

Bayesian Averaging, Prediction and Nonnested Model Selection

Han Hong and Bruce Preston¹

Previous version: January 2006

This version: November 2009

Abstract

This paper studies the asymptotic relationship between Bayesian model averaging and post-selection frequentist predictors in both nested and nonnested models. We derive conditions under which their difference is of a smaller order of magnitude than the inverse of the square root of the sample size in large samples. This result depends crucially on the relation between posterior odds and frequentist model selection criteria. Weak conditions are given under which consistent model selection is feasible, regardless of whether models are nested or nonnested and regardless of whether models are correctly specified or not, in the sense that they select the best model with the least number of parameters with probability converging to 1. Under these conditions, Bayesian posterior odds and BICs are consistent for selecting among nested models, but are not consistent for selecting among nonnested models and possibly overlapping models. These findings have important bearing for applied researchers who are frequent users of model selection tools for empirical investigation of model predictions.

JEL Classification: C14; C52

Keywords: Model selection criteria, Nonnested, Posterior odds, BIC

1 Introduction

Bayesian methods are becoming increasingly popular, both as a framework of model selection and also as a tool of forecasting — see, among others, Fernandez-Villaverde and Rubio-Ramirez (2004a), Schorfheide (2000), Stock and Watson (2001), Timmermann (2005), Clark and McCracken (2006) and Wright (2003). These methods are often used to summarize statistical properties of data, identify parameters of interest, and conduct policy evaluation. While empirical applications of these methods are abundant, less is understood about their theoretical sampling properties. This

¹Department of Economics, Stanford University and Department of Economics, Columbia University. We thank Raffaella Giacomini, Jin Hahn, Bernard Salanié, Barbara Rossi, Frank Schorfheide and Chris Sims, the editors, two anonymous referees, and participants of the SETA conference in Seoul, Korea for comments. We also thank the NSF and the Sloan Foundation for generous research support. The usual caveat applies.

paper provides a starting point for understanding the relation between Bayesian forecast averaging and frequentist model selection and prediction in a general framework that incorporates both nested and nonnested models.

We study the large sample properties of Bayesian prediction and model averaging for both nested and nonnested models. We first show that, for a single model, the difference between Bayesian and frequentist predictors are of smaller order of magnitude than the inverse of the square root of the sample size in large samples, regardless of the expected loss function used in forming the Bayesian predictors. This contrasts with the difference between MLE and Bayesian estimators, formed using a variety of loss functions, which is of the order $O_p\left(1/\sqrt{T}\right)$.

For multiple models, we derive general conditions under which the Bayesian posterior odds place asymptotically unit weight on the best model with the most parsimonious parameterization. Under these conditions, the Bayesian average model forecast is equivalent to the frequentist post-selection forecast up to a term that is of smaller order of magnitude than the inverse of the square root of the sample size in large samples. These findings essentially combine Schwarz' original contribution regarding BIC — that it is an asymptotic approximation to posterior odds — with the insights by Sin and White (1996) who demonstrate the inconsistency of BIC for selecting among nonnested models. The conditions we derive are weaker, more general, and allow for a much wider class of models.

An immediate consequence of multiple model comparison is that both the BIC and Bayesian posterior odds comparison are inconsistent in choosing a true and parsimonious model for selecting among nonnested models and some overlapping models. While this procedure will select one of the best fitting models, it does not necessarily choose the most parsimonious model with probability converging to 1 in large samples. Consistent selection among possibly nonnested models is feasible using nonnested model selection criteria in the spirit of Sin and White (1996).

These findings have important bearing for applied researchers who are frequent users of model selection tools for empirical investigation of model predictions. In addition, empirical analyses frequently find that forecasts generated from averages of a number of models typically perform better than forecasts of any one of the underlying models – see, for instance, Stock and Watson (2001). Our theoretical findings suggest that this can be because the models under consideration are close to each other and are all misspecified. As long as the posterior weights are non-degenerate among the set of models under comparison, it is possible for model averaging to outperform each model as long as all models are misspecified. Indeed, it is shown that for nonnested models, posterior weights will be non-degenerate, as long as the models are sufficiently close to each other. The case

of nested models is more interesting. It turns out that when the two models are sufficiently far from each other or sufficiently close from each other, the posterior weights will be degenerate. However, when the two models are “just close enough” but “not too close”, the posterior weights can be non-degenerate, and, as a consequence, model averaging can outperform each individual model.

Our results are of interest given the burgeoning use of Bayesian methods in the estimation of dynamic stochastic general equilibrium models in modern macroeconomics. See, *inter alia*, Fernandez-Villaverde and Rubio-Ramirez (2004a), Schorfheide (2000), Smets and Wouters (2002), Lubik and Schorfheide (2003) and Justiniano and Preston (2004) for examples of estimation in both closed and open economy settings. These papers all appeal to posterior odds ratios as a criterion for model selection. By giving a classical interpretation to the posterior odds ratio, the present paper intends to provide useful information regarding the conditions under which such selection procedures ensure consistency. The analysis contributes to understanding the practical limitations of standard model selection procedures given a finite amount of data.

The paper proceeds as follows. Section 2 describes model assumptions and derives their implications on the large sample behavior of the likelihood function. Section 3 demonstrates the asymptotic equivalence between Bayesian and frequentist predictors for a single model under weak conditions. The rest of the paper generalizes this result to multiple models. Section 4 first derives weak conditions under which the generalized posterior odds ratio is equivalent to BIC up to a term that is asymptotically negligible, and under which alternative model selection criteria are feasible to select consistently between both nested and nonnested models. Section 5 makes use of the asymptotic equivalence between posterior odds ratio and BIC to derive the relation between Bayesian model averaging and frequentist post-selection prediction. Finally, section 6 generalizes the implications for our results for Bayesian type model selection methods for non-likelihood-based objective functions as considered by Kim (2005), and section 7 concludes.

2 Model assumptions and implications

For clarity of exposition, in the rest of the paper we identify a model with the likelihood function that is being used to estimate model parameters. All results extend to general random distance functions that satisfy the stochastic equicontinuity assumptions stated below.

A parameter β is often estimated by maximizing a random log-likelihood function $\hat{Q}(\beta)$ associated with some model $f(y_t, \beta)$ that depends on observed data y_t and parameterized by the vector β :

$$\hat{Q}(\beta) \equiv Q(y_t, t = 1 \dots, T; \beta).$$

For example, under i.i.d sampling of the data, as in Vuong (1989) and Sin and White (1996), the

log-likelihood function takes the form of

$$\hat{Q}(\beta) = \sum_{t=1}^T \log f(y_t; \beta),$$

which minimizes the Kullback-Leibler distance between the parametric model and the data.

Under standard assumptions, the random objective function converges to a population limit when the sample size increases without bound. It is therefore assumed that there exists a function $Q(\beta)$, uniquely maximized at β_0 , which is the uniform limit of the random sample analog

$$\sup_{\beta \in \mathcal{B}} \left| \frac{1}{T} \hat{Q}(\beta) - Q(\beta) \right| \xrightarrow{p} 0.$$

Typically, the following decomposition holds for $\hat{Q}(\beta)$:

$$\hat{Q}(\hat{\beta}) = \underbrace{\hat{Q}(\hat{\beta}) - \hat{Q}(\beta_0)}_{(Qa)} + \underbrace{\hat{Q}(\beta_0) - TQ(\beta_0)}_{(Qb)} + \underbrace{TQ(\beta_0)}_{(Qc)}.$$

Under suitable regularity conditions, the following are true:

$$(Qa) = O_p(1), \quad (Qb) = O_p(\sqrt{T}), \quad (Qc) = O(T).$$

The regularity conditions under which the first equality holds are formally given below. They are the same as those in Chernozhukov and Hong (2003). They do not require the objective function to be smoothly differentiable, and permit complex nonlinear or simulation-based estimation methods. In particular, conditions that require smoothness of the objective function are typically violated in simulation-based estimation methods and in percentile-based non-smooth moment conditions. Even for simulation-based estimation methods, it can be difficult for researchers to insure that the simulated objective functions are smooth in model parameters.

ASSUMPTION 1 *The true parameter vector β_0 belongs to the interior of a compact convex subset \mathcal{B} of $R^{\dim(\beta)}$.*

ASSUMPTION 2 *For any $\delta > 0$, there exists $\epsilon > 0$, such that*

$$\liminf_{T \rightarrow \infty} P \left\{ \sup_{|\beta - \beta_0| \geq \delta} \frac{1}{T} \left(\hat{Q}(\beta) - \hat{Q}(\beta_0) \right) \leq -\epsilon \right\} = 1.$$

ASSUMPTION 3 *There exist quantities Δ_T , J_T , Ω_T , where $J_T \xrightarrow{p} -\mathcal{A}_\beta$, $\Omega_T = O(1)$,*

$$\frac{1}{\sqrt{T}} \Omega_T^{-1/2} \Delta_T \xrightarrow{d} N(0, I),$$

such that if we write

$$R_T(\beta) = \hat{Q}(\beta) - \hat{Q}(\beta_0) - (\beta - \beta_0)' \Delta_T + \frac{1}{2} (\beta - \beta_0)' (T J_T) (\beta - \beta_0)$$

then it holds that for any sequence of $\delta_T \rightarrow 0$

$$\sup_{|\beta - \beta_0| \leq \delta_T} \frac{R_T(\beta)}{1 + T|\beta - \beta_0|^2} = o_p(1).$$

THEOREM 1 *Under assumptions (1), (2) and (3), $\hat{Q}(\hat{\beta}) - \hat{Q}(\beta_0) = O_p(1)$.*

Given Pakes and Pollard (1989), Newey and McFadden (1994) and Andrews (1994), the result of Theorem 1 is rather straightforward. Its proof is incorporated in the beginning of the proof for Theorem 2 and is provided in the appendix.

The asymptotic distribution of $\hat{Q}(\hat{\beta}) - \hat{Q}(\beta_0)$ is also easy to derive in many situations. This distribution is useful for model selection tests but is not directly used in model selection criteria developed in section 4. In particular, satisfaction of the information matrix equality, $\Omega_T = -\mathcal{A}_\beta + o_p(1)$, is not necessary for our discussion of model selection criteria. This is especially relevant under potential model misspecification.

ASSUMPTION 4 $\hat{Q}(\beta_0) - TQ(\beta_0) = O_p(\sqrt{T})$.

Assumption 4 typically follows from an application of the central limit theorem

$$\frac{1}{\sqrt{T}} \left(\hat{Q}(\beta_0) - TQ(\beta_0) \right) \xrightarrow{d} N(0, \Sigma_Q), \quad (2.1)$$

where

$$\Sigma_Q = \lim Var \left(\frac{1}{\sqrt{T}} \left(\hat{Q}(\beta_0) - TQ(\beta_0) \right) \right).$$

For example, in the case of the log-likelihood function for i.i.d. observations, where

$$\hat{Q}(\beta) = \sum_{t=1}^T \log f(y_t; \beta) \quad \text{and} \quad Q(\beta) = E \log f(y; \beta),$$

such convergence follows immediately from the central limit theorem: $\Sigma = \text{Var}(\log f(y_t; \beta))$.

