# Posterior probability of the sire's genotype at a major locus based on progeny-test results for discrete characters

J.L. FOULLEY and J.M. ELSEN

*Institut National de la Recherche Agronomique, Station de Génétique Quantitative et Appliquée, Centre de Recherches de Jouy, 78350 Jouy-en-Josas, France*
*Institut National de la Recherche Agronomique, Station d'Amélioration Génétique des Animaux, Centre de Recherches de Toulouse, BP 27, 31326 Castanet-Tolosan Cedex, France*

## Summary

This paper presents the expression of the posterior probability of the different possible genotypes of a sire at a major locus based on progeny performance for a discrete trait. Binomial, multinomial and Poisson distributions are considered fort this trait. Sires and dams are assumed unrelated. Polygenic as well a major gene effects are supposed to influence the trait recorded. An approximation to the computation of the posterior probability is suggested by expressing the latter, conditionally to estimated values of location and dispersion parameters (population, environmental and polygenic effects) which influence progeny performance.

*Key words : major gene, Bayes' theorem, discrete characters.*

## Résumé

*Probabilité* a posteriori *du génotype paternel à un locus majeur à partir de l'observation des descendants : cas d'un caractère discret*

Cet article présente l'expression de la probabilité *a posteriori* des différents génotypes possibles d'un père à un locus majeur à partir des performances obtenues en contrôle de descendance pour un caractère discret. Des distributions binomiale, multinomiale et de Poisson sont considérées pour ce dernier. Les parents sont supposés non apparentés entre eux. Le caractère mesuré est supposé sous la dépendance d'un gène majeur et de polygènes. Un calcul approché de la probabilité est proposé basé sur l'expression de celle-ci, conditionnellement à des valeurs estimées des paramètres de position et de dispersion relatifs aux effets « population », environnementaux et polygéniques résiduels influençant les performances.

*Mots clés : gène majeur, théorème de Bayes, caractère discret.*

## I. Introduction

For a trait with a mixed model of inheritance (a single locus and polygenes), rules to assign genotype at a major locus as a function of observed performance are often empirical. A good example for this is given by the *Booroola* gene with major effect on ovulation rate and litter size in sheep (PIPER & BINDON, 1982). In that case, rules to assign genotype to females are usually those defined by DAVIS *et al.* (1982) ; *i.e.* : ewes are said carriers of the *F* gene if they have at least one observation where the ovulation rate is higher than 3. They are declared homozygote if one observation at least is higher than 5. These values 3 and 5 were derived from performance achieved in a relatively low prolific breed, the *Merino*. Obviously, these thresholds should be revised for more prolific breeds in which this gene is now introduced. Another example is muscular hypertrophy in cattle as affected by a major recessive gene (VISSAC, 1972 ; ROLLINS *et al.*, 1972 ; HANSET & MICHAUX, 1985 *a* & *b*).

For such sex limited traits, inference about genotypes of males require information from related females, usually groups of paternal half sibs. In that situation, the former rule can be applied to female progeny individually and genotype inference based upon agreement between expected and observed numbers of the different genotypes in the female progeny of a given sire. Again, such criteria are open to criticism due to ignoring polygenic variation and for other methodological reasons (FOULLEY & FREBLING, 1985).

To improve these empirical criteria, ELSEN *et al.* (1988) dealing with progeny-testing of males, suggested to use computation of posterior probabilities. Classically, the trait observed in the progeny was assumed normally distributed. The purpose of this study is to extend this approach to discrete phenotypic distributions as often encountered in reproductive and productive performance.

## II. Methodology

The following design is considered. Each sire whose genotype at a major locus is investigated, is randomly mated to several females having a given genotype (usually but not necessarily so), and progeny performance for a trait influenced by the major locus is recorded. Let us define the following symbols :

$G_i$          : a random variable corresponding to the genotype of individual i,

$g_i$          : a realized value of this random variable among r possible values coded by the integers 1, 2, ..., r ; G(g) will be used for sires and $G^*(g^*)$ for progeny,

**G** and **g** : column vectors made up of former components,

**y**          : a data vector with elements $y_{ij}$ where i = 1, 2, ..., q refers to sire and j = 1, 2, ..., $n_i$ to progeny within sire such as

$$\mathbf{y} = (\mathbf{y}'_1, \mathbf{y}'_2, ..., \mathbf{y}'_q)' \quad \text{with } \mathbf{y}_i = \{y_{ij}\} \text{ for } j = 1, 2, ..., n_i$$

$\mathbf{Y}$, $\mathbf{Y}_{ij}$, $\mathbf{Y}_i$ symbolize the corresponding random variables.

