

Minimum Description Length Model Selection Criteria for Generalized Linear Models

Mark H. Hansen and Bin Yu

Abstract

This paper derives several model selection criteria for generalized linear models (GLMs) following the principle of Minimum Description Length (MDL). We focus our attention on the mixture form of MDL. Normal or normal-inverse gamma distributions are used to construct the mixtures, depending on whether or not we choose to account for possible over-dispersion in the data. In the latter case, we apply Efron's [6] double exponential family characterization of GLMs. Standard Laplace approximations are then employed to derive computationally tractable selection rules. Each of the criteria we construct have adaptive penalties on model complexity, either explicitly or implicitly. Theoretical results for the normal linear model, and a set of simulations for logistic regression illustrate that mixture MDL can "bridge" the selection "extremes" AIC and BIC in the sense that it can mimic the performance of either criterion, depending on which is best for the situation at hand.

Keywords. AIC; Bayesian Methods; BIC; Code Length; Information Theory; Minimum Description Length; Model Selection; Generalized Linear Models.

1 Introduction

Statistical model selection attempts to decide between competing model classes for a data set. As a principle, maximum likelihood is not well suited for this problem as it suggests choosing the largest model under consideration. Following this strategy, we tend to overfit the data and choose models that have poor prediction power. Model selection emerged as a field in the 70's, introducing procedures that "corrected" the maximum likelihood approach. The most famous and widely used criteria are *A Information Criterion* (AIC) of Akaike [1, 2] and the *Bayesian Information Criterion* (BIC) of Schwarz [15]. They both take the form of a penalized maximized likelihood, but with different penalties: AIC adds 1 for each additional variable included in a model, while BIC adds $\log n/2$, where n is the sample size. Theoretical and simulation studies (cf. Shibata [16], Speed and Yu [18], and references therein), mostly in the

regression case, have revealed that when the underlying model is finite-dimensional (specified by a finite number of parameters), BIC is preferred; but when it is infinite-dimensional, AIC performs best. Unfortunately, in practical applications we rarely have this level of information about how the data were generated, and it is desirable to have selection criteria which perform well independent of the form of the underlying model. That is, we seek criteria which adapt automatically to the situation at hand. In this paper, we derive such adaptive model selection criteria for generalized linear models (GLMs) under the Minimum Description Length (MDL) framework. With MDL we find several generic prescriptions or “forms” for constructing such selection criteria. In this paper, we focus on one MDL form that is based on mixtures.

The MDL approach began with Kolmogorov’s theory of algorithmic complexity, matured in the literature on information theory, and has recently received renewed interest within the statistics community. By viewing statistical modeling as a means of generating *descriptions* of observed data, the MDL framework (cf. Rissanen [13], Barron et al, [3], and Hansen and Yu [8]) discriminates between competing model classes based on the *complexity* of each description. Precisely, the Minimum Description Length (MDL) Principle recommends that we

Choose the model that gives the shortest description of data.

While there are many kinds of descriptions and many ways to evaluate their complexity, we follow Rissanen [13] and use a *code length* formulation based on the candidate model.

To make this more precise, we first recall that for each probability distribution Q on a finite set \mathcal{A} there is an associated *code* that prepares elements of \mathcal{A} for transmission across some (noiseless) communication channel. We will consider binary codes, meaning that each codeword is a string of 0’s and 1’s. It is possible to find a code so that the number of bits (the number of 0’s and 1’s in a codeword) used to encode each symbol of $a \in \mathcal{A}$ is essentially $-\log_2 Q(a)$; that is, $-\log_2 Q$ can be thought of as a *code length function*. Huffman’s algorithm [5] takes a distribution Q and produces a so-called *prefix code* with the right length function.¹ Conversely, any integer-valued function L corresponds to the code length of some binary prefix code if and only if it satisfies Kraft’s inequality

$$\sum_{a \in \mathcal{A}} 2^{-L(a)} \leq 1, \tag{1}$$

see Cover and Thomas [5] for a proof. Therefore, given a prefix code on \mathcal{A} with length function L , we can define a distribution on \mathcal{A} as follows,

$$Q(a) = \frac{2^{-L(a)}}{\sum_{z \in \mathcal{A}} 2^{-L(z)}} \quad \text{for any } a \in \mathcal{A}.$$

¹While the details are beyond the scope of this short paper, the interested reader is referred to Hansen and Yu [8] and Cover and Thomas [5].

With Kraft's inequality, we find a correspondence between codes and probability distributions. In what follows, we will work with natural logs and take $-\log Q$ to be an idealized code length.

One of the early problems in information theory involves transmitting symbols that are randomly generated from a probability distribution P defined on \mathcal{A} . Let $A \in \mathcal{A}$ denote a random variable with this distribution. Now, from the discussion in the previous paragraph, any code defined on \mathcal{A} can be associated with an idealized code length function $-\log Q$ for some distribution Q . With this setup, the expected code length for symbols generated from P is given by $-E \log Q(A) = -\sum_a P(a) \log Q(a)$. By Jensen's inequality, we see that the shortest code length is achieved by a code that has $-\log P$ as its idealized length function. That is, the expected code length is bounded from below by $-E \log P(a) = -\sum_a P(a) \log P(a)$, the entropy of P . In the literature on information theory, this fact is known as Shannon's Inequality.

In this paper, we focus on descriptions of data that consist of probability models, and compare them based on the efficiency of the corresponding code in terms of improvements in code length relative to the entropy of the data generating process. When the competing models are members of a parametric family, using MDL to select a model, or rather, to estimate a parameter, is equivalent to maximum likelihood estimation (when the cost of transmitting the parameter estimate is fixed). To compare different model classes, different parametric families, or carry out model selection from among several candidate model classes, efficient codes for each class need to *fairly represent* its members. We will not elaborate on this idea, but instead comment it is possible to demonstrate rigorously that several coding schemes achieve this fairness and hence provide valid selection criteria (for say, iid or time series observations). We refer readers to Barron et al [3] and Hansen and Yu [8].

Among the schemes that yield valid selection criteria, the best known is the so-called *two-stage code*, in which we first encode the maximum likelihood estimate (MLE) of the parameters in the model, and then use the model with the MLE to encode the data (say, via Huffman's algorithm described above). Hence this form is a penalized likelihood, and to first order is exactly the same as BIC. Other forms of MDL include predictive, mixture and normalized maximum likelihood (NML). The predictive form makes the most sense when the data come in sequentially and has a close connection to prequential inference; the mixture codes will be described in more detail in the next section; and the NML form is new and evolving, and we know the code length expressions only in a few special cases including binomial model and Gaussian linear regression (cf. Rissanen [14], Barron et al [3], and Hansen and Yu [8]).

The rest of the paper is organized as follows. Section 2 gives the details of a mixture code in the context of regression-type models. Section 3 covers the gMDL model selection criterion (so named because of its use of Zellner's g -prior [20]) from Hansen and Yu [8] in the variance known and unknown cases to prepare the reader for the new results in Section 4. The criterion when σ^2 is known appeared originally in George and Foster [7] in the context of a Bayesian analysis. Section 3.3 contains

a new theorem to show the bridging effect of the gMDL criterion between AIC and BIC in a normal linear regression model.

Section 4 derives a version of the mixture form gMDL for GLMs. In this case, normal or normal-inverse gamma distributions are used to construct a mixture model, depending on whether or not we choose to account for possible over-dispersion in the data. When the dispersion parameter is known, the resulting criterion appeared first in Peterson [11] in the context of a Bayesian analysis. To account for dispersion effects, we use Efron’s [6] double exponential family characterization of GLMs as the likelihood. Standard Laplace approximations are employed to derive computationally tractable selection rules. Each of the criteria we construct have adaptive penalties on model complexity, either explicitly or implicitly. The last section of the paper contains a set of simulations for logistic regression to illustrate that mixture MDL can “bridge” AIC and BIC in the sense that it can mimic the performance of either criterion, depending on which is best for the situation at hand. The performance measures include the probability of selecting the correct model and test-error based on a selected model. The latter is found to be much less sensitive to the model selection criterion than the former due to the robustness of 0-1 loss in classification.

