

1 **Vet. Res. (2009) 40:30**
 DOI: 10.1051/vetres/2009013

www.vetres.org

2 © INRA, EDP Sciences, 2009

Original article

3 **Use of posterior predictive assessments to evaluate model** 4 **fit in multilevel logistic regression**

5 **Martin J. GREEN^{1,2*}, Graham F. MEDLEY³, William J. BROWNE⁴**

6
 7 ¹ School of Veterinary Medicine and Science, University of Nottingham, Sutton Bonington Campus,
 8 Sutton Bonington, LE12 5RD, United Kingdom

9 ² School of Mathematical Sciences, University of Nottingham, Nottingham, NG7 2RD, United Kingdom

10 ³ Department of Biological Sciences, University of Warwick, Coventry, CV4 7AL, United Kingdom

11 ⁴ Department of Clinical Veterinary Science, University of Bristol, Langford House, Langford, Bristol,
 12 BS40 5DT, United Kingdom

13 (Received 5 November 2008; accepted 24 March 2009)

14 **Abstract** – Assessing the fit of a model is an important final step in any statistical analysis, but this is
 15 not straightforward when complex discrete response models are used. Cross validation and posterior
 16 predictions have been suggested as methods to aid model criticism. In this paper a comparison is
 17 made between four methods of model predictive assessment in the context of a three level logistic
 18 regression model for clinical mastitis in dairy cattle; cross validation, a prediction using the full
 19 posterior predictive distribution and two ‘mixed’ predictive methods that incorporate higher level
 20 random effects simulated from the underlying model distribution. Cross validation is considered a
 21 gold standard method but is computationally intensive and thus a comparison is made between
 22 posterior predictive assessments and cross validation. The analyses revealed that mixed prediction
 23 methods produced results close to cross validation whilst the full posterior predictive assessment gave
 24 predictions that were over-optimistic (closer to the observed disease rates) compared with cross
 25 validation. A mixed prediction method that simulated random effects from both higher levels was
 26 best at identifying the outlying level two (farm-year) units of interest. It is concluded that this mixed
 27 prediction method, simulating random effects from both higher levels, is straightforward and may be
 28 of value in model criticism of multilevel logistic regression, a technique commonly used for animal
 29 health data with a hierarchical structure.

30 **model fit / posterior predictive assessment / mixed predictive assessment / cross validation / Bayesian**
 31 **multilevel model**

32

33 **1. INTRODUCTION**

34 Random effect statistical models are being
 35 increasingly used in veterinary sciences within
 36 both frequentist and Bayesian frameworks.
 37 Models are commonly specified with a binary
 38 outcome to represent, for example, ‘diseased’

or ‘non-diseased’ states and therefore take the 39
 form of multilevel logistic regression [5]. An 40
 important element of constructing and finalising 41
 a statistical model is to critically assess the fit 42
 and performance of the model [8]. However, 43
 model checking with discrete data regressions 44
 is problematic because usual methods, such as 45
 residual plots, have complicated reference dis- 46
 tributions that depend on the parameters in the 47
 model [7, 4]. Thus, these traditional methods 48

* Corresponding author: martin.green@nottingham.ac.uk

1 are considered to be of limited value in discrete
2 outcome, random effects models [2]. It may be
3 because of this that, in the applied literature,
4 particularly when complex discrete response
5 models are specified, attention to model fit is
6 often cursory.

7 In this research, a recently reported method
8 of mixed predictive model assessment [10] is
9 examined and illustrated in the context of an
10 example from veterinary epidemiology. The
11 concept is extended from the two level Poisson
12 regression originally reported, to a three logistic
13 regression setting with the focus of interest on
14 prediction of bovine clinical mastitis on dairy
15 farms in a specific year [6].

16 Posterior prediction is a general term used
17 when data are generated under a proposed
18 model, often so that comparisons can be made
19 between specific features of the observed and
20 generated data [3]. The approach provides a
21 useful means for model assessment and cross
22 validatory posterior predictive distributions are
23 generally considered a 'gold standard' [10,
24 13]. Using cross validation, the data are parti-
25 tioned ' k ' times into subsets and an analysis is
26 initially performed on the 'training' subset.
27 The other 'testing' subset(s) are retained to vali-
28 date the initial analysis by making predictions
29 from the data. Data predictions are compared
30 with the observed data. The procedure is
31 repeated k times and k may equal the total num-
32 ber of data points in the dataset or may repre-
33 sent groups of data within the full set. An
34 important element of cross validation is that
35 predictions made on each subset of testing data
36 are independent of the observed outcome for
37 that subset. The comparisons are used to iden-
38 tify discrepancies between model and data.

