

Bayesian Wavelet Regression on Curves With Application to a Spectroscopic Calibration Problem

P. J. BROWN, T. FEARN, and M. VANNUCCI

Motivated by calibration problems in near-infrared (NIR) spectroscopy, we consider the linear regression setting in which the many predictor variables arise from sampling an essentially continuous curve at equally spaced points and there may be multiple predictands. We tackle this regression problem by calculating the wavelet transforms of the discretized curves, then applying a Bayesian variable selection method using mixture priors to the multivariate regression of predictands on wavelet coefficients. For prediction purposes, we average over a set of likely models. Applied to a particular problem in NIR spectroscopy, this approach was able to find subsets of the wavelet coefficients with overall better predictive performance than the more usual approaches. In the application, the available predictors are measurements of the NIR reflectance spectrum of biscuit dough pieces at 256 equally spaced wavelengths. The aim is to predict the composition (i.e., the fat, flour, sugar, and water content) of the dough pieces using the spectral variables. Thus we have a multivariate regression of four predictands on 256 predictors with quite high intercorrelation among the predictors. A training set of 39 samples is available to fit this regression. Applying a wavelet transform replaces the 256 measurements on each spectrum with 256 wavelet coefficients that carry the same information. The variable selection method could use subsets of these coefficients that gave good predictions for all four compositional variables on a separate test set of samples. Selecting in the wavelet domain rather than from the original spectral variables is appealing in this application, because a single wavelet coefficient can carry information from a band of wavelengths in the original spectrum. This band can be narrow or wide, depending on the scale of the wavelet selected.

KEY WORDS: Markov chain Monte Carlo; Mixture prior; Model averaging; Multivariate regression; Near-infrared spectroscopy; Variable selection.

1. INTRODUCTION

This article presents a new way of tackling linear regression problems in which the predictor variables arise from sampling an essentially continuous curve at equally spaced points. The work was motivated by calibration problems in near-infrared (NIR) spectroscopy, of which the following example is typical.

1.1 Near-Infrared Spectroscopy of Biscuit Doughs

Quantitative NIR spectroscopy is used to analyze such diverse materials as food and drink, pharmaceutical products, and petrochemicals. The NIR spectrum of a sample of, say, wheat flour is a continuous curve measured by modern scanning instruments at hundreds of equally spaced wavelengths. The information contained in this curve can be used to predict the chemical composition of the sample. The problem lies in extracting the relevant information from possibly thousands of overlapping peaks. Osborne, Fearn, and Hindle (1993) described applications in food analysis and reviewed some of the standard approaches to the calibration problem.

The example studied in detail here arises from an experiment done to test the feasibility of NIR spectroscopy to

measure the composition of biscuit dough pieces (formed but unbaked biscuits), for possible on-line implementation. (For a full description of the experiment, see Osborne, Fearn, Miller, and Douglas 1984.) Briefly, two similar sample sets were made up, with the standard recipe varied to provide a large range for each of the four constituents under investigation: fat, sucrose, dry flour, and water. The calculated percentages of these four ingredients represent the $q = 4$ responses. There were $n = 39$ samples in the calibration or training set, with sample 23 excluded from the original 40 as an outlier, and a further $m = 39$ in the separate prediction or validation set, again after one outlier was excluded. Thus \mathbf{Y} and \mathbf{Y}_f , the matrices of compositional data for the training and validation sets, are both of dimension 39×4 .

An NIR reflectance spectrum is available for each dough piece. The original spectral data consist of 700 points measured from 1100 to 2498 nanometers (nm) in steps of 2 nm. For our analyses using wavelets, we have chosen to reduce the number of spectral points to save computational time. The first 140 and last 49 wavelengths, which were thought to contain little useful information, were removed, leaving a wavelength range from 1380 nm to 2400 nm, over which we took every other point, thus increasing the gap to 4 nm and reducing the number of points to $p = 256$. The matrices \mathbf{X} and \mathbf{X}_f of spectral data are then 39×256 . Samples of three centered spectra are given on the left side of Figure 1.

The aim is to derive an equation that will predict the response values \mathbf{Y} from the spectral data \mathbf{X} for future samples where \mathbf{Y} is unknown but \mathbf{X} can be measured cheaply and rapidly.

P. J. Brown is Pfizer Professor of Medical Statistics, Institute of Mathematics and Statistics, University of Kent, Canterbury, Kent, CT2 7NF, U.K. (E-mail: philip.j.brown@ukc.ac.uk). T. Fearn is Professor of Applied Statistics, Department of Statistical Science, University College London, WC1E 6BT, U.K. (E-mail: tom@stats.ucl.ac.uk). M. Vannucci is Assistant Professor of Statistics, Department of Statistics, Texas A&M University, College Station, TX 77843 (E-mail: mvannucci@stat.tamu.edu). This work was supported by the U.K. Engineering and Physical Sciences Research Council under the Stochastic Modelling in Science and Technology Initiative, grant GK/K73343. M. Vannucci also acknowledges support from the Texas Higher Advanced Research Program, grant 010366-0075 from Texas A&M International Research Travel Assistant Program and from the National Science Foundation CAREER award DMS-0093208. The spectroscopic calibration problem was provided by the Flour Milling and Baking Research Association. The authors thank the associate editor and a referee, as well as Mike West, Duke University, and Adrian Raftery, University of Washington, for suggestions that helped improve the article.

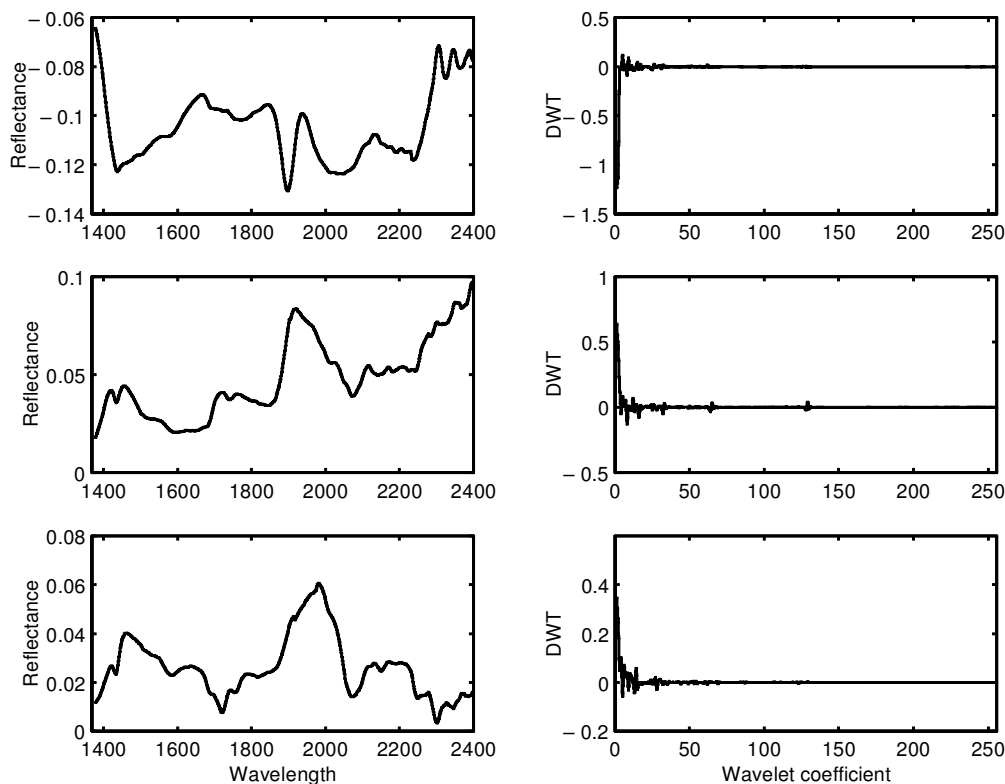


Figure 1. Original Spectra (left column) and Wavelet Transforms (right column).

1.2 Standard Analyses

The most commonly used approaches to this calibration problem regress \mathbf{Y} on \mathbf{X} , with the linear form

$$\mathbf{Y} = \mathbf{XB} + \mathbf{E}$$

being justified either by appeals to the Beer–Lambert law (Osborne et al. 1993) or on the grounds that it works in practice. In Section 6 we also investigate logistic transformations of the responses, showing that overall their impact on the prediction performance of the model is not beneficial.

