

# Forecasting with Model Uncertainty: Representations and Risk Reduction\*

Keisuke Hirano<sup>†</sup>

Jonathan H. Wright<sup>‡</sup>

First version: October 14, 2013

This version: April 5, 2016

## Abstract

We consider forecasting with uncertainty about the choice of predictor variables. The researcher wants to select a model, estimate the parameters, and use the parameter estimates for forecasting. We investigate the distributional properties of a number of different schemes for model choice and parameter estimation: in-sample model selection using the Akaike information criterion, out-of-sample model selection, and splitting the data into subsamples for model selection and parameter estimation. Using a weak-predictor local asymptotic scheme, we provide a representation result that facilitates comparison of the distributional properties of the procedures and their associated forecast risks. We develop a simulation procedure that improves the accuracy of the out-of-sample and split-sample methods uniformly over the local parameter space. We also examine how bootstrap aggregation (bagging) affects the local asymptotic risk of the estimators and their associated forecasts. Numerically, we find that for many values of the local parameter, the out-of-sample and split-sample schemes perform poorly if implemented in the conventional way. But they perform well, if implemented in conjunction with our risk-reduction method or bagging.

---

\*We are grateful to Don Andrews, Marine Carrasco, Russell Davidson, Gary Chamberlain, Sylvia Gonçalves, Bruce Hansen, Serena Ng, Peter Phillips, Jack Porter, four anonymous referees and the co-editor for very helpful discussions. The usual disclaimer applies.

<sup>†</sup>Department of Economics, University of Arizona, 1130 E.Helen St., Tucson AZ 85721. Email: [hirano@u.arizona.edu](mailto:hirano@u.arizona.edu)

<sup>‡</sup>Department of Economics, Johns Hopkins University, 3400 North Charles St., Baltimore MD 21218. Email: [wrightj@jhu.edu](mailto:wrightj@jhu.edu)

# 1 Introduction

In this paper, we reconsider the problem of forecasting when there is uncertainty about the forecasting model. As is well known, a model that fits well in sample may not be good for forecasting—a model may fit well in-sample, only to turn out to predict poorly. Consequently, it is common practice to select the model based on pseudo-out-of-sample fit from a sequence of recursive or rolling predictions. Parameters are then estimated over the whole sample period. The idea of using an out-of-sample criterion is discussed by Clark (2004) and West (2006). It is an idea with a long history, going back to Wilson (1934), Meese and Rogoff (1983), and Ashley, Granger, and Schmalensee (1980), and is very intuitive: it is what a researcher could have done at the time. Instead, one might select the model based on in-sample fit, but adjust for overfitting by using an information criterion, such as the Akaike Information Criterion (AIC) (Akaike, 1974), as advocated by Inoue and Kilian (2006).

We consider a pseudo-likelihood setting with a fixed number  $k$  of potential parameters to be estimated, each of which has a coefficient that is local to zero. The concept of model selection that we envision amounts to selecting a set of zero restrictions; in a regression setting, for example, this would indicate which predictors are excluded from the regression. Thus there are up to  $2^k$  possible models among which we can choose. Having chosen the model, we then have to estimate the parameters and use these for forecasting. Although some model will be best in terms of predictive accuracy, the local-to-zero nesting means that we can never consistently select that model. We consider various methods of model selection and forecasting, including: using in-sample fit with the AIC information criterion; selecting the model based on recursive pseudo-out-of-sample forecast accuracy and then using the whole dataset for parameter estimation; and splitting the sample into two parts, using one part for model selection and the other for parameter estimation. We call this last method the split-sample approach. Unlike the first two methods, it is not commonly used in practice. But it does ensure asymptotic independence between parameter estimates and model selection, unlike methods based on in-sample fit (Leeb and Pötscher, 2005; Hansen, 2009), and also unlike the standard out-of-sample approach.

We obtain asymptotic characterizations of these forecasting procedures under the local parameter sequence. A key step is to obtain an asymptotic representation of the partial sum process for the score function as the sum of a term that is directly informative about the local parameters, and another term

that is an independent Gaussian process. This allows us to provide a limit-experiment type representation of the procedures, from which we can calculate normalized local asymptotic mean square prediction errors up to  $O(T^{-1})$  terms. We show that the recursive pseudo-out-of-sample and split-sample procedures are inefficient, in the sense that their limit distributions depend on the ancillary Gaussian noise process.

Our characterizations also suggest ways to improve upon these procedures. The influence of the ancillary noise term in the limiting asymptotic representation can be eliminated by a conditioning argument. We can implement this noise reduction via a simulation-and-averaging scheme; doing this is shown to uniformly improve the out-of-sample and split-sample methods asymptotically for a wide variety of loss functions.

This method is related to bootstrap aggregating (bagging) (Breiman, 1996) in which the data are resampled, the forecasting method is applied to the resampled data, and the resulting forecasts are then averaged over all the bootstrap samples. Bagging has a smoothing effect that alters the risk properties of estimators, but averaging over bootstrap draws can also reduce the influence of the extraneous noise term in the out-of-sample and split-sample methods. Earlier theoretical work on bagging, notably Bühlmann and Yu (2002), emphasized its smoothing effect rather than the noise reduction effect, though Efron (2014) observes that it can reduce variability<sup>1</sup>.

We then numerically compare the various procedures in terms of their local asymptotic risk. In their standard forms the rank ordering among the in-sample, out-of-sample and split-sample methods depends on the localization parameter which is not consistently estimable (although one can form a confidence interval for it). Nonetheless, we find that for many values of the localization parameter, in-sample forecasting using the AIC gives the most accurate forecasts, out-of-sample prediction does worse, and the split-sample method does worst of all. This is intuitive because the out-of-sample and split-sample schemes are in some sense wasting data, and is essentially the argument of Inoue and Kilian (2004) and Inoue and Kilian (2006) for the use of in-sample rather than out-of-sample predictability tests. However, introducing the simulation or bagging step changes the rank ordering substantially, and entirely reverses it for many values of the localization parameter. Our simulation scheme has no effect on in-

---

<sup>1</sup>One other useful feature of the out-of-sample forecasting setup is that it can be constructed to use only real-time data which precisely mimics the data available to a researcher in the presence of data revisions. Unfortunately, adding our simulation scheme or bootstrap aggregation step destroys this feature.

sample forecasts, but reduces the local asymptotic mean square prediction error of the out-of-sample and split-sample forecasts *uniformly* in the localization parameter, and the reductions are generally numerically large. Bagging can modestly reduce the local asymptotic mean square prediction error of the in-sample forecasts over some parts of the parameter space, but it makes a more dramatic difference to the out-of-sample and split-sample forecasts. In short, the out-of-sample and split-sample forecasts are poor choices in their conventional form, but become very competitive when combined with a simulation or bagging step.

In the next section, we set up the model and introduce the various procedures we will evaluate. In Section 3, we derive asymptotic characterizations via our representation theorem for the partial sum process. Section 4 contains some extensions. Section 5 explores the local asymptotic risk properties of the procedures numerically. Section 6 examines some finite-sample simulation evidence. Section 7 contains an illustrative application, and Section 8 concludes.

## 2 Pseudo-Likelihood Framework

We observe  $(y_t, x_t)$  for  $t = 1, \dots, T$  and wish to forecast  $y_{T+1}$  given knowledge of  $x_{T+1}$ . Let the pseudo log (conditional) likelihood be

$$\ell(\beta) = \sum_{t=1}^T \ell_t(\beta) = \sum_{t=1}^T \ln f(y_t | x_t, \beta), \quad (2.1)$$

where  $f$  is a conditional density function and the parameter  $\beta$  is  $k \times 1$ . This framework could apply to cross-sectional regression of an outcome variable  $y_t$  on a  $k \times 1$  vector of predictors  $x_t$ ,  $h$ -step ahead forecasting regressions (where  $x_t$  are suitably lagged predictor variables), vector autoregression (where  $x_t$  contains lagged values of the vector  $y_t$ ), and nonlinear regression models. There may be unmodeled dependence, subject to the large-sample distributional assumptions imposed below.

Model selection amounts to setting some elements of  $\beta$  to zero, and estimating the others. Thus there are up to  $2^k$  possible models. Let  $m \subset \{1, \dots, k\}$  denote a model, with the interpretation that the elements of  $m$  indicate which coefficients of  $\beta$  are allowed to be nonzero. The set of possible models  $\mathcal{M}$  is a subset of the power set of  $\{1, \dots, k\}$ .

We consider a variety of strategies for model selection and parameter estimation. Each strategy will end

up giving us an estimator for  $\beta$ , some elements of which may be zero. We denote this overall estimator as  $\tilde{\beta}$ . The strategies that we consider are:

1. **Maximum Likelihood (MLE).** Set  $\tilde{\beta}$  to the unrestricted (pseudo) maximum likelihood estimator,  $\hat{\beta}$ , that maximizes  $\ell(\beta)$ . We assume that this is consistent for the pseudo-true value.
2. **James-Stein (JS).** The positive-part James-Stein estimator uses the unrestricted estimate  $\hat{\beta}$  and an estimate  $\hat{V}$  of its asymptotic variance-covariance matrix. The JS estimator for  $k > 2$  is:

$$\tilde{\beta} = \hat{\beta} \cdot \max \left\{ 1 - \frac{k-2}{T \hat{\beta}' \hat{V}^{-1} \hat{\beta}}, 0 \right\}. \quad (2.2)$$

3. **Laplace estimator.** Let  $t$  denote the vector of  $t$ -statistics consisting of the elements of  $\hat{\beta}$  each divided by the  $j$ th diagonal element of  $\hat{V}$ . The Laplace estimator of  $\beta$  is:

$$\tilde{\beta} = \hat{\beta} \circ \tilde{h}(t) \quad (2.3)$$

where  $\tilde{h}(x) = \left(1 - \frac{ch(x)}{x}\right)$  (element-by-element),  $h(x) = \frac{e^{-cx}\Phi(x-c) - e^{cx}\Phi(-x-c)}{e^{-cx}\Phi(x-c) + e^{cx}\Phi(-x-c)}$ ,  $c = \ln(2)$  and  $\Phi(\cdot)$  is the standard normal cdf. In the scalar case, this is equivalent to the Bayesian posterior mean corresponding to a certain Laplace prior distribution, and is therefore admissible (Magnus, 2002; Magnus, Powell, and Prüfer, 2010).

4. **LASSO Estimator.** Set  $\tilde{\beta}$  to maximize  $l(\beta) - \lambda_T \sum_{j=1}^k |\beta_j|$ , where  $\lambda_T$  is a penalty term. Penalized maximum likelihood estimators of this sort have been considered by Tibshirani (1996) and many others.
5. **Small Model.** Set  $\tilde{\beta} = 0$ .
6. **AIC (In-sample).** For each model  $m \in \mathcal{M}$ , let  $\hat{\beta}(m)$  denote the restricted pseudo-ML estimator that maximizes  $\ell(\beta)$  subject to the zero restrictions implied by  $m$ . (Thus,  $\hat{\beta}(m)$  is a  $k \times 1$  vector with zeros in the locations implied by  $m$ .) Let  $n(m)$  be the number of free parameters in model  $m$ . For each  $m \in \mathcal{M}$  we calculate the AIC objective function

$$AIC(m) = 2\ell(\hat{\beta}(m)) - 2n(m), \quad (2.4)$$

and choose the model  $m^*$  that maximizes  $AIC(m)$ . Then set  $\tilde{\beta} = \hat{\beta}(m^*)$ .

7. **Out of Sample.** For each model  $m$ , we estimate the model recursively starting a fraction  $\pi \in (0, 1)$  of the way through the sample, and calculate its one-period-ahead predictive density to obtain a pseudo out-of-sample estimate of predictive performance. Let  $\hat{\beta}_{1,t-1}(m)$  denote the pseudo maximum likelihood estimate for model  $m$  using observations 1 to  $t - 1$ . For each  $m$ , we calculate

$$\sum_{t=[T\pi]+1}^T \ell_t(\hat{\beta}_{1,t-1}(m)). \quad (2.5)$$

We then choose the model  $m$  that maximizes this predictive likelihood, and use the full sample for estimation of the model.<sup>2</sup>

8. **Split-sample.** For each model  $m$ , we calculate AIC using data up to a fraction  $\pi$  of the way through the sample:

$$AIC_{ss}(m) = 2\ell(\hat{\beta}_{1,[T\pi]}(m)) - 2n(m). \quad (2.6)$$

For  $m^* = \arg\max AIC_{ss}(m)$ , we use the second fraction of the sample to estimate the model parameters:

$$\tilde{\beta} = \hat{\beta}_{[T\pi]+1,T}(m^*). \quad (2.7)$$

Later in the paper, we also consider modifications of these procedures that attempt to improve their risk properties.

