

GENETIC EVALUATION AND SELECTION USING GENOTYPIC, PHENOTYPIC AND PEDIGREE INFORMATION

Rohan L. Fernando

Department of Animal Science, Iowa State University, Ames IA 50011, USA

SUMMARY

Genetic evaluation is simplest when genotypes are known at all QTL. Then, the BLUE of genetic differences can be obtained by modeling phenotypic values using a fixed linear model. Genetic evaluation is also straightforward when genotypic information is available only for some of the QTL and no genotypic information is available for the remaining QTL. Then, BLUP of genetic differences can be obtained by modeling phenotypic values using a mixed linear model. It is **expected** that genotypes at markers linked to QTL will be available for a considerable period before genotypes at QTL will be available for genetic evaluation. For this situation, BLUP of genetic differences can be obtained by using a mixed linear model that includes gametic effects at the marked QTL.

Keywords genetic evaluation, genetic markers, marker assisted selection

INTRODUCTION

Genetic variability for most economically important traits is assumed to be due to the segregation of alleles at several quantitative trait loci (QTL). Until recently, little information has been available on individual QTL, and genetic evaluation and selection have been based on phenotypic and pedigree information. Research in molecular genetics during the last decade, however, has resulted in the identification of a large number of polymorphic marker loci and a few candidate genes (Bishop et al. 1994; Archibald et al. 1995; Rothschild et al. 1996). Several experimental procedures and statistical techniques are being used to identify markers that are closely linked to QTL (Weller 1986; Weller and Fernando 1991; Haley and Knott 1992; Hoeschele 1993a; van Arendonk et al. 1994; Xu and Atchley 1995; Stricker et al. 1996; Uimari et al. 1996; Darvasi 1997; Gringola et al. 1997). It is expected that closely linked markers will be used to identify the QTL themselves.

It has been proposed that genotypic information at markers and QTL can be combined with phenotypic information to improve genetic evaluation and selection (Soller 1978; Smith and Webb 1981; Smith and Simpson 1986). Statistical methods for combining genotypic and phenotypic information across large pedigrees are discussed in this paper.

STATISTICAL METHODS

Assumptions and Notation. Consider a trait for which alleles are segregating at N

QTL. For simplicity, we assume two alleles per QTL and additive inheritance. Then, the vector of additive genotypic values for n animals can be written as

$$\mathbf{a} = \sum_{i=1}^N \mathbf{a}_i = \sum_{i=1}^N \mathbf{Q}_i \boldsymbol{\alpha}_i \quad (1)$$

where \mathbf{a}_i is the vector of additive genetic values at locus i , $\boldsymbol{\alpha}_i$ is 2×1 vector of additive effects for the two alleles at locus i , and \mathbf{Q}_i is an $n \times 2$ incidence matrix relating the two additive effects at locus i to the animals; the j th row of \mathbf{Q}_i is $[2, 0]$ if animal j is homozygous for the first allele at loci i , $[1, 1]$ if it is heterozygous, or $[0, 2]$ if it is homozygous for the second allele. For convenience, let the QTL be numbered such that: $1, \dots, k$ are loci for which genotypic information is available for the QTL, $k+1, \dots, k+m$ are loci for which genotypic information is not available for the QTL but is available for linked markers, and $k+m+1, \dots, N$ are loci for which genotypic information is not available for the QTL nor for linked marker loci. Also, let $\mathbf{Q}^K = [\mathbf{Q}_1, \dots, \mathbf{Q}_k]$, $\mathbf{Q}^M = [\mathbf{Q}_{k+1}, \dots, \mathbf{Q}_{k+m}]$ and $\mathbf{Q}^U = [\mathbf{Q}_{k+m+1}, \dots, \mathbf{Q}_N]$, and $\boldsymbol{\alpha}^K$, $\boldsymbol{\alpha}^M$ and $\boldsymbol{\alpha}^U$ denote the vectors of additive effects that correspond to \mathbf{Q}^K , \mathbf{Q}^M and \mathbf{Q}^U , respectively. Phenotypic values on animals can be modeled as

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{Q}\boldsymbol{\alpha} + \mathbf{e} \quad (2)$$

where $\boldsymbol{\beta}$ is a vector of fixed effects, \mathbf{X} is an incidence matrix relating $\boldsymbol{\beta}$ to \mathbf{y} , \mathbf{Z} is an incidence matrix relating \mathbf{a} to \mathbf{y} , $\mathbf{Q} = [\mathbf{Q}^K, \mathbf{Q}^M, \mathbf{Q}^U]$, $\boldsymbol{\alpha}$ is the vector of additive effects for all N loci, and \mathbf{e} is a residual vector with $E(\mathbf{e}) = 0$ and $\text{Var}(\mathbf{e}) = \mathbf{I}\sigma_e^2$.