2 Mixture MDL

In this paper we will consider regression-type models; that is, we would like to characterize the dependence of a random variable $Y \in \mathcal{Y} \subset \mathbb{R}$ on a vector of potential covariates $(X_1, \dots, X_K) \in \mathbb{R}^K$. We consider various parametric model classes (or conditional densities) for Y , indexed by a 0-1 binary vector $\boldsymbol{\gamma} = (\gamma_1, \dots, \gamma_K)$; each model depends on a subset of the covariates corresponding to 1’s in the model index vector $\boldsymbol{\gamma}$. Generically, we will let $\mathcal{M}_\boldsymbol{\gamma}$ denote a simple model class with dimension $k_\boldsymbol{\gamma} = \sum_{j=1}^K \gamma_j$, which depends on the predictors (X_1, \dots, X_K) through the linear combination

$$\sum_{j:\gamma_j=1} \beta_j X_j, \quad (2)$$

where $\boldsymbol{\beta}_\boldsymbol{\gamma} = (\beta_j)_{\{j:\gamma_j=1\}}$ is a vector of parameters. To fit this relationship, our basic data are observations of the form (Y_i, \mathbf{X}_i) , $i = 1, \dots, n$, where $\mathbf{X}_i = (X_{i1}, \dots, X_{iK})$. In observational studies it makes sense to consider \mathbf{X}_i as being random, whereas in designed experiments the values of the covariates are specified. Let $\mathbf{Y} = (Y_1, \dots, Y_n) \in \mathcal{Y}^n$ denote the vector of responses and let \mathbf{X}_K be the $n \times K$ full design matrix, $[\mathbf{X}_K]_{ij} = X_{ij}$. By $\mathbf{X}_\boldsymbol{\gamma}$ we mean a submatrix of \mathbf{X} consisting of those columns j for which $\gamma_j = 1$. We connect the data to the model (2) via the conditional density functions $f_{\boldsymbol{\theta}_\boldsymbol{\gamma}}(\mathbf{y}|\mathbf{X}_\boldsymbol{\gamma})$, $\mathbf{y} \in \mathcal{Y}^n$, for some set of parameters $\boldsymbol{\theta}_\boldsymbol{\gamma} \in \Theta$. (Typically, $\boldsymbol{\theta}_\boldsymbol{\gamma}$ will include regression parameters $\boldsymbol{\beta}_\boldsymbol{\gamma}$ and possibly a dispersion effect.) In order to assess the suitability of $\mathcal{M}_\boldsymbol{\gamma}$, we derive a description length for \mathbf{Y} based on $\mathcal{M}_\boldsymbol{\gamma}$.

For simplicity, we now drop the subscript γ except in places where we feel the need to remind the reader. The reader should interpret the model class \mathcal{M} , its dimension k , the design matrix \mathbf{X} , and the parameters $\boldsymbol{\theta}$ and $\boldsymbol{\beta}$ as all depending on some subset of the available predictors. We then judge the appropriateness of this model based on the so-called *mixture form* of MDL. As its name suggests, this criterion starts with a mixture distribution that combines all the members in the class \mathcal{M}

$$m(\mathbf{y}|\mathbf{X}) = \int f_{\boldsymbol{\theta}}(\mathbf{y}|\mathbf{X})w(\boldsymbol{\theta}|\mathbf{X})d\boldsymbol{\theta}, \quad \mathbf{y} \in \mathcal{Y}^n, \quad (3)$$

where w is a probability density function on $\boldsymbol{\theta}$. This integral has a closed form expression when $f_{\boldsymbol{\theta}}(\cdot|\mathbf{X})$ is an exponential family and w is a conjugate distribution.

If \mathcal{Y} is a finite set of values, we can use the distribution (3) to directly form a *mixture code* for strings $\mathbf{y} \in \mathcal{Y}^n$. In this setting, we assume that both sender and receiver know about the covariates \mathbf{X} , and we only have to transmit \mathbf{y} . As an example, suppose $\mathcal{Y} = \{0, 1\}$ so that \mathbf{y} is a binary string of length n . We use the model class \mathcal{M} and the distribution (3) to construct a mixture code for all 2^n strings $\mathbf{y} \in \{0, 1\}^n$. From our discussion in Section 1, we know that we can apply Huffman's algorithm to build a code that has the (idealized) length function $L(\mathbf{y}) = -\log_2 m(\mathbf{y}|\mathbf{X})$ for all $\mathbf{y} \in \mathcal{Y}^n$. This means that the number of bits required to transmit any $\mathbf{y} \in \{0, 1\}^n$ is essentially $-\log_2 m(\mathbf{y}|\mathbf{X})$. The MDL principle then distinguishes between candidate model classes based on the associated length function $L(\mathbf{Y})$, the number of bits required to transmit the observed data \mathbf{Y} . As mentioned earlier, we have chosen to use base e in the log for our derivations.

In Section 1, we only considered building codes for finite sets of symbols. When $Y_i \in \mathcal{Y} \subset \mathbb{R}$, $i = 1, \dots, n$, is a continuous response, we form an approximate length function by first discretizing the set \mathcal{Y} . That is, given a precision δ we obtain the description length

$$-\log \int f_{\boldsymbol{\theta}}(\mathbf{y}|\mathbf{X})w(\boldsymbol{\theta}|\mathbf{X})d\boldsymbol{\theta} + n \log \delta. \quad (4)$$

Assuming that the precision used for this approximation is the same regardless of model class \mathcal{M} , we again arrive at the expression

$$-\log \int f_{\boldsymbol{\theta}}(\mathbf{y}|\mathbf{X})w(\boldsymbol{\theta}|\mathbf{X})d\boldsymbol{\theta} \quad (5)$$

as a suitable length function. In the next section, we present a brief review of mixture MDL for the simple linear model. A full derivation of these results can be found in Hansen and Yu [8].

When choosing between two model classes, the mixture form of MDL (with fixed hyperparameters) is equivalent to a Bayes factor (Kass and Raftery [9]) based on the same distributions on the parameters spaces. As we will see in the next section, MDL allows for a natural, principled mechanism for dealing with hyperparameters

that distinguishes it from classical Bayesian analysis. Also, keep in mind that w is not introduced as a prior in the Bayesian sense, but rather as a device for creating a distribution for the data \mathbf{Y} from \mathcal{M} . This distinction also allows us more freedom in choosing w , and has spawned a number of novel applications in the engineering literature.

3 Regression

We begin with the simplest GLM, namely the normal linear model \mathcal{M}_γ :

$$Y_i = \sum_{j:\gamma_j=1} \beta_j X_{ij} + \epsilon_i. \quad (6)$$

where the ϵ_i are normally distributed with mean zero and variance σ^2 . To remind the reader that our basic model classes will consist of various subsets of the predictor variables (X_1, \dots, X_K) , we restored the γ notation in the above equation. For simplicity, however, from this point on, we will drop it and consider derivations with respect to a single model class, a single choice of γ . Technically, we do not need to assume that the relationship in (6) holds for some collection of predictors \mathbf{X}_K , but instead we will entertain model classes because they are capable of capturing the major features observed in the observed data string \mathbf{Y} . For comparison with more general GLMs later, we treat separately the case in which σ^2 is known and unknown. In the former case, the parameter vector $\boldsymbol{\theta}$ in the mixture (3) consists only of the coefficients $\boldsymbol{\beta}$; while in the latter, $\boldsymbol{\theta}$ involves both $\boldsymbol{\beta}$ and σ^2 .

