

Bayesian Model Averaging for Linear Regression Models*

Adrian E. Raftery David Madigan
University of Washington University of Washington

Jennifer A. Hoeting
Colorado State University

April 20, 1998

Journal of the American Statistical Association (1997) 92,
179-191

Abstract

We consider the problem of accounting for model uncertainty in linear regression models. Conditioning on a single selected model ignores model uncertainty, and thus leads to the underestimation of uncertainty when making inferences about quantities of interest. A Bayesian solution to this problem involves averaging over all possible models (i.e., combinations of predictors) when making inferences about quantities of

*Adrian E. Raftery is Professor of Statistics and Sociology, David Madigan is Assistant Professor of Statistics, both at the Department of Statistics, University of Washington, Box 354322, Seattle, WA 98195-4322. Jennifer Hoeting is Assistant Professor of Statistics at the Department of Statistics, Colorado State University, Fort Collins, CO 80523. The research of Raftery and Hoeting was partially supported by ONR Contract N-00014-91-J-1074. Madigan's research was partially supported by NSF grant no. DMS 92111627. The authors are grateful to Danika Lew for research assistance and the Editor, the Associate Editor, two anonymous referees and David Draper for very helpful comments that greatly improved the article.

interest. This approach is often not practical. In this paper we offer two alternative approaches. First we describe an ad hoc procedure called “Occam’s Window” which indicates a small set of models over which a model average can be computed. Second, we describe a Markov chain Monte Carlo approach which directly approximates the exact solution. In the presence of model uncertainty, both these model averaging procedures provide better predictive performance than any single model which might reasonably have been selected.

In the extreme case where there are many candidate predictors but no relationship between any of them and the response, standard variable selection procedures often choose some subset of variables that yields a high R^2 and a highly significant overall F value. In this situation, Occam’s Window usually indicates the null model as the only one to be considered, or else a small number of models including the null model, thus largely resolving the problem of selecting significant models when there is no signal in the data.

Software to implement our methods is available from StatLib.

Key Words: Bayes factor; Markov chain Monte Carlo model composition; Model uncertainty; Occam’s Window; Posterior model probability.

1 Introduction

The selection of subsets of predictor variables is a basic part of building a linear regression model. The objective of variable selection is typically stated as follows: given a dependent variable Y and a set of a candidate predictors X_1, X_2, \dots, X_k , find the “best” model of the form

$$Y = \beta_0 + \sum_{j=1}^p \beta_{i_j} X_{i_j} + \epsilon,$$

where $X_{i_1}, X_{i_2}, \dots, X_{i_p}$ is a subset of X_1, X_2, \dots, X_k . Here “best” may have any of several meanings, e.g., the model providing the most accurate predictions for new cases exchangeable with those used to fit the model.

A typical approach to data analysis is to carry out a model selection exercise leading to a single “best” model and then to make inference as if the selected model were the true model. However, this ignores a major component of uncertainty, namely uncertainty about the model itself (Leamer 1978; Hodges 1987; Raftery 1988, 1996; Moulton 1991; Draper 1995). As a consequence, uncertainty about quantities of interest can be underestimated. For striking examples of this see Miller (1984), Regal and Hook (1991), Madigan and York (1995), Raftery (1996), and Kass and Raftery (1995), and Draper (1995). A complete Bayesian solution to this problem involves averaging over *all* possible combinations of predictors when making inferences about quantities of interest. Indeed, this approach provides optimal predictive ability (Madigan and Raftery 1994). In many applications however, this averaging will not be a practical proposition and here we present two alternative approaches.

First we extend the Bayesian graphical model selection algorithm of Madigan and Raftery (1994) to linear regression models. We refer to this algorithm as “Occam’s Window.” This approach involves averaging over a reduced set of models. Second, we directly approximate the complete solution by applying the Markov chain Monte Carlo model composition (MC³) approach of Madigan and York (1995) to linear regression models. In this approach the posterior distribution of a quantity of interest is approximated by a Markov chain Monte Carlo method which generates a process that moves through model space. We show in an example that both of these model averaging approaches provide better predictive performance than any single model which might reasonably have been selected.

Freedman (1983) pointed out that when there are many predictors and there is no relationship between the predictors and the response, variable selection techniques can lead to a model with a high R^2 and a highly significant overall F value. By contrast, when a data set is generated with no relationship between the predictors and the response, Occam’s Window typically indicates the null model as the “best” model or as one of a small set of “best” models, thus largely resolving the problem of selecting a significant model for a null relationship.

The background literature for our approach includes several areas of research, namely the selection of subsets of predictor variables in linear regression models (Hocking 1976, Draper and Smith 1981, Shibata 1981, Linhart and Zucchini 1986, Miller 1990, Breiman 1992, Breiman and Spector 1992, Breiman 1995), Bayesian approaches to the selection of subsets of predictor variables in linear regression models (Mitchell and Beauchamp 1988, Schwarz 1978, George and McCulloch 1993, Laud and Ibrahim 1995), and model uncertainty (Leamer 1978, Freedman *et al.* 1986, Stewart and Davis 1986, Stewart 1987, Madigan and Raftery 1994).

In the next section we outline the philosophy underlying our approach. In Section 3 we describe how we selected prior distributions, and we outline the two model averaging approaches in Section 4. In Section 5 we provide an example and describe our assessment of predictive performance. In Section 6 we compare the performance of Occam’s Window to that of standard variable selection methods when there is no relationship between the predictors and the response. In Section 7 we discuss related work and suggest future directions.

2 Accounting for Model Uncertainty using BMA

As described above, basing inferences on a single “best” model as if the single selected model were true ignores model uncertainty which can result in underestimation of uncertainty about quantities of interest. There is a standard Bayesian solution to this problem, proposed by Leamer (1978). If $\mathcal{M} = \{M_1, \dots, M_K\}$ denotes the set of all models being considered and if Δ is the quantity of interest such as a future observation or the utility of a course of action, then the posterior distribution of Δ given the data D is

$$\text{pr}(\Delta | D) = \sum_{k=1}^K \text{pr}(\Delta | M_k, D) \text{pr}(M_k | D). \quad (1)$$

This is an average of the posterior distributions under each model weighted by the corresponding posterior model probabilities. We call this Bayesian Model Averaging (BMA). In equation (1), the posterior probability of model M_k is given by

$$\text{pr}(M_k | D) = \frac{\text{pr}(D | M_k) \text{pr}(M_k)}{\sum_{l=1}^K \text{pr}(D | M_l) \text{pr}(M_l)}, \quad (2)$$

where

$$\text{pr}(D | M_k) = \int \text{pr}(D | \theta_k, M_k) \text{pr}(\theta_k | M_k) d\theta_k \quad (3)$$

is the marginal likelihood of model M_k , θ_k is the vector of parameters of model M_k , $\text{pr}(\theta_k | M_k)$ is the prior density of θ_k under model M_k , $\text{pr}(D | \theta_k, M_k)$ is the likelihood, and $\text{pr}(M_k)$ is the prior probability that M_k is the true model. All probabilities are implicitly conditional on \mathcal{M} , the set of all models being considered. In this paper, we consider \mathcal{M} to be equal to the set of all possible combinations of predictors.

Averaging over *all* the models in this fashion provides better predictive ability, as measured by a logarithmic scoring rule, than using any single model M_j :

$$-E \left[\log \left\{ \sum_{k=1}^K \text{pr}(\Delta | M_k, D) \text{pr}(M_k | D) \right\} \right] \leq -E [\log \{ \text{pr}(\Delta | M_j, D) \}] \quad (j = 1, \dots, K),$$

where Δ is the observable to be predicted and the expectation is with respect to $\sum_{k=1}^K \text{pr}(\Delta | M_k, D)\text{pr}(M_k | D)$. This follows from the non-negativity of the Kullback-Leibler information divergence.