Genetic evaluation. For most of the current evaluations, $k=0$ and $m=0$. Then, at each locus i , $E(\mathbf{a}_i) = 0$ and $\text{Var}(\mathbf{a}_i) = \mathbf{A}\sigma_{a_i}^2$, where \mathbf{A} is the additive relationship matrix and $\sigma_{a_i}^2$ is the additive genetic variance for the i th QTL. Thus, $E(\mathbf{a}) = 0$ and, if the QTL are in gametic equilibrium, $\text{Var}(\mathbf{a}) = \mathbf{A}\sigma_a^2$, where $\sigma_a^2 = \sum \sigma_{a_i}^2$. For this situation, the best linear unbiased predictor (BLUP) of \mathbf{a} is obtained by solving Henderson's mixed model equations (Henderson 1973).

Genetic evaluation is even more straightforward when $k=N$. Then, the incidence matrix \mathbf{Q} is known, and, conditional on \mathbf{Q} , (2) is a fixed linear model with fixed effects $\boldsymbol{\beta}$ and $\boldsymbol{\alpha}$. For this situation, best linear unbiased estimates (BLUE) of genetic differences between animals can be easily obtained using least squares theory (Searle 1971), provided that genetic differences are not confounded with other fixed effects in the model.

Further, dominance and epistasis can be accommodated by including within locus and between locus allelic interactions. The method of analysis would be the same for linked or unlinked QTL. In principle, the same method can be used whether loci are in equilibrium or not, but including loci that are highly dependent can cause numerical problems due to collinearity of the incidence matrix. With repeated cycles of selection, genetic

evaluations will be biased because selection will be across fixed effects (Henderson 1975). This bias is due to the estimation of fixed genotypic effects, and as information accumulates the bias would disappear. Although genetic evaluation is straightforward for this situation, developing selection rules that maximize long-term genetic progress may not be as straightforward.

Nejati-Javaremi et al. (1997) have considered an alternative approach for the above situation. In their approach, when QTL genotypes are known, genetic evaluations are obtained by solving the following mixed model equations:

$$\begin{bmatrix} \mathbf{X}'\mathbf{X} & \mathbf{X}'\mathbf{Z} \\ \mathbf{T}\mathbf{Z}'\mathbf{X} & \mathbf{T}\mathbf{Z}'\mathbf{Z} + \mathbf{I}\sigma_e^2/\sigma_\alpha^2 \end{bmatrix} \begin{bmatrix} \tilde{\boldsymbol{\beta}} \\ \hat{\mathbf{a}} \end{bmatrix} = \begin{bmatrix} \mathbf{X}'\mathbf{y} \\ \mathbf{T}\mathbf{Z}'\mathbf{y} \end{bmatrix} \quad (3)$$

where \mathbf{T} is the total allelic identity matrix. The total allelic identity between individuals \mathbf{x} and \mathbf{y} was defined as

$$T_{xy} = \frac{\sum_{i=1}^N T_{xy,i}}{N}$$

where $T_{xy,i}$ is the total allelic identity between \mathbf{x} and \mathbf{y} for locus i . Their definition of total allelic identity between \mathbf{x} and \mathbf{y} for locus i is equivalent to twice the probability that a randomly picked allele from locus i in \mathbf{x} is identical in state to a randomly picked allele from locus i in \mathbf{y} .

