

THE CHOICE OF VARIABLES FOR PREDICTION IN CURVILINEAR MULTIPLE REGRESSION

By R. J. BROOKS

University College, London

A Bayesian formulation of the problem of analysing data from a curvilinear regression of y on x_1, x_2, \dots, x_r in order to predict a future value of y is considered. The problem is to obtain a criterion to decide which is the best subset of x_1, x_2, \dots, x_r to perform this prediction. Under very strict assumptions the criterion obtained is shown to use the same statistic as the orthodox (least squares) approach.

1. Introduction. The type of problem to be discussed in this paper can be illustrated by an example from Johnson and Leone ((1964) page 313). A rubber manufacturing company performed an experiment in which interest was focused on the relationship of the tear strength of a particular type of rubber with three variables (i) percentage of component A in rubber, (ii) percentage of component A in resin and (iii) percentage modifier. If these three variables are denoted by x_1, x_2 and x_3 respectively and the tear strength by y , then a quadratic regression model of y on x_1, x_2, x_3 was investigated. It was found that the values of the coefficients of x_3, x_3^2, x_1x_3 and x_2x_3 were all not significantly different from zero (using a 5% significance level). This could be taken to suggest that x_3 need not be included in the model and if we wanted to predict a future value of y we need only use (x_1, x_2) . In this paper, this type of problem is investigated, but from a Bayesian viewpoint. If the model of interest is a curvilinear regression of y on x_1, x_2, \dots, x_r and the object in mind is to predict a future value of y , we have to decide which is the best subset of x_1, x_2, \dots, x_r to use for this prediction.

We use the decision-theoretic approach of Lindley (1968), who discussed the same problem for linear regression of y on x_1, x_2, \dots, x_r . Consequently, the first four assumptions stated below are generalizations of those made by Lindley.

Let \mathcal{E} denote the results of an experiment to investigate the regression plus \mathcal{E}_0 , which denotes any knowledge prior to the experiment.

ASSUMPTION 1. θ and \mathbf{x} are two random variables in \mathcal{R}^u and \mathcal{R}^r respectively (Euclidean spaces of u and r dimensions) which are independent, given \mathcal{E} .

ASSUMPTION 2. y is a random variable in \mathcal{R}^1 with density $p(y|\theta, \mathbf{x}, \mathcal{E}_0)$ whose form does not depend on the results of the experiment, such that

$$(1) \quad E(y|\theta, \mathbf{x}, \mathcal{E}_0) = \theta^T \phi(\mathbf{x})$$

and

$$\text{Var}(y|\theta, \mathbf{x}, \mathcal{E}_0) = \sigma^2,$$

a known constant. Here, if $\mathbf{x} = (x_1, x_2, \dots, x_r)^T$, ξ is a known integer (> 0),

Received March 1971; revised May 1972.

$u = (\xi^{+r})$ and $A = \{\text{all sets } (\alpha_1, \alpha_2, \dots, \alpha_r): \alpha_i \text{ is integer } \geq 0, 0 \leq \sum_{i=1}^r \alpha_i \leq \xi\}$, $\phi(\mathbf{x})$ is a $u \times 1$ vector with typical element $x_1^{\alpha_1} x_2^{\alpha_2} \dots x_r^{\alpha_r}$ where $(\alpha_1, \alpha_2, \dots, \alpha_r) \in A$. The corresponding element of θ is defined to be $\theta_{\alpha_1 \alpha_2 \dots \alpha_r}$.

Equation (1) can then be written as

$$E(y | \theta, \mathbf{x}, \mathcal{E}_0) = \sum_A \theta_{\alpha_1 \alpha_2 \dots \alpha_r} x_1^{\alpha_1} x_2^{\alpha_2} \dots x_r^{\alpha_r}.$$

An example would be the two variable quadratic regression (i.e. $r = \xi = 2$),

$$(2) \quad E(y | \theta, \mathbf{x}, \mathcal{E}_0) = \theta_{00} + \theta_{10} x_1 + \theta_{01} x_2 + \theta_{20} x_1^2 + \theta_{11} x_1 x_2 + \theta_{02} x_2^2$$

where $\theta = (\theta_{00}, \theta_{10}, \theta_{01}, \theta_{20}, \theta_{11}, \theta_{02})^T$ and $\phi(\mathbf{x}) = (1, x_1, x_2, x_1^2, x_1 x_2, x_2^2)^T$.

In these two assumptions, \mathbf{x} and y refer to the future values of the independent and dependent variables, respectively.

Lindley (1968) discussed the special case $\xi = 1$, which in his notation can be written for the u variable linear regression as $E(y | \theta, \mathbf{x}, \mathcal{E}_0) = \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_u x_u$, where we put $x_1 = 1$ to allow for the constant term. If we replace the elements of the $u \times 1$ vector $(x_1, x_2, \dots, x_u)^T$ by the elements of the $u \times 1$ vector $\phi(\mathbf{x})$, then the initial part of Lindley's analysis, up to his expression (11), holds for the regression model we are considering. Lindley then assumes (Assumption 5) that $p(\mathbf{x} | \mathcal{E})$, which with our model will be $p(\phi(\mathbf{x}) | \mathcal{E})$, has linear regressions. Clearly this does not hold for $\xi > 1$, so the remainder of his prediction analysis is not applicable to curvilinear regression.