By a straightforward application of Bayes' theorem, one has :

$$\Pr(\mathbf{G} = \mathbf{g} \mid \mathbf{y}) \propto \Pr(\mathbf{G} = \mathbf{g}) \; p(\mathbf{y}|\mathbf{G} = \mathbf{g}) \qquad (1)$$

where $\Pr(\cdot)$ and $p(\cdot)$ refer to a probability and a density function respectively.

If sires and dams are unrelated, the random variables $(\mathbf{G}_i, \mathbf{Y}_i')$ are independent and (1) reduces just to :

$$\Pr(\mathbf{G} = \mathbf{g}|\mathbf{y}) = \prod_i \Pr(\mathbf{G}_i = \mathbf{g}_i|\mathbf{y}_i) \qquad (2)$$

and letting $\mathbf{G}_i' = \{g_{ij}'\}$ for $j = 1, 2, ..., n_i$, it can be shown that (ELSEN *et al.*, 1988)

$$\Pr(\mathbf{G}_i = k|\mathbf{y}_i) \propto \Pr(\mathbf{G}_i = k) \sum_{g_i'} \Pr(\mathbf{G}_i' = \mathbf{g}_i'|\mathbf{G}_i = k) \; p(\mathbf{y}_i|\mathbf{g}_i') \qquad (3)$$

with

$$\Pr(\mathbf{G}_i' = \mathbf{g}_i'|\mathbf{G}_i = k) = \prod_{j=1}^{n_i} \Pr(\mathbf{G}_{ij}' = g_{ij}'|\mathbf{G}_i = k) \qquad (4)$$

$$p(\mathbf{y}_i|\mathbf{g}_i') = \int_{R_\alpha} p(\mathbf{y}_i|\mathbf{g}_i', \; \alpha) \; dF(\alpha) \qquad (5a)$$

$$\int_{R_\alpha}[\prod_j p(y_{ij}|g_{ij}', \; \alpha)]dF(\alpha) \qquad (5b)$$

where $F(\alpha)$ is the cumulative density function pertaining to the prior distribution of a vector $\alpha$ of parameters $(\alpha \in R_\alpha)$ involved in environmental and residual polygenic influences and which will be discussed more in detail in the next chapter.

Using (4) and (5b) in (3) and changing the order of summation and product operators, one obtains the general expression of the posterior probability that sire i has genotype k knowing its progeny performance $y_i$ (ELSEN *et al.*, 1988) :

$$\Pr(\mathbf{G}_i = k|\mathbf{y}_i) \; \alpha \; \Pr(\mathbf{G}_i = k) \times \int_{R_\alpha}[\prod_{j=1}^{n_i} \sum_{\ell=1}^{r} \Pr(\mathbf{G}_{ij}' = \ell|\mathbf{G}_i = k)$$

$$p(y_{ij}|\mathbf{G}_{ij}' = \ell,\alpha)]dF(\alpha). \qquad (6)$$

Notice that (6) involves the following elements :

i) $\Pr(\mathbf{G}_i = k)$ is the prior probability that sire i has genotype k. These probabilities depend on the mating structure out of which tested males are born. For instance, when introducing the *Booroola* gene into a foreign population by successive backcrossings, one will usually have heterozygote $(F+)$ and homozygote normal $(++)$ rams in equal proportions at each generation. In human genetics, these prior probabilities are generally assumed to be those of a Hardy-Weinberg population.

ii) $\Pr(\mathbf{G}_{ij}' = \ell|\mathbf{G}_i = k)$ are elements of transition matrix such as the T matrix developed by GEPPERT & KOLLER (1938) and LI & SACKS (1954) when dams are chosen at random in a Hardy-Weinberg population.

iii) $p(y_{ij}|\mathbf{G}_{ij}' = \ell, \; \alpha)$ is the likelihood function for the observed performance of progeny j born out of sire i conditionally to its genotype and which will be discussed more in detail in the next section.

### III. Likelihood functions

Several cases may be considered according to which distribution is hypothesized for the performance trait recorded.