Beyond the likelihood setting, assumption 4 can in general be easily verified for extreme estimators based on optimizing random objective functions. These include M estimator, generalized method of moment estimators and their recent information theoretic variants. It can be shown to hold for generalized method of moment estimators regardless of whether the model is correctly specified or misspecified.

The property of the final term (Qc) is immediate.

3 Bayesian predictive analysis

Consider the Bayesian inference problem of predicting y_{T+1} given a data set Y_T with observations up to T . Typically we need to calculate the predictive density of y_{T+1} given Y_T , and for this purpose need to average over the posterior distribution of the model parameters β given the data Y_T :

$$\begin{aligned} f(y_{T+1}|Y_T) &= \int f(y_{T+1}|Y_T, \beta) f(\beta|Y_T) d\beta \\ &= \int f(y_{T+1}|Y_T, \beta) \frac{e^{\hat{Q}(\beta)} \pi(\beta)}{\int e^{\hat{Q}(\beta)} \pi(\beta) d\beta} d\beta. \end{aligned}$$

We will assume a markov structure of the model such that $f(y_{T+1}|Y_T; \beta) = f(y_{T+1}|y_T; \beta)$, where y_T is the most recent observation in Y_T . It is well understood that the first-order randomness in the prediction is driven by $f(y_{T+1}|y_T, \beta)$. The length of the prediction interval comprises two parts: the first is due to $f(y_{T+1}|y_T, \beta)$ and the second is due to the uncertainty from estimating β . While the second part will decrease to zero as the sample size T increases, the first part remains constant. It is straightforward to generalize y_T to a finite dimensional statistic of the sample.

We are interested in the second-order uncertainty in the prediction that is due to the estimation of the parameter β , and therefore will consider a fixed value \bar{y} of the random component y_T , and consider

$$f(y_{T+1}|\bar{y}, Y_T) = \int f(y_{T+1}|\bar{y}, \beta) f(\beta|Y_T) d\beta, \quad (3.2)$$

where \bar{y} can potentially differ from the observed realization of y_T in the sample. For example, one might consider out of sample predictions where \bar{y} does not take the realized value of y_T .

Point predictions can be constructed as functionals of the posterior predictive density $f(y_{T+1}|\bar{y}, Y_T)$. While researchers are most often interested in mean predictions and predictive intervals, a general

class of nonlinear predictors $\hat{\lambda}(\bar{y}, Y_T)$ can be defined as the solution to minimizing an expected loss function $\rho(\cdot)$, which is assumed to be convex:

$$\begin{aligned}\hat{\lambda}(\bar{y}, Y_T) &= \arg \min_{\lambda} E(\rho(y_{T+1}, \lambda) | \bar{y}, Y_T) \\ &= \arg \min_{\lambda} \int \rho(y_{T+1}, \lambda) f(y_{T+1} | \bar{y}, Y_T) dy_{T+1},\end{aligned}\tag{3.3}$$

where the predictive density $f(y_{T+1} | \bar{y}, Y_T)$ is defined in equation (3.2).

For example, a square loss function in which $\rho(x) = x^2$ results in the mean prediction

$$E(y_{T+1} | \bar{y}, Y_T) = \int E(y_{T+1} | \bar{y}, Y_T; \beta) f(\beta | Y_T) d\beta.$$

As another example, an absolute deviation loss function results in the median prediction

$$\text{med}(f_{y_{T+1}}(\cdot | \bar{y}, Y)) = \inf \left\{ x : \int^x f(y_{T+1} | \bar{y}, Y) dy_{T+1} \geq \frac{1}{2} \right\}.$$

As a more general class of examples, to construct a one-sided τ th predictive interval, we can take

$$\rho(y_{T+1}, \lambda) \equiv \rho_{\tau}(y_{T+1} - \lambda) = (\tau - 1(y_{T+1} \leq \lambda))(y_{T+1} - \lambda).\tag{3.4}$$

The following focuses on this loss function.

We are interested in comparing $\hat{\lambda}(\bar{y}, Y_T)$ to the nonlinear frequentist prediction, defined as

$$\tilde{\lambda}(\bar{y}, Y_T) = \arg \min_{\lambda} \int \rho(y_{T+1}, \lambda) f(y_{T+1} | \bar{y}, \hat{\beta}) dy_{T+1}.\tag{3.5}$$

Also define the infeasible loss function, where the uncertainty from estimation of the unknown parameters is absent, as

$$\begin{aligned}\bar{\rho}(\bar{y}, \lambda; \beta) &= \int \rho(y_{T+1}, \lambda) f(y_{T+1} | \bar{y}, \beta) dy_{T+1} \\ &= E(\rho(y_{T+1}, \lambda) | \bar{y}, \beta).\end{aligned}$$

Then the Bayesian predictor and the frequentist predictor can be written as

$$\hat{\lambda}(\bar{y}, Y_T) = \arg \min_{\lambda} \int \bar{\rho}(\bar{y}, \lambda; \beta) f(\beta | Y_T) d\beta,$$

and

$$\tilde{\lambda}(\bar{y}, Y_T) = \arg \min_{\lambda} \bar{\rho}(\bar{y}, \lambda; \hat{\beta}).$$

The next theorem establishes the asymptotic relation between $\hat{\lambda}(\bar{y}, Y_T)$ and $\tilde{\lambda}(\bar{y}, Y_T)$.

THEOREM 2 *Under assumptions (1), (2) and (3), and assuming that for each \bar{y} , $\bar{\rho}(\bar{y}, \lambda; \beta)$ is three times continuously differentiable with bounded derivatives in β and λ in a small neighborhood of β_0 and $\lambda_0 \equiv \arg \min_{\lambda} \bar{\rho}(\bar{y}, \lambda, \beta_0)$, with uniformly integrable derivatives, then*

$$\sqrt{T} \left(\hat{\lambda}(\bar{y}, Y_T) - \tilde{\lambda}(\bar{y}, Y_T) \right) \xrightarrow{p} 0.$$

Note that the conditions of this theorem only require the integrated loss function $\bar{\rho}(\bar{y}, \lambda; \beta)$ to be smoothly differentiable in λ and β , and does not impose smoothness conditions directly on $\rho(y_{T+1}, \lambda)$. Therefore, the results cover predictive intervals, as long as the predictive density given β is smoothly differentiable around the percentiles to be predicted.

The intuition of this theorem is that under the given regularity conditions, $\hat{\beta}$ is an asymptotic sufficient statistic for the random posterior distribution and $f(\beta|Y_T)$ is approximately normal with mean $\hat{\beta}$ and variance $-\frac{1}{T}(\mathcal{A}_{\beta})^{-1}$:

$$f(\beta|Y_T) \stackrel{A}{\approx} N\left(\hat{\beta}, -\frac{1}{T}(\mathcal{A}_{\beta})^{-1}\right).$$

This asymptotic approximation of the posterior distribution of β by a normal distribution forms the basis of the asymptotic approximation of the predictive distribution, and is also the key element of the results formally developed in section 4.

Mean prediction follows from the use of a quadratic loss function, and is convenient to analyze because of its linearity property. An immediate consequence of Theorem 2 is that

$$E(y_{T+1}|\bar{y}, Y_T) = E(y_{T+1}|\bar{y}; \hat{\beta}) + o_p\left(\frac{1}{\sqrt{T}}\right).$$

The same result can be obtained by directly applying a first-order Taylor expansion of $E(y_{T+1}|\bar{y}, \beta)$ around $\hat{\beta}$.

Given the generic notation of y_{T+1} , these arguments also apply without change to more general functions of y_{T+1} , in the sense that for a general function $t(\cdot)$ of y_{T+1} ,

$$E(t(y_{T+1})|\bar{y}, Y_T) = E\left(t(y_{T+1})|\bar{y}, \hat{\beta}\right) + o_p\left(\frac{1}{\sqrt{T}}\right).$$

Theorem 2 generalizes these results to nonlinear predictions that can be expressed as nonlinear functions of the predictive density $f(y_{T+1}|\bar{y}, Y_T)$ using the loss function setup in (3.3).

The conditions required for Theorem 2 are fairly weak and accomodate a wide range of models. These assumptions are automatically satisfied in a normal AR(1) model, in which $y_t = \beta y_{t-1} + \epsilon_t$,

where $\epsilon_t \sim N(0, \sigma^2)$. To illustrate, suppose that σ^2 is known, then in assumption 3, $\Delta_T = -\frac{1}{\sigma^2} \sum_{t=1}^T y_{t-1} \epsilon_t$, $\Omega_T = J_T = \frac{1}{\sigma^2} \frac{1}{T} \sum_{t=1}^T y_{t-1}^2$. In this model, $R_T(\beta) \equiv 0$ so that the last condition in assumption 3 is automatically satisfied. The identification assumption 2 also holds because of a nonsingular Hessian, as $E y_{t-1}^2 = \frac{\sigma^2}{1-\beta^2} > 0$. The differentiability assumption on $\bar{\rho}(\bar{y}, \lambda; \beta)$ is satisfied for the prediction interval as well. This is because while $\rho(y_{T+1}, \lambda)$ as defined in (3.4) is not smooth, its expectation conditional on \bar{y} and β defined as $\bar{\rho}(\bar{y}, \lambda; \beta)$ typically is. For example, when $\rho(y_{T+1}, \lambda) = 1(y_{T+1} \leq \lambda)$, $\bar{\rho}(\bar{y}, \lambda; \beta) = P(y_{T+1} \leq \lambda | \bar{y}, \beta)$ is often differentiable multiple times in both λ and β . In the AR(1) example, $P(y_{T+1} \leq \lambda | \bar{y}, \beta) = \Phi\left(\frac{\lambda - \beta' \bar{y}}{\sigma}\right)$.

Different loss functions can be used to construct a variety of Bayesian point estimators from the posterior density. Unless the loss function is convex and symmetric around 0, the corresponding Bayes estimator is typically different from the frequentist maximum likelihood estimator at the order of $O_p(1/\sqrt{T})$. In contrast, we found that when different loss functions are used to define different predictors, Bayesian and frequentist predictors coincide with each other up to the order $o_p(1/\sqrt{T})$. This is probably a more relevant result concerning loss functions because researchers are typically more interested in using loss functions to define the properties of predictions, rather than to define the estimator itself. Note also that we consider only prediction at a fixed value and have not yet considered out of sample variation. The recent work of Hansen (2008) discovered important negative correlation between in-sample and out-of-sample variations of predictive criterion functions.

4 Posterior odds and consistent model selection

Bayesian averaging is a popular method for making predictions in the context of multiple models. The weights used in Bayesian averaging are calculated through posterior odds ratios. To generalize the results in the previous section to multiple models, it is important to understand first the relation between posterior odds ratio calculation and consistent frequentist model selection criteria. This is of independent interest given the increasing use of Bayesian methods in economics; particularly the recent macroeconomics literature on estimation of dynamic stochastic general equilibrium models which makes use of the posterior odds ratio for model selection and prediction.

4.1 Large sample properties of Bayes factors

The posterior distribution is an integral transformation of the model (f, Q) , defined as

$$e^{\hat{Q}(\beta)} \pi(\beta) / \int e^{\hat{Q}(\beta')} \pi(\beta') d\beta'.$$

Associated with the posterior distribution is the Bayes factor

$$P_Q \times \int e^{\hat{Q}(\beta)} \pi(\beta) d\beta,$$

where P_Q is a prior probability weight of model (f, Q) .