We focus on obtaining limiting distributions for these estimators, which we denote generically as  $\tilde{\beta}$ . This may be of direct interest in itself as a parameter estimation problem. Alternatively, it may be useful for forecasting. Suppose that a forecaster creates a point forecast of  $y_{T+1}$  using a rule  $\hat{y}_{T+1}(\tilde{\beta})$ , where  $\tilde{\beta}$  denotes an estimator of  $\beta$ . The function  $\hat{y}_{T+1}(\cdot)$  may depend on  $x_{T+1}$ , and it may depend on the data up to time  $T$  in other ways. Given a forecasting loss function  $L(\hat{y}_{T+1}, y_{T+1})$ , let

$$\bar{L}(a, \beta) = E[L(\hat{y}_{T+1}(a), \beta)], \quad (2.8)$$

---

<sup>2</sup>Several authors test for the statistical significance of differences in out-of-sample forecasting performance with one model, typically a simple benchmark, as the null (see, for example Diebold and Mariano (1995) and Hansen and Timmermann (2013)). Here we are instead thinking of selecting the model based on the point estimate of its pseudo out-of-sample predictive performance.

where the expectation is with respect to the conditional distribution of  $y_{T+1}$  under parameter value  $\beta$ . Define the *regret* as

$$\bar{L}_r(a, \beta) = \bar{L}(a, \beta) - \inf_c \bar{L}(c, \beta). \quad (2.9)$$

This subtracts off the expected forecast loss associated with the infeasible optimal choice of the estimate used in the forecast rule. The regret risk is defined as

$$R(\tilde{\beta}, \beta) = E \left[ \bar{L}_r(\tilde{\beta}, \beta) \right], \quad (2.10)$$

where the expectation is with respect to the sampling distribution of  $\tilde{\beta}$ , again under the “true” parameter value  $\beta$ .

Typically, our estimators  $\tilde{\beta}$  will be  $\sqrt{T}$ -consistent for  $\beta$ , so a suitably normalized version of the forecast risk will depend on the local behavior of  $\bar{L}_r(a, \beta)$  as  $a \rightarrow \beta$ . For a number of standard forecast loss functions, and suitable choices for  $\hat{y}_{T+1}(\cdot)$ , we will have an expansion of the form

$$T \cdot \bar{L}_r(a, \beta) = G\left(\sqrt{T}(a - \beta)\right) + o(\|a - \beta\|), \quad (2.11)$$

for some convex function  $G$ . Then the normalized regret risk for the forecast rule  $\hat{y}_{T+1}(\tilde{\beta})$  will satisfy

$$T \cdot R(\tilde{\beta}, \beta) = E \left[ G\left(\sqrt{T}(\tilde{\beta} - \beta)\right) \right] + o_p(1). \quad (2.12)$$

Thus obtaining the limiting distribution of  $\sqrt{T}(\tilde{\beta} - \beta)$  will be key to characterizing the risk properties of the associated forecast rule. (The convexity of  $G$  will also prove to be useful for our risk reduction methods.) In the next section we give some examples of some cases where the representation in Equation (2.12) holds.

## 2.1 Example: Prediction in a Regression Model

To illustrate our approach in a simple setting, we consider prediction using a standard regression model:

$$y_t = \beta' x_t + u_t, \quad (2.13)$$

where the  $u_t$  are i.i.d. with mean 0, finite variance  $\sigma^2$ , and  $2+\delta$  finite moments for some  $\delta > 0$ . We assume  $x_t$  is a  $k \times 1$  i.i.d. vector that has been orthonormalized, so that  $E[x_t x_t'] = I_k$ . (The orthonormality and time-independence of  $x_t$  is not essential for the analysis, but simplifies the notation.)

This model fits into the general pseudo-likelihood framework of Section 2, using the standard Gaussian likelihood. Then  $\hat{\beta}$ , the unrestricted pseudo-ML estimator of  $\beta$ , is the OLS estimator; and  $\hat{\beta}(m)$ , the restricted pseudo-ML estimator under model  $m$ , is the restricted OLS estimator using only the regressors indicated by  $m$ . Each model corresponds to some subset of the  $k$  regressors that are to be used for forecasting.

Consider forecasting under squared error loss, with forecasts of the form  $\hat{y}_{T+1} = \tilde{\beta}' x_{T+1}$ , where  $\tilde{\beta}$  can again be any of the estimators of  $\beta$ . The mean squared prediction error is

$$\begin{aligned} MSPE &= E[(y_{T+1} - \tilde{\beta}' x_{T+1})^2] = E[(u_{T+1} - (\tilde{\beta} - \beta)' x_{T+1})^2] \\ &= \sigma^2 + E[(\tilde{\beta} - \beta)'(\tilde{\beta} - \beta)]. \end{aligned} \quad (2.14)$$

The first term on the right hand side of (2.14) is the risk of infeasible forecast rule that uses the true value of the parameter,  $\beta' x_{T+1}$ . The second term is  $O(T^{-1})$  and differs across forecasting methods. In the forecasting context, parameter estimation uncertainty is of small order relative to uncertainty in the distribution of the shocks. But uncertainty about the distribution of the shocks is common to all the models. We therefore subtract the term  $\sigma^2$  and scale the risk by  $T$  to obtain the normalized mean square prediction error criterion

$$NMSPE = T(MSPE - \sigma^2) = E\left[\sqrt{T}(\tilde{\beta} - \beta)' \sqrt{T}(\tilde{\beta} - \beta)\right]. \quad (2.15)$$

This is equivalent to the normalized regret risk, as defined above, for squared error loss.

Other loss functions lead to different expressions for the regret risk, but in some leading cases they will continue to satisfy Equation (2.12). For example suppose we work with linex forecast loss,

$$L(\hat{y}_{T+1}, y_{T+1}) = \exp(\theta(y_{T+1} - \hat{y}_{T+1})) - \theta(y_{T+1} - \hat{y}_{T+1}) - 1, \quad (2.16)$$

for some given  $\theta \neq 0$ , and suppose that  $u_{T+1}$  is Gaussian. In this case, the oracle point forecast using



knowledge of  $\beta$  and  $\sigma^2$  is  $y_{T+1}^* = \beta' x_{T+1} + \frac{\theta \sigma^2}{2}$ , so it is natural to consider rules of the form

$$\hat{y}_{T+1}(\tilde{\beta}) = \tilde{\beta}' x_{T+1} + \frac{\theta \hat{\sigma}^2}{2},$$

where  $\hat{\sigma}^2$  is a consistent estimator of  $\sigma^2$ . This forecast rule employs an estimate of the optimal “shift” of the naive forecast to account for the asymmetry of the linex loss function. Using a second-order Taylor expansion, it can be shown that the normalized regret risk for this forecast rule is

$$T \cdot R(\tilde{\beta}, \beta) = \frac{\theta^2}{2} E \left[ \sqrt{T}(\tilde{\beta} - \beta)' \sqrt{T}(\tilde{\beta} - \beta) \right] + o_P(1), \quad (2.17)$$

again satisfying (2.12).

To gain further intuition for our theoretical results in the next section, consider the special case where the  $u_t$  are i.i.d.  $N(0, 1)$  and the regressors are treated as fixed and satisfy  $\frac{1}{T} \sum_{t=1}^T x_t x_t' = I_k$ . Then the least squares estimator for the full set of parameters has an exact normal distribution:

$$\hat{\beta} = \left( \sum_{t=1}^T x_t x_t' \right)^{-1} \sum_{t=1}^T x_t y_t \sim N(\beta, \sigma^2 I_k / n) \quad (2.18)$$

and  $\hat{\beta}$  is a minimal sufficient statistic for  $\beta$ . If a procedure makes nontrivial use of information in the data other than that contained in  $\hat{\beta}$ , it is introducing an unnecessary source of randomness. In the next section we will obtain an asymptotic analog to this argument in the general pseudo-likelihood setting, and show how it applies to the various procedures we consider.

### 3 Local Asymptotics

In order to capture the role of parameter and model uncertainty in our analysis, the joint distribution of  $\{(y_1, x_1), \dots, (y_T, x_T)\}$  is assumed to be a triangular array with drifting parameters. Let  $\{(y_1, x_1), \dots, (y_T, x_T)\}$  have joint distribution  $P_T$ , and define the pseudo-true value of the parameter as

$$\beta_{0,T} = \arg \max_{\beta} \int \ell(\beta) dP_T. \quad (3.1)$$

We take the pseudo-true values (or functions of them) as our objects of interest. We will take limits as  $T \rightarrow \infty$  under the assumption that

$$\beta_{0,T} = \frac{b}{\sqrt{T}}, \quad b \in \mathbb{R}^k. \quad (3.2)$$

This type of drifting sequence has been used by Claeskens and Hjort (2008) and Inoue and Kilian (2004) and others to study the large-sample properties of model selection procedures. It preserves the role of parameter uncertainty in the asymptotic approximations, unlike fixed-alternative asymptotics in which model selection can determine which coefficients are nonzero with probability approaching one. Later in the paper, we will show some Monte-Carlo evidence indicating that this local parameterization provides a good approximation in small samples. The analysis could be extended to allow some components of  $\beta$  to be localized away from zero, corresponding to situations where some components of  $\beta$  are known to be nonzero. We use  $\rightarrow_d$  to denote weak convergence and  $\rightarrow_p$  to denote convergence in probability under the sequence of measures  $\{P_T\}_{T=1}^\infty$ . Our results below depend crucially on the convergence properties of the partial sums of the pseudo-likelihood function. We make the following high level assumptions.

**Assumption 3.1**

$$T^{-1/2} \sum_{t=1}^{[Tr]} \frac{\partial \ell_t(\beta_{0,T})}{\partial \beta} \rightarrow_d B(r),$$

where  $B(r)$  is a  $k$ -dimensional Brownian motion with positive definite covariance matrix  $\Lambda$ .

**Assumption 3.2** For all sequences  $\beta_T$  in a  $T^{-1/2}$ -neighborhood of zero,

$$-T^{-1} \sum_{t=1}^{[Tr]} \frac{\partial^2 \ell_t(\beta_T)}{\partial \beta \partial \beta'} \rightarrow_p r \Sigma.$$

where  $\Sigma$  is positive definite.

**Assumption 3.3** For any fixed  $k \times k$  matrix  $C$ ,

$$\sum_{t=1}^{[Tr]} \sum_{s=1}^{t-1} \frac{\partial l_s(\beta_{0,T})'}{\partial \beta} \frac{1}{t-1} C \frac{\partial l_t(\beta_{0,T})}{\partial \beta} \rightarrow_d \int_0^r \frac{1}{s} B(s) C dB(s).$$

These high-level assumptions would follow from conventional regularity conditions in correctly specified parametric models. In misspecified models, the assumptions require that the pseudo-true param-

ter sequence  $\beta_{0,T}$  is related to the distribution of the data in a smooth way.

To gain intuition for the results to follow, consider the case where the parametric model with conditional likelihood  $f(y_t|x_t, \beta)$  is correctly specified. Then, under standard regularity conditions, Assumptions 3.1, 3.2 and 3.3 will hold with  $\Lambda = \Sigma$ . Furthermore, the model will be locally asymptotically normal (LAN), and possess a limit experiment representation (see for example van der Vaart 1998, Chs. 7-9). In particular, consider any estimator sequence  $\tilde{\beta}$  with limiting distributions in the sense that

$$T^{1/2} \tilde{\beta} \rightarrow_d \mathcal{L}_b, \quad (3.3)$$

where the limit is taken under the drifting sequences of measures corresponding to  $\beta_{0,T} = T^{-1/2}b$ , and  $\mathcal{L}_b$  is a law that may depend on  $b$ . Then the estimator  $\tilde{\beta}$  has an asymptotic representation as a randomized estimator in a shifted normal model: if  $Y$  is a single draw from the  $N(\Sigma b, \Sigma)$  distribution, and  $U$  is random variable independent of  $Y$  (with sufficiently rich support<sup>3</sup>), there exists an estimator  $S(Y, U)$  with

$$S(Y, U) \sim \mathcal{L}_b \quad (3.4)$$

for all  $b$ . In other words, the sequence  $T^{1/2} \tilde{\beta}$  is asymptotically equivalent to the randomized estimator  $S$  under all values of the local parameter.

