

A simulation study of sample size demonstrated the importance of the number of events per variable to develop prediction models in clustered data

L. Wynants^{a,b,*}, W. Bouwmeester^{c,d}, K.G.M. Moons^c, M. Moerbeek^e, D. Timmerman^{f,g},
S. Van Huffel^{a,b}, B. Van Calster^{f,h}, Y. Vergouwe^h

^aKU Leuven Department of Electrical Engineering-ESAT, STADIUS Center for Dynamical Systems, Signal Processing and Data Analytics, Kasteelpark Arenberg 10, Box 2446, Leuven 3001, Belgium

^bKU Leuven iMinds Medical IT Department, Kasteelpark Arenberg 10, Box 2446, Leuven 3001, Belgium

^cJulius Center for Health Sciences and Primary Care, University Medical Center Utrecht, Heidelberglaan 100, 3584 CX Utrecht, The Netherlands

^dPharmerit B.V., Marten Meesweg 107, Rotterdam 3068 AV, The Netherlands

^eDepartment of Methodology and Statistics, Utrecht University, Padualaan 14, 3584 CH Utrecht, The Netherlands

^fKU Leuven Department of Development and Regeneration, Herestraat 49 Box 7003, Leuven 3000, Belgium

^gDepartment of Obstetrics and Gynaecology, University Hospitals Leuven, Herestraat 49, 3000 Leuven, Belgium

^hCenter for Medical Decision Sciences, Department of Public Health, Erasmus Medical Center, Wytemaweg 80, 3015 CN Rotterdam, The Netherlands

Accepted 9 February 2015; Published online 14 February 2015

Abstract

Objectives: This study aims to investigate the influence of the amount of clustering [intracluster correlation (ICC) = 0%, 5%, or 20%], the number of events per variable (EPV) or candidate predictor (EPV = 5, 10, 20, or 50), and backward variable selection on the performance of prediction models.

Study Design and Setting: Researchers frequently combine data from several centers to develop clinical prediction models. In our simulation study, we developed models from clustered training data using multilevel logistic regression and validated them in external data.

Results: The amount of clustering was not meaningfully associated with the models' predictive performance. The median calibration slope of models built in samples with EPV = 5 and strong clustering (ICC = 20%) was 0.71. With EPV = 5 and ICC = 0%, it was 0.72. A higher EPV related to an increased performance: the calibration slope was 0.85 at EPV = 10 and ICC = 20% and 0.96 at EPV = 50 and ICC = 20%. Variable selection sometimes led to a substantial relative bias in the estimated predictor effects (up to 118% at EPV = 5), but this had little influence on the model's performance in our simulations.

Conclusion: We recommend at least 10 EPV to fit prediction models in clustered data using logistic regression. Up to 50 EPV may be needed when variable selection is performed. © 2015 Elsevier Inc. All rights reserved.

Keywords: Clustered data; Multicenter study; Events per variable; Logistic model; Prediction model; Simulation study

1. Introduction

Clinical prediction models are useful aids to making a diagnosis or prognosis. They are often constructed using

multivariate logistic regression if the health outcome of interest is binary [1]. Researchers proposed to use a sample including at least 10 events per variable (EPV) or candidate predictor for the development of a prediction model [2,3].

Funding: This work was supported by a PhD fellowship from the Flanders' Agency for Innovation by Science and Technology (IWT Vlaanderen) to L.W.; the Research Foundation-Flanders (FWO) (a travel grant to L.W., a fundamental clinical research fellowship to D.T., and project grant G049312 N); the Netherlands Organization for Scientific Research (grant 917.11.383 to Y.V., projects 9120.8004 and 918.10.615 to K.G.M.M.); and ZonMw (grant 17088.25029 to W.B.). This work was further supported by the Research Council KUL [GOA/10/09 MaNet, CoE PFV/10/002 (OPTEC)], the Flemish Government [iMinds Medical

Information Technologies SBO 2014], and the Belgian Federal Science Policy Office [IUAP P7/19/(DYSCO, 'Dynamical systems, control and optimization', 2012–2017)]. The funding sources had no involvement in the study design, the data collection, analysis or interpretation, the writing of the article, or the decision to submit the article for publication.

Conflict of interest: None.

* Corresponding author. Tel.: (32)-16-3-21065; fax: +32 16 3 21970.
E-mail address: Laure.wynants@esat.kuleuven.be (L. Wynants).

What is new?

- The number of events per variable (EPV) can be used to guide sample size decisions in clustered data. There are no existing guidelines on the required number of events relative to the number of (candidate) predictors under study when data are clustered. We recommend to have at least 10 EPV for predefined models, although up to 50 may be needed when performing variable selection.
- Unlike previous studies, this study does not only investigate the influence of sample size and clustering on the bias in regression estimates, but also on the predictive performance of the regression model.
- This study further illustrates that besides the number of EPV, also the total number of observations contributes to the accuracy of regression coefficients and the performance of the prediction model.

The number of events is the number of observations in the smallest outcome category of the binary outcome. The number of variables, henceforward referred to as predictors, should be interpreted more broadly as the number of parameters to be considered. For example, more than one parameter per predictor must be estimated when polynomial terms are used to model a nonlinear effect or when dummy coding is used to model the effect of a qualitative predictor with more than two categories. Hence, in a data set of 500 observations with 50 events in total, only five parameters should be estimated to obtain an EPV of 10. Some researchers have proposed upward [4,5] and downward [6] adjustments for the EPV guideline, stating that the required EPV is influenced by the size of the predictor effects, the correlations among predictors, the prevalence of dichotomous predictors, and the predictor selection strategy. Predictor selection in particular may result in strongly biased estimated regression coefficients in small samples, which decreases the predictive performance of the prediction model when used in individuals other than those from which the model was developed [4].

Multicenter consortia are gaining popularity: recruiting patients from different sites produces a representative sample and reduces recruitment times [7,8]. They yield clustered data sets, that is, data sets with dependent observations because patients from one center may have more in common than patients from different centers [9]. These data sets can be analyzed using multilevel regression (also known as mixed- or random-effects regression, or hierarchical modeling), to build a prediction model [1,10,11].

