

# Genetic components of litter size variability in sheep

Magali SANCRISTOBAL-GAUDY<sup>a,\*</sup>, Loys BODIN<sup>b</sup>,  
Jean-Michel ELSÉN<sup>b</sup>, Claude CHEVALET<sup>a</sup>

<sup>a</sup> Laboratoire de génétique cellulaire, Institut national de la recherche agronomique,  
BP 27, 31326 Castanet-Tolosan, France

<sup>b</sup> Station d'amélioration génétique des animaux,  
Institut national de la recherche agronomique,  
BP 27, 31326 Castanet-Tolosan, France

(Received 6 June 2000; accepted 11 December 2000)

**Abstract** – Classical selection for increasing prolificacy in sheep leads to a concomitant increase in its variability, even though the objective of the breeder is to maximise the frequency of an intermediate litter size rather than the frequency of high litter sizes. For instance, in the Lacaune sheep breed raised in semi-intensive conditions, ewes lambing twins represent the economic optimum. Data for this breed, obtained from the national recording scheme, were analysed. Variance components were estimated in an infinitesimal model involving genes controlling the mean level as well as its environmental variability. Large heritability was found for the mean prolificacy, but a high potential for increasing the percentage of twins at lambing while reducing the environmental variability of prolificacy is also suspected. Quantification of the response to such a canalising selection was achieved.

**canalising selection / threshold trait / heterogeneous variances / litter size / sheep**

## 1. INTRODUCTION

Selection for increasing prolificacy in sheep, although leading to a better average litter size in selected lines, also leads to an increase in prolificacy variability. This phenomenon is well known for qualitative traits, where mean and variance are linked. Extreme litters are encountered in prolific ewes (Romanov; Finnish) with five or even more lambs per lambing, which is obviously unacceptable for ewe and lamb viability. Breeders would like to have litter sizes of two exactly – and not on average – or as often as possible. In many situations twins are the most profitable (Benoit, personal communication).

Based on the example of the French Lacaune breed, the aim of this work was to evaluate if sheep can be selected for the objective: “concentrating prolificacy

---

\* Correspondence and reprints  
E-mail: msc@toulouse.inra.fr

on 2". For that purpose, data consisting of litter size measurements on Lacaune sheep were analysed, using a direct adaptation to ordered categorical data of the quantitative genetic model described by SanCristobal-Gaudy *et al.* [22] relative to continuous traits. The hypothesis was stated that factors affect the underlying mean and/or the underlying environmental variability. These factors can be environmental, but also genetic. Variance components were estimated, giving the amount of genetic control on the mean and on the environmental variability, in a polygenic context. Prediction of the response to a selection for twins, based on the previous genetic parameter estimates, was derived using Monte Carlo simulation. Finally, this approach was compared with more traditional methods.

## 2. GENETIC MODEL

### 2.1. Threshold model for polytomous data – Likelihood approach

As Gianola and Foulley [10], Foulley and Gianola [8] or SanCristobal-Gaudy *et al.* [23] for example, we consider the threshold Wright model, based on an underlying Gaussian random variable. Thresholds transform this continuous variable into a multinomial variable with  $J$  ordered categories. Let us define  $I$  as cells indexed by  $i$  as combinations of levels of explanatory factors. Multinomial data are observed:

$$(N_{i1}, \dots, N_{ij}, \dots, N_{iJ}) \sim \mathcal{M}(n_{i+}; (\Pi_{i1}, \dots, \Pi_{ij}, \dots, \Pi_{iJ})) \quad (1)$$

with  $N_{ij}$  as the number of counts in cell  $i$  for the  $j$ th category, and  $\Pi_{ij}$  the probability that an unobservable Gaussian random variable  $Y_i \sim \mathcal{N}(\mu_i, \sigma_i^2)$  lies between two thresholds  $\tau_{j-1}$  and  $\tau_j$  (falls into the  $j^{\text{th}}$  ordered category). Setting  $\tau_0 = -\infty$  and  $\tau_J = +\infty$ , the following is obtained:

$$\begin{aligned} \Pi_{ij} &= P[\tau_{j-1} \leq Y_{ik} < \tau_j | Y_{ik} \sim \mathcal{N}(\mu_i, \sigma_i^2), k \in \{1, \dots, n_{i+}\}] \\ &= \Phi\left(\frac{\tau_j - \mu_i}{\sigma_i}\right) - \Phi\left(\frac{\tau_{j-1} - \mu_i}{\sigma_i}\right), \end{aligned} \quad (2)$$

where  $n_{i+}$  is the observed number of counts in cell  $i$  for all  $J$  categories:  $n_{i+} = \sum_j n_{ij}$ .

The underlying means  $\mu_i$  and variances  $\sigma_i^2$  are linear combinations of parameters to estimate:

$$\mu_i = \mathbf{x}'_i \boldsymbol{\beta}, \quad (3)$$

$$\ln \sigma_i^2 = \mathbf{p}'_i \boldsymbol{\delta}, \quad (4)$$

where  $\mathbf{x}'_i$  and  $\mathbf{p}'_i$  are incidence vectors,  $\boldsymbol{\beta}$  is a vector of location parameters, and  $\boldsymbol{\delta}$  is a vector of dispersion parameters.

### Estimation and hypothesis testing

The estimation procedure can simply be maximum likelihood, implementing for example a Fisher-scoring algorithm, exactly as in [8]. Moreover, the test of  $H_0 : \mathbf{K}'\boldsymbol{\delta} = 0$  vs.  $H_1 = \bar{H}_0$ , where  $\mathbf{K}$  is a full-rank matrix, is achieved with the log-likelihood ratio  $\lambda = -2(\mathcal{L}_1 - \mathcal{L}_0)$ , where  $\mathcal{L}_0$  (resp.  $\mathcal{L}_1$ ) is the log-likelihood of model  $\mathcal{M}_0$  (resp.  $\mathcal{M}_1$ ) corresponding to  $H_0$  (resp.  $H_1$ ). Asymptotically, the statistic  $\lambda$  follows a chi-square distribution under the null hypothesis  $H_0$ , with degrees of freedom equal to the difference in the number of estimated parameters between models  $\mathcal{M}_0$  and  $\mathcal{M}_1$ .

## 2.2. Bayesian approach

Furthermore, the Bayesian quantitative genetic model developed by SanCristobal-Gaudy *et al.* [22] is based upon the underlying continuous variable  $Y$  as follows:

$$\mu_i = \mathbf{t}_i'\boldsymbol{\theta} = \mathbf{x}_i'\boldsymbol{\beta} + \mathbf{z}_i'\mathbf{u}, \quad (5)$$

$$\ln \sigma_i^2 = \mathbf{w}_i'\boldsymbol{\gamma} = \mathbf{p}_i'\boldsymbol{\delta} + \mathbf{q}_i'\mathbf{v}, \quad (6)$$

where  $\mathbf{t}_i = (\mathbf{x}_i', \mathbf{z}_i)'$  and  $\mathbf{w}_i = (\mathbf{p}_i', \mathbf{q}_i)'$  are incidence vectors,  $\boldsymbol{\theta} = (\boldsymbol{\beta}', \mathbf{u}')'$  are location parameters, and  $\boldsymbol{\gamma} = (\boldsymbol{\delta}', \mathbf{v}')'$  are dispersion parameters. The parameters  $\boldsymbol{\beta}$  and  $\boldsymbol{\delta}$  have flat priors, in order to mimic a mixed model structure, while  $\mathbf{u}$  and  $\mathbf{v}$  represent genetic values, with a joint normal prior distribution:

$$\begin{pmatrix} \mathbf{u} \\ \mathbf{v} \end{pmatrix} | \sigma_u^2, \sigma_v^2, r \sim \mathcal{N} \left[ \mathbf{0}, \begin{pmatrix} \sigma_u^2 & r\sigma_u\sigma_v \\ r\sigma_u\sigma_v & \sigma_v^2 \end{pmatrix} \otimes \mathbf{A} \right], \quad (7)$$

where  $\otimes$  denotes the Kronecker product,  $\mathbf{A}$  is the relationship matrix between the animals present in the analysis,  $\sigma_u^2$  and  $\sigma_v^2$  are additive genetic variances relative to the location and log variance of the trait, respectively, and  $r$  is the correlation coefficient between genetic values  $\mathbf{u}$  and  $\mathbf{v}$ . Note that the continuous random variable  $Y$  is Gaussian conditional on  $(\mathbf{u}, \mathbf{v})$ . Using a now common incorrect terminology, the expressions “fixed effects” and “random effects” will sometimes be used in the following.

Here, focus is on the genetic aspect of the modelling of multinomial data, by the introduction of two (possibly) related groups of polygenes acting on the trait mean and log variance respectively.