The following theorem connects the properties of the Bayes posterior distribution to the extremum estimator and establishes the relation between the Bayes factor and the extremum estimator analog of BICs. Under the regularity conditions stated in assumptions 1 to 3, the Bayes factor is asymptotically equivalent to using BIC as a model selection criterion. While the relation between Bayes factor and BIC comes from the original contribution by Schwartz, the conditions in the following theorem are considerably weaker and more general.

THEOREM 3 *Under assumptions (1), (2) and (3), and suppose that $\hat{Q}(\hat{\beta}) = \inf_{\beta} \hat{Q}(\beta) + o_p(T^{-1/2})$, the Bayes factor satisfies the following relation*

$$T^{\frac{\dim(\beta)}{2}} \int e^{\hat{Q}(\beta) - \hat{Q}(\hat{\beta})} \pi(\beta) d\beta \xrightarrow{p} \pi(\beta_0) (2\pi)^{\frac{\dim(\beta)}{2}} \det(-\mathcal{A}_{\beta})^{-1/2}.$$

The formal details of the proof are relegated to an appendix. When $\hat{Q}(\beta)$ is smoothly differentiable, intuition can be gleaned by considering the expression:

$$\log \int e^{\hat{Q}(\beta)} \pi(\beta) d\beta - \hat{Q}(\hat{\beta}) = \log \int e^{\hat{Q}(\beta) - \hat{Q}(\hat{\beta})} \pi(\beta) d\beta.$$

It can be approximated up to an $o_p(1)$ term as follows. First, $\hat{Q}(\beta) - \hat{Q}(\hat{\beta})$ can be approximated by a quadratic function centered at $\hat{\beta}$, in a size $1/\sqrt{T}$ neighborhood around $\hat{\beta}$:

$$\hat{Q}(\beta) - \hat{Q}(\hat{\beta}) \approx \frac{1}{2} (\beta - \hat{\beta})' \frac{\partial^2 \hat{Q}(\hat{\beta})}{\partial \beta \partial \beta'} (\beta - \hat{\beta}).$$

Second, in this $1/\sqrt{T}$ size neighborhood, the prior density is approximately constant around β_0 as $\pi(\hat{\beta}) \xrightarrow{p} \pi(\beta_0)$. The impact of the prior density is negligible except for the value at $\pi(\beta_0)$. Outside the $1/\sqrt{T}$ size neighborhood around $\hat{\beta}$, the difference between $e^{\hat{Q}(\beta)}$ and $e^{\hat{Q}(\hat{\beta})}$ is exponentially small, and makes only asymptotically negligible construction to the overall integral.

The appendix proves formally the approximation:

$$\begin{aligned}
\log \int e^{\hat{Q}(\beta)} \pi(\beta) d\beta - \hat{Q}(\hat{\beta}) &= \log \pi(\beta_0) \int e^{\frac{1}{2}(\beta-\hat{\beta})' \frac{\partial^2 \hat{Q}(\hat{\beta})}{\partial \beta \partial \beta'} (\beta-\hat{\beta})} d\beta + o_p(1) \\
&= \log \pi(\beta_0) \int e^{\frac{1}{2}(\beta-\hat{\beta})' T \mathcal{A}_\beta (\beta-\hat{\beta})} d\beta + o_p(1) \\
&= \log \left(\pi(\beta_0) (2\pi)^{\frac{\dim(\beta)}{2}} \det(-T \mathcal{A}_\beta)^{-\frac{1}{2}} \right) + o_p(1) \\
&= \underbrace{\log \pi(\beta_0) + \frac{\dim(\beta)}{2} \log(2\pi) - \frac{1}{2} \det(-\mathcal{A}_\beta)}_{C(\mathcal{A}_\beta, \beta)} - \frac{1}{2} \dim(\beta) \log T + o_p(1),
\end{aligned}$$

where $C(\mathcal{A}_\beta, \beta)$ can also depend on the prior weight on the model itself.

4.2 Consistent model selection

Now introduce a competing model, $g(Y; \alpha)$, $\alpha \in \Lambda$ (generalizing the notations to multiple models is immediate), where α is estimated by the random log-likelihood function:

$$\hat{L}(\alpha) \equiv L(y_t, t = 1, \dots, T; \alpha).$$

Similarly to model (f, Q) , we assume that the decomposition of

$$\hat{L}(\hat{\alpha}) = \underbrace{\hat{L}(\hat{\alpha}) - \hat{L}(\alpha_0)}_{(La)} + \underbrace{\hat{L}(\alpha_0) - TL(\alpha_0)}_{(Lb)} + \underbrace{TL(\alpha_0)}_{(Lc)}$$

satisfies the stochastic order properties

$$(La) = O_p(1), \quad (Lb) = O_p(\sqrt{T}), \quad (Lc) = O(T).$$

We also assume that assumptions (1), (2), (3) and (4), and subsequently theorems 1, 2 and 3 hold for model (g, L) .

We follow the literature to define consistency in the context of model selection as a procedure which selects the most parsimonious model among the set of the models that generate the best fit with the data, cf Andrews and Lu (2001) and Hong, Preston, and Shum (2003). Formally, a consistent rule for model selection and parameter estimation is defined as the pair of parameter and model indicators $(\hat{\gamma}, \hat{\eta})$ where $\hat{\gamma} \in \{\hat{\beta}, \hat{\alpha}\}$ and $\hat{\eta} \in \{Q, L\}$, such that if

$$Q(\beta_0) > L(\alpha_0) \quad \text{or} \quad Q(\beta_0) = L(\alpha_0) \quad \text{and} \quad \dim(\beta) < \dim(\alpha), \quad (4.6)$$

then $P(\hat{\eta} \in \{Q, L\} = (\hat{\beta}, Q)) \rightarrow 1$. Otherwise $P(\hat{\eta} \in \{Q, L\} = (\hat{\alpha}, L)) \rightarrow 1$.

Obviously, this definition of consistency is specific to the context of model selection. In contrast, in hypothesis testing, consistency usually refers to a test procedure for which both the type 1 and type 2 errors go to zero as the sample size increases to infinity.

A conventional consistent model selection criterion takes the form of a comparison between

$$\hat{Q}(\hat{\beta}) - \dim(\beta) * C_T \quad \text{and} \quad \hat{L}(\hat{\alpha}) - \dim(\alpha) * C_T,$$

where C_T is a sequence of constants that tends to ∞ as T goes to ∞ , at a rate to be prescribed below, and $\dim(\cdot)$ indexes the parametric dimension of the model. $\hat{\beta}$ and $\hat{\alpha}$ are maximands of the corresponding objective functions. The second term in the model selection criteria penalizes the dimension of the parametric model. For instance, BIC takes $C_T = \log T$ and AIC adopts $C_T = 2$.

While the discussion so far has focused on parametric likelihood models, all results can be applied without modification to other non-likelihood-based models, including method of moment models. They are also applicable to non-log-likelihood random criterion functions that measure other notions of distance between the model and the data generating process, such as the Cressie-Read discrepancy statistic in generalized empirical likelihood methods.

For non-likelihood-based models, in addition to penalizing the parametric dimension of the model, the dimension of the estimation procedure, such as the number of moment conditions could also be penalized as in Andrews and Lu (2001). The goal of the latter penalization term is to preserve the parsimony of the model and the informativeness of the estimation procedure, thereby improving the precision of the model with a finite amount of data. To emphasize this possibility, we will adopt a general notation $\dim(f, Q)$ and $\dim(g, L)$, rather than $\dim(\beta)$ and $\dim(\alpha)$. The first argument refers to the parametric dimension of the model, while the second argument refers to the dimension of the information used in the inference procedure regarding model parameters.

For the specific case of the Bayesian (Schwartz) information criterion,

$$BIC = \hat{Q}(\hat{\beta}) - \hat{L}(\hat{\alpha}) - (\dim(f, Q) - \dim(g, L)) \times \log T$$

implicitly defining $C_T = \log T$. It now will be demonstrated that whether models are nested or nonnested has important consequences for the asymptotic properties of $\hat{Q}(\hat{\beta}) - \hat{L}(\hat{\alpha})$ and consistency of the BIC as model selection criterion. The two cases of nested and nonnested models are treated in turn.

4.2.1 *Nested Models*

There are two cases of interest: one model is better in the sense that $Q(\beta_0) \neq L(\alpha_0)$ and both models are equally good in the sense that $Q(\beta_0) = L(\alpha_0)$. Consider the former. Without loss of

generality, let $Q(\beta)$ be the larger nesting model and $L(\alpha)$ be the smaller model nested in $Q(\beta)$. When $\beta_0 \neq \alpha_0$, it is typically the case that $Q(\beta_0) > L(\alpha_0)$. The goal of BIC is to select the correct model, $Q(\beta)$, with probability converging to 1. In this case, this will be true because

$$BIC = \underbrace{(Qa) - (La)}_{O_p(1)} + \underbrace{(Qb) - (Lb)}_{O_p(\sqrt{T})} + \underbrace{T(Q(\beta_0) - L(\alpha_0))}_{O(T)} - (\dim(f, Q) - \dim(g, L)) \times \log T,$$

will be dominated by $+T(Q(\beta_0) - L(\alpha_0))$, which increases to $+\infty$ with probability converging to 1 as $T \rightarrow \infty$. In other words, for any $M > 0$,

$$P(BIC > M) \longrightarrow 1.$$

Hence, BIC selects the correct model with probability converging to 1 if there is one correct model. More generally, BIC selects one of the correct models with probability converging to 1.

Suppose now that both models are equally good, so that $Q(\beta_0) = L(\alpha_0)$. Because we are discussing nested models, this also means that $\beta_0 = \alpha_0$ (with the obvious abuse of the notation of equality with different dimensions), and that $\hat{Q}(\beta_0) = \hat{L}(\alpha_0)$ almost surely. In the case of likelihood models, this means that $f(Z_t; \beta_0) = g(Z_t; \alpha_0)$ almost surely, since $g(\cdot)$ is the smaller model nested inside $f(\cdot)$. The true model lies in the common subset of (β_0, α_0) and therefore has to be the same model.

In this case, the second term is identically equal to 0:

$$(Qb) - (Lb) = \hat{Q}(\beta_0) - \hat{L}(\alpha_0) - (TQ(\beta_0) - TL(\alpha_0)) \equiv 0.$$

Given that the last terms (Qc) and (Lc) disappear as a consequence of the equality $Q(\beta_0) = L(\alpha_0)$, the BIC comparison is reduced to

$$BIC = \underbrace{(Qa) - (La)}_{O_p(1)} - (\dim(f, Q) - \dim(g, L)) \times \log T.$$

The second term, which is of order $O(\log T)$, will dominate. So if $\dim(f, Q) > \dim(g, L)$, BIC will converge to $-\infty$ with probability converging to 1. In other words, for any $M > 0$,

$$P(BIC < -M) \longrightarrow 1.$$

Hence, given two equivalent models, in the sense of $Q(\beta_0) = L(\alpha_0)$, the BIC will choose the most parsimonious model (namely the one with the smallest dimension, either $\dim(f, Q)$ or (g, L)) with probability converging to 1.

