

SIMPLE REPARAMETERIZATION TO IMPROVE CONVERGENCE IN LINEAR MIXED MODELS

Gregor GORJANC¹, Tina FLISAR², Jose Carlos MARTÍNEZ-ÁVILA³, Luis Alberto GARCÍA-CORTÉS³

Received October 08, 2010; accepted December 01, 2010.
Delo je prispelo 08. oktobra 2010, sprejeto 01. decembra 2010.

Simple reparameterization to improve convergence in linear mixed models

Slow convergence and mixing are one of the main problems of Markov chain Monte Carlo (MCMC) algorithms applied to mixed models in animal breeding. Poor convergence is to a large extent caused by high posterior correlation between variance components and solutions for the levels of associated effects. A simple reparameterization of the conventional model for variance component estimation is presented which improves MCMC sampling and provides the same posterior distributions as the conventional model. Reparameterization is based on the rescaling of hierarchical (random) effects in a model, which alleviates posterior correlation. The developed model is compared against the conventional model using several simulated data sets. Results show that presented reparameterization has better behaviour of associated sampling methods and is several times more efficient for the low values of heritability.

Key words: statistics / mixed model / Bayesian analysis / MCMC / reparameterization / convergence

1 INTRODUCTION

Mixed models are abundantly used in the field of animal breeding and genetics with the aim to infer genetic values of animals given some phenotypic and pedigree information (Henderson, 1984). In its simplest form the mixed model can be written as:

$$\mathbf{y} = \mathbf{X}\mathbf{b} + \mathbf{Z}\mathbf{a} + \mathbf{e}, \quad (1)$$

where \mathbf{y} is a vector of phenotypes, \mathbf{b} is a vector of effects

Enostavna reparametrizacija za izboljšanje konvergence linearnih mešanih modelov

Počasna konvergenca je eden največjih problemov uporabe metode Monte Carlo z Markovimi verigami (MCMC) za mešane modele na področju genetike in selekcije domačih živali. Slaba konvergenca je v veliki meri posledica visoke posteriorne korelacije med komponentami variance in rešitvami za ravni pripadajočih vplivov. Predstavljamo enostavno reparametrizacijo običajnega modela, ki izboljša lastnosti metode MCMC in daje enake posteriorne porazdelitve parametrov modela kot standardni pristop. Reparametrizacija temelji na standardizaciji hierarhičnih (naključnih) vplivov v modelu, kar posledično spremeni posteriorne korelacije med parametri. Oba pristopa smo primerjali na večjem setu simuliranih podatkov. Rezultati kažejo, da reparametrizacija vodi do bolj učinkovitih metod MCMC vzorčenja in je nekajkrat bolj učinkovita za analizo lastnosti z nizko heritabiliteto.

Ključne besede: statistika / mešani model / bayesovska analiza / MCMC / reparametrizacija / konvergenca

like sex, breed, age, etc., \mathbf{a} is a vector of individual additive genetic effects and \mathbf{e} residual, $p(\mathbf{e} | \mathbf{a}^2) \sim \mathbf{N}(\mathbf{0}, \mathbf{I}\mathbf{a}^2)$, while \mathbf{X} and \mathbf{Z} are design matrices linking effects to phenotypic records. Pedigree information is included in the model hierarchically with prior distribution of individual additive genetic values, $p(\mathbf{a} | \mathbf{A}, \mathbf{a}^2) \sim \mathbf{N}(\mathbf{0}, \mathbf{A}\mathbf{a}^2)$. Henderson (1972) developed the so called mixed model equations (2) to efficiently obtain joint solutions for \mathbf{b} and \mathbf{a} , where $\mathbf{G} = \mathbf{A}\mathbf{a}^2$ and $\mathbf{R} = \mathbf{I}\mathbf{a}^2$:

¹ Univ. of Ljubljana, Biotechnical Fac., Dept. of Animal Science, Groblje 3, SI-1230 Domžale, Slovenia, Ph.D., e-mail: gregor.gorjanc@bf.uni-lj.si

² The same address as 1

³ Departamento de Mejora Genética, Instituto Nacional de Investigación Agraria, Carretera de La Coruña, km 7, 28040 Madrid, Spain, Ph.D.

