

Consistency of Bayesian Linear Model Selection With a Growing Number of Parameters

Zuofeng Shang and Murray K. Clayton

*Department of Statistics
University of Wisconsin, Madison
Madison, WI 53706*

Abstract

Linear models with a growing number of parameters have been widely used in modern statistics. One important problem about this kind of model is the variable selection issue. Bayesian approaches, which provide a stochastic search of informative variables, have gained popularity. In this paper, we will study the asymptotic properties related to Bayesian model selection when the model dimension p is growing with the sample size n . We consider $p \leq n$ and provide sufficient conditions under which: (1) with large probability, the posterior probability of the true model (from which samples are drawn) uniformly dominates the posterior probability of any incorrect models; and (2) with large probability, the posterior probability of the true model converges to one. Both (1) and (2) guarantee that the true model will be selected under a Bayesian framework. We also demonstrate several situations when (1) holds but (2) fails, which illustrates the difference between these two properties. Simulated examples are provided to illustrate the main results.

Keywords: Bayesian model selection; growing number of parameters; Posterior model consistency; consistency of Bayes factor; consistency of posterior odds ratio; Gibbs sampling.

1. Introduction

This work was motivated by efforts to analyze remotely sensed (satellite) data which consists of multiple spatial images. In the setting of interest, one image corresponds to a “response” while others correspond to covariates. To find the relationship between the response and covariate spatial images,

Zhang *et al.* (2010) proposed a functional concurrent linear model with varying coefficients and applied a wavelet approach to transform this model into a linear model (with a particular design matrix) which contains an n -vector of responses and a sparse p -vector of wavelet coefficients. Since the images contain thousands of pixels, the model dimension p , which is determined by the maximum decomposition level in the wavelet expansion, has to be large so that sufficiently fine details in the target images can be captured. On the other hand, p has an upper bound $p \leq (K+1)n$, where K is the total number of covariate images involved in the model. This is because each spatial image corresponds to a vector of wavelet coefficients which has dimension not exceeding n , and there are $K+1$ images in total with one of them representing the intercept and others the slopes. An important question is how to select the nonzero coefficients in the model, which is essentially a variable selection problem. Zhang *et al.* (2010) adopted a Lasso approach to address this.

The problem they handle relies on a specific design matrix induced by the wavelet structure. It is of interest, to frame the variable selection problem more broadly. More precisely, we suppose that data are drawn from the linear model

$$\mathbf{y} = X\beta + \epsilon, \tag{1.1}$$

where $\epsilon \sim N(\mathbf{0}, \sigma_0^2 I_n)$ is an n -vector of errors, $\mathbf{y} = (y_1, \dots, y_n)^T$ is an n -vector of responses, $\beta = (\beta_1, \dots, \beta_p)^T$ is a p -vector of parameters and $X = (X_1, \dots, X_p)$ is a $n \times p$ design matrix with X_j the j th column of X . It is also assumed that only a subset of X_1, \dots, X_p contribute to \mathbf{y} and we are interested in selecting the variables in this subset.

We consider a Bayesian variable selection (BVS) approach based on model (1.1). The Bayesian model to be considered is a variation of George and McCulloch (1993) and has been studied by Clyde *et al.* (1998), Clyde and George (2000), and Wolfe *et al.* (2004). Clearly, each subset of X_1, \dots, X_p defines a candidate model, so there are 2^p of them in total. According to George and McCulloch (1993), all the marginal posterior probabilities of these 2^p models can be calculated and the model with the largest posterior probability can be selected as the “best” model. This motivates the formal definition of posterior model consistency (PMC). We say that PMC holds if the true model, defined as the model from which samples are drawn, has a posterior probability approaching one. Since the sum of the posterior probabilities of all models equals one, when PMC holds, the posterior probability of any incorrect model will go to zero when n goes to infinity so that the true

model can be correctly selected.

PMC has been theoretically verified when p is fixed (see Fernández *et al.*, 2001; Moreno and Girón, 2005; Liang *et al.*, 2008; Casella *et al.*, 2009). However, fewer results have been derived when p is growing with n , an interesting and important scenario. For increasing p , Berger *et al.* (2003), Moreno *et al.* (2010) and Girón *et al.* (2010) proved consistency for Bayes factors. Although PMC and consistency of Bayes factors are equivalent for fixed p (see Liang *et al.*, 2008; Casella *et al.*, 2009), they are different for growing p . Actually, we will see below that consistency of the Bayes factor is equivalent to consistency of the posterior odds ratio under a general setting, but that the latter form of consistency is weaker than PMC. Therefore, it seems valuable to separately study PMC.

In this paper we will consider two classes of design matrix X , both with $p \leq n$, although our results can be generalized to $p \gg n$ when combined with certain dimension reduction approaches. In the first case, X is quite general. A representative situation is that the eigenvalues of $X^T X/n$ are uniformly bounded both above and below. Consistency is examined when p grows slower than n , say, $p \log n = o(n)$. We find that the posterior odds in favor of any incorrect model uniformly converges to zero, and the posterior probability of the true model converges to one. A second case we consider occurs when $X^T X/n$ is the identity matrix, i.e., $X^T X = nI_p$, and p grows as fast as n , say $p = n$. In that case, consistency of the posterior odds ratio and PMC are examined, i.e., the posterior odds ratio in favor of any incorrect model uniformly converges to zero, and the posterior probability of the true model converges to one. We also demonstrate how consistency of the posterior odds ratio can hold even though PMC fails.

The remainder of this paper is organized as follows. In Section 2, preliminaries and main results will be provided. In Section 3, a numerical example related to the results of Section 2 is displayed. Section 4 contains the conclusion. And technical arguments are included in Section 5.

2. Preliminaries and main results

Suppose the n dimensional response vector $\mathbf{y} = (y_1, \dots, y_n)^T$ and the n by p covariate matrix $X = (X_1, \dots, X_p)$ are linked by the model

$$\mathbf{y} = X\beta + \epsilon, \tag{2.1}$$

where the X_j 's are n -vectors, $\beta = (\beta_1, \dots, \beta_p)^T$ is an unknown p -vector and ϵ is a vector of random errors. Here, X is allowed to be either (1) random but independent of ϵ or (2) deterministic. For $1 \leq j \leq p$, define the state variable of β_j by $\gamma_j = I(\beta_j \neq 0)$ and $\gamma = (\gamma_1, \dots, \gamma_p)^T$, where $I(\cdot)$ is the indicator function. We call γ the state vector of β and denote the number of 1's in γ by $|\gamma|$. The state vector γ completely determines the inclusion or exclusion of β_j 's in model (2.1), and therefore, can define a model $\mathbf{y} = X_\gamma \beta_\gamma + \epsilon$, where X_γ is an $n \times |\gamma|$ submatrix of X whose columns are indexed by the nonzero components of γ , and β_γ is the subvector (with size $|\gamma|$) of β indexed by the nonzero components of γ . It is natural, therefore, to call each γ a model. Note that there are 2^p such γ 's representing 2^p different models. For any state vectors γ and γ' , let $(\gamma \setminus \gamma')_j = I(\gamma_j = 1, \gamma'_j = 0)$ denote the difference (which is also a state vector) between γ and γ' , i.e., the 0-1 vector indicating the variables that are present in γ but absent in γ' . We say that γ is nested in γ' (denoted by $\gamma \subset \gamma'$) if $\gamma \setminus \gamma' = 0$. Denote the true model coefficient vector by β^0 and the corresponding state vector by γ^0 , and let $s_n = |\gamma^0|$ denote the size of the true model.

In this paper we consider the following hierarchical Bayesian model which is a variation of the model used by George and McCulloch (1993)

$$\begin{aligned} \mathbf{y}|\beta, \sigma^2 &\sim N(X\beta, \sigma^2 I_n), \\ \beta_j|\gamma_j, \sigma^2 &\sim (1 - \gamma_j)\delta_0 + \gamma_j N(0, c_j \sigma^2), \\ 1/\sigma^2 &\sim \chi_\nu^2, \\ \gamma &\sim p(\gamma), \end{aligned} \tag{2.2}$$

where δ_0 is point mass measure concentrated at zero. Hereafter, ν will be fixed a priori. Let $\Sigma = \text{diag}(c)$ with $c = (c_j)_{1 \leq j \leq p}$ a p -vector of positive components, and let Σ_γ be the $|\gamma| \times |\gamma|$ sub-diagonal matrix of Σ corresponding to γ . Let $Z = (\mathbf{y}, X)$ denote the full data set. It follows by integrating out β and σ that the posterior distribution of γ is given by

$$p(\gamma|Z) \propto (2\pi)^{-n/2} \det(W_\gamma)^{-1/2} p(\gamma) \left\{ \frac{2}{1 + \mathbf{y}^T (I_n - X_\gamma U_\gamma^{-1} X_\gamma^T) \mathbf{y}} \right\}^{(n+\nu)/2}, \tag{2.3}$$

where $U_\gamma = \Sigma_\gamma^{-1} + X_\gamma^T X_\gamma$ and $W_\gamma = \Sigma_\gamma^{1/2} U_\gamma \Sigma_\gamma^{1/2}$. In particular, if $\gamma = \emptyset$ (the null model containing no covariate variables), (2.3) still holds if we adopt the conventions that $X_\emptyset = 0$ and $\Sigma_\emptyset = U_\emptyset = W_\emptyset = 1$.