39 There is an important difference between
40 conventional residual analysis and cross valida-
41 tion as a means of assessing outlying data
42 regions in the context of model assessment. In
43 conventional residual analysis, all data points
44 are included in the model fit and thus will have
45 a direct effect on model parameters and fitted
46 values, and hence the difference between
47 observed and fitted values. This is not the case
48 with cross validation when the data points or
49 groups have no influence at all on their cross
50 validatory predicted values, because they are

omitted during estimation, and in this respect,
classical residual plots are likely to be over-
optimistic in the assessment of model fit (i.e.
they may not identify all of the true outlying
regions) compared with cross validation. Outly-
ing units from cross validation are those for
which the other units do not provide sufficient
information for the model to fit; outliers from
residual analysis are those for which their
own influence is insufficient to provide a fit.
Therefore, regions of poor fit identified by cross
validation will not necessarily be identified by
residual analysis indicating the importance of
the former method.

A significant disadvantage of cross valida-
tion is that it is computationally intensive and
thus time consuming. A model has to be re-
estimated for each of k subsets and this may
include hundreds or thousands of data points
or regions. If Markov chain Monte Carlo
(MCMC) procedures are being used (as has
been recommended for random effects logistic
regression models [1]), and particularly with
large data sets, the timescale required means
that cross-validation may often become imprac-
tical (depending on the choice of k).

Alternative methods to cross-validatory pre-
dictions have been suggested that have the
advantage of being more straightforward to
compute and less computationally intensive.
Gelman et al. [3] proposed use of the full model
predictive distribution to make predictions on
any required aspect of the data. This method
may be over-optimistic in the context of model
checking (i.e. it may fail to identify true outly-
ing regions) compared to cross-validation
because, as for residual analysis, the prediction
of any data region tends to be strongly influ-
enced by the equivalent observed data for the
region. Marshall and Spiegelhalter [10] pro-
posed a method termed the 'mixed' predictive
check which they have illustrated in the context
of disease mapping, and which appeared to per-
form in a similar manner to cross validation.
The mixed predictive check incorporates simu-
lated random effects, generated from their
underlying distribution which is characterised
from fitting the initial model, rather than the
random effects estimated directly from the data.
Use of the mixed predictive distribution has

1 also been reported in the context of differential
 2 gene expression [9]. In that study, mixed pre-
 3 dictive Markov chain P values were used to
 4 evaluate hierarchical models [3, 10] but compar-
 5 isons were not made between different meth-
 6 ods of posterior predictions as a means to assess
 7 model fit. In this context, Markov chain P val-
 8 ues are an indicator of the probability that a pre-
 9 dicted data region is numerically higher
 10 (or lower) than the observed equivalent. If the
 11 probability is high (typically greater than 95%
 12 or 97.5%) or low (typically less than 5% or
 13 2.5%) then it suggests that the model is per-
 14 forming poorly in the data region.

15 The purpose of this paper is to illustrate and
 16 compare four methods of model predictive
 17 assessment in the context of a multilevel logis-
 18 tic regression model, in which the specific clin-
 19 ical interest was the prediction of disease in a
 20 higher level unit (in this example a farm-year).
 21 The methods are cross validation, a full poster-
 22 ior predictive assessment and two mixed predic-
 23 tive methods based on the approach proposed
 24 by Marshall and Spiegelhalter [10]. An exten-
 25 sion to the concept of the mixed prediction is
 26 described that is generalisable to three level
 27 hierarchical models.

28 **2. MATERIALS AND METHODS**

29 **2.1. The data and initial model**

30 The data for this analysis comprises clinical mas-
 31 titis and farm management information from fifty two
 32 commercial dairy herds, located throughout England
 33 and Wales, with a mean herd size of approximately
 34 150 cows and has been described in detail previously
 35 [6]. Data were collected over a two year period. The
 36 aim of the original research was to investigate the
 37 influence of cow characteristics, farm facilities and
 38 herd management strategies during the dry period,
 39 on the rate of clinical mastitis after calving. Interest
 40 was focussed on identifying determinants for clinical
 41 mastitis occurrence and to assess the extent to which
 42 these determinants could be used to predict the occur-
 43 rence of clinical mastitis in each year on each farm.
 44 The response variable was at the cow level; a cow
 45 either got a case of clinical mastitis (= 1) or not
 46 (= 0) within 30 days of calving and a cow could be
 47 at risk in both years of the study. Predictor variables