The problem is not straightforward, because there are many more predictor variables (256) than training samples (39) in our example. The most commonly used methods for overcoming this difficulty fall into two broad classes: variable selection and factor-based methods. When scanning NIR instruments first appeared, the standard approach was to select (typically using a stepwise procedure) predictors at a small number of wavelengths and use multiple linear regression with this subset (Hrushka 1987). Later, this approach was largely superseded by methods that reduce the p spectral variables to scores on a much smaller number of factors and then regress on these scores. Two variants—principal components regression (PCR; Cowe and McNicol 1985) and partial least squares regression (PLS; Wold, Martens, and Wold 1983)—are now widely used, with equal effectiveness, as the standard approaches. The increasing power of computers has triggered renewed research interest in wavelength selection, now using computer-intensive search methods.

1.3 Selecting Wavelet Coefficients

The approach that we investigate here involves selecting variables, but we select from derived variables. The idea is to transform each spectrum into a set of wavelet coefficients, the whole of which would suffice to reconstruct the spectrum, and select good predictors from among these. There are good reasons for thinking that this approach might have advantages over selecting from the original variables.

In previous work (Brown, Fearn, and Vannucci 1999; Brown, Vannucci, and Fearn 1998a,b) we explored Bayesian approaches to the problem of selecting predictor variables in this multivariate regression context. We apply this methodology to wavelet selection here.

We know that wavelets can be used successfully for compression of curves like the spectra in our example, in the sense that the curves can be accurately reconstructed from a fraction of the full set of wavelet coefficients (Trygg and Wold 1998; Walczak and Massart 1997). Furthermore, the wavelet decomposition of the curve is a local one, so that if the information relevant to our prediction problem is contained in a particular part or parts of the curve, as it typically is, then this information will be carried by a very small number of wavelet coefficients. Thus we may expect selection to work. The ability of wavelets to model the curve at different levels of resolution gives us the option of selecting from our curve at a range of bandwidths. In some situations it may be advantageous to select a sharp band, as we do when we select one of the original variables; in other situations a broad band, averaging over many adjacent points, may be preferable. Selecting from the wavelet coefficients gives us both of these options automatically and in a very computationally efficient framework.

Not surprisingly, Fourier analysis has also been successfully used for data compression and denoising of NIR spectral data (McClure, Hamid, Giesbrecht, and Weeks 1984). For this purpose, there is probably very little to choose between the Fourier and wavelet approaches. When it comes to selecting small numbers of coefficients for prediction, however, the local nature of wavelets makes them the obvious choice.

It is worth emphasizing that what we are doing here is not what is commonly described as wavelet regression. We are not fitting a smooth function to a single noisy curve by using either thresholding or shrinkage of wavelet coefficients (see Clyde and George 2000 for Bayesian approaches that use mixture modeling and Donoho, Johnstone, Kerkyacharian, and Picard 1995). Unlike those authors, we have several curves, the spectra of 39 dough pieces, each of which is transformed to wavelet coefficients. We then select some of these wavelet coefficients (the same ones for each spectrum), not because they give a good reconstruction of the curves (which they do not) or to remove noise (of which there is very little to start with) from the curves, but rather because the selected coefficients are useful for predicting some other quantity measured on the dough pieces. One consequence of this is that it is not necessarily the large wavelet coefficients that will be useful; small coefficients in critical regions of the spectrum also may carry important predictive information. Thus the standard thresholding or shrinkage approaches are just not relevant to this problem.

2. PRELIMINARIES

2.1 Wavelet Bases and Wavelet Transforms

Wavelets are families of functions that can accurately describe other functions in a parsimonious way. In $L^2(\mathbb{R})$, for example, an orthonormal wavelet basis is obtained as translations and dilations of a *mother wavelet* ψ as $\psi_{j,k}(x) = 2^{j/2}\psi(2^jx - k)$ with j, k integers. A function f is then represented by a wavelet series as

$$f(x) = \sum_{j,k \in \mathbb{Z}} d_{j,k} \psi_{j,k}(x), \quad (1)$$

with wavelet coefficients $d_{j,k} = \int f(x)\psi_{j,k}(x)dx$ describing features of the function f at the spatial location $2^{-j}k$ and frequency proportional to 2^j (or scale j).

Daubechies (1992) proposed a class of wavelet families that have compact support and a maximum number of vanishing moments for any given smoothness. These are used extensively in statistical applications.

Wavelets have been an extremely useful tool for the analysis and synthesis of discrete data. Let $\mathbf{Y} = (y_1, \dots, y_n)$, $n = 2^J$, be a sample of a function at equally spaced points. This vector of observations can be viewed as an approximation to the function at the fine scale J . A fast algorithm, the *discrete wavelet transform* (DWT), exists for decomposing Y into a set of wavelet coefficients (Mallat 1989) in only $O(n)$ operations. The DWT operates in practice by means of linear recursive filters. For illustration purposes, we can write it in matrix form as $\mathbf{Z} = \mathbf{W}\mathbf{Y}$, where \mathbf{W} is an orthogonal matrix corresponding to the discrete wavelet transform and \mathbf{Z} is a vector of wavelet coefficients describing features of the function at scales from

the fine $J - 1$ to a coarser one, say $J - r$. An algorithm for the inverse construction also exists.

Wavelet transforms can be computed very rapidly and have good compression properties. Because they are localized in both time and frequency, wavelets have the ability to represent many classes of functions in a sparse form by describing important features with few coefficients. (For a general exposition of wavelet theory see Daubechies 1992.)

2.2 Matrix-Variate Distributions

In what follows we use notation for matrix-variate distributions due to Dawid (1981). We write

$$\mathbf{V} - \mathbf{M} \sim \mathcal{N}(\mathbf{\Gamma}, \mathbf{\Sigma})$$

when the random matrix \mathbf{V} has a matrix-variate normal distribution with mean \mathbf{M} and covariance matrices $\gamma_{ij}\mathbf{\Sigma}$ and $\sigma_{jj}\mathbf{\Gamma}$ for its i th row and j th column. Such a \mathbf{V} could be generated as $\mathbf{V} = \mathbf{M} + \mathbf{A}'\mathbf{U}\mathbf{B}$, where \mathbf{M} , \mathbf{A} , and \mathbf{B} are fixed matrices such that $\mathbf{A}'\mathbf{A} = \mathbf{\Gamma}$ and $\mathbf{B}'\mathbf{B} = \mathbf{\Sigma}$ and \mathbf{U} is a random matrix with independent standard normal entries. This notation has the advantage of preserving the matrix structure instead of reshaping \mathbf{V} as a vector. It also makes for much easier formal Bayesian manipulation.

The other notation that we use is

$$\mathbf{W} \sim \mathcal{IW}(\delta; \mathbf{\Sigma})$$

for a random matrix \mathbf{W} with an inverse Wishart distribution with scale matrix $\mathbf{\Sigma}$ and shape parameter δ . The shape parameter differs from the more conventional degrees of freedom, again making for very easy Bayesian manipulations. With \mathbf{U} and \mathbf{B} defined as earlier, and with \mathbf{U} as $n \times p$ with $n > p$, $\mathbf{W} = \mathbf{B}'(\mathbf{U}'\mathbf{U})^{-1}\mathbf{B}$ has an inverse Wishart distribution with shape parameter $\delta = n - p + 1$ and scale matrix $\mathbf{\Sigma}$. The expectation of \mathbf{W} exists for $\delta > 2$ and is then $\mathbf{\Sigma}/(\delta - 2)$. (More details of these notations, and a corresponding form for the matrix-variate T , can be found in Brown 1993, App. A, or Dawid 1981.)

3. MODELING

3.1 Multivariate Regression Model

The basic setup that we consider is a multivariate linear regression model, with n observations on a q -variate response and p explanatory variables. Let \mathbf{Y} denote the $n \times q$ matrix of observed values of the responses and let \mathbf{X} be the $n \times p$ matrix of predictor variables. Our special concern is with functional predictor data; that is, the situation in which each row of X is a vector of observations of a curve $x(t)$ at p equally spaced points.

The standard multivariate normal regression model has, conditional on $\boldsymbol{\alpha}$, \mathbf{B} , $\mathbf{\Sigma}$, and \mathbf{X} ,

$$\mathbf{Y} - \mathbf{1}_n\boldsymbol{\alpha}' - \mathbf{X}\mathbf{B} \sim \mathcal{N}(\mathbf{0}_n, \mathbf{\Sigma}), \quad (2)$$

where $\mathbf{1}_n$ is an $n \times 1$ vector of 1's, $\boldsymbol{\alpha}$ is a $q \times 1$ vector of intercepts, and $\mathbf{B} = (\boldsymbol{\beta}_1, \dots, \boldsymbol{\beta}_q)$ is a $p \times q$ matrix of regression coefficients. Without loss of generality, we assume that the columns of \mathbf{X} have been centered by subtracting their means.