We review this material because relatively straightforward, direct analysis yield the MDL selection criteria. When we tackle the complete class of GLMs, things become more difficult, but the final forms are reminiscent of those derived in this section.

3.1 Known error variance σ^2

Here, we take $\boldsymbol{\theta} = \boldsymbol{\beta}$ and let $w(\boldsymbol{\beta}|\mathbf{X})$ be a normal distribution with mean zero and variance-covariance matrix $\sigma^2 V$. As $\mathcal{Y} = \mathbb{R}$, we have to appeal to the discretized form of MDL (4). By using a conjugate distribution, we are able to perform the integration in (5) exactly. This leads to a code length of the form

$$\begin{aligned} L(\mathbf{y}|V) &= -\log m(\mathbf{y}|\mathbf{X}, V) \\ &= \frac{1}{2} \log |V^{-1} + \mathbf{X}^t \mathbf{X}| + \frac{1}{2} \log |V| \\ &\quad + \frac{1}{2\sigma^2} \left(\mathbf{y}^t \mathbf{y} - \mathbf{y}^t \mathbf{X} (V^{-1} + \mathbf{X}^t \mathbf{X})^{-1} \mathbf{X}^t \mathbf{y} \right), \end{aligned} \quad (7)$$

where we have dropped terms that depend only on n . We have also made explicit the dependence of the length function on the variance-covariance matrix V . Clearly,

we can simplify this expression by taking $V = c(\mathbf{X}^t \mathbf{X})^{-1}$ so that

$$-\log m(\mathbf{y}|\mathbf{X}, c) = \frac{k}{2} \log(1+c) + \frac{1}{2\sigma^2} \left(\mathbf{y}^t \mathbf{y} - \frac{c}{1+c} FSS \right), \quad (8)$$

where $FSS = \mathbf{y}^t \mathbf{X}(\mathbf{X}^t \mathbf{X})^{-1} \mathbf{X}^t \mathbf{y}$ is the usual fitted sum of squares corresponding to the OLS estimate $\hat{\boldsymbol{\beta}}$,

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^t \mathbf{X})^{-1} \mathbf{X}^t \mathbf{Y}.$$

This particular choice of distribution is often attributed to Zellner [20] who christened it the g -prior. Because the mixture form reduces to a relatively simple expression, the g -prior has been used extensively to derive Bayesian model selection procedures for the normal linear model. Under this prior, it is not hard to show that the posterior mean of $\boldsymbol{\beta}$ is $\frac{c}{1+c} \hat{\boldsymbol{\beta}}$.

In (8) we have highlighted the dependence of the mixture on the scaling parameter c . George and Foster [7] studied various approaches to setting c , establishing that certain values lead to well-known selection criteria like AIC and BIC.² Ultimately, they propose an empirical Bayes approach, selecting an estimate \hat{c} via maximum likelihood. Hansen and Yu [8] take a similar approach to the hyperparameter c , but motivate it from a coding perspective. We review this approach here. Essentially, each choice of c produces a different mixture distribution and hence a different code. Therefore, to let c depend on the data, both sender and receiver need to agree on which value of c to use. Like c ; that is, c is transmitted first and then once each side knows which code to use, the data are sent. Of course, communicating c in this way will add to the code length, a charge that we make explicit by writing

$$L(\mathbf{y}) = L(\mathbf{y}|c) + L(c) = -\log m(\mathbf{y}|\mathbf{X}, c) + L(c). \quad (9)$$

Following Rissanen [13], the cost $L(c)$ is taken to be $\frac{1}{2} \log n$.³ Minimizing (9) with respect to c gives us

$$\hat{c} = \max \left(\frac{FSS}{k\sigma^2} - 1, 0 \right),$$

and substituting into (8) yields a code length (9) of the form

$$L(\mathbf{y}) = \begin{cases} \frac{\mathbf{y}^t \mathbf{y} - FSS}{2\sigma^2} + \frac{k}{2} [1 + \log(\frac{FSS}{\sigma^2 k})] + \frac{1}{2} \log n & \text{for } FSS > k\sigma^2 \\ \frac{\mathbf{y}^t \mathbf{y}}{2\sigma^2} & \text{otherwise} \end{cases}. \quad (10)$$

²A similar calibration between Bayesian methods and well-known selection criteria can also be found in Smith and Spiegelhalter [17].

³The cost $\frac{1}{2} \log n$ can be motivated as follows: for regular parametric families, an unknown parameter can be estimated at rate $1/\sqrt{n}$. Hence there is no need to code such a parameter with a precision finer than $1/\sqrt{n}$. Coding c with precision $1/\sqrt{n}$ gives a cost to the first order $-\log[1/\sqrt{n}] = \log n/2$.

When the minimizing value of c is zero, the prior on $\boldsymbol{\beta}$ becomes a point mass at zero, effectively producing the “null” model corresponding to all effects being zero. This accounts for the second case in the above expression. We should note that the extra $\frac{1}{2} \log n$ penalty is essential to guarantee consistency of the selection method when the null model is true. The Bayesian criterion of George and Foster [7] is basically the same, but leaves off this extra term.

3.2 Unknown error variance

We now consider the regression model (6) when σ^2 is unknown. While George and Foster [7] advocate estimating σ^2 and then applying the form (10), we prefer to assign a distribution to σ^2 and incorporate it into the mixture. Following Hansen and Yu [8], we employ a conjugate normal-inverse gamma distribution to form the mixture code; that is, $1/\sigma^2$ has a gamma distribution with shape parameter a ; and given σ^2 , $\boldsymbol{\beta}$ is normal with mean zero and variance $\sigma^2 V$. Setting $\tau = \sigma^2$, these densities are given by

$$w(\boldsymbol{\beta}, \tau) \propto \tau^{-\frac{(k+3)}{2}} \exp \left[\frac{-\boldsymbol{\beta}^t V^{-1} \boldsymbol{\beta} + a}{2\tau} \right], \quad (11)$$

where a and V are hyperparameters. Under this class of priors, the mixture distribution (3) has the form

$$\begin{aligned} -\log m(\mathbf{y}|\mathbf{X}, a, V) &= \frac{1}{2} \log |V^{-1} + \mathbf{X}^t \mathbf{X}| + \frac{1}{2} \log |V| - \frac{1}{2} \log a \\ &\quad + \frac{n+1}{2} \log \left(a + \mathbf{y}^t \mathbf{y} - \mathbf{y}^t \mathbf{X} (V^{-1} + \mathbf{X}^t \mathbf{X})^{-1} \mathbf{X}^t \mathbf{y} \right) \end{aligned} \quad (12)$$

where we have ignored terms that do not depend on our particular choice of model. The derivation of $m(\mathbf{y}|\mathbf{X}, a, V)$, the marginal or predictive distribution of \mathbf{y} , is standard and can be found in O’Hagan [10].

