



Available at

www.ElsevierMathematics.com

POWERED BY SCIENCE @ DIRECT®

Journal of Statistical Planning and
Inference 128 (2005) 149–164

journal of
statistical planning
and inference

www.elsevier.com/locate/jspi

A predictive density approach to predicting a future observable in multilevel models[☆]

David Afshartous^{a,*}, Jan de Leeuw^b

^a*School of Business Administration, University of Miami, Coral Gables, FL 33124-8237, USA*

^b*Department of Statistics, University of California, Los Angeles, CA 90095-1554, USA*

Received 25 November 2002; accepted 29 October 2003

Abstract

A predictive density function g^* is obtained for the multilevel model which is optimal in minimizing a criterion based on Kullback–Leibler divergence for a restricted class of predictive densities, thereby extending results for the normal linear model (J. Amer. Statist. Assoc. 81 (1986) 196). Based upon this predictive density approach, three prediction methods are examined: multilevel, prior, and OLS. The OLS prediction method corresponds to deriving a predictive density separately in each group, while the prior prediction method corresponds to deriving a predictive density for the entire model. The multilevel prediction method merely adjusts the prior prediction method by employing a well-known shrinkage estimator from multilevel model estimation. Multilevel data are simulated in order to assess the performance of these three methods. Both predictive intervals and predictive mean square error (PMSE) are used to assess the adequacy of prediction. The multilevel prediction method outperforms the OLS and prior prediction methods, somewhat surprising since the OLS and prior prediction methods are derived from the Kullback–Leibler divergence criterion. This suggests that the restricted class of predictive densities suggested by Levy and Perng for the normal linear model may need to be expanded for the multilevel model.

© 2003 Elsevier B.V. All rights reserved.

Keywords: Prediction; Predictive density; Multilevel model; Kullback–Leibler divergence; Predictive interval

[☆] This research was supported by a grant from the National Institute for Statistical Sciences.

* Corresponding author. Tel.: +1-305-2846595; fax: +1-305-2842321.

E-mail address: afshar@miami.edu (David Afshartous).

1. Introduction

A basic problem in predictive inference involves the prediction of a future observable Z based on the observed data Y in some passed experiment. Moreover, Z need not arise from the same stochastic model as Y . One approach to this problem is to attempt to “estimate” the stochastic model from which Z arises. Given such an estimate, there exist several options for predicting the future observable, e.g., the expected value of the stochastic process. Many authors have investigated this approach, often labeled the predictive density or predictive likelihood method (Levy and Perng, 1986; Butler, 1986; Geisser, 1971). Another approach is to forgo density estimation and seek to minimize some expected loss function, often within some prescribed class of predictors (Rao, 1987; Gotway and Cressie, 1993; Goldberger, 1962). Optimal predictors for both approaches have been derived for the general linear model. Moreover, there exist extensions to the multivariate case (Guttman and Hougaard, 1985; Keyes and Levy, 1996). The purpose of this paper is to extend the optimal predictive density results to the multilevel model. The outline of this paper is as follows: In Section 1.1 we review the notation of the multilevel model; in Section 2 we present the predictive density approach and the main result by Levy and Perng (1986) for the general linear model. In Sections 2.1–2.3 we develop and apply this result to the multilevel model, thereby obtaining three predictive densities with which to predict a future observation in a hierarchical dataset. In Section 2.4 we describe a simulation study to assess the predictive performance of these three densities; in Section 3 we present the results, and finally in Section 4 we provide a brief summary and directions for future research.

1.1. The multilevel model

Multilevel modeling is a tool often used when analyzing hierarchical data, e.g., students grouped within schools. In the multilevel model prediction problem, we seek to predict a future observable y_{*j} , i.e., a future case of the j th group. Restricting the discussion to the simple case of n_j primary units (level 1: students) grouped within J secondary units (level 2: schools), the basic multilevel model has the following level-1 model equation:

$$Y_j = X_j \beta_j + r_j. \quad (1.1)$$

Each X_j has dimensions $n_j \times p$, and $r_j \sim N(0, \sigma^2 \Psi_j)$, with Ψ_j usually taken as I_{n_j} . Some or all of the level-1 coefficients, β_j , are random variables, and may also be functions of level-2 (school) variables:

$$\beta_j = W_j \gamma + u_j. \quad (1.2)$$

The vector γ is of length q , each W_j has dimension $p \times q$ and is a matrix of background variables on the j th group, and $u_j \sim N(0, \tau)$. The elements of the random vector β_j are not independent as τ is not necessarily diagonal. For instance, for each regression equation there might exist a covariance between the slope and intercept.

The single equation model is obtained by combining Eqs. (1.1) and (1.2):

$$Y_j = X_j W_j \gamma + X_j u_j + r_j, \quad (1.3)$$

which may be viewed as a special case of the mixed linear model, with fixed effects γ and random effects u_j .¹ The expectation of y_j is $X_j W_j \gamma$ and its dispersion is $V_j = X_j \tau X_j' + \sigma^2 I$. Observations in the same group have correlated disturbances, and this correlation will be larger if their predictor profiles are more alike in the metric τ (de Leeuw and Kreft, 1995). Letting $d_j = Y_j - X_j W_j \gamma$, the full log-likelihood for the j th unit is

$$L_j(\sigma^2, \tau, \gamma) = -\frac{n_j}{2} \log(2\pi) - \frac{1}{2} \log |V_j| - \frac{1}{2} d_j' V_j^{-1} d_j. \quad (1.4)$$

Due to independence across level-2 units, we write the log-likelihood for the entire model as a sum of unit log-likelihoods, i.e.,

$$L(\sigma^2, \tau, \gamma) = \sum_{j=1}^J L_j(\sigma^2, \tau, \gamma). \quad (1.5)$$

Estimates of σ^2 , τ , and γ are obtained via full or restricted maximum likelihood. These estimates may in turn be employed in various approaches to estimate the level-1 coefficients β_j .² For a full review of estimation in multilevel models see Bryk and Raudenbush (1992). Although multilevel model estimation is an important topic, it is not the focus of this paper. The focus here lies in the prediction of a future observable y_{*j} and we shall employ a predictive density approach to this problem.