Following SanCristobal-Gaudy *et al.* [22,23], a sire model is written with

$$\mu_i = \mathbf{x}_i'\boldsymbol{\beta} + \frac{1}{2}\mathbf{z}_i'\mathbf{u}, \quad (8)$$

$$\sigma_i^2 = \frac{3}{4}\sigma_u^2 + \exp \left[ \mathbf{p}_i'\boldsymbol{\delta} + \frac{1}{2}\mathbf{q}_i'\mathbf{v} + \frac{3}{8}\sigma_v^2 \right] \quad (9)$$

replacing (5) and (6). Vectors  $\mathbf{u}$  and  $\mathbf{v}$  are genetic values of sires, and data are collected on their progeny.

*Model fitting*

Let us denote  $\mathbf{N} = (N_{ij})_{(i=1,\dots,I)(j=1,\dots,J)}$  as the observation,  $\boldsymbol{\sigma}^2 = (\sigma_u^2, \sigma_v^2, r)$  the set of variance component parameters, and  $\boldsymbol{\zeta} = (\boldsymbol{\tau}', \boldsymbol{\theta}', \boldsymbol{\gamma}')$  the other parameters with  $\boldsymbol{\tau} = (\tau_j)_{j=1,\dots,J}$  as the thresholds. The logarithm  $\mathcal{L}$  of the joint posterior distribution reads:

$$\mathcal{L} = \sum_{i=1}^I \sum_{j=1}^J n_{ij} \ln \Pi_{ij} - \frac{1}{2(1-r^2)} \left[ \frac{\mathbf{u}'\mathbf{A}^{-1}\mathbf{u}}{\sigma_u^2} - 2r \frac{\mathbf{u}'\mathbf{A}^{-1}\mathbf{v}}{\sigma_u\sigma_v} + \frac{\mathbf{v}'\mathbf{A}^{-1}\mathbf{v}}{\sigma_v^2} \right] - \frac{q}{2} \ln \sigma_u^2 - \frac{q}{2} \ln \sigma_v^2 - \frac{q}{2} \ln(1-r^2) + \text{const.} \quad (10)$$

where  $q$  denotes the number of elements in vector  $\mathbf{u}$  (or  $\mathbf{v}$ ).

Estimation of parameters  $\boldsymbol{\zeta}$  *via* the maximisation of  $\mathcal{L}$  with respect to  $\boldsymbol{\tau}$ ,  $\boldsymbol{\theta}$ ,  $\boldsymbol{\gamma}$  presents no theoretical difficulty when variance components are known. A Fisher-scoring algorithm leads to extended mixed-model equations (see Appendix).

When variance components have to be estimated, we chose to base the inference on the mode of the log marginal posterior distribution of variance components  $\boldsymbol{\sigma}^2$ :

$$\hat{\boldsymbol{\sigma}}^2 = \text{Argmax} \ln p(\boldsymbol{\sigma}^2 | \mathbf{N}), \quad (11)$$

by extension of the usual case ( $\sigma_v^2 = 0$ ) where the previous equation leads to REML estimates of variance components.

An EM-type algorithm was implemented as in [9,22], using an iterative algorithm where two systems are involved. The first system consists of BLUP-like mixed-model equations, where variance components are replaced by their current estimates. Solutions of these equations give current estimates of  $\boldsymbol{\zeta}$ . The second system updates the variance component estimates. When  $r$  is set to zero, equation (11) reduces to usual REML equations. However, numerical integration is required for multinomial data; details can be found in the Appendix.

At convergence, maximum *a posteriori* (MAP) estimates of  $\boldsymbol{\zeta}$  are obtained as a by-product:

$$\hat{\boldsymbol{\zeta}} = \text{Argmax} \ln p(\boldsymbol{\zeta} | \boldsymbol{\sigma}^2 = \hat{\boldsymbol{\sigma}}^2, \mathbf{N}). \quad (12)$$

### 3. ANALYSIS OF LITTER SIZE DATA

#### 3.1. Data

Data were collected from Lacaune ewe lambs born over 11 years as the result of inseminations made from 157 sires in 57 flocks. These flocks were a part of a selection scheme implemented in the Lacaune population since 1975 for

**Table I.** Significance effects of explanatory factors on the underlying mean. Reference model is  $YEAR + SEASON + AGE + HERD + SIRE$ .

Factor	Test statistics	df	<i>p</i> -value
– <i>YEAR</i>	15.8	10	0.1
– <i>SEASON</i>	10.4	1	0.001
– <i>AGE</i>	80.2	3	0
– <i>HERD</i>	557.2	56	0
– <i>SIRE</i>	788.2	156	0

increasing prolificacy and operating on farms through a sire progeny test, as described by Perret *et al.* [20]. In the experimental design, each ram offspring averaged 25 daughters spread among five different flocks (factor *HERD*) and each flock had ewe lambs of about eight different sires thus providing a suitable sample for the estimation of genetic values. The sample used in this study was limited to data for rams (factor *SIRE*) with at least 30 controlled daughters. It considered only the first lambing after natural oestrus in ewes of 4 age classes at mating (< 7, 7 to 11, 11 to 14, > 14 months of age, factor *AGE*), and obtained in two lambing seasons (November-December and March-April, factor *SEASON*). This sample involved the results of 11 723 litter sizes over 11 years (factor *YEAR*).

Litter sizes greater than 5 were grouped into the 5th and last category. The percentages of litters with 1, 2, 3, 4 and 5 or more lambs were 41.1, 47.5, 9.8, 1.5 and 0.1 respectively. The overall prolificacy of these ewes at their first lambing was 1.72.

### 3.2. Homoscedastic models

A usual homoscedastic threshold model is fitted, including the fixed effects *YEAR*, *HERD*, *SEASON*, *AGE* in an additive way, and a random sire effect ( $\mathbf{u}/2$ ), symbolically written as:

$$E(Y|\mathbf{u}) = YEAR + HERD + SEASON + AGE + \mathbf{u}/2 \quad (13)$$

on the underlying mean, where  $\mathbf{u} \sim \mathcal{N}_{157}(\mathbf{0}, \sigma_u^2 \mathbf{A})$  is the vector of sire genetic values and  $\mathbf{A}$  is the relationship matrix. Interactions were not taken into account in the model because of non-(or bad) estimability or statistical non-significance. The significance tests for the explanatory factors on the underlying mean are shown in Table I.

The estimation procedure of Gianola and Foulley [10] gave an estimate of heritability equal to  $\hat{h}_u^2 = 0.39$ .

**Table II.** Significance effects of explanatory factors on the underlying environmental log variance.

Reference model	Added factor	$n_{\min}$ <sup>(a)</sup>	$s_{\text{Max}}^2/s_{\text{min}}^2$ <sup>(b)</sup>	$\hat{\sigma}_{\text{Max}}^2/\hat{\sigma}_{\text{min}}^2$	Test statistics	df	p-value
const.	+ <i>YEAR</i>	156	1.38	1.6	20.4	10	0.026
	+ <i>SEASON</i>	5236	1.09	1.02	0.22	1	0.64
	+ <i>AGE</i>	619	1.25	1.22	3.6	3	0.31
	+ <i>HERD</i>	11	3.85	11.17	61.04	56	0.3
	+ <i>SIRE</i>	30	4.63	13.8	237.6	156	$3 \times 10^{-5}$
<i>SIRE</i>	+ <i>YEAR</i>			1.48	16	10	0.1
	+ <i>SEASON</i>			1.01	0.02	1	0.89
	+ <i>AGE</i>			1.28	4.5	3	0.21
	+ <i>HERD</i>			62.55	71.4	56	0.08

<sup>(a)</sup> Minimum number of observations among all levels of each factor.

<sup>(b)</sup> Observed ratio of highest variance over lowest variance among levels of each factor.

### 3.3. Heteroscedastic models

The previous additive model for the mean was used throughout the next analyses.