We extend this type of asymptotic representation, in terms of an asymptotically sufficient component and an independent randomization, to the pseudo-likelihood setup. We do this by establishing a large-sample representation of the partial sum process for the score function that corresponds to the  $(Y, U)$  limit experiments in parametric LAN models.

From Assumptions 3.1 and 3.2, it follows that:

$$T^{-1/2} \sum_{t=1}^{[Tr]} \frac{\partial \ell_t(0)}{\partial \beta} \rightarrow_d B(r) + r \Sigma b =: Y(r) \quad (3.5)$$

Thus the partial sums of the score function evaluated at  $\beta = 0$  converge to a Brownian motion with linear drift. By a standard argument, we can decompose this process into the sum of a normal random vector and a Brownian bridge:

---

<sup>3</sup>Typically, a representation  $S(Y, U)$  exists for  $U$  distributed uniform on  $[0, 1]$ , but for our results below, it is useful to allow  $U$  to have a more general form.

**Proposition 3.4** *Under the drifting sequence in Equation (3.2) and Assumptions 3.1 and 3.2,*

$$T^{-1/2} \sum_{t=1}^{\lfloor Tr \rfloor} \frac{\partial \ell_t(0)}{\partial \beta} \rightarrow_d Y(r) = rY + U_B(r), \quad (3.6)$$

where  $Y := Y(1) \sim N(\Sigma b, \Lambda)$ , and  $U_B(r)$  is a  $k$ -dimensional Brownian bridge with covariance matrix  $\Lambda$ , where  $U_B$  is independent of  $Y$ .

All proofs are given in Appendix A. This result decomposes the limit of the partial sums of the score function into two stochastic components, one of which depends on the local parameter  $b$  and one of which is ancillary.

Let  $\Sigma(m)$  denote the  $k \times k$  matrix that consists of the elements of  $\Sigma$  in the rows and columns indexed by  $m$  and zeros in all other locations, and let  $H(m)$  denote the Moore-Penrose generalized inverse of  $\Sigma(m)$ . Then  $T^{1/2} \hat{\beta} \rightarrow_d \Sigma^{-1} Y(1)$  and  $T^{1/2} \hat{\beta}(m) \rightarrow_d H(m) Y(1, m)$ , where  $Y(r, m)$  denotes the  $k \times 1$  vector with the elements of  $Y(r)$  in the locations indexed by  $m$  and zeros elsewhere. This leads to the following asymptotic characterizations of the procedures:

**Proposition 3.5** *Under the drifting sequence in Equation (3.2) and Assumptions 3.1, 3.2 and 3.3, we have the following limiting representations of the parameter estimation procedures:*

(i) *Using unrestricted MLE:*

$$T^{1/2} \hat{\beta} \rightarrow_d \Sigma^{-1} Y(1) \quad (3.7)$$

(ii) *Using the positive-part James-Stein estimator:*

$$T^{1/2} \tilde{\beta} \rightarrow_d \Sigma^{-1} Y(1) \max \left\{ 1 - \frac{k-2}{Y(1)' \Sigma^{-2} Y(1)}, 0 \right\} \quad (3.8)$$

(iii) *Using the Laplace estimator:*

$$T^{1/2} \tilde{\beta} \rightarrow_d Y(1) \circ \tilde{h}(WY(1)) \quad (3.9)$$

where  $W$  is a diagonal matrix with the  $j$ th diagonal element equal to the reciprocal of the square root of the  $j$ th diagonal element of  $\Sigma^{-1}$ .

(iv) *Selecting the model using the AIC:*

$$T^{1/2} \tilde{\beta} \rightarrow_d \sum_{m^*} H(m^*) Y(1, m^*) \mathbf{1}\{m^* = \arg\max_m [Y(1, m)' H(m) Y(1, m) - 2n(m)]\} \quad (3.10)$$

(v) *Selecting the model minimizing recursive out-of-sample error starting a fraction  $\pi$  of the way through the sample:*

$$T^{1/2} \tilde{\beta} \rightarrow_d \sum_{m^*} H(m^*) Y(1, m^*) \mathbf{1}\{m^* = \arg\max_m [Y(1)' H(m) Y(1) - \frac{1}{\pi} Y(\pi)' H(m) Y(\pi) + \text{tr}(H(m) \Lambda) \log(\pi)]\} \quad (3.11)$$

(vi) *Using the split-sample method, using the first fraction  $\pi$  of the sample for model selection and the rest for parameter estimation:*

$$T^{1/2} \tilde{\beta} \rightarrow_d \sum_{m^*} H(m^*) \frac{Y(1, m^*) - Y(\pi, m^*)}{1 - \pi} \mathbf{1}\{m^* = \arg\max_m [\frac{1}{\pi} Y(\pi, m)' H(m) Y(\pi, m) - 2n(m)]\} \quad (3.12)$$

where  $\sum_{m^*}$  denotes the summation over all the models in  $\mathcal{M}$ .

Of course, there are other criteria besides AIC that we could use for in-sample model selection. Some of these are asymptotically equivalent to AIC, such as Mallows'  $C_p$  criterion (Mallows, 1973) or leave-one-out cross-validation. Using any of these information criteria for in-sample model selection will give the same asymptotic distribution as in equation (3.10). Alternatively, one could use the Bayes information criterion (BIC). In the present setting, because the penalty term goes to zero at a rate slower than  $T^{-1}$ , the BIC will pick the small model ( $\beta = 0$ ) with probability converging to one. Part (v) of the proposition can immediately be adapted to selecting the model minimizing out-of-sample error with a rolling estimation window, as long as the estimation window contains a fixed fraction of the sample size, but not if it instead contains a fixed number of observations as in Giacomini and White (2006).

Inoue and Kilian (2004) considered the local power of some in-sample and out-of-sample tests of the hypothesis that  $\beta = 0$ . They derived equation (3.7) and a result very similar to equation (3.11). The expression in equation (3.11) involves asymptotically selecting the model by maximizing an expression

that is the difference between two quadratic forms plus a nonstochastic term. This derivation is also closely related to Hansen and Timmermann (2013) who showed that the difference between the out-of-sample predictive accuracy of two models is asymptotically equivalent to the difference between two Wald statistics plus a nonstochastic term.

### 3.1 Rao-Blackwellization

The estimators other than the out-of-sample and split-sample estimators can be viewed as generalized shrinkage estimators (Stock and Watson, 2012) as their limiting distributions are of the form:  $T^{1/2}\tilde{\beta} \rightarrow_d Yg(Y)$  for some nonlinear function  $g(Y)$ . In contrast, the limiting distributions in equations (3.11) and (3.12) are functions of both  $Y$  and an independent Brownian bridge,  $U_B(r)$ . Their dependence on the noise term  $U = U_B$  suggests a novel way to improve them.

In the statistical experiment of observing the pair  $(Y, U)$ , where  $Y \sim N(\Sigma b, \Lambda)$  and  $U$  is ancillary, the variable  $Y$  is sufficient. Thus, for any estimator  $S(Y, U)$ , consider its conditional expectation

$$\tilde{S}(Y) := E[S(Y, U)|Y]. \quad (3.13)$$

By the Rao-Blackwell theorem, the risk of  $\tilde{S}(Y)$  is less than or equal to that of  $S(Y, U)$  for all  $b$  for any convex risk function. This includes the regret functions in equations (??) and (??).

To implement the conditional estimators, we need consistent estimators  $\hat{\Lambda} \rightarrow_p \Lambda$  and  $\hat{\Sigma} \rightarrow_p \Sigma$ . Dependence in the scores poses no problem, so long as  $\hat{\Lambda}$  is a consistent estimate of the zero-frequency spectral density. Recall that  $T^{1/2}\hat{\beta}(m) \rightarrow_d H(m)Y(1, m)$ . Then take  $L$  independent artificially generated Brownian bridges  $\{U_B^i(r)\}_{i=1}^L$  with covariance matrix  $\hat{\Lambda}$ . For each  $i$ , consider the estimators:

$$\begin{aligned} \tilde{\beta}_{i,1} &= \sum_{m^*} \hat{\beta}(m^*) \mathbf{1}\{m^* = \arg\max_m [T\hat{\beta}(m)' \hat{\Sigma} \hat{\beta}(m) \\ &\quad - \frac{1}{\pi} [T^{1/2}\pi\hat{\beta}(m) + \hat{H}(m)U_B^i(\pi, m)]' \hat{\Sigma} [T^{1/2}\pi\hat{\beta}(m) + \hat{H}(m)U_B^i(\pi, m)] + \text{tr}(H(m)\hat{\Lambda})\log(\pi)]\} \end{aligned} \quad (3.14)$$

and

$$\begin{aligned} \tilde{\beta}_{i,2} &= \sum_{m^*} [\hat{\beta}(1, m^*) - T^{-1/2} \frac{U_b^i(\pi, m^*)}{1 - \pi}] \\ \mathbf{1}\{m^* = \arg \max_m [\frac{1}{\pi} [T^{1/2} \pi \hat{\beta}(m) + \hat{H}(m) U_B^i(\pi, m)]' \hat{\Sigma} [T^{1/2} \pi \hat{\beta}(m) + \hat{H}(m) U_B^i(\pi, m)] - 2n(m)]\} \end{aligned} \quad (3.15)$$

where  $\hat{H}(m)$  is the Moore-Penrose inverse of  $\hat{\Sigma}(m)$  and  $U_B^i(r, m)$  is the vector with the elements of  $U_B^i(r)$  in the locations indexed by  $m$  and zeros everywhere else. The next proposition gives their limiting distributions:

**Proposition 3.6** *Under the conditions for Proposition 3.5, for each  $i$ :*

$$\begin{aligned} T^{1/2} \tilde{\beta}_{i,1} &\rightarrow_d \sum_{m^*} H(m^*) \tilde{Y}_i(1, m^*) \mathbf{1}\{m^* = \arg \max_m [\tilde{Y}_i(1, m)' H(m) \tilde{Y}_i(1, m) \\ &\quad - \frac{1}{\pi} \tilde{Y}_i(\pi, m)' H(m) \tilde{Y}_i(\pi, m) + \text{tr}(H(m) \Lambda) \log(\pi)]\} \end{aligned} \quad (3.16)$$

and

$$T^{1/2} \tilde{\beta}_{i,2} \rightarrow_d \sum_{m^*} H(m^*) \frac{\tilde{Y}_i(1, m^*) - \tilde{Y}_i(\pi, m^*)}{1 - \pi} \mathbf{1}\{m^* = \arg \max_m [\frac{1}{\pi} \tilde{Y}_i(\pi, m)' H(m) \tilde{Y}_i(\pi, m) - 2n(m)]\} \quad (3.17)$$

where  $\tilde{Y}_i(r) = rY + U_B^i(r)$  and  $\tilde{Y}_i(r, m)$  is a  $k \times 1$  vector with the elements of  $\tilde{Y}_i(r)$  in the locations indexed by  $m$  and zeros elsewhere. These are the same distributions as in equations (3.11) and (3.12).

These estimators can then be averaged over  $i$ . After this step of averaging over different realizations of the Brownian bridge, the asymptotic distributions depend on  $Y$  alone and are asymptotically the expectations of the out-of-sample and split-sample estimators conditional on  $Y$ . Note that this Rao-Blackwellization (henceforth, RB) does applies only to the out-of-sample and split-sample estimators because it is only for these estimators that there is any ancillary noise process to eliminate.

In the special case of regression considered in Subsection 2.1, numerical calculations indicate that the limiting risk of the RB estimator is strictly lower than the original estimator for at least some values of  $b$ , implying that the out-of-sample and split-sample estimators are asymptotically inadmissible.