Multilevel regression enables the incorporation of center-specific intercepts and predictor effects [12]. Simulation studies have shown that the amount of clustering, the number of clusters, and cluster size influence the accuracy and precision of the parameter estimates, especially the random-effect variances and the predictor effects at the cluster level [13–15]. These studies, however, considered neither predictive performance nor the required number of EPV.

Here, we use simulated data and an empirical example of the classification of ovarian tumors to study the effect of the number of EPV for prediction modeling with clustered data when using multilevel logistic regression. We hypothesize that the number of EPV will influence both parameter estimates and the performance of the developed prediction model. The effective sample size [12] is smaller than the number of participants because they have not been sampled independently. The amount of clustering may therefore have a negative effect on the model's performance. However, this impact may be limited because only the estimated predictor effects are used in prediction, and previous research has shown that these are usually estimated with limited bias [13–15].

2. Design of the simulation study

We studied the influence of the number of EPV on the parameter estimates and the predictive performance of the logistic multilevel regression model. We have built models in samples with varying numbers of EPV. These samples were drawn from a source population with a certain degree of clustering. The predictive performance of the models was tested in the source population. In what follows, we describe the design of the simulation study (technical details in Appendix A at www.jclinepi.com, R code in Appendix B at www.jclinepi.com) and the performance measures we evaluated.

2.1. The source populations

We created source populations (of approximately 100,000 observations across 200 clusters) with an intraclass correlation (ICC) of 0%, 5%, or 20%. These values were chosen to reflect situations without clustering, with moderate clustering, and with extreme but realistic clustering in multicenter prediction research [16]. The source populations were generated according to a model with four uncorrelated continuous predictors ($X_1 \sim N(0,1)$, $X_2 \sim N(0,0.6)$, $X_3 \sim N(0,0.4)$, $X_4 \sim N(0,0.2)$), four uncorrelated dichotomous predictors (X_5 to X_8 , with prevalence 0.2, 0.3, 0.3, and 0.4, respectively), and a random intercept of which the variance was determined by the ICC. All regression coefficients were set at 0.8 to achieve a level of discrimination that is common in the applied literature. The overall intercept was set at -2.1 to obtain an outcome event rate of 0.3,

which is common for health outcomes in prediction research. For one population (with ICC = 20%), we introduced a correlation between the random intercepts and X_1 and X_5 , such that predictors were unequally distributed across clusters. The mean of X_1 ranged from -1.13 in the cluster with the lowest random intercept to 1.44 in the cluster with the highest random intercept, whereas the prevalence of X_5 ranged from 12% to 29%.

2.2. Sampling

We sampled data sets from the source populations. The number of EPV was set at 5, 10, 20, or 50. These values reflect popular choices and guidelines for prediction research. In clustered data, the EPV is determined as $EPV = N \times p/k$, where N is the sum of all individual cluster sizes n_j ($j = 1$ to J), that is, the total sample size, p is the sample's event rate, and k is the number of parameters in the model to be estimated, including the random intercept variance. We defined several simulation conditions, grouped into sets of simulations that are characterized by

the parameters that are varied (Table 1): the average number of observations per cluster ($\bar{n}_j = 5, 9, 10, 18, 20, 30, 36, 50$ or 89) in sets 1 to 3 and set 6, the number of clusters ($J = 5, 10, 20, 30,$ or 50) in set 4, and the sample's event rate ($P = 0.05, 0.1, 0.2, 0.3,$ or 0.5) in set 5. The event rate in the sample determines the total sample size when EPV and the number of predictors are fixed because $N = EPV \times k/p$. We drew 500 data sets per simulation condition.

2.3. Model building

Random intercept models were fitted in each sample. They were predefined models using all eight true predictors ($k = 9$, sets 1, 4, 5, 6), full models including the eight true predictors and eight noise variables ($k = 17$, set 2), or reduced models determined through backward variable selection ($\alpha = 0.1$) starting from all 16 candidate predictors ($k = 17$, set 3) (Table 1). The noise variables had the same distributions as the true predictors. To limit the amount of simulations, the presence of noise predictors and variable

Table 1. Simulation conditions

Set	Condition	Population		Sample				Model		EPV
		ICC (%)	Corr (X, u_j)	J	\bar{n}_j	N	p	k	Backward selection	
1	1.1	0	0	30	5	150	0.3	8 + 1	No	5
1	1.2	5	0	30	5	150	0.3	8 + 1	No	5
1	1.3	20	0	30	5	150	0.3	8 + 1	No	5
1	1.4	0	0	30	10	300	0.3	8 + 1	No	10
1	1.5	5	0	30	10	300	0.3	8 + 1	No	10
1	1.6	20	0	30	10	300	0.3	8 + 1	No	10
1	1.7	0	0	30	20	600	0.3	8 + 1	No	20
1	1.8	5	0	30	20	600	0.3	8 + 1	No	20
1	1.9	20	0	30	20	600	0.3	8 + 1	No	20
1	1.10	0	0	30	50	1,500	0.3	8 + 1	No	50
1	1.11	5	0	30	50	1,500	0.3	8 + 1	No	50
1	1.12	20	0	30	50	1,500	0.3	8 + 1	No	50
2	2.1	20	0	30	9	270	0.3	8 + 8 noise+1	No	5
2	2.2	20	0	30	18	540	0.3	8 + 8 noise+1	No	10
2	2.3	20	0	30	36	1,080	0.3	8 + 8 noise+1	No	20
2	2.4	20	0	30	89	2,670	0.3	8 + 8 noise+1	No	50
3	3.1	20	0	30	9	270	0.3	8 + 8 noise+1	Yes	5
3	3.2	20	0	30	18	540	0.3	8 + 8 noise+1	Yes	10
3	3.3	20	0	30	36	1,080	0.3	8 + 8 noise+1	Yes	20
3	3.4	20	0	30	89	2,670	0.3	8 + 8 noise+1	Yes	50
4	4.1	20	0	5	30	150	0.3	8 + 1	No	5
4	4.2	20	0	10	30	300	0.3	8 + 1	No	10
4	4.3	20	0	20	30	600	0.3	8 + 1	No	20
4	4.4	20	0	50	30	1,500	0.3	8 + 1	No	50
5	5.1	20	0	30	30	900	0.05	8 + 1	No	5
5	5.2	20	0	30	30	900	0.1	8 + 1	No	10
5	5.3	20	0	30	30	900	0.2	8 + 1	No	20
5	5.4	20	0	30	30	900	0.5	8 + 1	No	50
6	6.1	20	>0	30	5	150	0.3	8 + 1	No	5
6	6.2	20	>0	30	10	300	0.3	8 + 1	No	10
6	6.3	20	>0	30	20	600	0.3	8 + 1	No	20
6	6.4	20	>0	30	50	1,500	0.3	8 + 1	No	50

Abbreviations: ICC, intraclass correlation; corr(X, u_j), correlation between predictors and the random intercept; J , the number of clusters; \bar{n}_j , the average number of observations per cluster; N , the total number of observations; p , the event rate of the outcome in the sample; k , the number of parameters to be estimated, including the random intercept; EPV, the number of events per variable.