### A. *Normal distribution*

This situation has been described thoroughly by ELSEN *et al.* (1988) and will just be summarized here so as to introduce the parameterization. Let us now designate the $j^{th}$ progeny of the $i^{th}$ sire by the subscript m, the conditional distribution of the progeny performance may be written as :

$$Y_m | (G_{m=1}^{\cdot}), \; \alpha \sim N(w_{ml}^{'} \, \theta, \; \sigma_{e\ell}^2) \tag{7}$$

where $\alpha = (\theta', \; \gamma')$ is a concatenation of location ($\theta$) and dispersion ($\gamma$) parameters defined as follows :

$$\theta = (\beta', \; u')'. \tag{8}$$

The $\beta$ component usually represents « population » and systematic environmental effects such as herd $\times$ year $\times$ season, age of dam, etc.

To make the model flexible, these effects will be allowed to vary from one genotype to another. Therefore, we will write

$$\beta' = (\beta_1', \; \beta_2', \; ..., \; \beta_\ell', \; ..., \; \beta_r') \tag{9a}$$

and

$$\beta_\ell' = (v_\ell, \; b_{\ell 1}', \; b_{\ell 2}', \; ....,) \tag{9b}$$

where $v_\ell$ is the progeny mean for genotype $\ell$ and $b_{\ell 1}$, $b_{\ell 2}$, ... are effects of factors 1, 2, ... expressed in deviation from this mean.

Polygenic variations will be represented by the vector $u$ which may include different kinds of effects (FOULLEY *et al.*, 1987 c) *e.g.* sire and/or maternal grand sire transmitting abilities, additive genetic values and/or permanent influences. Again, these effects are formally allowed to vary according to the genotype of the progeny in which the trait is expressed

$$u' = (u_1', \; u_2', \; ..., \; u_\ell', \; ..., \; u_r'). \tag{10}$$

The vector $\gamma$ of dispersion parameters is decomposed into :

$$\gamma' = (\gamma_u', \; \gamma_e'). \tag{11}$$

The vector $\gamma_u$ component is formed by variances and covariances among elements in (10), say r variances $\rho_{u\ell}^2$ and $r(r-1)/2$ covariances $\sigma_{u\ell\ell'}$, when just one factor (*e.g.* additive genetic values or sire) is considered.

Conditionally to $G_m^{\cdot} = \ell$, $\beta$ and $u$, the distribution of progeny performance has a mean which is a specific linear combination of effects in $\theta$, say $w_{m\ell}^{'} \, \theta$, where $w_{m\ell}^{'}$   is

a known row incidence matrix for progreny m knowing its genotype is $\ell$ and has a variance which is the usual residual variance $\sigma_{e_\ell}^2$. This leads to

$$\gamma_e' = (\sigma_{e_1}^2, \ \sigma_{e_2}^2, \ ..., \ \sigma_{e_\ell}^2, \ ..., \ \sigma_{e_r}^2) \tag{12}$$

when allowance is made for possible heterogeneity of residual variances among progeny genotypes.

## B. *Binomial and multinomial distributions*

When the trait considered is an all-or-none trait, such as twinning in cattle for which a major gene effect has been suggested (MORRIS & DAY, 1986), one will use the « liability » concept originally developed by WRIGHT (1934). This model postulates an underlying normal distribution rendered dichotomous *via* an abrupt threshold.

Keeping the same notation as in the previous section, but with vector $\theta$ defined now as a vector of effects in the underlying scale, we will take for a binary variate, $y_m = 0$ or 1 for categories coded [0] to the left and [1] to the right of the threshold respectively.

$$\Pr[Y_m = y_m | (G_m^{\cdot} = \ell), \ \mu_{m\ell}] = [\Phi(\mu_{m\ell})]^{y_m}[1 - \Phi(\mu_{m\ell})]^{1 - y_m} \tag{13a}$$

with

$$\mu_{m\ell} = w_{m\ell}' \ \theta \tag{13b}$$

where $\Phi(\cdot)$ is the standardized normal cumulative density function and $\mu_{m\ell}$ the mean of the distribution of progeny performance conditionally to genotype $\ell$, the origin being taken at the threshold. It has been assumed in (13 a & b) that the residual variance in the underlying scale $\sigma_{e_\ell}^2$ is constant. This is not a very limiting constraint since the heterogeneity in residual variances in the underlying continuum results usually from a scale effect.