The unknown parameters are α , \mathbf{B} , and the $q \times q$ error covariance matrix Σ . A conjugate prior for this model is as follows. First, given Σ ,

$$\alpha' - \alpha'_0 \sim \mathcal{N}(h, \Sigma) \quad (3)$$

and, independently,

$$\mathbf{B} - \mathbf{B}_0 \sim \mathcal{N}(\mathbf{H}, \Sigma). \quad (4)$$

The marginal distribution of Σ is then

$$\Sigma \sim \mathcal{IW}(\delta; \mathbf{Q}). \quad (5)$$

Note that the priors on both α and \mathbf{B} have covariances dependent on Σ in a way that directly extends the univariate regression natural conjugate prior distributions.

In practice, we let $h \rightarrow \infty$ to represent vague prior knowledge about α and take $\mathbf{B}_0 = \mathbf{0}$, leaving the specification of \mathbf{H} , δ , and \mathbf{Q} to incorporate prior knowledge about our particular application.

3.2 Transformation to Wavelets

We now transform the predictor variables by applying to each row of \mathbf{X} a wavelet transform, as described in Section 2.1. In matrix form, multiplying each row of \mathbf{X} by the same matrix \mathbf{W} is equivalent to multiplying \mathbf{X} on the right side by \mathbf{W}' .

The wavelet transform is orthogonal (i.e., $\mathbf{W}'\mathbf{W} = \mathbf{I}$), and thus (2) can be written as

$$\mathbf{Y} - \mathbf{1}_n \alpha' - \mathbf{X}\mathbf{W}'\mathbf{W}\mathbf{B} \sim \mathcal{N}(\mathbf{I}_n, \Sigma). \quad (6)$$

We can now express the model in terms of wavelet transformations of the predictors as

$$\mathbf{Y} - \mathbf{1}_n \alpha' - \mathbf{Z}\tilde{\mathbf{B}} \sim \mathcal{N}(\mathbf{I}_n, \Sigma), \quad (7)$$

where $\mathbf{Z} = \mathbf{X}\mathbf{W}'$ is now a matrix of wavelet coefficients and $\tilde{\mathbf{B}} = \mathbf{W}\mathbf{B}$ is a matrix of regression coefficients. The transformed prior on $\tilde{\mathbf{B}}$, in the case of inclusion of all predictors, is

$$\tilde{\mathbf{B}} \sim \mathcal{N}(\tilde{\mathbf{H}}, \Sigma), \quad (8)$$

where $\tilde{\mathbf{H}} = \mathbf{W}\mathbf{H}\mathbf{W}'$ and the parameters α and Σ are unchanged by the orthogonal transformations, as are the priors (3) and (5).

In practice, wavelets exploit the recursive application of filters, and the \mathbf{W} -matrix notation is more useful for explanation than for computation. Vannucci and Corradi (1999) proposed a fast recursive algorithm for computing quantities such as $\mathbf{W}\mathbf{H}\mathbf{W}'$. Their algorithm has a useful link to the two-dimensional DWT (DWT2), making computations simple. The matrix $\mathbf{W}\mathbf{H}\mathbf{W}'$ can be computed from \mathbf{H} with an $O(n^2)$ algorithm. (For more details, see secs. 3.1 and 3.2 of Vannucci and Corradi 1999.)

3.3 A Framework for Variable Selection

To perform selection in the wavelet coefficient domain, we further elaborate the prior on $\tilde{\mathbf{B}}$ by introducing a latent binary p -vector γ . The j th element of γ , γ_j , may be either 1 or 0, depending on whether the j th column of \mathbf{Z} is or is not included in the model. When γ_j is unity, the covariance matrix of the corresponding row of $\tilde{\mathbf{B}}$ is “large,” and when γ_j is 0, the covariance matrix is a zero matrix. We have assumed that the prior expectation of $\tilde{\mathbf{B}}$ is $\mathbf{0}$, and so $\gamma_j = 0$ effectively deletes the j th explanatory variable (or wavelet coefficient) from the model. This gives, conditional on γ ,

$$\tilde{\mathbf{B}}_\gamma \sim \mathcal{N}(\tilde{\mathbf{H}}_\gamma, \Sigma), \quad (9)$$

where $\tilde{\mathbf{B}}_\gamma$ and $\tilde{\mathbf{H}}_\gamma$ are just $\tilde{\mathbf{B}}$ and $\tilde{\mathbf{H}}$ with the rows and, in the case of $\tilde{\mathbf{H}}$, columns for which $\gamma_j = 0$ deleted. Under this prior, each row of $\tilde{\mathbf{B}}$ is modeled as having a scale mixture of the type

$$[\tilde{\mathbf{B}}]_{[j]} \sim (1 - \gamma_j)I_0 + \gamma_j N(0, \tilde{h}_{jj}\Sigma), \quad (10)$$

with \tilde{h}_{jj} equal to the j th diagonal element of the matrix $\tilde{\mathbf{H}} = \mathbf{W}\mathbf{H}\mathbf{W}'$ and I_0 a distribution placing unit mass on the $1 \times q$ zero vector. Note that the rows of $\tilde{\mathbf{B}}$ are not independent.

A simple prior distribution $\pi(\gamma)$ for γ takes the γ_j to be independent with $\Pr(\gamma_j = 1) = w_j$ and $\Pr(\gamma_j = 0) = 1 - w_j$, with hyperparameters w_j to be specified, for $j = 1, \dots, p$. In our example we take all of the w_j equal to a common w , so that the nonzero elements of γ have a binomial distribution with expectation pw .

Mixture priors have been widely used for variable selection in the original model space, originally by Leamer (1978) and more recently by George and McCulloch (1997) and Mitchell and Beauchamp (1988) for the linear multiple regression case. Carlin and Chib (1995), Chipman (1996), and Geweke (1996), among others, concentrated on special features of these priors. Clyde, DeSimone, and Parmigiani (1996) used model mixing in prediction problems with correlated predictors when expressing the space of models in terms of an orthogonalization of the design matrix. Their methods are not directly applicable to our situation, because the wavelet transforms do not leave us with an orthogonal design. The use of mixture priors for selection in the multivariate regression setup has been investigated by Brown et al. (1998a,b).

4. SELECTING WAVELET COEFFICIENTS

4.1 Posterior Distribution of γ

The posterior distribution of γ given the data, $\pi(\gamma | \mathbf{Y}, \mathbf{Z})$, assigns a posterior probability to each γ -vector and thus to each possible subset of predictors (wavelet coefficients). This posterior arises from the combination of a likelihood that gives great weight to subsets explaining a high proportion of the variation in the responses \mathbf{Y} and a prior for γ that penalizes large subsets. It can be computed by integrating out α , \mathbf{B} , and Σ from the joint posterior distribution of these parameters and γ given the data. With the vague ($h \rightarrow \infty$) prior for α , this parameter is essentially estimated by the mean \mathbf{Y} in the calibration data (see Smith 1973), and to simplify the formulas

that follow, we now assume that the columns of \mathbf{Y} have been centered. (Full details of the prior to posterior analysis have been given by Brown et al. 1998b, who also considered other prior structures.) After some manipulation, we have

$$\pi(\boldsymbol{\gamma} \mid \mathbf{Y}, \mathbf{Z}) \propto g(\boldsymbol{\gamma}) = \{|\tilde{\mathbf{H}}_\gamma| |\mathbf{Z}'_\gamma \mathbf{Z}_\gamma + \tilde{\mathbf{H}}_\gamma^{-1}|\}^{-q/2} \times |\mathbf{Q}_\gamma|^{-(n+\delta+q-1)/2} \pi(\boldsymbol{\gamma}), \quad (11)$$

where $\mathbf{Q}_\gamma = \mathbf{Q} + \mathbf{Y}'\mathbf{Y} - \mathbf{Y}'\mathbf{Z}_\gamma(\mathbf{Z}'_\gamma \mathbf{Z}_\gamma + \tilde{\mathbf{H}}_\gamma^{-1})^{-1}\mathbf{Z}'_\gamma \mathbf{Y}$ and \mathbf{Z}_γ is \mathbf{Z} with the columns for which $\gamma_j = 0$ deleted. Care is needed in computing (11); the alternative forms discussed later may be useful.

