

FULLY BAYES FACTORS WITH A GENERALIZED g -PRIOR

BY YUZO MARUYAMA AND EDWARD I. GEORGE

University of Tokyo and University of Pennsylvania

For the normal linear model variable selection problem, we propose selection criteria based on a fully Bayes formulation with a generalization of Zellner's g -prior which allows for $p > n$. A special case of the prior formulation is seen to yield tractable closed forms for marginal densities and Bayes factors which reveal new model evaluation characteristics of potential interest.

1. Introduction. Suppose the normal linear regression model is used to relate y to the potential predictors x_1, \dots, x_p ,

$$(1.1) \quad \mathbf{y} \sim N_n(\alpha \mathbf{1}_n + \mathbf{X}_F \boldsymbol{\beta}_F, \sigma^2 \mathbf{I}_n),$$

where α is an unknown intercept parameter, $\mathbf{1}_n$ is an $n \times 1$ vector each component of which is one, $\mathbf{X}_F = (\mathbf{x}_1, \dots, \mathbf{x}_p)$ is an $n \times p$ design matrix, $\boldsymbol{\beta}_F$ is a $p \times 1$ vector of unknown regression coefficients, \mathbf{I}_n is an $n \times n$ identity matrix and σ^2 is an unknown positive scalar. (The subscript F denotes the full model.) We assume that the columns of \mathbf{X}_F have been standardized so that for $1 \leq i \leq p$, $\mathbf{x}'_i \mathbf{1}_n = 0$ and $\mathbf{x}'_i \mathbf{x}_i / n = 1$.

We shall be particularly interested in the variable selection problem where we would like to select an unknown subset of the important predictors. It will be convenient throughout to index each of these 2^p possible subset choices by the vector

$$\boldsymbol{\gamma} = (\gamma_1, \dots, \gamma_p)',$$

where $\gamma_i = 0$ or 1. We use $q_\gamma = \boldsymbol{\gamma}' \mathbf{1}_p$ to denote the size of the $\boldsymbol{\gamma}$ th subset. The problem then becomes that of selecting a submodel of (1.1) which has a density of the form

$$(1.2) \quad p(\mathbf{y} | \alpha, \boldsymbol{\beta}_\gamma, \sigma^2, \boldsymbol{\gamma}) = \phi_n(\mathbf{y}; \alpha \mathbf{1}_n + \mathbf{X}_\gamma \boldsymbol{\beta}_\gamma, \sigma^2 \mathbf{I}_n),$$

where $\phi_n(\mathbf{y}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$ denotes the n -variate normal density with mean vector $\boldsymbol{\mu}$ and covariance matrix $\boldsymbol{\Sigma}$. In (1.2), \mathbf{X}_γ is the $n \times q_\gamma$ matrix whose columns correspond to the $\boldsymbol{\gamma}$ th subset of x_1, \dots, x_p , and $\boldsymbol{\beta}_\gamma$ is a $q_\gamma \times 1$ vector of unknown regression coefficients. We assume throughout that \mathbf{X}_γ is of full rank denoted

$$r_\gamma = \min\{q_\gamma, n - 1\}.$$

Received September 2010; revised August 2011.

MSC2010 subject classifications. Primary 62F07, 62F15; secondary 62C10.

Key words and phrases. Bayes factor, model selection consistency, ridge regression, singular value decomposition, variable selection.

Last, let \mathcal{M}_γ denote the submodel given by (1.2).

A Bayesian approach to this problem entails the specification of prior distributions on the models $\pi_\gamma = \Pr(\mathcal{M}_\gamma)$, and on the parameters $p(\alpha, \beta_\gamma, \sigma^2)$ of each model. For each such specification, of key interest is the posterior probability of \mathcal{M}_γ given \mathbf{y} ,

$$(1.3) \quad \Pr(\mathcal{M}_\gamma|\mathbf{y}) = \frac{\pi_\gamma m_\gamma(\mathbf{y})}{\sum_\gamma \pi_\gamma m_\gamma(\mathbf{y})} = \frac{\pi_\gamma \text{BF}_{\gamma:N}}{\sum_\gamma \pi_\gamma \text{BF}_{\gamma:N}},$$

where $m_\gamma(\mathbf{y})$ is the marginal density of \mathbf{y} under \mathcal{M}_γ . In (1.3), $\text{BF}_{\gamma:N}$ is the so-called “null-based Bayes factor” for comparing each of \mathcal{M}_γ to the null model \mathcal{M}_N which is defined as

$$\text{BF}_{\gamma:N} = \frac{m_\gamma(\mathbf{y})}{m_N(\mathbf{y})},$$

where the null model \mathcal{M}_N is given by $\mathbf{y} \sim N_n(\alpha \mathbf{1}_n, \sigma^2 \mathbf{I}_n)$ and $m_N(\mathbf{y})$ is the marginal density of \mathbf{y} under the null model. For model selection, a popular strategy is to select the model for which $\Pr(\mathcal{M}_\gamma|\mathbf{y})$ or $\pi_\gamma \text{BF}_{\gamma:N}$ is largest.

Our main focus in this paper is to propose and study specifications for the parameter prior for each submodel \mathcal{M}_γ , which we will consider to be of the form

$$(1.4) \quad \begin{aligned} p(\alpha, \beta_\gamma, \sigma^2) &= p(\alpha)p(\sigma^2)p(\beta_\gamma|\sigma^2) \\ &= p(\alpha)p(\sigma^2) \int p(\beta_\gamma|\sigma^2, g)p(g) dg, \end{aligned}$$

where g is a hyperparameter. In Section 2 we explicitly describe our choices of prior forms for (1.4). Our key innovation there will be to use a generalization of

$$(1.5) \quad p(\beta_\gamma|\sigma^2, g) = \phi_{q_\gamma}(\beta; \mathbf{0}, g\sigma^2(\mathbf{X}'_\gamma \mathbf{X}_\gamma)^{-1}),$$

Zellner’s (1986) g -prior, a normal conjugate form which leads to tractable marginalization, for example, see George and Foster (2000), Fernández, Ley and Steel (2001), Liang et al. (2008). Under (1.5) and a flat prior on α , the marginal density of \mathbf{y} given g and σ^2 under \mathcal{M}_γ is given by

$$(1.6) \quad m_\gamma(\mathbf{y}|g, \sigma^2) \propto \exp\left(\frac{g}{g+1} \left\{ \max_{\alpha, \beta_\gamma} \log p(\mathbf{y}|\alpha, \beta_\gamma, \sigma^2) - q_\gamma H(g) \right\}\right),$$

where $H(g) = (2g)^{-1}(g+1)\log(g+1)$, a special case of the key relation in George and Foster (2000). As they point out, for particular values of g , when σ^2 is known, the Bayesian strategy of choosing \mathcal{M}_γ to maximize (1.6) corresponds to common fixed penalty selection criteria. For example, setting $H(g) = 2, \log n$ or $2 \log p$ (independently of \mathbf{y}) would correspond to AIC [Akaike (1974)], BIC [Schwarz (1978)] or RIC [Foster and George (1994)], BIC, or RIC, respectively. For a discussion of recommendations in the literature for choosing a fixed g depending on p and/or n , see Section 2.4 of Liang et al. (2008).

Although the correspondences to fixed penalty criteria are interesting, as a practical matter, it is necessary to deal with the uncertainty about g and σ^2 to obtain useful criteria. For this purpose, George and Foster (2000) proposed selecting the model maximizing $m_\gamma(\mathbf{y}|g, \sigma^2)$ based on an empirical Bayes estimate of g and the standard unbiased estimate of σ^2 . More recently, Cui and George (2008) proposed margining out g with respect to a prior, and Liang et al. (2008) proposed margining out g and σ^2 with respect to priors. It should be noted that the first paper to effectively use a prior integrating out g was Zellner and Siow (1980); they stated things in terms of multivariate Cauchy densities, which can always be expressed as a g -mixture of g -priors. All of these strategies lead to criteria that can be seen as adapting to the fixed penalty criterion which would be most suitable for the data at hand. In this paper, we shall similarly follow a fully Bayes approach, but with a generalization of the g -prior (1.5) and an extension of the considered class of priors on g .

After describing our prior forms in Section 2 and then calculating the marginals and Bayes factors in Section 3, we ultimately obtain our proposed g -prior Bayes factor (gBF), which is of the form (omitting the γ subscripts for clarity)

$$(1.7) \quad gBF_{\gamma:N} = \begin{cases} \left\{ \frac{\bar{d}}{d_q} \right\}^{-q} \frac{\{1 - R^2 + d_q^2 \|\hat{\boldsymbol{\beta}}_{LS}\|^2\}^{-1/4 - q/2}}{C_{n,q}(1 - R^2)^{(n-q)/2 - 3/4}}, & \text{if } q < n - 1, \\ \{\bar{d} \times \|\hat{\boldsymbol{\beta}}_{LS}^{MP}\|\}^{-n+1}, & \text{if } q \geq n - 1, \end{cases}$$

where $C_{n,q} \equiv \frac{B(1/4, (n-q)/2 - 3/4)}{B(q/2 + 1/4, (n-q)/2 - 3/4)}$ using the Beta function $B(\cdot, \cdot)$, R^2 is the familiar R -squared statistic under \mathcal{M}_γ , \bar{d} and d_r are, respectively, the geometric mean and minimum of the singular values of \mathbf{X}_γ , $\|\cdot\|$ is the L_2 norm, and finally, for the standardized response $(\mathbf{y} - \bar{y}\mathbf{1}_n)/\|\mathbf{y} - \bar{y}\mathbf{1}_n\|$, $\hat{\boldsymbol{\beta}}_{LS}$ is the usual least squares estimator, and $\hat{\boldsymbol{\beta}}_{LS}^{MP}$ is the least squares estimator using the Moore–Penrose inverse matrix.

Two immediately apparent features of (1.7) should be noted. First, in contrast to other fully Bayes factors for our selection problem, gBF is a closed form expression which allows for interpretation and straightforward calculation under any model. As will be seen in later sections, this transparency reveals that gBF not only rewards explained variation overall, but also rewards variation explained by the larger principal components of the design matrix. Second, gBF can be applied to all models even when the number of predictors p exceeds the number of observations n . This includes $p > n$ which is of increasing interest. This is not the case for (1.5) which requires $p \leq n - 1$ so that $\mathbf{X}'_\gamma \mathbf{X}_\gamma$ will be invertible for all q_γ , (recall that \mathbf{X}_γ has dimension at most $n - 1$ because its columns have been centered). Note also that when $p > n - 1$, penalized sum-of-squares criteria such as AIC, BIC and RIC will be unavailable for all submodels.

The organization of this paper is as follows. In Section 2 we propose prior forms including a generalized g -prior with a beta-prime prior for g . In Section 3 we

derive general Bayes factor expressions, and propose default hyperparameter settings which yield g BF above. In Section 4 we discuss appealing consequences of our default specifications. In Section 5 we describe conditional shrinkage estimation with the generalized g -prior. In Section 6 we show that g BF is consistent for model selection as $n \rightarrow \infty$. In Section 7 we provide a simulation evaluation of g BF performance.

2. A fully Bayes prior formulation. We now proceed to describe the prior components that form $p(\alpha, \beta_\gamma, \sigma^2)$ in (1.4). Throughout the remainder of the paper, we will omit the subscript γ for notational simplicity when there is no ambiguity. However, it is important to remember throughout that our formulations are to be applied to all of the 2^p possible submodels in (1.2).