Implementation of Bayesian model averaging is difficult for two reasons. First, the integrals in (3) can be hard to compute. Second, the number of terms in (1) can be enormous. In this paper, we present solutions to both of these problems.

3 Bayesian Framework

3.1 Modelling Framework

Each model we consider is of the form

$$Y = \beta_0 + \sum_{j=1}^p \beta_j X_j + \epsilon = X\beta + \epsilon, \quad (4)$$

where the observed data on p predictors are contained in the $n \times (p + 1)$ matrix X . The observed data on the dependent variable are contained in the n -vector Y . We assign to ϵ a normal distribution with mean 0 and variance σ^2 and assume that the ϵ 's in distinct cases are independent. We consider the $(p + 1)$ individual parameter vectors β and σ^2 to be unknown.

Where possible, informative prior distributions for β and σ^2 should be elicited and incorporated into the analysis—see Kadane *et al.* (1980) and Garthwaite and Dickey (1992). In the absence of expert opinion we seek to choose prior distributions which reflect uncertainty about the parameters and also embody reasonable *a priori* constraints. We use prior distributions that are proper but reasonably flat over the range of parameter values that could plausibly arise. These represent the common situation where there is some prior information, but rather little of it, and put us in the “stable estimation” case where results are relatively insensitive to changes in

the prior distribution (Edwards, Lindman, and Savage 1963). We use the standard normal-gamma conjugate class of priors,

$$\beta \sim N(\mu, \sigma^2 V),$$

$$\frac{\nu\lambda}{\sigma^2} \sim \chi_\nu^2.$$

Here ν , λ , the $(p+1) \times (p+1)$ matrix V , and the $(p+1)$ -vector μ are hyperparameters to be chosen.

The marginal likelihood for Y under a model M_i based on the proper priors described above is given by

$$p(Y|\mu_i, V_i, X_i, M_i) = \frac{\pi^{\frac{n}{2}}, \left(\frac{\nu+n}{2}\right)(\nu\lambda)^{\frac{\nu}{2}}}{\pi^{\frac{n}{2}}, \left(\frac{\nu}{2}\right)|I + X_i V_i X_i^t|^{\frac{1}{2}}} \cdot \left[\lambda\nu + (Y - X_i\mu_i)^t(I + X_i V_i X_i^t)^{-1}(Y - X_i\mu_i)\right]^{-\frac{(\nu+n)}{2}}, \quad (5)$$

where X_i is the design matrix and V_i is the covariance matrix for β corresponding to model M_i (Raiffa and Schlaifer 1961). The Bayes factor for M_0 versus M_1 , the ratio of equation (5) for $i = 0$ and $i = 1$, is then given by

$$B_{01} = \left(\frac{|I + X_1 V_1 X_1^t|}{|I + X_0 V_0 X_0^t|}\right)^{\frac{1}{2}} \left[\frac{\lambda\nu + (Y - X_0\mu_0)^t(I + X_0 V_0 X_0^t)^{-1}(Y - X_0\mu_0)}{\lambda\nu + (Y - X_1\mu_1)^t(I + X_1 V_1 X_1^t)^{-1}(Y - X_1\mu_1)}\right]^{-\frac{(\nu+n)}{2}}. \quad (6)$$

3.2 Selection of Prior Distributions

The Bayesian framework described above gives the user of the BMA approach the flexibility to modify the prior set-up as desired. In this section we describe the prior distribution set-up we adopt in our examples below.

reasonable desiderata and attempt to satisfy them. In what follows we assume that all the variables have been standardized to have zero mean and sample variance one. We would like:

1. The prior density $\text{pr}(\beta_1, \dots, \beta_p)$ to be reasonably flat over the unit hypercube $[-1, 1]^p$.
2. $\text{pr}(\sigma^2)$ to be reasonably flat over $(a, 1)$ for some small a .
3. $\text{Pr}(\sigma^2 \leq 1)$ to be large.

The order of importance of these desiderata is roughly the order in which they are listed. More formally, we maximize $\text{Pr}(\sigma^2 \leq 1)$ subject to:

a. $\frac{\text{pr}(\beta_1=0, \dots, \beta_p=0)}{\text{pr}(\beta_1=1, \dots, \beta_p=1)} \leq K_1$

Following Jeffreys (1961) we choose $K_1 = \sqrt{10}$.

b. $\frac{\max_{a < \sigma^2 < 1} \text{pr}(\sigma^2)}{\text{pr}(\sigma^2=a)} \leq K_2$

c. $\frac{\max_{a < \sigma^2 < 1} \text{pr}(\sigma^2)}{\text{pr}(\sigma^2=1)} \leq K_2$

Since desideratum 2 is less important than desideratum 1, we have chosen $K_2 = 10$.

For $a = 0.05$ this yields $\nu = 2.58$, $\lambda = 0.28$, and $\phi = 2.85$. For this set of hyperparameters $\text{Pr}(\sigma^2 \leq 1) = 0.81$. These settings of the hyperparameters were used in the examples below.

To compare our prior for β_i , $i = 1, \dots, p$, for a non-categorical predictor with the actual distribution of coefficients from real data, 13 data sets from several regression textbooks were collected (Appendix A). A histogram of the 100 coefficients from the

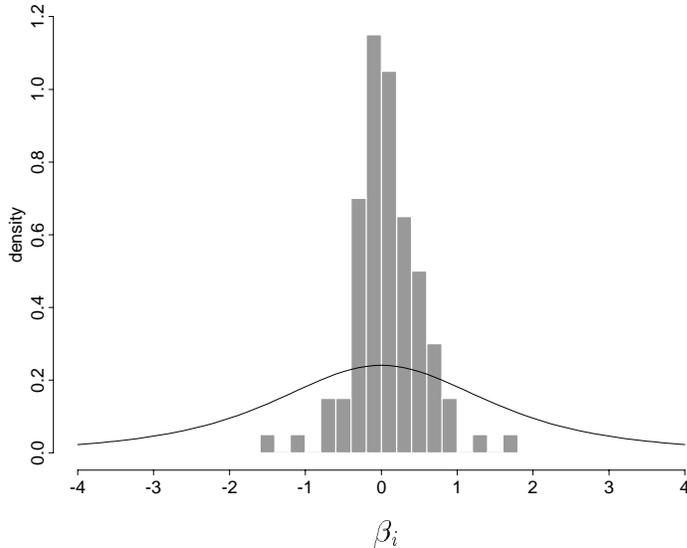


Figure 1: Histogram of 100 coefficients from standardized data, from 13 textbook data sets. The solid line is the prior density for β_i , $i = 1, \dots, p$.

standardized data plotted with the prior distribution resulting from the hyperparameters we use in this paper is shown in Figure 1. As desired, the prior density is relatively flat over the range of observed values.

4 Two approaches to Bayesian Model Averaging

4.1 Occam's Window

Our first method for accounting for model uncertainty starting from equation (1) involves applying the Occam's Window algorithm of Madigan and Raftery (1994) to linear regression models. Two basic principles underly this ad hoc approach.

First, if a model predicts the data far less well than the model which provides the best predictions, then it has effectively been discredited and should no longer be considered. Thus models not belonging to:

$$\mathcal{A}' = \left\{ M_k : \frac{\max_l \{\text{pr}(M_l | D)\}}{\text{pr}(M_k | D)} \leq C \right\}, \quad (7)$$

should be excluded from equation (1) where C is chosen by the data analyst and $\max_l \{\text{pr}(M_l | D)\}$ denotes the model with the highest posterior model probability.