A simplifying feature of this setup is that all of the computations can be formulated as least squares problems with modified \mathbf{Y} and \mathbf{Z} matrices. By writing

$$\tilde{\mathbf{Z}}_\gamma = \begin{pmatrix} \mathbf{Z}_\gamma \tilde{\mathbf{H}}_\gamma^{\frac{1}{2}} \\ \mathbf{I}_{p_\gamma} \end{pmatrix}, \quad \tilde{\mathbf{Y}} = \begin{pmatrix} \mathbf{Y} \\ \mathbf{0} \end{pmatrix},$$

where $\tilde{\mathbf{H}}_\gamma^{1/2}$ is a matrix square root of $\tilde{\mathbf{H}}_\gamma$ and p_γ is the number of 1's in $\boldsymbol{\gamma}$, the relevant quantities entering into (11) can be computed as

$$|\tilde{\mathbf{H}}_\gamma| |\mathbf{Z}'_\gamma \mathbf{Z}_\gamma + \tilde{\mathbf{H}}_\gamma^{-1}| = |\tilde{\mathbf{Z}}'_\gamma \tilde{\mathbf{Z}}_\gamma| \quad (12)$$

and

$$\mathbf{Q}_\gamma = \mathbf{Q} + \tilde{\mathbf{Y}}'\tilde{\mathbf{Y}} - \tilde{\mathbf{Y}}'\tilde{\mathbf{Z}}_\gamma(\tilde{\mathbf{Z}}'_\gamma \tilde{\mathbf{Z}}_\gamma)^{-1}\tilde{\mathbf{Z}}'_\gamma \tilde{\mathbf{Y}}; \quad (13)$$

that is, \mathbf{Q}_γ is given by \mathbf{Q} plus the residual sum of products matrix from the least squares regression of $\tilde{\mathbf{Y}}$ on $\tilde{\mathbf{Z}}_\gamma$. The QR decomposition can then be used (see, e.g., Seber 1984, chap. 10, sec. 1.1b), which avoids “squaring” as in (12) and (13).

4.2 Metropolis Search

Equation (11) gives the posterior probability of each of the 2^p different $\boldsymbol{\gamma}$ vectors, and thus of each choice of wavelet coefficient subsets. What remains to do is to look for “good” wavelet components by computing these posterior probabilities. When p is much greater than about 25, too many subsets exist for this to be feasible. Fortunately, we can use simulation methods that will find the $\boldsymbol{\gamma}$ vectors with relatively high posterior probabilities. We can then quickly identify useful coefficients that have high marginal probabilities of $\gamma_j = 1$.

Here we use a Metropolis search, as suggested for model selection by Madigan and York (1995) and applied to variable selection for regression by Brown et al. (1998a), George and McCulloch (1997), and Raftery, Madigan, and Hoeting (1997). The search starts from a randomly chosen $\boldsymbol{\gamma}^0$ and then moves through a sequence of further values of $\boldsymbol{\gamma}$. At each step, the algorithm generates a new candidate $\boldsymbol{\gamma}$ by randomly modifying the current one. Two types of moves are used:

- Add or delete a component by choosing at random one component in the current $\boldsymbol{\gamma}$ and changing its value. This move is chosen with probability ϕ .

- Swap two components by choosing independently at random a 0 and a 1 in the current $\boldsymbol{\gamma}$ and changing both of them. This move is chosen with probability $1 - \phi$.

The new candidate model, $\boldsymbol{\gamma}^*$, is accepted with probability

$$\min \left\{ \frac{g(\boldsymbol{\gamma}^*)}{g(\boldsymbol{\gamma})}, 1 \right\}. \quad (14)$$

Thus a more probable $\boldsymbol{\gamma}^*$ is always accepted, and a less probable one may be accepted. There is scope for further ingenuity in designing the sequence of random moves. For example, moves that add or subtract or swap two or three or more at a time, or a combination of these, may be useful.

The sequence of $\boldsymbol{\gamma}$'s generated by the search is a realization of a Markov chain, and the choice of acceptance probabilities ensures that the equilibrium distribution of this chain is the distribution given by (11). In typical uses of such schemes, the realizations are monitored after a suitable burn-in period to verify that they appear stationary. Here we have a closed form for the posterior distribution, and are using the chain simply to explore this distribution. Thus we have not been so concerned about strict convergence of the Markov chain. Following Brown et al. (1998a), we adopt a strategy of running the chain from a number of different starting points (four here) and looking at the four marginal distributions provided by the computed $g(\cdot)$ values of the visited $\boldsymbol{\gamma}$'s. We also look for good indication of mixing and explorations with returns. Because we know the relative probabilities, we do not need to worry about using a burn-in period.

Note that, given the form of the acceptance probability (14), $\boldsymbol{\gamma}$ -vectors with high posterior probability have a greater chance of appearing in the sequence. Thus we might expect that a long run of such a chain would visit many of the best subsets.

5. PREDICTION

Suppose now that we wish to predict \mathbf{Y}_f , an $m \times q$ matrix of further \mathbf{Y} -vectors given the corresponding \mathbf{X} -vectors, \mathbf{X}_f , ($m \times p$). First, we treat \mathbf{X}_f exactly as the training data have been treated, by subtracting the training data means and transforming to wavelet coefficients \mathbf{Z}_f , ($m \times p$). The model for \mathbf{Y}_f , following the model for the training data (6), is

$$\mathbf{Y}_f - \mathbf{1}_m \boldsymbol{\alpha}' - \mathbf{Z}_f \tilde{\mathbf{B}} \sim \mathcal{N}(\mathbf{I}_m, \boldsymbol{\Sigma}). \quad (15)$$

If we believe our Bayesian mixture model, then logically we should apply the same latent structure model to prediction as well as to training. This has the practical appeal of providing averaging over a range of likely models (Madigan and Raftery 1994).

Results of Brown et al. (1998b) demonstrate that with the columns of \mathbf{Y}_f centered using the mean \mathbf{Y} from the training set, the expectation of the predictive distribution $p(\mathbf{Y}_f \mid \boldsymbol{\gamma}, \mathbf{Z}, \mathbf{Y})$ is given by $\mathbf{Z}_f \hat{\mathbf{B}}_\gamma$ with

$$\hat{\mathbf{B}}_\gamma = (\mathbf{Z}'_\gamma \mathbf{Z}_\gamma + \tilde{\mathbf{H}}_\gamma^{-1})^{-1} \mathbf{Z}'_\gamma \mathbf{Y} = \tilde{\mathbf{H}}_\gamma^{\frac{1}{2}} (\tilde{\mathbf{Z}}'_\gamma \tilde{\mathbf{Z}}_\gamma)^{-1} \tilde{\mathbf{Z}}'_\gamma \tilde{\mathbf{Y}}. \quad (16)$$

Averaging over the posterior distribution of $\boldsymbol{\gamma}$ gives

$$\hat{\mathbf{Y}}_f = \sum_{\boldsymbol{\gamma}} \mathbf{Z}_f \hat{\mathbf{B}}_\gamma \pi(\boldsymbol{\gamma} \mid \mathbf{Z}, \mathbf{Y}), \quad (17)$$

and we might choose to approximate this by some restricted set of γ values, perhaps the r most likely values from the Metropolis search.

6. APPLICATION TO NEAR-INFRARED SPECTROSCOPY OF BISCUIT DOUGHS

We now apply the methodology developed earlier to the spectroscopic calibration problem described in Section 1.1. First, however, we report the results of some other analyses of these data.

6.1 Analysis by Standard Methods

For all of the analyses carried out here, both compositional and spectral data were centered by subtracting the training set means from the training and validation data. The responses, but not the spectral data, were also scaled, to give each of the four variables unit variance in the training set.

Mean squared prediction errors were converted back to the original scale by multiplying them by the training sample variances. This preprocessing of the data makes no difference to the standard analyses, which treat the response variables separately, but it simplifies the prior specification for our multivariate wavelet analysis.

Osborne et al. (1984) derived calibrations by multiple regression, using various stepwise procedures to select wavelengths for each constituent separately. The mean squared error of predictions on the 39 validation samples for their calibrations are reported in the first row of Table 1.

Brown et al. (1999) also selected small numbers of wavelengths to find calibration equations for this example. Their Bayesian decision theory approach differed from the approach of Osborne in being multivariate (i.e., in trying to find one small subset of wavelengths suitable for predicting all four constituents) and in using a more extensive search using simulated annealing. The results for this alternative wavelength selection approach are given in the second row of Table 1.

For the purpose of comparison, we carried out two other analyses using partial least squares regression (PLS) and principal components regression (PCR). These approaches, both of which construct factors from the full spectral data and then regress constituents on the factors, are very much the standard tools in NIR spectroscopy (see, e.g. Geladi and Martens 1996; Geladi, Martens, Hadjiiski, and Hopke 1996). For the computations, we used the PLS Toolbox 2.0 of Wise and Gallagher (Eigenvector Research, Manson, WA). Although there are multivariate versions of PLS, we took the usual approach of calibrating for each constituent separately. The number of factors used, selected in each case by cross-validation on the training set, was five for each of the PLS equations and six

for each of the PCR equations. The results given in rows three and four of Table 1 show that, as usual, there is little to choose between the two methods. These results are for PLS and PCR using the reduced 256-point spectrum that we used for our wavelet analysis. Repeating the analyses using the original 700-point spectrum yielded results very similar to those for PLS and somewhat worse than those reported for PCR.