2.1. *A generalized g -prior for β .* To motivate our proposed generalization of Zellner’s g -prior, we begin with a reconsideration of the original g -prior (1.5) for the case $p \leq n - 1$. The covariance matrix of the g -prior, $g\sigma^2(\mathbf{X}'\mathbf{X})^{-1}$, is proportional to the covariance matrix of the least squares estimator $\hat{\beta}_{LS}$. As a consequence of this choice, the marginal likelihood with respect to the g -prior appealingly becomes a function only of the residual sum-of-squares, RSS.

However, from the “matrix conditioning” viewpoint of Casella (1980, 1985) which advocates more shrinkage on higher variance estimates, the original g -prior may not be reasonable. To see why, let us rotate the problem by the $q \times q$ orthogonal matrix $\mathbf{W} = (\mathbf{w}_1, \dots, \mathbf{w}_q)$ which diagonalizes $\mathbf{X}'\mathbf{X}$ as

$$(2.1) \quad \mathbf{W}'(\mathbf{X}'\mathbf{X})\mathbf{W} = \mathbf{D}^2,$$

where $\mathbf{D} = \text{diag}(d_1, \dots, d_q)$ with

$$(2.2) \quad d_1 \geq \dots \geq d_q > 0.$$

Thus,

$$\mathbf{W}'\hat{\beta}_{LS} \sim N_q(\mathbf{W}'\beta, \sigma^2\mathbf{D}^{-2}).$$

Applying the g -prior (1.5) to these rotated coordinates would then induce the prior

$$\mathbf{W}'\beta \sim N_q(\mathbf{0}, g\sigma^2\mathbf{D}^{-2}),$$

which reveals the prior variances to be proportional to the sample variances of the elements of $\mathbf{W}'\hat{\beta}_{LS}$. This contradicts Casella (1980) who states, “if the sampling information is good, it is reasonable to downweight the prior guess.” To remedy this situation, we propose consideration of priors on β for which

$$\mathbf{W}'\beta \sim N_q(\mathbf{0}, \sigma^2\mathbf{\Psi}_q),$$

where the components of $\mathbf{\Psi}_q = \text{diag}(\psi_1, \dots, \psi_q)$ are in descending order, namely,

$$(2.3) \quad \psi_1 \geq \dots \geq \psi_q > 0.$$

Note that this would be satisfied for $\Psi_q \propto \mathbf{I}_q$, a consequence of the common assumption of exchangeable β components.

In fact, a slightly weaker ordering of the form

$$(2.4) \quad d_1^2 \psi_1 \geq \dots \geq d_q^2 \psi_q > 0$$

would still be reasonable because the resulting Bayes estimator of $\mathbf{w}'_i \beta$ would be of the form

$$(1 + \{d_i^2 \psi_i\}^{-1})^{-1} \mathbf{w}'_i \hat{\beta}_{LS},$$

so that under (2.4), the components of $\mathbf{W}' \hat{\beta}_{LS}$ with larger variance would be shrunk more. We note that the original g -prior (1.5), for which $\psi_i = g d_i^{-2}$, satisfies only the extreme boundary of (2.4), namely,

$$d_1^2 \psi_1 = \dots = d_q^2 \psi_q = g.$$

This violates (2.3) whenever $d_i > d_{i+1}$, in which case $\psi_i < \psi_{i+1}$.

An appealing general form for Ψ_q is $\Psi_q(g, \mathbf{v}) = \text{diag}(\psi_1(g, \mathbf{v}), \dots, \psi_q(g, \mathbf{v}))$, where

$$(2.5) \quad \psi_i(g, \mathbf{v}) = (1/d_i^2) \{v_i(1 + g) - 1\},$$

$\mathbf{v} = (v_1, \dots, v_q)'$ and $v_i \geq 1$ for any i , guaranteeing $\psi_i(g, \mathbf{v}) > 0$. Note that $\Psi_q(g, \mathbf{v})$, like the original g -prior, is controlled by a single hyperparameter $g > 0$. When $v_1 = \dots = v_q = 1$, $\sigma^2 \Psi_q(g, \mathbf{v})$ becomes $g \sigma^2 \mathbf{D}^{-2}$, yielding the covariance structure of the original g -prior. Although (2.4) will be satisfied whenever $v_1 \geq \dots \geq v_q \geq 1$, we shall ultimately be interested in a particular design dependent choice defined in Section 3.2. In summary, when $q \leq n - 1$, we propose a generalized g -prior for β of the form

$$(2.6) \quad p(\beta | \sigma^2, g) = \phi_q(\mathbf{W}' \beta; \mathbf{0}, \sigma^2 \Psi_q(g, \mathbf{v})),$$

where $v_1 \geq \dots \geq v_q \geq 1$.

When $q > n - 1$ and the rank of \mathbf{X} is $n - 1$, there exists a $q \times (n - 1)$ matrix $\mathbf{W} = (\mathbf{w}_1, \dots, \mathbf{w}_{n-1})$ which diagonalizes $\mathbf{X}'\mathbf{X}$ as

$$(2.7) \quad \mathbf{W}'(\mathbf{X}'\mathbf{X})\mathbf{W} = \mathbf{D}^2,$$

where $\mathbf{W}'\mathbf{W} = \mathbf{I}_{n-1}$ and $\mathbf{D} = \text{diag}(d_1, d_2, \dots, d_{n-1})$ with $d_1 \geq d_2 \geq \dots \geq d_{n-1} > 0$. For this case, we propose a generalized g -prior of the form

$$(2.8) \quad p(\beta | \sigma^2, g) = \phi_{n-1}(\mathbf{W}' \beta; \mathbf{0}, \sigma^2 \Psi_{n-1}(g, \mathbf{v})) p_{\#}(\mathbf{W}'_{\#} \beta),$$

where $\Psi_{n-1}(g, \mathbf{v}) = \text{diag}(\psi_1, \dots, \psi_{n-1})$ is again given by (2.5) and $v_1 \geq \dots \geq v_{n-1} \geq 1$. Here, $\mathbf{W}_{\#}$ is an arbitrary matrix which makes the $q \times q$ matrix $(\mathbf{W}, \mathbf{W}_{\#})$ orthogonal, and $p_{\#}(\cdot)$ is an arbitrary probability density on $\mathbf{W}'_{\#} \beta$, respectively. As will be seen, the choices of $\mathbf{W}_{\#}$ and $p_{\#}$ have no effect on the selection criteria we obtain, thus we leave them as arbitrary.

Combining the above two cases by letting

$$(2.9) \quad r = \min\{q, n - 1\},$$

our suggested generalized g -prior is of the form

$$(2.10) \quad p(\boldsymbol{\beta}|g, \sigma^2) = \phi_r(\mathbf{W}'\boldsymbol{\beta}; \mathbf{0}, \sigma^2\Psi_r(g, \mathbf{v})) \times \begin{cases} 1, & \text{if } q \leq n - 1, \\ p_{\#}(\mathbf{W}'_{\#}\boldsymbol{\beta}), & \text{if } q > n - 1, \end{cases}$$

where the $q \times r$ matrix \mathbf{W} satisfies both $\mathbf{W}'\mathbf{X}'\mathbf{X}\mathbf{W} = \text{diag}(d_1^2, \dots, d_r^2)$ and $\mathbf{W}'\mathbf{W} = \mathbf{I}_r$, and $\Psi_r(g, \mathbf{v}) = \text{diag}(\psi_1(g, \mathbf{v}), \dots, \psi_r(g, \mathbf{v}))$ with (2.5).

REMARK 2.1. In (2.1) and (2.7), let

$$(2.11) \quad \mathbf{U} = (\mathbf{u}_1, \dots, \mathbf{u}_r) = (\mathbf{X}\mathbf{w}_1/d_1, \dots, \mathbf{X}\mathbf{w}_r/d_r) = \mathbf{X}\mathbf{W}\mathbf{D}^{-1}.$$

Then $\mathbf{U}'\mathbf{U} = \mathbf{I}_r$ and

$$(2.12) \quad \mathbf{X} = \mathbf{U}\mathbf{D}\mathbf{W}' = \sum_{i=1}^r d_i \mathbf{u}_i \mathbf{w}'_i.$$

This is the nonnull part of the well-known singular value decomposition (SVD). The diagonal elements of $\mathbf{D} = \text{diag}(d_1, \dots, d_r)$ are the singular values of \mathbf{X} , and the columns of $\mathbf{U} = (\mathbf{u}_1, \dots, \mathbf{u}_r)$ are the normalized principal components of the column space of \mathbf{X} . Note that the components of the rotated vector $\mathbf{W}'\boldsymbol{\beta}$ are the coefficients for the principal component regression of \mathbf{y} on $\mathbf{U}\mathbf{D}$. From the definition of \mathbf{W} and \mathbf{U} by (2.1), (2.7) and (2.11), the signs of $\mathbf{u}_i \mathbf{w}'_i$ are determinate although the signs of \mathbf{w}_i and \mathbf{u}_i for $1 \leq i \leq r$ are indeterminate. These indeterminacies can safely be ignored in our development.

2.2. *Priors for g, α and σ^2 .* Turning to the prior for the hyperparameter g , we propose

$$(2.13) \quad p(g) = \frac{g^b (1 + g)^{-a-b-2}}{B(a + 1, b + 1)} I_{(0, \infty)}(g)$$

with $a > -1, b > -1$, a Pearson Type VI or *beta-prime* distribution under which $1/(1 + g)$ has a Beta distribution $\text{Be}(a + 1, b + 1)$. Choices for the hyperparameters a and b are discussed later.

Although Zellner and Siow (1980) did not explicitly use a g -prior formulation with a prior on g , their recommendation of a multivariate Cauchy form for $p(\boldsymbol{\beta}|\sigma^2)$ implicitly corresponds to using a g -prior with an inverse Gamma prior

$$(n/2)^{1/2} \{\Gamma(1/2)\}^{-1} g^{-3/2} e^{-n/(2g)}$$

on g . Both Cui and George (2008) and Liang et al. (2008) proposed using g -priors with priors of the form

$$(2.14) \quad p(g) = (a + 1)^{-1} (1 + g)^{-a-2},$$

the subclass of (2.13) with $b = 0$. Cases for which $b = O(n)$ will be of interest to us in what follows.

For the parameter α and σ^2 , we use the location invariant flat prior

$$(2.15) \quad p(\alpha) = I_{(-\infty, \infty)}(\alpha)$$

and the scale invariant prior

$$(2.16) \quad p(\sigma^2) = (\sigma^2)^{-1} I_{(0, \infty)}(\sigma^2),$$

respectively. Because α and σ^2 appear in every model, the use of these improper priors for Bayesian model selection is formally justified by Berger, Pericchi and Varshavsky (1998).

We note in passing that for the estimation of a multivariate normal mean, priors equivalent to (2.6), (2.13), (2.15) and (2.16) have been considered by Strawderman (1971) and extended by Maruyama and Strawderman (2005).

