

Predictive Model Selection

West and Harrison 1989)) and nonlinear models see

For the model selection problem in general, one can replace (1.2) and the accompanying descriptions of distributions for various quantities by

density to inferential use, adapting the philosophy advocated in Geisser (1991). The imagined replication makes y and \hat{y} comparable; in fact, exchangeable a priori. Moreover, the parameters in the model play a minimal role under replication. It seems clear that good models, among those under consideration, should make predictions close to what has been observed for an identical experiment. The criteria below are defined with this motivation.

For a given model m , consider

$$L_m^2 = E[(\hat{y} - y)'(\hat{y} - y)] ,$$

where the expectation is taken with respect to the PDRE f_m . The measure L_m^2 has the decomposition

$$L_m^2 = \sum_{i=1}^n \{ [E(\hat{y}_i) - y_i]^2 + V(\hat{y}_i) \} ,$$

as a sum of two components, one involving the means of the predictive distribution, and the other involving the variances. Thus a model's performance is measured by a combination of how close its predictions are to the observed data and the variability of the predictions. Good models will have small values of L_m^2 . It is often more convenient to use the measure

$$L_m = \sqrt{L_m^2}$$

since it is a distance on the response axis, measured in the same units as the response variable. We refer to L_m as the *L criterion*.

To define the second criterion, consider

$$M_m^* = f_m(y) .$$

This is the PDRE under model m , evaluated at the observed response y . A good model will have a large value of M_m^* . The ratio of M_m^* 's for two different models is an instance of what Pitkin (1991) calls the *posterior Bayes factor*. To gain, to facilitate interpretation, let

$$M_m = (M_m^*)^{-1/n}$$

which is in the units of the response variable, and small values of it indicate good models. We refer to M_m as the *M criterion*.

The third criterion we introduce for model selection is the Kullback-Leibler (KL) divergence between two predictive densities. Suppose f_1 and f_2 are two densities with respect to Lebesgue measure. Then, the KL divergence between f_1 and f_2 is defined by

$$D(f_1, f_2) = \int \log [f_1(x)/f_2(x)] f_1(x) dx .$$

In general, $D(f_1, f_2) \neq D(f_2, f_1)$, and $D(f, b) \geq 0$ with equality occurring only if $f = b$. The KL divergence has been used in the literature for a wide variety of statistical problems, and in connection with the Bayesian predictive distribution. For example, with its use (itchison 1975) shows that the predictive distribution best approximates the sampling distribution, Johnson and Geisser (1983) detect influential observations in linear regression, and (Culloch 1989) assesses the influence of model assumptions. Bhattacharjee and Dunsmore (1991) use the KL directed divergence to select variables in logistic regression.

For our purposes, suppose m_0 is a fixed model in \mathcal{M} from which we measure other models. In variable selection, for instance, a natural choice for m_0 might be the full model (1.1) with all of the k predictors. Using PDRE's of m_0 and m , we define

$$D(m) = D(m_0, m) + D(m, m)$$

2 Prior Distributions

with $c \geq 0$ quantifying, in multiples of the present experiment, the importance one wishes to attach to the prior guess η_0 . Thus under model m , T_m is a scalar multiple of the Fisher information matrix for $\beta^{(m)}$. Zellner's g-priors (Zellner, 1986) also have this structure for the precision matrix. It has the advantage of leading to analytically tractable and computationally feasible solutions.

Now, we take $\beta^{(m)}|\tau$ to be normally distributed, i.e.,

$$\beta^{(m)}|\tau \sim N_{k_m}(\mu^{(m)}, \tau T_m). \quad (2.5)$$

As a result of focusing on the observables, only a few easily interpreted quantities are needed to specify the prior. In particular, the prediction η_0 is turned into a prior for $\beta^{(m)}|\tau$ for each m in an automated fashion.

Finally, the prior distribution for τ is taken to be a gamma distribution with parameters $(\delta_0/2, \gamma_0/2)$, i.e., with density

$$\pi(\tau) d\tau \propto \tau^{\delta_0/2-1} \exp\{-\gamma_0\tau/2\} d\tau. \quad (2.6)$$

For a fixed model m , (2.5) and (2.6) result in the conjugate normal-gamma prior.

With this prior and the likelihood implied by (1.2) for each m , a straightforward derivation yields

$$Y \sim S_n \left(\eta_0, \eta_m, s_m^2 (I + (1-\gamma)P_m) \right), \quad (2.7)$$

where $\gamma = c/(1+c)$, $\eta_m = P_m(\gamma\eta_0 + (1-\gamma)y)$, $s_m^2 = (n+\delta_0)^{-1}(q_m + \gamma p_m + \gamma_0)$, $q_m = y'(I - P_m)y$, and $p_m = (y - \eta_0)'P_m(y - \eta_0)$. The PDRE in (2.2) for noninformative priors can be obtained from (2.7) by formally setting $\gamma = 0$, $\delta_0 = -k_m$, and $\gamma_0 = 0$. Moreover if X_m has rank $r_m < k_m$, replace k_m by r_m in (2.7) above. For brevity, the relevant expressions are given only for the case of conjugate priors in the remainder of this article.