Define $S_1 = \{\gamma | \gamma^0 \subset \gamma, \gamma \neq \gamma^0\}$ and $S_2 = \{\gamma | \gamma^0 \text{ is not nested in } \gamma\}$. It is clear that $S(n)$ defined by $S(n) = S_1 \cup S_2 \cup \{\gamma^0\}$ is the class of all state vectors. In particular, when $\gamma^0 = \emptyset$, S_2 is empty, and hence S_1 is the class of all state vectors excluding γ^0 . As was found by Liang *et al.* (2008), we will see later in this section that whether γ^0 is null or nonnull will result in some differences in the main results (especially in the assumptions that are needed to establish our main results); thus, we will treat these cases separately. When γ^0 is nonnull, we denote $\varphi_{\min}(n) = \min_{\gamma \in S_2} \lambda_- \left(\frac{1}{n} X_{\gamma^0 \setminus \gamma}^T (I_n - P_\gamma) X_{\gamma^0 \setminus \gamma} \right)$ and $\varphi_{\max}(n) = \max_{\gamma \in S_2} \lambda_+ \left(\frac{1}{n} X_{\gamma^0 \setminus \gamma}^T X_{\gamma^0 \setminus \gamma} \right)$, where $P_\gamma = X_\gamma (X_\gamma^T X_\gamma)^{-1} X_\gamma^T$ is a projection matrix, $\lambda_-(A)$ and $\lambda_+(A)$ are the minimal and maximal eigenvalues of the square matrix A . We also adopt the convention that $P_\emptyset = 0$. For the case that $\gamma^0 = \emptyset$, both φ_{\min} and φ_{\max} are meaningless, and S_1 will be focused on in this situation.

Before proceeding further, we introduce several types of consistency central to this work. Generally speaking, to make a correct model selection

$$\max_{\gamma \neq \gamma^0} p(\gamma|Z)/p(\gamma^0|Z) \rightarrow 0 \quad (2.4)$$

should hold as $n \rightarrow \infty$, which means that the posterior probability of the true model asymptotically dominates that of any incorrect model. Following a framework similar to that of Zellner (1978), the term $p(\gamma|Z)/p(\gamma^0|Z)$, which is called the posterior odds ratio in favor of γ , satisfies the relationship

$$p(\gamma|Z)/p(\gamma^0|Z) = BF(\gamma : \gamma^0) \frac{p(\gamma)}{p(\gamma^0)}, \quad (2.5)$$

where $BF(\gamma : \gamma^0) := p(Z|\gamma)/p(Z|\gamma^0)$ is the Bayes factor of γ versus γ^0 and $p(\gamma)/p(\gamma^0)$ is the prior odds ratio in favor of γ . The Bayes factor is consistent if for any $\gamma \neq \gamma^0$, $BF(\gamma : \gamma^0) \rightarrow 0$. The posterior odds ratio is consistent if for any $\gamma \neq \gamma^0$, $p(\gamma|Z)/p(\gamma^0|Z) \rightarrow 0$. It is easy to see that property (2.4) implies consistency of the posterior odds ratio. We say that posterior model consistency (PMC) holds if $p(\gamma^0|Z) \rightarrow 1$. These types of consistency all have been useful in Bayesian model selection. Representative references include (1) assessment of posterior odds ratio: Jeffreys (1967), Zellner (1971, 1978); (2) performance of Bayes factor: Berger and Pericchi (1996), Moreno *et al.* (1998, 2010), Casella *et al.* (2009); (3) PMC: Fernández *et al.* (2001), Liang *et al.* (2008).

It is easy to see that when

$$c_1^{-1} \leq \min_{\gamma} p(\gamma)/p(\gamma^0) \leq \max_{\gamma} p(\gamma)/p(\gamma^0) \leq c_1 \quad (2.6)$$

holds for some positive constant c_1 , consistency of the Bayes factor is equivalent to consistency of the posterior odds ratio, and that both are weaker than (2.4). A special case is that $p(\gamma) = 2^{-p}$ for all γ 's, which results in an indifference prior distribution for γ , see, e.g., Smith and Kohn (1996).

To illustrate the relationship between PMC and (2.4), note that

$$p(\gamma^0|Z) = \frac{1}{1 + \sum_{\gamma \neq \gamma^0} p(\gamma|Z)/p(\gamma^0|Z)}, \quad (2.7)$$

and thus $p(\gamma^0|Z) \rightarrow 1$ will imply (2.4). When p is fixed, it has been noted by Liang *et al.* (2008) that (2.4) implies PMC. However, when p grows with n , it will be shown later that this may not be true. This somewhat illustrates the difference between PMC and (2.4).

In what follows, we introduce some regularity conditions that are useful to establish our main results. We will also demonstrate some particular situations when these conditions are satisfied.

Assumption 2.1. There exists a constant $C_0 > 0$ such that for any n , $\max_{\gamma \in S(n)} p(\gamma)/p(\gamma^0) \leq C_0$.

Assumption 2.2. There exist positive constants C_1, C_2 such that with probability equal to one, $\liminf_n \varphi_{\min}(n) \geq C_1$ and $\limsup_n \varphi_{\max}(n) \leq C_2$.

Assumption 2.3. There exists a positive sequence ψ_n such that $\min_{j \in \gamma^0} |\beta_j^0| \geq \psi_n$ and, as $n \rightarrow \infty$, $\psi_n \sqrt{n} \rightarrow \infty$.

Assumption 2.4. $p_n \rightarrow \infty$, $s_n \leq p_n \leq n$ and $p_n \log n = o(n \log(1 + \min\{\psi_n^2, 1\}))$.

Assumption 2.5. $p_n \rightarrow \infty$, $s_n \leq p_n \leq n$ and $p_n \log p_n = o(n)$.

Hereafter, unless otherwise explicitly stated, we will drop the subscript from p_n .

Assumption 2.6. $\bar{\phi}_n = O(n^{\delta_0})$ for some $\delta_0 > 0$, where $\bar{\phi}_n = \max_{1 \leq j \leq p} c_j$.

Assumption 2.7. $k_n = O(\underline{\phi}_n)$, where $k_n = \|\beta_{\gamma^0}^0\|_2^2$ and $\underline{\phi}_n = \min_{1 \leq j \leq p} c_j$.

Assumption 2.8. There exist $C_3 > 0$ and $\delta \geq 0$ such that $n^{1-\delta} \underline{\phi}_n \rightarrow \infty$, and for any n , with probability equal to one,

$$\inf_{\gamma \in S_1} \lambda_- \left(\frac{1}{n} X'_{\gamma \setminus \gamma^0} (I_n - P_{\gamma^0}) X_{\gamma \setminus \gamma^0} \right) \geq C_3 n^{-\delta}. \quad (2.8)$$

Remark 2.1.

- (a). Assumption 2.1 is satisfied by some commonly used priors $p(\gamma)$, such as the flat prior $p(\gamma) = 2^{-p}$ (Smith and Kohn, 1996). More generally, if $p(\gamma_j = 1) = \theta_j$ is such that both $\prod_{j \in \gamma \setminus \gamma^0} \left(\frac{\theta_j}{1-\theta_j} \right)$ and $\prod_{j \in \gamma^0 \setminus \gamma} \left(\frac{1-\theta_j}{\theta_j} \right)$ are bounded, then Assumption 2.1 is satisfied.
- (b). We use Assumption 2.3 to prove consistency for a growing p . Fan and Peng (2004) introduced a similar assumption in the framework of smoothly clipped absolute deviation (SCAD) penalized optimization where \sqrt{n} in Assumption 2.3 was replaced by $1/\lambda_n$ with λ_n the penalty parameter. This condition requires the true parameters to be away from zero. Otherwise, it is impossible to distinguish between zero and nonzero parameters.
- (c). Assumptions 2.4 and 2.5 define a rate on the dimension p . In particular, when $\inf_n \psi_n > 0$, Assumption 2.4 is satisfied if $s_n \leq p$ and $p \log n = o(n)$. The results hold when s_n is either bounded or growing with n .
- (d). Assumption 2.6 excludes the possibility that $\bar{\phi}_n$ is extremely large, e.g., we exclude the situation that $\bar{\phi}_n = \exp(n^\omega)$ for some $\omega > 0$. Assumption 2.7 requires that $\underline{\phi}_n$ is not growing slower than $k_n = \|\beta_{\gamma^0}^0\|_2^2$. When the design matrix X is nonorthogonal, we use this assumption to facilitate the proof of consistency (see Theorem 2.2 below). But when X is orthogonal, this assumption is redundant and can be removed (see Corollary 2.5 below). \square

Assumptions 2.1, 2.3–2.7 are easily satisfied. The following proposition demonstrates that a broad class of design matrices X can satisfy Assumptions 2.2 and 2.8.

Proposition 2.1. If the $n \times p$ matrix X satisfies $\lambda_- \left(\frac{1}{n} X^T X \right) \geq c$, where $c > 0$ is constant, then for any $\gamma \subset \bar{\gamma}$ and $\gamma \neq \bar{\gamma}$,

$$\lambda_- \left(\frac{1}{n} X_{\bar{\gamma} \setminus \gamma}^T (I_n - P_\gamma) X_{\bar{\gamma} \setminus \gamma} \right) \geq c. \quad (2.9)$$

The proof of Proposition 2.1 can be found in Section 5 (Appendix).

Remark 2.2. Proposition 2.1 demonstrates that Assumptions 2.2 and 2.8 can hold under general classes of design matrices. One such class consists of matrices X satisfying

$$1/c_2 \leq \lambda_- \left(\frac{1}{n} X^T X \right) \leq \lambda_+ \left(\frac{1}{n} X^T X \right) \leq c_2, \quad (2.10)$$

where c_2 is some positive constant. For any $\gamma \in S_1$, we will have that $\gamma^0 \subset \gamma$ and $\gamma^0 \neq \gamma$. Thus, by Proposition 2.1, $\lambda_- \left(\frac{1}{n} X'_{\gamma \setminus \gamma^0} (I_n - P_{\gamma^0}) X_{\gamma \setminus \gamma^0} \right) \geq 1/c_2$, i.e., inequality (2.8) in Assumption 2.8 holds. Notice that when $\gamma \in S_2$, the relationship $\gamma \subset \gamma^0 \vee \gamma$ and $\gamma \neq \gamma^0 \vee \gamma$ holds, where $\gamma^0 \vee \gamma$ denotes the p -vector with j th component the larger of $(\gamma^0)_j$ and γ_j , then Assumption 2.2 follows by applying Proposition 2.1. \square

In the following text, we assume that data are generated from the true model $\mathbf{y} = X\beta^0 + \epsilon$ with $\epsilon \sim N(0, \sigma_0^2 I_n)$. Let γ^0 be the p -dimensional state vector corresponding to β^0 . Unless otherwise stated, the limits in our main results will be taken when $n \rightarrow \infty$.