were included at the cow, year and farm levels. The
 model hierarchical structure was cows within farm-
 years within farms, and can be summarised as:

$$\begin{aligned}
 CM_{ijk} &\sim \text{Bernoulli}(\pi_{ijk}) \\
 \text{Logit}(\pi_{ijk}) &= \beta_0 + \beta_1 X_{ijk}^{(1)} + \beta_2 X_{jk}^{(2)} \\
 &\quad + \beta_3 X_k^{(3)} + u_{jk} + v_{0k} + v_{1k} P_{ijk} \\
 u_{jk} &\sim N(0, \sigma_u^2), v_k = \begin{pmatrix} v_{0k} \\ v_{1k} \end{pmatrix} \sim \text{MVN}(0, \Omega_v)
 \end{aligned}
 \tag{1}$$

where the subscripts i, j and k denote the three
 model levels, π_{ijk} the fitted probability of clinical
 mastitis (CM) for cow i in year j on farm k , β_0
 the regression intercept, $X_{ijk}^{(1)}$ the vector of covari-
 ates at cow level, β_1 the coefficients for covariates
 $X_{ijk}^{(1)}$, $X_{jk}^{(2)}$ the vector of farm-year level covariates,
 β_2 the coefficients for covariates $X_{jk}^{(2)}$, $X_k^{(3)}$ the vec-
 tor of farm level covariates, β_3 the coefficients for
 covariates $X_k^{(3)}$, P_{ijk} is a covariate (within $X_{ijk}^{(1)}$) that
 identifies cows of parity one (after first calf), u_{jk} is a
 random effect to reflect residual variation between
 years within farms, and v_{0k} and v_{1k} are random
 effects to reflect residual variation between farms,
 and for the difference in rates for parity 1 cows
 between farms respectively.

Model selection was made from a rich dataset of
 more than 350 covariates. Model building has been
 described in detail previously [6] but briefly pro-
 ceeded as follows. Each of the covariates was exam-
 ined individually, within the specified model
 framework, to investigate individual associations with
 clinical mastitis whilst accounting for the data struc-
 ture. Initial covariate assessment was carried out using
 penalised quasi-likelihood for parameter estimation
 (MLwiN, [11]) and final models were selected using
 MCMC for parameter estimation in WinBUGS [12].
 A burn-in of at least 2 000 iterations was used for
 all MCMC runs during which time model conver-
 gence had occurred. Parameter estimates were based
 on a further 8 000 iterations. The final model included
 the following predictor variables; cow parity, cow his-
 toric infection status, whether the farm maintained a
 cow standing time of 30 min after administration of
 treatments at drying off (the end of the previous lacta-
 tion), whether farms reduced the milk yield of high
 yielding cows before drying off, whether cow bedding
 was disinfected during the early dry period, type of
 cow bedding during the late dry period, the time per-
 iod between sequential cleaning out of the calving
 pens, and the time between calving and the cows
 being first milked after calving.

2.2. Predictive assessments

Of particular clinical interest in the research was the prediction of the incidence rate of clinical mastitis (number of cases per cow at risk) for each of the $j = 1 \dots 103$ farm-years and thus the predictions of these rates were used to investigate methods of model assessment. Four methods of predictive assessment were compared; cross validation, a full posterior predictive check and two ‘mixed’ predictive assessments similar to that suggested by Marshall and Spiegelhalter [10]. After final model selection, each method of prediction was incorporated into the MCMC process. At each iteration after model convergence, a prediction was made for the occurrence of mastitis for each individual cow (y_{ijk}) by drawing from the appropriate conditional probability distribution (see below). Similarly, at each iteration, the number of predicted cases of clinical mastitis were summed over all cows in each farm-year and divided by the total cows at risk in each farm-year, to provide an MCMC estimate of the farm-year incidence rate of clinical mastitis. Predictions were made from 8 000 MCMC iterations after model convergence.