It is clear from the above arguments that instead of using $C_T = \log T$, we can choose any sequence of C_T such that $C_T \rightarrow \infty$ and $C_T = o(T)$.

4.2.2 *Nonnested Models*

Many model comparisons are performed among models that are not nested inside each other. A leading example is the choice between a nonlinear model and its linearized version. For instance, in an extensive literature on estimating consumption Euler equations, there has been debate on the appropriateness of using log-linear versus non-linear Euler equations to estimate household preference parameters. See, for instance, Carroll (2001), Paxson and Ludvigson (1999), and Attanasio and Low (forthcoming). More recently, Fernandez-Villaverde and Rubio-Ramirez (2003) and Fernandez-Villaverde and Rubio-Ramirez (2004b) show that nonlinear filtering methods render feasible the estimation of some classes of nonlinear dynamic stochastic general equilibrium models. Being equipped with model selection criteria to compare multiple non-linear non-nested models is clearly desirable.

In contrast to the case of nested models, the comparison of nonnested models imposes more stringent requirements on C_T for consistent model selection. Indeed, the further condition on C_T that $C_T/\sqrt{T} \rightarrow \infty$ is required in addition to $C_T = o(T)$. As an example, $C_T = \sqrt{T} \log T$ will satisfy both requirements, but $C_T = \log T$ will not satisfy the second requirement. Let's call the model selection criteria *NIC* when we choose $C_T = \sqrt{T} \log T$. To our knowledge these rate conditions are first due to Sin and White (1996) in the context of smooth likelihood models.

Suppose that $Q(\beta_0)$ is greater than $L(\alpha_0)$, implying model (f, β) is better than model (g, α) . Here the comparison of the quality of the two models is made under a specific context, in which f and g are the likelihoods of two competing models and the comparison is based on the Kullback Leibler information criteria between the true data generating process and the postulated models. The model selection tests proposed in Vuong (1989) make use of the same criteria for model comparison in nonnested models. In particular, this specific notion of model quality rules out the possibility of comparing models based on prediction fitting quality criteria that may differ from the criteria that are being used for parameter estimation.

As an example, consider a collection of regression functions $y_i = g_m(x_i, \gamma_m) + \epsilon_i^m$. The usual criterion for the fit of the regression functions $g_m(x_i, \gamma_m)$, at least for independent data, is the least square norm $E(y_i - g_m(x_i, \gamma_m))^2$. However, dependent on the parametric specification of the density of ϵ_m , the parameters γ_m can be estimated using a pseudo-likelihood method even when the parametric density of ϵ_m is potentially misspecified. When $m = 2$ and the densities for ϵ_1 and ϵ_2 are $f(\epsilon_i^1, \nu)$ and $h(\epsilon_i^2, \eta)$. The limit objective functions corresponding to the pseudo-likelihood estimators are $Q(\gamma_1^0) = E \log f(\epsilon_i^1, \nu)$ and $L(\gamma_2^0) = E \log h(\epsilon_i^2, \eta)$. When the two models are nonnested, even when $Q(\gamma_1^0) > L(\gamma_2^0)$, it is possible for $g_2(x_i, \gamma_2)$ to fit the dependent

data y_i better than $g_1(x_i, \gamma_1)$ by the mean square error criterion. Our assumption rules out the discrepancy between the objective function that is being used to obtain parameter estimates and that which is being used to comparing model qualities.²

Then, as before, NIC is dominated by $T(Q(\beta_0) - L(\alpha_0))$, which increases to $+\infty$ with probability converging to 1. This is true regardless of whether we choose $C_T = \log T$ or $C_T = \sqrt{T} \log T$, since both are of smaller order of magnitude than T . The behavior of NIC when one model is better than the other is essentially the same for both nested models and nonnested models.

When both models are equally good, so that $Q(\beta_0) = L(\alpha_0)$, NIC comprises the non-vanishing term

$$(Qb) - (Lb) = \hat{Q}(\beta_0) - \hat{L}(\alpha_0) - (TQ(\beta_0) - TL(\alpha_0))$$

which is of order $O(\sqrt{T})$. In contrast to the nested case, it is no longer true that $\hat{Q}(\beta_0) \equiv \hat{L}(\alpha_0)$ with probability one when the two models are equally good. The model selection criterion takes the form

$$NIC = \underbrace{(Qa) - (La)}_{O_p(1)} + \underbrace{(Qb) - (Lb)}_{O_p(\sqrt{T})} - (\dim(f, Q) - \dim(g, L)) \times C_T.$$

For choice of penalty function $C_T = \sqrt{T} \log T$, or any other sequence that increases to ∞ faster than \sqrt{T} , the last term dominates. So if $\dim(f, Q) > \dim(g, L)$, NIC converges to $-\infty$ with probability converging to 1, or

$$P(NIC < -M) \longrightarrow 1 \quad \text{for any } M > 0.$$

Thus, NIC will choose the most parsimonious model among the two models with probability converging to 1.

In contrast, if the BIC had been used, where $C_T = \log T$, then the final term fails to dominate. The second term, which is random, might instead dominate. It is immediate that model (f, Q) and model (g, L) will both be selected with strictly positive probabilities. Such model selection behavior does not have a clear interpretation when a unique ranking of models is desired. While an analyst may be content with identifying just one of the best fitting models, regardless of whether it is the most parsimonious or not, the definition of consistency we adopt from the literature seeks to uniquely rank models given a model selection criterion and a particular set of assumptions. As such, positive probability weights on multiple models fails our requirements.

²We thank a referee for pointing this out to us.

While the properties of AIC and BIC for selecting among nested parametric models are well understood (e.g. see Sin and White (1996)), the above discussion about the NIC is summarized in the following. Suppose that assumptions (1), (2), (3) and (4) hold for models (f, Q) and (g, L) , then the NIC model selection criterion, defined as

$$NIC = \hat{Q}(\hat{\beta}) - \hat{L}(\hat{\alpha}) - (\dim(f, Q) - \dim(g, L)) \times \sqrt{T} \log T$$

has the following properties for (f, Q) and (g, L) either nested or nonnested:

1. If $Q(\beta_0) > L(\alpha_0)$ then $P(NIC > M) \rightarrow 1$ for all $M > 0$;
2. If $Q(\beta_0) = L(\alpha_0)$ and $\dim(f, Q) - \dim(g, L) > 0$ then $P(NIC > M) \rightarrow 1$ for all $M > 0$.

These findings can be anticipated to some degree from Sin and White (1996) given Pakes and Pollard (1989), Newey and McFadden (1994) and Andrews (1994).

4.2.3 Overlapping Models

The competing models can also be overlapping (Vuong (1989) and Sin and White (1996)). These models can behave either like nested models or nonnested models, depending on the relation between the pseudo-true parameters of the competing models. In the first case, when $Q(\beta_0) = L(\alpha_0)$ implies $\hat{Q}(\beta_0) \equiv \hat{L}(\alpha_0)$, these models behave like nested models. In the second case, when $Q(\beta_0) = L(\alpha_0)$ does not imply that $\hat{Q}(\beta_0) \equiv \hat{L}(\alpha_0)$, these models behave like nonnested models. Because the conditions of consistent model selection for nonnested models are strictly stronger than those for nested models, when there is uncertainty regarding the behavior of overlapping models, we suggest following the rules for nonnested models, which require that $C_T/\sqrt{T} \rightarrow \infty$ and that $C_T = o(T)$.

4.3 Posterior odds and BIC

The posterior odds ratio between the two models is defined as

$$\log \frac{P_Q}{P_L} \times \frac{\int e^{\hat{Q}(\beta)} \pi(\beta) d\beta}{\int e^{\hat{L}(\alpha)} \gamma(\alpha) d\alpha} \quad (4.7)$$

when the two models have prior probability weights P_Q and $P_L = 1 - P_Q$. A classical result due to Schwarz (1978) shows that by exploiting the approximation in section 4.1, up to a term of order $o_p(1)$, the log posterior odds ratio can be written as

$$\hat{Q}(\hat{\beta}) - \hat{L}(\hat{\alpha}) - \left(\frac{1}{2} \dim(\beta) - \frac{1}{2} \dim(\alpha) \right) \log T + C(\mathcal{A}_\beta, \beta) - C(\mathcal{A}_\alpha, \alpha) + \log \frac{P_Q}{P_L}$$

implicitly defining a penalization term $C_T = \log T$. It is immediately clear that this expression is asymptotically equivalent, up to a constant that is asymptotically negligible, to BIC. As discussed previously, this choice of C_T gives a consistent model selection criterion only when comparing nested models. In the nonnested case, when the null hypothesis contends that two models are asymptotically equivalent in fit and therefore misspecified, it fails to select the most parsimonious model with probability converging to 1.

Fernandez-Villaverde and Rubio-Ramirez (2004a) also explore the large sample properties of the posterior odds ratio in the case of parametric likelihood. They demonstrate, assuming there exists a unique asymptotically superior model in the sense of $Q(\beta) > L(\alpha)$, that the posterior odds ratio will select model (f, Q) with probability converging to 1 as T goes to infinity. Theorem 3 serves to generalize this finding. First, the results presented here are established under very weak regularity conditions. Second, only by considering a null hypothesis that admits the possible equivalence of the two models, in the sense that $Q(\beta) = L(\alpha)$, can a complete classical interpretation properly be given to the posterior odds ratio as a model selection criterion. Statistically distinguishing models is fundamental to classical hypothesis testing. Given the absence of prior information in classical estimation it is necessary to entertain the possibility that two or more models are equally good in a testing framework. The results of this paper are therefore seen to be couched naturally in the classical paradigm.

4.4 Multiple Models

The previous theorem can be immediately extended to the case of multiple models. Suppose there are a total of M models denoted by $\Gamma_m, m = 1, \dots, M$ with corresponding sets of parameters $\gamma_m, m = 1, \dots, M$, of which k of them are equally good, in the sense of having the same limit objective function with magnitude $Q(\beta_0)$, but the other $M - k$ models have lower valued limit objective functions. For example, when $M = 2$, $\Gamma_1 = Q$ and $\Gamma_2 = L$, $\gamma_1 = \beta$ and $\gamma_2 = \alpha$.

The notion of consistent model selection defined in 4.6 can be generalized to allow for multiple models. Let \mathcal{BC} denote the space of (m, γ_m) vectors, which can be viewed as a “generalized parameter space” including a model choice indicator and the set of parameters corresponding to the chosen model. Define the set \mathcal{BCL}^0 as the set of (m, γ_m) for which $\Gamma_m(\gamma_m)$ achieves the minimum over the collection of the generalized parameter set. Furthermore, define

$$\mathcal{MBCL}^0 = \{(m, \gamma_m) \in \mathcal{BCL}^0 : \dim(\gamma_m) \geq \dim(\gamma_l), \forall l \in \mathcal{BCL}^0\}.$$

According to this definition, it can be shown analogously to the two model case that with probability converging to 1, both the BICs and the NICs for the k good models will be infinitely larger than

the those for the $M - k$ inferior models. In other words, with probability converging to 1, none of the $M - k$ inferior models will ever be selected by either BIC or NIC comparison.