The varying parameters within each set of conditions are indicated in bold.

selection (sets 2 and 3), the effects of the ICC (set 1), and the correlations between random intercepts and predictors (set 6) were only studied in samples with $J = 30$ and $P = 0.3$.

We used the following criteria for model convergence: 10 to 100 iterations to fit the model, a change of less than 10^{-5} in deviances of the models fitted in the last two iterations, and no outlying estimated regression coefficients and standard errors (visual inspection). All models fulfilled these criteria and no samples needed to be deleted.

2.4. Model evaluation

All models were evaluated in terms of the accuracy of estimated regression parameters and the predictive performance.

2.4.1. Bias in the estimated regression coefficients

We compared each estimated regression coefficient $\hat{\beta}$ to β_{sp} , the regression coefficient when the model was built in the source population. The percentage of relative bias in the estimated regression coefficients was defined as $100 \times (\hat{\beta} - \beta_{sp})/\beta_{sp}$. The use of β_{sp} ensures that random error originating from generating a source population is not included in the computation of the bias in the estimated regression coefficients. β_{sp} was between 0.770 and 0.857 for all predictors in all source populations. The relative bias of the estimated random intercept variance was computed analogously.

2.4.2. Predictive performance

The predictive performance of the resulting prediction models was tested in the corresponding source population. We used predictions for the average center, omitting the random intercept, to be able to make predictions in clusters that were not represented in the sample [11,17].

The validated C-index or concordance probability (C) [18] measured the discriminatory performance of the developed model. It was obtained by testing the fitted model in the source population. The calibration slope (b) [1,19] was used to evaluate the accuracy of predicted probabilities in the source population. It is obtained through logistic regression of the event indicator against the linear predictor of the prediction model. If b is smaller than one, there is overfitting, that is, the predicted probabilities are too extreme (too close to zero or one); if b is larger than one, there is underfitting, that is, the predicted probabilities are not extreme enough.

We also computed the within-cluster C-index (C_{within}) and the within-cluster calibration slope (b_{within}) to evaluate the performance at the cluster level instead of the population level. The former is a weighted combination of center-specific C-indices [20], and the latter is estimated using logistic regression with random cluster intercepts and a random cluster calibration slope [10].

The obtained C-indices and calibration slopes were compared with those of a model developed and evaluated in the source population (henceforward C_{sp} and b_{sp}), which serve as upper limits for discriminatory power and calibration in the given population. The relative C-index and the relative calibration slope were computed as $100 \times C/C_{sp}$ and $100 \times b/b_{sp}$ (or $100 \times C_{within}/C_{sp}$ within and $100 \times b_{within}/b_{sp}$ within for the within-cluster measures). Note that b_{sp} will deviate from one because predictions for the average center, omitting the random intercepts, are used. Even if random intercepts are used in prediction, the calibration slope will deviate from one because the cluster-specific random intercepts are shrunk to zero [12].

All simulations and calculations were performed in R version 2.14.0 (Vienna, Austria) [21]. The lmer function from the lme4 package was used to fit multilevel logistic regression models using Laplace approximation [22], and the rms package was used for model evaluation [18].

3. Results of the simulation study

3.1. Data clustering and the number of EPV

The amount of clustering (ICC) did not influence the relative bias of the estimated regression coefficients (Fig. 1A, representing results from simulation set 1). The median bias was close to 0% at each ICC. The bias of the estimated regression coefficients related to the EPV: the interquartile range (IQR) of the relative bias was largest for the lowest EPV values. At ICC = 20%, the median relative bias of $\hat{\beta}_4$ was -10% at EPV = 5 and -2% at EPV = 50. Similar patterns were observed for the other regression coefficients (Supplementary Table 1/Appendix at www.jclinepi.com). The random intercept variances were often underestimated, but a large number of EPV benefited estimation (median relative bias -15.2% at EPV = 5 to -5.2% at EPV = 50 for ICC = 20%).

Model performance also related to the EPV and minimally to the ICC. The within-cluster C-index of the model fitted and evaluated in the source population was 0.78 in each source population (ICC = 0%, 5%, and 20%). The relative within-cluster discrimination of models fitted in the samples was the lowest for the samples with EPV = 5, with median values of around 97% (Fig. 1B). The calibration slope of the models fitted in the source populations was 1.00 at each ICC, and the relative performance of the models fitted in the samples was similar for varying ICCs (Fig. 1C). At ICC = 20%, the median relative within-cluster calibration slope increased from 71.6% at EPV = 5 to 96.6% at EPV = 50. The IQR of the calibration slope also decreased with increasing EPV.