(i) First, factors that have a significant effect on the underlying trait environmental variability were sought. A likelihood ratio test was implemented. The reference model is the homoscedastic model with only fixed effects, including a sire fixed effect (model of the form (8)-(9), without  $\mathbf{u}$  nor  $\mathbf{v}$ ):

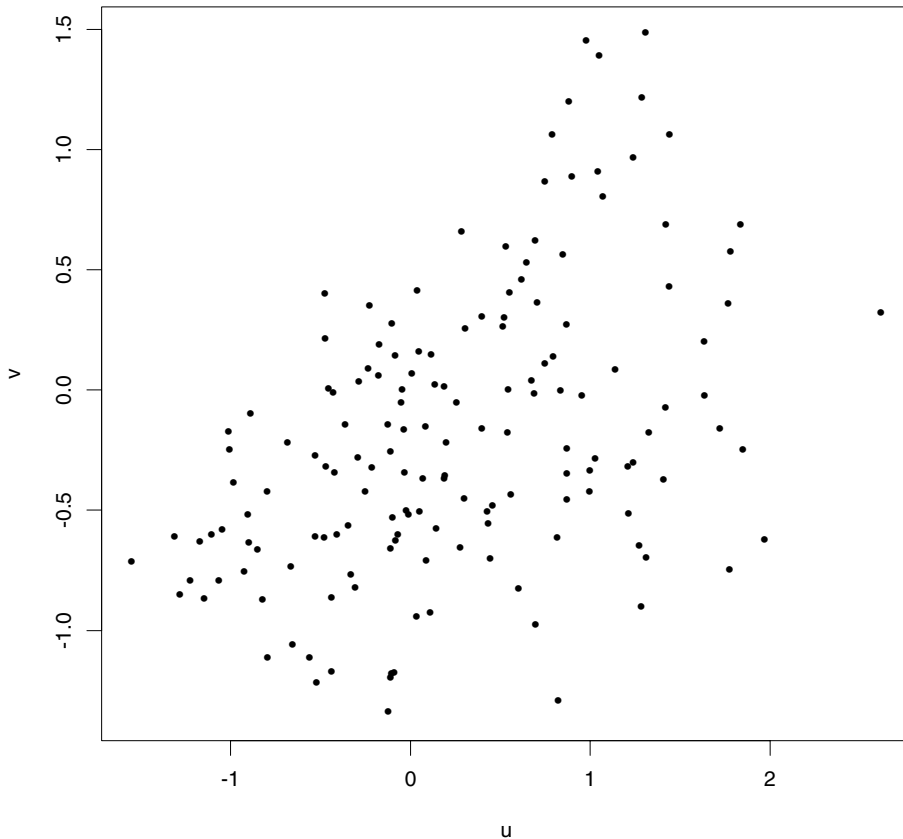
$$\mathcal{M}_0 : \begin{cases} E(Y) = \textit{YEAR} + \textit{HERD} + \textit{SEASON} + \textit{AGE} + \textit{SIRE} \\ \ln \text{Var}(Y) = \text{const.} \end{cases} \quad (14)$$

The current model for the significance test for, say, the *YEAR* factor, is for example:

$$\mathcal{M}_1 : \begin{cases} E(Y) = \textit{YEAR} + \textit{HERD} + \textit{SEASON} + \textit{AGE} + \textit{SIRE} \\ \ln \text{Var}(Y) = \textit{YEAR}. \end{cases} \quad (15)$$

Table II gives the results of a forward selection procedure for the model on log variances. It shows that only the sire (considered as a fixed effect) has a significant effect.

(ii) Then a mixed sire model (8)-(9), with  $\beta = (\textit{YEAR}, \textit{HERD}, \textit{SEASON}, \textit{AGE})$ ,  $\mathbf{u} = \textit{SIRE}$  and  $\mathbf{v} = \textit{SIRE}$ , is fitted in order to estimate the variance components. This gives  $\hat{h}_u^2 = 0.34$  (s.e. = 0.037),  $\hat{\sigma}_v^2 = 0.23$  (s.e. = 0.027)



**Figure 1.** Plot of estimated  $u$  and  $v$  genetic values of the 157 numbered sires, in genetic standard deviation units.

and  $\hat{r} = 0.19$  (s.e. = 0.092). These variance component estimates are approximately the same when the correlation  $r$  between the two sets of breeding values is arbitrarily set to 0 ( $\hat{\sigma}_v^2 = 0.25$  and  $\hat{h}_u^2 = 0.36$ , see also [23]).

The fixed effects and breeding value estimates are compared with those obtained with the mixed homoscedastic threshold model. They are close to each other, although the ranking is not exactly the same (not shown).

A plot of estimated breeding values ( $\hat{u}, \hat{v}$ ) (Fig. 1) allows to apprehend the joint ability of the 157 sires to produce high or low litter size on average *and* with a high or low variability.

In Table III, two sires with a mean prolificacy of the same order of magnitude are compared. The former has a high dispersion while the latter is canalised. The heteroscedastic model detects these differences and predicts slightly better the probabilities for the five categories. The total number of parameters is higher in the heteroscedastic than in the homoscedastic model,

**Table III.** Comparison of two sires. Expected probabilities correspond to an environment with average effect.

Sire	Mean prol.	$\hat{u}$	$\hat{v}$	Model	$\Pi_1$	$\Pi_2$	$\Pi_3$	$\Pi_4$	$\Pi_5$
44	1.80	0.738	0.283	raw data	0.40	0.43	0.14	0.03	0.00
				homosc. mod.	0.48	0.42	0.08	0.01	0.00
				hetero. mod.	0.46	0.36	0.13	0.04	0.01
83	1.73	0.621	-0.625	raw data	0.34	0.59	0.07	0.00	0.00
				homosc. mod.	0.49	0.47	0.04	0.00	0.00
				hetero. mod.	0.45	0.48	0.06	0.01	0.00

but the likelihood ratio test infers that the former better fits the Lacaune data, accounting for the extra number of parameters ( $p$ -value =  $3 \times 10^{-5}$ , see Tab. II).

The high estimate of genetic variance ( $\hat{\sigma}_v^2 = 0.23$ ) and of heritability ( $\hat{h}_u^2 = 0.34$ ) can be viewed as a great potential for the population to be canalised toward the phenotypic optimum of two (twins are economically the best), with a reduction of the environmental variability. The next section is a first attempt to quantify the expected response to such a selection, as was done for continuous traits [22].

#### 4. PREDICTION OF THE RESPONSE TO CANALISING SELECTION OF PROLIFICACY IN THE LACAUNE BREED

##### 4.1. Objective

One of the general objectives is the minimisation of discrepancies from an optimum

$$\Pi_0 = (\Pi_{0,1}, \dots, \Pi_{0,j}, \dots, \Pi_{0,j})$$

of the descendance performances.

The simple example of sheep breeders who wish to maximise the proportion of twins, first prompted this work. A single lamb and more than three lambs are economically undesirable. The optimum is then  $\Pi_0 = (0, 1, 0, \dots, 0)$ . In the remainder of the text, the focus will be on this particular target. Obviously, generalisations are straightforward without any conceptual addition.

##### 4.2. Selection schemes

Simulated selection schemes were run 1 000 times in order to have accurate empirical responses to canalising selection. A fixed number ( $n_p$ ) of unrelated sires were mated to  $n$  unrelated dams each, producing  $n$  daughters per sire family. Each daughter had one record (litter size), and the set of  $n$  performances



in a sire family was used to evaluate this sire. Different indices were compared and are detailed later. For the likelihood-based indices, animals were treated as if they were unrelated. True variance components were used (otherwise mentioned). After sire ranking,  $n_s$  sires were selected and produce  $n_p$  males for the next generation. The selection scheme was hence the same as in SanCristobal-Gaudy *et al.* [22], except that the phenotype was not directly

$$y = \mu + u + \exp\left(\frac{\eta + v}{2}\right) \varepsilon$$

but was set to  $j$  if  $y$  lied in the interval  $[\tau_{j-1}, \tau_j]$ .

Let us denote by  $i$  the sire,  $j$  the category,  $\Pi_{ij}$  the probability that father  $i$  has daughters with a litter size equal to  $j$  for  $j$  in the  $\{1, 2, 3, 4, 5\}$  set,  $n_{ij}$  the number of daughters of sire  $i$  that have a  $j$  litter size,  $I(\mathbf{n}_i)$  the index of sire  $i$  with  $\mathbf{n}_i = (n_{i1}, \dots, n_{i5})$ ,  $\sum_{j=1}^5 n_{ij} = n$ .

Two phenotypic selection indices were considered:

$$I_{PO}(\mathbf{n}_i) = \frac{n_{i2}}{n} \quad (16)$$

the empirical estimate of  $\Pi_{i2}$ , where the index  $P$  stands for phenotypic and  $O$  denotes on the observed scale;

if the discrete trait is treated as continuous, as in [22], the index is:

$$I_{PC}(\mathbf{n}_i) = (\bar{n}_i - y_0)^2 + S_i^2, \quad (17)$$

where  $C$  stands for continuous (data are considered as such),  $\bar{n}_i$  and  $S_i^2$  are the empirical mean and variance, respectively, of  $\mathbf{n}_i$  and  $y_0 = 2$ .