Because shortly we need to specify a prior distribution on regression coefficients, it is interesting to examine those resulting from a factor-type approach. Combining the coefficients for the regression of constituent on factor scores with the loadings that produce scores from the original spectral variables gives the coefficient vector that would be applied to a measured spectrum to give a prediction. Figure 2 plots these vectors for the PLS equations for the four constituents, showing the smoothness in the 256 coefficients that we attempt to reflect in our prior distribution for \mathbf{B} .

6.2 Wavelet Transforms of Spectra

To each spectrum we apply a wavelet transform, converting it to a set of 256 wavelet coefficients. We used the MATLAB toolbox Wavbox 4.3 (Taswell 1995) for this step. Using spectra with 2^m (m integer) points is not a real restriction here. Methods exist to overcome the limitation, allowing the DWT to be applied on any length of data. We used MP(4) (Daubechies 1992, p. 194), wavelets with four vanishing moments. The Daubechies wavelets have compact support, important for good localization, and a maximum number of vanishing moments for a given smoothness. A large number of vanishing moments leads to high compressibility, because the fine-scale wavelet coefficients are essentially 0 where the functions are smooth. On the other hand, support of the wavelets increases with an increasing number of vanishing moments, so there is a trade-off with the localization properties. Some rather limited exploration suggested that the chosen wavelet family is a good compromise for these data.

The graphs on the right side of Figure 1 show the wavelet transforms corresponding to the three NIR spectra in the left column. Coefficients are ordered from coarsest to finest.

6.3 Prior Settings

We need to specify the values of \mathbf{H} , δ , and \mathbf{Q} in (3), (4), and (5) and the prior probability w that an element of γ is 1. We wish to put in weak but proper prior information about Σ . We choose $\delta = 3$, because this is the smallest integer value such that the expectation of Σ , $E(\Sigma) = \mathbf{Q}/(\delta - 2)$, exists. The scale matrix \mathbf{Q} is chosen as $\mathbf{Q} = k \mathbf{I}_q$ with $k = .05$, comparable in size to the expected error variances of the standardized \mathbf{Y} given \mathbf{X} . With δ small, the choice of \mathbf{Q} is unlikely to be critical.

Much more likely to be influential are the choices of \mathbf{H} and w in the priors for \mathbf{B} and γ . To reflect the smoothness in the coefficients B , as exemplified in Figure 2, while keeping the form of \mathbf{H} simple, we have taken \mathbf{H} to be the variance matrix of a first-order autoregressive process, with $h_{ij} = \sigma^2 \rho^{|i-j|}$. We derived the values $\sigma^2 = 254$ and $\rho = .32$ by maximizing a type II likelihood (Good 1965). Integrating α , \mathbf{B} , and Σ from the joint distribution given by (2), (3), (4), and (5) for the

Table 1. Mean Squared Errors of Prediction on the 39 Biscuit Dough Pieces in the Validation Set Using Four Calibration Methods

Method	Fat	Sugar	Flour	Water
Stepwise MLR	.044	1.188	.722	.221
Decision theory	.076	.566	.265	.176
PLS	.151	.583	.375	.105
PCR	.160	.614	.388	.106

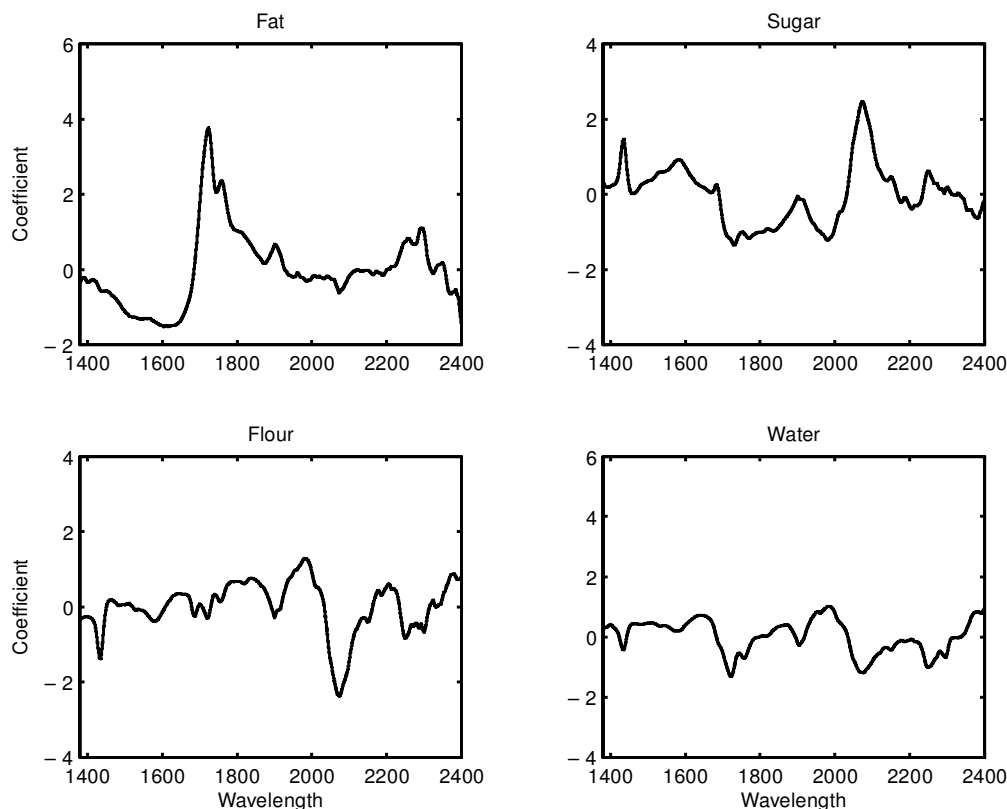


Figure 2. Coefficient Vectors From 5-Factor PLS Equations.

regression on the full untransformed spectra, with $h \rightarrow \infty$ and $\mathbf{B}_0 = \mathbf{0}$, we get

$$f \propto |\mathbf{K}|^{-q/2} |\mathbf{Q}|^{(\delta+q-1)/2} |\mathbf{Q} + \mathbf{Y}'\mathbf{K}^{-1}\mathbf{Y}|^{-(\delta+n+q-1)/2}, \quad (18)$$

where

$$\mathbf{K} = \mathbf{I}_n + \mathbf{X}\mathbf{H}\mathbf{X}'$$

and the columns of \mathbf{Y} are centered, as in Section 4.1. With $k = .05$ and $\delta = 3$ already fixed, (18) is a function, via \mathbf{H} , of σ^2 and ρ . We used for our prior the values of these hyperparameters that maximize (18). Possible underestimation due to the use of the full spectra was taken into account by multiplying the estimate of σ^2 by the inflation factor 256/20, reflecting our prior belief as to the expected number of included coefficients.

Figure 3 shows the diagonal elements of the matrix $\tilde{\mathbf{H}} = \mathbf{W}\mathbf{H}\mathbf{W}'$ implied by our choice of \mathbf{H} . The variance matrix of the i th column of \mathbf{B} in (8) is $\sigma_{ii}\tilde{\mathbf{H}}$, so this plot shows the pattern in the prior variances of the regression coefficients when the predictors are wavelet coefficients. The wavelet coefficients are ordered from coarsest to finest, so the decreasing prior variance means that there will be more shrinkage at the finer levels. This is a logical consequence of the smoothness that we have tried to express in the prior distribution. The spikes in the plot at the level transitions are from the boundary condition problems of the discrete wavelet transform.

We know from experience that good predictions can usually be obtained with 10 or so selected spectral points in examples of this type. Having no previous experience in selecting

wavelet coefficients, and wanting to induce a similarly “small” model without constraining the possibilities too severely, we chose w in the prior for $\boldsymbol{\gamma}$ so that the expected model size was $pw = 20$. We have given equal prior probability here to coefficients at the different levels. Although we considered the possibility of varying w in blocks, we had no strong prior opinions about which levels were likely to provide the most useful coefficients, apart from a suspicion that neither the coarsest nor the finest levels would feature strongly.

6.4 Implementing the Metropolis Search

We chose widely different starting points for the four Metropolis chains, by setting to 1 the first 1, the first 20, the first 128 (i.e., half), and all elements of $\boldsymbol{\gamma}$. There were 100,000 iterations in each run, where an iteration comprised either adding/deleting or swapping, as described in Section 4.2. The two moves were chosen with equal probability $\phi = 1/2$. Acceptance of the possible move by (14) was by generation of a Bernoulli random variable. Computation of $g(\boldsymbol{\gamma})$ and $g(\boldsymbol{\gamma}^*)$ was done using the QR decomposition of MATLAB.