However, the behavior of BIC and NIC can be very different among the k best models depending on whether they are nested or nonnested. If the k best models are all nested inside each other, both BICs and NICs will select the most parsimonious model with probability converging to 1, and are consistent in the sense defined above. In contrast, if there exist two models that are not nested inside each other, then BIC will put random weights on at least two models. But NIC will still choose the most parsimonious model among all the k best models.

5 Bayesian model averaging and post-selection prediction

The results in the previous section can now be applied to study prediction using the average of two or more models. It is well known (see for example Pötscher (1991)) that the frequentist asymptotic distribution properties of post-selection estimation and prediction are not affected by the model selection step, as long as the model selection criterion is consistent. Here we demonstrate conditions under which the property of consistent model selection is possessed by the Bayesian prediction procedure (in the sense of asymptotically placing unitary probability weight on a single model).

With two models (f, Q) and (g, L) , we can write the predictive density as

$$f(y_{T+1}|\bar{y}, Y_T) = \frac{BF_Q}{BF_Q + BF_L} f(y_{T+1}|\bar{y}, Y_T, Q) + \frac{BF_L}{BF_Q + BF_L} f(y_{T+1}|\bar{y}, Y_T, L)$$

where the posterior probability weights on each model are defined as

$$BF_Q = P_Q \int e^{\hat{Q}(\beta)} \pi(\beta) d\beta \quad \text{and} \quad BF_L = P_L \int e^{\hat{L}(\alpha)} \gamma(\alpha) d\alpha.$$

The prior weights P_Q and P_L are defined in (4.7).

In the case of likelihood models $e^{\hat{Q}(\beta)} = f(Y_T|\beta)$ and $e^{\hat{L}(\alpha)} = f(Y_T|\alpha)$. In particular, note

$$f(Y_T) = BF_Q + BF_L$$

is the marginal density of the data Y_T . The model specific predictive densities are

$$\begin{aligned} f(y_{T+1}|\bar{y}, Y_T, Q) &= \frac{\int e^{\hat{Q}(\beta)} \pi(\beta) f(y_{T+1}|\bar{y}, \beta) d\beta}{\int e^{\hat{Q}(\beta)} \pi(\beta) d\beta} \\ &= \int f(y_{T+1}|\bar{y}, \beta) f(\beta|Y_T, Q) d\beta \end{aligned}$$

and $f(y_{T+1}|\bar{y}, Y_T, L) = \int f(y_{T+1}|\bar{y}, \alpha) f(\alpha|Y_T, L) d\alpha$. As before, a general class of predictions can be defined by

$$\hat{\lambda}(\bar{y}, Y_T) = \arg \min_{\lambda} \int \rho(y_{T+1}, \lambda) f(y_{T+1}|\bar{y}, Y_T) dy_{T+1}.$$

The following theorem establishes the asymptotic properties of the Bayesian predictor formed by averaging the two models.

THEOREM 4 *Suppose the assumptions stated in Theorem 2 hold for both models (f, Q) and (g, L) . Also assume that one of the following two conditions holds:*

1. $Q(\beta_0) > L(\alpha_0)$, (f, Q) and (g, L) can be either nested or nonnested.
2. $Q(\beta_0) = L(\alpha_0)$ but (f, Q) is nested inside (L, g) . In addition, $\dim(\alpha) - \dim(\beta) > 1$.

Then $\sqrt{T} \left(\hat{\lambda}(\bar{y}, Y_T) - \tilde{\lambda}(\bar{y}, Y_T) \right) \xrightarrow{p} 0$, where $\tilde{\lambda}(\bar{y}, Y_T) = \arg \min_{\lambda} \bar{\rho}(\bar{y}, \lambda; \hat{\beta})$ is formed from (f, Q) .

It is clear from previous discussion that if $Q(\beta_0) = L(\alpha_0)$ and (f, Q) and (g, L) are not nested, then the weight on neither BF_Q nor BF_L will converge to 1, and the behavior of $f(y_{T+1}|\bar{y}, Y_T)$ will be a random average between the two models.

These theoretical results contrast with the conventional wisdom regarding forecasting: that averages (or some combination) of available forecasting models perform better than any one model taken in isolation. Recent studies by Stock and Watson (2001), Stock and Watson (2003) and Wright (2003) adduce evidence consistent with this view. The former papers show that a range of combination forecasts outperform benchmark autoregressive models when forecasting output growth of seven industrialized economies and similarly that the predictive content of asset prices in forecasting inflation and output growth when models are averaged improves upon univariate benchmark models. The latter paper shows that Bayesian model averaging can improve upon the random walk model of exchange rate forecasts. Understanding the fundamental cause of this disjunction between theoretical and empirical findings is left for future work. However, we offer the following remarks based on our findings

For forecasts generated from averages of a number of models to perform better than forecasts of any one of the underlying models, the Bayesian posterior weights have to be asymptotically non-degenerate among at least two of the models that are being averaged. From the previous analysis,

we have shown that this is possible among nonnested models when two of the models are “equally good” ($Q(\beta_0) = L(\alpha_0)$).

Our analysis also suggests that among nested models, unless their dimensions are essentially the same ($\dim(f, Q) = \dim(g, L)$), it is most likely that posterior weights will concentrate on one of the models, either the “best” model or the most parsimonious model. It turns out that it is still possible for the posterior weights to be non-degenerate among nested model when the two models are “just close enough” but “not too close”. To see this, suppose that $Q(\beta_0) - L(\alpha_0) = h(T)$, where $h(T)$ is a function of the sample size T . Recall the decomposition for nested models:

$$\underbrace{(Qa) - (La)}_{O_p(1)} + \underbrace{T(Q(\beta_0) - L(\alpha_0))}_{T \times O(h(T))} - (\dim(f, Q) - \dim(g, L)) \times C_T.$$

As long as $T \times h(T)$ is approximately the same order as C_T , the penalization term will not dominate the second term, and it is then possible for posterior weights to be non-degenerate if the last two terms happen to approximately offset each other.

6 Generalized nonnested model selection criteria

While results in the previous sections are stated with the parametric log-likelihood function in mind, they apply without modification to a wide class of random non-likelihood-based objective functions assuming direct comparison between $Q(\beta)$ and $L(\alpha)$ is interpretable. Such interpretation is possible, for example, when both are constructed as information-theoretic empirical likelihoods corresponding to a different set of model and moment conditions. Empirical likelihood-based Bayesian inference methods are suggested in Kim (2005). This section generalizes the implications of our results to this class of models.

In the case where $\hat{Q}(\beta)$ takes the form of a quadratic norm such as the GMM estimator, it can be similarly shown that when $Q(\beta_0) \neq 0$, or when the GMM model is misspecified, (2.1) continues to hold. On the other hand, when $Q(\beta_0) = 0$, or when the GMM model is correctly specified, $\hat{Q}(\beta_0) - TQ(\beta_0) = \hat{Q}(\beta_0)$ typically converges in distribution to the negative of the quadratic norm of a normal distribution, a special case of which is the χ^2 distribution when an optimal weighting matrix is being used. Regardless $\hat{Q}(\beta_0) - TQ(\beta_0) = O_p(1)$ implies $\hat{Q}(\beta_0) - TQ(\beta_0) = O_p(\sqrt{T})$, therefore the statement that $(Qb) = O_p(\sqrt{T})$ is valid in both cases.

For non-likelihood Q and L , the integral transformations

$$e^{\hat{Q}(\beta)} \pi(\beta) / \int e^{\hat{Q}(\beta')} \pi(\beta') d\beta' \quad \text{and} \quad e^{\hat{L}(\alpha)} \gamma(\alpha) / \int e^{\hat{L}(\alpha')} \gamma(\alpha') d\alpha',$$

can be properly interpreted as distributions. Chernozhukov and Hong (2003) show that estimation can proceed using simulation methods from Bayesian statistics, such as Markov Chain Monte

Carlo methods, using various location measures to identify the parameters of interest. These so-called Laplace-type estimators are defined analogously to Bayesian estimators but use general statistical criterion functions in place of the parametric likelihood function. By considering integral transformations of these statistical functions to give quasi-posterior distributions, this approach provides a useful alternative consistent estimation method to handle intensive computation of many classical extremum estimators such as the GMM estimators of Hansen (1982), Powell (1986)'s censored median regression, nonlinear IV regression such as Berry, Levinsohn, and Pakes (1995) and instrumental quantile regression as in Chernozhukov and Hansen (2005).

The following provides two specific examples of generalized posterior odds ratios constructed from non-likelihood random distance functions. They implicitly deliver a penalization term for the dimension of the information used in the estimation procedure, though all penalize the parametric dimension of the model.

It is worth underscoring that while penalization for parameterization is a natural choice for parametric likelihood models, in general it is not obvious this is necessarily the most desirable form of penalty function outside the likelihood framework. In the context of generalized Bayesian inference, there may be grounds to consider alternative penalty functions that also penalize the dimension of the estimation procedure. For instance, Andrews and Lu (2001) and Hong, Preston, and Shum (2003) consider such penalty functions that involve both the number of parameters and the number of moment conditions. The use of additional moments in GMM and GEL contexts is desirable on efficiency grounds.

Andrews and Lu (2001) proposed consistent model and moment selection criteria for GMM estimation. Interestingly, such selection criteria, which award the addition of moment conditions and penalize the addition of parameters, can not be achieved using a generalized Bayes factor constructed from the GMM objective function:

$$\hat{Q}(\beta) = -Tg_T(\beta)'W_Tg_T(\beta), \quad \text{where} \quad g_T(\beta) = \frac{1}{T} \sum_{t=1}^T m(y_t, \beta),$$

for $m(y_t, \beta)$ a vector of moment conditions and $W_T \xrightarrow{p} W$ positive definite.

With the objective function $\hat{Q}(\beta)$ and associated prior P_Q , the volume

$$P_Q \int e^{\hat{Q}(\beta)} \pi(\beta) d\beta = e^{\hat{Q}(\hat{\beta})} e^{C(\mathcal{A}_{\beta, \hat{\beta}})} T^{-\frac{1}{2} \dim(\beta)} \times e^{o_p(1)},$$

only shrinks at a rate related to the number of parameters and not the number of moment conditions. Generalized Bayes factors using the quadratic GMM objective functions do not encompass

the model selection criteria proposed by Andrews and Lu (2001). This is not surprising given the lack of likelihood interpretation of conventional two-step GMM estimators based on a quadratic norm of the moments.