Let I denote a subset of the integers $1, 2, \dots, r$ containing s members ($0 \leq s \leq r$), and J its complement. Define \mathbf{x}_I as the vector with elements x_i , $i \in I$, and similarly define \mathbf{x}_J . Define $\phi_I(\mathbf{x}_I)$ as the $v \times 1$ vector of all the elements of $\phi(\mathbf{x})$ which do not contain any elements or function of elements belonging to \mathbf{x}_J . Let the remaining elements of $\phi(\mathbf{x})$ be contained in $\phi_J(\mathbf{x})$. Define vectors θ_I and θ_J corresponding to $\phi_I(\mathbf{x}_I)$ and $\phi_J(\mathbf{x})$ respectively.

In the above example when $r = \xi = 2$, these definitions would mean that when $I = (1)$, $\phi_I(\mathbf{x}_I) = (1, x_1, x_1^2)^T$ and $\phi_J(\mathbf{x}) = (x_2, x_1 x_2, x_2^2)^T$ so that $\theta_I = (\theta_{00}, \theta_{10}, \theta_{20})^T$ and $\theta_J = (\theta_{01}, \theta_{11}, \theta_{02})^T$.

ASSUMPTION 3. The decision space consists of elements $(I, f(\cdot))$, where $f(\cdot)$ is a function from \mathcal{R}^s to \mathcal{R}^1 .

ASSUMPTION 4. The loss in predicting y using the subset of variables \mathbf{x}_I giving the prediction $f(\mathbf{x}_I)$ is

$$[y - f(\mathbf{x}_I)]^2 + c_I,$$

where $c_I (\geq 0)$ is the cost of observing the variables in \mathbf{x}_I per cost of unit error in prediction.

The loss function is assumed to consist of additive terms arising from the cost of the error $|y - f(\mathbf{x}_I)|$ in prediction and the cost of observing those variables in \mathbf{x}_I . The cost of error in prediction is assumed to be proportional to $[y - f(\mathbf{x}_I)]^2$.

Using results obtained by Lindley ((1968) (8) and (11)), we have:

LEMMA 1. Under Assumptions 1 to 4, (a) the best set of independent variables to use to predict y are those x_i with $i \in I$, where I is chosen to be that subset of $(1, 2, \dots, r)$ which minimizes

$$(3) \quad E(\boldsymbol{\theta})^T \mathbf{V}[\boldsymbol{\phi}_J(\mathbf{x})] E(\boldsymbol{\theta}) + c_I$$

where

$$\mathbf{V}[\boldsymbol{\phi}_J(\mathbf{x})] = E([\boldsymbol{\phi}(\mathbf{x}) - E\{\boldsymbol{\phi}(\mathbf{x}) | \mathbf{x}_I\}][\boldsymbol{\phi}(\mathbf{x}) - E\{\boldsymbol{\phi}(\mathbf{x}) | \mathbf{x}_I\}]^T);$$

(b) the optimum predictor of y is $E(\boldsymbol{\theta})^T E[\boldsymbol{\phi}(\mathbf{x}) | \mathbf{x}_I]$.

In (a) and (b), all expectations and dispersion matrices are conditional upon \mathcal{E} .

2. The regression experiment. To choose the best set of variables, it will clearly be useful to simplify (3). To do this we need to know the density $p(\mathbf{x} | \mathcal{E})$. Consequently, we now discuss the nature of the experiment we conduct leading to the information denoted by \mathcal{E} .

The situations described by Assumptions 5a and 5b both lead to $p(\mathbf{x} | \mathcal{E})$ belonging to the same family of densities.

ASSUMPTION 5a. In the experiment, n independent observations y_1, y_2, \dots, y_n are obtained at values $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ of the independent variables, where y_k has density $p(y_k | \boldsymbol{\theta}, \mathbf{x}_k, \mathcal{E}_0)$ given in Assumption 2. Here $\mathbf{x}_k = (x_{k1}, x_{k2}, \dots, x_{kr})^T$, ($k = 1, 2, \dots, n$). The \mathbf{x}_k are independent random variables from a common multinormal distribution with unknown mean vector $\boldsymbol{\mu}$ and unknown dispersion matrix $\boldsymbol{\Sigma}$. Furthermore, the future value $\mathbf{x} = (x_1, x_2, \dots, x_r)^T$ (see Assumption 1) is another independent value from the same distribution. Prior to the experiment, the parameters $(\boldsymbol{\mu}, \boldsymbol{\Sigma}^{-1})$ are independent of $\boldsymbol{\theta}$ and have a Normal-Wishart distribution (as defined by Ando and Kaufman ((1965) Section 1.3)).

ASSUMPTION 5b. This is as Assumption 5a, but the \mathbf{x}_k are selected. (\mathbf{x} still has the distribution given in Assumption 5a.)

The experiment discussed in detail by Lindley ((1968) Section 4) is a special case of that described by Assumption 5a, but with a diffuse prior for $(\boldsymbol{\mu}, \boldsymbol{\Sigma}^{-1})$, i.e. $p(\boldsymbol{\mu}, \boldsymbol{\Sigma}^{-1} | \mathcal{E}_0) \propto |\boldsymbol{\Sigma}|$.

LEMMA 2. Under Assumption 5a or 5b, $p(\mathbf{x} | \mathcal{E})$ is a Student density on \mathcal{R}^r (the parameters in the two cases are different).

(Note: throughout this paper, the definition of the Student density is that given by Raiffa and Schlaifer (1961) page 256.)

PROOF. Consider the situation of Assumption 5b. This is a designed experiment which does not give us any information about the distribution of \mathbf{x} ; thus $p(\mathbf{x} | \mathcal{E}) = p(\mathbf{x} | \mathcal{E}_0)$. Now

$$(4) \quad p(\mathbf{x} | \mathcal{E}_0) = \int \int p(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}^{-1}, \mathcal{E}_0) p(\boldsymbol{\mu}, \boldsymbol{\Sigma}^{-1} | \mathcal{E}_0) d\boldsymbol{\mu} d\boldsymbol{\Sigma}^{-1},$$

where $p(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}^{-1}, \mathcal{E}_0)$ is the density of a multivariate normal distribution and $p(\boldsymbol{\mu}, \boldsymbol{\Sigma}^{-1} | \mathcal{E}_0)$ is the density of a Normal-Wishart distribution. Ando and

Kaufman ((1965) Section 2.4) show that $p(\mathbf{x} | \mathcal{E}_0)$ is the density of a multivariate Student distribution.