The same patterns were observed for the relative overall C-index and the relative overall calibration slope (Supplementary Fig. 1A and B/Appendix at www.jclinepi.com), but the overall C-indices and calibration slopes of

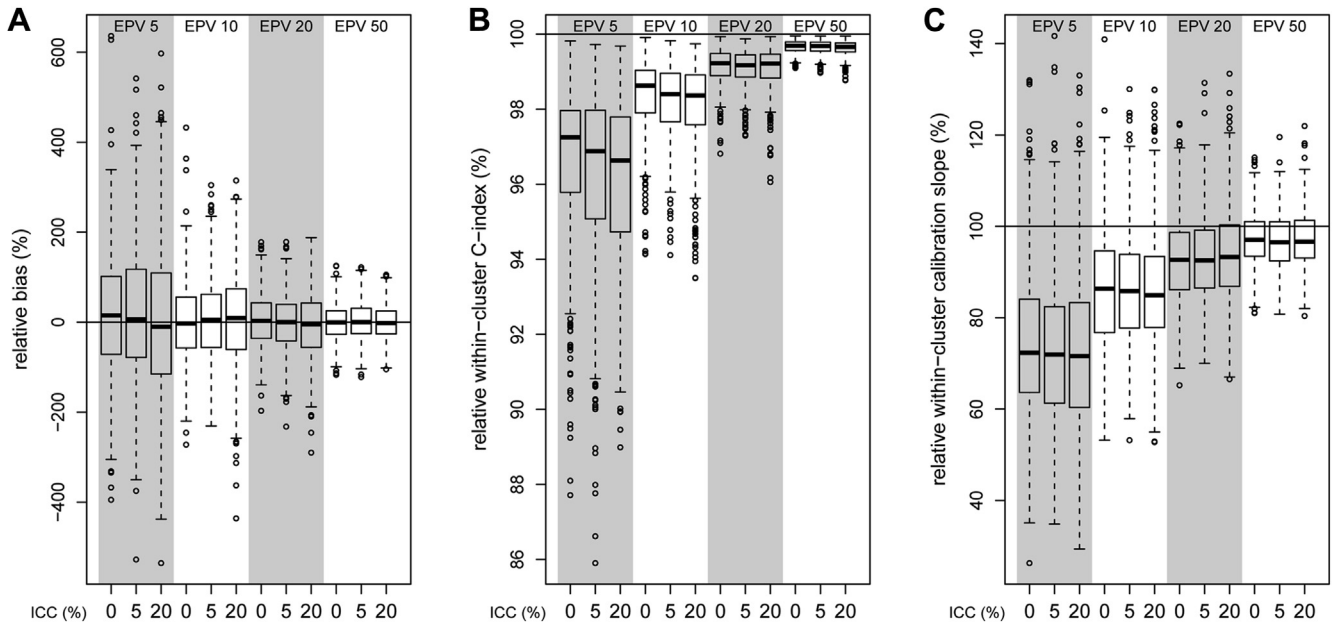


Fig. 1. The relative bias in estimated regression coefficients and the predictive performance in relation to the number of events per variable (EPV) and the amount of clustering (ICC). (A) the relative bias (%) in the estimated regression coefficient $\hat{\beta}_4$; (B) the relative within-cluster discrimination (%), computed as $100 \times (C_{\text{within}}/C_{\text{sp within}})$, where $C_{\text{sp within}}$ is the within-cluster C-index of a model fitted and evaluated in the source population; (C) the relative within-cluster calibration slope (%), defined as $100 \times (b_{\text{within}}/b_{\text{sp within}})$, where $b_{\text{sp within}}$ is the within-cluster calibration slope of a model fitted and evaluated in the source population. The box indicates the interquartile range (IQR), and the fat horizontal line within the box indicates the median. Whiskers extend to the lowest and highest data points still within 1.5 IQR of the box. Outliers beyond these points are represented by dots.

the models fitted in the source populations did decrease with increasing ICC (C_{sp} 0.78, 0.78, and 0.76 and b_{sp} 1.00, 0.97, and 0.88 for ICC = 0%, 5%, and 20%, respectively).

3.2. Variable selection

A high EPV was required to ensure the inclusion of important predictors and to prevent bias in the estimates when using backward variable selection ($\alpha = 0.1$). At EPV = 5, X_4 was selected in only 24% of the samples. To ensure the selection of X_4 in >90% of the samples, 50 EPV were required (Table 2). The median relative bias of $\hat{\beta}_4$ (conditional on selection) was 118% at EPV = 5 and disappeared at EPV = 50 (0.05% relative bias) (Fig. 2A, representing results from simulation conditions 1.3, 1.6, 1.9, and 1.12, and sets 2–3). The selection bias for other predictors was negligible. Predefined models that

included all candidate predictors showed biases similar to predefined models that included only the eight true predictors for the regression coefficients of the eight true predictors. The smaller IQRs in models with 16 variables reflect the larger sample sizes required to obtain the related EPV values with 16 rather than 8 predictors.

Models containing all 16 candidate predictors had a slightly better median relative within-cluster discrimination and a lower median relative within-cluster calibration slope than models after variable selection at EPV = 5 (95.9% vs. 95.3% of $C_{\text{sp within}} = 0.78$ and 67.3% vs. 72.1% of $b_{\text{sp within}} = 1.00$ Fig. 2B and C). This is in accordance with earlier findings [23]. The differences in performance reduced as EPV increased. At EPV = 10, 20, and 50, the predictive performance of the models after variable selection was comparable with the model including only true predictors. The lower relative within-cluster C-index after variable selection at EPV = 5 may be explained by the frequent exclusion of relevant predictors (Table 2). The same pattern was observed for the overall C-index and calibration slope (Supplementary Fig. 1C and D/Appendix at www.jclinepi.com).

Table 2. The selection frequency of predictors

Predictor	EPV = 5 (%)	EPV = 10 (%)	EPV = 20 (%)	EPV = 50 (%)
X_1	100	100	100	100
X_2	88	99	100	100
X_3	60	85	99	100
X_4	24	43	65	93
X_5	64	87	100	100
X_6	72	95	100	100
X_7	76	95	100	100
X_8	79	95	100	100

Abbreviation: EPV, events per variable.

3.3. Sample size

The IQR of the relative bias in the estimated regression coefficients and the model performance related to the total sample size. Samples with EPV = 5 and a total sample size of 900 showed a smaller range of relative bias than samples with EPV = 5 and a total sample size of 150 (Fig. 3A, representing