Our approach to handling the hyperparameter a will be the same as that in the previous section. Minimizing (12) with respect to a we find that $\hat{a} = (\mathbf{y}^t \mathbf{y} - \mathbf{y}^t \mathbf{X} (V^{-1} + \mathbf{X}^t \mathbf{X})^{-1} \mathbf{X}^t \mathbf{y})/n$ which leaves

$$\begin{aligned} -\log m(\mathbf{y}|\mathbf{X}, \hat{a}, V) &= \frac{1}{2} \log |V^{-1} + \mathbf{X}^t \mathbf{X}| + \frac{1}{2} \log |V| \\ &\quad + \frac{n}{2} \log \left(\mathbf{y}^t \mathbf{y} - \mathbf{y}^t \mathbf{X} (V^{-1} + \mathbf{X}^t \mathbf{X})^{-1} \mathbf{X}^t \mathbf{y} \right). \end{aligned} \quad (13)$$

As was true in the known-variance case, we can achieve a simplification in computing the mixture distribution if we again make Zellner’s choice of $V = c(\mathbf{X}^t \mathbf{X})^{-1}$. This leaves

$$-\log m(\mathbf{y}|\mathbf{X}, \hat{a}, c) = \frac{k}{2} \log(1+c) + \frac{n}{2} \log \left(\mathbf{y}^t \mathbf{y} - \frac{c}{1+c} FSS \right), \quad (14)$$

To settle the hyperparameter c , we again minimize the overall code length to find

$$\hat{c} = \max(F - 1, 0) \quad \text{with} \quad F = \frac{FSS}{kS}, \quad (15)$$

where F is the usual F -ratio for testing the hypothesis that each element of $\boldsymbol{\beta}$ is zero, and $S = RSS/(n - k)$. The truncation at zero in (15) rules out negative values of the prior variance. Rewriting (15), we find that \hat{c} is zero unless $R^2 > k/n$, where R^2 is the usual squared multiple correlation coefficient. When the value of \hat{c} is zero, the prior on $\boldsymbol{\beta}$ becomes a point mass at zero, effectively producing the “null” mixture model⁴ corresponding to zero regression effects. Substituting the optimal value of \hat{c} into (14), we arrive at a final mixture form

$$\text{gMDL} = \begin{cases} \frac{n}{2} \log S + \frac{k}{2} \log F + \log n, & R^2 \geq k/n \\ \frac{n}{2} \log \left(\frac{\mathbf{y}^t \mathbf{y}}{n} \right) + \frac{1}{2} \log n, & \text{otherwise.} \end{cases} \quad (16)$$

Note that we have added the cost to code the hyperparameters a and c , producing an extra $\log n$ and $(1/2) \log n$ in the upper and lower expressions, respectively.

3.3 Comparison

As alluded to in the introduction, two widely used model selection criteria are AIC and BIC. In the case of regression with an unknown variance, they take forms

$$\text{AIC} = \frac{n}{2} \log RSS + k \quad \text{and} \quad \text{BIC} = \frac{n}{2} \log RSS + \frac{k}{2} \log n. \quad (17)$$

Comparing these with (16), we see that the essential difference is in the penalty. Both AIC and BIC have data independent penalties, while gMDL has a data-dependent $\log F/2$ for each additional dimension.

By charging less for each new variable, AIC tends to include more terms. When the underlying model consists of many effects, or more precisely the model is infinite-dimensional, AIC tends to perform better. If we take our figure of merit to be prediction error, then AIC has been shown to be optimal in this setting both through theoretical and simulation studies. When the true, data generating mechanism is finite-dimensional (and is included among the candidates we are comparing), the stronger penalty of BIC tends to perform better. For this kind of problem, we can also judge selection criteria based on consistency (which leads to prediction optimality); that is, whether or not they ultimately select the correct model as the number of samples tends to infinity. BIC has been shown to perform optimally in this setting.

We will now demonstrate that gMDL with its adaptive penalty enjoys the advantages of both AIC and BIC in the regression context. We focus on the simple linear

⁴The null model is a scale mixture of normals, each $N(0, \tau)$ and τ having an inverse-gamma prior.

model because the expressions are easy to work with, although we expect the same kind of result will hold for GLMs. To simplify our analysis, we assume the regressors are ordered as X_{i1}, X_{i2}, \dots . Following Breiman and Freedman [4], we assume that $X_i = (X_{i1}, X_{i2}, \dots, X_{ij}, \dots)$ are Gaussian, zero-mean random vectors and let

$$\sigma_k^2 = \text{var} \left(\sum_{j=k+1}^{\infty} \beta_j X_{ij} \mid X_{i1}, \dots, X_{ik} \right).$$

Then the finite-dimensional model assumption implies that $\sigma_{k_0}^2 = 0$ for some $k_0 > 0$. Using similar arguments as those used to prove Theorem 1.4 in Breiman and Freedman [4] and the fact that $\frac{1}{n} \|y\|^2 = (\sigma^2 + \sigma_0^2)(1 + o_p(1))$, it is straightforward to establish the following two results.

Theorem 1

The quantity F in (15) satisfies $F = \left[\frac{n}{k} \frac{\sigma_0^2 - \sigma_k^2}{\sigma_k^2 + \sigma^2} + 1 \right] (1 + o_p(1))$ where $o_p(1) \rightarrow 0$ in probability uniformly over $0 \leq k \leq n/2$.

Corollary 1

If the model is finite-dimensional and the maximum dimension of the models examined $K = K_n = o(n)$, then gMDL is consistent and is also prediction-optimal.

The above theorem presents an expansion of the data dependent-penalty of gMDL, and the corollary establishes that gMDL enjoys the same optimality as BIC when the model is finite-dimensional. When $\sigma_k^2 > 0$ for all k , the underlying model is infinite-dimensional. In this case, the quantity F/n can be viewed as the average signal to noise ratio for the fitted model. Adjusting the penalty with $(k/2) \log F/n$, gMDL is able to adapt to perform well in terms of prediction in both domains, finite- or infinite-dimensional. The simulation studies in Hansen and Yu (2001) support this adaptivity of gMDL since there gMDL has an overall prediction performance better than AIC or BIC.

In the next section, we will show that the newly derived MDL-based criteria for GLMs are also adaptive.

4 Generalized Linear Models

The characterization of a GLM starts with an exponential family of the form

$$f(y) = \exp \left(\frac{y\psi - b_1(\psi)}{b_2(\phi)} + b_3(y, \phi) \right), \quad y \in \mathcal{Y}, \quad (18)$$

where b_1 , b_2 and b_3 are known functions. We refer to ψ as the canonical parameter for the family. Typically, we take $b_2(\phi) = \phi$, and refer to ϕ as the dispersion parameter. It plays the role of the noise-variance in the ordinary regression setup of the previous

section. The family (18) contains many practically important cases, including the normal, binomial, Poisson, exponential, gamma and inverse Gaussian distributions. With this model, it is not hard to show that if Y has distribution (18),

$$\begin{aligned} E(Y) &= \mu = b_1'(\psi) \\ \text{var}(Y) &= \sigma^2 = b_1''(\psi)b_3(\phi) \end{aligned} \quad (19)$$

As with the normal case above, the GLM framework allows us to study the dependence of a response variable $Y \in \mathcal{Y}$ on a vector of covariates (X_1, \dots, X_K) . Each model class corresponds to some value of the binary vector $\boldsymbol{\gamma} = (\gamma_1, \dots, \gamma_K)$, and we relate the mean μ of Y to a subset of the covariates via the linear predictor

$$\eta = g(\mu), \quad \text{for } \eta = \sum_{j:\gamma_j=1} \beta_j X_j, \quad (20)$$

where g is a one-to-one, continuously differentiable transformation known as the link function. Using (19) and (20) we see that $\eta = g(b_1'(\psi))$.⁵ Again we let $\boldsymbol{\beta}_\gamma = (\beta_j)_{j:\gamma_j=1}$ denote the vector of regression coefficients and $k_\gamma = \sum \gamma_j$ its dimension. The unknown parameters associated with this model are denoted $\boldsymbol{\theta}_\gamma$ and include both $\boldsymbol{\beta}_\gamma$ as well as a possible dispersion effect ϕ . We observe data of the form (Y_i, \mathbf{X}_i) for $i = 1, \dots, n$ where $\mathbf{X}_i = (X_{i1}, \dots, X_{iK})$ and again \mathbf{X}_K is the $n \times K$ full design matrix $[\mathbf{X}_K]_{ij} = X_{ij}$. We let \mathbf{X}_γ refer to a submatrix of \mathbf{X}_K consisting of only those columns j for which $\gamma_j = 1$. Let $f_{\boldsymbol{\theta}_\gamma}(\mathbf{y}|\mathbf{X}_\gamma)$ denote the density for \mathbf{Y} based on model class γ .