$$\begin{matrix} X^T R^{-1} X & X^T R^{-1} Z & \mathbf{b} & \mathbf{I} & \mathbf{I} X^T R^{-1} \mathbf{y} \\ Z^T R^{-1} X & Z^T R^{-1} Z + G^{-1} & \mathbf{a} & \mathbf{I} & Z^T R^{-1} \mathbf{y} \end{matrix} \quad (2)$$

Use of mixed model equations assumes known variance components c_a and c_e . Standard procedure is to estimate these variance components using restricted maximum likelihood method (REML; Patterson and Thompson, 1971) and to use these estimates in mixed model equations (2) ignoring the error of estimation in variance components.

Another approach to statistical inference, Bayesian approach, treats inference of all model parameters jointly. Although conceptually very appealing, Bayesian approach leads to formulas that are computationally intractable. This can be avoided by sampling methods such as Markov chain Monte Carlo (McMC; e.g., Gelman, *et al.*, 2004). Wang *et al.* (1993) showed how McMC methods can be used with linear mixed models in animal breeding applications. In the case of linear mixed models all McMC computations follow from the posterior distribution (3):

$$p(\mathbf{b}, \mathbf{a}, c_a, c_e) \propto |R|^{-1} \exp(-\frac{1}{2}(\mathbf{y} - \mathbf{Xb} - \mathbf{Za})^T R^{-1} (\mathbf{y} - \mathbf{Xb} - \mathbf{Za})) \times \exp(-\frac{1}{2} \mathbf{a}^T G^{-1} \mathbf{a}) \quad (3)$$

where prior distributions for c_a and c_e were assumed uniform (e.g., Gelman *et al.*, 2004). Given that c_a and \mathbf{a} are *a priori* correlated due to the prior definition of \mathbf{a} , the *a posteriori* correlation between them is expected to be high. This leads to high autocorrelation between consecutive samples, making McMC method inefficient. Autocorrelations can be really problematic with low or near zero values for some variance components (e.g. additive genetic variance). This is caused by the shrinkage of \mathbf{a} towards zero and in a next round of sampling variance component will again be close to zero, which can make the sampler stuck for quite some time at the values near zero (Gelman *et al.*, 2004).

Chib and Carlin (1999) proposed block sampling of some parameters in (2) to improve convergence. Autocorrelation has also been alleviated by the use of centered models (Gelfand *et al.*, 1995), parameter expanded models (Liu and Wu, 1999; Gelman *et al.*, 2003; Gelman, 2004) and data augmentation based models (Meng and van Dyk, 1997; van Dyk and Meng, 2001). These methods have been applied both to accelerate the Expectation-Maximization (EM) algorithm and to alleviate the autocorrelation of McMC algorithms. In this work a reparameterization will be employed where additive genetic values will be *a priori* uncorrelated with c_a . This approach will be compared against the conventional model of Wang *et al.* (1993).

2 METHOD

Let us consider a simple animal model $\mathbf{y} = \mathbf{Xb} + \mathbf{Za} + \mathbf{e}$ with the following distributional assumptions:

$$\begin{aligned} p(\mathbf{y} | \mathbf{b}, \mathbf{a}, \sigma^2) &\sim N(\mathbf{Xb} + \mathbf{Za}, \mathbf{I}\sigma^2) \\ p(\mathbf{a} | \mathbf{A}, \sigma^2) &\sim N(\mathbf{o}, \mathbf{A}\sigma^2) \\ p(\sigma^2) &\sim N(\mathbf{o}, \mathbf{I}\sigma^2) \end{aligned} \quad (4)$$

For this particular case and assuming uniform priors for \mathbf{b} and both variance components, $p(\mathbf{b}) \propto const.$, $p(c_a) \propto const.$, and $p(c_e) \propto const.$, the equation (3) becomes:

$$p(\mathbf{b}, \mathbf{a}, c_a, c_e) \propto \exp(-\frac{1}{2}(\mathbf{y} - \mathbf{Xb} - \mathbf{Za})^T R^{-1} (\mathbf{y} - \mathbf{Xb} - \mathbf{Za})) \times \exp(-\frac{1}{2} \mathbf{a}^T G^{-1} \mathbf{a}) \quad (5)$$

where n is the number of records and q the number of animals. Full conditionals of the posterior (5) can be sampled using the coefficient (left hand side) matrix of the mixed model equations (2), sums of squares, normal and scaled central deviates (Wang *et al.*, 1993).