### 3.2 Linear Regression Model and Bagging

In the special case of the regression model with orthonormal regressors, considered in subsection 2.1, we have  $\Lambda = \Sigma = \sigma^{-2} I_k$ . In this model, all of the estimators depend crucially on the partial sum process  $T^{-1/2} \sum_{t=1}^{\lfloor Tr \rfloor} x_t y_t$  and it follows from Proposition 3.4 that:

$$T^{-1/2} \sigma^{-2} \sum_{t=1}^{\lfloor Tr \rfloor} x_t y_t \rightarrow_d Y(r) \quad (3.18)$$

and Proposition 3.5 will immediately apply. In this case, moreover, the results of Knight and Fu (2000) apply to the LASSO estimator—they show that for the LASSO estimator if  $\max_t T^{-1} x_t' x_t \rightarrow 0$  and  $T^{-1/2} \lambda_T \rightarrow \lambda_0 \geq 0$  then

$$T^{1/2} \hat{\beta} \rightarrow_d \arg \min_v v' \Sigma v - 2v' Y(1) + \lambda_0 \sum_{j=1}^k |v_j|. \quad (3.19)$$

In the linear regression model (subsection 2.1), we can also consider adding a *bagging* step to each of the procedures. Bagging, or bootstrap aggregation, was proposed by Breiman (1996) as a way to smooth predictive procedures. Bühlmann and Yu (2002) study the large-sample properties of bagging. The  $i^{\text{th}}$  bagging step resamples from the pairs  $\{(x_t, y_t), t = 1, \dots, T\}$  with replacement to form a pseudo-sample  $\{x_t^*(i), y_t^*(i), t = 1, \dots, T\}$ . The full model-selection and estimation procedure is then applied to the  $i^{\text{th}}$  bootstrap sample. This is repeated  $L$  times, and the  $L$  estimates are averaged to obtain the bagged estimate that can be used for forecasting. The following proposition provides a key result for obtaining the limiting distribution with bagging.

**Proposition 3.7** *Let  $\{x_t^*, y_t^*, t = 1, \dots, T\}$  denote a bootstrap sample and  $g(\cdot)$  be an uniformly integrable  $\mathbb{R}^k$ -valued functional. Then*

$$E^* \left[ g \left( \frac{1}{\sigma^2 \sqrt{T}} \sum_{t=1}^{\lfloor Tr \rfloor} x_t^* y_t^* \right) \right] \rightarrow E^* \left[ g \left( \frac{r}{\sigma^2} \sqrt{T} \hat{\beta} + V(r) \right) \right] \quad a.s.,$$

where  $E^*$  represents the expectation with respect to the bootstrap conditional distribution and  $V(r)$  is a  $k \times 1$  Brownian motion with covariance matrix  $\sigma^{-2} I$ .

Thus the limiting distribution of a single bootstrap draw for the partial sums process mimics the result in



Proposition 3.4, except that the Brownian bridge  $U_B(r)$  is replaced with a Brownian motion  $V(r)$ . Define  $Y^*(r) = rY + V(r)$ . Using Proposition 3.7, we can obtain asymptotic representations for the different procedures incorporating a bagging step in analogy with (3.7)-(3.12):

**Proposition 3.8** *In large samples, the distributions of the alternative parameter estimates including a bagging step are as follows:*

(i) *Using unrestricted MLE:*

$$T^{1/2} \tilde{\beta} \rightarrow_d \Sigma^{-1} Y(1) \quad (3.20)$$

(ii) *Using the positive-part James-Stein estimator:*

$$T^{1/2} \tilde{\beta} \rightarrow_d E^* \left[ \Sigma^{-1} Y^*(1) \max \left\{ 1 - \frac{k-2}{Y^*(1)' \Sigma^{-2} Y^*(1)}, 0 \right\} \right] \quad (3.21)$$

(iii) *Using the Laplace estimator:*

$$T^{1/2} \tilde{\beta} \rightarrow_d E^* [Y^*(1) \circ \tilde{h}(WY^*(1))] \quad (3.22)$$

(iv) *Selecting the model using the AIC:*

$$T^{1/2} \tilde{\beta} \rightarrow_d E^* \left[ \sum_{m^*} H(m^*) Y^*(1, m^*) \mathbf{1}\{m^* = \arg \max_m [Y^*(1, m)' H(m) Y^*(1, m) - 2n(m)]\} \right] \quad (3.23)$$

(v) *Selecting the model minimizing out-of-sample error:*

$$T^{1/2} \tilde{\beta} \rightarrow_d E^* \left[ \sum_{m^*} H(m^*) Y^*(1, m^*) \mathbf{1}\{m^* = \arg \max_m [Y^*(1, m)' H(m) Y^*(1, m) - \frac{1}{\pi} Y^*(\pi, m)' H(m) Y^*(\pi, m) + \text{tr}(H(m) \Lambda) \log(\pi)]\} \right] \quad (3.24)$$

(vi) *Using the split-sample method:*

$$T^{1/2} \tilde{\beta} \rightarrow_d E^* \left[ \sum_{m^*} H(m^*) \frac{Y^*(1, m^*) - Y^*(\pi, m^*)}{1 - \pi} \mathbf{1}\{m^* = \arg \max_m [\frac{1}{\pi} Y^*(\pi, m)' H(m) Y^*(\pi, m) - 2n(m)]\} \right] \quad (3.25)$$

where  $\sum_{m^*}$  denotes the summation over all the models in  $\mathcal{M}$  and  $Y^*(r, m)$  is a  $k \times 1$  vector with the elements of  $Y^*(r)$  in the locations indexed by  $m$  and zeros elsewhere.

In Appendix B, we also provide more concrete expressions for the in-sample and split-sample procedures, in their standard form, with RB, and with bagging, in the special case where  $k = 1$ .

Bagging and our proposed RB procedure are closely related. RB uses simulation to integrate out the Brownian bridge  $U_B(r)$ . Bagging has an asymptotic representation that has a Brownian motion instead of  $U_B(r)$ , and then integrates out that Brownian motion. But RB has certain advantages. Bagging applies only in the case of independent data, whereas our RB approach can be used in any setting where we have consistent estimators of  $\Lambda$  and  $\Sigma$ . Also RB does not require resampling the data and reestimating the model. This may make RB especially attractive when the data are dependent or in circumstances where model estimation is computationally costly.

Breiman (1996) gave a heuristic argument for why bagging weakly reduces mean square error, but in fact bagging can increase mean square error. The calculations of Bühlmann and Yu (2002) showed this for the case of estimation with AIC model selection. See also Andreas and Stuetzle (2000) and Friedman and Hall (2007). On the other hand RB does indeed weakly reduce the local asymptotic risk for any convex loss function.

Since RB involves integrating out  $U_B(r)$ , it does not affect any procedure that does not depend on this ancillary noise. Because  $U_B(1) = 0$  full sample procedures won't depend on this noise. In particular, RB does not affect in-sample model selection with AIC or the Laplace estimator. But bagging will affect the limiting distribution of all of the procedures that we consider, except for the unrestricted MLE. In the case of in-sample model selection with AIC, bagging can be thought of as replacing hard thresholding with soft thresholding (see Appendix B for more discussion).

## 4 Extensions

In this section, we consider two extensions of the basic framework of our analysis, in the context of the linear regression model.

## 4.1 Unmodeled Structural Change

A variant of our basic regression model specifies that  $y_t = \beta'_t x_t + u_t$  where  $T^{1/2} \beta_{[Tr]} = W(r)$ , where  $r$  may be either a stochastic or nonstochastic process. This allows various forms of structural breaks, and is similar to specifications used by Andrews (1993) and Elliott and Mueller (2014). For example, if  $\beta_t = T^{-1/2} b + T^{-1/2} \tilde{b} 1(t > [Ts])$ , then  $W(r) = b + \tilde{b} 1(r > s)$ . Or, if  $\beta_t = T^{-1} \sum_{s=1}^t \eta_s$  with Gaussian shocks, then  $W(r)$  is a Brownian motion. Proposition 4.1 gives the asymptotic distribution of the partial sum process  $T^{-1/2} \sigma^{-2} \sum_{t=1}^{[Tr]} x_t y_t$  in this variant of our basic model:

**Proposition 4.1** *As  $T \rightarrow \infty$ , the partial sum process*

$$T^{-1/2} \sigma^{-2} \sum_{t=1}^{[Tr]} x_t y_t \rightarrow_d Z(r) \quad (4.1)$$

where  $Z(r) \stackrel{d}{=} \Sigma \int_0^r W(s) ds + r \xi + U_B(r)$ ,  $\xi \sim N(0, \Sigma)$  and  $U_B(r)$  is an independent  $k$ -dimensional Brownian bridge with covariance matrix  $\Sigma = \sigma^{-2} I$ .

Suppose that the researcher ignores the possibility of structural change, and simply uses the available estimators for forecasting. The limiting distributions of the estimators will be as in Propositions 3.5, with  $Y(r)$  replaced by  $Z(r)$  everywhere. Alternatively, the researcher might be aware of the possibility of structural change, and might choose to select among models and estimate parameters using a rolling window. The estimators will then have limiting distributions that are simple extensions of those in Propositions 3.5 and 3.8. Other approaches for dealing with the possibility of parameter instability might be considered, but we leave this topic for future research.

Structural instability may be part of the motivation for considering out-of-sample forecasting methods. It is true that RB is no longer guaranteed to reduce risk in the above model with structural change. Nonetheless, in Monte-Carlo simulations documented in the web appendix, we find that RB does in practice improve the risk of out-of-sample and split-sample forecasting approaches.

## 4.2 Model Combination

It may also be appealing to combine forecasts made from multiple models, instead of selecting a single model (Bates and Granger (1969) and Timmermann (2006)). Recalling that  $\hat{\beta}(1, m)$  denotes the parameter estimate from the model containing the variables indexed by  $m$  (with zeros in other locations), then we could estimate the parameter vector as  $\sum_m w(m) \hat{\beta}(1, m)$ , where  $\sum_m$  denotes the sum over all the models in  $\mathcal{M}$  and the weights sum to 1. As examples of weighting schemes, we could set the weight for model  $m$  to  $w(m) = \frac{\exp(AIC(m)/2)}{\sum_{m^*} \exp(AIC(m^*)/2)}$  (Buckland, Burnham, and Augustin, 1997) where  $AIC(m)$  denotes the Akaike Information Criterion in the model indexed by  $m$ , or we could weight models by out-of-sample predictive performance setting:

$$w(m) = \frac{\exp[\sum_{t=[T\pi]+1}^T \ell_t(\hat{\beta}_{1,t-1}(m))]}{\sum_{m^*} \exp[\sum_{t=[T\pi]+1}^T \ell_t(\hat{\beta}_{1,t-1}(m^*))]} \quad (4.2)$$

Alternatively, to do a combination version of the split-sample scheme, we could estimate the parameter vector as  $\sum_m w(m) \hat{\beta}^*(\pi, m)$  where  $w(m) = \frac{\exp(AIC(\pi, m)/2)}{\sum_m \exp(AIC(\pi, m)/2)}$  and  $AIC(\pi, m)$  denotes the Akaike Information Criterion for the model indexed by  $m$  computed only over the first fraction  $\pi$  of the sample.

**Proposition 4.2** *If the parameter vector is estimated by  $\sum_m w(m) \hat{\beta}(1, m)$  then in large samples, the distributions of the alternative parameter estimates will be:*

$$E\{\sum_m w(m) H(m) Y(1, m)\} \quad (4.3)$$

where

$$w(m) \propto \exp([Y(1, m)' H(m) Y(1, m) - 2n(m)]/2) \quad (4.4)$$

or

$$w(m) \propto \exp(Y(1, m)' H(m) Y(1, m) - \frac{1}{\pi} Y(1, m)' H(m) Y(1, m) + \text{tr}(H(m) \Lambda) \log(\pi)) \quad (4.5)$$

for exponential AIC and out-of-sample prediction error weights, respectively. Meanwhile, if the parameter vector is instead estimated by  $\sum_m w(m) \hat{\beta}^*(\pi, m)$  with exponential AIC weights, then in large samples, the

distribution of the estimator will be:

$$E\{\Sigma_{m^*} w(m) H(m) \frac{Y(1, m) - Y(\pi, m)}{1 - \pi}\} \quad (4.6)$$

where

$$w(m) \propto \exp\left(\left[\frac{1}{\pi} Y(\pi, m)' H(m) Y(\pi, m) - 2n(m)\right]/2\right) \quad (4.7)$$

The standard bagging step can be added to any of these methods for forecast combination and the resulting limiting distributions are also given by Proposition 4.2, except with  $Y(\cdot)$  and  $Y(\cdot, m)$  replaced by  $Y^*(\cdot)$  and  $Y^*(\cdot, m)$  everywhere. Or RB can be added, to integrate out  $U_B(r)$ .

An alternative and more standard way to obtain combination weights for the out-of-sample forecasting scheme would be to weight the forecasts by the inverse mean square error (Bates and Granger (1969) and Timmermann (2006)). Under our local asymptotics, this will give each model equal weight in large samples.