To describe the four methods of predictive assessment, we condense the model terms, such that the disease status for each cow (y_{ijk}) is conditional on a set of model fixed effect parameters β , covariates (at various levels) X_{ijk} , and random effects v_k , and u_{jk} :

$$y_{ijk} \sim p(y_{ijk} | \beta, X_{ijk}, V_k, U_{jk})$$

The random effects have parameters represented by σ_u^2 and Ω_v .

$$U_{jk} \sim p(U_{jk} | \sigma_u^2)$$

$$V_k \sim p(V_k | \Omega_v)$$

The four methods of predictive assessment employed were:

- A. Cross validation (“xval”). Each of the 103 farm-years was removed from the analysis in turn and the model fitted to a reduced data set excluding the jk th farm-year (denoted $(-jk)$), from which new model parameters were estimated ($\beta^{(-jk)}$, $v_k^{(-jk)}$, $u_k^{(-jk)}$, $\sigma_u^2^{(-jk)}$, $\Omega_v^{(-jk)}$). A replicate observation for the omitted data, y_{ijk}^{xval} was simulated from the conditional distribution;

$$y_{ijk}^{xval} \sim p(y_{ijk}^{xval} | \beta^{(-jk)}, X_{ijk}, u_{jk}^{xval}, v_k^{xval})$$

$$u_{jk}^{xval} \sim p(u_{jk}^{xval} | \sigma_u^2^{(-jk)})$$

$$v_k^{xval} \sim p(v_k^{xval} | \Omega_v^{(-jk)})$$

- B. Posterior predictive assessment from the full data (“full”). The predictive distribution was conditional on all fixed effect and random effect parameters estimated in the final model and a replicate observation y_{ijk}^{full} generated from the conditional distribution;

$$y_{ijk}^{full} \sim p(y_{ijk}^{full} | \beta, X_{ijk}, v_k, u_{jk})$$

- C. Mixed prediction 1 (“mix1”). This predictive distribution was conditional on the fixed effect parameters and the random effect distributions from which new random effects, u_{jk}^{mix1} and v_k^{mix1} , were simulated to make the prediction. Thus a replicate observation y_{ijk}^{mix1} was generated from the conditional distribution;

$$y_{jk}^{mix1} \sim p(y_{jk}^{mix1} | \beta, X_{ijk}, u_{jk}^{mix1}, v_k^{mix1})$$

$$u_j^{mix1} \sim p(u_j^{mix1} | \sigma_u^2)$$

$$v_k^{mix1} \sim p(v_k^{mix1} | \Omega_v)$$

- D. Mixed prediction 2 (“mix2”). This predictive distribution was conditional on the fixed effect parameters, the random effects distribution at level 2, (from which new random effects, u_{jk}^{mix2} were simulated), and the level 3 random effects from the model, v_k . Thus a replicate observation y_{ijk}^{mix2} was simulated from the conditional distribution;

$$y_{ijk}^{mix2} \sim p(y_{ijk}^{mix2} | \beta, X_{ijk}, u_{jk}^{mix2}, v_k)$$

$$u_{jk}^{mix2} \sim p(u_{jk}^{mix2} | \sigma_u^2)$$

2.3. Comparisons between methods of predictive assessments

In each case, predictions of farm-year incidence rates of clinical mastitis were compared with observed rates. Predictions from cross validation (taken as a gold standard) were also compared to the other methods of prediction to assess which best mimicked this procedure. To assess the degree of discrepancy between observed and predicted farm-year incidence rate of mastitis, the predicted distributions, y_{jk}^{pred} were compared to the observed values using Monte Carlo predictive P values. At each iteration of the MCMC procedure, an indicator variable was set to 1 when $y_{jk}^{pred} > y_{jk}$, to 0.5 if $y_{jk}^{pred} = y_{jk}$ and to 0 if $y_{jk}^{pred} < y_{jk}$; the Monte Carlo P value was estimated as the mean of this indicator variable.