The above approach has a Bayesian interpretation. To see this, let $\boldsymbol{\beta}$ have a flat prior distribution and let $\boldsymbol{\alpha}$ have a normal prior distribution with null mean and variance $\mathbf{I}\sigma_\alpha^2$. The ratio of variance components $\sigma_e^2/\sigma_\alpha^2$ is assumed to be known. Further, note that the total allelic identity matrix can be expressed as $\mathbf{T} = \frac{\mathbf{Q}\mathbf{Q}'}{2N}$. Now,

$$\begin{aligned} \text{Var}(\mathbf{a}|\text{QTL genotypes}, \sigma_\alpha^2) &= \mathbf{Q}\text{Var}(\boldsymbol{\alpha})\mathbf{Q}' \\ &= \mathbf{Q}\mathbf{Q}'\sigma_\alpha^2 \\ &= \mathbf{T}2N\sigma_\alpha^2 \end{aligned} \quad (4)$$

This conditional variance (4), with $\sigma_\alpha^2 = \frac{\sigma_e^2}{2N}$, is implied by the mixed model equations (3). Thus, from a Bayesian point of view, $\hat{\mathbf{a}}$ from (3) is the posterior mean of \mathbf{a} . It would be more straightforward to first get the posterior mean of $\boldsymbol{\alpha}$ as the solution to:

$$\begin{bmatrix} \mathbf{X}'\mathbf{X} & \mathbf{X}'\mathbf{Z}\mathbf{Q} \\ \mathbf{Q}'\mathbf{Z}'\mathbf{X} & \mathbf{Q}'\mathbf{Z}'\mathbf{Z}\mathbf{Q} + \mathbf{I}\sigma_e^2/\sigma_\alpha^2 \end{bmatrix} \begin{bmatrix} \tilde{\boldsymbol{\beta}} \\ \hat{\boldsymbol{\alpha}} \end{bmatrix} = \begin{bmatrix} \mathbf{X}'\mathbf{y} \\ \mathbf{Q}\mathbf{Z}'\mathbf{y} \end{bmatrix} \quad (5)$$

and then obtain the posterior mean of \mathbf{a} as $\hat{\mathbf{a}} = \mathbf{Q}\hat{\boldsymbol{\alpha}}$.

When $0 < k < N$ and $m = 0$, model (2) becomes

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{Q}^K\boldsymbol{\alpha}^K + \mathbf{Z}\mathbf{Q}^U\boldsymbol{\alpha}^U + \mathbf{e} \quad (6)$$

In this model, β, α^K , and α^U are fixed effects. Given the observed genotypes, Q^K is a fixed matrix but Q^U is random, and therefore, $u = Q^U \alpha^U$ is random. Suppose QTL $k+1, \dots, N$ are in gametic equilibrium with each other and with QTL 1, \dots, k . Then, $E(u) = 0$ and $\text{Var}(u) = A\sigma_U^2$, where $\sigma_U^2 = \sum_{i=k+1}^N \sigma_{a_i}^2$. Thus, genetic evaluations can be obtained from Henderson's mixed model equations for (6), where β and α^K are fixed effects and $u = Q^U \alpha^U$ is random (Kennedy et al. 1992).

Now consider the situation where genotypic information is missing for some of the first k loci, on some of the animals. Suppose, for example, genotype $i < k$ is missing for animal j . Then, the elements of row j of Q_i are not observed. The simplest approach for this situation would be to replace these unobserved row elements in Q_i by their expected values conditional on the observed genotypes at locus i , and proceed with the analysis for complete genotypic information (Clamp et al. 1992; Kinghorn and Kerr 1995). It is not difficult to show that genetic evaluations obtained by this approach are unbiased. Other properties are not known. More accurate evaluations can be obtained by combined segregation and linkage analysis (Stricker et al. 1996; Uimari et al. 1996). However, if genotypes are missing at many of the QTL, this approach may not be computationally feasible.