2. Predictive density approach

Let $f(y; \theta)$ denote the density function for Y and $g(z | y, \theta)$ denote the density function of Z conditioned upon having observed Y . The forms of f and g are assumed known, they are not necessarily the same, and they share the common parameter θ which belongs to some parameter space Θ . Hence the past experiment is informative for the future. A prediction function $s(z; y)$ for z is an estimator of $g(z | y, \theta)$, and if s is a density we call it a predictive density. Levy and Perng (1986) discuss this problem in the context of the general linear model: Consider an n -dimensional random vector Y and the m -dimensional random vector Z , where $Y = X\beta + \varepsilon$ and $Z = W\beta + \tau$, with the usual independence and constant variance assumptions for the error terms ($\varepsilon \sim N(0, \sigma^2 I_n)$ and $\tau \sim N(0, \sigma^2 I_m)$). Here $\beta \in \Omega_\beta \subset R^p$ is an unknown $p \times 1$ vector of regression coefficients while σ^2 is an unknown but positive scalar. It is further assumed that ε and τ are independent. Letting $p_n(y; X, \beta, \sigma^2)$ and $p_m(z; W, \beta, \sigma^2)$ denote the multivariate normal density functions of Y and Z , respectively, we have:

$$p_n(y; X, \beta, \sigma^2) = N(X\beta, \sigma^2 I_n),$$

$$p_m(z; W, \beta, \sigma^2) = N(W\beta, \sigma^2 I_m),$$

¹ For an excellent review of the estimation of fixed and random effects in the general mixed model see Robinson (1991).

² The term “estimation” is being used somewhat loosely when speaking of an estimate of β_j since β_j is a random variable. One may consider an estimate of the random variable β_j as an estimate of the mean of its distribution.

where $N(\mu, \sigma^2 I_k)$ represents the k -dimensional multivariate density function with mean vector μ and dispersion $\sigma^2 I_k$.

Under Kullback–Leibler (Kullback and Leibler, 1951) information loss,³ Levy and Perng (1986) derive an optimal estimator for the density of Z within a prescribed class of density estimators. Levy and Perng restrict the collection of possible density estimators to a subset of predictive densities, Ψ , and the density within this subset which minimizes the Kullback–Leibler measure is selected. Specifically, they consider the statistics defined by

$$t = t(y, z) = (z - W \hat{\beta}) / (n^{1/2} \hat{\sigma}), \tag{2.1}$$

where $\hat{\beta} = (X'X)^{-1} X'y$ and $\hat{\sigma}^2 = (y - X \hat{\beta})'(y - X \hat{\beta}) / n$ are the maximum likelihood estimates of β and σ^2 , respectively. Then Ψ is defined as the collection of all predictive densities that are functions of the statistic t , i.e.,

$$\Psi = \{s(z; W, y, X) : s(z; W, y, X) = g(t(y, z))\}, \tag{2.2}$$

where g is any probability density function. Two reasons for restricting attention to this class are provided: (1) It contains several commonly used predictors, and (2) the statistic $t(y, z)$ used to define Ψ results from a sequence of data reductions by applying the invariance principle under reasonable groups of transformations. They elaborate by demonstrating maximal invariance with respect to specific groups of transformations; see Levy and Perng (1986, p.197) for further details.

Recall, if $s(z; W, y, X)$ is a predictive density estimate for $p_m(z; W, \beta, \sigma^2)$, then the Kullback–Leibler divergence is defined as

$$\begin{aligned} D_{\beta, \sigma^2}(p_m, s) &= \int_{R^n} p_n(y; X, \beta, \sigma^2) \int_{R^m} p_m(z; W, \beta, \sigma^2) \\ &\quad \times \log p_m(z; W, \beta, \sigma^2) / s(z; W, y, X) \, dz \, dy \\ &= E_{Y, Z} \log [p_m(Z; W, \beta, \sigma^2) / s(Z; W, Y, X)]. \end{aligned}$$

Thus, a predictive density s is considered optimal with respect to Kullback–Leibler loss if s minimizes D_{β, σ^2} among all possible predictive densities uniformly with respect to β and σ^2 . Their main result is expressed as follows:

Theorem 1. *Let Ψ , D_{β, σ^2} and t be defined as above. The prediction density*

$$\begin{aligned} g^*(z; W, y, X) &= g^*(t(y, z)) \\ &= st_m(n - p, W \hat{\beta}, n \hat{\sigma}^2 A / (n - p)), \end{aligned}$$

where $A = I_m + W(X'X)^{-1} W'$, provides the unique minimum of D_{β, σ^2} among all s in Ψ uniformly in β and σ^2 .

³ Kullback–Leibler information measure was proposed by Atchinson (1975) as a general prediction measure; a discussion of the motivations and properties of this criterion may be found in Larimore (1983).

The notation $st_d(k, b, C)$ denotes a multivariate Student- t density function, where d is the dimension, k the degrees of freedom, b the location parameter, and C the dispersion matrix, and density function as follows:

$$st_d(k, b, C) = \Gamma[(k + d)/2] / [\pi^{d/2} \Gamma[k/2]] \\ \times [\det(kC)]^{1/2} \{1 + (z - b)'(kC)^{-1}(z - b)\}^{-(k+d)/2}.$$

As noted by Levy and Perng (1986), assuming a non-informative diffuse prior for (β, σ^2) , $f(\beta, \sigma^2) \propto 1/\sigma^2$, the predictive density g^* may be interpreted as a Bayesian predictive density.