Theorem 2.2. Suppose that γ^0 is nonnull and Assumptions 2.1–2.4, 2.6–2.8 are satisfied. Let $\delta \geq 0$ satisfy Assumption 2.8. If $p^{\alpha_0} = o(n^{1-\delta} \underline{\phi}_n)$ for some $\alpha_0 > 2$, then $\max_{\gamma \neq \gamma^0} p(\gamma|Z)/p(\gamma^0|Z) \rightarrow_p 0$. If $p^{\alpha_0+2} = o(n^{1-\delta} \underline{\phi}_n)$ for some $\alpha_0 > 2$, then $\sum_{\gamma \neq \gamma^0} p(\gamma|Z) \rightarrow_p 0$, and consequently, $p(\gamma^0|Z) \rightarrow_p 1$.

The proof of Theorem 2.2 follows by first deriving asymptotic approximations of the posterior odds ratios $p(\gamma|Z)/p(\gamma^0|Z)$ for any $\gamma \neq \gamma^0$, and then using these approximations to show that $\sum_{\gamma \neq \gamma^0} p(\gamma|Z)/p(\gamma^0|Z) \rightarrow_p 0$. The limit $p(\gamma^0|Z) \rightarrow_p 1$ thus immediately follows from (2.7). Details are in the Appendix.

Remark 2.3. Theorem 2.2 provides sufficient conditions under which (2.4) and PMC are satisfied. It asserts that, with large probability, $p(\gamma^0|Z)$

uniformly dominates $p(\gamma|Z)$ for any $\gamma \neq \gamma^0$, and with large probability, $p(\gamma^0|Z)$ approaches one. Thus, with large probability, the true model γ^0 can be selected from a Bayesian perspective. \square

Remark 2.4. A natural but interesting question is that, when the growth rate for p changes, should there be any change in the hyperparameters c_j 's so that property (2.4) or PMC still holds? Theorem 2.2 partly and heuristically answers this question. To see this, let us consider a special case that $p = n^a$ with $a \in (0, 1]$. Here the factor a controls how fast p grows. For instance, a larger a corresponds to a faster growth rate. By Theorem 2.2, to satisfy (2.4), one sufficient condition is $p^{\alpha_0} = o(n^{1-\delta}\underline{\phi}_n)$ for some fixed $\alpha_0 > 2$, which, in the special case of interest, becomes

$$n^{a\alpha_0+\delta-1} = o(\underline{\phi}_n). \quad (2.11)$$

The interpretation of (2.11) is that $\underline{\phi}_n$ controls the growth of $n^{a\alpha_0+\delta-1}$, and thus, $n^{a\alpha_0+\delta-1}$ serves as a lower bound for $\underline{\phi}_n$. When a increases (which corresponds to p growing faster), this lower bound should become larger. For instance, $a = 0.5$ corresponds to a lower bound $n^{0.5\alpha_0+\delta-1}$; however, when $a = 0.75$, this lower bound has increased to $n^{0.75\alpha_0+\delta-1}$. This heuristically shows that, in order to satisfy property (2.4) or PMC, the lower bound for $\underline{\phi}_n$ should generally increase when p grows faster. \square

Remark 2.5. When combined with certain dimension reduction techniques such as sure independence screening (SIS) proposed by Fan and Lv (2008), one can generalize Theorem 2.2 to the ultra-high dimensional setting, i.e., $p \gg n$. This framework has been explored by many authors from non-Bayesian perspectives (see, e.g., Meinshausen and Bühlmann, 2006; Meinshausen and Yu, 2009; Zhang and Huang, 2010; Bühlmann and Kalisch, 2010). Here, we explore it by a Bayesian way. The basic idea is to first reduce the high-dimensional linear model so that the model dimension is below n , and then apply Bayesian model (2.2) to this reduced linear model. Under suitable conditions and using the arguments similar to the proof of Theorem 2.2, one can show that the posterior probability of the true model based on the reduced linear model converges in probability to 1. We refer to Supplement A for the description of this result and details of the proof. \square

The following result is an application of Theorem 2.2 in a special setting, which allows the growth rate of p to be $p \log n = o(n)$.

Corollary 2.3. Suppose that γ^0 is nonnull and Assumptions 2.1, 2.2 and inequality (2.8) are satisfied. Assume that $\min_{j \in \gamma^0} |\beta_j^0| \geq \psi_n$ with $\inf_n \psi_n > 0$, and

p satisfies $p \log n = o(n)$. Suppose there exists a constant δ_0 with $\delta_0 > 3 + \delta$ such that $k_n = O(n^{\delta_0})$, where $\delta \geq 0$ is specified in inequality (2.8). Then with the selection $\bar{\phi}_n = O(n^{\delta_0})$ and $n^{\delta_0} = O(\underline{\phi}_n)$, we have $p(\gamma^0|Z) \rightarrow_p 1$.

The proof of Corollary 2.3 can be finished by choosing $\alpha_0 \in (2, \delta_0 - \delta - 1)$ and verifying the assumptions in Theorem 2.2.

Theorem 2.2 deals with the case when the true model is nonnull. If the true model is null, then the response vector \mathbf{y} will have a zero mean. The corresponding result is summarized below.

Theorem 2.4. Suppose γ^0 is null, i.e., $\mathbf{y} = \epsilon \sim N(0, \sigma_0^2 I_n)$, and that Assumptions 2.1, 2.5, and 2.8 are satisfied. If $p^{\alpha_0} = o(n^{1-\delta} \underline{\phi}_n)$ for some $\alpha_0 > 2$, then $\max_{\gamma \neq \gamma^0} p(\gamma|Z)/p(\gamma^0|Z) \rightarrow_p 0$. If $p^{\alpha_0+2} = o(n^{1-\delta} \underline{\phi}_n)$ for some $\alpha_0 > 2$, then $\sum_{\gamma \neq \gamma^0} p(\gamma|Z) \rightarrow_p 0$, and consequently, $p(\gamma^0|Z) \rightarrow_p 1$.

The proof of Theorem 2.4 is similar to Theorem 2.2 and can be found in Appendix.

Remark 2.6. Liang *et al.* (2008) applied mixture of g-priors in their Bayesian model, which is different from the priors used in this work, and obtained PMC in the case that p is fixed. However, since their model induces a non-analytical expression for $p(\gamma|Z)$, when p is growing with n , the theoretical derivation would become complicated. To overcome this difficulty, we consider conjugate priors in Bayesian model (2.2), which induces an analytical expression for $p(\gamma|Z)$. Under this framework, the derivation of both (2.4) and PMC becomes easier. \square

Although it is valid for a general type of design matrix, Theorem 2.2 requires that p grows slower than n . More precisely, if the Assumptions in Theorem 2.2 are satisfied, then $p = o(n)$. To see this, we notice that Assumptions 2.6, 2.7 and the fact that $\psi_n \leq k_n^{1/2}$ lead to $\psi_n = O(n^{\delta_0})$ for some $\delta_0 > 0$. Therefore, $p = o(n)$ follows from Assumption 2.4. In order to obtain consistency when p may grow as fast as n , one idea, but not the weakest possible, is to assume orthogonality of X , i.e., $X^T X = nI_p$, and to relax Assumption 2.7. To simplify the technical proof, we assume in the following Corollaries 2.5 and 2.6 that all c_j 's in model (2.2) are equal, and thus, $\bar{\phi}_n = \underline{\phi}_n$. We denote $\phi_n = \bar{\phi}_n$ (or $= \underline{\phi}_n$). Moreover, we need the following assumption about the growth rates of s_n and p to replace Assumptions 2.4 and 2.5.

Assumption 2.9. Let $a_n = n + \sigma_0^{-2}k_n/(n^{-1} + \phi_n)$ and $\zeta \in (1, \infty)$ be a constant such that $n\psi_n^2 > \sigma_0^2\zeta a_n$ as $n \rightarrow \infty$. The numbers p and s_n with $p \rightarrow \infty$ and $s_n \leq p \leq n$ satisfy

- (i). $s_n = o\left(\min\left\{\frac{(n+\nu)\log(n\psi_n^2/(\sigma_0^2\zeta a_n))}{\log(1+n\phi_n)}, n\psi_n^2, n\right\}\right)$.
- (ii). $p \log p = o(a_n)$.

Assumption 2.9 potentially allows the case $p = n$. To see this, suppose $s_n = O(1)$ and we choose ϕ_n such that $(n + \nu)/\log(1 + n\phi_n) \rightarrow \infty$. When a_n grows faster than $n \log n$ and $n\psi_n^2/a_n \rightarrow \infty$, $p = n$ will satisfy Assumption 2.9. However, this requires ψ_n^2 to grow at least faster than $\log n$. This extra requirement on ψ_n^2 has not been imposed by Theorems 2.2 and 2.4, and can be treated as the price which we pay to relax the growth rate for p . Under Assumption 2.9 and assuming orthogonality on X , we have the following consistency result which allows a faster growth rate for the dimension p .