## 5 Numerical Work

In this section we numerically explore the root mean squared error

$$\sqrt{E[(T^{1/2}\tilde{\beta} - b)'(T^{1/2}\tilde{\beta} - b)]}, \quad (5.1)$$

the square of which is asymptotically equivalent to the NMSPE in the regression model example. Given the expressions in Propositions 3.5 and 3.8, we can simulate the asymptotic risk of different methods in their standard form, with RB, and with bagging<sup>4</sup> for different choices of the localization parameter  $b$  and the number of potential predictors  $k$ . None of the methods gives the lowest risk uniformly in  $b$ . Always using the big model is minmax, but due to the Stein phenomenon, it may be dominated by shrinkage estimators. In all cases, RB and bagging are implemented using 100 replications, the out-of-sample and split-sample methods both set  $\pi = 0.5$ , and we set  $\Sigma = \Lambda = I_k$ . The asymptotic risk is symmetric in  $b$  and is consequently shown only for non-negative  $b$ . The bagging results from Proposition 3.8 apply only in the

---

<sup>4</sup>The results with bagging are based on Proposition 3.8, which applies only in the case of the linear regression model.

special case of the linear regression model, but RB applies in the general pseudo-likelihood framework. Figure 1 plots the asymptotic risk of the standard MLE, JS, in-sample, out-of-sample and split-sample methods, for the case  $k = 1$  against  $b$ . Results with RB and bagging are also included. Results with LASSO and Laplace estimators are not shown, but are in the web appendix.

Among the standard methods, selecting the model in-sample by AIC does better than the out-of-sample scheme for most values of  $b$ , which in turn dominates the split-sample method. But RB changes this ordering. RB reduces the risk of the out-of-sample and split-sample methods for all values of  $b$ , and makes them much more competitive. Bagging accomplishes much the same thing. The fact that bagging improves the out-of-sample and split-sample methods uniformly in  $b$  is just a numerical result, but it is also a theoretical result for RB. Neither bagging nor RB dominates the other in terms of risk. Bagging also helps with the in-sample method for some but not all values of  $b$ —this was also shown by Bühlmann and Yu (2002). Recall that RB does nothing to the in-sample method.

Among all the prediction methods represented in Figure 1, which one the researcher would ultimately would want to use depends on  $b$ , which is in turn not consistently estimable. But the split-sample and out-of-sample methods do best for many values of the localization parameter, as long as the bagging or RB step is included. Indeed, for all  $b$ , the best forecast is some method combined with bagging or RB.

We next consider the case where the number of potential predictors  $k$  is larger, but only one parameter actually takes on a nonzero value. (Of course, the researcher does not know this.) Without loss of generality, we let the nonzero element of  $b$  be the first element and so specify that  $b = (b_1, 0, \dots, 0)'$ . To keep the model space manageable, the set of possible models  $\mathcal{M}$  consists of models with the first  $l$  predictors, for  $l \in \{0, 1, \dots, k\}$ . Figures 2 and 3 plot the risk for  $k = 3$  and  $k = 20$  against  $b_1$  for in-sample, out-of-sample and split-sample methods in the standard form, with RB, and with bagging. The positive-part James-Stein estimator is also included. The split-sample method with either RB or bagging compares very favorably with the other alternatives. Other numerical results with multiple predictors are contained in the web appendix, and give similar conclusions.

Figure 4 plots the risk for  $k = 1$  against  $b$  for the in-sample, out-of-sample and split-sample forecast combination methods, in their standard form, with RB and with bagging. These are based on simulating the distributions in Proposition 4.2. The combination forecasts are generally better than forecasts based

on selecting an individual model. Nonetheless, with combined forecasts as with individual forecasts, in the absence of a randomization step, using in-sample AIC weights does best for most values of  $b$ . Adding in RB/bagging allows better predictions to be made. RB/bagging reduces the risk of the combination forecasts with out-of-sample or split-sample weights uniformly in  $b$ . Once an RB/bagging step is added in, there is no clear winner among the in-sample, out-of-sample and split-sample forecast combination methods.

Although our numerical work considers the quadratic loss function, it should be noted that the weak reduction in risk from RB applies with any convex loss function.

## 6 Monte Carlo Simulations

The results in the previous section are based on a local asymptotic sequence. The motivation for this is to provide a good approximation to the finite sample properties of different forecasting methods while retaining some assurance that they are not an artifact of a specific simulation design. As some check that the local asymptotics are indeed relevant to small samples, we did a small simulation consisting of equation (2.13) with  $t(5)$  errors<sup>5</sup> scaled to have unit variance, independent standard normal regressors, a sample size  $T = 100$ , and different values of  $k$ . In each simulation we drew  $T + 1$  observations on  $y_t$  and  $x_t$ , used the first  $T$  for model selection and parameter estimation according to one of the methods discussed above. Then given  $x_{T+1}$ , we worked out the prediction for  $y_{T+1}$ , and computed the mean square prediction error (MSPE).

Figure 5 plots the simulated root normalized mean square prediction errors ( $\sqrt{T \cdot (MSPE - 1)}$ ) against  $\beta$  for  $k = 1$  for the MLE, JS, in-sample, out-of-sample and split-sample methods, in their standard form, with RB and with bagging. Our simulations included results using leave-one-out cross-validation, the Laplace and LASSO estimators. These are omitted from Figure 5, but shown in the web appendix. Not surprisingly, leave-one-out cross-validation gave very similar results to AIC. The web appendix also contains results for higher values of  $k$ .

The simulations in Figure 5 and in the web appendix give very similar conclusions to the local asymptotic calculations (although MSPEs for all methods tend to be a bit higher than would be predicted by the local

---

<sup>5</sup>Results with normal errors were very similar.

asymptotics for  $k = 20$ , because of estimation of variance-covariance matrices). Without RB or bagging, the in-sample scheme generally gives the best forecasts, followed by out-of-sample, with the split-sample doing the worst. RB or bagging substantially improve the performance of the out-of-sample and split sample methods.

The web appendix also contains some simulations in which the parameters follow a random walk, but the econometrician treats them as though they were fixed. Although there is no theoretical result that RB has to reduce risk in this case, we find that risk is in fact lowered by RB as long as the parameter variation is not too great.

## 7 Application

We finally consider a small empirical illustration to the classic problem of forecasting stock returns. We follow the setup of Goyal and Welch (2008) in forecasting annual excess stock returns using data from 1926 to 2014 with 13 possible predictors: book-to-market ratio, Treasury bill yields, long-term yields, net equity expansion, inflation, percent equity issuing, long term returns, stock variance, default yield spread, default return spread, dividend-price ratio, dividend yield and earnings-price ratio. These are all the predictors considered by Goyal and Welch (2008) for which data are available over the full sample period, excepting some that are perfectly multicollinear with the included predictors. We then consider one-year-ahead forecasting of stock returns using all possible subsets of these 13 predictors including the empty set, for a total of 8,192 models. Definitions of the predictors and sources are in Goyal and Welch (2008). Each model includes an intercept.

AIC chooses a single predictor: the book-to-market ratio. The same is true selecting the model with the split-sample approach. Selecting the model by out-of-sample performance, percent equity issuing and the dividend price ratio are the two chosen predictors. Table 1 shows the probability of the chosen model having no predictors whatsoever (other than the intercept), using AIC, out-of-sample and split-sample schemes with bagging and out-of-sample and split-sample schemes with RB. Having no predictors other than the intercept amounts to forecasting excess stock returns using the unconditional mean. Table 1 also shows the expected number of predictors (other than the intercept) for all five simulation schemes. AIC with bagging, AIC with RB and the split-sample approach with RB all put most weight on models



with 0 or 1 predictors, and thereby implying heavy shrinkage. The out-of-sample approach with either bagging or RB chooses larger models. Table 1 also reports the fraction of draws for which each of the 13 possible predictors is included in the model, for each of the five simulation schemes. Finally, Table 1 gives the forecast for excess stock returns in 2015 using each of the methods. These are very close to the unconditional mean, indicating a considerable degree of shrinkage.

## 8 Conclusion

When forecasting using  $k$  potential predictors, each of which has a coefficient that is local to zero, there are several competing methods, none of which is most accurate uniformly in the localization parameter, which is in turn not consistently estimable. Optimizing the in-sample fit, as measured by the Akaike information criterion, generally does better than out-of-sample or split-sample methods. However, the out-of-sample and split-sample methods can be improved substantially by removing the impact of an ancillary noise term that appears in their limit representations, either through Rao-Blackwellization or bagging. Rao-Blackwellization uniformly lowers asymptotic risk of the out-of-sample and split-sample methods and can be implemented without having to resample the data. For important ranges of the local parameters, these modified procedures are very competitive with in-sample methods.

## A Appendix: Proof of Propositions

**Proof of Proposition 3.4:** We have

$$\begin{aligned} T^{-1/2} \sum_{t=1}^{[Tr]} \frac{\partial \ell_t(0)}{\partial \beta} &= T^{-1/2} \sum_{t=1}^{[Tr]} \frac{\partial \ell_t(\beta_{0,T})}{\partial \beta} - T^{-1/2} \sum_{t=1}^{[Tr]} \frac{\partial^2 \ell_t(\beta_{0,T})}{\partial \beta \partial \beta'} \beta_{0,T} + o_p(1) \\ &\rightarrow_d r \Sigma b + B(r), \end{aligned}$$

where  $B(r)$  denotes a Brownian motion with covariance matrix  $\Lambda$ . Let  $Y(r) = r \Sigma b + B(r)$ , and let  $Y = Y(1)$  which is  $N(\Sigma b, \Lambda)$ . Define  $U_B(r) = B(r) - rB(1)$ . Then  $U_B$  is a Brownian bridge by standard arguments, and by calculating the covariance between  $U_B(r)$  and  $B(1)$  it can be verified that  $U_B$  is independent of  $B(1)$ . We can therefore write

$$\begin{aligned} Y(r) &= r \Sigma b + B(r) \\ &= r \Sigma b + rB(1) + U_B(r) \\ &= rY + U_B(r). \end{aligned}$$

and  $Y = \Sigma b + B(1)$  is uncorrelated with  $U_B$  and hence independent of  $U_B$ . ■

**Proof of Proposition 3.5:** Let  $\hat{\beta}$  denote the unrestricted estimator and let  $\hat{\beta}(m)$  denote the restricted estimators as defined in Section 2. Equations (3.7) and (3.8) immediately follow because  $T^{1/2} \hat{\beta} \rightarrow_d Y$ .

**[note: I think we need an assumption either of identification or directly on limits of  $\hat{\beta}(m)$ ]**

The AIC objective function (to be maximized) is:

$$2 \sum_{t=1}^T l_t(\hat{\beta}(m)) - 2n(m) = 2 \sum_{t=1}^T l_t(0) + 2\hat{\beta}(m)' \sum_{t=1}^T \frac{\partial l_t(0)}{\partial \beta} + \hat{\beta}(m)' \sum_{t=1}^T \frac{\partial^2 l_t(0)}{\partial \beta \partial \beta'} \hat{\beta}(m) - 2n(m) + o_p(1),$$

which is asymptotically the same, up to the same affine transformation across all models, as

$$Y(1, m)' H(m) Y(1, m) - 2n(m),$$

noting that  $H(m) \Sigma H(m) = H(m)$ .

The OOS objective function (to be maximized) is:

$$\sum_{t=[T\pi]+1}^T l_t(\hat{\beta}_{1,t-1}(m)) = \sum_{t=[T\pi]+1}^T \left[ l_t(0) + \hat{\beta}_{1,t-1}(m)' \frac{\partial l_t(0)}{\partial \beta} + \frac{1}{2} \hat{\beta}_{1,t-1}(m)' \frac{\partial^2 l_t(0)}{\partial \beta \partial \beta'} \hat{\beta}_{1,t-1}(m) \right] + o_p(1).$$

Let  $S(m)$  be the diagonal matrix with ones in the diagonal elements indexed by  $m$  and zeros elsewhere.