The recent literature on generalized empirical likelihood (GEL) estimators proposes an information-theoretic alternative to efficient two-step method of moment estimators that minimizes a likelihood distance between the data and the model moments. A GEL estimator is defined as the saddle point of a GEL function,

$$\left(\hat{\beta}, \hat{\lambda}\right) = \arg \max_{\beta \in \mathcal{B}} \arg \min_{\lambda \in \Lambda} \hat{Q}(\beta, \lambda).$$

For example, in the case of exponential tilting

$$\hat{Q}(\beta, \lambda) = \sum_{t=1}^T \rho(\lambda' m(y_t, \beta)) \quad \text{where} \quad \rho(x) = e^x.$$

Given the connection between GEL and parametric likelihood models, it is interesting to ask whether one can define a generalized Bayes factor based on the GELs that mimics the model and moment selection criteria of Andrews and Lu (2001) and others. We define such a GEL Bayes factor as

$$\text{GELBF} = \int \frac{1}{\int e^{-\hat{Q}(\beta, \lambda)} \phi(\lambda) d\lambda} \pi(\beta) d\beta$$

where $\phi(\lambda)$ is a prior density on the lagrange multiplier λ . Intuitively, a large volume of the integral

$$\int e^{-\hat{Q}(\beta, \lambda)} \phi(\lambda) d\lambda$$

indicates that λ tends to be large, and therefore that the GMM model (f, Q) is more likely to be incorrect, or misspecified. Hence, we use its inverse to indicate the strength of the moments involved in the GMM model. The asymptotic equivalence between this GEL Bayes factor and a model and moment selection criteria is given in the following proposition.

Proposition 1 *Suppose assumptions (1), (2) and (3) hold in the parameter space of (β, λ) , assume also that there are interior points in the parameter space λ_0 and β_0 such that $\hat{Q}(\hat{\beta}, \hat{\lambda}) - TQ(\beta_0, \lambda_0) = O_p(1)$, where $Q(\beta, \lambda)$ is the uniform probably limit of $\hat{Q}(\beta, \lambda)/T$. Then*

$$\log \text{GELBF} = \hat{Q}(\hat{\beta}, \hat{\lambda}) + \frac{1}{2} (\dim(\lambda) - \dim(\beta)) \log T + O_p(1).$$

As a consequence of this proposition, asymptotically the GEL Bayes factor puts all weight on models with the best fit $Q(\beta_0, \lambda_0)$, and among models with equal best fit, the one with the largest number of over-identifying moment conditions $\dim(\lambda) - \dim(\beta)$. Model selection behaviors are undetermined among models with equal fit and equal number of overidentifying moment conditions.

7 Conclusion

This paper exploits the connections between Bayesian and classical predictions. For predictions based on minimization of a general class of nonlinear loss functions, we demonstrated conditions under which the asymptotic distribution properties of prediction intervals are not affected by model averaging and the posterior odds place asymptotically unitary probability weight on a single model. This establishes an analogue to the well-known classical result: that asymptotical distribution properties of post-selection estimation and prediction are not affected by first stage model selection so long as the model selection criterion is consistent.

Of course, when confronted with multiple misspecified models it is clear that there is no unique approach to model selection. Importantly, while the number of parameters might appropriately capture model parsimony in the nested case, this may not necessarily be true in the nonnested case. There may be better ways to capture model complexity in this instance. Exploring possible criteria is left for future work.

This paper should not be interpreted as claiming that model selection criteria can be used to the exclusion of standard measures of model fit. For example, there are a number of single model specification tests available — see, among others, Newey (1985), Fan (1994) and Hansen (1982) — which should be applied in assessing model fit. Even the best model chosen by model comparison methods may be rejected by such specification tests indicating that all models are far from the true data generating process (see Sims (2003)). Indeed, as emphasized by Gourieroux and Monfort (1995), hypothesis testing procedures may lead to the simultaneous acceptance or rejection of two nonnested hypothesis. The former may reflect a lack a data while the latter may suggest the testing framework is misspecified.

It is also worth noting Bayesian inference often advocates the use of more general classes of loss function for model evaluation – see Schorfheide (2000) for a recent discussion and promotion of such an approach. For instance, one potential criticism of ranking models based only on statistical fit and parameter parsimony, is that such criterion could well give rise to perverse policy recommendations if the selected model fails to capture important components for the transmission mechanism in the case of monetary policy. However, this will in general be true for any proposed criterion for ranking models, whether decision theoretic or purely statistical, and only serves to emphasize that analysis of this kind is never meant to supersede sound economic reasoning.

What this paper has attempted to treat carefully, given a plausible set of models, finite sample and a particular set of assumptions, is whether models can be statistically distinguished. It remains as further work to analyze the statistical properties of general classes of decision theoretic approaches

to model selection and whether said approaches resolve the inconsistency of posterior odds ratios under our assumptions.

References

- ANDREWS, D. (1994): “Empirical Process Methods in Econometrics,” in *Handbook of Econometrics, Vol. 4*, ed. by R. Engle, and D. McFadden, pp. 2248–2292. North Holland.
- ANDREWS, D., AND B. LU (2001): “Consistent model and moment selection procedures for GMM estimation with application to dynamic panel data models,” *Journal of Econometrics*, 101, 123–164.
- ATTANASIO, O., AND H. LOW (forthcoming): “Estimating Euler Equations,” *Review of Economic Dynamics*.
- BERRY, S., J. LEVINSOHN, AND A. PAKES (1995): “Automobile Prices in Market Equilibrium,” *Econometrica*, 63, 841–890.
- CARROLL, C. D. (2001): “Death to the Log-Linearized Consumption Euler Equation! (And Very Poor Health to the Second-Order Approximation),” *Advances in Macroeconomics*.
- CHERNOZHUKOV, V., AND C. HANSEN (2005): “An IV Model of Quantile Treatment Effects,” *Econometrica*, 73(1), 245–261.
- CHERNOZHUKOV, V., AND H. HONG (2003): “A MCMC Approach to Classical Estimation,” *Journal of Econometrics*, 115(2), 293–346.
- CLARK, T. E., AND M. W. MCCracken (2006): “Combining Forecasts from Nested Models,” working paper, Federal Reserve Bank of Kansas.
- FAN, Y. (1994): “Testing the Goodness of Fit of a Parametric Density Function by Kernel Method,” *Econometric Theory*, 10, 316–356.
- FERNANDEZ-VILLaverDE, J., AND J. F. RUBIO-RAMIREZ (2003): “Estimating Nonlinear Dynamic Equilibrium Economies: Linear versus Nonlinear Likelihood,” unpublished, University of Pennsylvania.
- (2004a): “Comparing Dynamic Equilibrium Models to Data: A Bayesian Approach,” *Journal of Econometrics*, 123, 153–187.
- (2004b): “Estimating Nonlinear Dynamic Equilibrium Economies: A Likelihood Approach,” University of Pennsylvania, PIER Working Paper 04-001.
- GOURIEROUX, C., AND A. MONFORT (1995): *Statistics and Econometric Models*. Cambridge University Press, Cambridge, UK.
- HANSEN, L. (1982): “Large Sample Properties of Generalized Method of Moments Estimators,” *Econometrica*, 50(4), 1029–1054.

- HANSEN, P. (2008): “In-Sample Fit and Out-of-Sample Fit: Their Joint Distribution and Its Implications for Model Selection,” manuscript, Department of Economics, Stanford University.
- HONG, H., B. PRESTON, AND M. SHUM (2003): “Generalized Empirical Likelihood-Based Model Selection Criteria for Moment Condition Models,” *Econometric Theory*, 19, 923–943.
- JUSTINIANO, A., AND B. PRESTON (2004): “Small Open-Economy DSGE Models: Specification, Estimation and Model Fit,” unpublished, Columbia University and International Monetary Fund.
- KIM, J. (2005): “Limited Information Likelihood in Models of Moment Restrictions and Data-based Models/Moments Determination,” manuscript, Seoul National University.
- LUBIK, T. A., AND F. SCHORFHEIDE (2003): “Do Central Banks Respond to Exchange Rate Movements? A Structural Investigation,” unpublished, Johns Hopkins University and University of Pennsylvania.
- NEWBY, W. (1985): “Maximum Likelihood Specification Testing and Conditional Moment Tests,” *Econometrica*, pp. 1047–1070.
- NEWBY, W., AND D. MCFADDEN (1994): “Large Sample Estimation and Hypothesis Testing,” in *Handbook of Econometrics*, Vol. 4, ed. by R. Engle, and D. McFadden, pp. 2113–2241. North Holland.
- PAKES, A., AND D. POLLARD (1989): “Simulation and the Asymptotics of Optimization Estimators,” *Econometrica*, 57(5), 1027–1057.
- PAXSON, C. H., AND S. LUDVIGSON (1999): “Approximation Bias in Euler Equation Estimation,” NBER Working Paper No. T0236.
- POLLARD, D. (1991): “Asymptotics for Least Absolute Deviation Regression Estimator,” *Econometric Theory*, 7, 186–199.
- PÖTSCHER, B. (1991): “Effects of model selection on inference,” *Econometric Theory*, 7(2), 163–185.
- POWELL, J. L. (1986): “Censored Regression Quantiles,” *Journal of Econometrics*, 32, 143–155.
- SCHORFHEIDE, F. (2000): “Loss Function-Based Evaluation of DSGE Models,” *Journal of Applied Econometrics*, 15, 645–670.
- SCHWARZ, G. (1978): “Estimating the dimension of a model,” *The annals of statistics*, 6(2), 461–464.
- SIMS, C. (2003): “Probability Models for Monetary Policy Decisions,” unpublished, Princeton University.
- SIN, C. Y., AND H. WHITE (1996): “Information Criteria for Selecting possibly misspecified parametric models,” *Journal of Econometrics*, 71, 207–225.
- SMETS, F., AND R. WOUTERS (2002): “An Estimated Dynamics Stochastic General Equilibrium Model of the Economy,” National Bank of Belgium, Working Paper No. 35.

- STOCK, J. H., AND M. W. WATSON (2001): “Forecasting Output and Inflation: The Role of Asset Prices,” NBER Working Paper 8180.
- (2003): “Combination Forecasts of Output Growth in a Seven-country Data Set,” unpublished, Princeton University.
- TIMMERMANN, A. (2005): “Forecast Combinations,” in *Handbook of Economic Forecasting*, ed. by C. W. G. Graham Elliott, and A. Timmermann. North Holland.
- VUONG, Q. (1989): “Likelihood-ratio tests for model selection and non-nested hypotheses,” *Econometrica*, pp. 307–333.
- WRIGHT, J. H. (2003): “Bayesian Model Averaging and Exchange Rate Forecasts,” Board of Governors of the Federal Reserve System, International Finance Discussion Papers, No. 779.

A Proof of Theorems 1 and 2

First consider theorem 1. It has been shown (e.g. Pakes and Pollard (1989), Newey and McFadden (1994) and Andrews (1994)) that under assumptions 1, 2 and 3,

$$\sqrt{T} \left(\hat{\beta} - \beta_0 \right) = -\mathcal{A}_\beta^{-1} \frac{\Delta_T}{\sqrt{T}} + o_p(1).$$

Combined with assumption 3, this implies that

$$\hat{Q} \left(\hat{\beta} \right) - \hat{Q} \left(\beta_0 \right) = \sqrt{T} \left(\hat{\beta} - \beta_0 \right)' \frac{\Delta_T}{\sqrt{T}} - \frac{1}{2} \sqrt{T} \left(\hat{\beta} - \beta_0 \right)' (J_T) \sqrt{T} \left(\hat{\beta} - \beta_0 \right) + o_p(1),$$

because of $R_T \left(\hat{\beta} \right) = o_p(1)$, we can write

$$\hat{Q} \left(\hat{\beta} \right) - \hat{Q} \left(\beta_0 \right) = -\frac{1}{2} \frac{\Delta_T'}{\sqrt{T}} \mathcal{A}_\beta^{-1} \frac{\Delta_T}{\sqrt{T}} + o_p(1).$$

The conclusion follows from the assumptions that $\frac{\Delta_T}{\sqrt{T}} = O_p(1)$ and that $-\mathcal{A}_\beta$ is positive definite.