Next, we consider genetic evaluation when $m > 0$. First, let $m = 1$, and suppose genotypes are observed at one marker that is in equilibrium with QTL $k+1$. Now, model (2) can be written as

$$y = X\beta + ZQ^K \alpha^K + ZQ^M \alpha^M + ZQ^U \alpha^U + e \quad (7)$$

In this model, $\beta, \alpha^K, \alpha^M$, and α^U are fixed effects. But, the matrix Q^M is random, therefore, $w = Q^M \alpha^M$ is random. Because the marker locus and QTL $k+1$ are in equilibrium, $E(w|M) = 0$, where M denotes the marker information. But, $\text{Var}(w|M) \neq A\sigma_{k+1}^2$, in general (Fernando and Grossman 1989). For efficient computations, it is more convenient to work with gametic values v than with additive genetic values w , and the model (7) can be written in terms of u and v as

$$y = X\beta + ZQ^K \alpha^K + ZKv + Zu + e \quad (8)$$

where K is an incidence matrix relating gametes to genotypic values, and $w = Kv$. If parental origin of the marker alleles can be inferred from marker genotypes, recursive formulae given by Fernando and Grossman (1989) can be used to compute $G = \text{Var}(v|M)$ and its inverse. The matrix G is a function of the recombination rate between the marker and the linked QTL and of the additive variance at the QTL. These parameters can be estimated by maximum likelihood (Weller and Fernando 1991; van Arendonk et al. 1994b; Gringola et al. 1997).

Formulae have been developed to accommodate the situation where parental origin of the marker alleles is not known (Wang et al. 1991; Hoeschele 1993b; van Arendonk et al.

Table 1. Pedigree with missing marker genotype for animal 3

Animal	Genotype	Sire	Dam
1	A_1A_2	0	0
2	A_1A_3	0	0
3	missing	1	2
4	A_1A_2	0	0
5	A_1A_2	3	4
6	A_1A_2	3	4

1994a; Wang et al. 1995). Let \mathbf{G}_i be the covariance matrix for gametic values of animals $1, \dots, i$. Then, the tabular method to compute \mathbf{G} can be expressed in matrix notation as

$$\mathbf{G}_{i+1} = \begin{bmatrix} \mathbf{G}_i & \mathbf{G}_i \mathbf{L} \\ \mathbf{L}' \mathbf{G}_i & \mathbf{C}_{i+1} \end{bmatrix} \quad (9)$$

where \mathbf{C}_{i+1} is a 2×2 matrix with $\sigma_{k+1}^2/2$ on the diagonal and $f_{i+1}\sigma_{k+1}^2/2$ on the off-diagonals, f_{i+1} is the conditional inbreeding coefficient for animal $i+1$ given the marker information, and \mathbf{L} is a $2i \times 2$ matrix (Wang et al. 1995). If the marker genotypes of the parents are not missing, each column of \mathbf{L} will contain at most four non-zero elements that are located in rows corresponding to the gametes of the parents of animal $i+1$. This leads to an efficient method to invert \mathbf{G} (van Arendonk et al. 1994a; Wang et al. 1995). Thus, genetic evaluations can be obtained from Henderson's mixed model equations for (8).

When genotypes of parents are missing, however, \mathbf{L} may not have a simple form, and inverting \mathbf{G} may not be simple. For example, consider the pedigree in Table 1, where the genotype M_3 for animal 3 is missing. Then, \mathbf{G} can be computed as

$$\mathbf{G} = \sum_{M_3} \mathbf{G}_{|M_3} \Pr(M_3|M) \quad (10)$$

where the summation is over all possible values for M_3 and $\mathbf{G}_{|M_3}$ is the value of \mathbf{G} computed with animal 3 having genotype M_3 (Hoeschele 1993b; Wang et al. 1995). When the marker genotypes are complete, \mathbf{L} has a simple, known form and $\mathbf{G}_{|M_3}$ can be computed easily using (9). Once \mathbf{G} is computed for a pedigree with missing genotypes using (10), one can solve for \mathbf{L} of any animal. For the pedigree in Table 1, \mathbf{L} for animal

6 is

$$\mathbf{L} = \begin{bmatrix} 0.073 & -0.040 \\ -0.069 & 0.080 \\ 0.073 & -0.043 \\ -0.028 & -0.012 \\ 0.252 & 0.190 \\ 0.247 & 0.078 \\ 0.207 & 0.066 \\ 0.079 & 0.441 \\ 0.255 & -0.062 \\ -0.090 & 0.302 \end{bmatrix} \quad (11)$$

When \mathbf{L} is dense as above, the inverse of \mathbf{G} cannot be obtained efficiently, and further, \mathbf{G}^{-1} will not be sparse. Then, obtaining genetic evaluations using Henderson's mixed model equations may not be feasible. It has been proposed to use a simple \mathbf{L} matrix to approximate \mathbf{G} when marker genotype information is not complete (Hoeschele 1993b; Wang et al. 1995). The consequences of such approximations need to be studied.