3. Marginal densities and Bayes factors.

3.1. *General forms.* The marginal densities of \mathbf{y} under $\mathcal{M}_\gamma (\neq \mathcal{M}_N)$ and \mathcal{M}_N are, by definition,

$$(3.1) \quad \begin{aligned} m_\gamma(\mathbf{y}) &= \int_{-\infty}^{\infty} \int_{R^q} \int_0^{\infty} p(\mathbf{y}|\alpha, \boldsymbol{\beta}_\gamma, \sigma^2) p(\alpha, \boldsymbol{\beta}_\gamma, \sigma^2) d\alpha d\boldsymbol{\beta}_\gamma d\sigma^2, \\ m_N(\mathbf{y}) &= \int_{-\infty}^{\infty} \int_0^{\infty} p(\mathbf{y}|\alpha, \sigma^2) p(\alpha, \sigma^2) d\alpha d\sigma^2, \end{aligned}$$

respectively. Under the priors

$$p(\alpha, \boldsymbol{\beta}_\gamma, \sigma^2) = p(\alpha) p(\sigma^2) \int_0^{\infty} p(\boldsymbol{\beta}_\gamma|\sigma^2, g) p(g) dg \quad \text{for } \mathcal{M}_\gamma (\neq \mathcal{M}_N)$$

and

$$p(\alpha, \sigma^2) = p(\alpha) p(\sigma^2) \quad \text{for } \mathcal{M}_N,$$

where $p(\boldsymbol{\beta}|\sigma^2, g)$, $p(\alpha)$ and $p(\sigma^2)$ are given by (2.10), (2.15) and (2.16), and $p(g)$ when $q < n - 1$ is given by (2.13) with $-1 < a < -1/2$ and $b = (n - 5)/2 - q/2 - a$ [$p(g)$ is arbitrary when $q \geq n - 1$], we have a following theorem about the Bayes factor ratio of the marginal densities under each of \mathcal{M}_γ and \mathcal{M}_N .

THEOREM 3.1. *The Bayes factor for comparing each of \mathcal{M}_Y to \mathcal{M}_N is*

$$(3.2) \quad \text{BF}_{Y:N}(a, \mathbf{v}) = \frac{m_Y(\mathbf{y})}{m_N(\mathbf{y})} = \begin{cases} \prod_{i=1}^q v_i^{-1/2} \frac{B(q/2 + a + 1, (n - q - 3)/2 - a)}{B(a + 1, (n - q - 3)/2 - a)} \\ \quad \times \frac{(1 - Q^2)^{-q/2 - a - 1}}{(1 - R^2)^{(n - q - 3)/2 - a}}, & \text{if } q < n - 1, \\ \prod_{i=1}^{n-1} v_i^{-1/2} (1 - Q^2)^{-(n-1)/2}, & \text{if } q \geq n - 1, \end{cases}$$

where $v_1 \geq \dots \geq v_r \geq 1$, R^2 and Q^2 are given by

$$(3.3) \quad R^2 = \sum_{i=1}^r \{\text{cor}(\mathbf{u}_i, \mathbf{y})\}^2, \quad Q^2 = \sum_{i=1}^r (1 - v_i^{-1}) \{\text{cor}(\mathbf{u}_i, \mathbf{y})\}^2.$$

Note that R^2 and Q^2 are the usual and a modified version of the R -squared statistics and $\text{cor}(\mathbf{u}_i, \mathbf{y})$ is the correlation of the response \mathbf{y} and the i th principal component of \mathbf{X} .

PROOF OF THEOREM 3.1. Defining $\mathbf{v} = \mathbf{y} - \bar{y}\mathbf{1}_n$, where \bar{y} is the mean of \mathbf{y} , so that

$$\|\mathbf{y} - \alpha\mathbf{1}_n - \mathbf{X}\boldsymbol{\beta}\|^2 = n(-\alpha + \bar{y})^2 + \|\mathbf{v} - \mathbf{X}\boldsymbol{\beta}\|^2,$$

we obtain

$$(3.4) \quad \int_{-\infty}^{\infty} p(\mathbf{y}|\alpha, \boldsymbol{\beta}, \sigma^2) d\alpha = \frac{n^{1/2}}{(2\pi\sigma^2)^{(n-1)/2}} \exp\left(-\frac{\|\mathbf{v} - \mathbf{X}\boldsymbol{\beta}\|^2}{2\sigma^2}\right).$$

We make the following orthogonal transformation when integration with respect to $\boldsymbol{\beta}$ is considered:

$$(3.5) \quad \boldsymbol{\beta} \rightarrow \begin{cases} \mathbf{W}'\boldsymbol{\beta} \equiv \boldsymbol{\beta}_*, & \text{if } q \leq n - 1, \\ \begin{pmatrix} \mathbf{W}'\boldsymbol{\beta} \\ \mathbf{W}'_{\#}\boldsymbol{\beta} \end{pmatrix} \equiv \begin{pmatrix} \boldsymbol{\beta}_* \\ \boldsymbol{\beta}_{\#} \end{pmatrix}, & \text{if } q > n - 1, \end{cases}$$

so that

$$\begin{aligned} & \int_{-\infty}^{\infty} \int_{R^q} p(\mathbf{y}|\alpha, \boldsymbol{\beta}, \sigma^2) p(\boldsymbol{\beta}|\sigma^2, g) d\alpha d\boldsymbol{\beta} \\ &= \frac{n^{1/2}}{(2\pi\sigma^2)^{(n-1)/2}} \frac{|\boldsymbol{\Psi}|^{-1/2}}{(2\pi\sigma^2)^{r/2}} \int_{R^r} \exp\left(-\frac{\|\mathbf{v} - \mathbf{U}\mathbf{D}\boldsymbol{\beta}_*\|^2}{2\sigma^2} - \frac{\boldsymbol{\beta}'_{*}\boldsymbol{\Psi}^{-1}\boldsymbol{\beta}_*}{2\sigma^2}\right) d\boldsymbol{\beta}_* \\ & \quad \times \begin{cases} 1, & \text{if } q \leq n - 1, \\ \int_{R^{q-n+1}} p_{\#}(\boldsymbol{\beta}_{\#}) d\boldsymbol{\beta}_{\#} (=1), & \text{if } q > n - 1. \end{cases} \end{aligned}$$

Completing the square $\|\mathbf{v} - \mathbf{UD}\boldsymbol{\beta}_*\|^2 + \boldsymbol{\beta}'_*\boldsymbol{\Psi}^{-1}\boldsymbol{\beta}_*$ with respect to $\boldsymbol{\beta}_*$, we have

$$\begin{aligned}
 & \|\mathbf{v} - \mathbf{UD}\boldsymbol{\beta}_*\|^2 + \boldsymbol{\beta}'_*\boldsymbol{\Psi}^{-1}\boldsymbol{\beta}_* \\
 &= \{\boldsymbol{\beta}_* - (\mathbf{D}^2 + \boldsymbol{\Psi}^{-1})^{-1}\mathbf{D}'\mathbf{U}'\mathbf{v}\}'(\mathbf{D}^2 + \boldsymbol{\Psi}^{-1}) \\
 (3.6) \quad & \times \{\boldsymbol{\beta}_* - (\mathbf{D}^2 + \boldsymbol{\Psi}^{-1})^{-1}\mathbf{D}'\mathbf{U}'\mathbf{v}\} \\
 & - \mathbf{v}'\mathbf{UD}(\mathbf{D}^2 + \boldsymbol{\Psi}^{-1})^{-1}\mathbf{D}'\mathbf{U}'\mathbf{v} + \mathbf{v}'\mathbf{v},
 \end{aligned}$$

where the residual term is rewritten as

$$\begin{aligned}
 & -\mathbf{v}'\mathbf{UD}(\mathbf{D}^2 + \boldsymbol{\Psi}^{-1})^{-1}\mathbf{D}'\mathbf{U}'\mathbf{v} + \mathbf{v}'\mathbf{v} \\
 &= -\mathbf{v}'\left(\sum_{i=1}^r \mathbf{u}_i\mathbf{u}'_i \frac{d_i^2}{d_i^2 + \psi_i^{-1}}\right)\mathbf{v} + \mathbf{v}'\mathbf{v} \\
 &= \frac{g\|\mathbf{v}\|^2}{g+1}\left\{1 - \sum_{i=1}^r \frac{(\mathbf{u}'_i\mathbf{v})^2}{\|\mathbf{v}\|^2}\right\} + \frac{\|\mathbf{v}\|^2}{1+g}\left\{1 - \sum_{i=1}^r \left(1 - \frac{1}{v_i}\right) \frac{(\mathbf{u}'_i\mathbf{v})^2}{\|\mathbf{v}\|^2}\right\}.
 \end{aligned}$$

Hence, by

$$|\boldsymbol{\Psi}| = \prod_{i=1}^r \frac{v_i + v_i g - 1}{d_i^2}, \quad |\mathbf{D}^2 + \boldsymbol{\Psi}^{-1}| = \prod_{i=1}^r \frac{d_i^2 v_i (1 + g)}{v_i + v_i g - 1},$$

we have

$$\begin{aligned}
 (3.7) \quad & \int_{-\infty}^{\infty} \int_{R^q} p(\mathbf{y}|\boldsymbol{\alpha}, \boldsymbol{\beta}, \sigma^2) p(\boldsymbol{\beta}|g, \sigma^2) d\boldsymbol{\alpha} d\boldsymbol{\beta} \\
 &= \frac{n^{1/2}}{(2\pi\sigma^2)^{(n-1)/2}} \frac{(1+g)^{-r/2}}{\prod_{i=1}^r v_i^{1/2}} \exp\left(-\frac{\|\mathbf{v}\|^2\{g(1-R^2) + 1 - Q^2\}}{2\sigma^2(g+1)}\right),
 \end{aligned}$$

where R^2 and Q^2 are given by (3.3).

Next we consider the integration with respect to σ^2 . By (3.7), we have

$$\begin{aligned}
 (3.8) \quad & \int_{-\infty}^{\infty} \int_{R^q} \int_0^{\infty} p(\mathbf{y}|\boldsymbol{\alpha}, \boldsymbol{\beta}, \sigma^2) p(\boldsymbol{\beta}|g, \sigma^2) \frac{1}{\sigma^2} d\boldsymbol{\alpha} d\boldsymbol{\beta} d\sigma^2 \\
 &= \int_0^{\infty} \frac{n^{1/2}}{(2\pi\sigma^2)^{(n-1)/2}} \frac{(1+g)^{-r/2}}{\prod_{i=1}^r v_i^{1/2}} \\
 & \quad \times \exp\left(-\frac{\|\mathbf{v}\|^2\{g(1-R^2) + 1 - Q^2\}}{2\sigma^2(g+1)}\right) \frac{1}{\sigma^2} d\sigma^2 \\
 &= \frac{K(n, \mathbf{y})}{\prod_{i=1}^r v_i^{1/2}} (1+g)^{-r/2+(n-1)/2} \{g(1-R^2) + 1 - Q^2\}^{-(n-1)/2},
 \end{aligned}$$

where

$$K(n, \mathbf{y}) = \frac{n^{1/2}\Gamma(\{n-1\}/2)}{\pi^{(n-1)/2}\|\mathbf{y} - \bar{y}\mathbf{1}_n\|^{n-1}}.$$

When $q \geq n - 1$, $R^2 = 1$ and $r = n - 1$ so that

$$(3.9) \quad \int_{-\infty}^{\infty} \int_{R^q} \int_0^{\infty} p(\mathbf{y}|\alpha, \boldsymbol{\beta}, \sigma^2) p(\boldsymbol{\beta}|g, \sigma^2) \frac{1}{\sigma^2} d\alpha d\boldsymbol{\beta} d\sigma^2 = \frac{K(n, \mathbf{y})}{\prod_{i=1}^{n-1} v_i^{1/2}} \{1 - Q^2\}^{-(n-1)/2},$$

which does not depend on g . Hence, in this case, $m_\gamma(\mathbf{y})$ does not depend on the prior density of g .