Now  $T^{1/2} \hat{\beta}_{1,[Tr]}(m) \rightarrow_d H(m) \frac{Y(r)}{r}$  and

$$\begin{aligned} \sum_{t=[T\pi]+1}^T \hat{\beta}_{1,t-1}(m)' \frac{\partial l_t(0)}{\partial \beta} &= \sum_{t=[T\pi]+1}^T \sum_{s=1}^{t-1} \frac{\partial l_s(\beta_{0,T})'}{\partial \beta} \frac{1}{t-1} H(m) \frac{\partial l_t(\beta_{0,T})}{\partial \beta} \\ &\quad + b' \Sigma H(m) T^{-1/2} \sum_{t=[T\pi]+1}^T \frac{\partial l_t(\beta_{0,T})}{\partial \beta} + \sum_{t=[T\pi]+1}^T \hat{\beta}_{1,t-1}(m)' \frac{1}{T} \Sigma b + o_p(1) \\ &\rightarrow_d \int_{\pi}^1 \frac{B(r)'}{r} H(m) dB(r) + b' \Sigma H(m) (B(1) - B(\pi)) + \int_{\pi}^1 \frac{Y(r)'}{r} H(m) \Sigma b dr \\ &= \int_{\pi}^1 \frac{Y(r)'}{r} H(m) dB(r) + \int_{\pi}^1 \frac{Y(r)'}{r} H(m) \Sigma b dr \\ &= \int_{\pi}^1 \frac{Y(r)'}{r} H(m) dY(r). \end{aligned}$$

Consequently, the OOS objective function is asymptotically the same, up to the same affine transformation across all models, as

$$2 \int_{\pi}^1 \frac{Y(r)'}{r} H(m) dY(r) - \int_{\pi}^1 \frac{Y(r)'}{r} H(m) \frac{Y(r)}{r} dr.$$

Using an argument similar to Hansen and Timmermann (2013), define

$$F(Y, r) = \frac{1}{r} Y(r)' H(m) Y(r) - \text{tr}(H(m) \Lambda) \log(r).$$

By Ito's lemma

$$\begin{aligned} dF(Y, r) &= \left\{ -\frac{1}{r^2} Y(r)' H(m) Y(r) - \frac{1}{r} \text{tr}(H(m) \Lambda) + 2b' \Sigma H(m) Y(r) + \frac{1}{r} \text{tr}(H(m) \Lambda) \right\} dr + \frac{2}{r} Y(r)' H(m) dB(r) \\ &= \left\{ -\frac{1}{r^2} Y(r)' H(m) Y(r) + 2b' \Sigma H(m) Y(r) \right\} dr + \frac{2}{r} Y(r)' H(m) dB(r) \\ &= \left\{ -\frac{1}{r^2} Y(r)' H(m) Y(r) + 2b' \Sigma H(m) Y(r) \right\} dr + \frac{2}{r} Y(r)' H(m) dY(r) - 2Y(r)' H(m) \Sigma b dr \\ &= \frac{2}{r} Y(r)' H(m) dY(r) - \frac{1}{r^2} Y(r)' H(m) Y(r) dr. \end{aligned}$$

So

$$\begin{aligned} F(Y, 1) - F(Y, \pi) &= \int_{\pi}^1 dF(Y, r) = 2 \int_{\pi}^1 \frac{Y(r)'}{r} H(m) dY(r) - \int_{\pi}^1 \frac{Y(r)'}{r} H(m) \frac{Y(r)}{r} dr \\ &= Y(1)' H(m) Y(1) - \frac{1}{\pi} Y(\pi)' H(m) Y(\pi) + tr(H(m) \Lambda) \log(\pi). \end{aligned}$$

Consequently, the OOS objective function is asymptotically the same, up to the same affine transformation across all models, as

$$Y(1)' H(m) Y(1) - \frac{1}{\pi} Y(\pi)' H(m) Y(\pi) + tr(H(m) \Lambda) \log(\pi).$$

The AIC estimated over the first fraction  $\pi$  of the sample is:

$$\begin{aligned} 2 \sum_{t=1}^{[T\pi]} l_t(\hat{\beta}_{1, [T\pi]}(m)) - 2n(m) &= 2 \sum_{t=1}^{[T\pi]} l_t(0) + 2\hat{\beta}_{1, [T\pi]}(m)' \sum_{t=1}^{[T\pi]} \frac{\partial l_t(y_t, 0)}{\partial \beta} \\ &\quad + \hat{\beta}_{1, [T\pi]}(m)' \sum_{t=1}^{[T\pi]} \frac{\partial^2 l_t(y_t, 0)}{\partial \beta \partial \beta'} \hat{\beta}_{1, [T\pi]}(m) - 2n(m) + o_p(1) \end{aligned}$$

which is asymptotically the same, up to the same affine transformation across all models, as

$$\frac{1}{\pi} Y(\pi, m)' H(m) Y(\pi, m) - 2n(m)$$

The limiting distributions in Proposition 3.5 all follow from these results and the facts that  $T^{1/2} \hat{\beta} \rightarrow_d \Sigma^{-1} Y(1)$  and  $T^{1/2} \hat{\beta}(m) \rightarrow_d H(m) Y(1, m)$ . ■

**Proof of Proposition 3.6:** We know that  $T^{1/2} \hat{\beta}(m) \rightarrow_d H(m) \tilde{Y}_i(1, m)$  and  $\hat{H}(m) \rightarrow_p H(m)$ . Hence  $T^{1/2} \hat{\beta}(m) + \hat{H}(m) \frac{U_B^i(r, m)}{r} \rightarrow_d H(m) \frac{\tilde{Y}_i(r, m)}{r}$ . The result follows immediately. ■

**Proof of Proposition 3.7:** Let  $\{x_t^*, y_t^*\}$  denote a bootstrap sample and let  $u_t^* = y_t^* - \beta_T' x_t^*$ ,  $t = 1, \dots, T$ .  $\hat{\beta}$  is the OLS estimate using the original data. Using the assumption that the  $x_t$  are orthonormal, we can take  $\hat{\beta} = T^{-1} \sum_{t=1}^T x_t y_t$ . Then, in the bootstrap conditional probability space (conditioning on the realization

of  $\{x_t, y_t\}$ ,

$$\begin{aligned}
\frac{1}{\sigma^2 \sqrt{T}} \sum_{t=1}^{[Tr]} x_t^* y_t^* - \frac{r}{\sigma^2} \sqrt{T} \hat{\beta} &= \frac{1}{\sigma^2 \sqrt{T}} \sum_{t=1}^{[Tr]} x_t^* y_t^* - \frac{r}{\sigma^2 \sqrt{T}} \sum_{t=1}^T x_t y_t \\
&= \frac{1}{\sigma^2 \sqrt{T}} \sum_{t=1}^{[Tr]} x_t^* (x_t^{*'} \beta_T + u_t^*) - \frac{r}{\sigma^2 \sqrt{T}} \sum_{t=1}^T x_t (x_t' \beta_T + u_t) \\
&= \left( \frac{1}{\sigma^2 \sqrt{T}} \sum_{t=1}^{[Tr]} \left[ x_t^* x_t^{*'} - \frac{1}{T} \sum_{s=1}^T x_s x_s' \right] \right) \beta_T + \frac{1}{\sigma^2 \sqrt{T}} \sum_{t=1}^{[Tr]} \left( x_t^* u_t^* - \frac{1}{T} \sum_{s=1}^T x_s u_s \right) \\
&= \left( \frac{1}{\sigma^2 T} \sum_{t=1}^{[Tr]} \left[ x_t^* x_t^{*'} - \frac{1}{T} \sum_{s=1}^T x_s x_s' \right] \right) b + \frac{1}{\sigma^2 \sqrt{T}} \sum_{t=1}^{[Tr]} \left( x_t^* u_t^* - \frac{1}{T} \sum_{s=1}^T x_s u_s \right).
\end{aligned}$$

The first term is  $o_{p^*}(1)$  by a Uniform Law of Large Numbers for partial sums (e.g., Gaenssler and Ziegler (1994)), and the second converges almost surely to  $V(r)$ , from Theorem 2.2 of Park (2002). Therefore

$$\frac{1}{\sigma^2 \sqrt{T}} \sum_{t=1}^{[Tr]} x_t^* y_t^* = \frac{r}{\sigma^2} \sqrt{T} \hat{\beta} + V(r) + o_p(1) \quad a.s.,$$

Given that  $g$  is uniformly integrable,

$$E^* \left[ g \left( \frac{1}{\sigma^2 \sqrt{T}} \sum_{t=1}^{[Tr]} x_t^* y_t^* \right) \right] \rightarrow E \left[ g \left( \frac{r}{\sigma^2} \sqrt{T} \hat{\beta} + V(r) \right) \right] \quad a.s.,$$

as required. ■

The proofs of Propositions 3.8 and 4.2 involve exactly the same calculations as in Proposition 3.5 and are hence omitted.

**Proof of Proposition 4.1:** We have

$$\begin{aligned}
T^{-1/2} \sigma^{-2} \sum_{t=1}^{[Tr]} x_t y_t &= T^{-1/2} \sigma^{-2} \sum_{t=1}^{[Tr]} x_t x_t' \beta_t + T^{-1/2} \sigma^{-2} \sum_{t=1}^{[Tr]} x_t u_t \\
&= T^{-3/2} \sigma^{-2} \sum_{t=1}^{[Tr]} x_t x_t' \sum_{s=1}^t \eta_s + T^{-1/2} \sigma^{-2} \sum_{t=1}^{[Tr]} x_t u_t \\
&\rightarrow_d \sigma_\eta \Sigma \int_0^r W(s) ds + B(r) = \sigma_\eta \Sigma \int_0^r W(s) ds + r \xi + U_B(r).
\end{aligned}$$

where  $B(r)$  is a Brownian motion with covariance matrix  $\Lambda$ . ■

## B Appendix: Shrinkage Representations in the case $k = 1$

In the case  $k = 1$ , and with  $\Sigma = \Lambda$ , some of the expressions in Propositions 3.5 and 3.8 can be simplified.

For the AIC estimator in its standard form we have:

$$T^{1/2} \tilde{\beta} \rightarrow_d \Sigma^{-1} Y 1(|Y| > \sqrt{2\Sigma}). \quad (\text{B1})$$

Rao-Blackwellization makes no difference to the AIC estimator, and equation (B1) continues to apply.

For the AIC estimator, with bagging, we have:

$$T^{1/2} \tilde{\beta} \rightarrow_d \Sigma^{-1} \{Y - Y\Phi(\sqrt{2} - \kappa Y) + \kappa\phi(\sqrt{2} - \kappa Y) + Y\Phi(-\sqrt{2} - \kappa Y) - \kappa\phi(-\sqrt{2} - \kappa Y)\},$$

where  $\kappa = \Sigma^{-1/2}$ , shown in proposition 2.2 of Bühlmann and Yu (2002).<sup>6</sup> Comparing this to equation (B1), in the context of the AIC estimator, bagging is effectively replacing a hard thresholding procedure with a soft thresholding counterpart.

For the split-sample estimator, we have:

$$T^{1/2} \tilde{\beta} \rightarrow_d \Sigma^{-1} z_1 1(|z_2| > \sqrt{\frac{2\Sigma}{\pi}}),$$

where  $z_1 = Y - \frac{U_B(\pi)}{1-\pi}$  and  $z_2 = Y + \frac{U_B(\pi)}{\pi}$ . By direct calculations,  $z_1$  is  $N(\Sigma b, \frac{1}{1-\pi}\Sigma)$ ,  $z_2$  is  $N(\Sigma b, \frac{1}{\pi}\Sigma)$  and  $z_1$  and  $z_2$  are mutually independent.