This approach can also be applied to ordered polytomies. In that case, the probability that progeny m responds in category q (q = 1, 2, ..., c + 1) can be written as

$$\Pr[Y_{mq = 1} | (G_m^{\cdot} = \ell), \ \eta_{m\ell}] = \Phi(t_q - \eta_{m\ell}) - \Phi(t_{q - 1} - \eta_{m\ell}) \tag{14a}$$

where $t_1, \ t_2, \ ..., \ t_q, \ ..., \ t_c$ are parameters locating a reference population from the c thresholds (GIANOLA & FOULLEY, 1983) with $t_0 = -\infty$ and $t_{c + 1} = +\infty$

and

$$\eta_{m\ell} = w_{m\ell}' \ \theta \tag{14b}$$

is a location parameter similar to (13b). Hence, the likelihood is product binomial or multinomial according to the distribution considered.

## C. *Poisson distribution*

Some traits such as ovulation rate or litter size might be better described with a Poisson than a multinomial distribution. Then, following FOULLEY *et al.*, (1987 *c*), the conditional distribution of progeny performance is given by

$$\Pr[Y_m = y_m | (G_m^{\cdot} = \ell), \ \lambda_{m\ell}] = \lambda_{m\ell}^{y_m} \ \exp(-\lambda_{m\ell})/y_m \ ! \tag{15a}$$

with

$$\lambda_{m\ell} = \exp(\mu_{m\ell}) \tag{15b}$$

$$\mu_{m\ell} = w'_{m\ell} \theta \tag{15c}$$

As pointed out by these authors, one may also envision truncated Poisson distributions. For zero excluded, (15a) must be replaced by

$$\Pr[Y_m = y_m > 0 \mid (G_m^* = \ell), \lambda_{m\ell}] = \lambda_{m\ell}^{Y_m} / \{y_m ! [\exp(\lambda_{m\ell}) - 1]\}. \tag{16}$$

## IV. Computations

Computing the exact posterior probability that sire i has genotype $\ell$ according to formula (6) is formidable task which requires integrating out both location ($\theta$) and dispersion ($\gamma$) parameters. Even when $\gamma$ is known, this integration involves the calculation or $\Sigma r^{n_i}$ terms which can be expressed analytically only in the normal case (ELSEN et al., 1988).

Therefore some approximations are necessary especially with discrete distributions. As suggested by ELSEN et al., (1988) with normal traits and GIANOLA et al. (1986) and FOULLEY et al. (1987 a, b) in genetic evaluation problems, one can evaluate (1) conditionally to $\theta$ and $\gamma$ values estimated from the data i.e. compute

$$\Pr(G = g|y, \theta^* (\gamma^*)] \tag{17}$$

$$\theta^*(\gamma) : \text{Max}_\theta \log p(\theta|y, \gamma) \tag{18a}$$

$$\gamma^* : \quad : \text{Max}_\gamma \log p(\gamma|y) \tag{18b}$$

where

$\theta^* (\gamma)$ is the mode of the posterior distribution of $\theta$ given $\gamma$

$\gamma^*$ is the mode of the marginal posterior distribution of $\gamma$ which corresponds to the marginal maximum likelihood (ML) estimator when a flat prior for $\gamma$ is used.

### A. *Estimation of location parameters*

In consideration of the threshold-liability model, we will take as prior density (GIANOLA & FOULLEY, 1983 ; FOULLEY et al., 1987 c)

$$p(\theta|\Gamma, \theta_0) \propto \exp\left[ -\frac{1}{2} (\theta - \theta_0)' \Gamma^{-1}(\theta - \theta_0)\right]. \tag{19}$$

Usually, but no necessarily so, one takes $\theta'_0 = (\delta', 0)$ and the $\beta$ component of $\Gamma$, say $\Gamma_\beta \to \infty$, so as to mimic a mixed model structure with $\Gamma$ depending on $\gamma_u$ only. Now, with the same assumptions as before, the likelihood can be written as

$$p(y|\theta) = \prod_{i=1}^{q} \prod_{j=1}^{n_i} \sum_{\ell=1}^{r} \Pr(G_{ij}^* = \ell)p [y_{ij}|(G_{ij}^* = \ell), \theta] \tag{20}$$

$$\Pr(G_{ij}^* = \ell) = \sum_{k=1}^{r} \Pr(G_i = k) \Pr(G_{ij}^* = \ell|G_i = k). \tag{21}$$