Here another approach is proposed, which alleviates the autocorrelation of samples from (5). It is based on the reparameterization of the model in the terms of a new augmented variable \mathbf{u} , $\mathbf{a} = \mathbf{u}c_a$. Such a model has been already proposed by Foulley and Quaas (1995) in a heterogeneous variance EM-REML context. To simplify

the notation, c_a is used instead of $\mathbf{u}c_a$, but the model is still considered written in terms of c_a . The model is now $\mathbf{y} = \mathbf{Xb} + \mathbf{Zu}c_a + \mathbf{e}$, with the following distributional assumptions:

$$\begin{aligned} p(\mathbf{y} | \mathbf{b}, \mathbf{u}, c_a) &\sim N(\mathbf{Xb} + \mathbf{Zu}c_a, \mathbf{I}K^2) \\ p(\mathbf{u} | \mathbf{A}) &\sim N(\mathbf{0}, \mathbf{A}), \\ p(c_a) &\sim N(\mathbf{o}, \mathbf{I}) \end{aligned} \quad (6)$$

The joint posterior distribution, assuming again uniform priors on \mathbf{b} and both variance components, is:

$$p(\mathbf{b}, \mathbf{u}, c_a) \propto \exp(-\frac{1}{2}(\mathbf{y} - \mathbf{Xb} - \mathbf{Zu}c_a)^T R^{-1} (\mathbf{y} - \mathbf{Xb} - \mathbf{Zu}c_a)) \times \exp(-\frac{1}{2} \mathbf{u}^T \mathbf{A}^{-1} \mathbf{u}) \quad (7)$$

Note that in (7) variance component c_a drops out from the last part, but c_a comes in the sum of squares of residuals. The full conditional distributions for the levels of both \mathbf{b} and \mathbf{a} are univariate normal distributions as in the conventional model, but considering $\mathbf{a} = \mathbf{u}c_a$:

$$p(\mathbf{b}, \mathbf{u}, c_a) \propto \exp(-\frac{1}{2}(\mathbf{y} - \mathbf{Xb} - \mathbf{Zu}c_a)^T R^{-1} (\mathbf{y} - \mathbf{Xb} - \mathbf{Zu}c_a)) \times \exp(-\frac{1}{2} \mathbf{u}^T \mathbf{A}^{-1} \mathbf{u}) \quad (8)$$

$$y \sim N(\mathbf{X}\mathbf{b} + \mathbf{Z}\mathbf{u}, \mathbf{C}) \quad (9)$$

where both \mathbf{C}_i and \mathbf{C}_j are closely related with the conventional mixed model (2) but modified as:

$$\mathbf{C} = \begin{pmatrix} \mathbf{X}^T \mathbf{X} & \mathbf{X}^T \mathbf{Z} \mathbf{a}_a \\ \mathbf{Z}^T \mathbf{X} \mathbf{a}_a & \mathbf{Z}^T \mathbf{Z} + A^{-1} \mathbf{a}_e^2 \end{pmatrix} \quad \mathbf{S} = \begin{pmatrix} \mathbf{X}^T \mathbf{y} \\ \mathbf{Z}^T \mathbf{y} + \mathbf{a}_a \end{pmatrix} \quad (10)$$

The full conditional distribution of \mathbf{a}_e can be sampled from scaled inverted chi-square distribution with $n - 2$ degrees of freedom as in the conventional model:

$$4r^2 | \mathbf{b}, \mathbf{u}, \mathbf{a}_e^2, \mathbf{y} \sim (\mathbf{y} - \mathbf{X}\mathbf{b} - \mathbf{Z}\mathbf{u})^T (\mathbf{y} - \mathbf{X}\mathbf{b} - \mathbf{Z}\mathbf{u}) / (n-2) \quad (11)$$