Let us examine this predictive density further by using it to create a predictive interval for z in the case where z has dimension one. For large values of n , our predictive interval is centered around the mean of the predictive density, $\hat{z} = W\hat{\beta}$, with margin of error taken as $1.96n\hat{\sigma}^2 A/(n - p)$, in this case a scalar since we have $A = 1 + W(X'X)^{-1}W'$. Upon closer examination, however, we see that this interval is very close to the standard exact prediction interval in linear regression. Specifically, the predictive density variance can be written as follows:

$$n\hat{\sigma}^2 A/(n - p) = \frac{n}{n - p} [\hat{\sigma}^2 + \hat{\sigma}^2 W(X'X)^{-1}W'] \\ = \frac{n}{n - p} [\hat{\text{var}}(z) + \hat{\text{var}}(W\hat{\beta})] \\ = \frac{n}{n - p} [\hat{\text{var}}(z) + \hat{\text{var}}(\hat{z})].$$

Thus, recalling that our usual exact predictive interval in linear regression has margin of error $t_{n-p, \alpha/2} [\hat{\text{var}}(z) + \hat{\text{var}}(\hat{z})]$, the only difference we obtain by using this predictive density to form a prediction interval is the adjustment of the term $n/(n - p)$ in the expression above. Hence, the resulting interval based on this optimal predictive density would be wider than the exact predictive interval.

We would like to extend this result to the multilevel model. We shall do this in three different ways. First, we extend the theorem above to each of the J groups in the multilevel model as independent OLS regression equations. Thus, in each of the J groups, the prediction problem is identical to the presentation above. This method is referred to as the OLS prediction method. Second, we do not ignore the multilevel structure and write the network of J models as one large model and derive the corresponding predictive density for this model. For reasons that will become clear later, this is called the prior prediction method. Finally, we alter the prior prediction method by utilizing a well-known result for the multilevel model to yield the multilevel prediction method.

2.1. OLS prediction method

Here, the level-1 β_j coefficients are not random variables regressed on level-2 variables. Instead, we simply have J separate regressions:

$$Y_j = X_j \beta_j + r_j, \tag{2.3}$$

and desire to predict the future observation in the j th group, y_{*j} :

$$y_{*j} = X_{*j}\beta_j + r_{*j}. \tag{2.4}$$

If y_{*j} were observed, X_{*j} would represent a row of the X_j design matrix and we have $r_{*j} \sim N(0, \sigma^2)$. Thus, we may immediately apply Theorem 1 above to yield the predictive density for y_{*j} :

$$g^*(y_{*j}; X_{*j}, Y_j, X_j) = st_1(n - p, X_{*j}\hat{\beta}_j, n\hat{\sigma}_j^2 A_j / (n - p)), \tag{2.5}$$

where $A_j = 1 + X_{*j}(X_j'X_j)^{-1}X_{*j}'$; $\hat{\beta}_j$ and $\hat{\sigma}_j^2$ are the usual OLS estimates for slope and residual variance.

We employ this predictive density to construct a predictive interval by taking its expected value, $X_{*j}\hat{\beta}_j$, as our point predictor for y_{*j} and use its variance to form our margin of error. Formally, we have the following predictive interval:

$$X_{*j}\hat{\beta}_j \pm t_{n-p, 0.975}\hat{\sigma}_j[nA_j/(n - p)]^{1/2}, \tag{2.6}$$

where $t_{n-p, 0.975}$ is the 0.975 critical value for a t distribution with $n - 2$ degrees of freedom.

2.2. Prior prediction method

In this case the structure of the data is incorporated via the multilevel model as discussed earlier. However, we first need to do some re-arranging. We shall manipulate the notation in the multilevel model such that it is presented as a special case of the general linear model. By appropriately stacking the data for each of the J level-2 units, we may write the model for the entire data without subscripts. Thus, we have

$$Y = X\beta + r, \tag{2.7}$$

with r normally distributed with mean 0 and dispersion Ψ where

$$Y = (Y'_1, Y'_2, \dots, Y'_J)',$$

$$\beta = (\beta'_1, \beta'_2, \dots, \beta'_J)',$$

$$r = (r'_1, r'_2, \dots, r'_J)',$$

$$X = \text{diag}(X_1, \dots, X_J),$$

$$\Psi = \text{diag}(\Psi_1, \dots, \Psi_J),$$

where Ψ_j is usually $\sigma^2 I_{n_j}$. The level-2 equation may also be written in no-subscript form through similar stacking manipulations:

$$\beta = W\gamma + u, \tag{2.8}$$

where u is normally distributed with mean 0 and covariance matrix T where

$$W = (W'_1, W'_2, \dots, W'_J)',$$

$$u = (u'_1, u'_2, \dots, u'_J)',$$

$$T = \text{diag}(\tau, \dots, \tau).$$

The entire model may now be written as

$$Y = XW\gamma + Xu + r, \tag{2.9}$$

where we note that $E(y) = XW\gamma$ and $\text{Var}(y) = XTX' + \Psi$.