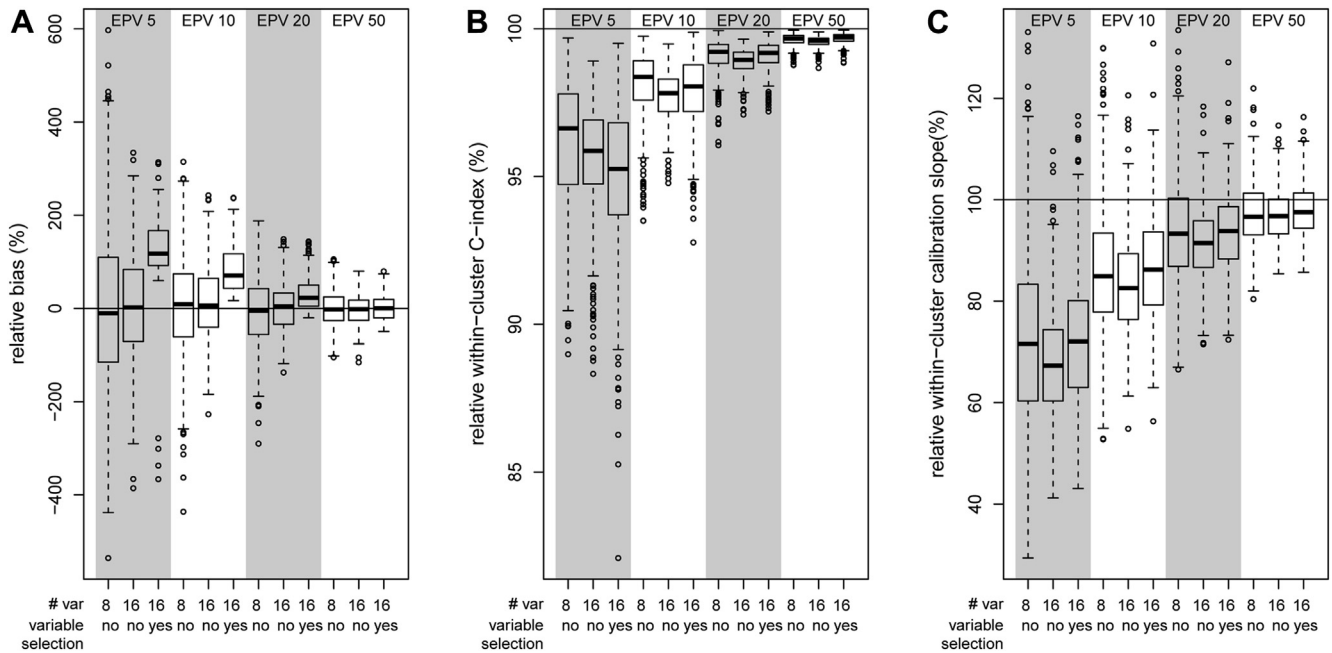


Fig. 2. The relative bias in estimated regression coefficients and the predictive performance in relation to the number of events per variable (EPV) and backward variable selection (ICC = 20%). #var: number of candidate predictors. (A) the relative bias (%) in the estimated regression coefficient $\hat{\beta}_4$; (B) the relative within-cluster discrimination (%), computed as $100 \times (C_{\text{within}}/C_{\text{sp within}})$, where $C_{\text{sp within}}$ is the within-cluster C-index of a model fitted and evaluated in the source population; (C) the relative within-cluster calibration slope (%), defined as $100 \times (b_{\text{within}}/b_{\text{sp within}})$, where $b_{\text{sp within}}$ is the within-cluster calibration slope of a model fitted and evaluated in the source population. The box indicates the interquartile range (IQR), and the fat horizontal line within the box indicates the median. Whiskers extend to the lowest and highest data points still within 1.5 IQR of the box. Outliers beyond these points are represented by dots.

results of simulation conditions 1.3, 1.6, 1.9, 1.12, and sets 4–5). A large total sample size and a large number of clusters reduced the relative bias in the estimated random intercept variance (Supplementary Fig. 2/Appendix at www.jclinepi.com). The relative within-cluster C-indices and calibration slopes of models fitted in the larger samples were higher (Fig. 3B and C). For a given total sample size and number of EPV, samples with many small clusters yielded a predictive performance comparable with samples with a few large clusters. Note that at EPV = 5, the IQRs of the bias in the estimated regression coefficients, the within-cluster C-indices, and the within-cluster calibration slopes were slightly smaller in samples with a few large clusters. The overall performance measures showed similar patterns (Supplementary Fig. 1E and F/Appendix at www.jclinepi.com).

3.4. Random cluster effects correlated with predictors

When the assumption of independence between random intercepts and predictors was violated, $\hat{\beta}_1$ was positively biased (Supplementary Fig. 3A/Appendix at www.jclinepi.com, representing the results of simulation conditions 1.3, 1.6, 1.9, 1.12, and set 6). The same holds for $\hat{\beta}_5$. These regression coefficients accounted for the cluster-level association of X_1 and X_5 with the random intercept, yielding more severely underestimated random intercept variances (Supplementary Fig. 3B/Appendix at www.jclinepi.com). Because $\hat{\beta}_1$ and $\hat{\beta}_5$ contributed to explaining differences

between clusters, the relative overall C-index was increased compared with the situation in which predictors and random intercepts were independent (Supplementary Fig. 3D/Appendix at www.jclinepi.com). This could not be observed for the within-cluster C-index, as $\hat{\beta}_1$ and $\hat{\beta}_5$ did not enable a better discrimination within clusters (Supplementary Fig. 3C/Appendix at www.jclinepi.com). Finally, because $\hat{\beta}_1$ and $\hat{\beta}_5$ were positively biased, overfitting was more problematic (Supplementary Fig. 3E and F/Appendix at www.jclinepi.com).

4. Empirical example

We developed clinical prediction models to diagnose ovarian cancer using data from the International Ovarian Tumor Analysis group [24,25]. We analyzed clinical and ultrasound information on 5,912 patients with ovarian masses from 24 hospitals, collected between 1999 and 2012. The data collected up until 2005 [$n = 1,571$, 9 centers, 409 (26%) malignant tumors] were used for the development of prediction models for tumor malignancy that included random intercepts for hospitals. We drew 100 samples of 409 events and 1,162 nonevents from the training set, with replacement. We considered seven predictors (Table 3). Together with the random intercept variance, this yielded eight parameters to estimate (EPV = 51). We further drew 100 random subsets of 154 patients (40