When $q < n - 1$, we consider the prior (2.13) of g with $-1 < a < -1/2$ and $b = (n - 5)/2 - q/2 - a$, where b is guaranteed to be strictly greater than -1 for $q < n - 1$. Then we have

$$(3.10) \quad \begin{aligned} m_\gamma(\mathbf{y}) &= \frac{K(n, \mathbf{y})}{\prod_{i=1}^q v_i^{1/2} B(a + 1, b + 1)} \\ &\quad \times \int_0^{\infty} \frac{g^b}{(1 + g)^{a+b+2}} \frac{\{g(1 - R^2) + 1 - Q^2\}^{-(n-1)/2}}{(1 + g)^{q/2 - (n-1)/2}} dg \\ &= \frac{K(n, \mathbf{y})(1 - Q^2)^{-(n-1)/2}}{\prod_{i=1}^q v_i^{1/2} B(a + 1, b + 1)} \int_0^{\infty} g^b \left(\frac{1 - R^2}{1 - Q^2} g + 1 \right)^{-(n-1)/2} dg \\ &= \frac{K(n, \mathbf{y})(1 - Q^2)^{-(n-1)/2+b+1}}{\prod_{i=1}^q v_i^{1/2} \{1 - R^2\}^{b+1}} \frac{B(q/2 + a + 1, b + 1)}{B(a + 1, b + 1)} \\ &= \frac{K(n, \mathbf{y})(1 - Q^2)^{-q/2-a-1}}{\prod_{i=1}^q v_i^{1/2} \{1 - R^2\}^{(n-q-3)/2-a}} \frac{B(q/2 + a + 1, (n - q - 3)/2 - a)}{B(a + 1, (n - q - 3)/2 - a)}. \end{aligned}$$

In the same way, $m_N(\mathbf{y})$ for the null model is obtained as

$$(3.11) \quad m_N(\mathbf{y}) = K(n, \mathbf{y}).$$

From (3.9), (3.10) and (3.11), the theorem follows. \square

REMARK 3.1. R^2 and Q^2 given by (3.3) are the usual and a modified form of the R -squared measure for multiple regression. They are here expressed in terms of $\{\text{cor}(\mathbf{u}_1, \mathbf{y})\}^2, \dots, \{\text{cor}(\mathbf{u}_r, \mathbf{y})\}^2$, the squared correlations of the response \mathbf{y} and the principal components $\mathbf{u}_1, \dots, \mathbf{u}_r$ of \mathbf{X} . For fixed q and \mathbf{v} , the BF criterion is increasing in both R^2 and Q^2 . The former is definitely reasonable. Larger Q^2 would also be reasonable when $v_1 \geq \dots \geq v_r$ so that Q^2 would put more weight on those components of $\mathbf{W}'\boldsymbol{\beta}$ for which d_i is larger and are consequently better estimated. In this sense, Q^2 would reward those models which are more stably estimated.

Beyond their influence through Q^2 , the choice of v_1, \dots, v_r plays a further influential role in $\text{BF}_{\gamma:N}$ through the $\prod_{i=1}^r v_i^{-1/2}$ terms in (3.2). In Section 3.2 below,

a default choice is proposed which, through these terms, rewards stable estimation. Note that if $v_i = 1$ for all i (i.e., the original g -prior), Q^2 becomes zero, $\prod_{i=1}^r v_i^{-1/2} \equiv 1$, and $\text{BF}_{\gamma:N}$ becomes a function of just R^2 and q . In this case, $\text{BF}_{\gamma:N}$ will not distinguish between models for which $q \geq n - 1$.

REMARK 3.2. The analytical simplification in (3.10) is a consequence of the choice $b = (n - 5)/2 - q/2 - a$, and results in a convenient closed form for our Bayes factor. Such a reduction is unavailable for other choices of b . For example, Liang et al. (2008) use Laplace approximations to avoid the evaluation of the special functions that arise in the resulting Bayes factor when $b = 0$. Another attractive feature of the choice $b = (n - 5)/2 - q/2 - a$ will be discussed in Section 4.2.

3.2. *Default choices.* At this point, we are ready to consider default choices for a and \mathbf{v} . For a , we recommend

$$(3.12) \quad a = -3/4,$$

the median of the range of values $(-1, -1/2)$ for which the marginal density is well defined for any choices of $q < n - 1$. In Section 4 we will explicitly see the appealing consequence of this choice on the asymptotic tail behavior of $p(\boldsymbol{\beta}|\sigma^2)$.

For \mathbf{v} , we recommend

$$(3.13) \quad \mathbf{v} = (d_1^2/d_r^2, d_2^2/d_r^2, \dots, 1)',$$

which coupled with (2.5) satisfies (2.4) since $v_1 \geq \dots \geq v_q \geq 1$ for this choice. Inserting this \mathbf{v} into (3.3) yields

$$(3.14) \quad \begin{aligned} Q^2 &= R^2 - d_r^2 \sum_{i=1}^r \frac{(\mathbf{u}'_i \mathbf{v})^2}{d_i^2 \mathbf{v}' \mathbf{v}} \\ &= R^2 - d_r^2 \|\mathbf{D}^{-1} \mathbf{U}' \{\mathbf{v}/\|\mathbf{v}\|\}\|^2 \\ &= \begin{cases} R^2 - d_q^2 \|\hat{\boldsymbol{\beta}}_{\text{LS}}\|^2, & \text{if } q < n - 1, \\ 1 - d_{n-1}^2 \|\hat{\boldsymbol{\beta}}_{\text{LS}}^{\text{MP}}\|^2, & \text{if } q \geq n - 1, \end{cases} \end{aligned}$$

where, for the standardized response $\mathbf{v}/\|\mathbf{v}\|$ for $\mathbf{v} = \mathbf{y} - \bar{y}\mathbf{1}_n$, $\hat{\boldsymbol{\beta}}_{\text{LS}}$ is the usual LS estimator for $q < n - 1$, and $\hat{\boldsymbol{\beta}}_{\text{LS}}^{\text{MP}}$ is the LS estimator based on the Moore–Penrose inverse matrix. The third equality in (3.14) follows from the fact that both $\hat{\boldsymbol{\beta}}_{\text{LS}}$ and $\hat{\boldsymbol{\beta}}_{\text{LS}}^{\text{MP}}$ for the response $\mathbf{v}/\|\mathbf{v}\|$ can be expressed as

$$\hat{\boldsymbol{\beta}} = \mathbf{W} \mathbf{D}^{-1} \mathbf{U}' \{\mathbf{v}/\|\mathbf{v}\|\},$$

and from the orthogonality of \mathbf{W} ,

$$\|\hat{\boldsymbol{\beta}}\|^2 = \|\mathbf{D}^{-1} \mathbf{U}' \{\mathbf{v}/\|\mathbf{v}\|\}\|^2.$$

It will also be useful to define

$$(3.15) \quad \bar{d} = \left(\prod_{i=1}^r d_i \right)^{1/r},$$

the geometric mean of the singular values d_1, \dots, d_r . Inserting our default choices for a and \mathbf{v} into $\text{BF}_{\gamma:N}(a, \mathbf{v})$ in (3.2), and noting that

$$(3.16) \quad \prod_{i=1}^r v_i^{-1/2} = (\bar{d}/d_r)^{-r},$$

we obtain our recommended Bayes factor in (1.7) which we denote by $g\text{BF}$ (g -prior Bayes factor):

$$(3.17) \quad g\text{BF}_{\gamma:N} = \begin{cases} \left\{ \frac{\bar{d}}{d_q} \right\}^{-q} \frac{B(q/2 + 1/4, (n - q)/2 - 3/4)}{B(1/4, (n - q)/2 - 3/4)} \\ \quad \times \frac{(1 - R^2 + d_q^2 \|\hat{\boldsymbol{\beta}}_{\text{LS}}\|^2)^{-1/4 - q/2}}{(1 - R^2)^{(n - q)/2 - 3/4}}, & \text{if } q < n - 1, \\ \{\bar{d} \times \|\hat{\boldsymbol{\beta}}_{\text{LS}}^{\text{MP}}\|\}^{-(n-1)}, & \text{if } q \geq n - 1, \end{cases}$$

which is a function of the key quantities q , R^2 , the LS estimators and the singular values of the design matrix.

REMARK 3.3. Like traditional selection criteria such as AIC, BIC and RIC, the $g\text{BF}$ criterion (3.17) rewards models for explained variation through R^2 . However, $g\text{BF}$ also rewards models for stability of estimation through smaller values of \bar{d}/d_q and $d_q \|\hat{\boldsymbol{\beta}}_{\text{LS}}\|$ for $q < n - 1$, and through smaller values of the product \bar{d}/d_{n-1} and $d_{n-1} \|\hat{\boldsymbol{\beta}}_{\text{LS}}^{\text{MP}}\|$ for $q \geq n - 1$, the case where R^2 is unavailable.

To see how these various quantities bear on stable estimation, note first that

$$(3.18) \quad \bar{d}/d_r = \left\{ \prod_{i=1}^r (d_i/d_r) \right\}^{1/r},$$

which gets smaller as the d_i/d_r ratios get smaller. Like the well-known condition number d_1/d_r , smaller values of (3.18) indicate a more stable design matrix \mathbf{X}_γ .

For $d_q \|\hat{\boldsymbol{\beta}}_{\text{LS}}\|$ and $d_{n-1} \|\hat{\boldsymbol{\beta}}_{\text{LS}}^{\text{MP}}\|$, note that each of these can be expressed as

$$(3.19) \quad d_r^2 \|\hat{\boldsymbol{\beta}}\|^2 = \sum_{i=1}^r \left(\frac{d_r}{d_i} \right)^2 \left\{ \frac{(\mathbf{u}'_i \mathbf{v})}{\|\mathbf{u}_i\| \|\mathbf{v}\|} \right\}^2 = \sum_{i=1}^r \left(\frac{d_r}{d_i} \right)^2 \{\text{cor}(\mathbf{u}_i, \mathbf{y})\}^2.$$

Thus, for a given set of d_i/d_r ratios, (3.19) gets smaller if the larger correlations $\text{cor}(\mathbf{u}_i, \mathbf{y})$ correspond to the larger d_i . Again, this is a measure of stability, as the largest principal components $d_i \mathbf{u}_i$ are the ones which are most stably estimated.

REMARK 3.4. The choice of ν in (3.13) will be especially sensitive to small values of d_r which would lead to large prior variances in (2.10). Thus, one bad x_i predictor variable could spoil the model. From an estimation point of view, this perhaps would be unwise. However, from a model selection point of view, the effect of a small d_r would have the effect of downweighting the model, through the stability measures discussed in Remark 3.3, in favor of models which left out the offending x_i . Thus, any unstable submodel with at least one such x_i , but possibly more, would be downweighted.

4. The effect of the default choices of a and b . In Section 3 we proposed the prior form $p(g)$ given by (2.13) with hyperparameters a and b , recommending the choices $a = -3/4$ and $b = (n - q - 5)/2 - a$ for the case $q < n - 1$ where the prior on g matters. In the following subsections, we show some appealing consequences of these choices.