Then, four selection indices were defined, using estimated breeding values  $\hat{u}_i$  and  $\hat{v}_i$  (when an heteroscedastic model is used) of sire  $i$ , on the observed ( $O$ ) or underlying ( $U$ ) scale. The estimates  $\hat{u}_i$  and  $\hat{v}_i$  are MAP estimates of breeding values (see paragraph 2.2), *i.e.* likelihood-based estimates (index  $L$ ):

$$I_{LhomO}(\mathbf{n}_i) = \Phi\left(\frac{\tau_2 - \mu - \hat{u}_i/2}{\sigma_e}\right) - \Phi\left(\frac{\tau_1 - \mu - \hat{u}_i/2}{\sigma_e}\right) \quad (18)$$

and  $\sigma_e = \sqrt{3\sigma_u^2/4 + \exp(\eta + \sigma_v^2/2)}$ , where *hom* means that the model is homoscedastic;

$$I_{LhetO}(\mathbf{n}_i) = \hat{\Pi}_{i2} = \Phi\left(\frac{\tau_2 - \mu - \hat{u}_i/2}{\hat{\sigma}_{e,i}}\right) - \Phi\left(\frac{\tau_1 - \mu - \hat{u}_i/2}{\hat{\sigma}_{e,i}}\right) \quad (19)$$

and  $\hat{\sigma}_{e,i} = \sqrt{3\sigma_u^2/4 + \exp(\eta + \hat{v}_i/2 + 3\sigma_v^2/8)}$ , where *het* means that the model is heteroscedastic;

$$I_{LhomU}(\mathbf{n}_i) = (\mu + \hat{u}_i/2 - y_0)^2, \quad (20)$$

with  $y_0 = \frac{\tau_1 + \tau_2}{2}$ ; and

$$I_{LhetU}(\mathbf{n}_i) = (\mu + \hat{u}_i/2 - y_0)^2 + (3\sigma_u^2 + \exp(\eta + \hat{v}_i/2 + 3\sigma_v^2/8)), \quad (21)$$

with  $y_0 = \frac{\tau_1 + \tau_2}{2}$ .

Particular parameters were chosen in order to mimic the Lacaune population analysed in the previous section:  $n_p = 30$ ,  $n_s = 5$ ,  $n = 30$  or  $100$ ,  $r = 0$ ,  $\sigma_u^2 = 0.64$ ,  $\sigma_v^2 = 0.25$ ,  $\mu$  and  $\eta$  such that the mean prolificacy equals 1.7 and the phenotypic variance equals 0.71,  $\tau_1 = 0.311$ ,  $\tau_2 = 2.193$ ,  $\tau_3 = 3.456$ , and  $\tau_4 = 4.637$ .

Data were also generated with  $\sigma_v^2 = 0.001$  and likelihood calculations were performed with  $\sigma_v^2 = 0.25$  and vice versa, to apprehend the impact of using a wrong model on selection efficiency.

Moreover, the model was slightly complicated by adding a fixed effect with two levels, say a *HERD* factor. Each sire  $i$  was given at generation  $t$  a proportion  $\alpha_{it}$  (resp.  $1 - \alpha_{it}$ ) of daughters in herd 1 (resp. 2), with  $\alpha_{it}$  drawn from a uniform distribution  $\mathcal{U}(0, 1)$ . The following parameterisation was adopted: the two levels had effects equal to  $a$  and  $-a$ , respectively. The particular value  $2a = 1.5$  was used in the simulations. It corresponds to a large effect encountered in the analysis of the Lacaune data.

At this point the following question arises: how can one introduce fixed effects in the index of selection when the relation between breeding values and phenotype (or index) is nonlinear? In the traditional linear case, let us denote  $\hat{\mu}_k + \hat{u}_i$  the estimated index of animal  $i$  in environment  $k$ . Evidently, the ranks of these indices do not depend on the environments. This is not the case in the threshold model since the ranks of

$$\hat{\Pi}_{2,i,k} = \Phi\left(\frac{\tau_2 - \hat{\mu}_k - \hat{u}_i}{\hat{\sigma}_{ik}}\right) - \Phi\left(\frac{\tau_1 - \hat{\mu}_k - \hat{u}_i}{\hat{\sigma}_{i,k}}\right) \quad (22)$$

do depend on environment  $k$ . In our particular case, the aim was to select sires giving the maximum of twins whatever the herd. The chosen index was

$$I_{LhetO} = \frac{1}{2}\Pi_{2,i,k=1} + \frac{1}{2}\Pi_{2,i,k=2} \quad (23)$$

since each sire has a probability of 1/2 of having a daughter in herd 1, by construction. More generally, each likelihood-based index  $I_{L*}$  of equations (18), (19), (20), and (21) is replaced by

$$\frac{1}{2}I_{L*,k=1} + \frac{1}{2}I_{L*,k=2}. \quad (24)$$

The effect of the herd was not taken into account in the phenotypic indices *PO* and *PC*.

### 4.3. Results

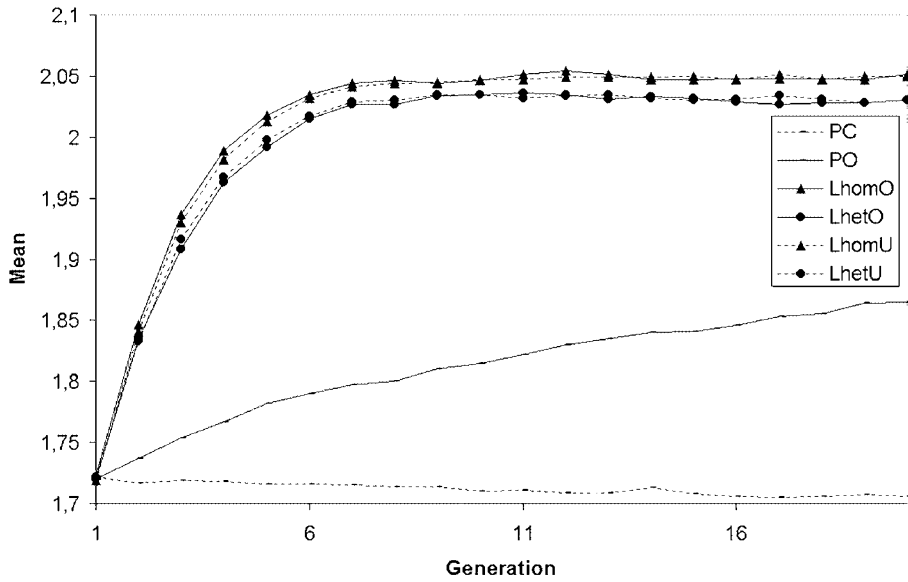
The six selection indices are compared in terms of mean prolificacy (Fig. 2), phenotypic standard deviation (Fig. 3) with the corresponding genetic progress for  $v$  (Fig. 4), and percentage of twins (Fig. 5) during 20 generations of selection, and  $n = 100$  daughters per sire. The shape of the  $u$  genetic progress is the same as the shape of the phenotypic mean in Figure 2 (not shown). Similarly, the percentage of quintuplets (not shown) behaves like the phenotypic standard deviation (Fig. 3). More importantly, the equivalence of indices corresponding to the same model, no matter the scale in which it is calculated (Observed or Underlying), is to be mentioned: *LhomO* behaves like *LhomU*, and *LhetO* like *LhetU*.

The phenotypic variance and the percentage of quintuplets are stabilised by the *PO* index, while the phenotypic mean tends very slowly towards the optimum. The *PC* index shows no progress in the mean prolificacy. This can be explained by the fact that the strong effect of the environment is not taken into account; this omission increases the residual variance and hence drastically decreases the heritability. The selection is consequently quite inefficient in moving the mean towards the target. The selection is nevertheless very efficient in decreasing the variance. In contrast the likelihood-based indices show a fast increase in the main criterion, that is the twin percentage and consequently the mean prolificacy. Because of the discrete nature of the data, the strong increase in the mean is accompanied by an increase in phenotypic variance. As soon as the population has reached the optimum on average, the phenotypic variance decreases provided that a heteroscedastic model is used (indices *LhetO* and *LhetU*). If not, the variance and the percentage of quintuplets are maintained at a high and constant level. Note that the *PC* index, also leading to a high genetic progress for  $v$  but with a lower mean than the *LhetO* and *LhetU* indices, shows a reduction in phenotypic variance.

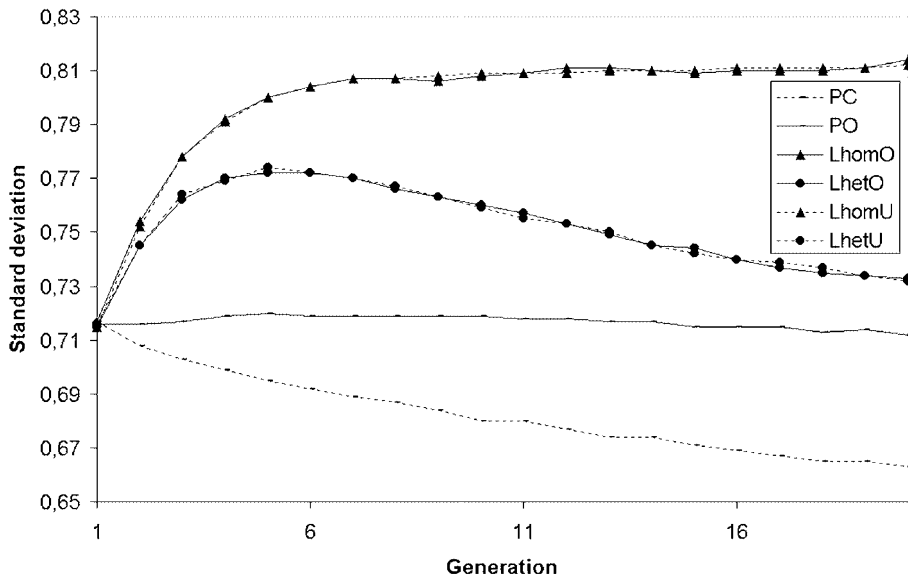
Since data are discrete, the link between the mean and variance is so strong that the underlying genetic progress in  $v$ , which is indeed high for the *LhetO* and *LhetU* indices (one genetic standard deviation gain in 10 generations of selection), is not visible on the phenotypic scale until the mean stops increasing. It is however possible to slow down the genetic progress of  $u$  in order to privilege the genetic progress of  $v$  and its phenotypic expression. This can be achieved by putting different weights in the index, like:

$$I_{LhetU}(\mathbf{n}_i) = w_1(\mu + \hat{u}_i/2 - y_0)^2 + w_2(3\sigma_u^2 + \exp(\eta + \hat{v}_i/2 + 3\sigma_v^2/8)). \quad (25)$$