For each chain, we recorded the visited $\boldsymbol{\gamma}$'s and their corresponding relative probability $g(\boldsymbol{\gamma})$. No burn-in was necessary, as relatively unlikely $\boldsymbol{\gamma}$'s would automatically be down-weighted in our analysis. There were approximately 38,000–40,000 successful moves for each of the four runs. Of these moves, around 95% were swaps. The relative probabilities of the set of distinct visited $\boldsymbol{\gamma}$ were then normalized to 1 over this set. Figure 4 plots the marginal probabilities for components of $\boldsymbol{\gamma}$, $P(\gamma_j = 1)$, $j = 1, \dots, 256$. The spikes show where regressor variables have been included in subsets with

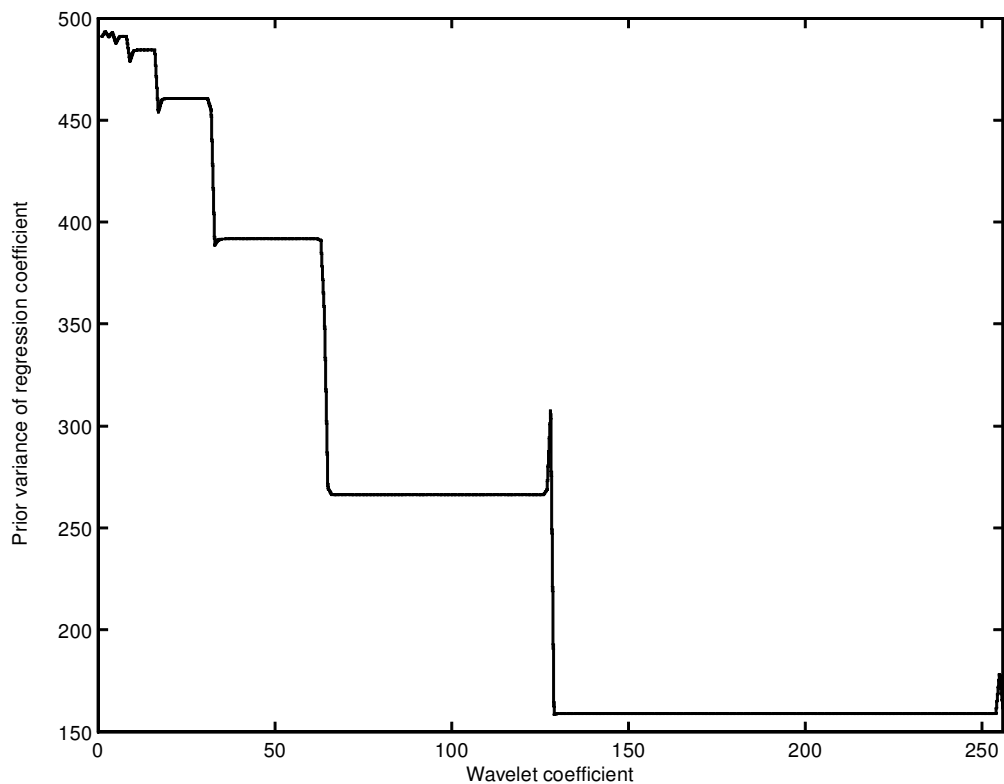


Figure 3. Diagonal Elements of $\tilde{H} = \tilde{W}\tilde{W}'$, With \tilde{h}_{ii} Plotted Against i .

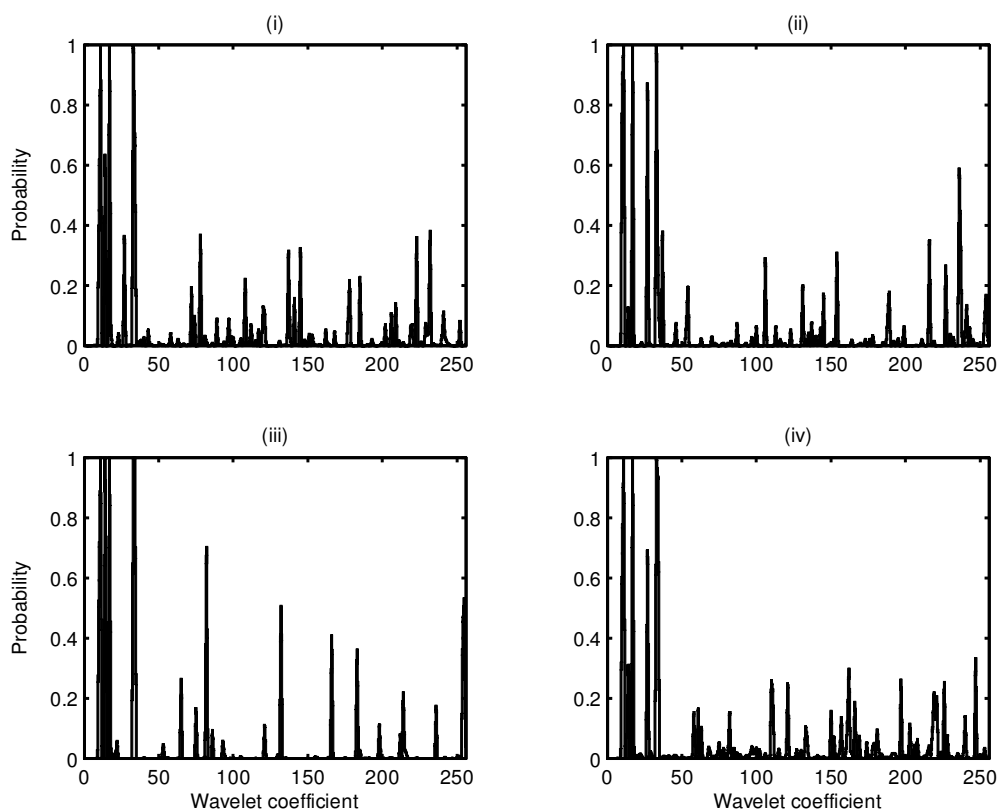


Figure 4. Marginal Probabilities of Components of γ for Four Runs.

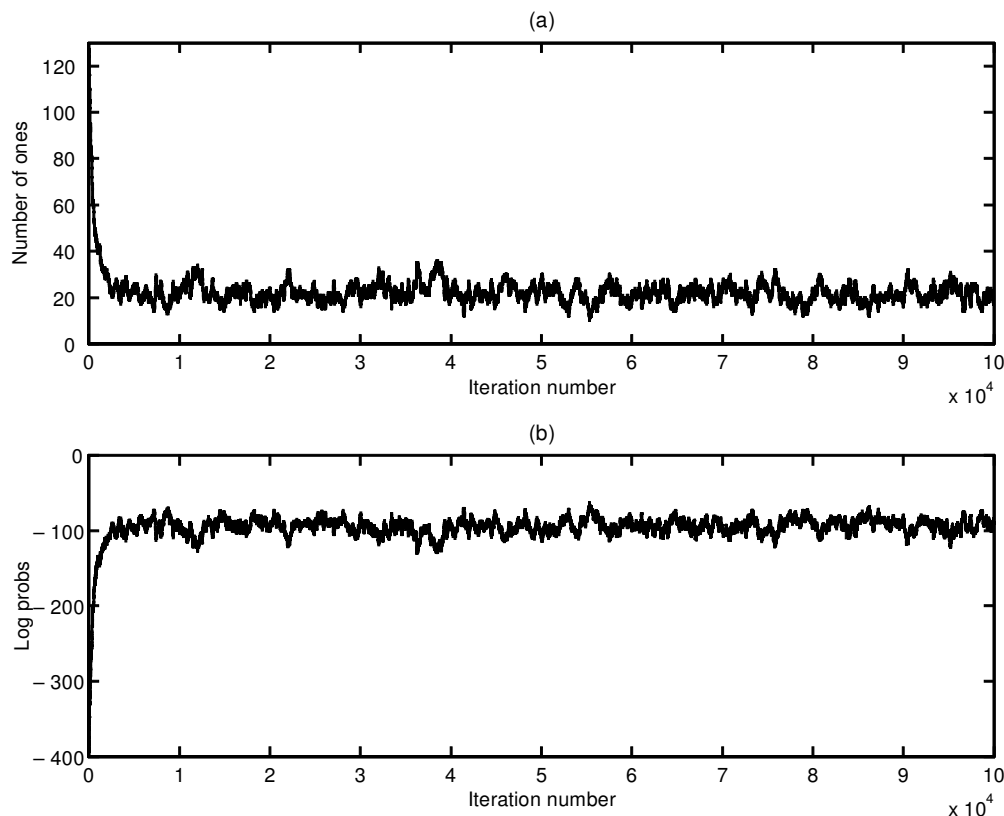


Figure 5. Plots in Sequence Order for Run (iii). (a) The number of 1's; (b) log relative probabilities.

high probability. For one of the runs, (iii), Figure 5 gives two more plots: the number of 1's, and the log-relative probabilities, $\log(g(\gamma))$, of the visited γ , plotted over the 100,000 iterations. The other runs produced very similar plots, quickly moving toward models of similar dimensions and posterior probability values.