Hence, the posterior distribution is

$$p(\theta|y, \Gamma, \theta_0) \propto (19) \times (20). \qquad (22)$$

Maximizing (22) with respect to $\theta$ involves solving a nonlinear system of equations, using for instance the Newton-Raphson algorithm (FOULLEY et al., 1987 b). Let $L(\theta)$ be the logarithm of the posterior density defined in (22) and, again m being a single subscript for the combination ij, the first and second partial derivatives of L with respect to $\theta$ are :

$$\overset{\bullet}{L}(\theta) = -\Gamma^{-1}(\theta - \theta_0) + \sum_m \sum_\ell q_{m\ell} \, v_{m\ell} \, w_{m\ell} \qquad (23a)$$

$$\overset{\bullet\bullet}{L}(\theta) = -\Gamma^{-1} - \sum_m \sum_\ell r_{m\ell} \, w_{m\ell} \, w'_{m\ell}. \qquad (23b)$$

Hence, putting

$$s_\ell = \{q_{m\ell} \, v_{m\ell}\} \text{ for } m = 1, 2, ..., N \qquad (24a)$$

$$R_\ell = \text{Diag } \{r_{m\ell}\} \text{ for } m = 1, 2, ..., N \qquad (24b)$$

$$W_\ell = (w_{1\ell}, \, w_{2\ell}, \, ..., \, w_{m\ell}, \, ..., \, w_{N\ell})' \qquad (24c)$$

$$z_\ell = W_\ell \theta + R_\ell^{-1} s_\ell \qquad (24d)$$

the Newton-Raphson algorithm consists in iterating from round t to t + 1 with

$$[(\sum_\ell W_\ell \, R_\ell^{[t]} \, W_\ell) + \Gamma^{-1}] \, \theta^{[t+1]} = \sum_\ell W'_\ell \, R_\ell^{[t]} \, z_\ell^{[t]} + \Gamma^{-1} \theta_0 \qquad (25)$$

where

$$q_{m\ell} = \frac{\Pr(G_m^* = \ell) \, p[y_m|(G_m^* = \ell), \, \theta, \, \Gamma]}{\sum_\ell \Pr(G_m^* = \ell) \, p[y_m|(G_m^* = \ell), \, \theta, \, \Gamma]} \qquad (26a)$$

$$v_{m\ell} = \partial \log p[y_m|(G_m^* = \ell), \, \theta, \, \Gamma]/\partial \mu_{m\ell} \qquad (26b)$$

$$r_{m\ell} = -q_{m\ell}(\partial v_{m\ell}/\partial \mu_{m\ell}) - q_{m\ell}(1 - q_{m\ell})v_{m\ell}^2 \qquad (26c)$$

Analytical expressions of these coefficients can be derived explicitly for the different discrete distributions considered previously, i.e.

i) for a Bernouilli variate

$$v_{m\ell} = \varphi(\mu_{m\ell}) \, [y_m - \Phi(\mu_{m\ell})]/\{\Phi(\mu_{m\ell}) \, [1 - \Phi(\mu_{m\ell})]\} \qquad (27a)$$

$$r_{m\ell} = q_{m\ell}v_{m\ell}(q_{m\ell}v_{m\ell} + \mu_{m\ell}) \qquad (27b)$$

where $\varphi$ is the standardized normal density function. For several ordered categories, explicit formulae for $v_{m\ell}$ and $E(\partial v_{m\ell}/\partial \eta_{m\ell})$ are shown in GIANOLA & FOULLEY (1983) ; moreover, the system in (25) must be augmented by sectors pertaining to the thresholds.

ii) for a Poisson variate

As shown by FOULLEY et al. (1987 c)

$$v_{m\ell} = y_m - \lambda_{m\ell} \qquad (28a)$$

$$r_{m\ell} = q_{m\ell} \, \lambda_{m\ell} - q_{m\ell} \, (1 - q_{m\ell}) \, (y_m - \lambda_{m\ell})^2. \qquad (28b)$$