Now, consider a future observable y_{*j} , i.e., a future observation in the j th unit. As before, let the level-1 data corresponding to this observation be denoted as X_{*j} , a $1 \times p$ row vector. To make the analogy with Levy and Perng (1986) explicit, recall that for the general linear model we had the following distributions:

$$p_n(y; X, \beta, \sigma^2) = N(X\beta, \sigma^2 I_n),$$

$$p_m(z; W, \beta, \sigma^2) = N(W\beta, \sigma^2 I_m).$$

Similarly, for the multilevel model in stacked form, we now have

$$p_N(y; X, \gamma, \sigma^2) = N(XW\gamma, XTX' + \Psi),$$

$$p_1(y_{*j}; X_{*j}, \gamma, \sigma^2) = N(X_{*j}W_j\gamma, V_*), \tag{2.10}$$

where $V_* = X_{*j}\tau X_{*j}' + \sigma^2$ and $N = \sum_{j=1}^J n_j$ equals the total number of cases (students) across all the units or groups (schools). The corresponding parameter estimate for the multilevel model is now γ instead of β and may be estimated as follows:

$$\hat{\gamma} = \left(\sum_{j=1}^J W_j' X_j' \hat{V}_j^{-1} X_j W_j \right)^{-1} \sum_{j=1}^J W_j' X_j' \hat{V}_j^{-1} y_j$$

$$\hat{V}_j = \text{v}\hat{\text{a}}\text{r}(y_j) = X_j \hat{\tau} X_j' + \hat{\sigma}^2 I, \tag{2.11}$$

where $\hat{\tau}$ and $\hat{\sigma}^2$ must be estimated iteratively via full or restricted maximum likelihood. The estimate above for the fixed effects γ may be interpreted as a generalized linear model (GLM) estimator. If the dispersion matrix in the multilevel model was diagonal as in the normal linear model, we could directly apply the result of Levy and Perng to obtain the corresponding predictive density for the multilevel model. However, as can readily be seen from Eq. (2.10), the dispersion matrix has a complicated structure. This problem may be solved, however, by making use of some transformations. In order to simplify presentation, we shall first extend Levy and Perng’s assumptions in the case of the normal linear model, and then directly apply this result to the multilevel model. Formally, let us generalize Levy and Perng’s case to that of the following:

$$p_n(y; X, \beta, \sigma^2) = N(X\beta, \sigma^2 \Sigma),$$

$$p_m(z; W, \beta, \sigma^2) = N(W\beta, \sigma^2 \Delta).$$

Assume that Σ and Δ are known matrices of rank n and m , respectively. Let G be an $n \times n$ matrix of rank n such that $\Sigma = G'G$. Similarly, let H be an $m \times m$ matrix of rank m such that $\Delta = H'H$. Let $\tilde{y} = G'^{-1}y$ and let $\tilde{z} = H'^{-1}z$. Similarly, let $\tilde{X} = G'^{-1}X$ and let $\tilde{W} = H'^{-1}W$. Thus, the models above have now been transformed

as follows:

$$p_n(\tilde{y}; \tilde{X}, \beta, \sigma^2) = N(\tilde{X}\beta, \sigma^2 I_n),$$

$$p_m(\tilde{z}; \tilde{W}, \beta, \sigma^2) = N(\tilde{W}\beta, \sigma^2 I_m).$$

Hence, we are back in the original format of Levy and Perng’s problem and may apply the main result to the transformed variables \tilde{y} and \tilde{z} in order to produce the optimal predictive density for \tilde{z} . Note that our corresponding maximum likelihood estimates in the transformed model are now $\hat{\beta} = (\tilde{X}'\tilde{X})^{-1}\tilde{X}'\tilde{y}$ and $\hat{\sigma}^2 = (\tilde{y} - \tilde{X}\hat{\beta})'(\tilde{y} - \tilde{X}\hat{\beta})/n$. However, it can be readily shown that $\hat{\beta} = (\tilde{X}'\tilde{X})^{-1}\tilde{X}'\tilde{y} = (X'\Sigma^{-1}X)^{-1}X'\Sigma^{-1}y$ (Graybill, 1976, p. 207). Following along similar lines as Levy and Perng’s original development, define $\tilde{t} = (\tilde{z} - \tilde{W}\hat{\beta})/(n^{1/2}\hat{\sigma})$ and restrict the set of candidate optimal predictive densities for \tilde{z} to $\Psi = \{s(\tilde{z}; \tilde{W}, \tilde{y}, \tilde{X}) : s(\tilde{z}; \tilde{W}, \tilde{y}, \tilde{X}) = g(\tilde{t}(\tilde{y}, \tilde{z}))\}$. Thus, applying Theorem 1, we have the optimal predictive density g^* for \tilde{z} over the restricted set Ψ as follows:

$$g^*(\tilde{z}; \tilde{W}, \tilde{y}, \tilde{X}) = g^*(\tilde{t}(\tilde{y}, \tilde{z}))$$

$$= st_m(n - p, \tilde{W}\hat{\beta}, n\hat{\sigma}^2 A/(n - p)),$$

where $A = I_m + \tilde{W}(\tilde{X}'\tilde{X})^{-1}\tilde{W}'$, provides the unique minimum of D_{β, σ^2} among all s in Ψ uniformly in β and σ^2 .