In the examples below we use $C = 20$. The number of models in Occam’s Window increases as the value of C decreases.

Second, appealing to Occam’s razor, we exclude models which receive less support from the data than any of their simpler submodels. More formally we also exclude from (1) models belonging to:

$$\mathcal{B} = \left\{ M_k : \exists M_l \in \mathcal{M}, M_l \subset M_k, \frac{\text{pr}(M_l | D)}{\text{pr}(M_k | D)} > 1 \right\}. \quad (8)$$

Equation (1) is then replaced by

$$\text{pr}(\Delta | D) = \frac{\sum_{M_k \in \mathcal{A}} \text{pr}(\Delta | M_k, D) \text{pr}(D | M_k) \text{pr}(M_k)}{\sum_{M_k \in \mathcal{A}} \text{pr}(D | M_k) \text{pr}(M_k)}, \quad (9)$$

where

$$\mathcal{A} = \mathcal{M} \setminus \mathcal{B}. \quad (10)$$

This greatly reduces the number of models in the sum in equation (1) and now all that is required is a search strategy to identify the models in \mathcal{A} . Two further principles underly the search strategy. The first principle — “Occam’s Window” — concerns the interpretation of the ratio of posterior model probabilities $\text{pr}(M_1 | D)/\text{pr}(M_0 | D)$. Here M_0 is a model with one less predictor than M_1 . The essential idea is shown in Figure 2. If there is evidence for M_0 then M_1 is rejected, but to reject M_0 we require strong evidence *for* the larger model, M_1 . If the evidence is inconclusive (falling in Occam’s Window) neither model is rejected. The second principle is that if M_0 is rejected, then so are all of the models nested within it.

These principles fully define the strategy. Typically, in our experience, the number of terms in (1) is reduced to fewer than 25, and often to as few as one or two. Madigan and Raftery (1994) provide a detailed description of the algorithm and show how averaging over the selected models provides better predictive performance than basing inference on a single model in each of the examples they consider.

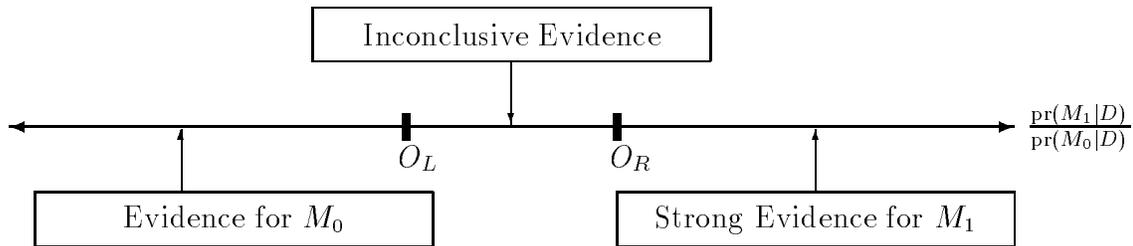


Figure 2: Occam’s Window: Interpreting the posterior odds for nested models.

4.2 Markov Chain Monte Carlo Model Composition

Our second approach is to approximate (1) using a Markov chain Monte Carlo (MCMC) approach (see, for example, Smith and Roberts 1993). For our application, we adopt the Markov chain Monte Carlo model composition (MC³) methodology of Madigan and York (1995) which generates a stochastic process which moves through model space. We can construct a Markov chain $\{M(t), t = 1, 2, \dots\}$ with state space \mathcal{M} and equilibrium distribution $\text{pr}(M_i | D)$. If we simulate this Markov chain for $t = 1, \dots, N$, then under certain regularity conditions, for any function $g(M_i)$ defined on \mathcal{M} , the average

$$\hat{G} = \frac{1}{N} \sum_{t=1}^N g(M(t)) \quad (11)$$

converges almost surely to $E(g(M))$ as $N \rightarrow \infty$ (Smith and Roberts 1993). To compute (1) in this fashion set $g(M) = \text{pr}(\Delta | M, D)$.

To construct the Markov chain we define a neighborhood $\text{nbd}(M)$ for each $M \in \mathcal{M}$ which consists of the model M itself and the set of models with either one variable more or one variable fewer than M . Define a transition matrix q by setting $q(M \rightarrow M') = 0$ for all $M' \notin \text{nbd}(M)$ and $q(M \rightarrow M')$ constant for all $M' \in \text{nbd}(M)$. If the chain is currently in state M , we proceed by drawing M' from $q(M \rightarrow M')$. It is

then accepted with probability

$$\min \left\{ 1, \frac{\text{pr}(M' | D)}{\text{pr}(M | D)} \right\}.$$

Otherwise the state stays in state M . Madigan and York (1995) described MC³ for discrete graphical models.

Software for implementing the MC³ algorithm is described in the Appendix.

5 Model uncertainty and prediction

5.1 Example: Crime and Punishment

5.1.1 Crime and Punishment: Overview

Up to the 1960s, criminal behavior was traditionally viewed as deviant and linked to the offender’s presumed exceptional psychological, social or family circumstances (Taft and England 1964). Becker (1968) and Stigler (1970) argued, on the contrary, that the decision to engage in criminal activity is a rational choice determined by its costs and benefits relative to other (legitimate) opportunities.

In an influential article, Ehrlich (1973) developed this argument theoretically, specified it mathematically, and tested it empirically using aggregate data from 47 U.S. states in 1960. Errors in Ehrlich’s empirical analysis were corrected by Vandaele (1978) who gave the corrected data, which we use here; see also Cox and Snell (1982)¹.

Ehrlich’s theory goes as follows. The costs of crime are related to the probability of imprisonment and the average time served in prison, which in turn are influenced by police expenditures, which may themselves have an independent deterrent effect.

The benefits of crime are related to both the aggregate wealth and income inequality

¹Ehrlich’s study has been much criticized (e.g. Brier and Fienberg 1980) and we use it here for purely illustrative purposes. For economy of expression, we use causal language and speak of “effects”, even though the validity of this language for these data is dubious. Since people, not states, commit crimes, these data may reflect aggregation bias.

in the surrounding community. The expected net payoff from alternative legitimate activities is related to educational level and the availability of employment, the latter being measured by the unemployment and labor force participation rates. The payoff from legitimate activities was expected to be lower (in 1960) for nonwhites and for young males than for others, so that states with high proportions of these were expected also to have higher crime rates. Vandaele (1978) also included an indicator variable for southern states, the sex ratio, and the state population as control variables, but the theoretical rationale for inclusion of these predictors is unclear.

We thus have 15 candidate predictors of crime rate (Table 4), and so potentially $2^{15} = 32,768$ different models. As in the original analyses, all data were transformed logarithmically. Standard diagnostic checking (e.g. Draper and Smith 1981) did not reveal any gross violations of the assumptions underlying normal linear regression.

Ehrlich's analysis concentrated on the relationship between crime rate and predictors 14 and 15 (probability of imprisonment and average time served in state prisons). In his original analysis, Ehrlich (1973) focused on two regression models, consisting of the predictors (9, 12, 13, 14, 15) and (1, 6, 9, 10, 12, 13, 14, 15), respectively, which were chosen in advance based on theoretical grounds.