If the Poisson model truncated at zero is employed, $v_{m\ell}$ becomes

$$v_{m\ell} = y_m - \{\lambda_{m\ell}/[1 - \exp(-\lambda_{m\ell})]\} \qquad (29a)$$

and $r_{m\ell}$ is calculated according to (26c) with

$$\partial v_{m\ell}/\partial \mu_{m\ell} = - \frac{\lambda_{m\ell}}{1 - \exp(- \lambda_{m\ell})} \left[ 1 - \frac{\lambda_{m\ell} \exp(- \lambda_{m\ell})}{1 - \exp(- \lambda_{m\ell})} \right].$$ (29b)

The system in (25) with formulae (26a, b & c) are similar to those given by FOULLEY et al. (1987 a) for genetic evaluation with uncertain parternity. The $q_{m\ell}$ coefficient gives rise to an interesting interpretation since (26a) can be viewed as the posterior probability that progeny m has genotype $\ell$. This expression would also occur naturally if the problem was set up in terms of incomplete data and solve accordingly via the EM algorithm (DEMPSTER et al., 1977, p. 16). As a matter of fact, letting $q_{m\ell} = 1$ in the expressions of $v_{m\ell}$ and $r_{m\ell}$, one obtains the usual coefficients encountered in genetic evaluation for binary (FOULLEY et al., 1983 ; FOULLEY et al., 1987 d) and Poisson variates (FOULLEY et al., 1987 c). Finally, the form of the system in (25) indicates that the analysis is carried out conditionally to the different possible genotypes of progeny with appropriate weighting factors on the left and right handsides depending on the posterior probabilities of genotypes. This generates two sources of nonlinearity as shown clearly in the formula for $r_{m\ell}$ (26c), one due to the form of the distribution and the second to uncertainty about genotype of progeny.

## B. *Estimation of dispersion parameters*

The value of $\theta$ calculated from (25) is the mode of the condition distribution of $\theta$ given $\Gamma$ and the data. It remains to replace $\Gamma$ by its marginal ML estimator. For the sake of simplicity we will consider the case of just one random factor u such as sire transmitting abilities or individual additive genetic values. Then, the unknown is the vector $\gamma_u$ (see 11) formed by the $r(r + 1)/2$ different elements of the $G = \{g_{k\ell}\}$ matrix of « u » components of variance and covariance. The general procedure for discrete variates has been presented by HARVILLE & MEE (1984), FOULLEY et al. (1987 b & c). These authors have shown that maximization of the logposterior density of $\gamma_u$ with respect to $\gamma_u$ using a diffuse prior, results in the equation

$$E_C \left\{ \frac{\partial}{\partial \gamma_u} \log[p(u|\gamma_u)] \right\} = 0$$ (30)

here $E_C$ indicates expectation with respect to the density $p(u|y, \gamma_u)$.

Within the framework of a normal distribution for $u|\gamma_u$, the formal solution to (30) is very general whatever form of the likelihood. Provided some approximations are made about the first two moments of $u|y, \gamma_u$, computations amount to iterate with

$$g_{k\ell}^{[s + 1]} = [u_k'^{[s]} u_\ell^{[s]} + \text{tr} (C_{k\ell}^{[s]})]/q$$ (31)

where

$u_k^{[s]}$, $u_\ell^{[s]}$ are solutions to (25) in $u_k$ and $u_\ell$ given $\gamma_u = \gamma_u^{[s]}$

$C_{k\ell}^{[s]}$ is the $(q \times q)$ submatrix pertaining to genotypes k and $\ell$ in the u part of the inverse of the coefficient matrix in (25).

## C. *Computation of the posterior probabilities*

Under the same assumptions as in II (no genetic relationships among parents), formula (17) may be simply written as

$$(17) = \prod_i \Pr[G_i = g_i | y_i, \theta^*(\gamma^*)].$$

Hence, after having estimated $\theta$ and $\gamma$, the posterior probability that sire i has genotype k will be calculated as

$$\Pr[G_i = k | y_i, \theta^*(\gamma^*)] = N_{ik}/D_i \qquad (32a)$$

with

$$N_{ik} = \Pr(G_i = k) \prod_{j=1}^{n_i} \sum_{\ell=1}^{r} \Pr(G_{ij}^* = \ell | G_i = k) \, p[y_{ij} | (G_{ij}^* = \ell), \theta^*(\gamma^*)] \qquad (32b)$$

$$D_i = \sum_{k=1}^{r} N_{ik}. \qquad (32c)$$

The last terms in (32 b) are computed from (13 a & b), (14 a & b) and (15 a, b & c) replacing $\theta$ by $\theta^*(\gamma^*)$.