Of course, we are interested in the optimal predictive density of z , not \tilde{z} , so we must transform back to the original units. Since we have $\tilde{z} = H'^{-1}z$, this implies that $z = H'\tilde{z}$. Thus, our predictive density for z is as follows:

$$g^*(z; W, y, X) = g^*(t(y, z))$$

$$= st_m(n - p, W\hat{\beta}, n\hat{\sigma}^2 \Delta A/(n - p)), \tag{2.12}$$

where once again we emphasize that the estimates of $\hat{\beta}$ and $\hat{\sigma}^2$ above are not the same as that in the original development. Recall that previously we showed that the dispersion term derived by Levy and Perng for the normal linear model may be written as an adjusted exact predictive interval, where the adjustment factor was $n/(n - p)$. Now, with a little bit of algebra, we demonstrate a similar result for our more general case:

$$\frac{n\hat{\sigma}^2 \Delta A}{n - p} = \frac{n\hat{\sigma}^2 \Delta}{n - p} [I + \tilde{W}(\tilde{X}'\tilde{X})^{-1}\tilde{W}']$$

$$= \frac{n}{n - p} [\hat{\sigma}^2 \Delta + \hat{\sigma}^2 \Delta \tilde{W}(\tilde{X}'\tilde{X})^{-1}\tilde{W}']$$

$$= \frac{n}{n - p} [\text{vâr}(z) + \Delta \text{vâr}(\hat{z})]$$

$$= \frac{n}{n - p} [\text{vâr}(z) + \Delta \text{vâr}(H'^{-1}\hat{z})]$$

$$\begin{aligned}
 &= \frac{n}{n-p} [\text{vâr}(z) + \Delta \Delta^{-1} \text{vâr}(\hat{z})] \\
 &= \frac{n}{n-p} [\text{vâr}(z) + \text{vâr}(\hat{z})].
 \end{aligned}
 \tag{2.13}$$

Thus, as in the general case, this illustrates the dispersion of the predictive density with respect to an adjustment to the margin of error in an exact predictive interval. Let us apply this result to our multilevel model such that we may produce a predictive density for y_{*j} . The only difference in assumptions is that certain components (σ^2 and τ) of the variances for y and y_{*j} are unknown and must be estimated. Directly applying the results from Eqs. (2.12) and (2.13) yields the following predictive density:

$$t(y_{*j}; X_{*j}, y, X, W) = st_1 \left(N - q, X_{*j} W_j \hat{\gamma}, \frac{N}{N - q} B_j \right),
 \tag{2.14}$$

where q is the length of γ and $B_j = \text{vâr}(y_{*j}) + \text{vâr}(\hat{y}_{*j})$ and may be written as follows:

$$\begin{aligned}
 B_j &= X_{*j} \hat{\tau} X_{*j}' + \hat{\sigma}_j^2 + \text{vâr}(X_{*j} W_j \hat{\gamma}) \\
 &= X_{*j} \hat{\tau} X_{*j}' + \hat{\sigma}_j^2 + X_{*j} W_j \text{vâr}(\hat{\gamma}) W_j' X_{*j}' \\
 &= \hat{V}_* + \hat{\sigma}_j^2 + X_{*j} W_j \left(\sum_{j=1}^J W_j' X_j V_j^{-1} X_j W_j \right)^{-1} W_j' X_{*j}'.
 \end{aligned}$$

Comparing the main result here with that from the previous section, the center of the prediction density is now $X_{*j} W_j \hat{\gamma}$ instead of $X_{*j} \hat{\beta}_j$. One may view $\hat{\gamma}$ as analogous to $\hat{\beta}_j$ with respect to the application of the theorem, noting that maximum likelihood is satisfied via generalized least squares in the former and OLS in the latter. As before, we employ this predictive density by taking its expected value, $X_{*j} W_j \hat{\gamma}$, as our point predictor for y_{*j} and use the variance to form our margin of error. Formally, we have the following predictive interval:

$$X_{*j} W_j \hat{\gamma} \pm t_{N-q, 0.975} [N B_j / (N - q)]^{1/2}.
 \tag{2.15}$$

Readers familiar with multilevel models will recognize that this corresponds to employing the prior estimate of β_j , $\hat{\beta}_j^{\text{Prior}} = W_j \hat{\gamma}$, in forming $\hat{y}_{*j} = X_{*j} \hat{\beta}_j^{\text{Prior}}$; Hence the term prior prediction method. Similarly, in the previous section we employed the OLS estimate for β_j and obtained the OLS prediction method. Although the predictive density above corresponding to the prior prediction method is the optimal predictive density in the sense of Levy and Perng (1986), it behooves the researcher to investigate the effect of using the popular multilevel estimate of β_j in place of either the OLS or prior estimate.

2.3. Multilevel prediction method

One of the main results in the multilevel model literature is the shrinkage estimator for β_j , which may be expressed as a weighted combination of the OLS and prior estimate. Intuitively, the higher the reliability of the OLS estimate the larger the weight

attached to the OLS estimate, and vice versa. Formally, the multilevel model estimate $\hat{\beta}_j^*$ may be written as follows:

$$\hat{\beta}_j^* = \delta_j \hat{\beta}_j + (I - \delta_j)W_j\hat{\gamma}, \tag{2.16}$$

where

$$\delta_j = \hat{\tau}[\hat{\tau} + \hat{\sigma}^2(X_j'X_j)^{-1}]^{-1}, \tag{2.17}$$

is the ratio of the parameter variance for β_j (τ) relative to the variance of the OLS estimator for β_j ($\sigma^2(X_j'X_j)^{-1}$) plus this parameter variance matrix. Thus, if the OLS estimate is unreliable, $\hat{\beta}_j^*$ will pull $\hat{\beta}_j$ towards $W_j\hat{\gamma}$, the prior estimate. Once again, the variance components τ and σ^2 must be estimated iteratively and γ is estimated via the generalized least squares as in the previous section. The shrinkage estimator above for β_j which employs the estimator of Eq. (2.11) for γ yields the minimum mean square linear unbiased estimator (MMSLUE) of β_j (Harville, 1976).⁴

The multilevel estimator may also be expressed as $\hat{\beta}_j^* = W_j\hat{\gamma} + \hat{u}_j$, where u_j may be interpreted in the mixed model sense as the random effect of the j th group. With respect to the prediction of y_{*j} , we will now take our predicted value of y_{*j} to be $X_{*j}\hat{\beta}_j^*$, which may also be written as $\hat{y}_{*j} = X_{*j}W_j\hat{\gamma} + X_{*j}\hat{u}_j$. Harville (1976) showed that this may also be written as follows:

$$\hat{y}_{*j} = X_{*j}W_j\hat{\gamma} + \hat{V}_{*j}\hat{V}_j^{-1}(y_j - X_jW_j\hat{\gamma}), \tag{2.18}$$

where $\hat{V}_{*j} = \text{cov}(y_{*j}, y_j) = X_{*j}\hat{\tau}X_j' + \hat{\sigma}^2$. This last representation illustrates our prediction as the conditional expectation of y_{*j} given the data Y . Furthermore, Rao (1973, p. 522) showed that \hat{y}_{*j} is the best predictor of y_{*j} with respect to the minimum mean square error criterion.

