3387</b>	11.519	0.8175	<b>0.1607</b>	9.1310	0.8254	<b>0.3520</b>	8.2753	0.8903
$T_{L3}$	<b>0.2718</b>	8.3617	0.7382	<b>0.2199</b>	9.3498	0.8352	<b>0.4600</b>	10.759	0.9792
$T_{L4}$	<b>0.1638</b>	4.2507	0.4721	<b>0.1532</b>	5.1157	0.6335	<b>0.4300</b>	5.7437	0.8624

a wide range of simulation experiments, where we considered not only the accuracy of the predictor (measured by the root mean squared prediction error), but also the accuracy of the prediction variance. The WALS predictor generally produces the lowest mean squared error. The estimated variance of the WALS predictor, while typically larger than the variance of the pretesting and ridge predictors, is more accurate, and when model uncertainty increases the dominance of WALS becomes more pronounced. These results, together with the fact that the WALS predictor is easy to compute, suggest that the WALS predictor is a serious candidate in economic prediction and forecasting.

## References

- Aiolfi, M. and A. Timmermann (2006). Persistence of forecasting performance and conditional combination strategies, *Journal of Econometrics*, 135, 31–53.
- Bates, J.M. and C.W.J. Granger (1969). The combination of forecasts, *Operational Research Quarterly*, 20, 451–468.
- Beran, R. (1995). Prediction in random coefficient regression, *Journal of Statistical Planning and Inference*, 43, 205–213.
- Bjørnland, H.C., K. Gerdrup, A.S. Jore, C. Smith, and L.A. Thorsrud (2012). Does forecast combination improve Norges Bank inflation forecasts?, *Oxford Bulletin of Economics and Statistics*, 74, 163–179.
- Bondeson, J. (1990). Prediction in random coefficient regression models, *Biometrical Journal*, 32, 387–405.
- Buckland, S.T., K.P. Burnham, and N.H. Augustin (1997). Model selection: An integral part of inference, *Biometrics*, 53, 603–618.
- Danilov, D. and J.R. Magnus (2004a). On the harm that ignoring pretesting can cause, *Journal of Econometrics*, 122, 27–46.
- Danilov, D. and J.R. Magnus (2004b). Forecast accuracy after pretesting with an application to the stock market, *Journal of Forecasting*, 23, 251–274.
- Elliott, G. and A. Timmermann (2004). Optimal forecast combinations under general loss functions and forecast error distributions, *Journal of Econometrics*, 122, 47–79.

- Froehlich, B.R. (1973). Some estimators for a random coefficient regression model, *Journal of the American Statistical Association*, 68, 329–335.
- Giles, J.A. and D.E.A. Giles (1993). Pre-test estimation and testing in econometrics: Recent developments, *Journal of Economic Surveys*, 7, 145–197.
- Golub, G.H., M. Heath, and G. Wahba (1979). Generalized cross-validation as a method for choosing a good ridge parameter, *Technometrics*, 21, 215–223.
- Granger, C.W.J. (2003). Time series concepts for conditional distributions, *Oxford Bulletin of Economics and Statistics*, 65, 689–701.
- Hansen, B.E. (2007). Least squares model averaging, *Econometrica*, 75, 1175–1189.
- Hansen, B.E. (2008). Least-squares forecast averaging, *Journal of Econometrics*, 146, 342–350.
- Hendry, D.F. and M.P. Clements (2004). Pooling of forecasts, *Econometrics Journal*, 7, 1–31.
- Hildreth, C. and J.P. Houck (1968). Some estimators for a linear model with random coefficients, *Journal of the American Statistical Association*, 63, 584–595.
- Hjort, N.L. and G. Claeskens (2003). Frequentist model average estimators, *Journal of the American Statistical Association*, 98, 879–899.
- Hoerl, A.E. and R.W. Kennard (1970). Ridge regression: Biased estimation for nonorthogonal problems, *Technometrics*, 12, 55–67.
- Jackson, T. and S. Karlsson (2004). Finding good predictors for inflation: A Bayesian model averaging approach, *Journal of Forecasting*, 23, 479–496.
- Johnston, J. and J. DiNardo (1997). *Econometric Methods*, Fourth Edition, McGraw-Hill, New York.
- Magnus, J.R. (1999). The traditional pretest estimator, *Theory of Probability and Its Applications*, 44, 293–308.
- Magnus, J.R. (2002). Estimation of the mean of a univariate normal distribution with known variance, *Econometrics Journal*, 5, 225–236.

- Magnus, J.R., O. Powell, and P. Prüfer (2010). A comparison of two model averaging techniques with an application to growth empirics, *Journal of Econometrics*, 154, 139–153.
- Magnus, J.R., A.T.K. Wan, and X. Zhang (2011). Weighted average least squares estimation with nonspherical disturbances and an application to the Hong Kong housing market, *Computational Statistics and Data Analysis*, 55, 1331–1341.
- Ohtani, K. and T. Toyoda (1980). Estimation of regression coefficients after a preliminary test for homoscedasticity, *Journal of Econometrics*, 12, 151–159.
- Poirier, D.J. (1995). *Intermediate Statistics and Econometrics: A Comparative Approach*, MIT Press, Cambridge, MA.
- Rubin, H. (1950). Note on random coefficients, in: T.C. Koopmans (Ed.), *Statistical Inference in Dynamic Economic Models*, Cowles Commission Monograph No. 10, pp. 419–421.
- Ruud, P.A. (2000). *An Introduction to Classical Econometric Theory*, Oxford University Press, New York/Oxford.
- Stock, J.H. and M.W. Watson (2004). Combination forecasts of output growth in a seven-country data set, *Journal of Forecasting*, 23, 405–430.
- Swamy, P.A.V.B. (1970). Efficient inference in a random coefficient regression model, *Econometrica*, 38, 311–323.
- Timmermann, A. (2006). Forecast combinations, in: G. Elliott, C.W.J. Granger, and A. Timmermann (Eds.), *Handbook of Economic Forecasting*, Vol. 1, Elsevier, Amsterdam, pp. 135–196.
- Toyoda, T. and T.D. Wallace (1976). Optimal critical values for pre-testing in regression, *Econometrica*, 44, 365–375.
- Wan, A.T.K., X. Zhang, and G. Zou (2010). Least squares model averaging by Mallows criterion, *Journal of Econometrics*, 156, 277–283.
- Wan, A.T.K. and G. Zou (2003). Optimal critical values of pre-tests when estimating the regression error variance: Analytical findings under a general loss structure, *Journal of Econometrics*, 114, 165–196.

- Whittle, P. (1963). *Prediction and Regulation by Linear Least-Square Methods*, The English Universities Press Ltd, London.
- Yang, Y. (2004). Combining forecasting procedures: Some theoretical results, *Econometric Theory*, 20, 176–222.