In the situation of Assumption 5a (a random experiment), we have

$$(5) \quad p(\mathbf{x} | \mathcal{E}) = \int \int p(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}^{-1}, \mathcal{E}_0) p(\boldsymbol{\mu}, \boldsymbol{\Sigma}^{-1} | \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n, \mathcal{E}_0) d\boldsymbol{\mu} d\boldsymbol{\Sigma}^{-1}$$

where $p(\boldsymbol{\mu}, \boldsymbol{\Sigma}^{-1} | \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n, \mathcal{E}_0)$ is the density of the posterior distribution of $(\boldsymbol{\mu}, \boldsymbol{\Sigma}^{-1})$ which is Normal-Wishart (Ando and Kaufman (1965) Section 1.3). By comparing (5) with (4), $p(\mathbf{x} | \mathcal{E})$ will be the density of a multivariate Student distribution.

The parameters of the densities (4) and (5) can be found by comparison with the results in Ando and Kaufman (1965) and should be such that the moments necessary for calculating $V[\boldsymbol{\phi}_J(\mathbf{x}) | \mathcal{E}]$ exist.

3. Choice of variables for prediction. As a result of the discussion of Section 2, we assume $p(\mathbf{x} | \mathcal{E})$ is a Student density on \mathcal{R}^r with ν degrees of freedom. We make use of the following result.

LEMMA 3. *If $p(\mathbf{x} | \mathcal{E})$ is a Student density on \mathcal{R}^r , then for fixed I ,*

$$(6) \quad E[\boldsymbol{\phi}_J(\mathbf{x}) | \mathbf{x}_I, \mathcal{E}] = \mathbf{B}_{JI} \boldsymbol{\phi}_I(\mathbf{x}_I)$$

where \mathbf{B}_{JI} is a $(u - v) \times v$ matrix with elements depending on \mathcal{E} .

The proof of Lemma 3 is found in the appendix.

By analogy with Lindley ((1968) Section 3), we find that

$$(7) \quad V[\boldsymbol{\phi}_J(\mathbf{x}) | \mathcal{E}] = \mathbf{M}_{JJ} - \mathbf{M}_{JI} \mathbf{M}_{II}^{-1} \mathbf{M}_{IJ}$$

where $\mathbf{M}_{IJ} = E[\boldsymbol{\phi}_I(\mathbf{x}_I) \boldsymbol{\phi}_J(\mathbf{x})^T | \mathcal{E}]$ and $\mathbf{M}_{II}, \mathbf{M}_{JI}, \mathbf{M}_{JJ}$ are defined similarly. Here, in $V[\boldsymbol{\phi}_J(\mathbf{x}) | \mathcal{E}]$, the zero rows and columns corresponding to the elements of $\boldsymbol{\phi}_I(\mathbf{x}_I)$ have been omitted. Thus, from (3), we have:

LEMMA 4. *Under Assumptions 1 to 5, to predict y the optimum set I is chosen to minimize*

$$(8) \quad E(\boldsymbol{\theta}_J | \mathcal{E})^T (\mathbf{M}_{JJ} - \mathbf{M}_{JI} \mathbf{M}_{II}^{-1} \mathbf{M}_{IJ}) E(\boldsymbol{\theta}_J | \mathcal{E}) + c_I.$$

It is evident that we require moments of the form

$$(9) \quad E(x_1^{\eta_1} x_2^{\eta_2} \dots x_r^{\eta_r} | \mathcal{E})$$

for integers $\eta_i \geq 0$ such that $0 \leq \sum_{i=1}^r \eta_i \leq 2\xi$. These moments can be expressed in terms of the elements of $E(\mathbf{x} | \mathcal{E})$ and $V(\mathbf{x} | \mathcal{E})$. (As an example, see Lemma 5.) In the case of the random experiment these elements depend on $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n, \mathcal{E}_0$, but in the case of the designed experiment they only depend on \mathcal{E}_0 .

4. An approximate method. In general, it is tedious to evaluate (9) in terms of the elements of $E(\mathbf{x} | \mathcal{E})$ and $V(\mathbf{x} | \mathcal{E})$. However, in the case of the random experiment, we can find an approximation to (8) which is analogous to a result by Lindley ((1968) Section 5, (31)) for multiple linear regression. In addition to Assumptions 1 to 5a we need the following:

ASSUMPTION 6. The prior distribution of θ is uniform over \mathcal{R}^r . Also, in Assumption 5a, the prior distribution of (μ, Σ^{-1}) is the limiting case in which $p(\mu, \Sigma^{-1} | \mathcal{E}_0) \propto |\Sigma|$.

In this situation it is well known that $E(\theta | \mathcal{E})$ is the least squares estimator of θ .

Furthermore, Lindley ((1968) Section 4) gives $\nu = n - 1$,

$$(10) \quad E(\mathbf{x} | \mathcal{E}) = n^{-1} \sum_{k=1}^n \mathbf{x}_k$$

and

$$(11) \quad \mathbf{V}(\mathbf{x} | \mathcal{E}) = n^{-1} \mathbf{W} \mathbf{W}^T$$

to order n^{-1} , where the (k, i) th element of \mathbf{W} is $x_{ki} - \bar{x}_i$, and $\bar{x}_i = n^{-1} \sum_{k=1}^n x_{ki}$.