Getting back to our predictive density, we center this predictive density around $X_{*j}\hat{\beta}_j^*$ and form its dispersion in a manner analogous to the previous sections, yielding:

$$t(y_{*j}; X_{*j}, y, X, W) = st_1 \left(N - q, X_{*j}\hat{\beta}_j^*, \frac{N}{N - q}C_j \right), \tag{2.19}$$

where q is the length of γ and $C_j = \text{var}(y_{*j}) + \text{var}(\hat{y}_{*j})$ and may be written as

$$C_j = \Omega_{*j} + M_j(\text{var}(\hat{\gamma}))M_j',$$

where $\Omega_{*j} = V_* + V_{*j}V_j^{-1}V_{j*}$ and $M_j = X_{*j}W_j - V_{*j}V_j^{-1}X_jW_j$.⁵ As before, we will form our predictive interval for y_{*j} by centering it around the distribution's mean and using its variance to form the margin of error. Formally, we have the following predictive interval:

$$X_{*j}\hat{\beta}_j^* \pm t_{N-q, 0.975}[NC_j/(N - q)]^{1/2}. \tag{2.20}$$

⁴ One must restrict oneself to the class of unbiased estimators since a MMSLE does not exist for the unknown γ case (Pfefferman, 1984).

⁵ This expression is derived in Liski and Nummi, 1996.

Table 1
Simulation specification

$\sigma^2 = 0.25$
$\tau_{00} = \tau_{11} = 0.125$
$\tau_{01} = 0.03$
$X_{ij} \sim N(0, 1)$
$W_j \sim N(0, 1)$
$J = (25, 50, 100)$
$n = (10, 25, 50)$

We investigate the difference between the OLS, prior, and multilevel prediction methods mentioned above through a simulation study. The design of the simulation study is explained in the next section.

2.4. Simulation study

The design for simulating multilevel data is based on the design of Busing (1993). As a simplification, we consider a simple 2-level multilevel model with one explanatory variable at each level and equal numbers of units in each group. Data are simulated in two stages. At stage one, the level-1 random coefficients are simulated as follows:⁶

$$\beta_{0j} = \gamma_{00} + \gamma_{01}W_j + u_{0j},$$

$$\beta_{1j} = \gamma_{10} + \gamma_{11}W_j + u_{1j}.$$

The γ 's are the fixed effects and are set to a predetermined value; we set them all equal to one as in Busing (1993). The scalar W_j is a standard normal random variable, while the error components, u_{0j} and u_{1j} , have a bivariate normal distribution with mean $(0, 0)$ and a 2×2 covariance matrix τ . We set the two diagonal elements of τ , τ_{00} and τ_{11} , to be 0.125 and the off-diagonal covariance term τ_{01} at 0.03, following one of Busing's major design conditions. This yields in intraclass correlation ρ of 0.33.⁷

At stage two, the level-1 observations are generated according to the following equation:

$$Y_{ij} = \beta_{0j} + \beta_{1j}X_{ij} + \varepsilon_{ij}. \quad (2.21)$$

The level-1 explanatory variable, X_{ij} , is simulated as a standard normal random variable, while the level-1 error ε_{ij} is a normal random variable with mean 0 and variance σ^2 equal to 0.25. The parameter specification for simulating multilevel data is summarized in Table 1.

⁶ There is a slight abuse of notation here. Previously W_j represented a matrix while here it represents a scalar.

⁷ The intraclass correlation is defined as follows: $\rho = \tau_{00}/(\tau_{00} + \sigma^2)$ and thus measures the degree to which units within the same unit are related.

Table 2
Mean fractional coverage for multilevel, prior, and OLS prediction intervals

J	$n = 10$	$n = 25$	$n = 50$
25	0.949, 0.946, 0.984	0.961, 0.951, 0.973	0.952, 0.950, 0.952
50	0.956, 0.952, 0.989	0.956, 0.953, 0.96	0.951, 0.951, 0.955
100	0.960, 0.948, 0.987	0.953, 0.953, 0.966	0.956, 0.947, 0.959

The multilevel data are simulated under a variety of specifications for the number of groups (J) and number of units per group (n). Once again following Busing (1993), the number of groups studied are 25, 50, and 100, while the number of units per group are 10, 25, and 50. Moreover, one additional “future” observation is generated for each of the J groups. Thus, for the $J = 100$, $n = 25$ design specification, 100 additional observations are generated and set aside. These are the observations that will be predicted; they are not used for estimative purposes.