Corollary 2.5. Assume that $X^T X = nI_p$ and $\Sigma = \phi_n I_p$ with $n\phi_n \rightarrow \infty$. Suppose γ^0 is nonnull and that Assumptions 2.1, 2.6 and 2.9 are satisfied. If $p^{\alpha_0(n+\nu)/a_n} = o(n\phi_n)$ for some $\alpha_0 > 2$, then $\max_{\gamma \neq \gamma^0} p(\gamma|Z)/p(\gamma^0|Z) \rightarrow_p 0$. If $p = o\left((n + \nu) \log\left(\frac{n\psi_n^2}{\sigma_0^2\zeta a_n}\right)\right)$ with ζ specified in Assumption 2.9, and $p^{2+\alpha_0(n+\nu)/a_n} = o(n\phi_n)$ for some $\alpha_0 > 2$, then $\sum_{\gamma \neq \gamma^0} p(\gamma|Z) \rightarrow_p 0$, and consequently, $p(\gamma^0|Z) \rightarrow_p 1$.

The proof of Corollary 2.5 is similar to those for Theorems 2.2 and 2.4 and is given in Supplement B. The following result, which requires a special model set-up, demonstrates that PMC and consistency of the posterior odds ratio may hold in some situations but fail in others.

Corollary 2.6. Assume $p = n$, $X^T X = nI_n$ and $\Sigma = \phi_n I_n$. Suppose $\min_{j \in \gamma^0} |\beta_j^0| \geq \psi_n$ with $\psi_n^2 = c_1 n^{1+\delta_1} (\log n)^2$ for some constants $\delta_1 > 1$ and $c_1 > 0$, $k_n = O(\psi_n^2)$ and $p(\gamma) = \text{constant}$ for all γ . Assume that $s_n = s$ with $s > 0$ a fixed integer, i.e., the true parameter vector β^0 contains exactly s nonzero components.

- (a). Suppose $\phi_n = c_2 n^{\delta_2}$ for some constants $c_2 > 0$ and δ_2 .

- i. If $-1 < \delta_2 \leq 1$, then $\max_{\gamma \neq \gamma^0} p(\gamma|Z)/p(\gamma^0|Z) \rightarrow_p 0$, but PMC does not hold. Specifically, when $-1 < \delta_2 < 1$, $p(\gamma^0|Z) \rightarrow 0$, a.s.; when $\delta_2 = 1$, then there exists a constant c_0 with $0 < c_0 < 1$ such that $\limsup_n p(\gamma^0|Z) \leq c_0$, a.s.
- ii. If $1 < \delta_2 \leq \delta_1$, then $p(\gamma^0|Z) \rightarrow_p 1$.
- (b). If $n^{\log n} = O(\phi_n)$, then $p(\emptyset|Z)/p(\gamma^0|Z) \rightarrow_p \infty$, where \emptyset represents the null model. Therefore, $p(\gamma^0|Z) \rightarrow_p 0$.
- (c). If $n\phi_n \rightarrow \eta \in [0, \infty)$, then almost surely, $\liminf_n \max_{\gamma \neq \gamma^0} p(\gamma|Z)/p(\gamma^0|Z) \geq (1 + \eta)^{-1/2}$ and $\lim_n p(\gamma^0|Z) = 0$.

The proof of Corollary 2.6 is given in Supplement B.

Remark 2.7. The main contribution of Corollary 2.6 is to demonstrate the difference between PMC and (2.4), and provide example growth rates for ϕ_n under which the two forms of consistency fail. Although this is obtained in a special situation, similar results should be still true under a more general setting, for instance, where $p < n$ or $X^T X$ is not diagonal, but we do not consider those circumstances here.

Corollary 2.6 (a) demonstrates that (2.4) does not necessarily imply PMC. This means that, although the posterior probability of the true model might not be approaching one, the ratio of the posterior probabilities of any “incorrect” model and the true model can still converge to zero. This phenomenon will not occur when p is fixed. In practice, (2.4) is sufficient to make a correct model selection even if PMC might fail.

Corollary 2.6 (b) and (c) demonstrate that in order to make a correct model selection, ϕ_n cannot be either too small or too large. Specifically, when $\phi_n = o(n^{-1})$, it follows by Corollary 2.6 (c) that almost surely $\liminf_n \max_{\gamma \neq \gamma^0} p(\gamma|Z)/p(\gamma^0|Z) \geq 1$. Thus, with probability one, for any $\varepsilon > 0$, there exists an integer N such that for any $n \geq N$

$$\max_{\gamma \neq \gamma^0} p(\gamma|Z)/p(\gamma^0|Z) \geq 1 - \varepsilon.$$

This implies that there exists a model, say γ^* , such that $p(\gamma^*|Z) \geq (1 - \varepsilon)p(\gamma^0|Z)$. Thus, when ε is small, either $p(\gamma^*|Z) > p(\gamma^0|Z)$, or $p(\gamma^*|Z)$ is very close to $p(\gamma^0|Z)$, which will both affect the selection result. On the

other hand, when ϕ_n is growing faster than $n^{n \log n}$, it follows from (b) that the null model will be preferred in favor of γ^0 .

Corollary 2.6 (b) and (c) can be also understood intuitively. When ϕ_n is too small, the two distribution components in the mixture prior of β tend to be indistinguishable so that it is difficult to separate the true model from some incorrect model; when ϕ_n approaches infinity, by (2.3), the posterior probability of any nonnull model approaches zero, and thus, all β_j 's are forced to be zero. This conclusion has been empirically obtained by Smith and Kohn (1996) under spline regression models. \square

Remark 2.8. Using arguments similar to the proofs of Theorems 2.2 and 2.4, and by the Borel-Cantelli lemma of Shao (2003), one can show the almost sure convergence of $p(\gamma^0|Z)$. We refer to Supplement C for details. \square

To conclude this section, let us look at an example which demonstrates that, when $\bar{\phi}_n = \bar{\phi}$ and $\underline{\phi}_n = \underline{\phi}$ with $\bar{\phi}$ and $\underline{\phi}$ unrelated to n , consistency might still hold under certain circumstances. This is motivated by a full Bayesian framework which requires all hyperparameters to be fixed.

Example 2.1. If a full Bayesian approach is desired, then we have to preselect the hyperparameters c_j 's, and so $\bar{\phi}_n = \bar{\phi}$ and $\underline{\phi}_n = \underline{\phi}$ must be fixed. Assume that $k_n = O(1)$, which is a slightly weaker assumption than that in Jiang (2007). Note that Assumptions 2.6 and 2.7 follow immediately. Suppose $\min_{j \in \gamma^0} |\beta_j^0| \geq \psi_n$ with $\psi_n = n^{-1/4} \sqrt{\log n}$, the prior distribution of model γ satisfies Assumption 2.1. Assume that $s_n = s$ with $s > 0$ a fixed integer (thus, the true model is nonnull), and design matrix X satisfies (2.10). Therefore, by Proposition 2.1 and Remark 2.2, Assumptions 2.2 and 2.8 both hold. We also notice that Assumption 2.3 is well satisfied. It follows from Theorem 2.2 that if $p = n^r$ for some $0 < r < 1/2$, then with probability approaching one, (2.4) holds, i.e., the true model can be correctly selected; if $p = n^r$ for some $0 < r < 1/4$, then PMC holds in probability.

3. Numerical results

In Section 2 we saw (e.g., Corollary 2.6) that the limiting behavior of $p(\gamma^0|Z)$ may rely on the hyperparameters c_j 's. In this section, simulated examples are given to numerically illustrate the relationship between $p(\gamma^0|Z)$ and the c_j 's.

To construct the random design matrix X , we generated *iid* p -dimensional row vectors $U_1, \dots, U_n \sim N(\mathbf{0}, I_p)$ and let U be an $n \times p$ matrix with i th row

U_i for $i = 1, \dots, n$. Then we let $X = \sqrt{n}U (U^T U)^{-1/2}$. Thus, $X^T X = nI_p$. (We choose X to be orthonormal for purposes of illustration, although, as we saw in the preceding section, results can be derived for general X .) To explore the dimension effect, we have considered two growth rates for p with respect to n : (1) $p = n^{1/2}$ and (2) $p = n^{3/4}$. Data were simulated from model (2.1) with $\sigma = 1$, $s_n = 5$ and the true model coefficients $(\beta_1^0, \dots, \beta_5^0) = (2, \dots, 2)$ and $(\beta_6^0, \dots, \beta_p^0) = (0, \dots, 0)$. We considered sample sizes $n = 70$ and 200 respectively.

The hierarchical Bayesian model (2.2) was fitted and the prior distributions on σ^2 and γ were assumed to be $1/\sigma^2 \sim \chi_4^2$ and $p(\gamma_j = 1) = p(\gamma_j = 0) = 1/2$, for any $j = 1, \dots, p$. For simplicity, we considered the case that $c_1 = \dots = c_p = \phi_n$. The values of ϕ_n were chosen to be $\phi_n = 10^q$ with $q = 1, 2, 3, 4, 5, 6, 8, 10, 20, 30, 60$. In particular, $q = 1, 2$ represent small ϕ_n 's and $q = 20, 30, 60$ represent extremely large ϕ_n 's. After 20,000 samples of (β, γ, σ) were drawn from the posterior distribution $p(\beta, \gamma, \sigma | Z)$ using a sub-blockwise Gibbs sampler developed by Godsill and Rayner (1998), we recorded the last 10,000 samples and treated the previous 10,000 samples as burnins. Convergence has been assessed by applying Gelman-Rubin's statistic to 5 parallel Markov chains for each ϕ_n . We denote $\gamma^{(1)}, \dots, \gamma^{(10000)}$ to be the last 10,000 samples of γ . Then $p(\gamma^0 | Z)$ is approximated by $p(\gamma^0 | Z) \approx \sum_{t=1}^{10000} I(\gamma^{(t)} = \gamma^0) / 10000$.