Let $\Phi(\mathbf{X})$ be the matrix whose k th row is $\phi(\mathbf{x}_k)$, where \mathbf{X} denotes the matrix with elements $x_{ki} (k = 1, 2, \dots, n; i = 1, 2, \dots, r)$. Partition \mathbf{X} into two matrices \mathbf{X}_I and \mathbf{X}_J corresponding to the variables in \mathbf{x}_I and \mathbf{x}_J respectively. Similarly, partition $\Phi(\mathbf{X})$ into $\Phi_I(\mathbf{X}_I)$ and $\Phi_J(\mathbf{X}_J)$ corresponding to $\phi_I(\mathbf{x}_I)$ and $\phi_J(\mathbf{x}_J)$, respectively. For clarity, we write $\mathbf{Z} = \Phi(\mathbf{X})$, $\mathbf{Z}_I = \Phi_I(\mathbf{X}_I)$ and $\mathbf{Z}_J = \Phi_J(\mathbf{X}_J)$.

For n sufficiently large, from (10) and (11) we see that to order n^{-1} , $E(x_i | \mathcal{E})$ and $E(x_i x_j | \mathcal{E})$ equal

$$n^{-1} \sum_{k=1}^n x_{ki} \quad \text{and} \quad n^{-1} \sum_{k=1}^n x_{ki} x_{kj},$$

respectively. Also, $p(\mathbf{x} | \mathcal{E})$ is approximately a multinormal density, and if we approximate $E(x_1^{\eta_1} x_2^{\eta_2} \dots x_r^{\eta_r} | \mathcal{E})$ (see (9)) by

$$n^{-1} \sum_{k=1}^n x_{k1}^{\eta_1} x_{k2}^{\eta_2} \dots x_{kr}^{\eta_r},$$

then we replace \mathbf{M}_{IJ} by $n^{-1} \mathbf{Z}_I^T \mathbf{Z}_J$, and replace \mathbf{M}_{II} , \mathbf{M}_{JI} , \mathbf{M}_{JJ} similarly. Expression (7) is then approximately

$$(12) \quad n^{-1} [\mathbf{Z}_J^T \mathbf{Z}_J - (\mathbf{Z}_I^T \mathbf{Z}_J)^T (\mathbf{Z}_I^T \mathbf{Z}_I)^{-1} (\mathbf{Z}_I^T \mathbf{Z}_J)] .$$

In an orthodox (least squares) approach to this analysis we would calculate $R(I)$, the reduction in the residual sum of squares due to extending the model which includes only those independent variables in \mathbf{x}_I , (i.e. $E(\mathbf{y} | \theta_I, \mathbf{X}_I, \mathcal{E}_0) = \mathbf{Z}_I \theta_I$) to include all the independent variables (i.e. $E(\mathbf{y} | \theta, \mathbf{X}, \mathcal{E}_0) = \mathbf{Z} \theta$). Kendall and Stuart ((1966) (35.122)) give

$$(13) \quad R(I) = E(\theta_J | \mathcal{E})^T [\mathbf{Z}_J^T \mathbf{Z}_J - (\mathbf{Z}_I^T \mathbf{Z}_J)^T (\mathbf{Z}_I^T \mathbf{Z}_I)^{-1} (\mathbf{Z}_I^T \mathbf{Z}_J)] E(\theta_J | \mathcal{E})$$

where $E(\theta_J | \mathcal{E})$ is the least squares estimator of θ_J .

Using the approximation to (7) as given by (12), we have:

THEOREM 1. Under Assumptions 1 to 5a and 6, for large n , to predict y the optimum set I is chosen to minimize

$$(14) \quad n^{-1} R(I) + c_I .$$

This result is comparable to that obtained by Lindley ((1968) (31)) for multiple linear regression. To determine the best subset of independent variables to use to predict a future value of y , for the curvilinear regression defined in (1), both the orthodox and Bayesian approaches calculate the statistic $R(I)$ for sufficiently large sample size. However, this has been done with three restrictions:

- (a) we have the full model as defined in (1),
- (b) \mathbf{x} has a multinormal distribution,
- (c) \mathbf{x} and $(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n)$ have a common distribution.

If we have, for example, the model

$$E(y | \boldsymbol{\theta}, \mathbf{x}, \mathcal{E}_0) = \theta_{00} + \theta_{10}x_1 + \theta_{01}x_2 + \theta_{02}x_2^2$$

then we can still use (8) by considering the model given by (2) and assuming prior to the experiment that $\theta_{20} = \theta_{11} = 0$. However since the k th row of \mathbf{Z} is then $(1, x_{k1}, x_{k2}, x_{k1}^2, x_{k1}x_{k2}, x_{k2}^2)$, this means that $E(\boldsymbol{\theta}_j | \mathcal{E})$ is not given as in (13) and the approximation (14) will no longer hold, so that the Bayesian method does not, in this case, use the same statistic $R(I)$ as the orthodox method. Hence, we impose the restriction (a).

An easier example to demonstrate the restriction (a) would be to consider simple linear regression through the origin, where we write $E(y | \boldsymbol{\theta}, x_1, \mathcal{E}_0) = \theta_0 + \theta_1x_1$ assuming prior to the experiment $\theta_0 = 0$ and θ_1 is uniformly distributed over \mathcal{R}^1 .