4.1. *The effect of a on the tail behavior of $p(\beta|\sigma^2)$.* Combining $p(\beta|g, \sigma^2)$ in (2.10) with $p(g)$ in (2.13), the probability density of β given σ^2 is given by

$$(4.1) \quad p(\beta|\sigma^2) = \int_0^\infty \frac{\phi_q(\mathbf{W}'\beta; \mathbf{0}, \sigma^2\Psi_q(g, \nu))}{B(a + 1, b + 1)} \frac{g^b}{(1 + g)^{a+b+2}} dg.$$

To examine the asymptotic behavior of the density $p(\beta|\sigma^2)$ as $\|\beta\| \rightarrow \infty$, we appeal to the Tauberian theorem for the Laplace transform [see Geluk and de Haan (1987)], which tells us that the contribution of the integral (4.2) around zero becomes negligible as $\|\beta\| \rightarrow \infty$. Thus, we have only to consider the integration between ν_1 and ∞ (the major term).

Since $d_1 \geq \dots \geq d_q$, and assuming $\nu_1 \geq \dots \geq \nu_q$, we have

$$(4.2) \quad \frac{d_q^2}{(\nu_1 + 1)g} \leq \frac{d_i^2}{\nu_i + \nu_i g - 1} \leq \frac{d_1^2}{\nu_q g}$$

for $g \geq \nu_1$ and any i , which implies

$$\begin{aligned} & C \frac{d_q^q}{(\nu_1 + 1)^{q/2}} \int_{\nu_1}^\infty \left(\frac{g}{g + 1}\right)^{a+b+2} \left(\frac{1}{g}\right)^{q/2+a+2} \exp\left(-\frac{1}{g} \frac{d_1^2 \|\mathbf{W}'\beta\|^2}{2\nu_q \sigma^2}\right) dg \\ & \leq \text{the major term of } p(\beta|\sigma^2) \\ & \leq C \frac{d_1^q}{\nu_q^{q/2}} \int_{\nu_1}^\infty \left(\frac{g}{g + 1}\right)^{a+b+2} \left(\frac{1}{g}\right)^{q/2+a+2} \exp\left(-\frac{1}{g} \frac{d_q^2 \|\mathbf{W}'\beta\|^2}{2(\nu_1 + 1)\sigma^2}\right) dg, \end{aligned}$$

where $C = \{B(a + 1, b + 1)\}^{-1} (2\pi\sigma^2)^{-q/2}$. Thus, by the Tauberian theorem, there exist $C_1 < C_2$ such that

$$(4.3) \quad C_1 < \frac{\|\beta\|^{q+2a+2}}{(\sigma^2)^{a+1}} p(\beta|\sigma^2) < C_2$$

for sufficiently large $\|\beta\|$.

From (4.3), we see that the asymptotic tail behavior of $p(\beta|\sigma^2)$ is determined by a and unaffected by b . Smaller a yields flatter tail behavior, thereby diminishing the prior influence of $p(\beta|\sigma^2)$. For $a = -1/2$ the asymptotic tail behavior of $p(\beta|\sigma^2)$, $\|\beta\|^{-q-1}$, corresponds to that of multivariate Cauchy distribution recommended by Zellner and Siow (1980). In contrast, the asymptotic tail behavior of our choice $a = -3/4$, $\|\beta\|^{-q-1/2}$, is even flatter than that of the multivariate Cauchy distribution.

4.2. *The effect of b on the implicit $O(n)$ choice of g .* For implementations of the original g -prior (1.5), Zellner (1986) and others have recommended choices for which $g = O(n)$. This prevents the g -prior from asymptotically dominating the likelihood which would occur if g was unchanged as n increased. The recommendation of choosing $g = O(n)$ also applies to the choice of a fixed g for the generalized g -prior (2.10) where

$$\text{tr}\{\text{Var}(\beta|g, \sigma^2)\} = \sigma^2 \sum_{i=1}^q \frac{v_i + v_i g - 1}{d_i^2}.$$

Since $d_i^2 = O(n)$ for $1 \leq i \leq q$ by Lemma B.1, $\text{tr}\{\text{Var}(\beta|g, \sigma^2)\} = gO(n^{-1})$ if v_i is bounded. Therefore, the choice $g = O(n)$ will also prevent the generalized g -prior from asymptotically dominating the likelihood, and stabilize it in the sense that $\text{tr}\{\text{Var}(\beta|g, \sigma^2)\} = O(1)$ when $g = O(n)$.

For our fully Bayes case, where g is treated as a random variable, our choice of b , in addition to yielding a closed form for the marginal density in (3.10), also yields an implicit $O(n)$ choice of g , in the sense that

$$\begin{aligned} [\text{mode of } g] &= \frac{b}{a+2} = \frac{2(n-q)-7}{5}, \\ \frac{1}{E[g^{-1}]} &= \frac{b}{a+1} = 2(n-q)-7 \end{aligned}$$

for our recommended choices $a = -3/4$ and $b = (n - q - 5)/2 - a$. (Note that $E[g]$ does not exist under the choice $a = -3/4$.)

5. Shrinkage estimation conditionally on a model. In this section we consider estimation conditionally on a model \mathcal{M}_γ . Because β is not identifiable when $q > n - 1$, and hence not estimable, we instead focus on estimation of $\mathbf{X}\beta$, which is always estimable. For this purpose, we consider estimation of $\mathbf{X}\beta$ under scaled quadratic loss $(\delta - \mathbf{X}\beta)' \mathbf{Q}(\delta - \mathbf{X}\beta)/\sigma^2$ for positive-definite \mathbf{Q} . The Bayes estimator under this loss for any \mathbf{Q} is of the form

$$(5.1) \quad \mathbf{X}\hat{\beta}_B = \mathbf{X}E[\sigma^{-2}\beta|\mathbf{y}]/E[\sigma^{-2}|\mathbf{y}].$$

From calculations similar to those in Section 3, under our priors given in Section 2, a simple closed form can be obtained for this estimator as follows. In contrast, such a simple closed form is not available for the usual Bayes estimator, $\mathbf{X}E[\boldsymbol{\beta}_\gamma | \mathbf{y}]$, the posterior mean under $(\boldsymbol{\delta} - \mathbf{X}\boldsymbol{\beta})' \mathbf{Q}(\boldsymbol{\delta} - \mathbf{X}\boldsymbol{\beta})$ which does not scale for the variance σ^2 .

THEOREM 5.1. *The Bayes estimator under scaled quadratic loss is given by*

$$(5.2) \quad \mathbf{X}\hat{\boldsymbol{\beta}}_B = \sum_{i=1}^r (1 - H(\mathbf{y})/v_i)(\mathbf{u}'_i \mathbf{y}) \mathbf{u}_i,$$

where

$$H(\mathbf{y}) = \begin{cases} \left(1 + \frac{1 - R^2}{1 - R^2} \frac{(n - q - 3)/2 - a}{q/2 + a + 1}\right)^{-1}, & q < n - 1, \\ \{1 + E[g]\}^{-1}, & q \geq n - 1. \end{cases}$$

PROOF. See the [Appendix](#). \square

Thus, when $q \geq n - 1$, we must specify the mean of prior density of g , although no such specification was needed for model selection. A reasonable specification may be $E[g] = d_{n-1}^2/d_1^2$, a function of the condition number d_1/d_{n-1} of the linear equation. For extremely large values of d_1/d_{n-1} , the coefficients of the first and the last terms in (5.2) become nearly 1 and 0, respectively. See [Casella \(1985\)](#) and [Maruyama and Strawderman \(2005\)](#) for further discussion of the condition number.

Thus, for our recommended choices of hyperparameters $a = -3/4$ and $v_i = d_i^2/d_r^2$ for $1 \leq i \leq r$, our recommended estimator of $\mathbf{X}\boldsymbol{\beta}$ for a given model \mathcal{M}_γ is

$$(5.3) \quad \mathbf{X}\hat{\boldsymbol{\beta}}_B = \sum_{i=1}^r (1 - \{d_r^2/d_i^2\}H(\mathbf{y}))(\mathbf{u}'_i \mathbf{y}) \mathbf{u}_i,$$

where

$$(5.4) \quad H(\mathbf{y}) = \begin{cases} \left(1 + \frac{1 - R^2 + d_q^2 \|\hat{\boldsymbol{\beta}}_{LS}\|^2}{1 - R^2} \frac{n/2 - q/2 - 3/4}{q/2 + 1/4}\right)^{-1}, & \text{if } q < n - 1, \\ (1 + d_{n-1}^2/d_1^2)^{-1}, & \text{if } q \geq n - 1. \end{cases}$$

REMARK 5.1. As mentioned in [Remark 3.4](#), a small value of d_r could be problematic for estimation. This is reflected in (5.3) where a small d_r would diminish overall shrinkage. However, the probability of such a model would be severely downweighted in the model selection context, and so this diminished shrinkage would be of little consequence.

6. Model selection consistency. In this section we consider the model selection consistency in the case where p is fixed and n approaches infinity. Posterior consistency for model choice means

$$\text{plim}_{n \rightarrow \infty} \Pr(\mathcal{M}_T | y) = 1 \quad \text{when } \mathcal{M}_T \text{ is the true model,}$$

where plim denotes convergence in probability under the true model \mathcal{M}_T , namely, $\mathbf{y} = \alpha_T \mathbf{1}_n + \mathbf{X}_T \boldsymbol{\beta}_T + \boldsymbol{\varepsilon}$, where \mathbf{X}_T is the $n \times q_T$ true design matrix and $\boldsymbol{\beta}_T$ is the true $(q_T \times 1)$ coefficient vector and $\boldsymbol{\varepsilon}_n \sim N_n(\mathbf{0}, \sigma^2 \mathbf{I}_n)$.

Let us show that our general criterion, $\text{BF}_{\gamma:N}(a, \mathbf{v})$ given by (3.2) with bounded v_1 , is model selection consistent. This is clearly equivalent to

$$(6.1) \quad \text{plim}_{n \rightarrow \infty} \frac{\text{BF}_{\gamma:N}(a, \mathbf{v})}{\text{BF}_{T:N}(a, \mathbf{v})} = 0 \quad \forall \mathcal{M}_\gamma \neq \mathcal{M}_T.$$

Recall that we have already assumed that $\mathbf{x}'_i \mathbf{1}_n = 0$ and $\mathbf{x}'_i \mathbf{x}_i / n = 1$ for any $1 \leq i \leq p$. To obtain model selection consistency, we also assume the following:

- (A1) The correlation between x_i and x_j , $\mathbf{x}'_i \mathbf{x}_j / n$, has a limit as $n \rightarrow \infty$.
- (A2) The limit of the correlation matrix of x_1, \dots, x_p , $\lim_{n \rightarrow \infty} \mathbf{X}'_F \mathbf{X}_F / n$, is positive definite.

Assumption (A1) is the standard assumption which also appears in Knight and Fu (2000) and Zou (2006). Assumption (A2) is natural because the columns of \mathbf{X}_F are assumed to be linearly independent.

Our main consistency theorem is as follows. Note that our recommended choice $v_1 = d_1^2 / d_q^2$ is bounded by Lemma B.1 in the Appendix.