To study the frequentist property of $p(\gamma^0 | Z)$, we have generated 100 data sets Z_1, \dots, Z_{100} independently from model (2.1), and for each ϕ_n calculated the corresponding 100 posterior probabilities $p(\gamma^0 | Z_m)$, $m = 1, \dots, 100$. This idea was inspired from Fernández *et al.* (2001) who studied the Bayesian selection problem when p is fixed. For any $\alpha \in (0, 1)$, we define $RF(\alpha)$ to be the relative frequency of $p(\gamma^0 | Z_m)$'s greater than α , i.e., $RF(\alpha) = \sum_{m=1}^{100} I(p(\gamma^0 | Z_m) \geq \alpha) / 100$. Note that $RF(\alpha)$ is an estimate of $pr(p(\gamma^0 | Z) \geq \alpha)$, where $pr(\cdot)$ denotes the probability measure associated with the underlying probability space. The values of ϕ_n which satisfy PMC will therefore make $RF(\alpha)$ close to one, and those ϕ_n for which PMC fails will make $RF(\alpha)$ deviate from one. Next, we explore how $RF(\alpha)$ changes with ϕ_n for various values of α .

Figure 1 displays the relationship between $RF(\alpha)$ and $\kappa = \log \log \phi_n$ for $\alpha = 0.99, 0.90$ and growth rates (1) and (2). We have observed that there is a

large range for ϕ_n , which we might call a “feasible” region, such that $RF(\alpha)$ is close to one; while $RF(\alpha)$ approaches zero for either small or large ϕ_n . We also observed that, for the faster growth rate setting (2), the left boundary of the feasible region for ϕ_n lies more to the right, which has been expected based on Remark 2.3. Furthermore, the simulation results demonstrate that, when n increases, the right boundary of the feasible region for ϕ_n moves to the right.

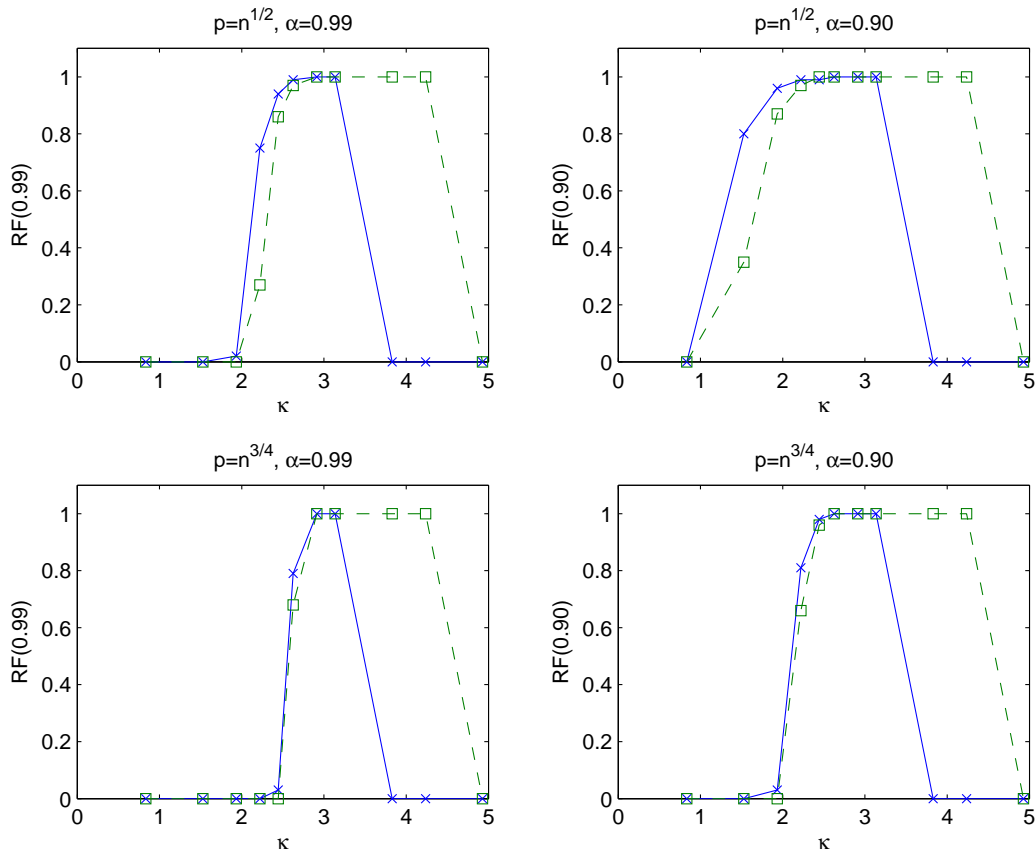


Figure 1: Plot of $RF(\alpha)$ versus $\kappa = \log \log \phi_n$. The vertical axis represents $RF(\alpha)$ with $\alpha = 0.99$ and 0.90 , and the horizontal axis represents $\kappa = \log \log \phi_n$ with $\phi_n = 10^q$ for $q = 1, 2, 3, 4, 5, 6, 8, 10, 20, 30, 60$. Two cases $n = 70$ (solid line with “ \times ” marks) and $n = 200$ (dashed line with “ \square ” marks) have been considered.

4. Conclusion

Previous work about posterior model consistency (PMC) includes Fernández *et al.* (2001) and Liang *et al.* (2008) when the number of parameters p is fixed. In this paper, we have studied PMC when the model dimension p grows with sample size n . Specifically, we have shown that, under a variation of the Bayesian model proposed by George and McCulloch (1993), the posterior probability of the true model converges to one, i.e., PMC holds. We have obtained this result in two situations: (i) design matrix X is general while p grows slower than n , e.g., $p \log n = o(n)$; (ii) $X^T X/n$ is the identity matrix and p may grow as fast as n , e.g., $p = n$. Furthermore, we have demonstrated under a special framework that the consistency results may fail if ϕ_n is too small or too large, where ϕ_n is the hyperparameter controlling the prior variance of the nonzero model coefficients. Precisely, when $\phi_n = o(n^{-1})$ (an example of small order) or when $n^{n \log n} = O(\phi_n)$ (an example of large order), both PMC and consistency of the posterior odds ratio fail. Besides that, our results do not require that the candidate models are pairwise nested.

Berger *et al.* (2003), Moreno *et al.* (2010) and Girón *et al.* (2010) have proved the consistency of Bayes factor when p is growing with n . This form of consistency, under our framework, is equivalent to the consistency of the posterior odds ratio if the prior odds ratio are uniformly bounded from above and below, so it is of interest to illustrate the relationship between PMC and consistency of posterior odds ratio. We have considered a special framework and shown that PMC implies consistency of the posterior odds ratio but the reverse may not be true. This is different from the finding by Liang *et al.* (2008) who demonstrate the equivalence of PMC and consistency of the Bayes factor when p is fixed. When combined with dimension reduction procedures such as SIS (Fan and Lv, 2008), our results can be also extended to ultrahigh-dimensional situations.

Two extensions of the current work are noteworthy. First, Assumption 2.7 is a technical assumption used to facilitate the proof and may not be the weakest possible. We leave it to future work to determine whether this condition can be further weakened or even removed. Second, model (2.2) results in a closed form for $p(\gamma|Z)$, which substantially reduces the complexity of the theory in this paper. However, there exist other useful Bayesian models which result in non-analytical forms for $p(\gamma|Z)$, such as a model with mixture g-priors (Liang *et al.*, 2008). The proof of consistency becomes complicated for such models, especially when p grows with n . We conjecture that the

asymptotics developed in this paper can be extended to models of this form and intend to explore these separately.

5. Appendix: proofs

In this section, we prove the main results in Section 2. We also prove some lemmas which are useful to establish the main results. Let $pr(\cdot)$ denote the probability measure associated with the underlying probability space.

Proof of Proposition 2.1. It follows by assumption that $\frac{1}{n}X_{\bar{\gamma}}^T X_{\bar{\gamma}} \geq cI_{|\bar{\gamma}|}$. Letting $X_{\bar{\gamma}} = (X_{\gamma}, X_{\bar{\gamma} \setminus \gamma})$, we can write $\frac{1}{n}X_{\bar{\gamma}}^T X_{\bar{\gamma}} = \begin{pmatrix} A & B \\ B^T & C \end{pmatrix}$, where $A = X_{\gamma}^T X_{\gamma}/n$, $B = X_{\gamma}^T X_{\bar{\gamma} \setminus \gamma}/n$ and $C = X_{\bar{\gamma} \setminus \gamma}^T X_{\bar{\gamma} \setminus \gamma}/n$. By formula for the inverse of blocked matrix (Seber and Lee, 2003, page 466), the lower right corner of $(\frac{1}{n}X_{\bar{\gamma}}^T X_{\bar{\gamma}})^{-1}$ is B_{22}^{-1} with $B_{22} = C - B^T A^{-1} B = \frac{1}{n}X_{\bar{\gamma} \setminus \gamma}^T (I_n - P_{\gamma}) X_{\bar{\gamma} \setminus \gamma}$. Then $B_{22}^{-1} \leq c^{-1}I$, which implies $\lambda_{-}(B_{22}) \geq c$. \square

Lemma 5.1. Suppose $\epsilon \sim N(0, \sigma_0^2 I_n)$. Then:

- (a). Let $v_{\gamma} = (I_n - P_{\gamma})X_{\gamma^0 \setminus \gamma} \beta_{\gamma^0 \setminus \gamma}^0$. If S_2 is nonnull, then $\max_{\gamma \in S_2} |v_{\gamma}^T \epsilon| / \|v_{\gamma}\|_2 = O_p(\sqrt{p})$, where we adopt the convention that $|v_{\gamma}^T \epsilon| / \|v_{\gamma}\|_2 = 0$ when $v_{\gamma} = 0$.
- (b). If S_1 is nonnull, then for any $\alpha > 2$, with probability approaching one, $\max_{\gamma \in S_1} \epsilon^T (P_{\gamma} - P_{\gamma^0}) \epsilon / (|\gamma| - s_n) \leq \alpha \sigma_0^2 \log p$.
- (c). If S_2 is nonnull, and we adopt the convention that $\epsilon^T P_{\gamma} \epsilon / |\gamma| = 0$ when γ is null, then for any $\alpha > 2$, with probability approaching one, $\max_{\gamma \in S_2} \epsilon^T P_{\gamma} \epsilon / |\gamma| \leq \alpha \sigma_0^2 \log p$.