Despite the very different starting points, the regions explored by the four chains have clear similarities in that plots of marginals are overall broadly similar. However, there are also clear differences, with some chains making frequent use of variables not picked up by others. Although we would not claim convergence, all four chains arrive at some similarly "good," albeit different, subsets. With mutually correlated predictor variables, as we have here, there will always tend to be many solutions to the problem of finding the best predictors. We adopt the pragmatic stance that all we are trying to do is identify some of the good solutions. If we happen to miss some other good ones, then this is unfortunate but not disastrous.

6.5 Results

We pooled the distinct γ 's visited by the four chains, normalized the relative probabilities, and ordered them according to probability. Then we predicted the further 39 unseen samples using Bayes model averaging. Mean squared prediction errors converted to the original scale were .063, .449, .348, and .050. These results use the best 500 models, accounting for almost 99% of the total visited probability and using 219 wavelet coefficients. They improve considerably on all of the standard methods reported in Table 1. The single best subset

among the visited ones had 10 coefficients, accounted for 9% of the total visited probability, and least squares predictions gave mean squared errors of .059, .466, .351 and .047.

Examining the scales of the selected coefficients is interesting. The model making the best predictions used coefficients (10, 11, 14, 17, 33, 34, 82, 132, 166, 255), which include (0%, 0%, 0%, 37%, 6%, 6%, 1%, 2%) of all of the coefficients at the eight levels from the coarsest to the finest scales. The most useful coefficients seem to be in the middle of the scale, with some more toward the finer end. Note that the very coarsest coefficients, which would be essential to any reconstruction of the spectra, are not used, despite the smoothness implicit in \mathbf{H} and the consequent increased shrinkage associated with finer scales, as seen in Figure 3.

Some idea of the locations of the wavelet coefficients selected by the modal model can be obtained from Figure 6. Here we used the linearity of the wavelet transform and the linearity of the prediction equation to express the prediction equation as a vector of coefficients to be applied to the original spectral data. This vector is obtained by applying the inverse wavelet transform to the columns of the matrix of the least squares estimates of the regression coefficients. Because selection has discarded unnecessary details, most of the coefficients are very close to 0. We thus display only the 1600–1800 nm range, which permits a better view of the features of coefficients in the range of interest. These coefficient vectors can be compared directly with those shown in Figure 2. They are applied to the same spectral data to produce predictions, and despite the different scales, some comparisons are possible. For example, the two coefficient vectors for fat can be eas-

ily interpreted. Fats and oils have absorption bands at around 1730 and 1765 nm (Osborne et al. 1993), and strong positive coefficients in this region are the major features of both plots. The wavelet-based equation (Fig. 6) is simpler, the selection having discarded much unnecessary detail. The other coefficient vectors show less resemblance and are also harder to interpret. One thing to bear in mind in interpreting these plots is that the four constituents add up to 100%. Thus it should not be surprising to find the fat measurement peak in all eight plots nor that the coefficient vectors for sugar and flour are so strongly inversely related.

Finally, we comment briefly on results that we obtained by investigating logistic transformations of the data. Our response variables are in fact percentages and are constrained to sum up to 100; thus they lie on a simplex, rather than in the full q -dimensional space. Sample ranges are 18 ± 3 for fat, 17 ± 7 for sugar, 51 ± 6 for flour, and 14 ± 3 for water.

Following Aitchison (1986), we transformed the original Y 's into log ratios of the form

$$\begin{aligned} Z_1 &= \ln(Y_1/Y_3), & Z_2 &= \ln(Y_2/Y_3), & \text{and} \\ Z_3 &= \ln(Y_4/Y_3). \end{aligned} \quad (19)$$

The choice of the third ingredient for the denominator was the most natural, in that flour is the major constituent and also because ingredients in recipes often are expressed as a ratio to flour content. We centered and scaled the Z variables and recomputed empirical Bayes estimates for σ^2 and ρ . (The other hyperparameters were not affected by the logistic transformation.) We then ran four Metropolis chains using the start-

ing points used previously with the data in the original scale. Diagnostic plots and plots of the marginals appeared very similar to those of Figures 4 and 5. The four chains visited 151,183 distinct models. We finally computed Bayes model averaging and least squares predictions with the best model, unscaling the predicted values and transforming them back to the original scale as

$$\begin{aligned} Y_i &= \frac{100 \exp(Z_i)}{\sum_{j=1}^3 \exp(Z_j) + 1}, \quad i = 1, 2, \\ Y_3 &= \frac{100}{\sum_{j=1}^3 \exp(Z_j) + 1}, \quad \text{and} \quad Y_4 = \frac{100 \exp(Z_3)}{\sum_{j=1}^3 \exp(Z_j) + 1}. \end{aligned}$$

The best 500 visited models accounted for 99.4% of the total visited probability, used 214 wavelet coefficients, and gave Bayes mean squared prediction errors of .058, .819, .457, and .080. The single best subset among the visited ones had 10 coefficients, accounted for 15.8% of the total visited probability, and gave least squares prediction errors of .091, .793, .496, and .119. The logistic transformation does not seem to have a positive impact on the predictive performance, and overall, the simpler linear model seems to be adequate. Our approach gives Bayes predictions on the original scale satisfying the constraint of summing exactly to 100. This stems from the conjugate prior, linearity in Y , zero mean for the prior distribution of the regression coefficients, and vague prior for the intercept. It is easily seen that with \mathbf{X} centered, $\hat{\boldsymbol{\alpha}}' \mathbf{1} = 100$ and $\hat{\mathbf{B}} \mathbf{1} = 0$ for either least squares or Bayes estimates, and hence all predictions sum to 100. In addition, had we imposed a singular distribution to deal with the four responses summing

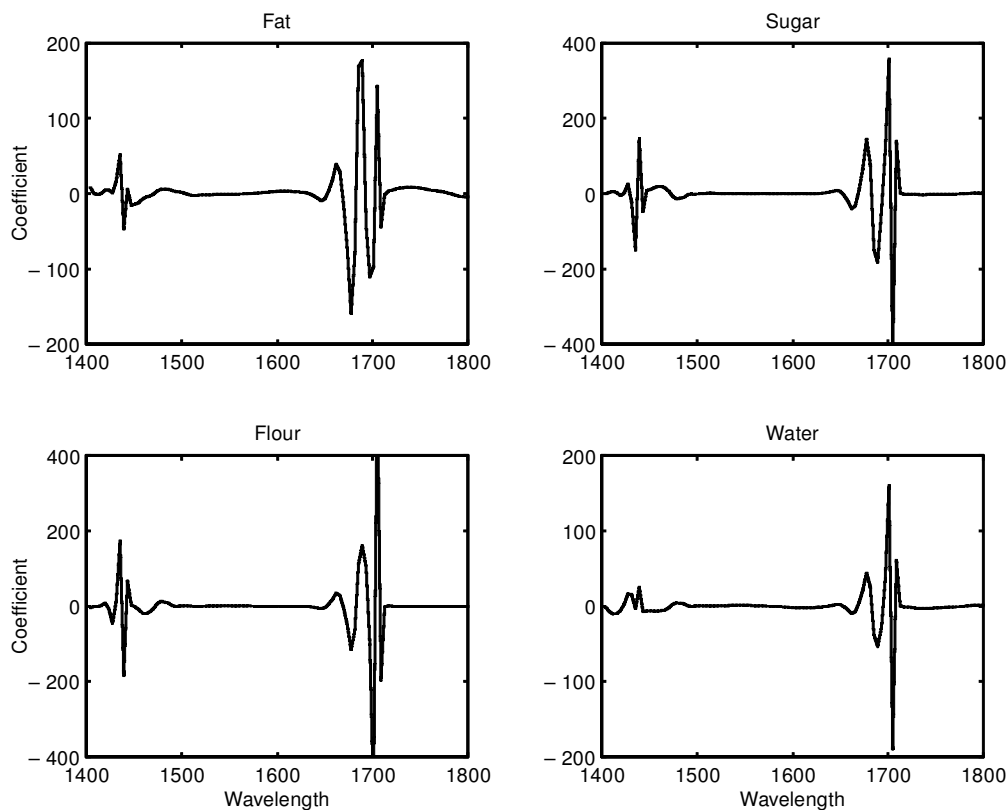


Figure 6. Coefficient Vectors From the "Best" Wavelet Equations.

to 100, then a proper analysis would suggest eliminating one component, but then the predictions for the remaining three components would be as we have derived. Thus our analysis is Bayes for the singular problem, even though superficially it ignores this aspect. This does not address the positivity constraint, but the composition variables are all so far from the boundaries compared to residual error that this is not an issue. Our desire to stick with the original scale is supported by the Beer–Lambert law, which linearly relates absorbance to composition. The logistic transform distorts this linear relationship.