To compare Ehrlich's results with models that might be selected using standard techniques, we chose three popular variable selection techniques, Efroymson's stepwise method (Miller 1990), minimum Mallows's C_p , and maximum adjusted R^2 (Weisberg 1985). Efroymson's stepwise method is like forward selection except that when a new variable is added to the subset, partial correlations are considered to see if any of the variables currently in the subset should be dropped. Similar hybrid methods are found in most standard statistical computer packages. Problems with stepwise regression, Mallows's C_p , and adjusted R^2 are well known (see, for example, Weisberg

Table 1: Models selected for crime data. For the stepwise procedure, $F=3.84$ was used for the F-to-enter and F-to-delete value. This corresponds approximately to the 5% level.

#	Method	Variables	R^2 (%)	# vars.	$\hat{\beta}_{14}$	$\hat{\beta}_{15}$	P_{15}
1	Full model	All	87	15	-.30	-.27	.133
2	Stepwise regression	1 3 4 9 11 13 14	83	7	-.19	—	—
3	Mallows' C_p	1 3 4 9 11 12 13 14 15	85	9	-.30	-.30	.050
4	Adjusted R^2	1 3 4 7 8 9 11 12 13 14 15	86	11	-.30	-.25	.129
5	Ehrlich model 1	9 12 13 14 15	66	5	-.45	-.55	.009
6	Ehrlich model 2	1 6 9 10 12 13 14 15	70	8	-.43	-.53	.011

Note: P_{15} is the p -value from a two-sided t -test for testing $\beta_{15} = 0$.

1985).

Table 1 displays the results from the full model with all 15 predictors, three models selected using standard variable selection techniques, and the two models chosen by Ehrlich on theoretical grounds. The three models chosen using variable selection techniques (models 2, 3, 4) share many of the same variables and have high values of R^2 . Ehrlich's theoretically chosen models fit the data less well. There are striking differences, indeed conflicts between the results from the different models. Even the models chosen using statistical techniques lead to conflicting conclusions about the main questions of interest, in spite of the models' superficial similarity.

Consider first the predictor for probability of imprisonment, X_{14} . This predictor is significant in all six models, so interest focuses on estimating the size of its effect. To aid interpretation, recall that all variables have been transformed logarithmically, so that, when all other predictors are held fixed, $\beta_{14} = -.30$ means roughly that a 10% increase in the probability of imprisonment produces a 3% reduction in the crime rate.

The estimates of β_{14} fluctuate wildly between models. The stepwise regression model gives an estimate that is about one-third lower in absolute value than the full model, enough to be of policy importance; this difference is equal to about 1.7 standard errors. The Ehrlich models give estimates that are about one-half higher than the full model, and more than twice as big as those from stepwise regression (in absolute value). There is clearly considerable model uncertainty about this parameter.

Now consider β_{15} , the effect of the average time served in state prisons. Whether this is significant at all is not clear, and t -tests based on different models lead to conflicting conclusions. In the full model, β_{15} has a non-significant p -value of .133, while stepwise regression leads to a model that does not include this variable. On the other hand, Mallows' C_p leads to a model in which the p -value for β_{15} is significant at the .05 level, while with adjusted R^2 it is again not significant. In contrast, in Ehrlich's models it is highly significant.

Together these results paint a confused picture about β_{14} and β_{15} . Below we will argue that the confusion can be resolved by taking explicit account of model uncertainty.

5.1.2 Crime and Punishment: Model Averaging

For the model averaging strategies, all possible combinations of predictors were assumed to be equally likely *a priori*. To implement Occam's Window, we started from the null model and used the "Up" algorithm only (see Madigan and Raftery 1994). The selected models and their posterior model probabilities are shown in Table 2. The models with posterior model probabilities of 1.2% or larger as indicated by MC^3 are shown in Table 3. In total, 1772 different models were visited during 30,000 iterations of MC^3 . Occam's Window chose 22 models in this example, clearly indicating

model uncertainty. Choosing any one model and making inferences as if it were the “true” model ignores model uncertainty. The consequences of basing inferences on a single model will be explored further in the next section.

The top models indicated by the two methods (Tables 2 and 3) are quite similar. The posterior probabilities are normalized over all selected models for Occam’s Window and over all possible combinations of the 15 predictors for MC³. So, the posterior probabilities for the same models differ across the model averaging method, but this has little effect on the relationships between the models as measured by the Bayes factor.

Table 4 shows the posterior probability that the coefficient for each predictor does not equal 0, i.e., $\Pr(\beta_i \neq 0|D)$, obtained by summing the posterior model probabilities across models for each predictor. The results from Occam’s Window and MC³ are fairly close for most of the predictors. There are several predictors with high $\Pr(\beta_i \neq 0|D)$ including the proportion of young males, mean years of schooling, police expenditure, income inequality, and probability of imprisonment.

Comparing the two models analyzed by Ehrlich (1973), consisting of the predictors (9, 12, 13, 14, 15) and (1, 6, 9, 10, 12, 13, 14, 15), with the results in Table 4, we see that there are several predictors included in Ehrlich’s analysis that receive little support from the data. The estimated $\Pr(\beta_i \neq 0|D)$ is quite small for predictors 6, 10, 12, and 15. There are also variables for which there is empirical support but which Ehrlich did not include (3 and 4). Indeed, Ehrlich’s two selected models have very low posterior probabilities.

Ehrlich’s work attracted attention primarily because of his conclusion that both the probability of imprisonment (predictor 14) and the average prison term (predictor 15) reduced the crime rate. The posterior distributions for the coefficients of these

Table 2: Crime data: Occam's Window Posterior Model Probabilities.

Model					Posterior model probability %
1	3 4	9 11	13 14		12.6
1	3 4	11	13 14		9.0
1	3 4	9	13 14		8.4
1	3 5	9 11	13 14		8.0
	3 4	8 9	13 14		7.6
1	3 4		13 14		6.3
1	3 4	11	13		5.8
1	3 5	11	13 14		5.7
1	3 4		13		4.9
1	3 5	9	13 14		4.8
	3 5 8 9		13 14		4.4
	3 4	9	13 14		4.1
	3 5 9		13 14		3.6
1	3 5		13 14		3.5
	2 3 4		13 14		2.0
1	3 5	11	13		1.9
	3 4		13 14		1.6
	3 5		13 14		1.6
	3 4		13		1.4
1	3 5		13		1.4
	3 5		13		0.7
1	4		12 13		0.7

Table 4: Crime data: $\Pr(\beta_i \neq 0|D)$, expressed as a percentage

Predictor number	Predictor	Occam's Window	MC ³	Ehrlich's models	
1	percent of males 14-24	73	79	★	
2	indicator variable for southern state	2	17		
3	mean years of schooling	99	98		
4	police expenditure in 1960	64	72		
5	police expenditure in 1959	36	50		
6	labor force participation rate	0	6	★	
7	number of males per 1000 females	0	7		
8	state population	12	23		
9	number of nonwhites per 1000 people	53	62	*	★
10	unemployment rate of urban males 14-24	0	11	★	
11	unemployment rate of urban males 35-39	43	45		
12	wealth	1	30	*	★
13	income inequality	100	100	*	★
14	probability of imprisonment	83	83	*	★
15	average time served in state prisons	0	22	*	★

predictors, based on the model averaging results of MC³, are shown in Figures 3 and 4. The MC³ posterior distribution for β_{14} is indeed centered away from 0 with a small spike at 0. The posterior distribution for β_{14} based on Occam's Window is quite similar. The spike corresponds to $P(\beta_{14} = 0|D)$. This is an artifact of our approach in which it is possible to consider models with a predictor fully removed from the model. This is in contrast to the practice of setting the predictor close to 0 with high probability as in George and McCulloch (1993). In contrast to Figure 3, the MC³ posterior distribution for the coefficient corresponding to average prison term is centered close to 0 and has a large spike at 0 (Figure 4). Occam's Window indicates a spike at 0 only, or no support for inclusion of this predictor. By averaging over all models, our results indicate support for a relationship between crime rate and