Next consider theorem 2. Define $\hat{\psi}(\bar{y}, Y_T) = \sqrt{T} \left(\hat{\lambda}(\bar{y}, Y_T) - \tilde{\lambda}(\bar{y}, Y_T) \right)$ so that $\hat{\lambda}(\bar{y}, Y_T) = \tilde{\lambda}(\bar{y}, Y_T) + \hat{\psi}(\bar{y}, Y_T) / \sqrt{T}$. By definition, $\hat{\psi}(\bar{y}, Y_T)$ minimizes the following loss function with respect to ψ ,

$$\int \bar{\rho} \left(\bar{y}, \tilde{\lambda}(\bar{y}, Y_T) + \frac{\psi}{\sqrt{T}}; \beta \right) f(\beta | Y_T) d\beta.$$

Also define $h \equiv \sqrt{T} \left(\beta - \hat{\beta} \right)$ as the localized parameter space around $\hat{\beta}$. The implied density for localized parameter h is given by

$$\xi(h) = \left(\frac{1}{\sqrt{T}} \right)^{\dim(\beta)} f \left(\hat{\beta} + \frac{h}{\sqrt{T}} | Y_T \right).$$

Then $\hat{\psi}(\bar{y}, Y_T)$ also minimizes the equivalent loss function of

$$Q_T(\psi) = \int \bar{\rho}\left(\bar{y}, \tilde{\lambda} + \frac{\psi}{\sqrt{T}}; \hat{\beta} + \frac{h}{\sqrt{T}}\right) \xi(h) dh$$

where we are using the shorthand notations $\hat{\lambda} = \hat{\lambda}(\bar{y}, Y_T)$ and $\tilde{\lambda} = \tilde{\lambda}(\bar{y}, Y_T)$. For a given ψ , we are interested in the asymptotic behavior of $Q_T(\psi)$ as $T \rightarrow \infty$. Define

$$\bar{Q}_T(\psi) = \bar{\rho}\left(\bar{y}, \tilde{\lambda} + \frac{\psi}{\sqrt{T}}; \hat{\beta}\right) + \int \bar{\rho}\left(\bar{y}, \tilde{\lambda}; \hat{\beta} + \frac{h}{\sqrt{T}}\right) \xi(h) dh - \bar{\rho}\left(\bar{y}, \tilde{\lambda}; \hat{\beta}\right). \quad (\text{A.8})$$

Essentially, $\bar{Q}_T(\psi)$ is a first order approximation to $Q_T(\psi)$. Under the assumptions stated in Theorem 2, it can be shown that for each ψ ,

$$T [Q_T(\psi) - \bar{Q}_T(\psi)] \xrightarrow{p} 0. \quad (\text{A.9})$$

Because $Q_T(\psi)$ and $\bar{Q}_T(\psi)$ are both convex in ψ , and since $\bar{Q}_T(\psi)$ is uniquely minimized at $\psi \equiv 0$, the convexity lemma (e.g. Pollard (1991)) is used to deliver the desired result that

$$\hat{\psi} = \sqrt{T} \left(\hat{\lambda}(\bar{y}, Y_T) - \tilde{\lambda}(\bar{y}, Y_T) \right) \xrightarrow{p} 0.$$

Proof of (A.9): As $T \rightarrow \infty$, $\xi(h)$ converges in a strong total variation norm in probability to $\phi\left(h; -\mathcal{A}_\beta^{-1}\right)$, the multivariate normal density with mean 0 and variance $-\mathcal{A}_\beta^{-1}$. In fact, the proof of Theorem 3 shows that

$$\int |h|^\alpha \xi(h) dh = O_p(1) \quad \text{and} \quad \int |h|^\alpha |\xi(h) - \phi\left(h; -\mathcal{A}_\beta^{-1}\right)| dh = o_p(1),$$

for all $\alpha \geq 0$. This, combined with the stated assumption that the differentiability of $\bar{\rho}(\bar{Y}, \lambda; \beta)$ with respect to λ and β , will imply the stated convergence in probability. Note that

$$\begin{aligned} T [Q_T(\psi) - \bar{Q}_T(\psi)] &= \int \sqrt{T} \left[\bar{\rho}^{(\beta)}\left(\tilde{\lambda} + \frac{\psi}{\sqrt{T}}, \beta^{**}(h)\right) - \bar{\rho}^{(\beta)}\left(\tilde{\lambda}, \beta^*(h)\right) \right] h \xi(h) dh \\ &= \int \bar{\rho}^{(\beta, \lambda)}(\lambda^{**}(h), \beta^{**}(h)) \psi h \xi(h) dh + \sqrt{T} \int \left[\bar{\rho}^{(\beta)}\left(\hat{\lambda}, \beta^{**}(h)\right) - \bar{\rho}^{(\beta)}\left(\hat{\lambda}, \beta^*(h)\right) \right] h \xi(h) dh. \end{aligned}$$

Because both $\beta^{**}(h)$ and $\beta^*(h)$ are $o_p\left(\frac{h}{\sqrt{T}}\right)$ apart, the second term is clearly of order $o_p(1)$. The first term is up to $o_p(1)$ different from

$$\int \bar{\rho}^{(\beta, \lambda)}\left(\tilde{\lambda}, \hat{\beta}\right) \psi h \xi(h) dh = \int \bar{\rho}^{(\beta, \lambda)}\left(\tilde{\lambda}, \hat{\beta}\right) \psi h \phi\left(h; -\mathcal{A}_\beta^{-1}\right) dh + o_p(1).$$

End of proof of theorem 2. ■

An alternative proof can be based on a standard Taylor expansion of the first order conditions that define $\hat{\lambda}(\bar{y}, Y_T)$ and $\tilde{\lambda}(\bar{y}, Y_T)$, but the notations will be more complicated.

B Proof of Theorem 3

Define $\tilde{\beta}_T = \beta_0 - \frac{1}{T} \mathcal{A}_\beta^{-1} \Delta_T$, and define $h = \sqrt{T} (\beta - \tilde{\beta}_T)$. Then through a change of variables, we can write

$$T^{\frac{\dim(\beta)}{2}} \int e^{\hat{Q}(\beta)} \pi(\beta) d\beta = \int e^{\hat{Q}\left(\frac{h}{\sqrt{T}} + \tilde{\beta}_T\right)} \pi\left(\frac{h}{\sqrt{T}} + \tilde{\beta}_T\right) dh.$$

Chernozhukov and Hong (2003) has shown that (equation A5 of p326) under the same set of assumptions:

$$\begin{aligned} & \int e^{\hat{Q}\left(\frac{h}{\sqrt{T}} + \tilde{\beta}_T\right) - \hat{Q}(\beta_0) + \frac{1}{2T} \Delta_T' \mathcal{A}_\beta^{-1} \Delta_T} \pi\left(\frac{h}{\sqrt{T}} + \tilde{\beta}_T\right) dh \\ & \xrightarrow{p} \pi(\beta_0) (2\pi)^{\frac{\dim(\beta)}{2}} \det(-\mathcal{A}_\beta)^{-1/2}. \end{aligned}$$

Therefore the proof for theorem 3 will be completed if one can show that

$$\hat{Q}(\hat{\beta}) - \left(\hat{Q}(\beta_0) - \frac{1}{2T} \Delta_T' \mathcal{A}_\beta^{-1} \Delta_T \right) \xrightarrow{p} 0,$$

where $\hat{\beta}$ is the conventional M estimator, defined as (see Pakes and Pollard (1989)),

$$\hat{Q}(\hat{\beta}) = \inf_{\beta} \hat{Q}(\beta) + o_p(T^{-1/2}).$$

This is indeed the conclusion of Theorem 1. ■

C Proof of Theorem 4

It is clear from Theorem 3 and its following discussions that under either one of the stated conditions,

$$w_Q \equiv \frac{BF_Q}{BF_Q + BF_L} \xrightarrow{p} 1 \quad \text{as } T \rightarrow \infty.$$

It is because from Theorem 3, for constants

$$C_Q = P_Q \pi(\beta_0) (2\pi)^{\dim(\beta)/2} \det(-\mathcal{A}_\beta)^{-1/2}$$

and

$$C_L = P_L \pi(\alpha_0) (2\pi)^{\dim(\alpha)/2} \det(-\mathcal{A}_\alpha)^{-1/2}$$

we can write

$$\begin{aligned} 1 - w_Q &= \frac{C_L e^{\hat{L}(\hat{\alpha})} T^{-\dim(\alpha)/2} (1 + o_p(1))}{C_Q e^{\hat{Q}(\hat{\beta})} T^{-\dim(\beta)/2} (1 + o_p(1)) + C_L e^{\hat{L}(\hat{\alpha})} T^{-\dim(\alpha)/2} (1 + o_p(1))} \\ &= \frac{(1 + o_p(1)) \frac{C_L}{C_Q} e^{\hat{L}(\hat{\alpha}) - \hat{Q}(\hat{\beta})} T^{\frac{d_\beta - d_\alpha}{2}}}{(1 + o_p(1)) \left(1 + \frac{C_L}{C_Q} e^{\hat{L}(\hat{\alpha}) - \hat{Q}(\hat{\beta})} T^{\frac{d_\beta - d_\alpha}{2}} \right)}. \end{aligned}$$

While $w_Q \xrightarrow{p} 1$ under the stated conditions, the specific rate of convergence depends on the specific condition stated in Theorem 3. Under condition 1, It is clear that $\exists \delta > 0$ such that with probability converging to 1, for all T large enough,

$$1 - w_Q < e^{-T\delta}. \quad (\text{C.10})$$

To show (C.10), note that $\hat{L}(\hat{\alpha}) - \hat{Q}(\hat{\beta}) = T(L(\alpha_0) - Q(\beta_0))(1 + o_p(1))$. Under condition 1, $\exists \delta$ such that $L(\alpha_0) - Q(\beta_0) < -2\delta$. Hence the numerator in (C.10), propotional to $e^{-2T\delta(1+o_p(1))}$, will be smaller than $e^{-T\delta}$ with probability converging to 1. Similar, the denominator converges to 1. Hence (C.10) holds.

On the other hand, when condition 2 holds, it can then be shown that

$$T^{\frac{d_\alpha - d_\beta}{2}} (1 - w_Q) \xrightarrow{d} \frac{C_L}{C_Q} \exp(\bar{\chi}^2), \quad (\text{C.11})$$

where $\bar{\chi}^2$ is typically distributed as the quadratic form of a random vector explained in the following. This is because now that $\hat{L}(\hat{\alpha}) - \hat{Q}(\hat{\beta}) = \hat{L}(\hat{\alpha}) - \hat{L}(\alpha_0) - \hat{Q}(\hat{\beta}) + \hat{Q}(\beta_0)$, it can be written as $T(\hat{\alpha} - \alpha_0)' \mathcal{A}_\alpha (\hat{\alpha} - \alpha_0) - T(\hat{\beta} - \beta_0)' \mathcal{A}_\beta (\hat{\beta} - \beta_0)$. It converges to a quadratic norm of the normal asymptotic distribution of the two sets of parameter estimates. The weighting matrix of the quadratic norm is not necessarily the asymptotic variance of the parameters. Therefore the quadratic norm does not necessarily follow from a chi-square distribution.