The restriction (b) can most easily be seen by considering simple linear regression, i.e. $\xi = r = 1$. Suppose that instead of (b), we assume x_1 has a uniform distribution over $(0, \lambda)$, where given \mathcal{E}_0 , λ is independent of $\boldsymbol{\theta}$ and $p(\lambda | \mathcal{E}_0) \propto \lambda^{-(\alpha+1)}(\lambda > 0; \alpha \geq 0$ is known). (This is a conjugate prior density for λ .) According to (3), we have to find the minimum of $c_{(1)}$ and

$$E(\theta_1 | \mathcal{E})^2 \text{Var}(x_1 | \mathcal{E}).$$

With (c), $\text{Var}(x_1 | \mathcal{E})$ depends on the observations $x_{11}, x_{21}, \dots, x_{n1}$ only through $M = \max_{k=1,2,\dots,n} \{x_{k1}\}$. This can easily be seen, because

$$p(x_1 | \mathcal{E}) = \int p(x_1 | \lambda, \mathcal{E}_0)p(\lambda | x_{11}, x_{21}, \dots, x_{n1}, \mathcal{E}_0) d\lambda$$

where $p(\lambda | x_{11}, x_{21}, \dots, x_{n1}, \mathcal{E}_0)$ is the posterior density of λ , which depends on $x_{11}, x_{21}, \dots, x_{n1}$ only through M (see Raiffa and Schlaifer (1961) page 54). Thus it is obvious that when $I = \phi$, $E(\theta_1 | \mathcal{E})^2 \text{Var}(x_1 | \mathcal{E})$ is not approximated by $n^{-1}R(I)$, where $R(I)$ is obtained from (13), and so (14) does not hold for all I .

If instead of (c) we assume the experiment is designed (as in Assumption 5b), then $V[\boldsymbol{\phi}_j(\mathbf{x}) | \mathcal{E}]$ does not depend on $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$. Consequently, (7) cannot be approximated by (12) and it is obvious that we do not use $R(I)$.

A further point of interest that should be mentioned while we are comparing the approach of this paper and orthodox methods, is that, for example, in the model given by (2), we are not investigating whether we should use a quadratic

model rather than a linear model. By our method, if we include x_1 , we would include x_1^2 , as we assume the extra cost of using x_1^2 is a computing cost, which may be negligible. This procedure may seem unreasonable if say ξ is greater than 2. However, many problems of this nature just require the use of a response surface for which $\xi = 2$, and it is towards this class of problems that the procedure is particularly directed.

5. An example. Consider the data from a regression experiment described in Williams ((1959) pages 42–45). The regression model is quadratic in the independent variable, thus $r = 1$, $\xi = 2$ and $E(y | \theta, x_1, \mathcal{E}_0) = \theta_0 + \theta_1 x_1 + \theta_2 x_1^2$.

With the assumptions of Section 4, we examine the criteria for selecting the best subset of independent variables, both in the original form given by (3) or equivalently (8), and in the approximate form given by (14). The choice of variables is simply between observing x_1 or not observing x_1 in order to predict a future value of y , i.e. $I = (1)$ or ϕ respectively.

When $I = \phi$, (3), or (8), will equal S where

$$S = e_1^2 \text{Var}(x_1 | \mathcal{E}) + 2e_1 e_2 \text{Cov}(x_1, x_1^2 | \mathcal{E}) + e_2^2 \text{Var}(x_1^2 | \mathcal{E})$$

and $e_i = E(\theta_i | \mathcal{E})$, $i = 1, 2$. The quantities $\mu = E(x_1 | \mathcal{E})$ and $\tau^2 = \text{Var}(x_1 | \mathcal{E})$ are obtained from (10) and (11) respectively, and in order to calculate S , we express $E(x_1^3 | \mathcal{E})$ and $E(x_1^4 | \mathcal{E})$ in terms of μ and τ^2 using the following lemma.

LEMMA 5. *If x has a Student distribution with ν degrees of freedom, mean μ and variance τ^2 , then*

$$\begin{aligned} E(x^3) &= 3\tau^2\mu + \mu^3, \\ E(x^4) &= 3k\tau^4 + 6\tau^2\mu^2 + \mu^4 \end{aligned}$$

where $k = (\nu - 2)/(\nu - 4)$. Here we assume $\nu > 4$.

The proof of Lemma 5 is given in the appendix.

In this example, $\nu = n - 1$, and consequently

$$S = (3k - 1)e_2^2\tau^4 + (e_1 + 2\mu e_2)^2\tau^2$$

where $k = (n - 3)/(n - 5)$.

Alternatively, if we proceed as in Section 4, S is approximately $n^{-1}R(\phi)$.

When $I = (1)$, (3) and (14) both yield c , the cost of observing x per cost of unit error in prediction. Consequently, we choose $I = (1)$ if and only if $c < S$ or $c < n^{-1}R(\phi)$, depending on whether we approximate (3) or not.

The values of $\sum x_{k1}$, $\sum x_{k1}^2$, e_1 , e_2 and $R(\phi)$ are given in Williams as 1646.4, 81747, 9.4743, 0.50863 and 21621500, respectively. The number of observations n is 36; hence $S = 580110$ and $n^{-1}R(\phi) = 600597$. The error in using $n^{-1}R(\phi)$ instead of S is just below $3\frac{1}{2}\%$ and therefore negligible.

Comparing the Bayesian and orthodox methods used on the above example, we see that a large value of S has corresponded to a highly significant value of the statistic used for testing the inclusion of (θ_1, θ_2) in the regression (see

Williams, Table 3.10). Intuitively, this is what we would expect, since a large S means that the cost of observing x would have to be very much larger than the cost of unit error in prediction before we would find it too expensive to observe x_1 in order to predict y .