As with our treatment of the regression context, maintaining the model index γ needlessly complicates our derivations. From this point on, we again drop it, reminding the reader that terms like \mathcal{M} , \mathbf{X} , k , and $\boldsymbol{\beta}$ all refer to a specific subset of covariates. For all the GLM cases, we will begin with a Laplace approximation to the mixture form which will be exact for the normal linear model. That is, we start with

$$m(\mathbf{y}|\mathbf{X}) \approx (2\pi)^{\frac{k}{2}} | -H^{-1}(\tilde{\boldsymbol{\beta}}) |^{\frac{1}{2}} f_{\tilde{\boldsymbol{\beta}}}(\mathbf{y}|\mathbf{X}) w(\tilde{\boldsymbol{\beta}}) \quad (21)$$

where H is the Hessian of $h(\boldsymbol{\beta}) = \log f_{\boldsymbol{\beta}}(\mathbf{y}|\mathbf{X}) + \log w(\boldsymbol{\beta})$ and $\tilde{\boldsymbol{\beta}}$ is the posterior mode of $\boldsymbol{\beta}$. In working with this form, we will repeatedly make use of the Fisher information matrix

$$I(\boldsymbol{\beta}) = \mathbf{X}^t \mathbf{W}(\boldsymbol{\beta}) \mathbf{X},$$

where \mathbf{W} is a diagonal weight matrix. Note that for GLMs, the observed Fisher information is the same as the Fisher information when we use the canonical parameterization.

Form (21) is still difficult to work with in practice because there is typically no closed-form expression for the posterior mode. We will now consider several criteria that make sensible choices for f and w that lead to computationally tractable criteria.

⁵Taking $b' = g^{-1}$ means that the canonical parameter ψ and the linear predictor η are the same. This choice of g is known as the canonical link function.

4.1 Direct approach

In this section, we will derive a criterion that first appeared in Peterson [11]. As with the regression context, the original motivation for this form was not MDL, but rather an approximation to a full Bayesian approach. Our analysis will follow closely the case of σ^2 known for regression. Let $\boldsymbol{\beta}$ be the MLE of $\boldsymbol{\beta}$, and assume that the prior $w(\boldsymbol{\beta})$ is normal with mean zero and variance-covariance V . Then, we can approximate $\tilde{\boldsymbol{\beta}}$ via a single Newton step

$$\begin{aligned}\tilde{\boldsymbol{\beta}} &\approx \boldsymbol{\beta} - H(\boldsymbol{\beta})^{-1}h'(\boldsymbol{\beta}) \\ &\approx \boldsymbol{\beta} + (I(\boldsymbol{\beta}) + V^{-1})^{-1}V\boldsymbol{\beta}\end{aligned}$$

using the fact that $H(\boldsymbol{\beta}) = -(I(\boldsymbol{\beta}) + V^{-1})$, where $I(\hat{\boldsymbol{\beta}})$ is the Fisher information evaluated at $\hat{\boldsymbol{\beta}}$. We now focus on the case where the prior variance-covariance matrix for $\boldsymbol{\beta}$ is simply $cI(\hat{\boldsymbol{\beta}})^{-1}$. For the normal linear model, this leads us to Zellner's g -prior. Unfortunately, for the other important members of this family, the prior variance-covariance matrix will depend on $\hat{\boldsymbol{\beta}}$. From a strict coding perspective this is hard to accept; it would imply that sender and receiver both know the coefficient $\hat{\boldsymbol{\beta}}$ (or at least $I(\hat{\boldsymbol{\beta}})$). Nonetheless, it is instructive to follow this line of analysis and compare it with the results of the previous section. For $V = cI(\hat{\boldsymbol{\beta}})^{-1}$ we find that the one-step Newton-Raphson iteration gives

$$\tilde{\boldsymbol{\beta}} \approx \frac{c}{1+c}\hat{\boldsymbol{\beta}}$$

which agrees with our regression form of MDL when σ^2 is known.

Continuing with the expression (21), we find that

$$\begin{aligned}\log w(\tilde{\boldsymbol{\beta}}) &\approx \log w\left(\frac{c}{1+c}\hat{\boldsymbol{\beta}}\right) \\ &= -\frac{k}{2}\log 2\pi + \log |cI(\hat{\boldsymbol{\beta}})| - \frac{1}{2}\frac{c}{(1+c)^2}\hat{\boldsymbol{\beta}}^t I(\hat{\boldsymbol{\beta}})\hat{\boldsymbol{\beta}}\end{aligned}\quad (22)$$

and that after a Taylor expansion of $\log f_{\boldsymbol{\beta}}(\mathbf{y}|\mathbf{X})$ around $\hat{\boldsymbol{\beta}}$

$$\begin{aligned}\log f_{\tilde{\boldsymbol{\beta}}}(\mathbf{y}|\mathbf{X}) &\approx \log f_{\hat{\boldsymbol{\beta}}}(\mathbf{y}|\mathbf{X}) - \frac{1}{2}(\tilde{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}})^t I(\hat{\boldsymbol{\beta}})(\tilde{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}) \\ &= \log f_{\hat{\boldsymbol{\beta}}}(\mathbf{y}|\mathbf{X}) - \frac{1}{2}\frac{1}{(1+c)^2}\hat{\boldsymbol{\beta}}^t I(\hat{\boldsymbol{\beta}})\hat{\boldsymbol{\beta}}\end{aligned}\quad (23)$$

Combining (22) and (23) we arrive at the expression

$$\begin{aligned}\log f_{\tilde{\boldsymbol{\beta}}}(\mathbf{y}|\mathbf{X}) + \log w(\tilde{\boldsymbol{\beta}}) &\approx \log f_{\hat{\boldsymbol{\beta}}}(\mathbf{y}|\mathbf{X}) - \frac{1}{2}\frac{1}{1+c}\hat{\boldsymbol{\beta}}^t I(\hat{\boldsymbol{\beta}})\hat{\boldsymbol{\beta}} \\ &\quad - \frac{k}{2}\log 2\pi + \frac{k}{2}\log c + \frac{1}{2}\log |I(\hat{\boldsymbol{\beta}})|.\end{aligned}\quad (24)$$

Finally, collecting terms in (21) we find an expression for the code length given c

$$-\log m(\mathbf{y}|c, \mathbf{X}) \approx \frac{k}{2} \log(1+c) + \frac{1}{2} \frac{1}{1+c} \hat{\boldsymbol{\beta}}^t I(\hat{\boldsymbol{\beta}}) \hat{\boldsymbol{\beta}} - \log f_{\hat{\boldsymbol{\beta}}}(\mathbf{y}|\mathbf{X}).$$

We then eliminate the hyperparameter c using the same minimization approach in (9). This yields

$$\hat{c} = \max \left(\frac{\hat{\boldsymbol{\beta}}^t I(\hat{\boldsymbol{\beta}}) \hat{\boldsymbol{\beta}}}{k} - 1, 0 \right).$$

Substituting this in the mixture form, we find the final MDL criterion

$$L(\mathbf{y}) = \begin{cases} -\log f_{\hat{\boldsymbol{\beta}}}(\mathbf{y}|\mathbf{X}, \hat{\boldsymbol{\beta}}) + \frac{k}{2} \left[1 + \log \left(\frac{\hat{\boldsymbol{\beta}}^t I(\hat{\boldsymbol{\beta}}) \hat{\boldsymbol{\beta}}}{k} \right) \right] + \frac{1}{2} \log n & \text{for } \hat{\boldsymbol{\beta}}^t I(\hat{\boldsymbol{\beta}}) \hat{\boldsymbol{\beta}} > k \\ -\log f_0(\mathbf{y}|\mathbf{X}) & \text{otherwise} \end{cases}$$

The function $f_0(\mathbf{y}|\mathbf{X})$ represents the log-likelihood when all the regression effects are zero. Again, we have added an extra $\frac{1}{2} \log n$ term to the top expression to account for the coding cost of c . This corresponds exactly to the regression context when σ^2 is known.