Using the definition of $\bar{\rho}(\bar{y}, \lambda; \beta)$, and define $\bar{\rho}(\bar{y}, \lambda; \alpha)$ similarly, we can write $\hat{\lambda}(\bar{y}, Y_T)$ as the minimizer with respect to λ of

$$w_Q \int \bar{\rho}(\bar{y}, \lambda; \beta) f(\beta|Y_T, Q) d\beta + (1 - w_Q) \int \bar{\rho}(\bar{y}, \lambda; \alpha) f(\alpha|Y_T, L) d\alpha.$$

As before, define $\hat{\psi}(\bar{y}, Y_T) = \sqrt{T} \left(\hat{\lambda}(\bar{y}, Y_T) - \tilde{\lambda}(\bar{y}, Y_T) \right)$ and define $h \equiv \sqrt{T} (\beta - \hat{\beta})$. Then $\hat{\psi}(\bar{y}, Y_T)$ equivalently minimizes, with respect to ψ ,

$$Q_T(\psi) = w_Q Q_T^1(\psi) + (1 - w_Q) Q_T^2(\psi)$$

where, with $\tilde{\lambda} \equiv \tilde{\lambda}(\bar{y}, Y_T)$,

$$Q_T^1(\psi) = \int \bar{\rho}\left(\bar{y}, \tilde{\lambda} + \frac{\psi}{\sqrt{T}}; \hat{\beta} + \frac{h}{\sqrt{T}}\right) \xi(h) dh,$$

and

$$Q_T^2(\psi) = \int \bar{\rho}\left(\bar{y}, \tilde{\lambda} + \frac{\psi}{\sqrt{T}}; \alpha\right) f(\alpha|Y_T, L) d\alpha.$$

Now recall the definition of $\bar{Q}_T(\psi)$ in equation (A.8) in the proof of theorem 2. Also define

$$\tilde{Q}_T(\psi) = w_Q \bar{Q}_T(\psi) + (1 - w_Q) \bar{\rho}\left(\bar{y}, \tilde{\lambda}; \hat{\alpha}\right).$$

We are going to show that with this definition

$$T\left(Q_T(\psi) - \tilde{Q}_T(\psi)\right) \xrightarrow{p} 0. \quad (\text{C.12})$$

If (C.12) holds, it then follows again from the convexity lemma of Pollard (1991) and the fact that $\tilde{Q}_T(\psi)$ is uniquely optimized at $\psi = 0$ that

$$\hat{\psi}(\bar{y}, Y_T) = \sqrt{T}\left(\hat{\lambda}(\bar{y}, Y_T) - \tilde{\lambda}(\bar{y}, Y_T)\right) \xrightarrow{p} 0.$$

Finally, we will verify (C.12). With the definition of $\tilde{Q}_T(\psi)$, we can write

$$T\left(Q_T(\psi) - \tilde{Q}_T(\psi)\right) = Tw_Q\left(Q_T^1(\psi) - \bar{Q}_T(\psi)\right) + T(1 - w_Q)\left(Q_T^2(\psi) - \bar{\rho}\left(\bar{y}, \tilde{\lambda}; \hat{\alpha}\right)\right).$$

Because $w_Q \xrightarrow{p} 1$, it follows from the proof of Theorem 2 that $w_Q T\left(Q_T^1(\psi) - \bar{Q}_T(\psi)\right) \xrightarrow{p} 0$. As the sample size increases, $f(\alpha|Y_T, L)$ tends to concentrate on $\hat{\alpha}$, therefore it can also be shown that

$$Q_T^2(\psi) - \bar{\rho}\left(\bar{y}, \tilde{\lambda}; \hat{\alpha}\right) \xrightarrow{p} 0. \quad (\text{C.13})$$

To see (C.13), note that consistency of the Bayes posterior distribution is automatically implied by the asymptotic normality of the localized posterior distribution in Theorem 3. Therefore for any $\rho > 0$, $\delta > 0$ and $\epsilon > 0$,

$$\lim_{T \rightarrow \infty} P\left(\int_{|\alpha - \hat{\alpha}| > \delta} \alpha^\rho f(\alpha|Y_T, L) d\alpha > \epsilon\right) = 0.$$

Then the difference in (C.13) can be decomposed into the sum of

$$\int_{|\alpha - \hat{\alpha}| < \delta} \left(\bar{\rho}\left(\bar{y}, \tilde{\lambda} + \frac{\psi}{\sqrt{T}}; \alpha\right) - \bar{\rho}\left(\bar{y}, \tilde{\lambda}; \hat{\alpha}\right)\right) f(\alpha|Y_T, L) d\alpha$$

and

$$\int_{|\alpha - \hat{\alpha}| \geq \delta} \left(\bar{\rho} \left(\bar{y}, \tilde{\lambda} + \frac{\psi}{\sqrt{T}}; \alpha \right) - \bar{\rho} \left(\bar{y}, \tilde{\lambda}, \hat{\alpha} \right) \right) f(\alpha | Y_T, L) d\alpha.$$

Both terms can be made as small as desired probabilistically by choosing small values of δ .

Now if either condition 1 holds or if condition 2 holds (in which case $d_\alpha - d_\beta \geq 2$), then because of (C.10) and (C.11), $T(1 - w_Q) = O_p(1)$. Combine this with (C.13), we see that it is always the case that $T(1 - w_Q) \left(Q_T^2(\psi) - \bar{\rho} \left(\bar{y}, \tilde{\lambda}; \hat{\alpha} \right) \right) = o_p(1)$. Therefore (C.12) holds. ■

Remark: It also follows from the same arguments as above that the results of Theorem 4 does not hold when $d_\alpha = d_\beta + 1$. In fact, in this case, we can redefine

$$\tilde{Q}_T(\psi) = w_Q \bar{Q}_T(\psi) + (1 - w_Q) \left[\bar{\rho} \left(\bar{y}, \tilde{\lambda}; \hat{\alpha} \right) + \frac{\partial}{\partial \lambda} \bar{\rho} \left(\bar{y}, \tilde{\lambda}; \hat{\alpha} \right) \frac{\psi}{\sqrt{T}} \right].$$

We can then follow the same logic as before to show that

$$T \left(Q_T(\psi) - \tilde{Q}_T(\psi) \right) \xrightarrow{p} 0.$$

Note that in the definition of $\bar{Q}_T(\psi)$ in equation (A.8) of Theorem 2, using a second order Taylor expansion we can replace $\bar{\rho} \left(\bar{y}, \tilde{\lambda} + \frac{\psi}{\sqrt{T}}; \hat{\beta} \right)$ by

$$\frac{1}{T} \psi' \frac{1}{2} \bar{\rho}_{\lambda\lambda} \left(\bar{y}, \tilde{\lambda}; \hat{\beta} \right) \psi + o_p \left(\frac{1}{T} \right).$$

As such we can write

$$\begin{aligned} & T \left(\tilde{Q}_T(\psi) - w_Q \left(\int \bar{\rho} \left(\bar{y}, \tilde{\lambda}; \hat{\beta} + \frac{h}{\sqrt{T}} \right) \xi(h) dh - \bar{\rho} \left(\bar{y}, \tilde{\lambda}; \hat{\beta} \right) \right) - (1 - w_Q) \bar{\rho} \left(\bar{y}, \tilde{\lambda}; \hat{\alpha} \right) \right) \\ &= \frac{1}{2} \psi' \bar{\rho}_{\lambda\lambda} \left(\tilde{\lambda}; \hat{\beta} \right) \psi + \sqrt{T} (1 - w_Q) \frac{\partial}{\partial \lambda} \bar{\rho} \left(\bar{y}, \tilde{\lambda}; \hat{\alpha} \right) \psi + o_p(1). \end{aligned}$$

It follows from both (C.11) and the convergence of $\frac{\partial}{\partial \lambda} \bar{\rho} \left(\bar{y}, \tilde{\lambda}; \hat{\alpha} \right) \xrightarrow{p} \frac{\partial}{\partial \lambda} \bar{\rho}(\bar{y}, \lambda_0; \alpha_0)$ that

$$\sqrt{T} (1 - w_Q) \frac{\partial}{\partial \lambda} \bar{\rho} \left(\bar{y}, \tilde{\lambda}; \hat{\alpha} \right) \xrightarrow{d} \frac{C_L}{C_Q} \exp(\bar{\chi}^2) \frac{\partial}{\partial \lambda} \bar{\rho}(\bar{y}, \lambda_0; \alpha_0)'$$

where $\lambda_0 = \arg \min_{\lambda} \bar{\rho}(\bar{y}, \lambda; \beta_0)$. Hence again with convexity arguments for uniform convergence we can show that

$$\begin{aligned} \hat{\psi}(\bar{y}, Y_T) &\xrightarrow{d} \arg \min_{\psi} \frac{1}{2} \psi' \bar{\rho}_{\lambda\lambda}(\lambda_0; \beta_0) \psi + \frac{C_L}{C_Q} \exp(\bar{\chi}^2) \frac{\partial}{\partial \lambda} \bar{\rho}(\bar{y}, \lambda_0; \alpha_0)' \psi \\ &= -\bar{\rho}_{\lambda\lambda}(\lambda_0; \beta_0)^{-1} \frac{C_L}{C_Q} \exp(\bar{\chi}^2) \frac{\partial}{\partial \lambda} \bar{\rho}(\bar{y}, \lambda_0; \alpha_0)'. \end{aligned}$$

In other words, if $d_\beta = d_\alpha - 1$, $\sqrt{T} \left(\hat{\lambda}(\bar{y}, Y_T) - \tilde{\lambda}(\bar{y}, Y_T) \right)$ converges in distribution to a non degenerate random variable.

D Proof of Proposition 1

We will only outline the main steps because the arguments essentially mirror those of Theorem 3. Define $\hat{\lambda}(\beta) = \arg \min_{\lambda} \hat{Q}(\beta, \lambda)$. It can be shown that for some constant C_1 ,

$$\int e^{-\hat{Q}(\beta, \lambda)} \phi(\lambda) d\lambda = e^{-\hat{Q}(\beta, \hat{\lambda}(\beta))} \sqrt{T}^{-\dim(\lambda)} C_1 (1 + o_p(1)).$$

Then we can write

$$GELBF = \sqrt{T}^{+\dim(\lambda)} \frac{1}{C_1} (1 + o_p(1)) \int e^{\hat{Q}(\beta, \hat{\lambda}(\beta))} \pi(\beta) d\beta.$$

It can also be shown that for another constant C_2 :

$$\int e^{\hat{Q}(\beta, \hat{\lambda}(\beta))} \pi(\beta) d\beta = C_2 e^{\hat{Q}(\hat{\beta}, \hat{\lambda})} T^{-\frac{\dim(\beta)}{2}} (1 + o_p(1)).$$

The result then follows from combining the previous relations. ■