## V. Discussion-Conclusion

It has been implicitly assumed till now that one record per progeny is available. Taking into account several performances per animal can be easily achieved using a « repeatability » model as follows :

$$p(y_i | g_i, \alpha) = \prod_{jt} p[y_{ijt} | (G_{ij}^* = g_{ij}^*), \alpha]$$

where t is the subscript for the $t^{th}$ performance within progeny ij. The location parameter used in the likelihood is then, for progeny m given genotype $\ell$ :

$$\mu_{m\ell t} = w_{m\ell t}' \theta$$

where, as previously $\theta' = (\beta', u')$, and u can be parameterized as $u' = (s', p')$, s being a vector of sire effects and p a vector of permanent environmental effects for a given progeny within sire. The corresponding dispersion parameters become

$$\Gamma_u = \begin{pmatrix} I\sigma_s^2 & O \\ O & I\sigma_p^2 \end{pmatrix}$$

and the coefficient $\rho_k = (\sigma_{sk}^2 + \sigma_{pk}^2)/(\sigma_{sk}^2 + \sigma_{pk}^2 + \sigma_{\ell k}^2)$

designates the usual repeatability coefficient for genotype k which may be assumed constant as in ELSEN et al. (1988). The procedure described in the previous chapter IV is still valid especially the method for estimating dispersion parameters which can be easily extended to several sets of random factors (FOULLEY et al., 1987 c).

It is also worth mentioning that the approach followed in chapter IV provides as a by-product a genetic evaluation method for traits with a mixed model of inheritance (ELSTON & STEWART, 1971 ; MORTON & MCLEAN, 1974 ; LALOUEL et al., 1983) and with completely or partially unknown genotypic information. In that case, the coefficient matrix in (25) may be very large especially with field data sets. Some problems encountered in solving such large non linear systems (e.g. convergence properties, precision, computing costs) have been recently discussed by MISTZAL & GIANOLA (1987) for sire evaluation programs dealing with threshold traits. Anyhow, this is another example of how the Bayesian paradigm can be used to solve problems in that area which cannot be readily addressed via the BLUP machinery.

However with the method proposed, nuisance parameters have been averaged out and not exactly integrated out as it should be. Therefore, it is important to bear in mind that this procedure as others (HASSTEDT, 1982) is an approximation, the domain of validity of which should be more carefully addressed using for instance realistic examples with Monte-Carlo simulation techniques.

Some practical provisions can be suggested to apply this procedure shrewdly. First, genetic and phenotypic parameters must be taken as known to reduce the degree of nonlinearity of the problem especially when a limited number of sires is tested. Secondly, one has better to choose at start a simple parameterization with no specific effects and variances according to genotypes of progeny. Finally, when the distribution of performance is clearly multimodal, one may expect some of the coefficients q to have extreme values. In order to improve the assignment of genotypes to sires, one might use for instance some prior information about genotypic means in calculating the q's at the first round of iteration. Anyhow, we are truly conscious that the approximations proposed in (17) and (18 a & b) may severely limit the potential interest of this methodology as long as formula (6) cannot be calculated efficiently by numerical procedures of integration.

### Acknowledgements

### References

DAVIS G.H., MONTGOMERY G.W., ALLISON A.J., KELLY R.W., BRAY A.R., 1982. Segregation of a major gene influencing fecundity in progeny of Booroola sheep. *N.Z. J. Agr. Res.*, **25,** 525-529.

DEMPSTER A.P., LAIRD N.M., RUBIN R.B., 1977. Maximum likelihood estimation from incomplete data via the EM algorithm. *J. Roy. Statist. Soc.*, **B, 39,** 1-20.

ELSEN J.M., VU TIEN KHANG J., LE ROY P., 1988. A statistical model for genotype determination at a major locus in a progeny test design. *Génét. Sél. Evol.*, **20,** 211-226

ELSTON N.E., STEWART J., 1971. A general model for the genetic analysis of pedigree data. *Hum. Hered.*, **21,** 523-542.

FOULLEY J.L., FREBLING J., 1985. Critères de détection indirecte des taureaux porteurs de la translocation 1-29 à partir des caryotypes de leurs descendants. *Génét. Sél. Evol.*, **17,** 341-350.