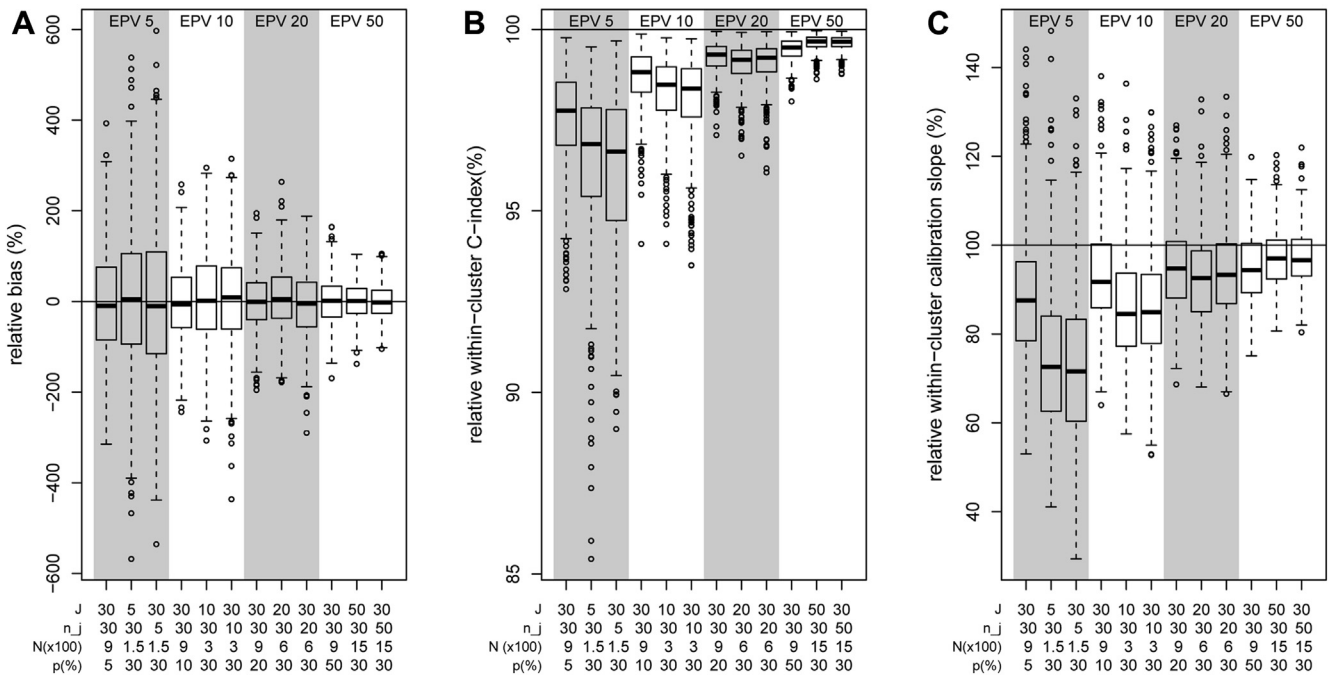


Fig. 3. The relative bias in estimated regression coefficients and the predictive performance in relation to the number of events per variable (EPV) and sample characteristics (ICC = 20%). *J*: the number of clusters, *n_j*: the average number of observations per cluster; *N* (×100): the total number of observations, to be multiplied by 100; *p*: prevalence. (A) the relative bias (%) in the estimated regression coefficient $\hat{\beta}_4$; (B) the relative within-cluster discrimination (%), computed as $100 \times (C_{\text{within}}/C_{\text{sp within}})$, where $C_{\text{sp within}}$ is the within-cluster C-index of a model fitted and evaluated in the source population; (C) the relative within-cluster calibration slope (%), defined as $100 \times (b_{\text{within}}/b_{\text{sp within}})$, where $b_{\text{sp within}}$ is the within-cluster calibration slope of a model fitted and evaluated in the source population. The box indicates the interquartile range (IQR), and the fat horizontal line within the box indicates the median. Whiskers extend to the lowest and highest data points still within 1.5 IQR of the box. Outliers beyond these points are represented by dots.

malignancies) to obtain development samples with EPV = 5. The regression coefficients of the prediction models are shown in Table 3. The most extreme estimates were obtained when variable selection was performed with EPV

= 5. Note that the random intercept variances estimated by the models fitted with EPV = 5 were lower than the estimates obtained with EPV = 51. This is in line with the findings of the simulation study, where the random

Table 3. The fitted models for ovarian tumor diagnosis in 100 samples, based on EPV = 5 vs. EPV = 51, without (model 1 and model 2) and with (model 3 and model 4) backward variable selection ($\alpha = 0.10$)

Predictor	Model 1 (EPV = 51) median (IQR)	Model 2 (EPV = 5) median (IQR)	Model 3 (EPV = 51) median (IQR) (selection frequency)	Model 4 (EPV = 5) median (IQR) (selection frequency)
Age (per 10 yr)	0.38 (0.35 to 0.42)	0.40 (0.27 to 0.55)	0.39 (0.36 to 0.42) (100)	0.48 (0.39 to 0.61) (76)
Maximum diameter of the lesion (per 10 mm)	0.12 (0.11 to 0.13)	0.13 (0.08 to 0.18)	0.12 (0.11 to 0.13) (100)	0.16 (0.11 to 0.19) (82)
Presence of solid tissue in the lesion (yes vs. no)	3.25 (3.13 to 3.40)	3.69 (2.98 to 4.25)	3.23 (3.08 to 3.39) (100)	3.49 (2.91 to 4.08) (91)
Family history of ovarian cancer (yes vs. no)	0.40 (0.20 to 0.61)	0.40 (−0.08 to 1.01)	0.69 (0.62 to 0.82) (30)	2.17 (1.99 to 2.26) (4)
Current use of hormonal therapy (yes vs. no)	−0.35 (−0.46 to −0.24)	−0.47 (−0.87 to −0.07)	−0.43 (−0.54 to −0.38) (60)	−1.24 (−1.39 to −1.15) (15)
Pelvic pain during examination (yes vs. no)	−0.13 (−0.25 to 0.01)	−0.32 (−0.62 to 0.03)	−0.39 (−0.43 to −0.33) (21)	−1.35 (−1.63 to −1.23) (8)
Presence of papillary structures (yes vs. no)	−0.07 (−0.20 to 0.08)	−0.15 (−0.57 to 0.28)	−0.28 (−0.32 to −0.26) (14)	1.03 (−1.14 to 1.17) (24)
Intercept	−6.50 (−6.80 to −6.25)	−7.20 (−8.49 to −5.79)	−6.56 (−6.81 to −6.28)	−6.68 (−7.85 to −5.04)
Random intercept variance	0.59 (0.45 to 0.75)	0.20 (0.00 to 0.67)	0.61 (0.45 to 0.75)	0.19 (0.00 to 0.67)

Abbreviations: EPV, events per variable; IQR, interquartile range.