4.2 Accounting for over-dispersion

In many families, like the Poisson and binomial models, the dispersion parameter is fixed $\phi = 1$. However, in practice it is often the case that the data do not support this, forcing us to consider over-dispersed models. There are a variety of ways to introduce extra variability into the form (18), many of which are primarily meant as computational devices. Efron [6] constructs a family to explicitly account for over-dispersion that admits an analysis for GLMs similar to that for ordinary regression in the σ^2 -unknown case. A related technique was independently derived by West [19].

To understand this form, we have to first rewrite the log-likelihood for a GLM in terms of its mean vector $l(\mathbf{y}|\boldsymbol{\mu})$, where $\boldsymbol{\mu} = (\mu_1, \dots, \mu_n)$. Now, using this notation, without the restriction (20) on the mean, the maximum value of the log-likelihood is simply $l(\mathbf{y}|\mathbf{y})$. We then define the deviance as the difference

$$D(\mathbf{y}|\boldsymbol{\beta}) = 2l(\mathbf{y}|\mathbf{y}) - 2l(\mathbf{y}|\boldsymbol{\mu}),$$

where $\boldsymbol{\beta}$ is the vector of regression coefficients that yield $\boldsymbol{\mu}$ through (20). To incorporate a dispersion parameter, Efron [6] motivates the use of

$$\tau^{-n/2} e^{l(\mathbf{y}|\boldsymbol{\mu})/\tau + (1-1/\tau)l(\mathbf{y}|\mathbf{y})} \tag{25}$$

as an (approximate) likelihood. Technically, this expression should include a normalizing constant $C(\tau, \boldsymbol{\beta})$. Following Efron [6], however, it can be shown that $C(\tau, \boldsymbol{\beta}) = 1 + O(n^{-1})$, and hence can be ignored for reasonable sample sizes. Rewriting (25), we will work with

$$\tau^{-n/2} e^{-\frac{(2l(\mathbf{y}|\mathbf{y}) - 2l(\mathbf{y}|\boldsymbol{\mu}))}{2\tau}} e^{l(\mathbf{y}|\mathbf{y})} = \tau^{-n/2} e^{-\frac{D(\mathbf{y}|\boldsymbol{\beta})}{2\tau}} e^{l(\mathbf{y}|\mathbf{y})}. \quad (26)$$

Then, arguing as we did for the σ^2 unknown case in regression, we use a normal-inverse gamma prior with variance-covariance matrix τV . The joint probability of $\boldsymbol{\beta}$, τ and \mathbf{y} is given by

$$\tau^{-1-(n+1)/2} \frac{e^{-a/2\tau}}{\sqrt{\pi}} \sqrt{\frac{a}{2}} (2\pi\tau)^{-k/2} |V^{-1}|^{1/2} e^{-\frac{\boldsymbol{\beta}^t V^{-1} \boldsymbol{\beta} - D(\mathbf{y}|\boldsymbol{\beta})}{2\tau}} \quad (27)$$

To integrate out $\boldsymbol{\beta}$, we use the Laplace method again which this time yields

$$\frac{\tau^{-1-(n+1)/2}}{\sqrt{\pi}} \sqrt{\frac{a}{2}} |V^{-1}|^{1/2} |V^{-1} + I(\tilde{\boldsymbol{\beta}})|^{-1/2} e^{-\frac{-a - \tilde{\boldsymbol{\beta}}^t V^{-1} \tilde{\boldsymbol{\beta}} - D(\mathbf{y}|\tilde{\boldsymbol{\beta}})}{2\tau}} \quad (28)$$

where $I(\tilde{\boldsymbol{\beta}})$ is the Fisher information matrix evaluated at the posterior mode $\tilde{\boldsymbol{\beta}}$. Integrating with respect to τ then yields

$$\begin{aligned} -\log m(\mathbf{y}|a, V) &= \frac{n+1}{2} \log \left(a + \tilde{\boldsymbol{\beta}}^t V^{-1} \tilde{\boldsymbol{\beta}} + D(\mathbf{y}|\tilde{\boldsymbol{\beta}}) \right) \\ &\quad - \frac{1}{2} \log a + \frac{1}{2} \log |V| + \frac{1}{2} \log |V^{-1} + I(\tilde{\boldsymbol{\beta}})| \end{aligned} \quad (29)$$

Following the prescription in the regression context, we eliminate the hyperparameter a by minimizing the overall code length. In this case we easily find that

$$\begin{aligned} -\log m(\mathbf{y}|\hat{a}, V) &= \frac{n}{2} \log \left(\tilde{\boldsymbol{\beta}}^t V^{-1} \tilde{\boldsymbol{\beta}} + D(\mathbf{y}|\tilde{\boldsymbol{\beta}}) \right) \\ &\quad + \frac{1}{2} \log |V| + \frac{1}{2} \log |V^{-1} + I(\tilde{\boldsymbol{\beta}})| \end{aligned}$$

We have now obtained a usable criterion for model selection. Specifying V , we can compute $\tilde{\boldsymbol{\beta}}$ with simple Newton-Raphson iterations. In our regression analysis, we used Zellner's g -prior for $\boldsymbol{\beta}$ which led to a closed-form selection criterion. The analog in this case is $V = cI^{-1}(\hat{\boldsymbol{\beta}})$. For a GLM, this choice is somewhat unsettling because $I(\hat{\boldsymbol{\beta}})$ is computed at the MLE. If we were to adhere to a strict MDL setting, it would not make sense; from a coding perspective, both sender and receiver would have to know about $\hat{\boldsymbol{\beta}}$, or at least $I(\hat{\boldsymbol{\beta}})$. Recall that for a GLM, the Fisher information matrix takes the form $\mathbf{X}^t \mathbf{W}(\hat{\boldsymbol{\beta}}) \mathbf{X}$ where \mathbf{W} is a diagonal weight matrix. One simple alternative is to take $V = c(\mathbf{X}^t \mathbf{X})^{-1}$, or $V = c\mathbf{1}$, where $\mathbf{1}$ is the identity matrix. In each of these cases, we must either approximate the $\tilde{\boldsymbol{\beta}}$ or iterate to find it. We will consider both kinds of selection criteria.

Following the approximation route, if we choose $V = cI^{-1}(\hat{\boldsymbol{\beta}})$, we get

$$\tilde{\boldsymbol{\beta}} \approx \frac{c}{1+c} \hat{\boldsymbol{\beta}} \quad (30)$$

and

$$\frac{k}{2} \log(1+c) + \frac{n}{2} \log \left(\frac{1}{1+c} \hat{\boldsymbol{\beta}}^t I(\hat{\boldsymbol{\beta}}) \hat{\boldsymbol{\beta}} + D(\mathbf{y}|\hat{\boldsymbol{\beta}}) \right) \quad (31)$$

Here we have substituted in the one-step Newton-Raphson approximation for $\tilde{\boldsymbol{\beta}}$ and have approximated the deviance $D(\mathbf{y}|\tilde{\boldsymbol{\beta}})$ by a Taylor expansion around $\hat{\boldsymbol{\beta}}$ and used a relation from Raftery [12]. Maximizing with respect to c yields

$$\hat{c} = \max(F - 1, 0) \quad (32)$$

where

$$F = \frac{(n-k) \hat{\boldsymbol{\beta}}^t I(\hat{\boldsymbol{\beta}}) \hat{\boldsymbol{\beta}}}{k D(\mathbf{y}|\hat{\boldsymbol{\beta}})}.$$

This then gives the form

$$L(\mathbf{y}) = \begin{cases} \frac{n}{2} \log \frac{D(\mathbf{y}|\hat{\boldsymbol{\beta}})}{n-k} + \frac{k}{2} \log F + \log n & \text{if } F > 1 \\ \frac{n}{2} \log \frac{D(\mathbf{y}|\mathbf{0})}{n} + \frac{1}{2} \log n & \text{otherwise} \end{cases},$$

where $D(\mathbf{y}|\mathbf{0})$ represents the deviance calculated under a model with zero regression effects.