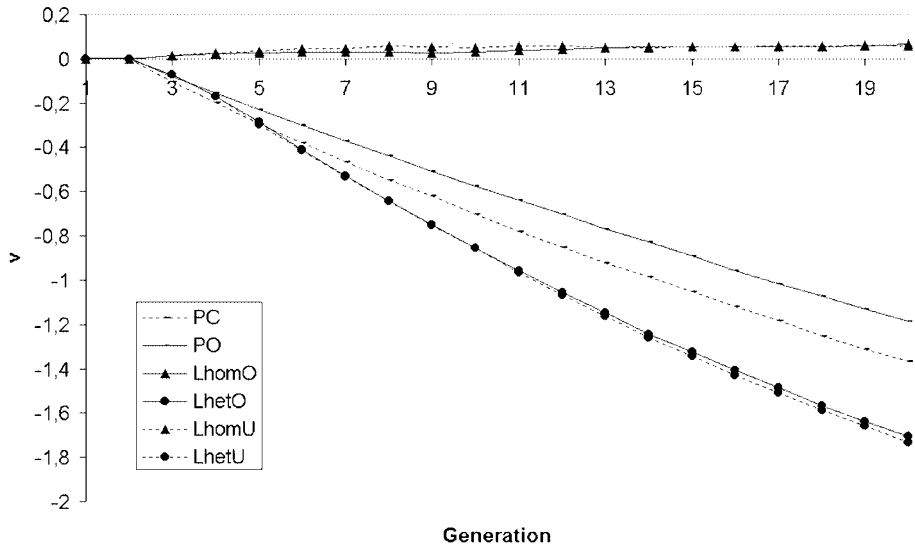
For Figure 6, the particular values  $w_1 = 1$  and  $w_2 = 50$  were chosen. Compared to the *PO* index (Fig. 6), the mean evolves faster towards the optimum, while the variance decreases, showing that the weighted index *LhetU* has the highest performances whatever the point of view (mean or variance evolution).



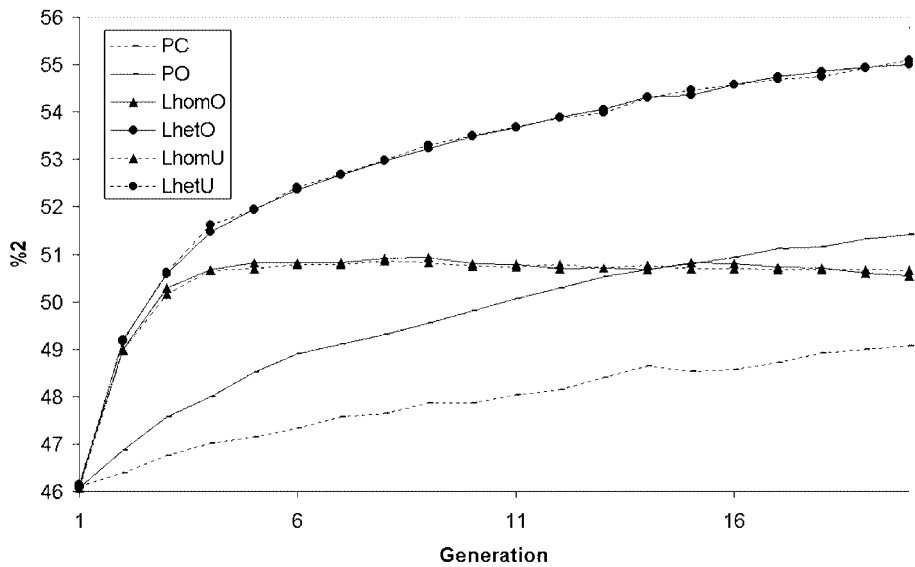
**Figure 2.** Evolution of phenotypic means for the six indices of selection. Simulations were performed with  $n_p = 30$ ,  $n_s = 5$ ,  $n = 100$ ,  $r = 0$ ,  $\sigma_u^2 = 0.64$ ,  $\sigma_v^2 = 0.25$ ,  $\mu = 0.61$ ,  $\eta = -0.6$ ,  $a = 1.5$ ,  $\tau_1 = 0.311$ ,  $\tau_2 = 2.193$ ,  $\tau_3 = 3.456$ , and  $\tau_4 = 4.637$ .



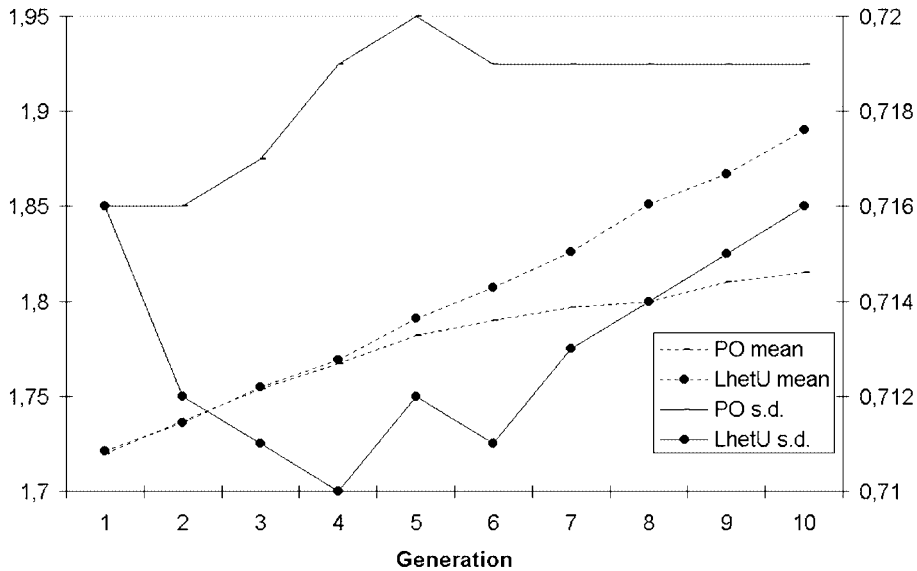
**Figure 3.** Evolution of phenotypic standard deviations for the six indices of selection. Simulation parameters as for Figure 2.



**Figure 4.** Genetic progress of  $v$  expressed in genetic standard deviation units. Simulation parameters as for Figure 2.



**Figure 5.** Evolution of twin percentages for the six indices of selection. Simulation parameters as for Figure 2.



**Figure 6.** Joint evolution of phenotypic mean and standard deviation. Indices *PO* and *LhetU* with weights 1 and 50 on mean and variance. Simulation parameters as for Figure 2.

When a parameter  $\sigma_v^2$  is set to 0.252 in the heteroscedastic model, while its true value is 0, then the selection based on the heteroscedastic indices *LhetO* or *LhetU* acts as if the genetic variance  $\sigma_v^2$  was already null, *i.e.* the indices *LhetO* or *LhetU* are quite equivalent to indices *LhomO* or *LhomU* in this case. For example, the mean prolificacy is only 3% lower with heteroscedastic than with homoscedastic models, while the phenotypic standard deviation is also 2% lower after three generations of selection. This means that the heteroscedastic approach does not slow down the efficiency of the selection if a higher genetic variance in  $v$  is wrongly put in the model.

The previous figures aimed at understanding the global long-term behaviour of some canalising selection indices. In practice, for the particular Lacaune breed, the short-term response to selection is given in Table IV in terms of mean prolificacy, phenotypic standard deviation, underlying genetic progress and percentages of single, twin, triplets, quadruplets and quintuplets or more. In this case,  $n = 30$  progeny per sire is assumed.