Note that if we observe x_1 , the optimal predictor of y is, by Lemma 2, $e_0 + e_1x_1 + e_2x_1^2$; otherwise it is $e_0 + e_1E(x_1|\mathcal{E}) + e_2E(x_1^2|\mathcal{E})$.

APPENDIX

PROOF OF LEMMA 3. That is, if $p(\mathbf{x}|\mathcal{E})$ is a non-degenerate Student density on \mathcal{R}^r with ν degrees of freedom, then for fixed I ,

$$(6) \quad E[\phi_J(\mathbf{x}) | \mathbf{x}_I, \mathcal{E}] = \mathbf{B}_{JI}\phi_I(\mathbf{x}_I)$$

where \mathbf{B}_{JI} is a matrix of dimension $(u - v) \times v$, whose elements depend on \mathcal{E} and are regarded here as constants. We omit the \mathcal{E} for clarity.

Suppose without loss of generality, that \mathbf{x} is partitioned so that $\mathbf{x}_I = (x_1, x_2, \dots, x_s)^T$ and $\mathbf{x}_J = (x_{s+1}, x_{s+2}, \dots, x_r)^T$. Then $\phi_I(\mathbf{x}_I)$ contains all variables of the form $x_1^{\alpha_1}x_2^{\alpha_2} \dots x_s^{\alpha_s}$, and $\phi_J(\mathbf{x})$ all variables of the form $x_1^{\alpha_1}x_2^{\alpha_2} \dots x_r^{\alpha_r}$ where at least one of $\alpha_{s+1}, \alpha_{s+2}, \dots, \alpha_r$ is nonzero. An individual element of $E[\phi_J(\mathbf{x}) | \mathbf{x}_I]$ can then be written as $x_1^{\alpha_1}x_2^{\alpha_2} \dots x_s^{\alpha_s}E[\phi(\mathbf{x}_J) | \mathbf{x}_I]$ where $\phi(\mathbf{x}_J) = x_{s+1}^{\alpha_{s+1}}x_{s+2}^{\alpha_{s+2}} \dots x_r^{\alpha_r}$ and at least one of $\alpha_{s+1}, \alpha_{s+2}, \dots, \alpha_r$ is nonzero.

The definition of the non-degenerate Student density function on \mathcal{R}^r is given by Raiffa and Schlaifer ((1961) page 256) as

$$(15) \quad p_s(\mathbf{x} | \boldsymbol{\mu}, \mathbf{H}, \nu) \equiv \int_0^\infty p_N(\mathbf{x} | \boldsymbol{\mu}, h\mathbf{H})p_r(h | 1, \nu) dh,$$

where \mathbf{H} is a positive semi-definite matrix, $p_N(\mathbf{x} | \boldsymbol{\mu}, h\mathbf{H})$ is the density of the multinormal distribution with mean $\boldsymbol{\mu}$ and dispersion matrix $h^{-1}\mathbf{H}^{-1}$, $p_r(h | 1, \nu) = \exp(-\frac{1}{2}\nu h)(\frac{1}{2}\nu h)^{\frac{1}{2}\nu-1}\frac{1}{2}\nu/\Gamma(\frac{1}{2}\nu)$ is the Gamma-2 density ($h \geq 0$), as defined by Raiffa and Schlaifer ((1961) page 226).

Corresponding to the partition of \mathbf{x} into \mathbf{x}_I and \mathbf{x}_J , suppose that $\boldsymbol{\mu}$ and \mathbf{H} are partitioned as

$$\begin{bmatrix} \boldsymbol{\mu}_I \\ \boldsymbol{\mu}_J \end{bmatrix}$$

and

$$\begin{bmatrix} \mathbf{H}_{II} & \mathbf{H}_{IJ} \\ \mathbf{H}_{JI} & \mathbf{H}_{JJ} \end{bmatrix}$$

respectively.

Using (15), it follows that

$$(16) \quad E_s[\phi(\mathbf{x}_J) | \mathbf{x}_I, \boldsymbol{\mu}, \mathbf{H}, \nu] = \frac{\int_0^\infty E_N[\phi(\mathbf{x}_J) | \mathbf{x}_I, \boldsymbol{\mu}, h\mathbf{H}]p_N(\mathbf{x}_I | \boldsymbol{\mu}_I, h\mathbf{H}_{II})p_r(h | 1, \nu) dh}{p_s(\mathbf{x}_I | \boldsymbol{\mu}_I, \mathbf{H}_{II}, \nu)}$$

where

$$p_s(\mathbf{x}_I | \boldsymbol{\mu}_I, \mathbf{H}_{II}, \nu) = \int_0^\infty p_N(\mathbf{x}_I | \boldsymbol{\mu}_I, h\mathbf{H}_{II})p_r(h | 1, \nu) dh.$$

Now

$$(17) \quad E_N[\psi(\mathbf{x}_J) | \mathbf{x}_I, \boldsymbol{\mu}, h\mathbf{H}] = \int x_{s+1}^{\alpha_{s+1}} x_{s+2}^{\alpha_{s+2}} \cdots x_r^{\alpha_r} p_N(\mathbf{x}_J | \mathbf{x}_I, \boldsymbol{\mu}, h\mathbf{H}) d\mathbf{x}_J$$

where $p_N(\mathbf{x}_J | \mathbf{x}_I, \boldsymbol{\mu}, h\mathbf{H}) = (2\pi)^{-\frac{1}{2}(r-s)} |h\mathbf{H}_{JJ}|^{\frac{1}{2}} \exp[-\frac{1}{2}h(\mathbf{x}_J - \mathbf{m}_J)^T \mathbf{H}_{JJ}(\mathbf{x}_J - \mathbf{m}_J)]$,

$$(18) \quad \begin{aligned} \mathbf{m}_J &= \boldsymbol{\mu}_J - (h\mathbf{H}_{JJ})^{-1} h\mathbf{H}_{JI}(\mathbf{x}_I - \boldsymbol{\mu}_I) \\ &= \boldsymbol{\mu}_J - \mathbf{H}_{JJ}^{-1} \mathbf{H}_{JI}(\mathbf{x}_I - \boldsymbol{\mu}_I), \end{aligned}$$

(see Raiffa and Schlaifer (1961) page 250).