Table 4. The performance of the fitted models for ovarian tumor diagnosis, based on EPV = 5 vs. EPV = 51 in 100 bootstrap samples, without (model 1 and model 2) and with (model 3 and model 4) backward variable selection ($\alpha = 0.10$)

Performance statistic	Model 1 (EPV = 51)	Model 2 (EPV = 5)	Model 3 (EPV = 51)	Model 4 (EPV = 5)
Within-cluster C-index (IQR)	0.862 (0.860 to 0.864)	0.854 (0.844 to 0.859)	0.862 (0.861 to 0.863)	0.851 (0.842 to 0.859)
Within-cluster calibration slope (IQR)	0.966 (0.934 to 1.011)	0.763 (0.509 to 0.855)	0.970 (0.939 to 1.015)	0.815 (0.735 to 0.933)
Overall C-index (IQR)	0.882 (0.881 to 0.884)	0.873 (0.865 to 0.878)	0.882 (0.881 to 0.883)	0.868 (0.858 to 0.878)
Overall calibration slope (IQR)	0.955 (0.923 to 1.002)	0.751 (0.477 to 0.850)	0.959 (0.925 to 1.00)	0.809 (0.732 to 0.925)

Abbreviations: EPV, events per variable; IQR, interquartile range.

intercept variance was underestimated in samples with EPV = 5 (Supplementary Fig. 2/Appendix at www.jclinepi.com).

The data collected after 2005 [$n = 4,341$, 22 centers, 1,522 (35%) malignancies] were used for the models' validation. The models' performance depended on the number of EPV. The full model fitted with EPV = 51 gave a median validated within-cluster calibration slope of 0.966 (IQR 0.934–1.011), whereas the full model fitted with EPV = 5 gave a much lower median calibration slope of 0.763 (IQR 0.509–0.855). The median validated within-cluster C-indices were 0.862 (IQR 0.860–0.864) and 0.854 (IQR 0.844–0.859), respectively. The effects of EPV were similar when backward variable selection ($\alpha = 0.10$) was used. The median within-cluster calibration slopes were 0.970 (IQR 0.939–1.015) and 0.815 (0.735–0.933) for the models fitted with EPV = 51 and EPV = 5, respectively. The median within-cluster C-indices were 0.862 (IQR 0.861–0.863) and 0.851 (IQR 0.842–0.859), respectively. Compared with a model developed with a high EPV (EPV = 51), a low EPV (EPV = 5) resulted in a similar discriminative ability of the prediction model but more overfitting (Table 4).

5. Discussion

Our simulation research showed that the number of EPV determines the bias in parameter estimates of prediction models developed in clustered data using multilevel logistic regression, as well as the resulting models' predictive performance (discrimination and calibration) in external data. At EPV = 5 and ICC = 20%, predictions were too extreme, but this overfitting disappeared at EPV = 50. Models built on samples with EPV = 50 were also slightly better at discriminating between events and nonevents. This was illustrated in our case study. The amount of clustering was not meaningfully associated with the models' predictive performance. Our simulation results also suggest that larger samples provide better models for a given EPV. This means that nonevents contribute to the stability of the prediction model, provided the number of events is sufficient. When variable selection was performed, a high number of EPV were needed to ensure the inclusion of relevant predictors and reduce estimation bias in the estimated predictor effects.

In accordance with earlier proposals for nonclustered data [2,3], we recommend to use at least 10 EPV to fit a predefined prediction model in clustered data, although

up to 50 EPV may be needed when stepwise variable selection is applied [4]. There were no negative effects of clustering in the simulation conditions we have considered if the number of EPV was sufficiently large. However, it must be noted that it is impossible to obtain an optimal overall calibration slope for a random intercept model, if predictions for the average center (omitting random intercept estimates) are used. Even if the model was fitted on the data of the entire source population, the median overall calibration slope was 0.88 rather than 1 when the ICC was 20%. The within-cluster calibration slope does not suffer from this issue.

Common formulas for power calculations for the design of experiments take into account the ICC because clustering reduces the effective sample size and necessitates larger samples [12]. Nonetheless, the estimation of the regression coefficients used for prediction purposes is little influenced by clustering [13–15], provided that the assumption of independence between predictors and random effects is not violated, and there is no interaction between predictor effects and cluster [9]. Hence, the existing EPV guidelines apply in clustered data if a random cluster intercept is added to the prediction model and the estimation of this additional parameter is accounted for in the EPV calculation.

Multilevel models are very useful tools when analyzing clustered data sets. Random slopes can be used to investigate interactions between predictor effects and cluster, and random intercepts can be used to model differences in outcome prevalence across clusters. It is difficult to reliably estimate the variance of the random intercept, but our results show that estimation improves when the number of clusters increases. This confirms earlier findings [13–15]. It has been recommended to collect data in at least 50 clusters [14], although random-effect variances may still be slightly underestimated with hundreds of clusters [15]. In reality, however, the number of centers in multicenter research is most often smaller than 50 and is determined by weighing benefits, such as the reduction of recruitment times, against practical concerns, to maintain the manageability of the study. Our findings also demonstrate that a high EPV benefits the estimation of the ICC, regardless of the total sample size or the number of clusters. This has not been studied previously.

We believe that the focus on the predictive performance of models, alongside the estimation of regression parameters, is a strength of this study because the purpose of clinical prediction modeling is to make reliable predictions for

new subjects [1]. For the same reason, other evaluation criteria such as confidence interval coverage, type I error rates, and the power of statistical tests are of lesser importance. The uncertainty inherent in model building was acknowledged by studying variable selection. We have used within-cluster performance measures to acknowledge the use of prediction models in separate centers [26]. Our simulation study, like all simulation studies, is restricted by our choice of simulation parameter settings. However, we have chosen practically relevant settings for our simulation parameters, such as the ICC, the number of EPV, and the number of clusters. Furthermore, we did not assume equal numbers of observations in clusters, as this hardly ever occurs in multicenter research [8]. The parameters that were not varied, such as the estimation method (Laplace approximation, [12]) and the variable selection method (backward variable selection [1,18]), were set at generally accepted or recommended choices. A final strength of our study is that we studied the effect of violating the assumption of independence between predictors and random effects [12]. Because two predictors were dependent on the random intercepts, they were also correlated with each other. We did not consider random predictor effects, assuming homogeneity of predictor effects across clusters. Center-level predictors were not included in the study. Reliable estimation of center-level effects will depend on the number of clusters.