2.5. Prediction results

The adequacy of prediction was checked in two ways: predictive intervals and predictive mean square error (PMSE). The predictive interval method is performed as follows. For each of the future observations to be predicted, a predictive interval is formed from the respective predictive distribution and we check whether or not the observation lies in this interval. Thus, for $J = 50$ we will have a possible range of 0–50 correct predictive intervals. Moreover, to check the variability of such coverage, each of the nine $J \times n$ design conditions are simulated 100 times, each time checking the percent of correct intervals. The data simulations were performed in XLISP-STAT while the multilevel model estimation was performed with TERRACE-TWO.⁸

The predictive interval results are given in Table 2, where we give the mean of the fraction of correct intervals over 100 simulations for each design specification. For instance, the entries in the top left cell shows that for the $J = 25$, $n = 10$ design condition the mean fractional coverage for the multilevel, prior, and OLS predictive intervals were 0.949, 0.946, and 0.984, respectively, over the 100 simulations.

For all three methods, the coverage rate is close to that expected from a theoretical 95% predictive interval. Moreover, there is not much difference between simple OLS and the multilevel intervals. (1) This could be a result of the wideness of our margin of error, and (2) it could be a result of the discreteness of the assessment approach we have employed. In order to get around this problem, we also examine the popular predictive mean square error (PMSE) approach to assessing predictive performance.

For the predictive mean square error (PMSE) approach, we employ the standard technique of taking the average of the sum of the squared errors (SSE) of the observed and predicted values. The predicted values are taken as the expected value of our

⁸ An XLISP-STAT program written by Hilden-Minton (1994, 1995), which incorporates both the EM algorithm and Fisher scoring for parameter estimation.

Table 3
Mean PMSE over 100 simulations for multilevel, prior, and OLS prediction

J	$n = 10$	$n = 25$	$n = 50$
25	0.2958, 0.4914, 0.3056	0.2591, 0.4691, 0.2610	0.2758, 0.489, 0.27666
50	0.2963, 0.4817, 0.3128	0.2644, 0.4785, 0.2674	0.2558, 0.4839, 0.2567
100	0.3005, 0.5048, 0.3188	0.2765, 0.5073, 0.2786	0.2677, 0.5056, 0.2682

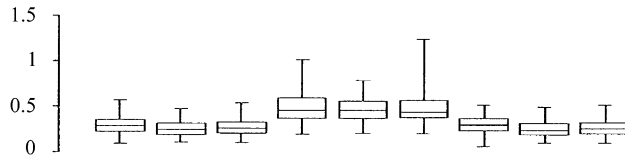


Fig. 1. $J = 25$; $n = 10, 25, 50$ for multilevel, prior and OLS PMSE over 100 simulations.

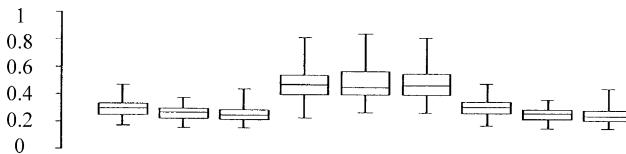


Fig. 2. $J = 50$; $n = 10, 25, 50$ for multilevel, prior and OLS PMSE over 100 simulations.

predictive density, varying according to our estimate of β_j . Once again, for each of the nine $J \times n$ design conditions, we calculate our result 100 times. The results are summarized in Table 3, where each entry is the average of PMSE over 100 simulations. The multilevel method is clearly the best, closely followed by the OLS method. As expected, the discrepancy between the multilevel and OLS method becomes less as n increases. Increasing J should have no effect on the OLS method since this method forms predictions separately for each group. On the other hand, an increase in J should decrease PMSE for the multilevel method since the multilevel method uses all of the data; however, this is not entirely confirmed in these simulations, possibly because the increase in J is not large enough to make a difference. Somewhat of a surprise, the prior prediction method performs the worst of the three methods, increasingly worse as J increases. Although the multilevel prediction rule is superior, the differential gain is not incredibly large and does not increase dramatically as the design tends towards smaller J and n , i.e., specifications where we would expect the multilevel prediction rule to further outperform the other methods.

These results are more clearly illustrated in Figs. 1–3, where we display boxplots of the distribution of PMSE for the three methods over the 100 simulations. Consider Fig. 1 where $J=25$ is fixed. Starting from the left, the first three boxplots correspond to the PMSE for the multilevel prediction method for $n=10, 25$, and 50 , respectively. The next three boxplots correspond to the PMSE for the prior prediction method for $n=10$,

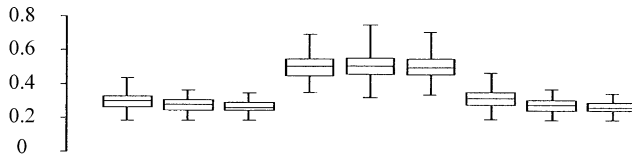


Fig. 3. $J = 100$; $n = 10, 25, 50$ for multilevel, prior and OLS PMSE over 100 simulations.

25, and 50, respectively. And, finally, the last three boxplots correspond to the PMSE for the OLS Prediction Method for $n = 10, 25$, and 50, respectively. Figs. 2 and 3 are arrayed similarly for $J = 50$ and 100, respectively. The effect of group size n is clear as PMSE decreases within each prediction method as n increases. An exception, however, occurs in Fig. 1 for the multilevel method. And, once again, the poor performance of the prior prediction method is apparent as the corresponding boxplots have higher medians in all design specifications. Moreover, as indicated by the boxplots, examining the standard deviation of PMSE over the 100 simulations confirms that the multilevel predictive approach is also the least variable.⁹

These results demonstrate that in spite of the fact that the OLS and prior prediction rules are based on predictive densities which are optimal in the sense of the Kullback–Leibler divergence criterion as employed by Levy and Perng (1986), the predictive performance of the multilevel prediction rule is superior. Part of the reason for this result may arise from the fact that we have restricted the collection of possible density estimators to a subset of prediction densities. This restriction, although possibly useful for theoretical purposes of density estimation, has clearly failed to produce the best density with respect to predicting future observations. Indeed, Levy and Perng (1986) employ this particular restriction in order to demonstrate their result with respect to several other commonly used predictive densities that also belong to this restricted set of predictive densities; whether this set is a reasonably large collection of predictive densities is not their main concern. For the multilevel model at least, our results indicate that this collection needs to be larger.