With RB, for the split-sample estimator in the  $i$ th simulated sample, we have:

$$T^{1/2} \tilde{\beta}_i \rightarrow_d \Sigma^{-1} (Y - \sqrt{\frac{\pi\Sigma}{1-\pi}} z(i)) 1(|Y + \sqrt{\frac{(1-\pi)\Sigma}{\pi}} z(i)| > \sqrt{\frac{2\Sigma}{\pi}}) = (\Sigma^{-1} Y - \gamma z(i)) 1(|\gamma Y + z(i)| > \sqrt{\frac{2}{1-\pi}}),$$

where  $z(i)$  is  $N(0, 1)$ , and is independent of  $Y$  and  $\gamma = \sqrt{\frac{\pi}{(1-\pi)\Sigma}}$ . Thus for the overall split-sample estima-

---

<sup>6</sup>Indeed, given the orthonormal setting, even if  $k > 1$ , if we sort the coefficient estimates by their absolute magnitude and apply AIC sequentially to these models, dropping variables one at a time as long as called for by the information criterion, then the above two expressions will apply to each element of  $\tilde{\beta} - \beta$  (Bühlmann and Yu, 2002; Stock and Watson, 2012). But the use of the AIC that we are considering in this paper is to select among all  $2^k$  possible models and so no such simplification is available in this case.

tor with RB, we have:

$$T^{1/2} \tilde{\beta} \rightarrow_d \Sigma^{-1} \{Y - Y\Phi(\sqrt{\frac{2}{1-\pi}} - \gamma Y) - \gamma\phi(\sqrt{\frac{2}{1-\pi}} - \gamma Y) + Y\Phi(-\sqrt{\frac{2}{1-\pi}} - \gamma Y) + \gamma\phi(-\sqrt{\frac{2}{1-\pi}} - \gamma Y)\}.$$

Meanwhile, for bagging the split-sample estimator in the  $i$ th bagging sample, we have:

$$T^{1/2} \tilde{\beta}_i \rightarrow_d z_1(i) 1(|z_2(i)| > \sqrt{\frac{2\Sigma}{\pi}})$$

where  $z_1(i) = Y + \frac{V_i(1) - V_i(\pi)}{1-\pi}$  and  $z_2(i) = Y + \frac{V_i(\pi)}{\pi}$ . By direct calculations,  $z_1(i)|Y$  is  $N(Y, \frac{\Sigma}{(1-\pi)})$ ,  $z_2(i)|Y$  is  $N(Y, \frac{\Sigma}{\pi})$  and the two are independent, conditional on  $Y$ . Thus for the overall split-sample with bagging estimator:

$$T^{1/2} \tilde{\beta} \rightarrow_d \Sigma^{-1} \{Y - Y\Phi(\sqrt{2} - \sqrt{\frac{\pi}{\Sigma}} Y) + Y\Phi(-\sqrt{2} - \sqrt{\frac{\pi}{\Sigma}} Y)\}.$$

For the out-of-sample estimator, we have:

$$T^{1/2} \tilde{\beta} \rightarrow_d \Sigma^{-1} Y 1(|z_3| < \sqrt{\max(Y^2 \pi + \pi \log(\pi) \Sigma, 0)}),$$

where  $z_3 = \pi Y + U_B(\pi)$  is  $N(\pi \Sigma b, \pi \Sigma)$

With RB, for the out-of-sample estimator, in the  $i$ th simulated sample, we have:

$$T^{1/2} \tilde{\beta}_i \rightarrow_d \Sigma^{-1} Y 1(|z(i)| < \sqrt{\max(\frac{Y^2}{\pi \Sigma^2} + \frac{\log(\pi)}{\pi \Sigma}, 0) - b}),$$

Thus for the overall out-of-sample estimator with RB:

$$T^{1/2} \tilde{\beta} \rightarrow_d \Sigma^{-1} Y \{\Phi(\sqrt{\max(\frac{Y^2}{\pi \Sigma^2} + \frac{\log(\pi)}{\pi \Sigma}, 0) - b}) - \Phi(-\sqrt{\max(\frac{Y^2}{\pi \Sigma^2} + \frac{\log(\pi)}{\pi \Sigma}, 0) - b})\}.$$

Lastly, for bagging the out-of-sample estimator in the  $i$ th bagging sample, we have:

$$T^{1/2} \tilde{\beta}_i \rightarrow_d \Sigma^{-1} (Y + V_i(1)) 1((Y + V_i(1))^2 - (\pi Y + V_i(\pi))^2 + \log(\pi) \Sigma > 0).$$

In this last case, we have no closed form expression for the overall out-of-sample estimator with bagging.

In all cases, we can write the limit of  $T^{1/2}\tilde{\beta}$  as  $Yg(Y)$  and we think of  $g(Y)$  as the implied local asymptotic shrinkage function. The web appendix plots the implied  $g(Y)$  functions for the AIC, out-of-sample and split sample estimators, in their standard form, and with RB and bagging.



## References

- AKAIKE, H. (1974): "A New Look at the Statistical Model Identification," *IEEE Transactions on Automatic Control*, 19, 716–723.
- ANDREAS, B., AND W. STUETZLE (2000): "Bagging does not always Decrease Mean Squared Error," mimeo.
- ANDREWS, D. W. K. (1993): "Tests for Parameter Instability and Structural Change with Unknown Change Point," *Econometrica*, 61(4), 821–856.
- ASHLEY, R., C. W. GRANGER, AND R. SCHMALENSEE (1980): "Advertising and Aggregate Consumption: An Analysis of Causality," *Econometrica*, 48, 1149–1167.
- BATES, J. M., AND C. W. GRANGER (1969): "The combination of forecasts," *Operations Research Quarterly*, 20, 451–468.
- BREIMAN, L. (1996): "Bagging Predictors," *Machine Learning*, 36, 105–139.
- BUCKLAND, S. T., K. P. BURNHAM, AND N. H. AUGUSTIN (1997): "Model Selection: An Integral Part of Inference," *Biometrics*, 53, 603–618.
- BÜHLMANN, P., AND B. YU (2002): "Analyzing Bagging," *Annals of Statistics*, 30, 927–961.
- CLAESKENS, G., AND N. L. HJORT (2008): *Model Selection and Model Averaging*. Cambridge University Press, Cambridge.
- CLARK, T. E. (2004): "Can Out-of-Sample Forecast Comparisons help Prevent Overfitting?," *Journal of Forecasting*, 23, 115–139.
- DIEBOLD, F. X., AND R. S. MARIANO (1995): "Comparing Predictive Accuracy," *Journal of Business and Economic Statistics*, 13, 253–263.
- EFRON, B. (2014): "Estimation and Accuracy After Model Selection (with discussion)," *Journal of the American Statistical Association*, 109, 991–1007.
- ELLIOTT, G., AND U. K. MUELLER (2014): "Pre and Post Break Parameter Inference," *Journal of Econometrics*, 180(2), 141–157.

- FRIEDMAN, J. H., AND P. HALL (2007): “On Bagging and Nonlinear Estimation,” *Journal of Statistical Planning and Inference*, 137, 669–683.
- GAENSSLER, P., AND K. ZIEGLER (1994): “A Uniform Law of Large Numbers for Set-Index Processes with Applications to Empirical and Partial-Sum Processes,” in *Probability in Banach Spaces*, 9, ed. by J. Hoffmann-Jorgensen, J. Kuelbs, and M. B. Marcus. Springer.
- GIACOMINI, R., AND H. WHITE (2006): “Tests of Conditional Predictive Ability,” *Econometrica*, 74, 1545–1578.
- GOYAL, A., AND I. WELCH (2008): “A Comprehensive Look at The Empirical Performance of Equity Premium Prediction,” *Review of Financial Studies*, 21, 1455–1508.
- HANSEN, P. R. (2009): “In-Sample Fit and Out-of-Sample Fit: Their Joint Distribution and its Implications for Model Selection,” mimeo.
- HANSEN, P. R., AND A. TIMMERMAN (2013): “Equivalence Between Out-of-Sample Forecast Comparisons and Wald Statistics,” mimeo.
- INOUE, A., AND L. KILIAN (2004): “In-Sample or Out-of-Sample Tests of Predictability: Which One Should We Use?,” *Econometric Reviews*, 23, 371–402.
- (2006): “On the Selection of Forecasting Models,” *Journal of Econometrics*, 130, 273–306.
- KNIGHT, K., AND W. FU (2000): “Asymptotics for LASSO-Type Estimators,” *Annals of Statistics*, 28, 1356–1378.
- LEEB, H., AND B. M. PÖTSCHER (2005): “Model selection and inference: Facts and Fiction,” *Econometric Theory*, 21, 21–59.
- 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.
- MALLOWS, C. L. (1973): “Some Comments on  $C_p$ ,” *Technometrics*, 15, 661–675.

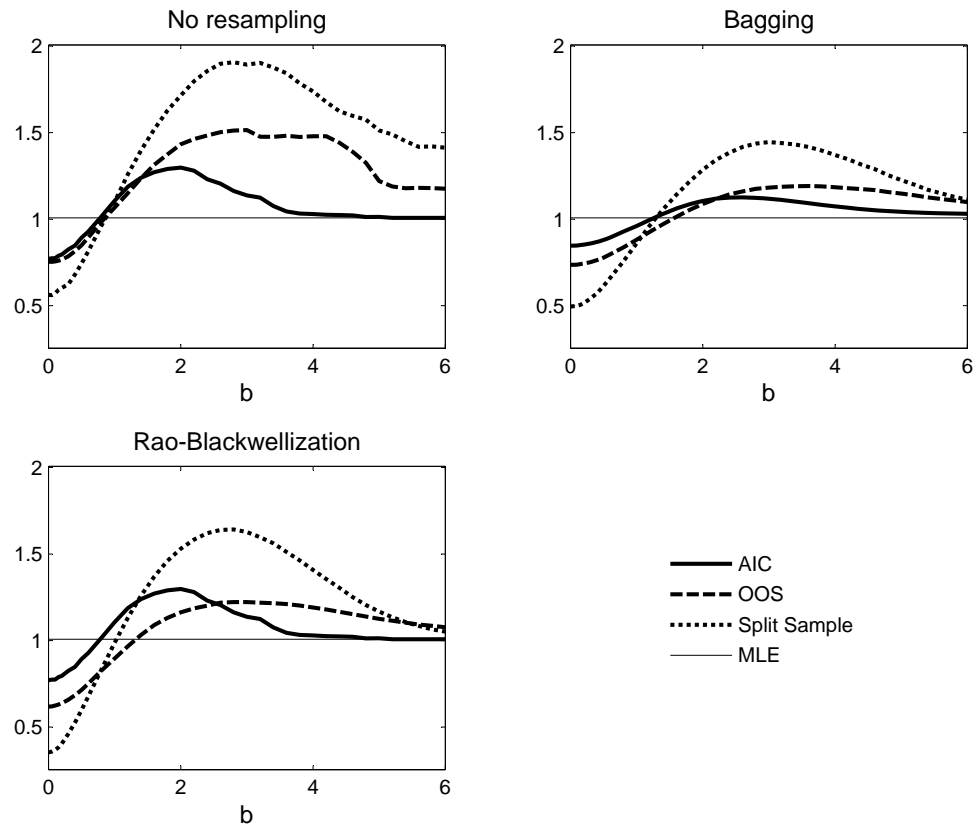
- MEESE, R. A., AND K. ROGOFF (1983): “Empirical Exchange Rate Models of the Seventies: Do They Fit Out of Sample?,” *Journal of International Economics*, 14, 3–24.
- PARK, J. (2002): “An Invariance Principle for Sieve Bootstrap in Time Series,” *Econometric Theory*, 18, 469–490.
- STOCK, J. H., AND M. W. WATSON (2012): “Generalized Shrinkage Methods for Forecasting Using Many Predictors,” *Journal of Business and Economic Statistics*, 30, 481–493.
- TIBSHIRANI, R. (1996): “Regression Shrinkage and Selection via the Lasso,” *Journal of the Royal Statistical Society, Series B*, 58, 267–288.
- TIMMERMANN, A. (2006): “Forecast Combination,” in *Handbook of Economic Forecasting*, ed. by C. W. Granger, G. Elliott, and A. Timmermann, Amsterdam. North Holland.
- VAN DER VAART, A. W. (1998): *Asymptotic Statistics*. Cambridge University Press, Cambridge.
- WEST, K. D. (2006): “Forecast Evaluation,” in *Handbook of Economic Forecasting*, ed. by C. W. Granger, G. Elliott, and A. Timmermann, Amsterdam. North Holland.
- WILSON, E. (1934): “The Periodogram of American Business Activity,” 34, 375–417.