We make the transformation $h^{\frac{1}{2}} \mathbf{U}_{JJ}(\mathbf{x}_J - \mathbf{m}_J) = \boldsymbol{\omega}_J$ where $\mathbf{U}_{JJ}^T \mathbf{I} \mathbf{U}_{JJ} = \mathbf{H}_{JJ}$ and \mathbf{I} is the unit matrix of order $r - s$; then

$$h(\mathbf{x}_J - \mathbf{m}_J)^T \mathbf{H}_{JJ}(\mathbf{x}_J - \mathbf{m}_J) = \boldsymbol{\omega}_J^T \mathbf{I} \boldsymbol{\omega}_J,$$

and the Jacobian of the transformation is $|h\mathbf{H}_{JJ}|^{-\frac{1}{2}}$. Each element of \mathbf{x}_J will be of the form

$$x_j = h^{-\frac{1}{2}} \sum_{l=s+1}^r u^{(jl)} \omega_l + m_j, \quad j = s + 1, \dots, r,$$

where each $u^{(jl)}$ is a constant and the m_j are the corresponding elements of \mathbf{m}_J ; the right-hand side of (17) then becomes

$$\begin{aligned} (2\pi)^{-\frac{1}{2}(r-s)} \int \prod_{j=s+1}^r [h^{-\frac{1}{2}} \sum_{l=s+1}^r u^{(jl)} \omega_l + m_j]^{\alpha_j} \exp(-\frac{1}{2} \boldsymbol{\omega}_J^T \mathbf{I} \boldsymbol{\omega}_J) d\boldsymbol{\omega}_J \\ = \sum_B K_{\beta_{s+1} \dots \beta_r} m_{s+1}^{\beta_{s+1}} \dots m_r^{\beta_r} h^{-\frac{1}{2} \sum_{j=s+1}^r (\alpha_j - \beta_j)} \end{aligned}$$

where $B = \{\text{all sets } (\beta_{s+1}, \dots, \beta_r) : \beta_j \text{ is integer } \geq 0, 0 \leq \beta_j \leq \alpha_j\}$ and each $K_{\beta_{s+1} \dots \beta_r}$ is a constant. The numerator of the right-hand side of (16) is now

$$\sum_B K_{\beta_{s+1} \dots \beta_r} m_{s+1}^{\beta_{s+1}} \dots m_r^{\beta_r} \int_0^\infty h^{-\frac{1}{2} \sum (\alpha_j - \beta_j)} p_N(\mathbf{x}_I | \boldsymbol{\mu}_I, h\mathbf{H}_{II}) p_r(h | 1, \nu) dh$$

where $p_N(\mathbf{x}_I | \boldsymbol{\mu}_I, h\mathbf{H}_{II}) = (2\pi)^{-\frac{1}{2}s} |h\mathbf{H}_{II}|^{\frac{1}{2}} \exp[-(h/2)(\mathbf{x}_I - \boldsymbol{\mu}_I)^T \mathbf{H}_{II}(\mathbf{x}_I - \boldsymbol{\mu}_I)]$, and hence the integral in this expression is proportional to

$$[\nu + (\mathbf{x}_I - \boldsymbol{\mu}_I)^T \mathbf{H}_{II}(\mathbf{x}_I - \boldsymbol{\mu}_I)]^{-\frac{1}{2}(\nu+s) + \frac{1}{2} \sum (\alpha_j - \beta_j)}.$$

Now

$$p_s(\mathbf{x}_I | \boldsymbol{\mu}_I, \mathbf{H}_{II}, \nu) \propto [\nu + (\mathbf{x}_I - \boldsymbol{\mu}_I)^T \mathbf{H}_{II}(\mathbf{x}_I - \boldsymbol{\mu}_I)]^{-\frac{1}{2}(\nu+s)},$$

and for the right-hand side of (16) we then obtain

$$\sum_B K'_{\beta_{s+1} \dots \beta_r} m_{s+1}^{\beta_{s+1}} \dots m_r^{\beta_r} [\nu + (\mathbf{x}_I - \boldsymbol{\mu}_I)^T \mathbf{H}_{II}(\mathbf{x}_I - \boldsymbol{\mu}_I)]^{\frac{1}{2} \sum (\alpha_j - \beta_j)}$$

where $K'_{\beta_{s+1} \dots \beta_r} = L_{\beta_{s+1} \dots \beta_r} K_{\beta_{s+1} \dots \beta_r}$, each $L_{\beta_{s+1} \dots \beta_r}$ being a constant.

We can write

$$(\mathbf{x}_I - \boldsymbol{\mu}_I)^T \mathbf{H}_{II}(\mathbf{x}_I - \boldsymbol{\mu}_I) = \sum_\Gamma M_{\gamma_1 \dots \gamma_s} x_1^{\gamma_1} \dots x_s^{\gamma_s}$$

where $\Gamma = \{\text{all sets } (\gamma_1, \dots, \gamma_s) : \gamma_i \text{ is integer } \geq 0, 0 \leq \sum_{i=1}^s \gamma_i \leq 2\}$; furthermore, from (18), $m_j = \sum_{i=1}^s Q_{ji} x_i + R_j$ for constants R_j and Q_{ji} ($j = s + 1, \dots, r$; $i = 1, 2, \dots, s$), and thus

$$\prod_{j=s+1}^r m_j^{\beta_j} = \sum_\Delta Q'_{\delta_1 \dots \delta_s} x_1^{\delta_1} \dots x_s^{\delta_s}$$

where $\Delta = \{\text{all sets } (\delta_1, \dots, \delta_s) : \delta_i \text{ integer } \geq 0, 0 \leq \sum_{i=1}^s \delta_i \leq \sum_{j=s+1}^r \beta_j\}$ and each $Q'_{\delta_1 \dots \delta_s}$ is a constant.