3. Summary

A predictive density for the multilevel model has been derived in order to facilitate the prediction of future observables in multilevel data. Based upon this predictive density, three prediction methods have been examined: multilevel, prior, and OLS prediction. The OLS prediction method corresponds to deriving a predictive density separately in each group, while the prior prediction method corresponds to deriving a predictive density for the entire model. The multilevel prediction method merely adjusts the prior prediction method by using a well-known result from multilevel model estimation. The adequacy of prediction has been assessed through both predictive intervals and

⁹ For the $J = 100, n = 10$ specification, the standard deviations of SSE for the multilevel, prior, and OLS methods are 3.678, 5.357, and 3.982, respectively.

predictive mean square error (PMSE). Based on simulated multilevel data, the multilevel method is superior. This indicates that for the multilevel model the restriction used by Levy and Perng (1986) in the context of the normal linear model is possibly overly conservative.

The differential gain in prediction for the multilevel method, however, is not incredibly large, nor does this differential gain increase appreciably as the design conditions tend towards smaller J and n , i.e., specifications where we would expect the multilevel method to outperform the OLS method. To be sure, our results might vary if we widen the $J \times n$ space or change other design parameters aside from J and n , e.g., the various parameters of Table 1. In the sequel we explore this enhanced design space and also present a decomposition of prediction error to assess the relative costs of missing data and parameter estimation.

4. For further reading

The following references could also be of interest to the reader. Afshartous, 1997; Afshartous and Hilden-Minton, 1996; Harville, 1985.

Acknowledgements

We greatly appreciate the comments of anonymous referees that have enhanced the quality of the paper.

References

- Afshartous, D., 1997. Prediction in Multilevel Models, Ph.D. Dissertation, UCLA, unpublished.
- Afshartous, D., Hilden-Minton, J., 1996. TERRACE-TWO: an XLISP-STAT software package for estimating multilevel models: user's guide, Technical Report, UCLA Department of Statistics.
- Atchinson, J., 1975. Goodness of prediction fit. *Biometrika* 62, 547–554.
- Bryk, A., Raudenbush, S., 1992. Hierarchical Linear Models. Sage Publications, Newbury Park.
- Busing, F., 1993. Distribution characteristics of variance estimates in two-level models. Technical Report PRM 93-04, Department of Psychometrics and Research Methodology, University of Leiden, Leiden, Netherlands.
- Butler, R.W., 1986. Predictive likelihood with applications. *J. Roy. Statist. Soc. Ser. B* 48, 1–38.
- de Leeuw, J., Kreft, I., 1995. Questioning multilevel models. *J. Educ. Behav. Statist.* 20, 171–189.
- Geisser, S., 1971. The inferential use of predictive distributions. In: Godambe, V.P., Sprott, D.A. (Eds.), *Foundations of Statistical Inference*. Holt, Rhinehart, and Winston, Toronto, pp. 456–469.
- Goldberger, A.S., 1962. Best linear unbiased prediction in the general linear model. *J. Amer. Statist. Assoc.* 57, 369–375.
- Gotway, C., Cressie, N., 1993. Improved multivariate prediction under a general linear model. *J. Multivariate Anal.* 45, 56–72.
- Graybill, F.A., 1976. *Theory and Application of the Linear Model*. Wadsworth, Belmont, CA.
- Guttman, I., Hougaard, P., 1985. Studentization and prediction problems in multivariate multiple regression. *Comm. in Statist., Part A* 14, 1251–1258.
- Harville, D.A., 1976. Extension of the Gauss–Markov theorem to include the estimation of random effects. *Ann. Statist.* 4, 384–396.

- Harville, D.A., 1985. Decomposition of prediction error. *J. Amer. Statist. Assoc.* 80, 132–138.
- Hilden-Minton, J., 1994. TERRACE-TWO: A new xliisp-stat package for multilevel modeling with diagnostics. UCLA Statistics Series. www.stat.ucla.edu.
- Hilden-Minton, J., 1995. Multilevel diagnostics for mixed and hierarchical linear models. Ph.D. Dissertation, UCLA, unpublished.
- Keyes, T., Levy, M., 1996. Goodness of prediction fit for multivariate linear models. *J. Amer. Statist. Assoc.* 91, 191–197.
- Kullback, S., Leibler, R.A., 1951. On information and sufficiency. *Ann. Math. Statist.* 22, 525–540.
- Larimore, W.E., 1983. Predictive inference, sufficiency, entropy and an asymptotic likelihood principle. *Biometrika* 70, 175–182.
- Levy, M.S., Perng, S.K., 1986. An optimal prediction function for the normal linear model. *J. Amer. Statist. Assoc.* 81, 196–198.
- Liski, E.P., Nummi, T., 1996. Prediction in repeated-measures models with engineering applications. *Technometrics* 38, 25–36.
- Pfefferman, D., 1984. On extensions of the Gauss–Markov theorem to the case of stochastic regression coefficients. *J. Roy. Statist. Soc. Ser. B* 46, 139–148.
- Rao, C.R., 1973. *Linear Statistical Inference and its Applications*, 2nd Edition. Wiley, New York.
- Rao, C.R., 1987. Prediction of future observations in growth curve models. *Statist. Sci.* 2, 434–471.
- Robinson, G.K., 1991. That BLUP is a good thing. *Statist. Sci.* 6, 15–51.