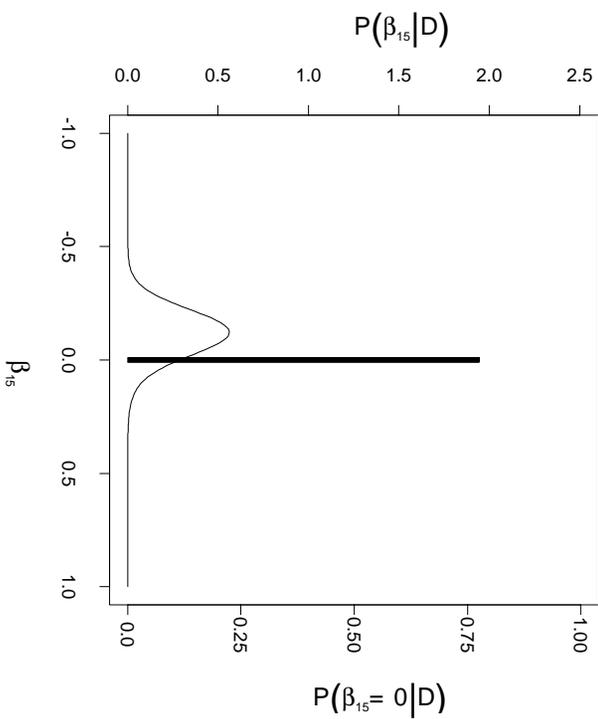


Figure 4: Posterior distribution for β_{15} , the coefficient for the predictor “average time served in state prisons”, based on the model average over a large set of models from MC³. See Figure 3.

predictor 14, but not predictor 15. Our model averaging results are consistent with those of Ehrlich for the probability of imprisonment, but not for the average prison term.

Among the variables that measure the expected benefits from crime, Ehrlich concluded that both wealth and income inequality had an effect; we found this to be true for income inequality but not for wealth. For the predictors that represent the payoff from legitimate activities, Ehrlich found the effects of variables 1, 6, 10 and 11 to be unclear; he did not include mean schooling in his model. We found strong evidence for the effect of some of these variables, notably the percent of young males and mean schooling, but the effects of unemployment and labor force participation are either improved or unlikely. Finally, the “control” variables that have no theoretical basis (2, 7, 8) turned out, satisfyingly, to have no empirical support either.

The model averaging results for the predictors for police expenditures lead to an interesting interpretation. Police expenditure was measured in two successive years and the measures are highly correlated ($r = .993$). The data show clearly that the 1960 crime rate is associated with police expenditures, and that only one of the two measures (X_4 and X_5) is needed, but they do not say for sure which measure should be used. Each model in Occam’s Window and each model visited by MC³ contains one predictor or the other, but not both. For both methods we have $\Pr[(\beta_4 \neq 0) \cup (\beta_5 \neq 0) | D] = 1$, so the data provide very strong evidence for an association with police expenditures.

In summary, we found strong support for some of Ehrlich’s conclusions but not for others. In particular, by averaging over all models, our results indicate support for a relationship between crime rate and probability of imprisonment, but not for average time served in state prisons.

5.1.3 Crime and Punishment: Assessment of Predictive Performance

We use the predictive ability of the selected models for future observations to measure the effectiveness of a model selection strategy. Our specific objective is to compare the quality of the predictions based on model averaging with the quality of predictions based on any single model that an analyst might reasonably have selected.

To measure performance we randomly split the complete data set into two subsets. Other percentage splits can be adopted. A 50/50 split was chosen here so that each portion would contain enough data to be a representative sample. We ran Occam’s Window and MC³ using half of the data. This set is called the training set, D^T . We evaluated performance using the prediction set made up of the remaining half of the data, $D^P = D \setminus D^T$. Within this framework, we assessed predictive performance

using numerical and graphical measures of performance.

Predictive coverage was measured using the proportion of observations in the performance set that fall in the corresponding 90% prediction interval. For both Occam’s Window and MC^3 , 80% of the observations in the performance set fell in the 90% prediction intervals over the averaged models (Table 5). David Draper (personal communication) suggested that BMA falls somewhat short of nominal coverage here because aspects of model uncertainty other than model selection have not been assessed. In Hoeting *et al.* (1995, 1996) we extend BMA to account for uncertainty in the selection of transformations and in the identification of outliers.

For comparison with other standard variable selection techniques, three popular variable selection procedures, discussed above, were used to select two or three “best” models. The models chosen using these methods are given in Table 5. All of the individual models chosen using standard techniques performed considerably worse than the model averaging approaches, with prediction coverage ranging from 58% to 67%. Thus the model averaging strategies improved predictive coverage substantially compared with any single model that might reasonably have been chosen.

A sensitivity analysis for priors chosen within the framework described in Section 3.2 indicates that the results for Occam’s Window and MC^3 are not highly sensitive to the choice of prior. The results for Occam’s Window and MC^3 using 3 different sets of priors were quite similar.

In an attempt to provide a graphical measure of predictive performance, a “calibration plot” was used to determine if the predictions were well calibrated. A model is well calibrated if, for example, 70% of the observations in the test data set are less than or equal to the 70th percentile of the posterior predictive distribution. The calibration plot shows the degree of calibration for different models with the pos-

Table 5: Crime data: Performance comparison. Predictive coverage % is the percentage of observations in the performance set that fall in the 90% prediction interval. Model numbers correspond to the i th model chosen using the given model selection method. For example, C_p (1) is the first model chosen using the C_p method. The percentage values shown for the stepwise procedures correspond to the significance levels for the F-to-enter and F-to-delete values. For example, F=3.84 corresponds approximately to the 5% level.

Method	Model	Predictive coverage %
MC ³	model averaging	80
Occam's Window	model averaging	80
Stepwise (5%)	3 4 9 13	67
Adjusted R ² (2)	1 2 3 4 5 8 11 12 13 15	67
Adjusted R ² (3)	1 2 3 4 5 6 8 11 12 13 15	67
Stepwise (15%)	3 4 8 9 13 15	63
C_p (2)	1 2 3 4 11 13	63
Adjusted R ² (1)	1 2 3 4 5 11 12 13 15	58
C_p (1)	1 2 3 4 11 13 15	58
C_p (3)	1 2 3 4 11 12 13 15	58

terior predictive probability on the x -axis and the percentage of observed data less than or equal to the posterior predictive probability on the y -axis. In a calibration plot, perfect calibration is the 45° line and so the closer the a model's calibration line is to the 45° line, the better calibrated it is. The calibration plot is similar to reliability diagrams used to assess probability forecasts (see, for example, Murphy and Winkler 1977). The calibration plot for the model chosen by stepwise selection and for model averaging using Occam's Window is shown in Figure 5. The shaded area in Figure 5 shows where the model averaging strategy produces predictions that are better calibrated than predictions from the model chosen by the stepwise model selection procedure. The calibration plot for MC³ is similar.