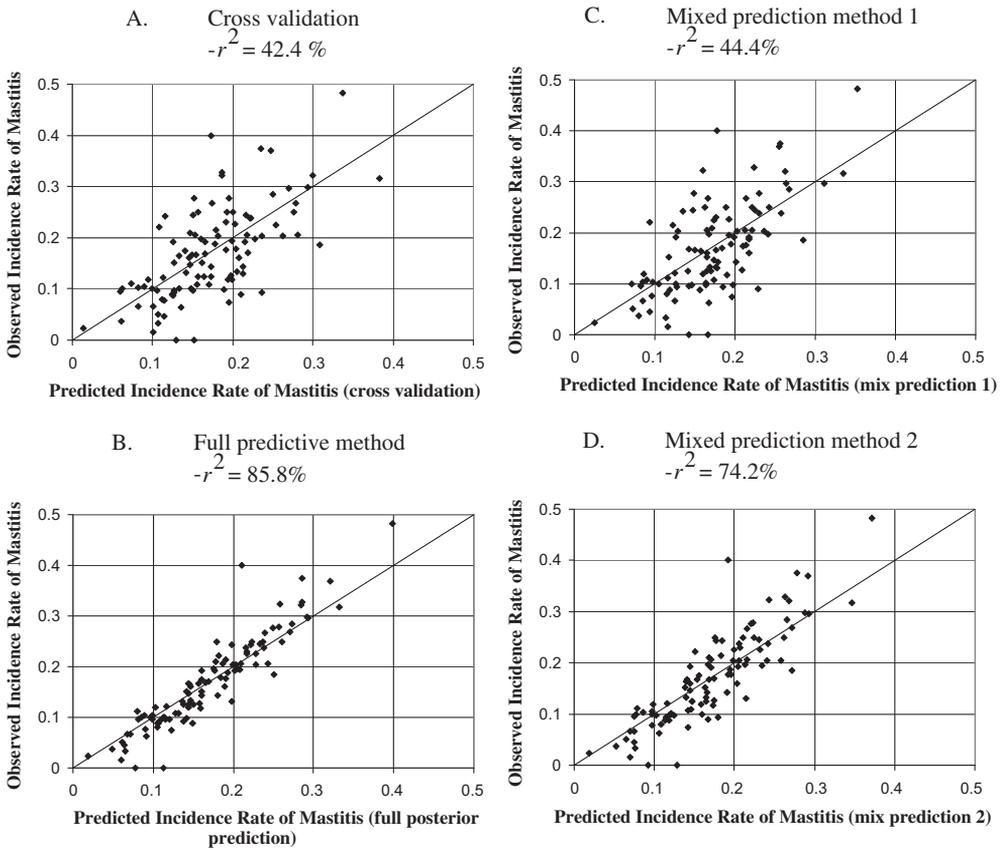


Figure 1. Plots of observed against predicted farm-year incidence rates of clinical mastitis (cases per cow at risk per year).

1 Therefore predictive P values > 0.975 or < 0.025
 2 indicated that the probability of the observed inci-
 3 dence rate of clinical mastitis being within the pre-
 4 dicted distribution was less than 5% and
 5 represented a relatively extreme result.

6 **3. RESULTS**

7 Figure 1 (A–D) illustrates the mean pre-
 8 dicted incidence rate of clinical mastitis for each
 9 method of posterior prediction, plotted against
 10 the observed incidence of clinical mastitis.
 11 The graphs illustrate that the full posterior pre-
 12 dictive method most closely resembled the

observed data and cross validation and the
 “mix1” method displayed considerably more
 variability. The “mix2” method provided an
 intermediate result. Figure 2 illustrates the com-
 parison between mixed and full predictive
 methods and cross validation. Both mixed pre-
 dictive methods yielded better estimates of the
 cross validity prediction than the full post-
 erior predictive method, and the “mix2” method
 produced estimates most similar to cross
 validation.

The median error for each predictive method
 was calculated as the median of the unsigned dif-
 ferences between predicted and cross validity
 farm-year incidence rates of clinical mastitis, as

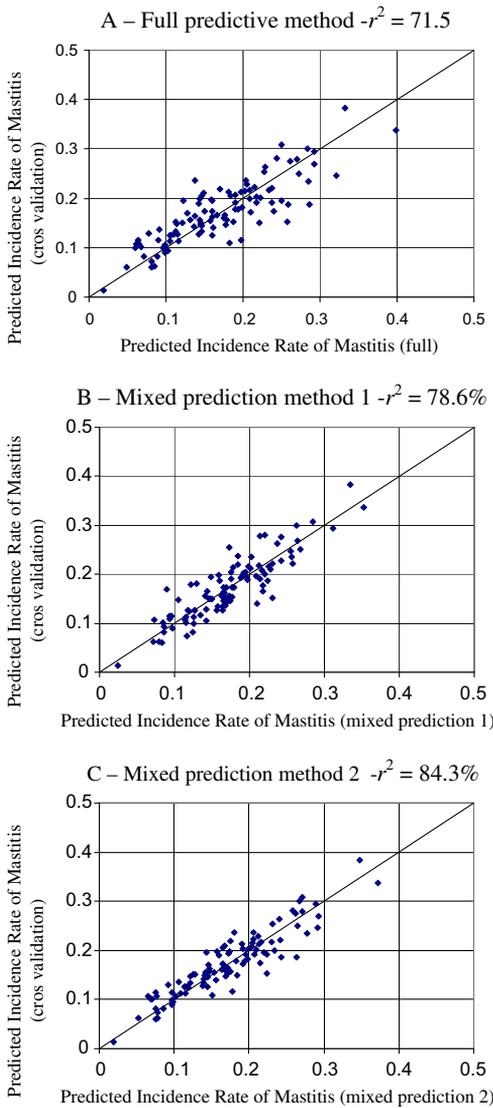


Figure 2. Plots of cross validity predictions of farm-year clinical mastitis incidence against full and mixed predictive methods of farm-year clinical mastitis incidence (cases per cow at risk per year).