For the other choices of $V^{-1} = c\Sigma$, we do not have a closed-form expression for the maximizing c . Instead, we can perform a search, but this is best done in conjunction with finding $\tilde{\boldsymbol{\beta}}$. It is also possible to use our approximate $\tilde{\boldsymbol{\beta}}$ (30) to derive a simple iteration to find c . In this case, we find

$$c = \frac{k R_c}{n \hat{\boldsymbol{\beta}}^t I(\hat{\boldsymbol{\beta}}) \left(I(\hat{\boldsymbol{\beta}}) + c\Sigma \right)^{-1} \Sigma \left(I(\hat{\boldsymbol{\beta}}) + c\Sigma \right)^{-1} I(\hat{\boldsymbol{\beta}}) \hat{\boldsymbol{\beta}} + R_c \text{trace} \left(c\mathbf{1} + I(\hat{\boldsymbol{\beta}}) \Sigma^{-1} \right)} \quad (33)$$

where

$$R_c = D(\mathbf{y}|\hat{\boldsymbol{\beta}}) + c \hat{\boldsymbol{\beta}}^t \Sigma \hat{\boldsymbol{\beta}} - c^2 \hat{\boldsymbol{\beta}}^t \Sigma (I(\hat{\boldsymbol{\beta}}) + c\Sigma)^{-1} \Sigma \hat{\boldsymbol{\beta}}. \quad (34)$$

Convergence of this algorithm is usually fairly fast, although as we will see, it can depend on the starting values.

Table 1: Classification errors for the different selection criteria.

| Coefficients | ρ | Bayes Rate | Mix ϕ | $\phi = 1$ | BIC | AIC | $\mathbf{1}$ Iter | $\mathbf{X}^t \mathbf{X}$ Iter | $\mathbf{1}$ search | $\mathbf{X}^t \mathbf{X}$ search |
|---------------------|--------|---------------|------------|------------|-------|-------|----------------------|-----------------------------------|------------------------|-------------------------------------|
| 3 2 2 0 0 | 0 | 0.125 | 0.138 | 0.138 | 0.135 | 0.137 | 0.137 | 0.138 | 0.137 | 0.137 |
| | 0.75 | 0.087 | 0.101 | 0.101 | 0.104 | 0.101 | 0.100 | 0.104 | 0.100 | 0.101 |
| 5 1 1 1 1 | 0 | 0.098 | 0.115 | 0.116 | 0.128 | 0.118 | 0.120 | 0.118 | 0.118 | 0.118 |
| | 0.75 | 0.072 | 0.087 | 0.087 | 0.092 | 0.089 | 0.087 | 0.087 | 0.087 | 0.089 |
| 2 2 2 2 2 | 0 | 0.115 | 0.130 | 0.130 | 0.131 | 0.130 | 0.130 | 0.131 | 0.130 | 0.130 |
| | 0.75 | 0.067 | 0.081 | 0.083 | 0.095 | 0.086 | 0.081 | 0.081 | 0.081 | 0.087 |
| 3.0 1.5 1.5 0.5 0.0 | 0 | 0.137 | 0.153 | 0.152 | 0.154 | 0.152 | 0.153 | 0.155 | 0.153 | 0.153 |
| | 0.75 | 0.093 | 0.108 | 0.108 | 0.112 | 0.109 | 0.108 | 0.111 | 0.108 | 0.109 |
| 5 0 0 0 0 | 0 | 0.103 | 0.116 | 0.114 | 0.110 | 0.114 | 0.113 | 0.113 | 0.113 | 0.113 |
| | 0.75 | 0.104 | 0.116 | 0.115 | 0.109 | 0.114 | 0.114 | 0.112 | 0.114 | 0.113 |
| 2 0 0 0 0 | 0 | 0.220 | 0.233 | 0.233 | 0.228 | 0.233 | 0.230 | 0.230 | 0.230 | 0.230 |
| | 0.75 | 0.221 | 0.231 | 0.231 | 0.227 | 0.232 | 0.230 | 0.229 | 0.230 | 0.229 |
| 1.5 1.5 1.5 1.5 0.0 | 0 | 0.163 | 0.177 | 0.177 | 0.177 | 0.177 | 0.177 | 0.180 | 0.177 | 0.177 |
| | 0.75 | 0.101 | 0.119 | 0.120 | 0.129 | 0.121 | 0.118 | 0.124 | 0.118 | 0.122 |
| 8 4 2 1 0 | 0 | 0.060 | 0.074 | 0.074 | 0.077 | 0.074 | 0.075 | 0.074 | 0.074 | 0.074 |
| | 0.75 | 0.040 | 0.057 | 0.058 | 0.060 | 0.058 | 0.057 | 0.057 | 0.057 | 0.058 |

5 Simulations

We have chosen 8 different simulation setups to compare AIC and BIC with the new MDL-based criteria derived in this section. We focus on logistic regression, and consider $K = 5$ potential covariates. We specify two distributions on \mathbf{X} . In the first, each column consists of $n = 100$ observations from a standard normal distribution and the different columns are independent. In the second case, we again use normal covariates, but now we consider a correlation structure of the form

$$\text{cov}(X_i, X_j) = \rho^{|i-j|}, \text{ for } i, j = 1, \dots, 5.$$

In the simulations below, we took $\rho = 0.75$. The data \mathbf{Y} was then generated using the standard logistic GLM using one of 8 different coefficient vectors. All $2^5 = 32$ possible models were fit and compared using the various selection criteria. In Table 1, we present the classification error rate for each procedure: Column 4 corresponds to mixture MDL with a normal-inverse gamma mixing distribution to capture dispersion effects and $V^{-1} = c^{-1}I(\hat{\boldsymbol{\beta}})$ (Section 4.2); Column 5 corresponds to mixture MDL with a fixed dispersion parameter ϕ and hence a normal mixing distribution again with $V^{-1} = c^{-1}I(\hat{\boldsymbol{\beta}})$ (Section 4.1); Columns 6 and 7 are BIC and AIC (17). Columns 8 through 11 also make use of the normal-inverse gamma distribution but with different choices of the variance-covariance matrix V^{-1} ; $c^{-1}\mathbf{1}$ for 8 and 10, and $c^{-1}\mathbf{X}^t\mathbf{X}$ for 9 and 11. Columns 8 and 10 differ only in how we estimate $\tilde{\boldsymbol{\beta}}$ and \hat{c} ; in the first case the iteration (33) is used, while in the second a full search is performed to identify both $\tilde{\boldsymbol{\beta}}$ and the appropriate value of c . The same holds for Columns 9 and 11, but with the different variance-covariance matrix.