The L criterion under model m is now given by

$$L_m = \{ (1 + \lambda_m)q_m + \gamma(\gamma + \lambda_m)p_m + \lambda_m\gamma_0 \}^{1/2}, \quad (2.8)$$

where $\lambda_m = \frac{n+(1-\gamma)k_m}{n+\delta_0-2}$. We see that L_m^2 above is a linear function of q_m and p_m . The quantity q_m is the squared length of the projection of the data onto the error space of model m , i.e., the error sum of squares for model m . The quantity p_m represents a penalty for a bad prior guess at Y . It is the squared length of the projection of the "guessing error" onto the model's column space. Under reference priors, (2.8) reduces to $L_m = (2n-1)(n-k_m-2)^{-1}q_m^{1/2}$. In this case, L_m is similar to the root mean square criterion.

To calculate the calibration number S , one can sample from the marginal distribution

$$Y \sim S_n \left(\eta_0, \eta_{m^*}, \gamma_0^{-1} (I + \gamma^{-1}(1-\gamma)P_{m^*}) \right),$$

and calculate L_{m^*} with each sample. Here m^* is the model that minimizes L_m .) The standard deviations of these values provides a Monte Carlo approximation to S . If one is using the reference priors in (2.1), however, it is well known that the marginal distribution of Y is improper. In this case, one could sample from the conditional distribution $Y|\tau \sim N(0, \tau(I - P_{m^*}))$ with τ replaced by $\tilde{\tau}$, the mode of the posterior distribution of τ using m^* . The standard deviation of the resulting samples of L_{m^*} can be viewed as an approximation to $\sigma = [V(L_{m^*}|\tau = \tilde{\tau})]^{1/2}$. For large n one can obtain the analytic approximation

$$\approx \tilde{\tau}^{-1}$$

where

$$= \frac{2}{\quad}$$

number S . Under the improper reference priors, however, one is not guaranteed such automatic protection and hence must be careful to not include in \mathcal{M} an model that can field a

3.2 Transformation Selection

In linear regression, transformations of the predictor variables can often lead to more accurate predictions and a model that better fits the data. Box and Cox (1964) discuss transformations with an emphasis on transforming the response variable. They also mention briefly a possible Bayesian approach. It appears, however, that the literature on Bayesian transformation methods is sparse at best.

Here, we show how two of the predictive criteria can be used to select a specific member of a suitably chosen parametric transformation family. The AIC criterion as defined in this article is not applicable to this problem. Consider equation (1.2) where a single model $m \in \mathcal{M}$

an intercept and the Box-Cox transformation on x . In addition, to denote the dependence of the criteria on α , we write $M_m \equiv M(\alpha)$ and $L_m \equiv L(\alpha)$. Under the noninformative prior (2.1), $M(\alpha)$ and $L(\alpha)$ are equivalent and we get the minimizer $\hat{\alpha} = -1.325$, with $L(\hat{\alpha}) = 0.160$ and $M(\hat{\alpha}) = 0.099$. Results for the coefficient of determination (R^2), residual sums of squares (RSS), and the criterion functions are given for three regression models in Table 2.

Table 2 - Comparison of Models, Vapor Pressure Data

Model	R^2	RSS	$L(\alpha)$
-------	-------	-----	-------------

4 Discussion

The minimizations of the criterion functions L and M for the transformation problems were carried out numerically since analytic methods are not readily available. The computations were greatly facilitated by LISP-STAT (Tierne, 1990), which made it possible to carry out the calculations with relatively few lines of code. The functions NEWTON-X and NEL-ED-X were used with good success. For the examples of this paper, the calculations proceeded quite fast on a SUN SPARC station. Starting values of $\alpha = (1, \dots, 1)'$ worked well. Other starting values were also used.

An important issue in an model selection procedure is that of model assumptions. It is well known that violations of the same can result in the addition or omission of variables in a variable selection procedure. AIC and BIC, for instance, are not robust to outliers or influential points. The criteria proposed in this article likely suffer from the same problems. Simultaneously checking and selecting models is difficult, and there are no definitive solutions to this problem. However, the proposed methods are robust to outliers and influential points.

given by

$$M_m = \pi^{1/2} \left(\frac{\frac{n+\delta_0}{2}}{n + \delta_0/2} (2 - \gamma)^{k_m/2} \right)^{1/n} \frac{1}{m} \left(1 + \frac{b_m}{m} \right)^{1 + \frac{\delta_0}{2n}},$$

where $m = q_m + \gamma p_m + \gamma_0$ and $b_m = q_m + \frac{\gamma^2}{2-\gamma} p_m$. Again, both m and b_m are linear combinations of the residual sum of squares and the “guessing error”.