After some algebra the full conditional of \mathbf{a}_a is

$$P(\mathbf{a}_a | \mathbf{b}, \mathbf{u}, \mathbf{y}) \propto \exp \left\{ -\frac{\mathbf{u}^T \mathbf{U} (\mathbf{y} - \mathbf{X}\mathbf{b})}{\mathbf{u}^T \mathbf{Z}^T \mathbf{Z} \mathbf{u}} \mathbf{a}_a + \frac{\mathbf{a}_a^2}{2 \mathbf{u}^T \mathbf{Z}^T \mathbf{Z} \mathbf{u}} \right\} \quad (12)$$

from which a truncated normal distribution can be recognized when presented in terms of \mathbf{a}_a with mean $\frac{\mathbf{u}^T \mathbf{Z}^T (\mathbf{y} - \mathbf{X}\mathbf{b})}{\mathbf{u}^T \mathbf{Z}^T \mathbf{Z} \mathbf{u}}$, variance $\frac{1}{\mathbf{u}^T \mathbf{Z}^T \mathbf{Z} \mathbf{u}}$, and truncation point at 0:

$$P(\mathbf{a}_a | \mathbf{b}, \mathbf{u}, \mathbf{y}) \sim TN \left(\frac{\mathbf{u}^T \mathbf{Z}^T (\mathbf{y} - \mathbf{X}\mathbf{b})}{\mathbf{u}^T \mathbf{Z}^T \mathbf{Z} \mathbf{u}}, \frac{1}{\mathbf{u}^T \mathbf{Z}^T \mathbf{Z} \mathbf{u}}, 0 \right) \quad (13)$$

When the full conditional distribution of \mathbf{a}_a does not involve the neighbourhood of zero, it is a scaled non-central χ^2 distribution with 1 degree of freedom, with a scale parameter $\frac{1}{\mathbf{u}^T \mathbf{Z}^T \mathbf{Z} \mathbf{u}}$ and noncentrality parameter

$$\lambda = \frac{\mathbf{u}^T \mathbf{Z}^T (\mathbf{y} - \mathbf{X}\mathbf{b}) (\mathbf{y} - \mathbf{X}\mathbf{b})^T \mathbf{Z} \mathbf{u}}{2 \mathbf{u}^T \mathbf{Z}^T \mathbf{Z} \mathbf{u} \mathbf{a}_a} \quad (14)$$

For cases where the posterior distribution of \mathbf{a}_a is close to zero, the Metropolis-Hastings algorithm with positive proposal can be implemented, where the natural logarithm of the conditional density derived from (12) is:

$$\ln p(\mathbf{a}_a | \mathbf{b}, \mathbf{u}, \mathbf{y}) = -\frac{1}{2} \mathbf{a}_a^2 \mathbf{u}^T \mathbf{Z}^T \mathbf{Z} \mathbf{u} + \frac{\mathbf{u}^T \mathbf{Z}^T (\mathbf{y} - \mathbf{X}\mathbf{b})}{\mathbf{u}^T \mathbf{Z}^T \mathbf{Z} \mathbf{u}} \mathbf{a}_a \quad (15)$$

where T represents mean and p variance from (13).

3 APPLICATION

Seven simulated datasets were used to compare the length of burn-in period and Monte Carlo variance of the model $\mathbf{y} = \mathbf{X}\mathbf{b} + \mathbf{Z}\mathbf{u} + \mathbf{e}$ against the conventional sire model $\mathbf{y} = \mathbf{X}\mathbf{b} + \mathbf{Z}\mathbf{s} + \mathbf{e}$. All datasets consisted of 10,000 records, 100 herds (\mathbf{b}) and 500 unrelated sires (\mathbf{s}). Records were randomly assigned to herds and sires, i.e., having on average 100 records per herd and 20 records per sire. True phenotypic variance was 100 and sire variances for each simulated case were: 0.25, 0.5, 1.25, 2, 3.75, 5, and 10.