In Table 1 we see that most of the selection criteria behave the same, at least in

Table 2: Summarizing the number of times different sized models were selected for a sample of simulation runs given in Table 1.

| Coefficients | Model | | | | $\mathbf{1}$ | | $\mathbf{X}^t \mathbf{X}$ | | $\mathbf{1}$ | | $\mathbf{X}^t \mathbf{X}$ | |
|---------------------------|---------|------------|------------|-----|--------------|------|---------------------------|--------|--------------|--------|---------------------------|-----|
| | Summary | Mix ϕ | $\phi = 1$ | BIC | AIC | Iter | Iter | search | search | search | search | |
| $\beta = (2, 0, 0, 0, 0)$ | 0-1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | 0-2 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | 0-3 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | 0-4 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | * | 1-0 | 131 | 134 | 215 | 121 | 176 | 179 | 183 | 179 | 183 | 179 |
| | 1-1 | 70 | 72 | 31 | 85 | 55 | 53 | 51 | 51 | 51 | 51 | 51 |
| | 1-2 | 37 | 37 | 4 | 40 | 17 | 18 | 15 | 19 | 15 | 19 | 19 |
| | 1-3 | 12 | 7 | 0 | 4 | 2 | 0 | 1 | 1 | 1 | 1 | 1 |
| | 1-4 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| $\beta = (3, 2, 2, 0, 0)$ | 0-1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | 0-2 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | 1-0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | 1-1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | 1-2 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | 2-0 | 0 | 0 | 0 | 0 | 0 | 3 | 0 | 0 | 0 | 0 | 0 |
| | 2-1 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | 2-2 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | * | 3-0 | 111 | 134 | 227 | 173 | 176 | 71 | 184 | 180 | 184 | 180 |
| | 3-1 | 103 | 93 | 20 | 71 | 49 | 31 | 61 | 64 | 61 | 64 | 64 |
| 3-2 | 36 | 23 | 2 | 6 | 25 | 145 | 5 | 6 | 5 | 6 | 6 | |
| $\beta = (5, 1, 1, 1, 1)$ | 1-0 | 0 | 0 | 4 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 |
| | 2-0 | 0 | 1 | 35 | 2 | 4 | 19 | 3 | 3 | 3 | 3 | 3 |
| | 3-0 | 9 | 12 | 64 | 24 | 33 | 13 | 23 | 23 | 23 | 23 | 23 |
| | 4-0 | 56 | 70 | 89 | 94 | 79 | 18 | 86 | 89 | 86 | 89 | 89 |
| | * | 5-0 | 185 | 167 | 58 | 130 | 134 | 199 | 138 | 135 | 199 | 138 |
| | 5-0 | 185 | 167 | 58 | 130 | 134 | 199 | 138 | 135 | 199 | 138 | 135 |

terms of classification error; this 0-1 error is very robust. In Table 2 we illustrate the types of models selected by each scheme. In the first column, we identify the simulations from Table 1. In the second column, we present a model summary of the form $x - y$ where x denotes the number of variables correctly included in the model and y denotes the number of excess variables. So, for the first panel of Table 2, the true model $(2, 0, 0, 0, 0)$ consists of only one effect. The heading “1-0” represents the correct model and is marked with a “*”, while the column “1-1” means that one extra term was included. From this table, we see that the three MDL criteria (Columns 9, 11 and 12) adapt to either AIC or BIC depending on which performs better in all 8 set-ups. Column 10 seemed to have some problems, and we believe this is because the iterations (33) failed to converge properly (possibly due to the approximations used to generate generate the form). Finally, we see that the columns using $I(\hat{\beta})$ can perform poorly (those denoted Mixture ϕ and $\phi = 1$). Recall that we derived these forms even though their reliance on $\hat{\beta}$ violates the basic coding ideas behind MDL.

To consider the cases in more depth, we start with the first panel of Table 2. Here “truth” is a small model, $(2, 0, 0, 0, 0)$, an ideal case for BIC. Clearly, this criterion selects the right model more often than the other procedures. The mixture MDL

procedures that use variance-covariance matrices other than $I(\hat{\beta})$ also perform quite well. In terms of test error, each of the procedures are about the same. Overall, we can recommend the MDL based criteria in terms of their ability to adapt and select concise models.

In the second panel of Table 2, the coefficient vector is $(3, 2, 2, 0, 0)$, a middle-ground case. The $I(\hat{\beta})$ criteria perform rather poorly, as does the $\mathbf{X}^t \mathbf{X}$ case with iterations (33) to find \hat{c} . In the latter case, the poor performance is even reflected in the prediction error. Again, we intend to examine whether it is the approximation that led to (33) that caused the problem, or if it was poor starting values for the iterations.

Finally, in the last panel of Table 2, we consider a “full” model with coefficient vector $(5, 1, 1, 1, 1)$, an ideal situation for AIC. Here we see that BIC fails to capture the correct model form, and the test error is slightly worse as a result. All the MDL criteria outperform even AIC in terms of identifying the correct model, although this does not translate into significant test error improvements.

6 Acknowledgments

B. Yu’s research is partially supported by grants from NSF (FD01-12731) and ARO (DAAD19-01-1-0643). M. Hansen would like to thank John Chambers, Diane Lambert, Daryl Pregibon and Duncan Temple Lang for helpful discussions.

Dedication

This paper is dedicated to Terry Speed, who introduced B. Yu to MDL 15 years ago. By encouraging the highest academic standards, Terry is an inspiration to all of Berkeley’s students.

Mark H. Hansen
Statistics Research, Bell Laboratories, Lucent Technologies
 cocteau@bell-labs.com.

Bin Yu
Department of Statistics, University of California, Berkeley
 binyu@stat.berkeley.edu.

References

- [1] H. Akaike. A new look at the statistical model identification. *IEEE Trans. AC*, 19:716–723, 1974.

- [2] H. Akaike. An objective use of bayesian models. *Ann. Inst. Statist. Math.*, 29:9–20, 1977.
- [3] A. Barron, J. Rissanen, and B. Yu. Minimum description length principle in coding and modeling. *IEEE Trans. Inform. Theory. (Special Commemorative Issue: Information Theory: 1948-1998)*, 44:2743–2760, 1998.
- [4] L. Breiman and D. Freedman. How many variables should be entered in a regression equation? *J. Amer. Statist. Assoc.*, 78:131–136, 1983.
- [5] T. M. Cover and J. A. Thomas. *Elements of Information Theory*. John Wiley & Sons, New York, 1991.
- [6] B. Efron. Double exponential families and their use in generalized linear regression. *J. Amer. Statist. Assoc.*, 81:709–721, 1986.
- [7] E. I. George and D. P. Foster. Calibration and empirical bayes variable selection. *Biometrika*, 87:731–747, 2000.
- [8] M. Hansen and B. Yu. Model selection and minimum description length principle. *J. Amer. Statist. Assoc.*, 96:746–774, 2001.
- [9] R. E. Kass and A. E. Raftery. Bayes factors. *J. Amer. Statist. Assoc.*, 90:773–795, 1995.
- [10] A. O’Hagan. *Kendall’s Advanced Theory of Statistics: Bayesian Inference. Vol 2B*. John Wiley & Sons, New York, 1994.
- [11] J. J. Peterson. A note on some model selection criteria. *Stat. and Prob. Letters*, 4:227–230, 1986.
- [12] A. E. Raftery. Approximate bayes factors and accounting for model uncertainty in generalised linear models. *Biometrika*, 83:251–266, 1996.
- [13] J. Rissanen. *Stochastic complexity and statistical inquiry*. World Scientific, Singapore, 1989.
- [14] J. Rissanen. Fisher information and stochastic complexity. *IEEE Trans. Inform. Theory*, 42:48–54, 1996.
- [15] G. Schwarz. Estimating the dimension of a model. *Ann. Statist.*, 6:4611–464, 1978.
- [16] R. Shibata. An optimal selection of regression variables. *Biometrika*, 68:45–54, 1981.
- [17] A. F. M. Smith and D. J. Spiegelhalter. Bayes factors and choice criteria for linear models. *Journal of the Royal Statistical Society, Series B*, 42:213–220, 1980.
- [18] T. Speed and B. Yu. Model selection and prediction: normal regression. *J. Inst. Statist. Math.*, 45:35–54, 1993.

- [19] M. West. Generalized linear models: scale parameters, outlier accomodation and prior distributions. In J. M. Bernardo, M. H. DeGroot, D. V. Lindley, and A. F. M. Smith, editors, *Bayesian Statistics 2*, pages 531–558, North-Holland, 1985. Elsevier Science Publishers B. V.
- [20] A. Zellner. On assessing prior distributions and bayesian regression analysis with g-prior distributions. In P.K. Goel and A. Zellner, editors, *Bayesian Inference and Decision Techniques: Essays in Honor of Bruno de Finetti*, pages 233–243, Amsterdam, 1986. North-Holland.