THEOREM 6.1. *Under assumptions (A1) and (A2), if v_1 is bounded, then $\text{BF}_{\gamma:N}(a, \mathbf{v})$ is consistent for model selection.*

7. Simulated performance evaluations. In this section we report on a number of simulated performance comparisons between our recommended Bayes factor $g\text{BF}_{\gamma:N}$ and the following selection criteria:

$$\text{ZE} = (1 - R^2)^{-(n-q)/2+3/4} \frac{B(q/2 + 1/4, (n - q)/2 - 3/4)}{B(1/4, (n - q)/2 - 3/4)},$$

$$\text{EB} = \max_g m_\gamma(\mathbf{y} | g, \hat{\sigma}^2),$$

$$\text{AIC} = -2 \times \text{maximum log likelihood} + 2(q + 2),$$

$$\text{AICc} = -2 \times \text{maximum log likelihood} + 2(q + 2) \frac{n}{n - q - 3},$$

$$\text{BIC} = -2 \times \text{maximum log likelihood} + q \log n.$$

Here, ZE is the special case of $\text{BF}_{\gamma:N}$ with $a = -3/4$ and $v_1 = \dots = v_q = 1$ (corresponding to Zellner's g -prior). Note that comparisons of $g\text{BF}$ with ZE should

reveal the effect of our choice of descending ν . EB is the empirical Bayes criterion of George and Foster (2000) in (1.6), also based on the original g -prior, with $\hat{\sigma}^2 = \text{RSS}_\gamma / (n - q_\gamma - 1)$ plugged in. Finally, AICc is the well-known correction of AIC proposed by Hurvich and Tsai (1989).

For these comparisons, we consider data generated by submodels (1.2) of (1.1) with $p = 16$ potential predictors for two different choices of the underlying design matrix \mathbf{X}_F . For the first choice, which we refer to as the correlated case, each row of the 16 predictors are generated as $x_1, \dots, x_{13} \sim N(0, 1)$, and $x_{14}, x_{15}, x_{16} \sim U(-1, 1)$ (the uniform distribution) with the following pairwise correlations:

$$(7.1) \quad \begin{array}{ccccc} \text{cor}=0.9 & & \text{cor}=0.5 & & \text{cor}=0.1 \\ \underbrace{x_1, x_2} & , & \underbrace{x_3, x_4} & , & \underbrace{x_5, x_6} & , & \underbrace{x_7, x_8} & , & \underbrace{x_9, x_{10}} \\ & & \text{cor}=-0.7 & & & & \text{cor}=-0.3 & & \end{array}$$

and independently otherwise. For the second choice, which we refer to as the simple case, each row of the 16 predictors are generated as x_1, \dots, x_{16} i.i.d. $\sim N(0, 1)$.

For our first set of comparisons, we set $n = 30$ (larger than $p = 16$) and considered 4 submodels where the true predictors are:

- $x_1, x_2, x_3, x_4, x_5, x_6, x_7, x_8, x_9, x_{10}, x_{11}, x_{12}, x_{13}, x_{14}, x_{15}, x_{16}$ ($q_T = 16$),
- $x_1, x_2, x_3, x_4, x_5, x_6, x_7, x_8, x_9, x_{10}, x_{11}, x_{14}$ ($q_T = 12$),
- $x_1, x_2, x_5, x_6, x_9, x_{10}, x_{11}, x_{14}$ ($q_T = 8$),
- x_1, x_2, x_5, x_6 ($q_T = 4$)

(where q_T denotes the number of true predictors) and the true model is given by

$$(7.2) \quad Y = 1 + 2 \sum_{i \in \{\text{true}\}} x_i + \{\text{normal error term } N(0, 1)\}.$$

In both cases, after generating pseudo random x_1, \dots, x_{16} , we centered and scaled them as noted in Section 1.

REMARK 7.1. With simulations of performance in Bayesian model selection, the answers primarily depend on the assumed prior. Here we have chosen all the $\beta_i = 2$, an extreme form of the assumption of exchangeability.

Table 1 compares the criteria by how often the true model was selected as best, or in the top 3, among the 2^{16} candidate models across the $N = 500$ replications. We note the following:

- In the correlated cases, EB, ZE and g BF were very similar for $q_T = 4, 8$, but g BF was much better for $q = 12, 16$.
- In the simple cases, g BF, ZE and EB were very similar, suggesting no effect of our extension of Zellner’s g -prior with descending ν .
- In both the correlated and simple cases, AIC and BIC were poor for all cases except $q_T = 16$.

TABLE 1
Rank of the true model

<i>q_T</i> :	16		12		8		4	
	Rank: 1st	1st-3rd	1st	1st-3rd	1st	1st-3rd	1st	1st-3rd
Correlated case								
<i>g</i> BF	0.71	0.91	0.73	0.94	0.69	0.87	0.66	0.86
ZE	0.40	0.70	0.63	0.89	0.68	0.89	0.67	0.87
EB	0.41	0.71	0.63	0.90	0.67	0.88	0.66	0.85
AIC	0.95	0.99	0.23	0.38	0.09	0.17	0.05	0.08
AICc	0.25	0.45	0.67	0.90	0.52	0.75	0.25	0.44
BIC	0.88	0.98	0.41	0.65	0.31	0.43	0.23	0.42
Simple case								
<i>g</i> BF	0.98	0.99	0.83	0.97	0.75	0.93	0.67	0.85
ZE	0.94	0.98	0.87	0.97	0.78	0.95	0.69	0.88
EB	0.95	0.98	0.87	0.98	0.76	0.95	0.65	0.87
AIC	1.00	1.00	0.22	0.37	0.08	0.13	0.05	0.08
AICc	0.82	0.87	0.85	0.97	0.55	0.80	0.24	0.46
BIC	0.99	1.00	0.41	0.65	0.27	0.46	0.22	0.39

- In both the correlated and simple cases, AICc was poor for *q_T* = 16 and 4 but good for *q_T* = 8, 12.

Overall, Table 1 suggests that *g*BF is stable and good for most cases, and that our generalization of Zellner’s *g*-prior is effective in the correlated case.

On data from the same setup with *n* = 30 and *N* = 500, Table 2 compares the models selected by each criterion based on their (in-sample) predictive error

$$\frac{(\hat{y}_* - \alpha_T 1_n - X_T \beta_T)' (\hat{y}_* - \alpha_T 1_n - X_T \beta_T)}{n\sigma^2},$$

where *X_T*, *α_T* and *β_T* are the true *n* × *q_T* design matrix, the true intercept and the true coefficients. The prediction *ŷ_{*}* for each selected model is given by *ȳ*1_{*n*} + *X_{γ*}**β̂_{γ*}*, where *X_{γ*}* is the selected design matrix, *β̂_{γ*}* is the Bayes estimator for *g*BF, ZE and EB, and is the least squares estimator for AIC, BIC and AICc. To aid in gauging these comparisons, we also included the “oracle” prediction error, namely, that based on the least squares estimate under the true model.

The summary statistics reported in Table 2 are the mean predictive error, and the lower quantile (LQ) and upper quantile (UQ) of the predictive errors. In terms of predictive performance, the comparisons are similar to those in Table 1. Overall, we see that *g*BF works well in this setting.

For our final evaluations, we use data again simulated from the simple form (7.2), but now with *x*₁, *x*₂, . . . , *x*₁₂, *x*₁₄, *x*₁₅ as the true predictors (*q_T* = 14) and a small sample size *n* = 12 (smaller than *p* = 16). Since *p* > *q_T* > *n*, the true model

TABLE 2
Prediction error comparisons

	16		12		8		4	
	Mean	(LQ, UQ)						
Correlated case								
Oracle	0.57	(0.43, 0.68)	0.43	(0.31, 0.53)	0.30	(0.20, 0.38)	0.17	(0.09, 0.22)
gBF	0.70	(0.44, 0.78)	0.52	(0.32, 0.61)	0.37	(0.22, 0.47)	0.26	(0.11, 0.35)
ZE	1.02	(0.53, 1.20)	0.59	(0.35, 0.71)	0.41	(0.23, 0.53)	0.27	(0.11, 0.37)
EB	1.00	(0.52, 1.16)	0.58	(0.35, 0.70)	0.41	(0.23, 0.53)	0.27	(0.11, 0.37)
AIC	0.56	(0.42, 0.67)	0.54	(0.40, 0.65)	0.51	(0.37, 0.62)	0.48	(0.33, 0.59)
AICc	1.29	(0.65, 1.65)	0.56	(0.34, 0.68)	0.42	(0.25, 0.52)	0.36	(0.22, 0.47)
BIC	0.58	(0.42, 0.69)	0.53	(0.38, 0.64)	0.46	(0.31, 0.58)	0.39	(0.23, 0.51)
Simple case								
Oracle	0.57	(0.43, 0.68)	0.43	(0.31, 0.53)	0.30	(0.20, 0.38)	0.17	(0.09, 0.22)
gBF	0.57	(0.41, 0.67)	0.45	(0.33, 0.56)	0.35	(0.21, 0.45)	0.25	(0.12, 0.33)
ZE	0.66	(0.42, 0.70)	0.45	(0.32, 0.56)	0.34	(0.21, 0.44)	0.24	(0.12, 0.32)
EB	0.65	(0.42, 0.69)	0.45	(0.32, 0.56)	0.35	(0.21, 0.45)	0.25	(0.12, 0.34)
AIC	0.56	(0.42, 0.67)	0.54	(0.39, 0.65)	0.51	(0.37, 0.63)	0.48	(0.32, 0.60)
AICc	0.98	(0.45, 0.83)	0.46	(0.33, 0.55)	0.39	(0.25, 0.50)	0.35	(0.20, 0.47)
BIC	0.56	(0.42, 0.67)	0.52	(0.37, 0.64)	0.45	(0.30, 0.57)	0.38	(0.21, 0.50)

is not identifiable here. Furthermore, AIC, BIC, AICc, ZE and EB cannot even be computed (because $p > n$) and so we confine our evaluations to gBF.

For this very difficult selection situation, gBF did not rank the complete true model of dimension $q_T = 14$ as best even once across the $N = 500$ iterations. In fact, as shown by the frequency of model sizes selected as best by gBF in Table 3, the top selected model was always of dimension less than $n = 12$, the dimension required for identifiability. However, if one instead considers the overall gBF rankings across all possible models, a different picture emerges.

As can be seen in Table 4, which summarizes the relative rank of the true model (rank/ 2^{16}) over the $N = 500$ iterations (smaller is better), gBF often ranked the true model relatively high. Indeed, the mean relative gBF rank of the true model was 0.035 in the correlated case and 0.039 in the simple structure case. Both of

TABLE 3
Model size frequencies in the many predictors case

	0-6	7	8	9	10	11	12-16
Correlated	0.10	0.11	0.22	0.34	0.16	0.07	0.00
Simple	0.11	0.15	0.21	0.33	0.14	0.06	0.00

TABLE 4
The relative rank of the true model

	Min	LQ	Median	Mean	UQ	Max
Correlated	0.001	0.012	0.023	0.035	0.042	0.518
Simple	0.001	0.013	0.023	0.039	0.043	0.555

these mean ranks were the highest mean ranks achieved by any of the $2^{16} = 65,536$ candidate models! The true model ranks were evidently more stable than the other model ranks which varied more from iteration to iteration. Rather than select a single top ranked model in this context, it would seem to be better to use gBF to restrict interest to a promising subset.

Further, it should be noted that gBF performed best among the larger unidentified models as shown by Table 5, which reports the frequencies with which the true model was ranked highly among the $(16 \times 15)/2 = 120$ candidate models with exactly 14 predictors. To our knowledge, we know of no other analytical selection criterion for choosing between models with $R^2 = 1$, which is the case here.

Finally, we call attention to Table 6 which reports the observed gBF predictor selection frequencies across the top ranked gBF models over the $N = 500$ iterations. These frequencies show that the top gBF models tended to at least be partially correct in the sense that, for the most part, the true individual predictors [designated by (T)] were selected more often than not.