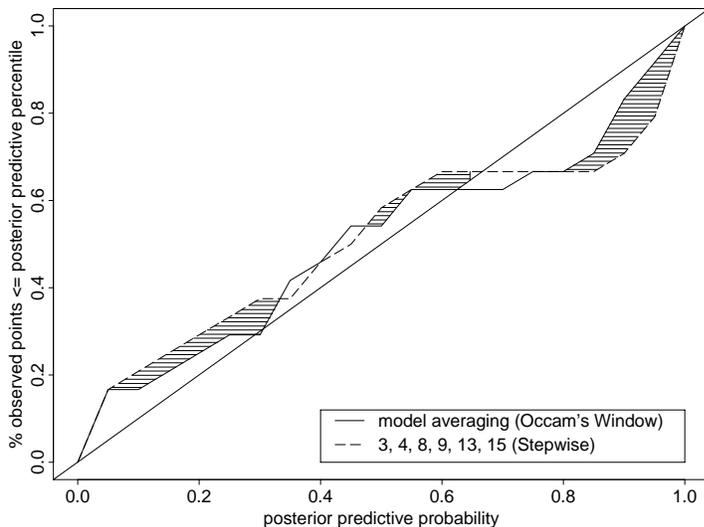


Figure 5: Crime data: Calibration plot

These performance measures support our claim that conditioning on a single selected model ignores model uncertainty which, in turn, leads to the underestimation of uncertainty when making inferences about quantities of interest. Model averaging leads to better calibrated predictive distributions.

5.2 Simulated Examples: Predictive Performance

In the example above, the true answer is unknown. To further demonstrate the usefulness of BMA, we use several simulated examples. In our examples below, we follow the format of George and McCulloch (1993).

Example 5.2.1 In this example we investigate the impact of model averaging on predictive performance when there is little model uncertainty. For the training set, we simulated $p = 15$ predictors and $n = 50$ observations as independent standard normal vectors. The response was generated using the model

$$\mathbf{Y} = \mathbf{X}_4 + \mathbf{X}_5 + \epsilon \quad (12)$$

where $\epsilon \sim N_{50}(0, \sigma^2)$ with $\sigma = 2.5$. Least squares estimates for these data are given in Table 6. There is little model uncertainty in this example; only the p -values for β_4 and β_5 were smaller than 0.1. Fifty additional observations were generated in the

Table 6: Least Squares Estimates for Example 5.2.1 ($\hat{\sigma} = 2.9$).

	β_0	β_1	β_2	β_3	β_4	β_5	β_6	β_7	β_8	β_9	β_{10}	β_{11}	β_{12}	β_{13}	β_{14}	β_{15}
β	0	0	0	0	1.00	1.00	0	0	0	0	0	0	0	0	0	0
$\hat{\beta}$.42	.21	.40	.07	.95	1.72	.20	.34	-.32	.24	-.15	.6	-.45	-.08	.20	.18
$\hat{\sigma}_\beta$.46	.55	.56	.36	.52	.47	.39	.58	.49	.45	.44	.55	.48	.52	.45	.47

Table 7: Performance comparison for Example 5.2.1. Predictive coverage for a 90% prediction interval. Predictive coverage for BMA (all models) is estimated using the 371 models with posterior model probabilities greater than 0.0001. See Table 5.

Method	Model	Predictive coverage %
BMA (estimated coverage)	model averaging	72
Occam's Window	model averaging	70
Adj R ² (3)	2 4 5 8 11	70
Cp (3)	4 5 11	70
True model & Stepwise (5%)	4 5	68
Stepwise (15%) & Cp (2)	2 4 5	68
Cp(1)	4 5	68
Adj R ² (2)	2 4 5 10 11	68
Adj R ² (1)	2 4 5 11	66

same manner to create the prediction set.

In this example the true model, the model averaging techniques, and models selected using standard techniques all have poor predictive coverage (Table 7). It is slightly encouraging that BMA performs better than the true model, but the improvement is too small to be significant. This and other similar examples simulated by the authors show that when there is very little model uncertainty, predictive performance is not significantly improved by model averaging.

Example 5.2.2 This example demonstrates the performance of BMA when a subset

Table 8: Least Squares Estimates for Example 5.2.2 ($\hat{\sigma} = 2.21$).

	β_0	β_1	β_2	β_3	β_4	β_5	β_6	β_7	β_8	β_9	β_{10}	β_{11}	β_{12}	β_{13}	β_{14}	β_{15}
β	0	1.00	1.00	1.00	1.00	1.00	0	0	0	0	0	0	0	0	0	0
$\hat{\beta}$	0.12	.80	1.07	1.03	-.18	.55	-.67	.28	-.11	.31	.29	.11	-.09	-.39	.73	-.96
$\hat{\sigma}_\beta$.38	.49	.41	.45	.53	.58	.37	.41	.49	.33	.34	.40	.32	.35	.37	.40

of the predictors is correlated. For the training set, we simulated $p = 15$ predictors and $n = 50$ observations. Predictors 1 through 10 were obtained as independent standard normal vectors, $\mathbf{X}_1, \dots, \mathbf{X}_{10}$ iid $\sim N(0, 1)$, and predictors 11 through 15 were generated using the framework

$$[\mathbf{X}_{11}, \dots, \mathbf{X}_{15}] = [\mathbf{X}_1, \dots, \mathbf{X}_5] \left([0.3 \ 0.5 \ 0.7 \ 0.9 \ 1.1]^T [1 \ 1 \ 1 \ 1 \ 1] \right) + \epsilon$$

where $\epsilon \sim N(0, 1)$. The response was generated using the model

$$\mathbf{Y} = \mathbf{X}_1 + \mathbf{X}_2 + \mathbf{X}_3 + \mathbf{X}_4 + \mathbf{X}_5 + \epsilon \quad (13)$$

where $\epsilon \sim N_{50}(0, \sigma^2)$ with $\sigma = 2.5$. Least squares estimates for these data are given in Table 8. The correlation structure resulted in moderate pairwise correlation between predictors 1 to 5 and 11 to 15 ($\text{corr}(X_1, X_{11})=0.39$, $\text{corr}(X_2, X_{12})=0.41$, $\text{corr}(X_3, X_{13})=0.56$, $\text{corr}(X_4, X_{14})=0.71$, $\text{corr}(X_5, X_{15})=0.69$) and small pairwise correlations elsewhere (median correlation equal to -0.02). Fifty additional observations were generated in the same manner to create the prediction set.

Table 9 shows that in this example model averaging has better predictive performance than any single model that might have been selected. In this example, the poor performance of the true model and the other single models selected using standard techniques demonstrate that model uncertainty can strongly influence predictive performance.

Table 9: Performance comparison for Example 5.2.2. Predictive coverage for a 90% prediction interval. Predictive coverage for BMA (all models) is estimated using the 1014 models with posterior model probabilities greater than 0.00005. See Table 5.

Method	Model	Predictive coverage %
MC ³	model averaging	92
Occam's Window	model averaging	86
Stepwise (5 & 15%)	1 2 3 4	80
True Model	1 2 3 4 5	78
C_p (2) & Adjusted R ² (1)	1 2 3 6 13 14 15	72
C_p (3)	1 2 3 6 10 14 15	72
Adjusted R ² (3)	1 2 3 6 7 13 14 15	72
C_p (1)	1 2 3 5 14 15	70
Adjusted R ² (2)	1 2 3 6 10 13 14 15	70

6 Successful identification of the null model

Linear regression models are frequently used even when little is known about the relationship between the predictors and the response. When there is a weak relationship between the predictors and the response, the overall F -statistic will be small and thus the null hypothesis that the null model is true fails to be rejected. However, many data analysts carry out model selection regardless of the F -statistic value for the overall model. Problems can then occur as subsequent model selection techniques often choose a model which includes a subset of the predictors. Freedman (1983) has shown that in the extreme case where there is *no* relationship between the predictors and the response variable, omitting the predictors with the smallest t-values (e.g., $p > 0.25$) can result in a model with a highly significant F statistic and high R². In contrast, if the response and predictors are independent, Occam's Window typically

indicates the null model only, or as one of a small number of “best” models.