Markov chain Monte Carlo method was implemented using Gibbs sampler for the full conditional distributions described in (8, 9, and 11), while Metropolis sampler was used for sampling from (15). The length of burn-in period was determined by the use of coupling argument (Johnson, 1996; García-Cortés *et al.*, 1998), where the tolerance of difference between two chains for the sire variance component was set to 10^{-4} . After the burn-in period, chains with 20,000 samples were produced. Monte Carlo error was calculated empirically after 50 replicates for each simulated dataset. Presented

Table 1: Average (\pm standard deviation obtained empirically from 50 replicates) burn-in length by model and true heritability (h^2)
Preglednica 1: Povprečna (\pm standardni odklon, pridobljen empirično iz 50 ponovitev) dolžina ogrevalne faze glede na model in dejansko vrednost heritabilite (h^2)

True h^2	Conventional model	Reparameterized model
0.01	569.6 \pm 266.1	9.8 \pm 6.4
0.02	332.7 \pm 165.2	8.4 \pm 3.9
0.05	173.9 \pm 37.1	7.8 \pm 2.6
0.10	162.4 \pm 41.2	7.8 \pm 2.9
0.15	55.1 \pm 5.8	6.8 \pm 2.4
0.20	42.6 \pm 2.7	7.4 \pm 2.2
0.40	25.2 \pm 3.6	8.4 \pm 3.5

results show the rate of convergence in the terms of burn-in period (Table 1) and after burn-in period (Table 2) for the conventional model (4) and the new reparameterized model (6).

Reparameterization of the model resulted in substantial reduction in burn-in phase of MCMC procedure (Table 1), especially with the low values of heritability. Inspection of trace plots (not shown) showed that in the case of low heritability values for additive genetic variance were very close to zero as well as individual additive genetic values, which is expected. However, conventional model was prone to stuck in that configuration, while reparameterized model more easily explored wider pa-

Table 2: Posterior mean (\pm standard deviation obtained empirically from 50 replicates) for the component of variance between sires by model and true heritability (h^2)

Preglednica 2: Posteriorno povprečje (\pm standardni odklon, pridobljen empirično iz 50 ponovitev) komponente variance med očeti glede na model in dejansko vrednost heritabilite (h^2)

True O^{n^2}	True h^2	Conventional model	Reparameterized model
0.25	0.01	0.39 \pm 0.03	0.38 \pm 0.01
0.50	0.02	0.91 \pm 0.03	0.98 \pm 0.01
1.25	0.05	1.45 \pm 0.02	1.44 \pm 0.01
2.50	0.10	1.69 \pm 0.02	1.69 \pm 0.01
3.75	0.15	4.39 \pm 0.01	4.39 \pm 0.01
5.00	0.20	6.02 \pm 0.01	6.03 \pm 0.01
10.00	0.40	13.03 \pm 0.01	13.05 \pm 0.02

parameter space, which in turn leads to faster convergence to stationary distribution (e.g., Gelman *et al.*, 2004).

Both models gave the same posterior mean on average (Table 2) for variance between sires. Only results for this effect are reported as this is one of the parameters that are hard to accurately estimate in linear mixed models (e.g., Gelman *et al.*, 2004). Posterior means for variance between sires were larger than the true value. This can be attributed to skewed posterior distributions for this effect. Monte Carlo variance obtained after 50 replicates of conventional analysis was sensitive to the value of the true heritability, while this was not the case for reparameterized model. In addition, Monte Carlo variance was higher with conventional model for heritabilities up to 0.1. More stable behaviour of reparameterized model was due to the possibility of easier escape from the neighbourhood of zero value for variance between sires. This means that reparameterized model is of a great value when traits with low heritability are analysed.

4 DISCUSSION

The new data augmentation scheme resulted in an algorithm faster than the conventional Gibbs sampler for linear mixed models. Estimates for variance components do not suffer from getting stuck when visiting values close to zero and then the rate of convergence does not depend on the true value of heritability. When new model was applied to data sets with small heritability, Monte Carlo variance was around five times smaller. Therefore, the new model needs about twenty five times shorter chains to get the same Monte Carlo variance as the conventional model of Wang *et al.* (1993). The new model can be easily implemented in existing programs for the conventional

model - slightly modifying the mixed model equations according to (10) and using the Metropolis algorithm to sample from the full conditional density of cr_a .