REMARK 7.2. The only variables that were under-selected by gBF in Table 6 were (x_3, x_4) and (x_{14}, x_{15}) in the correlated case. Although x_3 and x_4 are true predictors, their under-selection may be explained by the high negative correlation between them. Interestingly, the under-selection of x_{14} and x_{15} is not explained by correlation (as they are independent in both the correlated and simple cases). Rather, since all predictors have been standardized, it suggests that in this setting, selection of $U(-1, 1)$ predictors may be more difficult than $N(0, 1)$ predictors (they are uniform in the correlated case and normal in the simple case).

TABLE 5
Frequency that the true model was ranked highly among models with 14 predictors

	1st	1st-2nd	1st-3rd
Correlated	0.14	0.22	0.26
Simple	0.13	0.20	0.26

TABLE 6
Predictor frequencies in the many predictors case

	x_1 (T)	x_2 (T)	x_3 (T)	x_4 (T)	x_5 (T)	x_6 (T)
Correlated	0.65	0.63	0.44	0.46	0.62	0.60
Simple	0.54	0.54	0.54	0.54	0.54	0.57
	x_7 (T)	x_8 (T)	x_9 (T)	x_{10} (T)	x_{11} (T)	x_{12} (T)
Correlated	0.56	0.56	0.59	0.58	0.58	0.60
Simple	0.55	0.55	0.54	0.56	0.52	0.50
	x_{13} (F)	x_{14} (T)	x_{15} (T)	x_{16} (F)		
Correlated	0.40	0.43	0.45	0.40		
Simple	0.34	0.55	0.57	0.39		

APPENDIX A: PROOF OF THEOREM 5.1

We proceed by finding a simple closed form for $\hat{\beta}_B$ in (5.1). Making use of the transformation (3.5), and by the calculation in (3.6), $E[\beta_{\#}|y] = E[\beta_{\#}]$ (say, $\mu_{\#}$) and

$$\begin{aligned} \mathbf{W} \frac{E[\sigma^{-2}\beta_{*}|y]}{E[\sigma^{-2}|y]} &= \frac{1}{E[\sigma^{-2}|y]} E \left[\sigma^{-2} \sum_{i=1}^r \frac{\mathbf{u}'_i y}{d_i} \left\{ 1 - \frac{1}{v_i(1+g)} \right\} \mathbf{w}_i | y \right] \\ &= \sum_{i=1}^r \frac{\mathbf{u}'_i y}{d_i} \left\{ 1 - \frac{H(y)}{v_i} \right\} \mathbf{w}_i, \end{aligned}$$

where

$$(A.1) \quad H(y) = \frac{E[\sigma^{-2}(1+g)^{-1}|y]}{E[\sigma^{-2}|y]}.$$

Thus,

$$(A.2) \quad \hat{\beta}_B = \sum_{i=1}^r \frac{\mathbf{u}'_i y}{d_i} \left(1 - \frac{H(y)}{v_i} \right) \mathbf{w}_i + \begin{cases} \mathbf{0}, & \text{if } q \leq n - 1, \\ \mathbf{W}_{\#} \mu_{\#}, & \text{if } q > n - 1. \end{cases}$$

Since β is not identifiable when $q \geq n - 1$, it is not surprising that $\hat{\beta}_B$ is incompletely defined due to the arbitrariness of $\mathbf{W}_{\#} \mu_{\#}$. However, because $\mathbf{XW}_{\#} = \mathbf{0}$, this arbitrariness is not an issue for the estimation of $\mathbf{X}\beta$, for which we obtain

$$(A.3) \quad \mathbf{X}\hat{\beta}_B = \sum_{i=1}^r (\mathbf{u}'_i y) \mathbf{u}_i \left(1 - \frac{H(y)}{v_i} \right).$$

It now only remains to obtain a closed form for $H(\mathbf{y})$. As in (3.4), (3.7) and (3.8) in Section 3,

$$\begin{aligned}
 & \int_{-\infty}^{\infty} \int_{R^q} \int_0^{\infty} \frac{1}{\sigma^2} p(\mathbf{y}|\alpha, \boldsymbol{\beta}, \sigma^2) p(\boldsymbol{\beta}|g, \sigma^2) \frac{1}{\sigma^2} d\alpha d\boldsymbol{\beta} d\sigma^2 \\
 &= \int_0^{\infty} \{\sigma^2\}^{-(n+1)/2} \frac{n^{1/2}}{(2\pi)^{(n-1)/2}} \frac{(1+g)^{-r/2}}{\prod_{i=1}^r v_i^{1/2}} \\
 \text{(A.4)} \quad & \times \exp\left(-\frac{\|\mathbf{v}\|^2\{g(1-R^2)+1-Q^2\}}{2\sigma^2(g+1)}\right) \frac{1}{\sigma^2} d\sigma^2 \\
 &= \frac{2n^{1/2}\Gamma(\{n+1\}/2)}{\pi^{(n-1)/2}} \frac{\|\mathbf{v}\|^{-n-1}}{\prod_{i=1}^r v_i^{1/2}} (1+g)^{-r/2+(n+1)/2} \\
 & \times \{g(1-R^2)+1-Q^2\}^{-(n+1)/2},
 \end{aligned}$$

which differs slightly from (3.8) because of the extra $1/\sigma^2$ term in the first expression. Letting

$$\text{(A.5)} \quad L(\mathbf{y}|g) = (1+g)^{-r/2+(n+1)/2} \{g(1-R^2)+1-Q^2\}^{-(n+1)/2},$$

we have

$$\begin{aligned}
 H(\mathbf{y}) &= \frac{\int_0^{\infty} (1+g)^{-1} L(\mathbf{y}|g) p(g) dg}{\int_0^{\infty} L(\mathbf{y}|g) p(g) dg} \\
 &= \frac{\int_0^{\infty} (1+g)^{-r/2+(n-1)/2} \{g(1-R^2)+1-Q^2\}^{-(n+1)/2} p(g) dg}{\int_0^{\infty} (1+g)^{-r/2+(n+1)/2} \{g(1-R^2)+1-Q^2\}^{-(n+1)/2} p(g) dg}.
 \end{aligned}$$

When $q < n - 1$, under the prior (2.13) used in Section 3, namely,

$$p(g) = \frac{g^b(1+g)^{-a-b-2}}{B(a+1, b+1)} = \frac{g^b(1+g)^{-(n-r-1)/2}}{B(a+1, b+1)},$$

where $b = (n - 5)/2 - r/2 - a$, we have

$$\begin{aligned}
 H(\mathbf{y}) &= \frac{\int_0^{\infty} g^b \{g(1-R^2)+1-Q^2\}^{-(n+1)/2} dg}{\int_0^{\infty} g^b (1+g) \{g(1-R^2)+1-Q^2\}^{-(n+1)/2} dg} \\
 &= \left(1 + \frac{\int_0^{\infty} g^{b+1} \{g(1-R^2)+1-Q^2\}^{-(n+1)/2} dg}{\int_0^{\infty} g^b \{g(1-R^2)+1-Q^2\}^{-(n+1)/2} dg}\right)^{-1} \\
 &= \left(1 + \frac{1-Q^2}{1-R^2} \frac{B(q/2+a+1, b+2)}{B(q/2+a+2, b+1)}\right)^{-1} \\
 &= \left(1 + \frac{1-Q^2}{1-R^2} \frac{(n-q-3)/2-a}{q/2+a+1}\right)^{-1}.
 \end{aligned}$$

On the other hand, when $q \geq n - 1$, it follows that $R^2 = 1, r = n - 1, L(\mathbf{y}|g) = (1 + g)(1 - Q^2)^{-(n+1)/2}$ and, hence,

$$(A.6) \quad H(\mathbf{y}) = \frac{\int_0^\infty p(g) dg}{\int_0^\infty (1 + g)p(g) dg} = \{1 + E[g]\}^{-1}.$$

APPENDIX B: PROOF OF THEOREM 6.1

B.1. Some preliminary lemmas. Under the assumptions (A1) and (A2) in Section 6, we will give the following lemmas (Lemma B.1 on \mathbf{X}_T and \mathbf{X}_γ and Lemmas B.2, B.3 on R_T^2 and R_γ^2) for our main proof. See also Fernández, Ley and Steel (2001) and Liang et al. (2008). Note that (A2) implies that, for any model \mathcal{M}_γ , there exists a positive definite matrix \mathbf{H}_γ such that

$$(B.1) \quad \lim_{n \rightarrow \infty} \frac{1}{n} \mathbf{X}'_\gamma \mathbf{X}_\gamma = \mathbf{H}_\gamma.$$

LEMMA B.1. (1) Let $d_1[\gamma]$ and $d_q[\gamma]$ be the maximum and minimum of singular values of \mathbf{X}_γ . Then $\{d_1[\gamma]\}^2/n$ and $\{d_q[\gamma]\}^2/n$ approach the maximum and minimum eigenvalues of \mathbf{H}_γ , respectively.

(2) The $q_T \times q_T$ limit

$$(B.2) \quad \lim_{n \rightarrow \infty} n^{-1} \mathbf{X}'_T \mathbf{X}_\gamma (\mathbf{X}'_\gamma \mathbf{X}_\gamma)^{-1} \mathbf{X}'_\gamma \mathbf{X}_T = \mathbf{H}(T, \gamma)$$

exists.

(3) When $\gamma \not\supseteq T$, the rank of $\mathbf{H}_T - \mathbf{H}(T, \gamma)$ is given by the number of nonoverlapping predictors and $\beta'_T \mathbf{H}_T \beta_T > \beta'_T \mathbf{H}(T, \gamma) \beta_T$.

(4) $\mathbf{H}_T - \mathbf{H}(T, \gamma) = \mathbf{0}$ for $\gamma \supseteq T$.

LEMMA B.2. Let $\gamma \not\supseteq T$. Then

$$(B.3) \quad \text{plim}_{n \rightarrow \infty} R_\gamma^2 = \frac{\beta'_T \mathbf{H}(\gamma, T) \beta_T}{\sigma^2 + \beta'_T \mathbf{H}_T \beta_T} \left(< \frac{\beta'_T \mathbf{H}_T \beta_T}{\sigma^2 + \beta'_T \mathbf{H}_T \beta_T} \right).$$

PROOF. For the submodel $\mathcal{M}_\gamma, 1 - R_\gamma^2$ is given by

$$\|\mathbf{Q}_\gamma(\mathbf{y} - \bar{y}\mathbf{1}_n)\|^2 / \|\mathbf{y} - \bar{y}\mathbf{1}_n\|^2$$

with $\mathbf{Q}_\gamma = \mathbf{I} - \mathbf{X}_\gamma (\mathbf{X}'_\gamma \mathbf{X}_\gamma)^{-1} \mathbf{X}'_\gamma$. The numerator and denominator are rewritten as

$$(B.4) \quad \begin{aligned} \|\mathbf{Q}_\gamma(\mathbf{y} - \bar{y}\mathbf{1}_n)\|^2 &= \|\mathbf{Q}_\gamma \mathbf{X}_T \beta_T + \mathbf{Q}_\gamma \check{\boldsymbol{\varepsilon}}\|^2 \\ &= \beta'_T \mathbf{X}'_T \mathbf{Q}_\gamma \mathbf{X}_T \beta_T + 2\beta'_T \mathbf{X}'_T \mathbf{Q}_\gamma \boldsymbol{\varepsilon} + \check{\boldsymbol{\varepsilon}}' \mathbf{Q}_\gamma \check{\boldsymbol{\varepsilon}}, \end{aligned}$$

where $\check{\boldsymbol{\varepsilon}} = \boldsymbol{\varepsilon} - \bar{\boldsymbol{\varepsilon}}\mathbf{1}_n$ and, similarly,

$$\|\mathbf{y} - \bar{y}\mathbf{1}_n\|^2 = \beta'_T \mathbf{X}'_T \mathbf{X}_T \beta_T + 2\beta'_T \mathbf{X}'_T \boldsymbol{\varepsilon} + \|\boldsymbol{\varepsilon}\|^2.$$

Hence, $1 - R_\gamma^2$ can be rewritten as

$$(B.5) \quad \frac{\beta'_T \{X'_T Q_\gamma X_T / n\} \beta_T + 2\beta'_T \{X'_T Q_\gamma \epsilon / n\} + \|Q_\gamma \check{\epsilon}\|^2 / n}{\beta'_T \{X'_T X_T / n\} \beta_T + 2\beta'_T \{X'_T \epsilon / n\} + \|\check{\epsilon}\|^2 / n}.$$

In (B.5), $\beta'_T X'_T \epsilon / n$ approaches 0 in probability because $E[\epsilon] = \mathbf{0}$, $\text{var}[\epsilon] = \sigma^2 \mathbf{I}_n$, $E[X'_T \epsilon / n] = \mathbf{0}$ and

$$(B.6) \quad \text{var}(X'_T \epsilon / n) = n^{-1} \sigma^2 \{X'_T X_T / n\} \rightarrow \mathbf{0}.$$

Similarly $\beta'_T \{X'_T Q_\gamma \epsilon / n\} \rightarrow 0$ in probability. Further, both $\|\check{\epsilon}\|^2 / n$ and $\|Q_\gamma \check{\epsilon}\|^2 / n$ for any γ converge to σ^2 in probability.