FOULLEY J.L., GIANOLA D., THOMPSON R., 1983. Prediction of genetic merit from data on categorical and quantitative variates with an application to calving difficulty, birth weight and pelvic opening. *Génét. Sél. Evol.*, **15,** 407-424.

FOULLEY J.L., GIANOLA D., PLANCHENAULT D., 1987 *a*. Sire evaluation with uncertain paternity. *Génét. Sél. Evol.*, **19,** 83-102.

FOULLEY J.L., IM S., GIANOLA D., HÖSCHELE I., 1987 b. Empirical Bayes estimation of parameters for n polygenic binary traits. *Génét. Sél. Evol.,* **19,** 197-224.

FOULLEY J.L., GIANOLA D., IM S., 1987 c. Genetic evaluation for traits distributed as Poisson-binomial with reference to reproductive traits. *Theor. Appl. Genet.,* **73,** 870-877.

FOULLEY J.L., GIANOLA D., IM S., 1987 d. Genetic evaluation for discrete polygenic traits in animal beeding. *In : Advances in statistical methods for genetic improvement of animals,* Springer-Verlag, Berlin, forthcoming.

GEPPERT S., KOLLER S., 1938. *Erbmathematik : Theorie der Veterbung in Bevölkerung und Sippe.* 228 p., Quelle und Meyer, Leipzig.

GIANOLA D., FOULLEY J.L., 1983. Sire evaluation for ordered categorical data with a threshold model. *Génét. Sél. Evol.,* **15,** 201-224.

GIANOLA D., FOULLEY J.L., FERNANDO R.L., 1986. Prediction of breeding values when variances are not known. *Génét. Sél. Evol.,* **18,** 485-498.

HANSET R., MICHAUX C., 1985 a. On the genetic determinism of muscular hypertrophy in the Belgian White and Blue cattle breed. 1. Experimental data. *Génét. Sél. Evol.,* **17,** 359-368.

HANSET R., MICHAUX C., 1985 b. On the genetic determinism of muscular hypertrophy in the Belgian White and Blue cattle breed. 2. Population data. *Génét. Sél. Evol.,* **17,** 369-386.

HARVILLE D.A., MEE R.W., 1984. A mixed model procedure for analyzing ordered categorical data. *Biometrics,* **40,** 393-408.

HASSTEDT S.J., 1982. A mixed-model likelihood approximation on large pedigrees. *Comput. Biomed. Res.,* **15,** 295-307.

LALOUEL J.M., RAO D.C., MORTON N.E., ELSTON R.C., 1983. A unified model for complex segregation analysis. *Am. J. Hum. Genet.,* **35,** 816-826.

LI C.C., SACKS L., 1954. The derivation of joint distribution and correlation between relatives by the use of stochastic matrices. *Biometrics,* **10,** 347-360.

MISZTAL I., GIANOLA D., 1987. Computing aspects of a nonlinear method of sire evaluation for categorical data. *J. Dairy Sci.,* **70,** suppl. 1, 124 (abstr.).

MORRIS C.A., DAY A.M., 1986. Ovulation results from cattle herds with high twinning frequency. *In :* DICKERSON G.E., JOHNSON R.K. (ed.), *3rd World Congress on Genetics applied to Livestock Production, Lincoln, Nebraska, July 16-22, 1986,* **11,** 96-100, Editorial Garsi, Madrid.

MORTON N.E., McLEAN C.J., 1974. Analysis of family resemblance. 3. Complex segregation analysis of quantitative traits. *Am. J. Hum. Genet.,* **26,** 489-503.

PIPER L.R., BINDON B.M., 1982. Genetic segregation for fecundity in Booroola Merino sheep. *In :* BARTON R.A., SMITH W.C. (ed.), *Proceedings of the World Congress on Sheep and Beef Cattle Breeding,* vol. 1 (Technical), 395-400, The Dunmore Press, Palmerston North.

ROLLINS W.C., TANAKA M., NOTT C.F.G., THIESSEN R.B., 1972. On the mode of inheritance of double muscled conformation in bovines. *Hilgardia,* **41,** 433-456.

VISSAC B., 1972. L'hypertrophie musculaire d'origine génétique ou caractère culard. *Ann. Génét. Sél. Anim.,* **4,** 87-97.

WRIGHT S., 1934. An analysis of variability in number of digits in an inbred strain of Guinea pigs. *Genetics,* **19,** 506-536.