Following Freedman (1983), we generated 5100 independent observations from a standard normal distribution to create a matrix with 100 rows and 51 columns. The first column was taken to be the dependent variable in a regression equation and the other 50 columns were taken to be the predictors. Thus the predictors are independent of the response by construction. For the entire data set, the multiple regression results were as follows:

- $R^2 = .55$, $p = .29$;
- 18 coefficients out of 50 were significant at the .25 level;
- 4 coefficients out of 50 were significant at the .05 level.

Three different variable selection procedures were used on the simulated data. The first of these was the method used by Freedman (1983), in which all predictors with p -values of 0.25 or lower were included in a second pass over the data. The results from this method were as follows:

- $R^2 = .40$, $p = .0003$;
- 17 coefficients out of 18 were significant at the .25 level;
- 10 coefficients out of 18 were significant at the .05 level.

These results are highly misleading as they indicate a definite relationship between the response and the predictors, whereas, in fact, the data are all noise.

The second model selection method used on the full data set was Efron's stepwise method. This indicated a model with 15 predictors with the following results:

- $R^2 = .40$, $p = .0001$;
- all 15 were significant at the .25 level;

- 10 coefficients out of 15 were significant at the .05 level.

Again a model is chosen which misleadingly appears to have a great deal of explanatory power.

The third variable selection method used was Occam's Window. The only model chosen by this method was the null model.

The procedure described above was repeated 10 times with similar results. In 5 simulations, Occam's Window chose only the null model. For the remaining simulations 3 models or fewer were chosen along with the null model. For the non-null models that were chosen, all models had R^2 values less than 0.15. For all of the simulations the selection procedure used by Freedman (1983) and the stepwise method chose models with many predictors and highly significant R^2 values.

At best, Occam's Window correctly indicates that the null model is the only model that should be chosen when there is no signal in the data. At worst, Occam's Window chooses the null model along with several other models. The presence of the null model among those chosen by Occam's Window should indicate to a researcher that there may be evidence for a lack of signal in the data he or she is analyzing.

To examine the possibility that our Bayesian approach favors parsimony to the extent that Occam's Window finds no signal even when there is one, we did an additional simulation study. We generated 3000 observations from a standard normal distribution to create a data set with 100 observations and 30 candidate predictors. The response Y was allowed only to depend on X_1 , where $Y = 0.5X_1 + \epsilon$ with $\epsilon \sim N(0,0.75)$. Thus Y still has unit variance and the "true" R^2 for the model equals 0.20.

For this simulated data, Occam's Window contained one model only, the correct model with X_1 . In contrast, the screening method used by Freedman produced a

model with 6 predictors, including X_1 , with 4 of them significant at the 0.1 level. Stepwise regression indicated a model with 2 predictors, including X_1 , both of them significant at the 0.025 level. So the two standard variable selection methods indicated evidence for variables that were in fact not at all associated with the dependent variable while Occam's Window chose the correct model.

These examples provide evidence that Occam's Window overcomes the problem of selection of the null model when there is no signal in the data.

7 Discussion

7.1 Related Work

Draper (1995) has also addressed the problem of assessing model uncertainty. Draper's approach is based on the idea of *model expansion*, i.e., starting with a single reasonable model chosen by a data-analytic search, expanding model space to include those models which are suggested by context or other considerations, and then averaging over this model class. Draper does not directly address the problem of model uncertainty in variable selection. However, one could consider Occam's Window to be a practical implementation of model expansion.

George and McCulloch (1993) have developed the Stochastic Search Variable Selection (SSVS) method, which is similar in spirit to MC³. They define a Markov chain which moves through model space and parameter space at the same time. Their method never actually removes a predictor from the full model, but only sets it close to zero with high probability. Our approach avoids this by integrating analytically over parameter space.

We have focused here on Bayesian solutions to the model uncertainty problem. There has been very little written about frequentist solutions to the problem. Perhaps

the most obvious frequentist solution is to bootstrap the entire data analysis, including model selection. However, Freedman *et al.* (1986) have shown that this does not necessarily give a satisfactory solution to the problem.

7.2 Conclusions

The prior distribution of the covariance matrix for β described in Section 3.2 depends on the actual data, including both the dependent and independent variables. A similar data-dependent approach to the assessment of the priors was used by Raftery (1996). While this may appear at first sight to be contrary to the idea of a prior, our objective was to develop priors that lead to posteriors similar to those of a person with little prior information. Examples analyzed to date suggest that this objective was achieved. The priors for β lead to a reasonable prior variance and result in conclusions that are not highly sensitive to the choice of hyperparameters. Thus the data-dependence does not appear to be a drawback.

In a strict sense, our data dependent priors do not correspond to a Bayesian subjective prior. Our priors might be considered to be an approximation to a true Bayesian subjective prior and might be appropriate when little prior information is available. We have followed other authors, including Zellner (1986), George and McCullough (1993), and Laud and Ibrahim (1995), in referring to our approach as Bayesian.

The choice of which procedure to use — Occam’s Window or MC³ — will depend on the particular application. Occam’s Window will be most useful when one is interested in making inferences about the relationships between the variables. Occam’s Window also tends to be much faster computationally. MC³ is the better procedure to choose if the goal is good predictions or if the posterior distribution of some quan-

tity is of more interest than the nature of the “true” model and if computer time is not a critical consideration. However, each approach is flexible enough to be used successfully for both inference *and* prediction.

We have described two procedures that can be used to account for model uncertainty in variable selection for linear regression models. In addition to variable selection, there is also uncertainty involved in the identification of outliers and in the choice of transformations in regression. To broaden the flexibility of our current procedures as well as to improve our ability to account for model uncertainty, we have extended BMA to include transformation selection and outlier identification (in work reported elsewhere Hoeting *et al.* 1995, 1996).

Appendix A: Data for Figure 1

Data from selected textbooks used to make Figure 1.

Data set	Source	page number	number of observations	number of predictors
Attitude Survey	Chatterjee and Price (1991)	70	30	6
Equal Education Opportunity	Chatterjee and Price (1991)	176	70	3
Gasoline Mileage	Chatterjee and Price (1991)	261	30	10
Nuclear Power	Cox and Snell (1982)	81	32	10
Crime	Cox and Snell (1982)	170	47	13
Hald	Draper and Smith (1981)	630	13	4
Grades	Hamilton (1993)	83	118	3
Swiss Fertility	Mosteller and Tukey (1977)	550	47	5
Surgical Unit	Neter, Wasserman and Kutner (1990)	439, 468	108	4
Berkeley Study	Weisberg (1985)			
Girls		56	32	10
Boys		57	26	10
Housing	Weisberg (1985)	241	27	9
Highway	Weisberg (1985)	206	39	13

Appendix B: Software for Implementing MC³

BMA is a set of S-PLUS functions which can be obtained free of charge via the World Wide Web address <http://lib.stat.cmu.edu/S/bma> or by sending an e-mail message containing the text “send BMA from S” to the Internet address *statlib@stat.cmu.edu*.

The program MC3.REG performs Markov chain Monte Carlo model composition for linear regression. The set of programs fully implements the MC³ algorithm described in Section 4.2.