Now suppose marker information is available for $m > 1$ unlinked QTL. Genetic evaluations for this situation can be obtained by expanding model (8) to include the **gametic** effects of the m QTL with marker information. This will result in $2nm$ equations for **gametic** effects being included in the mixed model equations. An alternative is to combine the **gametic** effects for the m QTL with marker information and the u QTL without marker information as

$$\mathbf{c} = \sum_{i=k+1}^{k+m} \mathbf{K} \mathbf{v}_i + \mathbf{u} \quad (12)$$

where \mathbf{v}_i is the **gametic** effect for the i th QTL, and the variance of \mathbf{c} given the marker information is

$$\text{Var}(\mathbf{c}|\mathbf{M}) = \mathbf{K} \sum_{i=k+1}^{k+m} \text{Var}(\mathbf{v}_i|\mathbf{M}) \mathbf{K}' + \text{Var}(\mathbf{u}) \quad (13)$$

van Arendonk et al. (1994a) have proposed an algorithm to invert $\text{Var}(\mathbf{c}|\mathbf{M})$. Their algorithm may be useful in very particular situations, but for a general pedigree with marker information on most animals, in comparisons that we have made, their method was less efficient than inversion based on the Cholesky decomposition. Further, although the inverse matrix of each $\text{Var}(\mathbf{v}_i|\mathbf{M})$ and $\text{Var}(\mathbf{u})$ is very sparse, the inverse of $\text{Var}(\mathbf{c}|\mathbf{M})$ is a dense matrix. Thus, the mixed model equations for a model with \mathbf{c} in it will be very dense (van Arendonk et al. 1994a).

Hoeschele (1993) showed that the inverse of $\text{Var} \begin{bmatrix} \mathbf{c}|\mathbf{M} \\ \mathbf{v}|\mathbf{M} \end{bmatrix}$ is sparse and can be computed efficiently ($\mathbf{v} = [\mathbf{v}_{k+1}, \dots, \mathbf{v}_{k+m}]$). Further, she showed that this efficiency is not lost by

eliminating from \mathbf{v} : **gametic** effects of animals that do not have marker genotypes and are not in the pedigree path connecting **genotyped** descendants, or the **gametic** effect of a genotyped animal without offspring, if the **gametic** effect was inherited from an unknown parent or from a parent with both **gametic** effects eliminated from \mathbf{v} . Thus, equations for certain **gametic** effects could be eliminated from the mixed model equations without making them dense. Reduced animal models can also be employed to reduce the number of equations (Cantet and Smith 1991; Saito and Iwaisaki 1996).

When some of the QTL with marker information are linked, Goddard (1992) proposed to model each \mathbf{v}_i conditional on markers flanking the i th QTL. He gave formulae to compute $\text{Var}(\mathbf{v}_i|\mathbf{M})$ and its inverse. These results are also based on a simple, known form for \mathbf{L} in (9). When the linkage phase between the flanking markers is known, \mathbf{L} has a simple, known form. Even with complete marker genotype information, the linkage phase of some parents may not be known and approximations will have to be employed.

Genetic evaluations can also be undertaken by computing the conditional mean of \mathbf{y} given the genotypic, phenotypic, and pedigree information. This would require specifying completely the distribution of the random effects in the model and computing the mean by methods related to complex segregation and linkage analysis. Gibbs sampling can be used to reduce the computational burden. It has been shown, however, that the Gibbs sampler may not converge when more than two alleles are segregating at a locus.

ACKNOWLEDGMENTS

Discussions with Daniel Gianola, Liviu Totir and Bruce Southey are gratefully acknowledged. Research support was provided by award no. 94-37205-1307 of the National Research Initiative Competitive Grants Program of the USDA.