In conclusion, this study acknowledges the clustered nature of data sets collected in multicenter research and shows that the number of EPV is useful in guiding sample size decisions when the aim is to develop prediction models using clustered data.

Supplementary data

Supplementary data related to this article can be found at <http://dx.doi.org/10.1016/j.jclinepi.2015.02.002>.

References

- [1] Steyerberg EW. *Clinical prediction models: a practical approach to development, validation, and updating*. New York, NY: Springer US; 2009.
- [2] Peduzzi P, Concato J, Kemper E, Holford TR, Feinstein AR. A simulation study of the number of events per variable in logistic regression analysis. *J Clin Epidemiol* 1996;49:1373–9.
- [3] Harrell FE, Lee KL, Mark DB. Multivariable prognostic models: issues in developing models, evaluating assumptions and adequacy, and measuring and reducing errors. *Stat Med* 1996;15:361–87.
- [4] Steyerberg EW, Eijkemans MJ, Habbema JD. Stepwise selection in small data sets: a simulation study of bias in logistic regression analysis. *J Clin Epidemiol* 1999;52:935–42.
- [5] Courvoisier DS, Combescurie C, Agoritsas T, Gayet-Ageron A, Perneger TV. Performance of logistic regression modeling: beyond the number of events per variable, the role of data structure. *J Clin Epidemiol* 2011;64:993–1000.
- [6] Vittinghoff E, McCulloch CE. Relaxing the rule of ten events per variable in logistic and Cox regression. *Am J Epidemiol* 2007;165:710–8.
- [7] Sprague S, Matta JM, Bhandari M, Dodgin D, Clark CR, Kregor P, et al. Multicenter collaboration in observational research: improving generalizability and efficiency. *J Bone Joint Surg Am* 2009;91(Suppl 3):80–6.
- [8] Senn S. Some controversies in planning and analysing multi-centre trials. *Stat Med* 1998;17:1753–65. discussion 1799–1800.
- [9] Localio AR, Berlin JA, Ten Have TR, Kimmel SE. Adjustments for center in multicenter studies: an overview. *Ann Intern Med* 2001;135:112–23.
- [10] Bouwmeester W, Twisk J, Kappen T, Klei W, Moons K, Vergouwe Y. Prediction models for clustered data: comparison of a random intercept and standard regression model. *BMC Med Res Methodol* 2013;13:19.
- [11] Debray TPA, Moons KG, Ahmed I, Koffijberg H, Riley RD. A framework for developing, implementing, and evaluating clinical prediction models in an individual participant data meta-analysis. *Stat Med* 2013;32:3158–80.
- [12] Snijders TAB, Bosker RJ. In: Snijders TAB, Bosker RJ, editors. *Multilevel analysis: an introduction to basic and advanced multilevel modeling*. 2nd ed. London: Sage; 2012.
- [13] Maas CJM, Hox JJ. Sufficient sample sizes for multilevel modeling. *Methodol Eur J Res Methods Behav Soc Sci* 2005;1(3):86–92.
- [14] Moineddin R, Matheson FI, Glazier RH. A simulation study of sample size for multilevel logistic regression models. *BMC Med Res Methodol* 2007;7:34.
- [15] Paccagnella O. Sample size and accuracy of estimates in multilevel models. *Methodol Eur J Res Methods Behav Soc Sci* 2011;7(3):111–20.
- [16] Adams G, Gulliford MC, Ukoumunne OC, Eldridge S, Chinn S, Campbell MJ. Patterns of intra-cluster correlation from primary care research to inform study design and analysis. *J Clin Epidemiol* 2004;57:785–94.
- [17] Skrondal A, Rabe-Hesketh S. Prediction in multilevel generalized linear models. *J R Stat Soc Ser A Stat Soc* 2009;172(3):659–87.
- [18] Harrell FE. In: Harrell FE Jr, editor. *Regression modeling strategies: with applications to linear models, logistic regression, and survival analysis*. New York, NY: Springer; 2001.
- [19] Cox DR. Two further applications of a model for binary regression. *Biometrika* 1958;45:562–5.
- [20] Van Oirbeek R, Lesaffre E. Assessing the predictive ability of a multilevel binary regression model. *Comput Stat Data Anal* 2012;56(6):1966–80.
- [21] R Development Core Team. *R: a language and environment for statistical computing*. Vienna, Austria: R Foundation for Statistical Computing; 2011. [October 7, 2014]. Available at <http://www.R-project.org>. Accessed March 4, 2015.
- [22] Bates D, Maechler M, Bolker B. *lme4: linear mixed-effects models using Eigen and Eigen++*. R package version 0.999375-42 2011. [October 7, 2014]. Available at <http://CRAN.R-project.org/package=lme4>. Accessed March 4, 2015.
- [23] Steyerberg EW, Eijkemans MJ, Harrell FE Jr, Habbema JD. Prognostic modelling with logistic regression analysis: a comparison of selection and estimation methods in small data sets. *Stat Med* 2000;19:1059–79.
- [24] Timmerman D, Valentin L, Bourne TH, Collins WP, Verrelst H, Vergote I. Terms, definitions and measurements to describe the sonographic features of adnexal tumors: a consensus opinion from the International Ovarian Tumor Analysis (IOTA) group. *Ultrasound Obstetrics Gynecol* 2000;16:500–5.
- [25] Kaijser J, Bourne T, Valentin L, Sayasneh A, Van Holsbeke C, Vergote I, et al. Improving strategies for diagnosing ovarian cancer: a summary of the International Ovarian Tumor Analysis (IOTA) studies. *Ultrasound Obstet Gynecol* 2013;41:9.
- [26] van Klaveren D, Steyerberg E, Perel P, Vergouwe Y. Assessing discriminative ability of risk models in clustered data. *BMC Med Res Methodol* 2014;14:5.