1 a percentage of the cross validity farm-year
 2 incidence rate of clinical mastitis. The median
 3 errors were 13.7%, 11.5% and 9.4% for the full
 4 posterior prediction, the mixed prediction 1,
 5 and for mixed prediction 2 respectively.

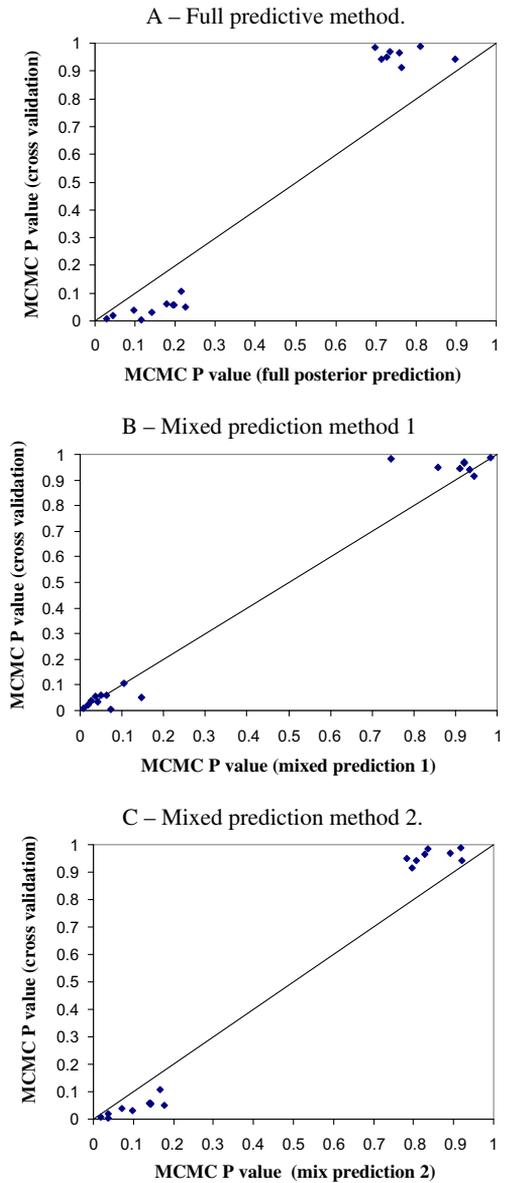


Figure 3. Comparison of MCMC P values from cross validation (for values > 0.80 and < 0.20) and from different methods of predictive assessment for farm-year incidence of clinical mastitis.

Figure 3 illustrates the MCMC P values 6
 obtained from the different predictive methods 7
 to compare with the most extreme P values 8

Table I. Sensitivity and specificity of MCMC *P* values for each prediction method (full = full posterior predictive method, mix 1 and mix 2 = mixed predictive methods 1 and 2 respectively) compared to MCMC *P* values for cross validation, at different *P* value thresholds (as specified).

		Cross validation		Total	Sens (%)	Spec (%)
		0	1			
<i>P</i> value > 0.90 or < 0.10						
full	0	86	14	100	17.6	100.0
	1	0	3	3		
	Total	86	17	103		
mix 1	0	84	3	87	82.4	97.7
	1	2	14	16		
	Total	86	17	103		
mix 2	0	86	10	96	41.2	100.0
	1	0	7	7		
	Total	86	17	103		
<i>P</i> value > 0.95 or < 0.05						
full	0	93	8	101	20.0	100.0
	1	0	2	2		
	Total	93	10	103		
mix 1	0	90	5	95	50.0	96.8
	1	3	5	8		
	Total	93	10	103		
mix 2	0	93	7	100	30.0	100.0
	1	0	3	3		
	Total	93	10	103		
<i>P</i> value > 0.975 or < 0.025						
full	0	98	5	103	0.0	100.0
	1	0	0	0		
	Total	98	5	103		
mix 1	0	98	2	100	60.0	100.0
	1	0	3	3		
	Total	98	5	103		
mix 2	0	98	4	102	20.0	100.0
	1	0	1	1		
	Total	98	5	103		

1 identified with cross validation, these being the
 2 most divergent regions eligible for identification
 3 and further investigation. At large and small *P*
 4 values (*P* < 0.20 or > 0.80) the mixed predic-
 5 tive methods performed more similarly to cross
 6 validation than the full posterior prediction with
 7 the “mix1” method most closely representing
 8 cross validatory MCMC *P* values. This is con-
 9 firmed in Table I that provide the sensitivity and
 10 specificity for each predictive method, taking
 11 cross validation MCMC *P* values as the “gold
 12 standard”, and different *P* value thresholds.