Our procedure is very similar to the parameter expanded models presented in (Liu and Wu, 1999; Gelman *et al.*, 2003; Gelman, 2004) among others for both the most frequent EM and Bayesian MCMC. Their approach also standardizes the additive genetic values, but in terms of $\mathbf{a} = \mathbf{ua}$, where a represents an extra augmented variables in the model, while our approach standardizes breeding values with its hyper-parameter, i.e., o_a . The data augmentation scheme presented here can be understood as a particular case of that presented in van Dyk and Meng (2001), which is based on linear transformations of random variables, such as $\mathbf{y} = \mathbf{Xb} + \mathbf{Zp} + \mathbf{e}$, where $\mathbf{p} = \mathbf{Yu} + y$. In our case $\mathbf{Y} = \mathbf{I}cr^{-1}$ and $y = 0$, is the simplest case having a significant reduction of the Monte Carlo variance.

Reparameterized model has been tested with a sire model example. Further research is necessary for animal models or multiple trait models (Henderson, 1984), where the amount of missing information may be higher causing more stringency in standard MCMC samplers. In such cases reparameterization in terms of \mathbf{u} is expected to provide even better results than presented here.

5 CONCLUSION

In summary, reparameterization of hierarchical effects resulted in a feasible Markov chain Monte Carlo algorithm that accelerates the convergence of the conventional sampling methods for Bayesian analysis of linear mixed models. This procedure requires a little programming effort for implementation by researchers who have experience with the conventional sampling methods.

6 REFERENCES

- Chib S., Carlin B.P. 1999. On MCMC sampling in hierarchical longitudinal models. *Statistics and Computing*, 9, 1: 17-26
- van Dyk D.A., Meng X.L. 2001. The art of data augmentation (with discussion). *Journal of Graphical and Computational Statistics*, 10, 1: 1-111
- Foulley J.L., Quaas R.L. 1995. Heterogeneous variances in Gaussian linear mixed models. *Genetics, Selection and Evolution*, 27,3: 211-228
- García-Cortés L.A., Rico M., Groeneveld E. 1998. Using coupling with the Gibbs sampler to assess convergence in animal models. *Journal of Animal Science*, 76, 2: 441-447
- Gelfand A.E., Sahu S.K., Carlin B.P. 1995. Efficient parameterizations for normal linear mixed models. *Biometrika*, 82: 479-488
- Gelman A. 2004. Parameterization and Bayesian model-

- ling. *Journal of American Statistical Association*, 99, 466: 537-545
- Gelman A., Carlin J.B., Stern H.S., Rubin D.B. 2004. *Bayesian data analysis*. Chapman & Hall / CRC, 2 edition
- Gelman A., Huang Z., van Dyk D.A., Boscardin W.J. 2003. Transformed and parameter-expanded Gibbs samplers for multilevel linear and generalized linear model. Technical report, Department of Statistics. Columbia University
- Henderson C.R. 1972. Sire evaluation and genetic trends. In: *Proceedings of the Animal Breeding and Genetics Symposium in Honor of Dr. J.L. Lush*, Champaign, 29 jul. 1972. ASAS, ADSA, PSA: 10-41
- Henderson C.R. 1984. *Applications of Linear Models in Animal Breeding*. Guelph, University of Guelph
- Johnson V.E. 1996. Studying convergence of Markov chain Monte Carlo algorithms using coupled sample paths. *Journal of American Statistical Association*, 91, 433: 154-166
- Liu J.S., Wu Y. 1999. Parameter expansion for data augmentation. *Journal of American Statistical Association*, 94, 448: 1264-1274
- Meng X.L., van Dyk D.A. 1997. The EM algorithm - an old folk-song sung to a fast new tune (with discussion). *Journal of Royal Statistical Society, B Statistical Methodology*, 59, 3: 511-567
- Patterson H.D., Thompson R. 1971. Recovery of inter-block information when block sizes are unequal. *Biometrics*, 58, 8: 545-554
- Wang C.S., Rutledge J.J., Gianola D. 1993. Marginal inferences about variance components in a mixed linear model using Gibbs sampler. *Genetics, Selection and Evolution*, 25, 1: 41-62