Proof of Lemma 5.1. We prove the result for the case where X is deterministic, and briefly talk about the proofs for the case where X is random and independent of ϵ .

(a) We first assume that X is deterministic. By inequality (9.3) in Durrett (2005), if $\xi \sim N(0, 1)$, then there exists a C_0 such that for any $t > 1$, $pr(|\xi| \geq t) \leq C_0 \exp(-t^2/2)$. Note that $|v_{\gamma}^T \epsilon| / (\sigma_0 \|v_{\gamma}\|_2) \sim N(0, 1)$, and therefore, by Bonferroni's inequality,

$$pr \left(\max_{\gamma \in S_2} \frac{|v_{\gamma}^T \epsilon|}{\|v_{\gamma}\|_2} \geq t \right) \leq \sum_{\gamma \in S_2} pr \left(\frac{|v_{\gamma}^T \epsilon|}{\|v_{\gamma}\|_2} \geq t \right) \leq C_0 2^p \exp \left(-\frac{t^2}{2\sigma_0^2} \right).$$

Then the result holds by setting $t = C\sigma_0\sqrt{2p}$ with large C . When X is random but independent of ϵ , note that the conditional distribution of $|v_\gamma^T \epsilon|/(\sigma_0 \|v_\gamma\|_2)$ given X is $N(0, 1)$. Thus, the proof can be finished by the above arguments.

(b) Suppose X is deterministic. First, if $\xi = \chi_\mu^2$, then by Chebyshev's inequality, for any $2 < \alpha' < \alpha$,

$$\begin{aligned} & pr(\xi \geq \alpha\mu \log p) \\ &= pr(\exp(\xi/\alpha') \geq \exp((\alpha/\alpha')\mu \log p)) \\ &\leq \exp(-(\alpha/\alpha')\mu \log p) E\{\exp(\xi/\alpha')\} \\ &= (1 - 2/\alpha')^{-\mu/2} \exp(-(\alpha/\alpha')\mu \log p). \end{aligned}$$

Using this inequality, Bonferroni's inequality, and the fact that when $\gamma \in S_1$, $\epsilon^T(P_\gamma - P_{\gamma^0})\epsilon \sim \sigma_0^2 \chi_{|\gamma| - s_n}^2$, we have

$$\begin{aligned} & pr\left(\max_{\gamma \in S_1} \frac{\epsilon^T(P_\gamma - P_{\gamma^0})\epsilon}{|\gamma| - s_n} \geq \alpha\sigma_0^2 \log p\right) \\ &\leq \sum_{\gamma \in S_1} pr(\epsilon^T(P_\gamma - P_{\gamma^0})\epsilon \geq \alpha\sigma_0^2(|\gamma| - s_n) \log p) \\ &\leq \sum_{\gamma \in S_1} (1 - 2/\alpha')^{-(|\gamma| - s_n)/2} \exp(-(\alpha/\alpha')(|\gamma| - s_n) \log p) \\ &= \sum_{r=1}^{p-s_n} \binom{p-s_n}{r} (1 - 2/\alpha')^{-r/2} \exp(-(\alpha/\alpha')r \log p) \\ &= \left(1 + (1 - 2/\alpha')^{-1/2} p^{-\alpha/\alpha'}\right)^{p-s_n} - 1 \rightarrow 0. \end{aligned}$$

When X is random and independent of ϵ , then conditioning on X , $\epsilon^T(P_\gamma - P_{\gamma^0})\epsilon \sim \sigma_0^2 \chi_{|\gamma| - s_n}^2$. Thus, the conclusion follows from the above arguments.

(c) We let X be deterministic. The case where X is random can be handled similarly. Assume that S_2 contains nonnull models, and note that when γ is nonnull, $\epsilon^T P_\gamma \epsilon \sim \sigma_0^2 \chi_{|\gamma|}^2$. Fix arbitrarily α' such that $2 < \alpha' < \alpha$.

Then by the proof of part (b) we have

$$\begin{aligned}
& pr \left(\max_{\gamma \in S_2} \frac{\epsilon^T P_\gamma \epsilon}{|\gamma|} \geq \alpha \sigma_0^2 \log p \right) \\
&= pr \left(\max_{\gamma \in S_2 \setminus \{\emptyset\}} \frac{\epsilon^T P_\gamma \epsilon}{|\gamma|} \geq \alpha \sigma_0^2 \log p \right) \\
&\leq \sum_{\gamma \in S_2 \setminus \{\emptyset\}} pr \left(\epsilon^T P_\gamma \epsilon \geq \alpha \sigma_0^2 |\gamma| \log p \right) \\
&\leq \sum_{\gamma \in S_2 \setminus \{\emptyset\}} (1 - \alpha'/2)^{-|\gamma|/2} \exp(-(\alpha/\alpha')|\gamma| \log p) \\
&\leq \sum_{r=1}^p \binom{p}{r} (1 - 2/\alpha')^{-r/2} p^{-(\alpha/\alpha')r} \\
&= \left(1 + (1 - 2/\alpha')^{-1/2} p^{-\alpha/\alpha'} \right)^p - 1 \rightarrow 0. \quad \square
\end{aligned}$$

Proof of Theorem 2.2. We have

$$\begin{aligned}
-\log(p(\gamma|Z)/p(\gamma^0|Z)) &= -\log\left(\frac{p(\gamma)}{p(\gamma^0)}\right) + \frac{1}{2} \log\left(\frac{\det(W_\gamma)}{\det(W_{\gamma^0})}\right) \\
&\quad + \frac{n+\nu}{2} \log\left(\frac{1 + \mathbf{y}^T (I_n - X_\gamma U_\gamma^{-1} X_\gamma^T) \mathbf{y}}{1 + \mathbf{y}^T (I_n - X_{\gamma^0} U_{\gamma^0}^{-1} X_{\gamma^0}^T) \mathbf{y}}\right) \\
&= -\log\left(\frac{p(\gamma)}{p(\gamma^0)}\right) + \frac{1}{2} \log\left(\frac{\det(W_\gamma)}{\det(W_{\gamma^0})}\right) \\
&\quad + \frac{n+\nu}{2} \log\left(\frac{1 + \mathbf{y}^T (I_n - X_\gamma U_\gamma^{-1} X_\gamma^T) \mathbf{y}}{1 + \mathbf{y}^T (I_n - P_\gamma) \mathbf{y}}\right) \\
&\quad - \frac{n+\nu}{2} \log\left(\frac{1 + \mathbf{y}^T (I_n - X_{\gamma^0} U_{\gamma^0}^{-1} X_{\gamma^0}^T) \mathbf{y}}{1 + \mathbf{y}^T (I_n - P_{\gamma^0}) \mathbf{y}}\right) \\
&\quad + \frac{n+\nu}{2} \log\left(\frac{1 + \mathbf{y}^T (I_n - P_\gamma) \mathbf{y}}{1 + \mathbf{y}^T (I_n - P_{\gamma^0}) \mathbf{y}}\right). \quad (5.1)
\end{aligned}$$

Denote the above summands by T_1, T_2, T_3, T_4, T_5 . By Assumption 2.6, T_1 is bounded below. Since $U_\gamma \geq X_\gamma^T X_\gamma$, we have $T_3 \geq 0$ for any n .

To approximate T_4 , let

$$\Delta = \mathbf{y}^T X_{\gamma^0} (X_{\gamma^0}^T X_{\gamma^0})^{-1} \left(\Sigma_{\gamma^0} + (X_{\gamma^0}^T X_{\gamma^0})^{-1} \right)^{-1} (X_{\gamma^0}^T X_{\gamma^0})^{-1} X_{\gamma^0}^T \mathbf{y}.$$

By the Sherman-Morrison-Woodbury matrix identity (Seber and Lee, 2003, page 467),

$$U_{\gamma^0}^{-1} - (X_{\gamma^0}^T X_{\gamma^0})^{-1} = - (X_{\gamma^0}^T X_{\gamma^0})^{-1} \left(\Sigma_{\gamma^0} + (X_{\gamma^0}^T X_{\gamma^0})^{-1} \right)^{-1} (X_{\gamma^0}^T X_{\gamma^0})^{-1}. \quad (5.2)$$

By (5.2) and the fact that $\left(\Sigma_{\gamma^0} + (X_{\gamma^0}^T X_{\gamma^0})^{-1} \right)^{-1} \leq \Sigma_{\gamma^0}^{-1}$, we have

$$\begin{aligned} & \frac{1 + \mathbf{y}^T (I_n - X_{\gamma^0} U_{\gamma^0}^{-1} X_{\gamma^0}^T) \mathbf{y}}{1 + \mathbf{y}^T (I_n - P_{\gamma^0}) \mathbf{y}} \\ &= 1 + \frac{\Delta}{1 + \mathbf{y}^T (I_n - P_{\gamma^0}) \mathbf{y}} \\ &\leq 1 + 2 \left(\frac{(\beta_{\gamma^0}^0)^T \Sigma_{\gamma^0}^{-1} \beta_{\gamma^0}^0 + \epsilon^T X_{\gamma^0} (X_{\gamma^0}^T X_{\gamma^0})^{-1} \Sigma_{\gamma^0}^{-1} (X_{\gamma^0}^T X_{\gamma^0})^{-1} X_{\gamma^0}^T \epsilon}{1 + \mathbf{y}^T (I_n - P_{\gamma^0}) \mathbf{y}} \right) \\ &\leq 1 + 2 \underline{\phi}_n^{-1} \left(\frac{\|\beta_{\gamma^0}^0\|_2^2 + \epsilon^T X_{\gamma^0} (X_{\gamma^0}^T X_{\gamma^0})^{-2} X_{\gamma^0}^T \epsilon}{1 + \mathbf{y}^T (I_n - P_{\gamma^0}) \mathbf{y}} \right). \end{aligned}$$