REFERENCES

- Archibald, A. L., Haley, C. S., Brown, J. F., Couperwhite, S., McQueen, H. A., Nicholson, D., Coppieters, W., Van de Weghe, A., Stratil, A., Wintero, A. K. and et al. (1995) **Mamm. Genome** 6:157–175.
- Bishop, M. D., Kappes, S. M., Keele, J. W., Stone, R. T., Sunden, S. L. F., Hawkins, G. A., Toldo, S. S., Fries, R., Grosz, M. D., Yoo, J. and Beattie, C. W. (1994) **Genetics** 136:619–639.
- Cantet, R. J. C. and Smith, C. (1991) **Genet. Sel. Evol.** 23:221–233.
- Clamp, P. A., Beever, J. E., Fernando, R. L., McLaren, D. G. and Schook, L. B. (1992) **J. Anim. Sci.** 70:2695–2706.
- Darvasi, A. (1997) **Mamm. Genome** 8:163–167.
- Fernando, R. L. and Grossman, M. (1989) **Genet. Sel. Evol.** 21:467–477.

- Goddard, M. E. (1992) *Theor. Appl. Genet.* 83:878–886.
- Gringola, F. E., Hoeschele, I. and Tier, B. (1997) *GSE* 28:479–490.
- Haley, C. S. and Knott, S. A. (1992) *Heredity* 69:315–324.
- Henderson, C. R. (1973) Proc. “Anim. Breed. Genet. Symp. in Honor of Dr. J. L. Lush”, pp. 1041. Amer. Soc. Anim. Sci. and Amer. Dairy Sci. Assoc., Champaign, IL.
- Henderson, C. R. (1975) *Biometrics* 31:423–449.
- Hoeschele, I. (1993) *Theor. Appl. Genet.* 85:946–952.
- Hoeschele, I. (1993) *J. Dairy Sci.* 76:1693–1713.
- Kennedy, B. W., Quiton, M. and Van Arendonk, J. A. M. (1992) *J. Anim. Sci.* 70:2000–2012.
- Kinghorn, B. P. and Kerr, R. J. (1995) Proc. “Proc. Eleventh Conf. Australian Assoc. of Anim. Breeding and Genetics”. University of New England, Armidale, Australia.
- Nejati-Javaremi, A., Smith, C. and Gibson, J. P. (1997) *J. Anim. Sci.* 75:1738–1745.
- Rothschild, M., Jacobson, C., Vaske, D., Tuggle, C., Wang, L., Short, T., Eckardt, G., Sasaki, S., Vincent, A., McLaren, D., Southwood, O., van der Steen, H., Mileham, A. and Plastow, G. (1996) *Proc. Natl. Acad. Sci. USA* 93:201–205.
- Saito, S. and Iwaisaki, H. (1996) *Genet. Sel. Evol.* 28:465–477.
- Searle, S. R. (1971) Linear Models. John Wiley and Sons, New York.
- Smith, C. and Simpson, S. P. (1986) *J. Anim. Breed. Genet.* 103:205–217.
- Smith, C. and Webb, A. J. (1981) *J. Anim. Breed. Genet.* 98:161.
- Soller, M. (1978) *Anim. Prod.* 27:133–139.
- Stricker, C., Fernando, R. L. and Elston, R. C. (1996) *Genetics* 141:1651–1656.
- Uimari, P., Thaller, G. and Hoeschele, I. (1996) *Genetics* 143:1831–1842.
- van Arendonk, J. A. M., Tier, B. and Kinghorn, B. (1994) *Genetics* 137:319–329.
- van Arendonk, J. A. M., Tier, B. and Kinghorn, B. P. (1994) Proc. “17th Int. Congr. Genet.”, p. 192. Birmingham, UK.
- Wang, T., Fernando, R. L., van der Beek, S., Grossman, M. and van Arendonk, J. A. M. (1995) *Genet. Sel. Evol.* 27:251–274.
- Wang, T., van der Beek, S., Fernando, R. L. and Grossman, M. (1991) *J. Anim. Sci.* 69(Suppl. 1):202.
- Weller, J. and Fernando, R. L. (1991) . In: “Gene Mapping: Strategies, Techniques and Applications”, pp. 305–328, Editors L. B. Schook, H. A. Lewin and D. G. McLaren, Marcel Dekker.
- Weller, J. I. (1986) *Biometrics* 42:627–640.
- Xu, S. and Atchley, W. R. (1995) *Genetics* 141:1189–1197.