References

- Becker, G.S. (1968), Crime and punishment: An economic approach. *Journal of Political Economy*, **76**, 169–217.
- Brier, S.S. and Fienberg, S.E. (1980), Recent econometric modeling of crime and punishment: Support for the deterrence hypothesis? *Evaluation Review*, **4**, 147–191.
- Breiman, L. (1968), *Probability*, Addison-Wesley, Reading.
- Breiman, L. (1992), The little bootstrap and other methods for dimensionality selection in regression: X-fixed prediction error, *Journal of the American Statistical Association*, **87**, 738–754.
- Breiman, L. and Spector, P. (1992), Submodel selection and evaluation in regression, *International Statistical Review*, **60**, 291–319.
- Breiman, L. (1995), Better subset regression using the nonnegative garrote, *Technometrics*, **37**, 373–384.
- Chatterjee, S. and Price, B. (1991), *Regression Analysis by Example*, 2nd edition, New York: Wiley.
- Cox, D. R. and Snell, E. J. (1982), *Applied Statistics : Principles and Examples*, New York: Chapman and Hall.
- Chung, Kai Lai (1967), *Markov Chains with Stationary Transition Probabilities* (2nd ed), Berlin: Springer-Verlag.
- Draper, D. (1995), Assessment and propagation of model uncertainty (with Discussion), *Journal of the Royal Statistical Society B*, **57**, 45–97.
- Draper, N.R. and Smith, H. (1981), *Applied Regression Analysis*, (2nd. edition), New York: Wiley.
- Edwards, W., Lindman, H. and Savage, L.J. (1963), Bayesian statistical inference for

- psychological research, *Psychological Review*, **70**, 193–242.
- Ehrlich, I. (1973), Participation in illegitimate activities: a theoretical and empirical investigation, *Journal of Political Economy*, **81**, 521–565.
- Freedman, D.A. (1983), A note on screening regression equations, *The American Statistician*, **37**, No. 2, 152–155.
- Freedman, D. A., Navidi, W.C. and Peters, S.C. (1986), On the impact of variable selection in fitting regression equations, In *On Model Uncertainty and its Statistical Implications* (T.K. Dijkstra, ed.), Berlin: Springer-Verlag, pp. 1–16.
- Garthwaite, P.H. and Dickey, J.M. (1992), Elicitation of prior distributions for variable-selection problems in regression, *Annals of Statistics*, **20**, No. 4, 1697–1719.
- Geisser, S. (1980), Discussion on sampling and Bayes' inference in scientific modelling and robustness (by G.E.P. Box), *Journal of the Royal Statistical Society A*, **143**, 416–417.
- George, E.I. and McCulloch, R.E. (1993), Variable selection via Gibbs sampling, *Journal of the American Statistical Association*, **88**, No. 423, 881–890.
- Good, I.J. (1952), Rational decisions, *Journal of the Royal Statistical Society B*, **14**, 107–114.
- Hamilton, L.C. (1993), *Statistics with Stata 3*, Belmont, CA: Duxbury Press.
- Hocking, R.R. (1976), The analysis and selection of variables in linear regression, *Biometrics*, **32**, 1–51.
- Hodges, J.S. (1987), Uncertainty, policy analysis and statistics, *Statistical Science*, **2**, 259–291.
- Hoeting, J.A., Raftery, A.E., and Madigan, D. (1996), “A Method for Simultaneous Variable Selection and Outlier Identification in Linear Regression”, to appear in *Journal of Computational Statistics and Data Analysis*.

- Hoeting, J.A., Raftery, A.E., and Madigan, D. (1995), “Simultaneous Variable and Transformation Selection in Linear Regression”, Technical Report 9506, Department of Statistics, Colorado State University.
- Jeffreys, H. (1961), *Theory of Probability*, (3rd ed.), Oxford University Press.
- Kadane, J.B., Dickey, J.M., Winkler, R.L., Smith, W.S. and Peters, S.C. (1980), Interactive elicitation of opinion for a normal linear model, *Journal of the American Statistical Association*, **75**, 845–854.
- Kass, R.E. and Raftery, A.E. (1995), Bayes factors, *Journal of the American Statistical Association*, **90**, 773–795.
- Laud, P.W. and Ibrahim, J.G. (1995), Predictive Model Selection, *Journal of the Royal Statistical Society - B*, **57**, 247-262.
- Leamer, E.E. (1978), *Specification Searches*, New York: Wiley.
- Linhart, H. and Zucchini, W. (1986), *Model Selection*, New York: Wiley.
- Madigan, D. and Raftery, A.E. (1994), Model selection and accounting for model uncertainty in graphical models using Occam’s Window, *Journal of the American Statistical Association*, **89**, 1535-1546.
- Madigan, D. and York, J. (1995), Bayesian graphical models for discrete data, *International Statistical Review*, **63**, 215-232.
- Miller, A.J. (1984), Selection of subsets of regression variables (with Discussion), *Journal of the Royal Statistical Society (Series A)*, **147**, 389–425.
- Miller, A.J. (1990), *Subset Selection in Regression*, New York: Chapman-Hall.
- Mitchell, T.J. and Beauchamp, J.J. (1988), Bayesian variable selection in linear regression (with discussion), *Journal of the American Statistical Association*, **83**, 1023–1036.
- Mosteller, F. and Tukey, J.W. (1977), *Data Analysis and Regression*, Reading, Mass.:

Addison–Wesley.

- Moulton, B.R. (1991), A Bayesian approach to regression selection and estimation with application to a price index for radio services, *Journal of Econometrics*, **49**, 169–193.
- Murphy, A.H. and Winkler R.L. (1977), Reliability of subjective probability forecasts of precipitation and temperature, *Applied Statistics*, **26**, 41–47.
- Neter, J., Wasserman, W., and Kutner, M. (1990), *Applied Linear Statistical Models*, Homewood, IL: Irwin.
- Raftery, A.E. (1988), Approximate Bayes factors for generalized linear models. *Technical Report no. 121*, Department of Statistics, University of Washington.
- Raftery, A.E. (1996), Approximate Bayes factors and accounting for model uncertainty in generalized linear models, *Biometrika*, to appear.
- Raiffa, H. and Schlaifer, R. (1961), *Applied Statistical Decision Theory*, Cambridge, MA: The MIT Press.
- Regal, R. and Hook, E.B. (1991), The effects of model selection on confidence intervals for the size of a closed population, *Statistics in Medicine*, **10**, 717–721.
- Schwarz, G. (1978), Estimating the dimension of a model, *The Annals of Statistics*, **6**, 461–464.
- Shibata, R. (1981), An optimal selection of regression variables, *Biometrika*, **68**, 45–54.
- Smith, A.F.M. and Roberts, G.O. (1993), Bayesian computation via Gibbs sampler and related Markov chain Monte Carlo methods, *Journal of the Royal Statistics Society B*, **55**, 3–24.
- Stewart, L. (1987), Hierarchical Bayesian analysis using Monte Carlo integration: Computing posterior distributions when there are many possible models, *The*

Statistician, 36, 211-219.

Stewart, L. and Davis, W.W. (1986), Bayesian posterior distributions over sets of possible models with inferences computed by Monte Carlo integration, *The Statistician*, 35, 175–182.

Stigler, G.J. (1970), The optimum enforcement of laws. *Journal of Political Economy*, 78, 526–536.

Taft, D.R. and England, R.W. (1964), *Criminology* (4th ed.). New York: Macmillan.

Vandaele, W. (1978), Participation in illegitimate activities; Ehrlich revisited, In *Deterrence and Incapacitation*, (eds. Blumstein, A., Cohen, J. and Nagin, D.). Washington, D.C.: National Academy of Sciences, 270–335.

Weisberg, S. (1985), *Applied Linear Regression*, 2nd ed., New York: Wiley.

Zellner, A. (1986), On assessing prior distributions and Bayesian regression analysis with g -prior distributions, In *Bayesian Inference and Decision Techniques-Essays in Honor of Bruno de Finetti*, (eds. P.K. Goel and A. Zellner), Amsterdam: North-Holland, 233–243.