Since $\mathbf{y}^T (I_n - P_{\gamma^0}) \mathbf{y} / n = \epsilon^T (I_n - P_{\gamma^0}) \epsilon / n \rightarrow_p \sigma_0^2$, and $E\{\epsilon^T X_{\gamma^0} (X_{\gamma^0}^T X_{\gamma^0})^{-2} X_{\gamma^0}^T \epsilon\} \leq \sigma_0^2 s_n (n \varphi_{\min}(n))^{-1}$, we have $\epsilon^T X_{\gamma^0} (X_{\gamma^0}^T X_{\gamma^0})^{-2} X_{\gamma^0}^T \epsilon = O_p(s_n (n \varphi_{\min}(n))^{-1})$. Therefore, by Assumptions 2.2 and 2.3, and the fact that $k_n \geq s_n \psi_n^2$, we can show that

$$\frac{1 + \mathbf{y}^T (I_n - X_{\gamma^0} U_{\gamma^0}^{-1} X_{\gamma^0}^T) \mathbf{y}}{1 + \mathbf{y}^T (I_n - P_{\gamma^0}) \mathbf{y}} \leq 1 + \frac{2k_n}{n \underline{\phi}_n \sigma_0^2} (1 + o_p(1)). \quad (5.3)$$

Consequently, $0 \leq -T_4 = O_p(1)$ follows from the condition that $k_n = O(\underline{\phi}_n)$ (Assumption 2.7).

Next we approximate T_2 and T_5 in the following Lemmas 5.2 and 5.3.

Lemma 5.2. Under Assumption 2.8, if $\gamma \in S_1$, then $T_2 \geq 2^{-1}(|\gamma| - s_n) \log(1 + C_3 n^{1-\delta} \underline{\phi}_n)$. Under Assumption 2.2, if $\gamma \in S_2$, $T_2 \geq -2^{-1} s_n \log(1 + C_2 n \bar{\phi}_n)$, where \bar{C}_2 and C_3 are constants given in Assumptions 2.2 and 2.8 respectively.

Proof of Lemma 5.2. If $\gamma \in S_1$, it follows from the determinant formula for block matrices (Seber and Lee, 2003, page 468), and Assumption

2.8 that

$$\begin{aligned}
\det(U_\gamma) &= \det(U_{\gamma^0}) \det\left(\Sigma_{\gamma\setminus\gamma^0}^{-1} + X_{\gamma\setminus\gamma^0}^T(I_n - X_{\gamma^0}U_{\gamma^0}^{-1}X_{\gamma^0}^T)X_{\gamma\setminus\gamma^0}\right) \\
&\geq \det(U_{\gamma^0}) \det\left(\Sigma_{\gamma\setminus\gamma^0}^{-1} + X_{\gamma\setminus\gamma^0}^T(I_n - P_{\gamma^0})X_{\gamma\setminus\gamma^0}\right) \\
&\geq \det(U_{\gamma^0}) \det\left(\Sigma_{\gamma\setminus\gamma^0}^{-1} + C_3n^{1-\delta}I_{|\gamma\setminus\gamma^0|}\right).
\end{aligned}$$

Therefore,

$$\begin{aligned}
\frac{\det(W_\gamma)}{\det(W_{\gamma^0})} &= \frac{\det(\Sigma_\gamma)}{\det(\Sigma_{\gamma^0})} \frac{\det(U_\gamma)}{\det(U_{\gamma^0})} \\
&\geq \det(\Sigma_{\gamma\setminus\gamma^0}) \det\left(\Sigma_{\gamma\setminus\gamma^0}^{-1} + C_3n^{1-\delta}I_{|\gamma\setminus\gamma^0|}\right) \\
&= \det\left(I_{|\gamma\setminus\gamma^0|} + C_3n^{1-\delta}\Sigma_{\gamma\setminus\gamma^0}\right) \\
&\geq \det\left((1 + C_3n^{1-\delta}\underline{\phi}_n)I_{|\gamma\setminus\gamma^0|}\right) = (1 + C_3n^{1-\delta}\underline{\phi}_n)^{|\gamma|-s_n}, \quad (5.4)
\end{aligned}$$

which shows that $T_2 \geq 2^{-1}(|\gamma| - s_n) \log(1 + C_3n^{1-\delta}\underline{\phi}_n)$. If $\gamma \in S_2$, note that $\det(W_\gamma) \geq 1$, and by Assumption 2.2

$$T_2 \geq -\frac{1}{2} \log(\det(W_{\gamma^0})) \geq -\frac{1}{2} \log(\det(I_{s_n} + C_2n\Sigma_{\gamma^0})) \geq -2^{-1}s_n \log(1 + C_2n\bar{\phi}_n),$$

which completes the proof of Lemma 5.2. \square

Lemma 5.3. Let $\alpha_0 > 2$. If either Assumption 2.4 or 2.5 is satisfied, when n is large, with large probability and uniformly for $\gamma \in S_1$, $T_5 \geq -2^{-1}(|\gamma| - s_n)\alpha_0 \log p$. If both Assumptions 2.2 and 2.4 are satisfied, there exists a constant C' such that when n is large, with large probability and uniformly for $\gamma \in S_2$, $T_5 \geq 2^{-1}(n + \nu) \log(1 + C'\psi_n^2)$.

Proof of Lemma 5.3. We consider $\gamma \in S_1$ and S_2 separately. Notice that Assumption 2.4 implies that $p \log p = o(n \log(1 + \psi_n^2))$, and therefore implies that $p \log p = o(n\psi_n^2)$. Let $v_\gamma = (I_n - P_\gamma)X_{\gamma^0\setminus\gamma}\beta_{\gamma^0\setminus\gamma}^0$. From Lemma 5.1 (a) and (c), there exists $C > 0$ such that when n is sufficiently large, with

large probability, for any $\gamma \in S_2$,

$$\begin{aligned}
\mathbf{y}^T(I_n - P_\gamma)\mathbf{y} &= \|v_\gamma\|_2^2 + 2v_\gamma^T\epsilon + \epsilon^T(I_n - P_\gamma)\epsilon \\
&\geq \|v_\gamma\|_2^2 - 2C\sqrt{p}\|v_\gamma\|_2 + \epsilon^T\epsilon - C|\gamma|\log p \\
&\geq \|v_\gamma\|_2^2 \left(1 - 2C\frac{\sqrt{p}}{\|v_\gamma\|_2} - C\frac{p\log p}{\|v_\gamma\|_2^2}\right) + \epsilon^T\epsilon \\
&\geq \|v_\gamma\|_2^2 \left(1 - 2C\sqrt{\frac{p}{n\varphi_{\min}(n)\psi_n^2}} - C\frac{p\log p}{n\varphi_{\min}(n)\psi_n^2}\right) + \epsilon^T\epsilon \\
&= \|v_\gamma\|_2^2(1 + o(1)) + \epsilon^T\epsilon \\
&\geq n\varphi_{\min}(n)\|\beta_{\gamma^0\setminus\gamma}^0\|_2^2(1 + o(1)) + \epsilon^T\epsilon \\
&\geq n\varphi_{\min}(n)\psi_n^2(1 + o(1)) + \epsilon^T\epsilon. \tag{5.5}
\end{aligned}$$

It is easy to see that Assumption 2.4 implies that $s_n = o(n)$, and therefore, $\epsilon^T(I_n - P_{\gamma^0})\epsilon = n\sigma_0^2(1 + o_p(1))$. Thus, by (5.5), there exists a C' such that for sufficiently large n , with large probability, uniformly for $\gamma \in S_2$,

$$T_5 \geq \frac{n + \nu}{2} \log \left(\frac{1 + n\varphi_{\min}(n)\psi_n^2(1 + o(1)) + \epsilon^T\epsilon}{1 + \epsilon^T(I_n - P_{\gamma^0})\epsilon} \right) \geq \frac{n + \nu}{2} \log(1 + C'\psi_n^2). \tag{5.6}$$

On the other hand, by properties of projection matrices and Lemma 5.1 (b), when n is sufficiently large, with large probability, we have uniformly for $\gamma \in S_1$,

$$\begin{aligned}
&\frac{1 + \mathbf{y}^T(I_n - P_\gamma)\mathbf{y}}{1 + \mathbf{y}^T(I_n - P_{\gamma^0})\mathbf{y}} \\
&= 1 - \frac{\mathbf{y}^T(P_\gamma - P_{\gamma^0})\mathbf{y}}{1 + \mathbf{y}^T(I_n - P_{\gamma^0})\mathbf{y}} \\
&= 1 - \frac{(\beta_{\gamma^0}^0)^T X_{\gamma^0}^T (P_\gamma - P_{\gamma^0}) X_{\gamma^0} \beta_{\gamma^0} + 2(\beta_{\gamma^0}^0)^T X_{\gamma^0}^T (P_\gamma - P_{\gamma^0})\epsilon + \epsilon^T (P_\gamma - P_{\gamma^0})\epsilon}{1 + \mathbf{y}^T(I_n - P_{\gamma^0})\mathbf{y}} \\
&= 1 - \frac{\epsilon^T (P_\gamma - P_{\gamma^0})\epsilon}{1 + \epsilon^T(I_n - P_{\gamma^0})\epsilon} \geq 1 - \frac{\alpha(|\gamma| - s_n) \log p}{n},
\end{aligned}$$

where we have temporarily fixed an α such that $2 < \alpha < \sqrt{2\alpha_0}$. It follows by the inequality that $\log(1 - x) \geq -(\alpha/2)x$ when $x \in (0, 1 - 2/\alpha)$, and by Assumption 2.4 or 2.5 (which both imply that $(|\gamma| - s_n) \log p/n$ approaches zero uniformly for $\gamma \in S_1$) that for sufficiently large n , with large probability

and uniformly for $\gamma \in S_1$,

$$T_5 \geq \frac{n+\nu}{2} \log \left(1 - \frac{\alpha(|\gamma| - s_n) \log p}{n} \right) \geq -2^{-1}(|\gamma| - s_n) \alpha_0 \log p, \quad (5.7)$$

which completes the proof of Lemma 5.3. \square

Now we are ready to finish the proof of Theorem 2.2. By (5.3), Lemma 5.2, Lemma 5.3, Assumption 2.4, and the fact that $p^{\alpha_0} = o(\rho_n)$ with $\rho_n \equiv n^{1-\delta} \underline{\phi}_n$, with large probability, uniformly for $\gamma \in S_1$,

$$\begin{aligned} p(\gamma|Z)/p(\gamma^0|Z) &\leq \tilde{C} \exp \left(-2^{-1}(|\gamma| - s_n) \log((1 + C_3 \rho_n)/p^{\alpha_0}) \right) \\ &= \tilde{C} \left(\frac{1 + C_3 \rho_n}{p^{\alpha_0}} \right)^{-2^{-1}(|\gamma| - s_n)} \rightarrow 0. \end{aligned} \quad (5.8)$$