**Table 1: Application Results**

	AICB	OOSB	SSB	OOSRB	SSRB
Probability of no predictors	0.72	0.00	0.40	0.00	0.34
Expected Number of predictors	0.28	4.86	0.60	5.00	0.66
Marginal probability of inclusion of each predictor:					
Book-to-market ratio	0.20	0.40	0.30	0.46	0.58
Treasury bill yields	0.08	0.36	0.14	0.68	0.02
Long-term yields	0.00	0.20	0.00	0.48	0.00
Net equity expansion	0.00	0.34	0.14	0.08	0.06
Inflation	0.00	0.32	0.00	0.22	0.00
Percent equity issuing	0.00	0.78	0.02	0.86	0.00
Long term returns	0.00	0.54	0.00	0.80	0.00
Stock variance	0.00	0.32	0.00	0.42	0.00
Default yield spread	0.00	0.40	0.00	0.52	0.00
Default return spread	0.00	0.26	0.00	0.00	0.00
Dividend-price ratio	0.00	0.20	0.00	0.12	0.00
Dividend yield	0.00	0.30	0.00	0.20	0.00
Earnings-price ratio	0.00	0.44	0.00	0.16	0.00
Forecast for Excess Returns in 2015	0.05	0.05	0.05	0.08	0.04
Memo: Unconditional Mean:1926-2014 : 0.06					

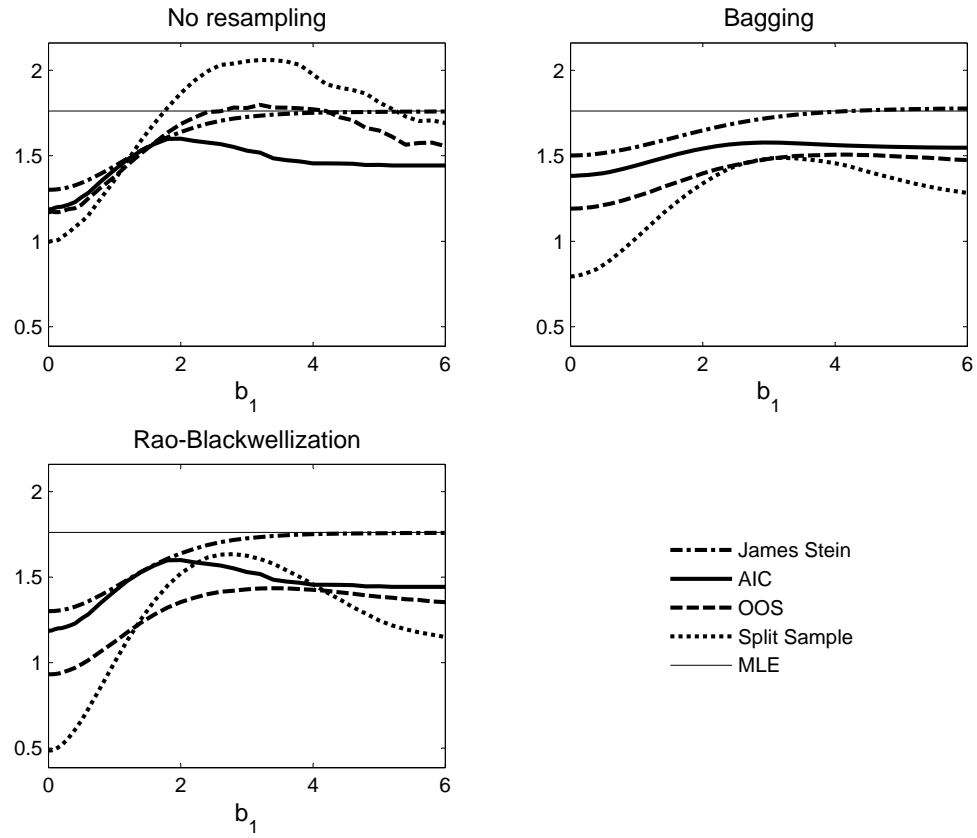
Notes: For the empirical application described in section 7, this table reports the fraction of draws with each simulation scheme that includes no predictors (other than the intercept), the average number of predictors (other than the intercept) chosen in each scheme, and the fraction of draws that includes each one of the predictors. The objective is forecasting year-ahead excess stock returns using annual data from 1926-2014. The forecast for excess stock returns in 2015 and the unconditional mean are also included. Model selection is over all 8,192 permutations of the 13 predictors. Definitions of the predictors and data sources are in Goyal and Welch (2008).

Figure 1: Local Asymptotic Risk ( $k = 1$ )



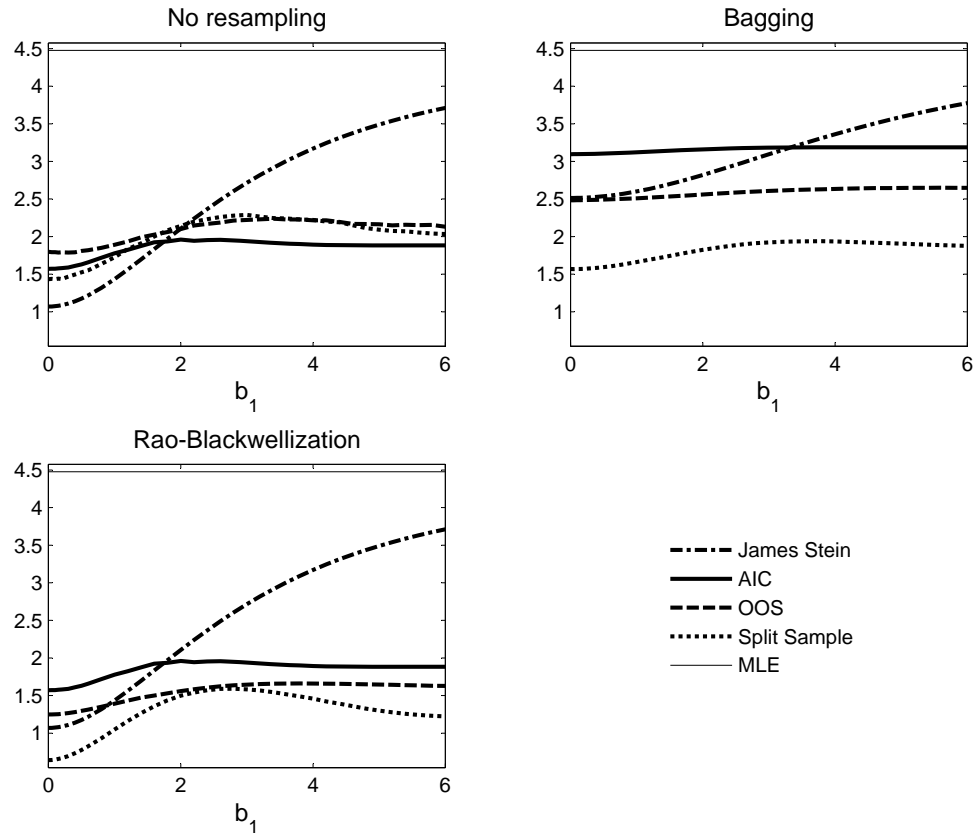
Notes: These are the simulated local asymptotic risk values, equation (5.1), for different procedures, plotted against  $b$ . Note that MLE is the same without any resampling, with bagging or with RB. AIC is the same without any resampling or with RB.

Figure 2: Local Asymptotic Risk ( $k = 3$ )



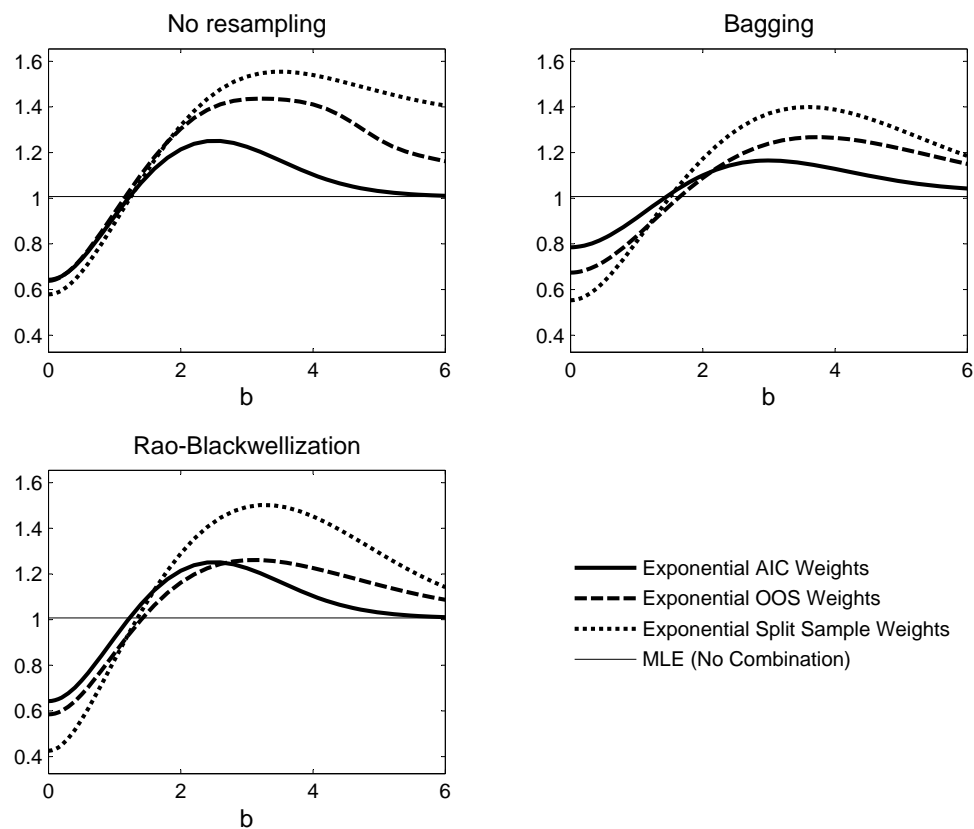
Notes: These are the simulated local asymptotic risk values, equation (5.1), for different procedures, plotted against  $b_1$ , where  $b = (b_1, 0, \dots, 0)'$ . Note that MLE is the same without any resampling, with bagging or with RB. AIC is the same without any resampling or with RB.

Figure 3: Local Asymptotic Risk ( $k = 20$ )



Notes: These are the simulated local asymptotic risk values, equation (5.1), for different procedures, plotted against  $b_1$ , where  $b = (b_1, 0, \dots, 0)'$ . Note that MLE is the same without any resampling, with bagging or with RB. AIC is the same without any resampling or with RB.

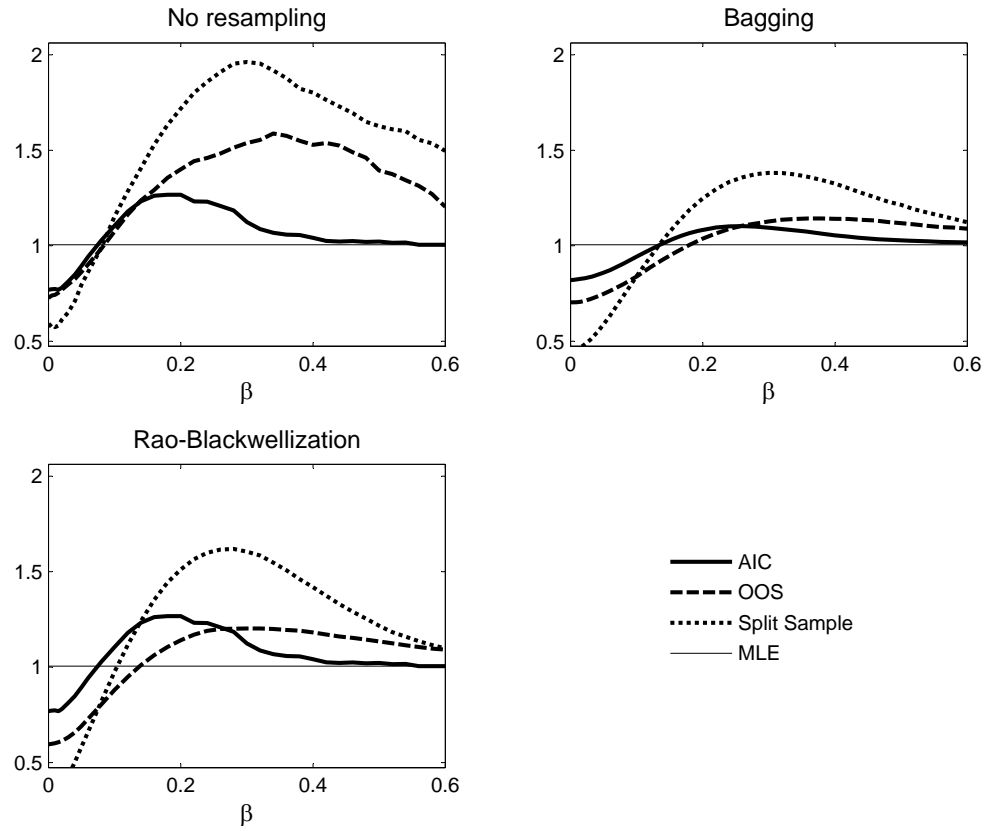
Figure 4: Local Asymptotic Risk: Combination Forecasts ( $k = 1$ )



Notes: These are the simulated local asymptotic risk values, equation (5.1), for different procedures, plotted against  $b$ . Note that MLE is the same without any resampling, with bagging or with RB. Exponential AIC forecast combination is the same without any resampling or with Rao-Blackwellization.



Figure 5: Root Normalized Mean Square Prediction Errors ( $k = 1$ )



Notes: These are the simulated root normalized mean square prediction errors using different procedures, plotted against  $\beta$ . There is one possible predictor and the sample size is  $T = 100$ . Note that MLE is the same without any resampling, with bagging or with RB. AIC is the same without any resampling or with RB.