From (16) we finally obtain that an element of $E_s[\phi_j(\mathbf{x}) | \mathbf{x}_T, \boldsymbol{\mu}, \mathbf{H}, \nu]$ has the form

$$(19) \quad x_1^{\alpha_1} \dots x_s^{\alpha_s} \sum_B K'_{\beta_{s+1} \dots \beta_r} [\nu + \sum_{\Gamma} M_{r_1 \dots r_s} x_1^{r_1} \dots x_s^{r_s}]^{\sum_{j=s+1}^r (\alpha_j - \beta_j)} \times \sum_{\Delta} Q'_{\delta_1 \dots \delta_s} x_1^{\alpha_1} \dots x_s^{\delta_s}.$$

The maximum power of any $x_i, i = 1, 2, \dots, s$, is

$$\alpha_i + 2 \cdot \frac{1}{2} \cdot \sum_{j=s+1}^r (\alpha_j - \beta_j) + \sum_{j=s+1}^r \beta_j = \alpha_i + \sum_{j=s+1}^r \alpha_j \leq \xi,$$

and the maximum sum of the powers of x_i is

$$\sum_{i=1}^s \alpha_i + 2 \cdot \frac{1}{2} \cdot \sum_{j=s+1}^r (\alpha_j - \beta_j) + \sum_{j=s+1}^r \beta_j = \sum_{i=1}^r \alpha_i \leq \xi.$$

Consequently (19) can be written as

$$\sum_Z b_{\zeta_1 \dots \zeta_s} x_1^{\zeta_1} \dots x_s^{\zeta_s},$$

where $Z = \{\text{all sets } (\zeta_1, \dots, \zeta_s) : \zeta_i \text{ is integer } \geq 0, 0 \leq \sum_{i=1}^s \zeta_i \leq \xi\}$ and each $b_{\zeta_1 \dots \zeta_s}$ is a constant. Written in matrix notation, this result means that (6) is true, since all possible forms of $x_1^{\zeta_1} \dots x_s^{\zeta_s}$ will be contained in the elements of $\phi_T(\mathbf{x}_T)$.

PROOF OF LEMMA 5. Using the definition of the non-degenerate Student density function on \mathcal{R}^1 as given by (15), we can show

$$(20) \quad E_s(x^\eta | \mu, H, \nu) = \int_0^\infty E_N(x^\eta | \mu, hH) p_T(h | 1, \nu) dh.$$

Since, from Raiffa and Schlaifer ((1961), page 257, (8.29)), we have

$$E_s(x | \mu, H, \nu) = \mu \quad \text{and} \quad \text{Var}_s(x | \mu, H, \nu) = \nu H^{-1} / (\nu - 2),$$

then μ is defined as in the statement of the lemma, and $H = \nu / (\nu - 2) \tau^2$.

It is easier to make the transformation $\omega = x - \mu$ and calculate $E_s(\omega^\eta | \mu, H, \nu)$ for $\eta = 3$ and 4 using (20), and then transform back to x to obtain $E_s(x^\eta | \mu, H, \nu)$. Since $E_N(\omega^3 | 0, hH) = 0$ and $E_N(\omega^4 | 0, hH) = 3(hH)^{-2} = 3(\nu - 2)^2 \tau^4 / h^2 \nu^2$, then from (20) we obtain $E_s(\omega^3 | 0, H, \nu) = 0$ and $E_s(\omega^4 | 0, H, \nu) = 3k\tau^4$, where $k = (\nu - 2) / (\nu - 4)$.

Transforming back to x , we obtain

$$E_s(x^3 | \mu, H, \nu) = 3\mu\tau^2 + \mu^3$$

and

$$E_s(x^4 | \mu, H, \nu) = 3k\tau^4 + 6\tau^2\mu^2 + \mu^4,$$

these being the quantities $E(x^3)$ and $E(x^4)$ required in Lemma 5.

Acknowledgments. The author is grateful to Professor D.V. Lindley and a referee for their comments. He is also indebted to the Science Research Council under whose financial support the work was undertaken.

REFERENCES

ANDO, A. and KAUFMAN, G. M. (1965). Bayesian analysis of the independent multinomial process—neither mean nor precision known. *J. Amer. Statist. Assoc.* **60** 347-358.

- JOHNSON, N. L. and LEONE, F. C. (1964). *Statistics and Experimental Design in Engineering and the Physical Sciences*, 2. Wiley, New York.
- KENDALL, M. G. and STUART, A. (1966). *The Advanced Theory of Statistics. Design and Analysis, and Time Series*. 3. Griffin, London.
- LINDLEY, D. V. (1968). The choice of variables in multiple regression. *J. Roy. Statist. Soc. Ser. B.* 30 31-53.
- RAIFFA, H. and SCHLAIFER, R. (1961). *Applied Statistical Decision Theory*. Harvard Graduate School of Business Administration, Boston.
- WILLIAMS, E. J. (1959). *Regression Analysis*. Wiley, New York.

UNIVERSITY COLLEGE LONDON
GOWER STREET
LONDON WC 1, ENGLAND