By Assumptions 2.4 and 2.6, it can be verified that $s_n \log(1 + C_2 n \bar{\phi}_n) \ll \frac{n+\nu}{2} \log(1 + C' \psi_n^2)$. So, with large probability, uniformly for $\gamma \in S_2$,

$$\begin{aligned} p(\gamma|Z)/p(\gamma^0|Z) &\leq \tilde{C} \exp \left(2^{-1} s_n \log(1 + C_2 n \bar{\phi}_n) - \frac{n+\nu}{2} \log(1 + C' \psi_n^2) \right) \\ &\leq \tilde{C} (1 + C' \psi_n^2)^{-\frac{n+\nu}{4}}. \end{aligned} \quad (5.9)$$

where \tilde{C} in (5.8) and (5.9) depends on the lower bounds of T_1 and T_4 . For the proof of PMC, we consider two cases. It is easy to see from (5.8) that

$$\begin{aligned} \sum_{\gamma \in S_1} p(\gamma|Z)/p(\gamma^0|Z) &\leq \tilde{C} \sum_{\gamma \in S_1} \left(\frac{1 + C_3 \rho_n}{p^{\alpha_0}} \right)^{-2^{-1}(|\gamma| - s_n)} \\ &= \tilde{C} \sum_{r=1}^{p-s_n} \binom{p-s_n}{r} \left(\frac{1 + C_3 \rho_n}{p^{\alpha_0}} \right)^{-\frac{r}{2}} \\ &= \tilde{C} \left[\left(1 + \left(\frac{1 + C_3 \rho_n}{p^{\alpha_0}} \right)^{-\frac{1}{2}} \right)^{p-s_n} - 1 \right] \rightarrow 0 \end{aligned} \quad (5.10)$$

where the last limit result follows from the assumption that $p^{\alpha_0+2} = o(\rho_n)$.

Similarly, by (5.9), and $p \log n = o(n \log(1 + \psi_n^2))$ (which follows from Assumption 2.4), we can show that

$$\sum_{\gamma \in S_2} p(\gamma|Z)/p(\gamma^0|Z) \leq \tilde{C} 2^p (1 + C' \psi_n^2)^{-(n+\nu)/4} \rightarrow 0. \quad (5.11)$$

This completes the proof of Theorem 2.2. \square

Proof of Theorem 2.4. The assumption that γ^0 is null implies that the model class S_2 is empty. Similar to the proof of Theorem 2.2, we need to approximate T_1 to T_5 in (5.1). This is easier when the true model is null since $T_4 = 0$, and by Lemma 5.2, when γ is nonnull, $T_2 \geq 2^{-1}|\gamma| \log(1 + C_3 n^{1-\delta} \underline{\phi}_n)$. Since T_1 and T_3 are still bounded below, the proof is reduced to approximate T_5 . By Lemma 5.3, Assumption 2.5, and that $s_n = 0$, when n is large, with large probability and uniformly for $\gamma \in S_1$, $T_5 \geq -2^{-1}|\gamma|\alpha_0 \log p$. Therefore, the remaining proofs can be finished by arguments similar to (5.8) and (5.10). \square

Supplement Materials

Supplements A–C are given in the authors' website <http://www.stat.wisc.edu/~shang/>

Acknowledgement The authors wish to thank Professor Jun Shao for suggestions that helped to improve the present work.

References

- Berger, J. O. and Pericchi, L. (1996). The intrinsic Bayes factor for model selection and prediction. *J. Amer. Statist. Assoc.* **91**, 109–122.
- Berger, J. O., Ghosh, J. K. and Mukhopadhyay, N. (2003). Approximations and consistency of Bayes factors as model dimension grows. *J. Statist. Planning. Inference.* **112**, 241–258.
- Bühlmann, P., and Kalisch, M. and Maathuis, M. H. (2010). Variable selection in high-dimensional linear models: partially faithful distributions and the PC-simple algorithm. *Biometrika* **97**, 261–278.
- Casella, C., Girón, F. J., Martínez, M. L. and Moreno, E. (2009). Consistency of Bayesian procedures for variable selection. *Ann. Statist.* **37**, 1207–1228.
- Clyde, M. and George, E. I. (2000). Flexible empirical Bayes estimation for wavelets. *J. R. Stat. Soc. Ser. B.* **62**, 681–698.

- Clyde, M., Parmigiani, G. and Vidakovic, B. (1998). Multiple shrinkage and subset selection in wavelets. *Biometrika* **85**, 391–401.
- Durrett, R. (2005). *Probability: Theory and Examples*. 3rd Ed. Wadsworth-Brooks/Cole, Pacific Grove.
- Fan, J. and Lv, J. (2008). Sure independence screening for ultrahigh dimensional feature space. *J. R. Stat. Soc. Ser. B.* **70**, 849–911.
- Fan, J. and Peng, H. (2004). Nonconcave penalized likelihood with a diverging number of parameters. *Ann. Statist.* **32**, 928–961.
- Fernández, C., Ley, E. and Steel, M. F. J. (2001). Benchmark priors for Bayesian model averaging. *J. Econometrics* **100**, 381–427.
- George, E. and McCulloch, R. (1993). Variable selection via Gibbs sampling. *J. Amer. Statist. Assoc.* **88**, 881–889.
- Godsill, J. S. and Rayner, P. J. W. (1998). Robust reconstruction and analysis of autoregressive signals in impulsive noise using the Gibbs sampler. *IEEE Trans. Speech Audio Process* **6**, 352–372.
- Jeffreys, H. (1967). *Theory of Probability*. 4th Ed. Oxford Univ. Press, Oxford.
- Jiang, W. (2007). Bayesian variable selection for high dimensional generalized linear models: Convergence rates of the fitted densities. *Ann. Statist.* **35**, 1487–1511.
- Liang, F., Paulo, R., Molina, G., Clyde, M. and Berger, J. O. (2008). Mixtures of g -priors for Bayesian variable selection. *J. Amer. Statist. Assoc.* **103**, 410–423.
- Meinshausen, N. and Bühlmann, P. (2006). High-dimensional graphs and variable selection with the lasso. *Ann. Statist.* **34**, 1436–1462.
- Meinshausen, N. and Yu, B. (2009). Lasso-type recovery of sparse representations for high-dimensional data. *Ann. Statist.* **37**, 246–270.
- Moreno, E., Bertolino, F. and Racugno, W. (1998). An intrinsic limiting procedure for model selection and hypotheses testing. *J. Amer. Statist. Assoc.* **93**, 1451–1460.

- Moreno, E. and Girón, F. J. (2005). Consistency of Bayes factors for intrinsic priors in normal linear models. *C. R. Math. Acad. Sci. Paris* **340**, 911–914.
- Moreno, E., Girón, F. J. and Casella, G. (2010). Consistency of objective Bayes factors as the model dimension grows. *Ann. Statist.* **38**, 1937–1952.
- Girón, F. J., Moreno, E., Casella, G. and Martínez, M. L. (2010). Consistency of objective Bayes factors for nonnested linear models and increasing model dimension. *Revista de la Real Academia de Ciencias Exactas, Físicas y Naturales. Serie A. Matemáticas* **104**, 57–67.
- Shao, J. (2003). *Mathematical Statistics*, 2nd Ed. Springer Texts in Statistics. Springer, New York.
- Seber, G. A. F. and Lee, A. J. (2003). *Linear Regression Analysis*, 2nd Ed. Wiley-Interscience [John Wiley & Sons], Hoboken, NJ.
- Smith, M. S. and Kohn, R. (1996). Nonparametric regression using Bayesian variable selection. *J. Econometrics* **75**, 317–344.
- Wolfe, P. J., Godsill, S. J. and Ng, W.-J. (2004). Bayesian variable selection and regularization for time-frequency surface estimation. *J. R. Stat. Soc. Ser. B.* **66**, 575–589.
- Zellner, A. (1971). *An Introduction to Bayesian Inference in Econometrics*. Wiley, New York.
- Zellner, A. (1978). Jeffreys-Bayes posterior odds ratio and the Akaike information criterion for discriminating between models. *Econom. Lett.* **1**, 337–342.
- Zhang, J., Clayton, M. K. and Townsend, P. (2010). Functional concurrent linear regression model for spatial images. *Journal of Agricultural, Biological and Environmental Statistics* In press.
- Zhang, C.-H. and Huang, J. (2008). The sparsity and bias of the LASSO selection in high-dimensional linear regression. *Ann. Statist.* **36**, 1567–1594.