An exact expression for the criterion is not available since the necessary integral is not tractable. However, for large n , we can approximate the distribution in (2.1) by a $N(\eta_m, \left(\frac{n+\delta_0}{n+\delta_0-2}\right)^{-1} s_m^{-2} (I + (1-\gamma)P_m)^{-1})$ distribution. Taking m_0 to be the full model, define

$$v = \frac{(n + m)(n + m_0 - 2)}{(n + m_0)(n + m - 2)},$$

where $m = m_0 = 0$ for the normal-gamma priors, and $m = -k_m$, $m_0 = -k_{m_0}$ under noninformative priors.

- [2] Pitchison, J., and Dunsmore, I. R. (1975), *Statistical Prediction Analysis*, New York : Cambridge University Press.
- [3] Pitkin, E. (1991), "Posterior Bayes Factors," (with discussion), *Journal of the Royal Statistical Society*, Ser. B, 53, 111-142.
- [4] Akaike, H. (1973), "Information Theory and an Extension of the Maximum Likelihood Principle," *International Symposium on Information Theory*, eds. B. N. Petrov and F. Csaki, pp.267-281. Budapest: Akademia Kiado.
- [5] Bates, D. M., and Watts, D. G. (1988) *Nonlinear Regression Analysis and its Applications*, New York : John Wiley & Sons.
- [6] Bernardo, J. M. (1985), Comment on "Outliers and Influential Observations in Linear Models", (with discussion), in *Bayesian Statistics 2*, eds. Bernardo, J. M., DeGroot, M. H., Lindle, D. V., and Smith, A. F. M., Amsterdam : North-Holland, p.492.
- [7] Bhattacharjee, S. K., and Dunsmore, I. R. (1991) "The Influence of Variables in Logistic Regression", *Biometrika*, 78, 851-856.
- [8] Box, G. E. P., and Cox, D. R. (1964), "The Analysis of Transformations" (with discussion), *Journal of the Royal Statistical Society*, Ser. B, 26, 211-252.
- [9] Box, G. E. P., and Jenkins, G. M. (1976). *Time Series Analysis : Forecasting and Control*, (2nd Ed.), San Francisco: Holden-Day.
- [10] Box, G. E. P., and Kanemasu, H. (1973), "Posterior Probabilities of Candidate Models in Model Discrimination," Technical Report 322, University of Wisconsin.
- [11] Box, G.E. P., and Jenkins, D. R. (1986), "Dispersion Effects From Fractional Designs," *Technometrics*, 28, 19-24.
- [12] Box, G. E. P., and Tiao, G. C. (1973), *Bayesian Inference in Statistical Analysis*. Reading, MA : Addison-Wesley.
- [13] Box, G. E. P., and Tidwell, P. W. (1962), "Transformation of the Independent Variables," *Technometrics*, 4, 531-550.
- [14] Carroll, R. J., and Ruppert, D. (1988), *Transformation and Weighting in Regression*, London : Chapman and Hall.
- [15] Christensen, R. (1990), *Log-Linear Models*, New York : Springer-Verlag.
- [16] Claiborn, K., Geisser, S., and Jennings, D. E. (1986), "A Comparison of Several Model Selection Procedures," in *Studies in Bayesian Econometrics and Statistics*, eds. P. K. Goel and A. Zellner, New York : Elsevier.

[1] Cook, R.

- [32] 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.
- [33] Pettit, L. I., and Smith, J. F. (1985), "Outliers and Influential Observations in Linear Models", (with discussion), in *Bayesian Statistics 2*, eds. Bernardo, J. M., DeGroot, M. H., Lindle, D. V., and Smith, J. F., Amsterdam : North-Holland, p.492.
- [34] San Martini, M., and Spezzaferri, F. (1984), "Bayesian Predictive Model Selection Criterion," *Journal of the Royal Statistical Society, Ser. B*, 46, 296-303.
- [35] San Martini, M., and Spezzaferri, F. (1986) "Selection of Variables in Multiple Regression for Prediction and Control," *Statistica*, 118-121.
- [36] Schwarz, G. (1978), "Estimating the Dimension of a Model", *Annals of Statistics*, 6, 461-464.
- [37] Smith, J. F., and Spiegelhalter, D. J. (1980), "Bayesian Factors and Choice Criteria for Linear Models," *Journal of the Royal Statistical Society, Ser. B.*, 42, 213-220.
- [38] Spiegelhalter, D. J., and Smith, J. F. (1982), "Bayesian Factors for Linear and Log-linear Models with Vague Prior Information", *Journal of the Royal Statistical Society, Ser. B.*, 44, 31-38.
- [39] Taguchi, G., and Wu, Y. (1980), *Introduction to Off-Line Quality Control*, Nagoya, Japan : Central Japan Quality Control Association.
- [40] Tierney, L. (1990), *Lisp-Stat : An Object-Oriented Environment for Statistical Computing and Dynamic Graphics*, New York: John Wiley.
- [41] Tuke, J. W. (1974), *Exploratory Data Analysis*, Reading, MA : Addison-Wesley. (19410TDF) (Nagoya)