## 5. DISCUSSION

The first aim of this work was the analysis of the genetic components of litter size in the Lacaune sheep breed. A liability model was chosen, as is often done for the analysis of polytomous data in animal genetics. A high

**Table IV.** Performances of six selection indices.  $n = 30$ ,  $\sigma_v^2 = 0.252$ .

Gen.	Index	Average prolificacy		Standard deviation		$\Pi_1$	$\Pi_2$	$\Pi_3$	$\Pi_4$	$\Pi_5$
		Phen.	$u$	Phen.	$v$					
0		1.71	0	0.71	0	42.4	45.7	10.3	1.4	0.12
1	<i>PC</i>	1.72	0	0.71	0	41.5	46.4	10.6	1.4	0.11
	<i>PO</i>	1.74	0	0.72	0	40.6	46.7	11.0	1.6	0.13
	<i>LhomO</i>	1.84	0	0.75	0	35.3	48.7	13.5	2.2	0.21
	<i>LhetO</i>	1.82	0	0.75	0	35.4	48.7	13.2	2.1	0.19
	<i>LhomU</i>	1.83	0	0.75	0	35.5	48.6	13.4	2.3	0.20
	<i>LhetU</i>	1.82	0	0.75	0	36.0	48.6	13.1	2.1	0.20
5	<i>PC</i>	1.76	0.09	0.71	-0.14	39.1	47.9	11.3	1.5	0.12
	<i>PO</i>	1.82	0.19	0.74	-0.10	35.9	48.9	13.1	2.0	0.17
	<i>LhomO</i>	2.02	0.58	0.80	0.02	26.0	50.8	18.8	4.0	0.45
	<i>LhetO</i>	2.00	0.55	0.78	-0.10	26.1	51.5	18.5	3.6	0.34
	<i>LhomU</i>	2.02	0.58	0.80	0.02	26.1	50.7	18.8	4.0	0.46
	<i>LhetU</i>	2.00	0.55	0.78	-0.09	26.1	51.5	18.5	3.6	0.35

heritability estimate ( $\hat{h}_u^2 = 0.34$  on the underlying scale) was found for mean prolificacy. This value is greater than estimates generally found in the literature but it was observed before in this particular sheep population by Bodin *et al.* [1]. Although the structure of the data seems suitable for giving unbiased heritability estimates, according to Engel *et al.* [5] and Engel and Buist [6], some authors like Matos *et al.* [15] remark higher heritability estimates with a sire model than with an animal model for litter size. Other estimation procedures could have been chosen such as the quasi-score used by Jaffrezic *et al.* [12], or MCMC techniques. The only advantages of an EM approach are the certainty of convergence of the algorithm to a local minimum of the function to optimise, and the slight modification of the traditional REML equations. But the need for a MC step in the EM algorithm leads to heavy computations, which may tell in favour of full MCMC techniques.

The infinitesimal model proposed by SanCristobal-Gaudy *et al.* [22] for continuous traits was extended here to polytomous traits *via* a continuous underlying variable, allowing the modelling of the environmental variability as is usually done for the mean. The year, herd, season and age have no significant effects on the variability of litter size in the Lacaune population, but the sire factor has an important influence. The inclusion of the relationship matrix allows the interpretation of the sire variance  $\sigma_v^2$  of the log residual variances in the underlying scale as an additive genetic variance. The estimate of this parameter was found equal to  $\hat{\sigma}_v^2 = 0.23$ ; it corresponds to a maximum value

of the ratio of sire variances on the underlying scale equal to  $\sigma_{\text{Max}}^2/\sigma_{\text{min}}^2 = \exp(v_{\text{Max}} - v_{\text{min}}) \approx \exp(6\sigma_v) \approx 18$ , which is pretty high. At present, this value, however, has no comparison in the literature.

The second aim of this work was the prediction of the response to a selection for homogenising litter size around the target of two lambs per lambing. This problem is already complicated in standard situations, due to nonlinearity. An immediate extension of the work of Im and Gianola [11] shows that the parent-offspring regression is nonlinear for polytomous data with more than two categories. Some of the heritability estimates proposed by Magnussen and Kremer [13] cannot be extended to multiple-category data. Analytical expressions for the selection response of a binary trait given by Foulley [7] are unfortunately not feasible when a multiplicative model is set on the underlying environmental variance. The simulations performed in the previous section were imposed by these analytical complications.

Quantitatively, canalising selection is less efficient here than for continuous traits, due to the relationship between phenotypic mean and variance for discrete traits. The Lacaune situation is particularly difficult since one aspect of the objective is the increase of mean prolificacy, whose consequence (the increase of phenotypic variance) has an opposite action on the other aspect of the objective (reduction of the environmental variance). Despite a high genetic progress on the underlying environmental variance, only a small part of this is reproduced on the observed scale.

In fact, the model assumes a constant genetic variance in the mean value of the underlying variable  $Y$  and fixed threshold values that define a limit to the possible reduction in phenotypic variance, corresponding to the case in which  $\text{Var}(Y) = \sigma_u^2$ . At the limit, the expected proportions of litter sizes are equal to 0.12, 0.76, 0.11, 0.003 and  $10^{-5}$ , in increasing order. No reduction in the genetic variance was envisaged for this theoretical limit. More flexible models, derived from a physiological analysis (as in the work of Mariana *et al.* [14]), or involving the effects of QTLs or major genes on mean prolificacy, might probably be required to make such mid- and long-term predictions of the response to canalising selection more realistic.

Qualitatively, the analysed indices can be ranked on the basis of their related selection responses. In every case, the indices based on a heteroscedastic model (*LhetO* and *LhetU*) gave the best results for this criterion. A gain in the selection of categorical traits based on a threshold model over a linear model was already pointed out by Meuwissen *et al.* [17]. Moreover, the omission of an environmental factor with large effect, like the *HERD* in the simulations, has disastrous consequences on the selection, stressed by the nonlinearity between breeding values and index. Long-term figures were given in order to understand the global dynamics of certain canalising selections. So far, the selection objective had been the increase of twin proportion for the next generation.



In practice however, short- or mid-term figures are interesting for breeders. Then, generation-dependent weights in the selection indices can be envisaged, generalising the use of weights as in index (25):

$$w_{1,t}(\mu + \hat{u}_i/2 - y_0)^2 + w_{2,t}(3\sigma_u^2 + \exp(\eta + \hat{v}_i/2 + 3\sigma_v^2/8)) \quad (26)$$

or

$$\sum_{j=1,J} c_{j,t} \hat{\Pi}_{j,t} \quad (27)$$

for generation  $t$ , these weights should be chosen optimally to maximise a selection objective over  $T$  generations:

$$\sum_{t=1,T} \sum_{j=1,J} c_{0,j,t} \Pi_{0,j,t}. \quad (28)$$

To be fully comprehensive, the quantity  $\Pi_{j,t}$  in equation 27 must be calculated over all the possible levels of environment  $k$  as in (23):

$$\sum_k p_{k,t} \Pi_{k,j,t}, \quad (29)$$

where  $p_{k,t}$  is the incidence of level  $k$  in the whole population. Economic studies will estimate weights  $c_{0,j,t}$  (Benoit, personal communication).

One must note that the Lacaune population analysed in this paper has been selected for increasing the mean litter size. The observed high heterogeneity in sire variances may be due to the presence of polygenes controlling the residual variance (sensitivity to the environment), as was done in this paper. Heteroscedasticity may also be due to a major gene controlling the mean and segregating in the population, with the progeny of homozygote sires being less variable than heterozygotes. A canalising selection will favour homozygotes by reducing the variability, and pertaining polygenes will move the population mean to the optimum. The existence of such a major gene is currently being tested by Bodin *et al.* [3]. However, the genetics of reproduction traits is difficult (see for example Bodin *et al.* [2]), and no tool is currently available for fully understanding the genetic determinism of litter size variability.

## ACKNOWLEDGEMENTS

We would like to thank Christèle Robert-Granié for kindly reading the manuscript, and two referees for helpful comments.