The “mix1” method had the highest sensitivity
 indicating that this method identified the largest
 proportion of “true” extreme values as deter-
 mined by cross validation. The “mix1” method
 identified 82.4% (14 out of 17) of extreme val-
 ues when a threshold of < 0.10 or > 0.90 was
 used and 60% (3 out of 5) of extreme values
 with a threshold set at < 0.025 or > 0.975.

The computing times to complete 10 000
 iterations (using an Intel Centrino 2.0 GHz Pro-
 cessor, 1.5GB RAM) for 103 cross validatory
 predictions and the “mix1” method were 334

1 h and 3.6 h respectively. This did not include the
 2 time required to format the data and set up each
 3 model and this took approximately the same
 4 time per model. Thus it took approximately
 5 103 times longer for the cross validatory predic-
 6 tions than the “mix1” method.

7 4. DISCUSSION

8 Identifying divergent data regions in statisti-
 9 cal modelling is important for two reasons.
 10 Firstly, numerous divergent regions could indi-
 11 cate that underlying statistical assumptions are
 12 incorrect, for example the model does not cap-
 13 ture the true data structure. Secondly, individual
 14 divergent units could represent those that are
 15 fundamentally different from other units in the
 16 dataset after accounting for predictor variables,
 17 and the possible absence of unknown but
 18 important explanatory covariates. In either case,
 19 further investigations would be warranted.
 20 Cross validation provides a useful method of
 21 accurately identifying divergent units in com-
 22 plex statistical models, but faster methods
 23 would be of practical value in model assess-
 24 ment and it was for this reason that the alterna-
 25 tive strategies were investigated in this research.

26 The predictions of clinical mastitis incidence
 27 rates obtained from the different methods show
 28 clear differences in results obtained, as shown
 29 in Figure 1. The full predictive method pro-
 30 vided predicted incidence rates of clinical mas-
 31 titis that most closely resembled the observed
 32 incidence rates, but these appeared to be over-
 33 optimistic in terms of model performance in
 34 comparison to cross validatory predictions. This
 35 is not surprising since the random effects from
 36 the initial model are directly incorporated into
 37 the prediction steps but it does highlight the dif-
 38 ference between this method and cross
 39 validation.

40 For the three level logistic regression models
 41 in this example, the mixed predictive methods
 42 provided a better approximation to cross-
 43 validation than the full posterior predictive
 44 assessment. This is concordant with the first
 45 study that used a mixed prediction for approxi-
 46 mating cross validation in a two level Poisson
 47 model for disease mapping [10]. In the current

48 study using a three level logistic regression 48
 49 model, the “mix2” method provided the closest 49
 50 overall approximation to cross validatory pre- 50
 51 dictions of farm-year incidence of clinical 51
 52 mastitis. However, the “mix1” method per- 52
 53 formed best for the more extreme outlying val- 53
 54 ues identified by cross validation and thus this 54
 55 method was more useful for identifying the 55
 56 most divergent higher level units in these data. 56
 57 The mixed predictive methods look promising 57
 58 as a means of practical model assessment for 58
 59 the relatively common statistical approach of 59
 60 multilevel logistic regression and as such, war- 60
 61 rant further investigations. 61

62 Importantly, the mixed predictive methods 62
 63 take considerably less time to implement 63
 64 (in this example approximately one hundredth 64
 65 of the time of cross validation) and therefore pro- 65
 66 vide a clear advantage in terms of practical use. 66
 67 The “mix2” method is essentially a compromise 67
 68 between the “mix1” method and a full posterior 68
 69 prediction. The method simulates a new random 69
 70 effect at level 2 but uses the estimated random 70
 71 effects from the model at level 3. In the current 71
 72 example there were only two level 2 units for 72
 73 each level 3 unit and it may be that if more level 73
 74 two units existed for each level 3 units, mixed 74
 75 prediction method 2 would tend to become simi- 75
 76 lar to mixed method 1 (the higher level unit hav- 76
 77 ing less influence on the predicted data). 77
 78 Similarly, the relative performance of the two 78
 79 mixed predictive methods may depend on the 79
 80 relative sizes of the higher level variances and 80
 81 more research into the importance of the relative 81
 82 size of higher level variances when using mixed 82
 83 predictive methods would be beneficial. In this 83
 84 example the variance at level two (farm-year) 84
 85 was 0.06 and at level three (farm) was 0.10 85
 86 (for cows greater than parity one) and 0.64 (for 86
 87 cows of parity one). If the level three variances 87
 88 had been very small in comparison to the level 88
 89 2 variance, it is possible that both mixed predic- 89
 90 tive methods used in this study would have 90
 91 yielded similar results. Further investigations 91
 92 of mixed predictive methods using different 92
 93 types of models, numbers of levels, units per 93
 94 level and relative sizes of higher unit variances 94
 95 would be worthwhile. 95