Therefore, by parts (2) and (3) of Lemma B.1, R_γ^2 for $\gamma \not\supseteq T$ approaches

$$\frac{\beta'_T \mathbf{H}(\gamma, T) \beta_T}{\sigma^2 + \beta'_T \mathbf{H}_T \beta_T} \left(< \frac{\beta'_T \mathbf{H}_T \beta_T}{\sigma^2 + \beta'_T \mathbf{H}_T \beta_T} \right)$$

in probability. \square

LEMMA B.3. *Let $\gamma \supseteq T$. Then:*

(1) $R_\gamma^2 \geq R_T^2$ for any n and

$$(B.7) \quad \text{plim}_{n \rightarrow \infty} R_T^2 = \text{plim}_{n \rightarrow \infty} R_\gamma^2 = \frac{\beta'_T \mathbf{H}_T \beta_T}{\sigma^2 + \beta'_T \mathbf{H}_T \beta_T}.$$

(2) $\{(1 - R_T^2)/(1 - R_\gamma^2)\}^n$ is bounded from above in probability.

PROOF. (1) When $\gamma \supseteq T$, $Q_\gamma X_T = \mathbf{0}$. Hence, as in (B.5), we have

$$(B.8) \quad 1 - R_\gamma^2 = \frac{\|Q_\gamma \check{\epsilon}\|^2 / n}{\beta'_T \{X'_T X_T / n\} \beta_T + 2\beta'_T \{X'_T \epsilon / n\} + \|\check{\epsilon}\|^2 / n},$$

$$1 - R_T^2 = \frac{\|Q_T \check{\epsilon}\|^2 / n}{\beta'_T \{X'_T X_T / n\} \beta_T + 2\beta'_T \{X'_T \epsilon / n\} + \|\check{\epsilon}\|^2 / n}.$$

Since $\|Q_T \check{\epsilon}\|^2 / n > \|Q_\gamma \check{\epsilon}\|^2 / n$ for any n and both approach σ^2 in probability, part (1) follows.

(2) By (B.8), $(1 - R_T^2)/(1 - R_\gamma^2)$ is given by $\|Q_T \check{\epsilon}\|^2 / \|Q_\gamma \check{\epsilon}\|^2$. Further, we have

$$1 \leq \frac{1 - R_T^2}{1 - R_\gamma^2} = \frac{\|Q_T \check{\epsilon}\|^2}{\|Q_\gamma \check{\epsilon}\|^2} \leq \frac{\|\check{\epsilon}\|^2}{\|Q_\gamma \check{\epsilon}\|^2} = \frac{1}{W_\gamma},$$

where $W_\gamma \sim (1 + \chi_{q_\gamma}^2 / \chi_{n-q_\gamma-1}^2)^{-1}$, for independent $\chi_{n-q_\gamma-1}^2$ and $\chi_{q_\gamma}^2$. Hence,

$$\begin{aligned} \{1 + \chi_{q_\gamma}^2 / \chi_{n-q_\gamma-1}^2\}^{-n} &= \{1 + \{n / \chi_{n-q_\gamma-1}^2\} \{\chi_{q_\gamma}^2 / n\}\}^{-n} \\ &\sim \exp(-\chi_{q_\gamma}^2) \quad \text{as } n \rightarrow \infty \end{aligned}$$

since $\chi_{n-q_\gamma-1}^2/n \rightarrow 1$ in probability. Therefore, W_γ^{-n} is bounded in probability from above and part (2) follows. \square

B.2. The proof of Theorem 6.1. Note that

$$v_1^{-1} \leq 1 - Q_\gamma^2 \leq 1$$

by (3.3),

$$v_1^{-q/2} \leq \prod_{i=1}^q v_i^{-1/2} \leq 1,$$

because the v_i 's are descending,

$$\frac{B(q/2 + a + 1, (n - q - 3)/2 - a)}{B(a + 1, (n - q - 3)/2 - a)} = \frac{\Gamma(q/2 + a + 1) \Gamma(\{n - q - 1\}/2)}{\Gamma(a + 1) \Gamma(\{n - 1\}/2)}$$

and

$$\lim_{n \rightarrow \infty} (n/2)^{q/2} \frac{\Gamma(\{n - q - 1\}/2)}{\Gamma(\{n - 1\}/2)} = 1$$

by Stirling's formula. Then, by (3.2), there exist $c_1(\gamma) < c_2(\gamma)$ (which do not depend on n) such that

$$c_1(\gamma) < \{n^{q_\gamma} (1 - R_\gamma^2)^n\}^{1/2} \frac{\text{BF}_{\gamma:N}(a, v)}{(1 - R_\gamma^2)^{(q_\gamma+3)/2+a}} < c_2(\gamma)$$

for sufficiently large n . By Lemmas B.2 and B.3, R_γ^2 goes to some constant in probability. Hence, to show consistency, it suffices to show that

$$(B.9) \quad \text{plim}_{n \rightarrow \infty} n^{q_T - q_\gamma} \left(\frac{1 - R_T^2}{1 - R_\gamma^2} \right)^n = 0.$$

Consider the following two situations:

(1) $\gamma \not\geq T$: by Lemmas B.2 and B.3, $(1 - R_T^2)/(1 - R_\gamma^2)$ is strictly less than 1 in probability. Hence, $\{(1 - R_T^2)/(1 - R_\gamma^2)\}^n$ converges to zero in probability exponentially fast with respect to n . Therefore, no matter what value $q_T - q_\gamma$ takes, (B.9) is satisfied.

(2) $\gamma \geq T$: by Lemma B.3, $\{(1 - R_T^2)/(1 - R_\gamma^2)\}^n$ is bounded in probability. Since $q_\gamma > q_T$, (B.9) is satisfied.

Acknowledgments. We are very grateful to a referee for wonderful insights which substantially helped us to strengthen this paper.

REFERENCES

- AKAIKE, H. (1974). A new look at the statistical model identification. *IEEE Trans. Automat. Control* **AC-19** 716–723. System identification and time-series analysis. [MR0423716](#)
- BERGER, J. O., PERICCHI, L. R. and VARSHAVSKY, J. A. (1998). Bayes factors and marginal distributions in invariant situations. *Sankhyā Ser. A* **60** 307–321. [MR1718789](#)
- CASELLA, G. (1980). Minimax ridge regression estimation. *Ann. Statist.* **8** 1036–1056. [MR0585702](#)
- CASELLA, G. (1985). Condition numbers and minimax ridge regression estimators. *J. Amer. Statist. Assoc.* **80** 753–758. [MR0803264](#)
- CUI, W. and GEORGE, E. I. (2008). Empirical Bayes vs. fully Bayes variable selection. *J. Statist. Plann. Inference* **138** 888–900. [MR2416869](#)
- FERNÁNDEZ, C., LEY, E. and STEEL, M. F. J. (2001). Benchmark priors for Bayesian model averaging. *J. Econometrics* **100** 381–427. [MR1820410](#)
- FOSTER, D. P. and GEORGE, E. I. (1994). The risk inflation criterion for multiple regression. *Ann. Statist.* **22** 1947–1975. [MR1329177](#)
- GELUK, J. L. and DE HAAN, L. (1987). *Regular Variation, Extensions and Tauberian Theorems. CWI Tract* **40**. Math. Centrum, Centrum Wisk. Inform., Amsterdam. [MR0906871](#)
- GEORGE, E. I. and FOSTER, D. P. (2000). Calibration and empirical Bayes variable selection. *Biometrika* **87** 731–747. [MR1813972](#)
- HURVICH, C. M. and TSAI, C.-L. (1989). Regression and time series model selection in small samples. *Biometrika* **76** 297–307. [MR1016020](#)
- KNIGHT, K. and FU, W. (2000). Asymptotics for lasso-type estimators. *Ann. Statist.* **28** 1356–1378. [MR1805787](#)
- LIANG, F., PAULO, R., MOLINA, G., CLYDE, M. A. and BERGER, J. O. (2008). Mixtures of g priors for Bayesian variable selection. *J. Amer. Statist. Assoc.* **103** 410–423. [MR2420243](#)
- MARUYAMA, Y. and STRAWDERMAN, W. E. (2005). A new class of generalized Bayes minimax ridge regression estimators. *Ann. Statist.* **33** 1753–1770. [MR2166561](#)
- SCHWARZ, G. (1978). Estimating the dimension of a model. *Ann. Statist.* **6** 461–464. [MR0468014](#)
- STRAWDERMAN, W. E. (1971). Proper Bayes minimax estimators of the multivariate normal mean. *Ann. Math. Statist.* **42** 385–388. [MR0397939](#)
- ZELLNER, A. (1986). On assessing prior distributions and Bayesian regression analysis with g -prior distributions. In *Bayesian Inference and Decision Techniques. Stud. Bayesian Econometrics Statist.* **6** 233–243. North-Holland, Amsterdam. [MR0881437](#)
- ZELLNER, A. and SIOW, A. (1980). Posterior odds ratios for selected regression hypotheses. In *Bayesian Statistics: Proceedings of the First International Meeting Held in Valencia (Spain)* (J. M. Bernardo, M. H. DeGroot, D. V. Lindley and A. F. M. Smith, eds.) 585–603. Univ. Valencia, Valencia.
- ZOU, H. (2006). The adaptive lasso and its oracle properties. *J. Amer. Statist. Assoc.* **101** 1418–1429. [MR2279469](#)

CENTER FOR SPATIAL INFORMATION SCIENCE
UNIVERSITY OF TOKYO
5-1-5 KASHIWANOHA, KASHIWA-SHI
CHIBA, 277-8568
JAPAN
E-MAIL: maruyama@csis.u-tokyo.ac.jp

DEPARTMENT OF STATISTICS
UNIVERSITY OF PENNSYLVANIA
400 JON M. HUNTSMAN HALL
3730 WALNUT STREET
PHILADELPHIA, PENNSYLVANIA 19104-6302
USA
E-MAIL: edgeorge@wharton.upenn.edu