**REFERENCES**

- [1] Bodin L., Bibé B., Blanc M.R., Ricordeau G., Genetic correlation relationship between prepuberal plasma FSH levels and reproductive performance in Lacaune ewe lambs, *Genet. Sel. Evol.* 20 (1988) 489–498.
- [2] Bodin L., Elsen J.M., Hanocq E., François D., Lajous D., Manfredi E., Mialon M.M., Boichard D., Foulley J.L., SanCristobal-Gaudy M., Teyssier J., Thimonier J., Chemineau P., Génétique de la reproduction chez les ruminants, *INRA Prod. Anim.* 12 (1999) 87–100.
- [3] Bodin L., Elsen J.M., Poivey J.P., SanCristobal-Gaudy M., Belloc J.P., Bibé B., Segregation of a major gene influencing ovulation in progeny of Lacaune meat sheep, in: 51st Annual Meeting of the European Association for Animal Production, 21–24 August 2000, Den Haag.
- [4] Bulmer M.G., *The mathematical theory of quantitative genetics*, Clarendon Press, Oxford, 1980.
- [5] Engel B., Buist W., Visscher A., Inference for threshold models with variance components from the generalized linear mixed model perspective, *Genet. Sel. Evol.* 27 (1995) 15–32.
- [6] Engel B., Buist W., Bias reduction of approximate maximum likelihood estimates for heritability in thresholds models, *Biometrics* 54 (1998) 1155–1164.
- [7] Foulley J.L., Prediction of selection response for threshold dichotomous traits, *Genetics* 132 (1992) 1187–1194.
- [8] Foulley J.L., Gianola D., Statistical analysis of ordered categorical data *via* a structural heteroskedastic threshold model, *Genet. Sel. Evol.* 28 (1996) 217–320.
- [9] Foulley J.L., Gianola D., San Cristobal M., Im S., A method for assessing extend and sources of heterogeneity of residual variances in mixed linear models, *J. Dairy Sci.* 73 (1990) 1612–1624.
- [10] Gianola D., Foulley J.L., Sire evaluation for ordered categorical data with a threshold model, *Genet. Sel. Evol.* 15 (1983) 201–224.
- [11] Im S., Gianola D., Offspring-parent regression for a binary trait, *Theor. Appl. Genet.* 75 (1988) 720–722.
- [12] Jaffrezic F., Robert-Granié C., Foulley J.L., A quasi-score approach to the analysis of ordered categorical data *via* a mixed heteroskedastic threshold model, *Genet. Sel. Evol.* 31 (1999) 301–318.
- [13] Magnussen S., Kremer A., The beta-binomial model for estimating heritabilities of binary traits, *Theor. Appl. Genet.* 91 (1995) 544–552.
- [14] Mariana J.C., Corpet F., Chevalet C., Lacker's model: control of follicular growth and ovulation in domestic species, *Acta Biotheoretica* 42 (1994) 245–262.
- [15] Matos C.A.P., Thomas D.L., Gianola D., Tempelman R.J., Young L.D., Genetic analysis of discrete reproductive traits in sheep using linear and nonlinear models: I. Estimation of genetic parameters, *J. Anim. Sci.* 75 (1997) 76–87.
- [16] Manfredi E., Foulley J.L., San Cristobal M., Gillard P., Genetic parameters for twinning in the Maine-Anjou breed, *Genet. Sel. Evol.* 23 (1991) 421–430.
- [17] Meuwissen T.H.E., Engel B., van der Werf J.H.J., Maximising selection efficiency for categorical traits, *J. Anim. Sci.* 73 (1995) 1933–1939.
- [18] Misztal I., Gianola D., Foulley J.L., Computing aspects of a nonlinear method of sire evaluation for categorical data, *J. Dairy Sci.* 72 (1989) 1557–1568.

- [19] Numerical Algorithms Group, The NAG Fortran Library Manual, NAG Ltd., Oxford, 1990.
- [20] Perret G., Bodin L., Mercadier M., Scheme for genetic improvement of reproductive abilities in Lacaune sheep, in: 43rd Annual Meeting of the EAAP, 1992, Madrid, Spain.
- [21] SanCristobal M., Foulley J.L., Manfredi E., Inference about multiplicative heteroskedastic components of variance in a mixed linear Gaussian model with an application to beef cattle breeding, *Genet. Sel. Evol.* 30 (1993) 423–451.
- [22] SanCristobal-Gaudy M., Elsen J.M., Bodin L., Chevalet C., Prediction of the response to a selection for canalisation of a continuous trait in animal breeding, *Genet. Sel. Evol.* 25 (1998) 3–30.
- [23] San Cristóbal-Gaudy M., Bodin L., Elsen J.M., Chevalet C., Selección para un óptimo: aplicación al tamaño de la camada en ovino, *ITEA 94A* (1998) 206–215.

## APPENDIX

This appendix is devoted to the parameter estimation for multinomial data. In order to shorten algebraic expressions, we define the following notations:

$$\alpha_{ij} = \frac{\tau_j - \mu_i}{\sigma_i},$$

$$\phi_{ij} = \phi(\alpha_{ij}),$$

$$\xi_i = \begin{cases} \frac{\exp\left(\mathbf{w}'_i \boldsymbol{\gamma} + \frac{3}{8} \sigma_v^2\right)}{\sigma_i^2} & \text{for a sire model} \\ 1 & \text{for an individual model} \end{cases} \quad (30)$$

$$\mathbf{t}'_i = \begin{cases} \left(\mathbf{x}'_i, \frac{1}{2} \mathbf{z}'_i\right) & \text{for a sire model} \\ (\mathbf{x}'_i, \mathbf{z}'_i) & \text{for an individual model} \end{cases} \quad (31)$$

$$\mathbf{w}'_i = \begin{cases} \left(\mathbf{p}'_i, \frac{1}{2} \mathbf{q}'_i\right) & \text{for a sire model} \\ (\mathbf{p}'_i, \mathbf{q}'_i) & \text{for an individual model} \end{cases} \quad (32)$$

where  $\phi$  is the density function of the standardised normal variable.

The maximisation of  $\mathcal{L}$  with respect to  $\boldsymbol{\zeta}$  can be achieved *via* a Fisher-scoring iterative algorithm. Each iteration  $t$  consists in solving a linear system:

$$-\left[E \frac{\partial^2 \mathcal{L}}{\partial \boldsymbol{\zeta}^2}\right]^{[t-1]} \left(\hat{\boldsymbol{\zeta}}^{[t]} - \hat{\boldsymbol{\zeta}}^{[t-1]}\right) = \left[\frac{\partial \mathcal{L}}{\partial \boldsymbol{\zeta}}\right]^{[t-1]}, \quad (33)$$

where  $E$  denotes expectation.

Here and in the following,  $\alpha_{i0}\phi_{i0}$  and  $\alpha_{iJ}\phi_{iJ}$  are replaced by their limit in  $\tau_0 \rightarrow -\infty$  and  $\tau_J \rightarrow +\infty$  respectively, *i.e.* by 0.

The Fisher-scoring algorithm requires the information matrix, which can be obtained from the Hessian matrix and the fact that (equation (1))

$$EN_{ij} = n_{i+}\Pi_{ij}. \quad (34)$$

Elements of the gradient of  $\mathcal{L}$  are equal to:

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \tau_j} &= \sum_{i=1}^I \frac{\phi_{ij}}{\sigma_i} \left( \frac{n_{ij}}{\Pi_{ij}} - \frac{n_{i,j+1}}{\Pi_{i,j+1}} \right), \text{ for } j = 1, \dots, J-1, \\ \frac{\partial \mathcal{L}}{\partial \theta} &= - \sum_{i=1}^I t_i \frac{1}{\sigma_i} \sum_{j=1}^J n_{ij} \frac{\phi_{ij} - \phi_{i,j-1}}{\Pi_{ij}} - \frac{1}{1-r^2} \left[ \Omega_{\theta}^{-} \theta - r \frac{\sigma_v}{\sigma_u} \Omega_{\gamma}^{-} \gamma \right], \\ \frac{\partial \mathcal{L}}{\partial \gamma} &= - \frac{1}{2} \sum_{i=1}^I w_i \xi_i \sum_{j=1}^J n_{ij} \frac{\alpha_{ij}\phi_{ij} - \alpha_{i,j-1}\phi_{i,j-1}}{\Pi_{ij}} - \frac{1}{1-r^2} \left[ \Omega_{\gamma}^{-} \gamma - r \frac{\sigma_u}{\sigma_v} \Omega_{\theta}^{-} \theta \right], \end{aligned} \quad (35)$$

where  $\Omega^{-}$  denotes a generalised inverse of  $\Omega$ , with

$$\Omega_{\theta} = \begin{pmatrix} 0 & 0 \\ 0 & \sigma_u^2 A \end{pmatrix} \quad (36)$$

and

$$\Omega_{\gamma} = \begin{pmatrix} 0 & 0 \\ 0 & \sigma_v^2 A \end{pmatrix}. \quad (37)$$

The results presented in [8] are a special case of these equations with  $\xi_i = 1$  and  $r = 0$ .