96 From our results, it would appear that, out of 96
 97 the methods examined, the “mix1” method is 97

likely to provide the closest representation of cross validation for potentially divergent data regions in multilevel logistic regression. However, it is important to note that these results apply only to one dataset and whilst in agreement with a previous study [10], need to be viewed with this perspective. It may be possible to generalise this approach to logistic regression and other multilevel models, but more research in this area is required.

Our results indicate that whilst mixed predictions provide a reasonable approximation to cross validation, they do not provide precise replication of the results. Therefore, a pragmatic approach for implementation of mixed predictive assessments may be for an initial highlighting of possible divergent data regions on which to undertake further model checking using cross validation. Thus, instead of undertaking cross validation on all possible regions an intermediate step could be to first use a mixed prediction approach and then to use cross validation for data regions that are potentially divergent based on the mixed prediction. A reduced mixed prediction MCMC *P* value threshold could be used to improve the likelihood that all ‘true’ outliers are identified, possibly the central 80 percentile region and cross validation then carried out on regions that fall outside this interval. This would increase the sensitivity of identifying “true” divergent regions using the mixed methods but would reduce the computing time required compared to using cross validation for all regions.

Assessment of model performance is important and problematic particularly when large datasets and complex model structures are used. Posterior predictions are recognised as a useful method to investigate model fit and more research on mixed posterior predictions may be useful to facilitate straightforward, fast assessments for these types of model.

Acknowledgements. Martin Green is funded by a Wellcome Trust Intermediate Clinical Fellowship.

REFERENCES

[1] Browne W.J., Draper D., A comparison of Bayesian and likelihood-based methods for fitting multilevel models, *Bayesian Analysis* (2006) 1: 473–514.

[2] Dohoo I.R., Martin W., Stryhn H., Veterinary epidemiologic research, Atlantic Veterinary College Inc., Prince Edward Island, Canada, 2003.

[3] Gelman A., Meng X., Stern H., Posterior predictive assessment of model fitness via realized discrepancies, *Statistica Sinica* (1996) 6:733–807.

[4] Gelman A., Goegebeur Y., Tuerlinckx F., van Mechelen I., Diagnostic checks for discrete data regression models using posterior predictive simulations, *Appl. Stat.* (2000) 49:247–268.

[5] Goldstein H., *Multilevel Statistical Models*, London, Edward Arnold, 1995.

[6] Green M.J., Bradley A.J., Medley G.F., Browne W.J., Cow, farm and management factors during the dry period that determine the rate of clinical mastitis after calving, *J. Dairy Sci.* (2007) 90:3764–3776.

[7] Landwehr J.M., Pregibon D., Shoemaker A.C., Graphical methods for assessing logistic regression models (with discussion), *J. Am. Stat. Assoc.* (1984) 79:61–83.

[8] Langford I., Lewis T., Outliers in multilevel data, *J. R. Stat. Soc. Ser. A* (1998) 161:121–160.

[9] Lewin A., Richardson S., Marshall C., Glazier A., Aitman T., Bayesian modelling of differential gene expression, *Biometrics* (2006) 62:1–9.

[10] Marshall E.C., Spiegelhalter D.J., Approximate cross-validators predictive checks in disease mapping, *Stat. Med.* (2003) 22:1649–1660.

[11] Rasbash J., Browne W.J., Healy M., Cameron B., Charlton C., MLwiN Version 2.02, Multilevel Models Project, Centre for Multilevel Modelling, Bristol, UK, 2005.

[12] Spiegelhalter D.J., Thomas A., Best N., WinBUGS Version 1.4.1., Imperial College and MRC, UK, 2004.

[13] Stern H.H., Cressie N., Posterior predictive model checks for disease mapping models, *Stat. Med.* (2000) 19:2377–2397.

45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88