7. DISCUSSION

In specifying the prior parameters for our example, we made a number of arbitrary choices. In particular, the choice of H as the variance matrix of an autoregressive process is a rather crude representation of the smoothness in the coefficients. It might be interesting to try to model this in more detail. However, although there may be room for improvement, the simple structure used here does appear to work well.

Another area for future investigation is the use of more sophisticated wavelet systems in this context. The additional flexibility of wavelet packets (Coifman, Meyer, and Wickerhauser 1992) or m -band wavelets (Mallet, Coomans, Kautsky, and De Vel 1997) might lead to improved predictions, or might be just an unnecessary complication.

[Received October 1998. Revised November 2000.]

REFERENCES

- Aitchison, J. (1986), *The Statistical Analysis of Compositional Data*, London: Chapman and Hall.
- Brown, P. J. (1993), *Measurement, Regression, and Calibration*, Oxford, U.K.: Clarendon Press.
- Brown, P. J., Fearn, T., and Vannucci, M. (1999), "The Choice of Variables in Multivariate Regression: A Bayesian Non-Conjugate Decision Theory Approach," *Biometrika*, 86, 635–648.
- Brown, P. J., Vannucci, M., and Fearn, T. (1998a), "Bayesian Wavelength Selection in Multicomponent Analysis," *Journal of Chemometrics*, 12, 173–182.
- (1998b), "Multivariate Bayesian Variable Selection and Prediction," *Journal of the Royal Statistical Society, Ser. B*, 60, 627–641.
- Carlin, B. P., and Chib, S. (1995), "Bayesian Model Choice via Markov Chain Monte Carlo," *Journal of the Royal Statistical Society, Ser. B*, 57, 473–484.
- Chipman, H. (1996), "Bayesian Variable Selection With Related Predictors," *Canadian Journal of Statistics*, 24, 17–36.
- Clyde, M., DeSimone, H., and Parmigiani, G. (1996), "Prediction via Orthogonalized Model Mixing," *Journal of the American Statistical Association*, 91, 1197–1208.
- Clyde, M., and George, E. I. (2000), "Flexible Empirical Bayes Estimation for Wavelets," *Journal of the Royal Statistical Society, Ser. B*, 62, 681–698.
- Coifman, R. R., Meyer, Y., and Wickerhauser, M. V. (1992), "Wavelet Analysis and Signal Processing," in *Wavelets and Their Applications*, eds M. B. Ruskai, G. Beylkin, R. Coifman, I. Daubechies, S. Mallat, Y. Meyer, and L. Raphael, Boston: Jones and Bartlett, pp. 153–178.
- Cowe, I. A., and McNicol, J. W. (1985), "The Use of Principal Components in the Analysis of Near-Infrared Spectra," *Applied Spectroscopy*, 39, 257–266.
- Daubechies, I. (1992), *Ten Lectures on Wavelets* (Vol. 61, CBMS-NSF Regional Conference Series in Applied Mathematics), Philadelphia: Society for Industrial and Applied Mathematics.
- Dawid, A. P. (1981), "Some Matrix-Variate Distribution Theory: Notational Considerations and a Bayesian Application," *Biometrika*, 68, 265–274.
- Donoho, D., Johnstone, I., Kerkyacharian, G., and Picard, D. (1995), "Wavelet Shrinkage: Asymptopia?" (with discussion), *Journal of the Royal Statistical Society, Ser. B*, 57, 301–369.
- Geladi, P., and Martens, H. (1996), "A Calibration Tutorial for Spectral Data. Part 1. Data Pretreatment and Principal Component Regression Using Matlab," *Journal of Near Infrared Spectroscopy*, 4, 225–242.
- Geladi, P., Martens, H., Hadjiiski, L., and Hopke, P. (1996), "A Calibration Tutorial for Spectral Data. Part 2. Partial Least Squares Regression Using Matlab and Some Neural Network Results," *Journal of Near Infrared Spectroscopy*, 4, 243–255.
- George, E. I., and McCulloch, R. E. (1997), "Approaches for Bayesian Variable Selection," *Statistica Sinica*, 7, 339–373.
- Geweke, J. (1996), "Variable Selection and Model Comparison in Regression," in *Bayesian Statistics 5*, eds J. M. Bernardo, J. O. Berger, A. P. Dawid, and A. F. M. Smith, Oxford, U.K.: Clarendon Press, pp. 609–620.
- Good, I. J. (1965), *The Estimation of Probabilities. An Essay on Modern Bayesian Methods*, Cambridge, MA: MIT Press.
- Hrushka, W. R. (1987), "Data Analysis: Wavelength Selection Methods," in *Near-Infrared Technology in the Agricultural and Food Industries*, eds P. Williams and K. Norris, St Paul, MN: American Association of Cereal Chemists, pp. 35–55.
- Leamer, E. E. (1978), "Regression Selection Strategies and Revealed Priors," *Journal of the American Statistical Association*, 73, 580–587.
- Madigan, D., and Raftery, A. E. (1994), "Model Selection and Accounting for Model Uncertainty in Graphical Models Using Occam's Window," *Journal of the American Statistical Association*, 89, 1535–1546.
- Madigan, D., and York, J. (1995), "Bayesian Graphical Models for Discrete Data," *International Statistical Review*, 63, 215–232.
- Mallat, S. G. (1989), "Multiresolution Approximations and Wavelet Orthonormal Bases of $L_2(\mathbb{R})$," *Transactions of the American Mathematical Society*, 315, 69–87.
- Mallet, Y., Coomans, D., Kautsky, J., and De Vel, O. (1997), "Classification Using Adaptive Wavelets for Feature Extraction," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 19, 1058–1066.
- McClure, W. F., Hamid, A., Giesbrecht, F. G., and Weeks, W. W. (1984), "Fourier Analysis Enhances NIR Diffuse Reflectance Spectroscopy," *Applied Spectroscopy*, 38, 322–329.
- Mitchell, T. J., and Beauchamp, J. J. (1988), "Bayesian Variable Selection in Linear Regression," *Journal of the American Statistical Association*, 83, 1023–1036.
- Osborne, B. G., Fearn, T., and Hindle, P. H. (1993), *Practical NIR Spectroscopy*, Harlow, U.K.: Longman.
- Osborne, B. G., Fearn, T., Miller, A. R., and Douglas, S. (1984), "Application of Near-Infrared Reflectance Spectroscopy to Compositional Analysis of Biscuits and Biscuit Doughs," *Journal of the Science of Food and Agriculture*, 35, 99–105.
- Raftery, A. E., Madigan, D., and Hoeting, J. A. (1997), "Bayesian Model Averaging for Linear Regression Models," *Journal of the American Statistical Association*, 92, 179–191.
- Seber, G. A. F. (1984), *Multivariate Observations*, New York: Wiley.
- Smith, A. F. M. (1973), "A General Bayesian Linear Model," *Journal of the Royal Statistical Society, Ser. B*, 35, 67–75.
- Taswell, C. (1995), "Wavbox 4: A Software Toolbox for Wavelet Transforms and Adaptive Wavelet Packet Decompositions," in *Wavelets and Statistics*, eds A. Antoniadis and G. Oppenheim, New York: Springer-Verlag, pp. 361–375.
- Trygg, J., and Wold, S. (1998), "PLS Regression on Wavelet-Compressed NIR Spectra," *Chemometrics and Intelligent Laboratory Systems*, 42, 209–220.
- Vannucci, M., and Corradi, F. (1999), "Covariance Structure of Wavelet Coefficients: Theory and Models in a Bayesian Perspective," *Journal of the Royal Statistical Society, Ser. B*, 61, 971–986.
- Walczak, B., and Massart, D. L. (1997), "Noise Suppression and Signal Compression Using the Wavelet Packet Transform," *Chemometrics and Intelligent Laboratory Systems*, 36, 81–94.
- Wold, S., Martens, H., and Wold, H. (1983), "The Multivariate Calibration Problem in Chemistry Solved by PLS," in *Matrix Pencils*, eds A. Ruhe and B. Kagstrom, Heidelberg: Springer, pp. 286–293.