We present hereafter the elements of the inverse of the Fisher information matrix:

$$\begin{aligned} -E \frac{\partial^2 \mathcal{L}}{\partial \tau_j^2} &= \sum_{i=1}^I n_{i+} \frac{\phi_{ij}^2}{\sigma_i^2} \left( \frac{1}{\Pi_{ij}} + \frac{1}{\Pi_{i,j+1}} \right), \\ -E \frac{\partial^2 \mathcal{L}}{\partial \tau_j \partial \tau_{j-1}} &= - \sum_{i=1}^I n_{i+} \frac{\phi_{ij}\phi_{i,j-1}}{\Pi_{ij}\sigma_i^2}, \end{aligned}$$

$$\begin{aligned}
-E \frac{\partial^2 \mathcal{L}}{\partial \tau_j \partial \tau_k} &= 0 \text{ for } j \neq k-1, k, k+1, \\
-E \frac{\partial^2 \mathcal{L}}{\partial \tau_j \partial \theta} &= \sum_{i=1}^I t_i n_{i+} \frac{\phi_{ij}}{\sigma_i^2} \left( \frac{\phi_{i,j+1} - \phi_{ij}}{\Pi_{i,j+1}} - \frac{\phi_{ij} - \phi_{i,j-1}}{\Pi_{ij}} \right), \\
-E \frac{\partial^2 \mathcal{L}}{\partial \tau_j \partial \boldsymbol{\gamma}} &= \frac{1}{2} \sum_{i=1}^I w_i n_{i+} \xi_i \frac{\phi_{ij}}{\sigma_i} \\
&\quad \times \left( \frac{\alpha_{i,j+1} \phi_{i,j+1} - \alpha_{ij} \phi_{ij}}{\Pi_{i,j+1}} - \frac{\alpha_{ij} \phi_{ij} - \alpha_{i,j-1} \phi_{i,j-1}}{\Pi_{ij}} \right), \\
-E \frac{\partial^2 \mathcal{L}}{\partial \theta^2} &= \sum_{i=1}^I t_i t'_i \frac{1}{\sigma_i^2} n_{i+} \sum_{j=1}^J \frac{(\phi_{ij} - \phi_{i,j-1})^2}{\Pi_{ij}} + \frac{1}{1-r^2} \Omega_{\theta}^-, \\
-E \frac{\partial^2 \mathcal{L}}{\partial \boldsymbol{\gamma}^2} &= \frac{1}{4} \sum_{i=1}^I w_i w'_i n_{i+} \xi_i^2 \sum_{j=1}^J \frac{(\alpha_{ij} \phi_{ij} - \alpha_{i,j-1} \phi_{i,j-1})^2}{\Pi_{ij}} + \frac{1}{1-r^2} \Omega_{\boldsymbol{\gamma}}^-, \\
-E \frac{\partial^2 \mathcal{L}}{\partial \theta \partial \boldsymbol{\gamma}} &= \frac{1}{2} \sum_{i=1}^I t_i w'_i \frac{1}{\sigma_i} n_{i+} \xi_i \sum_{j=1}^J \frac{(\alpha_{ij} \phi_{ij} - \alpha_{i,j-1} \phi_{i,j-1})(\phi_{ij} - \phi_{i,j-1})}{\Pi_{ij}}.
\end{aligned} \tag{38}$$

### Link to the Gaussian case

As in Gianola and Foulley [10], terms appearing in the derivatives of log-likelihood  $\mathcal{L}$  have some link to the terms of the Gaussian case. For example, the parallel between

$$\left( \frac{y_i - \mu_i}{\sigma_i} \right)^2 - n_{i+}$$

(equation (14b) in Foulley *et al.* [9]) and

$$\begin{aligned}
& - \sum_{j=1}^J n_{ij} \frac{\alpha_{ij} \phi_{ij} - \alpha_{i,j-1} \phi_{i,j-1}}{\Pi_{ij}} \\
&= \sum_j n_{ij} E \left[ \left( \frac{Y_{ik} - \mu_i}{\sigma_i} \right)^2 \mid \tau_{j-1} < Y_{ik} < \tau_j \right] - n_{i+}
\end{aligned}$$

in  $\partial \mathcal{L} / \partial \boldsymbol{\gamma}$  is interesting to highlight.

Similarly, in  $\partial^2 \mathcal{L} / \partial \theta^2$ ,

$$\begin{aligned} \frac{1}{\sigma_i^2} \sum_{j=1}^J n_{ij} \left[ \frac{(\phi_{ij} - \phi_{i,j-1})^2}{\Pi_{ij}^2} + \frac{\alpha_{ij} \phi_{ij} - \alpha_{i,j-1} \phi_{i,j-1}}{\Pi_{ij}} \right] \\ = \sum_j \frac{n_{ij}}{\sigma_i^2} \left\{ 1 + E^2 \left[ \frac{Y_{ik} - \mu_i}{\sigma_i} \mid \tau_{j-1} < Y_{ik} < \tau_j \right] \right. \\ \left. - E \left[ \left( \frac{Y_{ik} - \mu_i}{\sigma_i} \right)^2 \mid \tau_{j-1} < Y_{ik} < \tau_j \right] \right\} \end{aligned}$$

corresponds to  $\frac{n_{i+}}{\sigma_i^2}$  in the continuous case, and

$$\begin{aligned} \frac{1}{4} \sum_{j=1}^J n_{ij} \left[ \frac{(\alpha_{ij} \phi_{ij} - \alpha_{i,j-1} \phi_{i,j-1})^2}{\Pi_{ij}^2} - \frac{\alpha_{ij} \phi_{ij} - \alpha_{i,j-1} \phi_{i,j-1}}{\Pi_{ij}} + \frac{\alpha_{ij}^3 \phi_{ij} - \alpha_{i,j-1}^3 \phi_{i,j-1}}{\Pi_{ij}} \right] \\ = \frac{1}{4} \sum_j n_{ij} \left\{ 2E \left[ \left( \frac{Y_{ik} - \mu_i}{\sigma_i} \right)^2 \mid \tau_{j-1} < Y_{ik} < \tau_j \right] \right. \\ \left. + E^2 \left[ \left( \frac{Y_{ik} - \mu_i}{\sigma_i} \right)^2 \mid \tau_{j-1} < Y_{ik} < \tau_j \right] - E \left[ \left( \frac{Y_{ik} - \mu_i}{\sigma_i} \right)^4 \mid \tau_{j-1} < Y_{ik} < \tau_j \right] \right\} \end{aligned}$$

to the simpler expression  $\frac{(y_i - \mu_i)^2}{2\sigma_i^2}$  in the  $\partial^2 \mathcal{L} / \partial \gamma^2$  equation for the continuous case (equation (14d) in [9]).

#### Variance component estimation

The first system (33) gives updated location parameters to solve the Fisher-scoring equations.

The second system is relative to the dispersion parameters. Newton-Raphson equations are:

$$- \left[ \frac{\partial^2 \ln p(\sigma^2 | N)}{\partial (\sigma^2)^2} \right]^{[t-1]} \left( \hat{\sigma}^{2[t]} - \hat{\sigma}^{2[t-1]} \right) = \left[ \frac{\partial \ln p(\sigma^2 | N)}{\partial \sigma^2} \right]^{[t-1]}. \quad (39)$$

It can be proven, as in [9], that the previous system can be written as

$$- \left[ E_c \frac{\partial^2 \mathcal{L}}{\partial (\sigma^2)^2} + \text{Var}_c \frac{\partial \mathcal{L}}{\partial \sigma^2} \right]^{[t-1]} \left( \hat{\sigma}^{2[t]} - \hat{\sigma}^{2[t-1]} \right) = \left[ E_c \frac{\partial \mathcal{L}}{\partial \sigma^2} \right]^{[t-1]}, \quad (40)$$

where  $E_c$  and  $\text{Var}_c$  denote expectation and variance respectively, relative to the distribution of  $\zeta | N, \hat{\sigma}^{2[t-1]}$ . A usual large sample approximation of this

distribution is given by

$$\zeta|N, \hat{\sigma}^{2[t-1]} \sim \mathcal{N} \left( \hat{\zeta}^{[t]}, \hat{C}_\zeta^{[t]} \right), \quad (41)$$

where  $\hat{\zeta}^{[t]}$  is the solution of the system (33) and  $\hat{C}_\zeta^{[t]}$  the inverse of the coefficient matrix of the same system.

The first order derivative and the second order derivative of (40) have already been calculated (see (35) and (38)). However, their conditional expectation and variance have no explicit expressions, so that numerical integration is needed to calculate the right-hand side and the coefficient matrix of the  $\zeta$  equations (40), and is clarified in the following.

$S$  values are randomly drawn from the normal distribution

$$\zeta_s \sim \mathcal{N} \left( \hat{\zeta}^{[t]}, \hat{C}_\zeta^{[t]} \right) \quad s = 1, \dots, S, \quad (42)$$

and used to get approximations

$$\begin{aligned} E_c \frac{\partial \mathcal{L}}{\partial \sigma^2} &\doteq \frac{1}{S} \sum_s \frac{\partial \mathcal{L}}{\partial \sigma^2}(\zeta_s) \\ E_c \frac{\partial^2 \mathcal{L}}{\partial (\sigma^2)^2} &\doteq \frac{1}{S} \sum_s \frac{\partial^2 \mathcal{L}}{\partial (\sigma^2)^2}(\zeta_s) \\ \text{Var}_c \frac{\partial \mathcal{L}}{\partial \sigma^2} &\doteq \frac{1}{S} \sum_s \left[ \frac{\partial \mathcal{L}}{\partial \sigma^2}(\zeta_s) \right]^2 - \left[ \frac{1}{S} \sum_s \frac{\partial \mathcal{L}}{\partial \sigma^2}(\zeta_s) \right]^2. \end{aligned} \quad (43)$$

Another possible and simpler system in  $\sigma^2$  takes only account of

$$E_c \frac{\partial^2 \mathcal{L}}{\partial (\sigma^2)^2}$$

in the coefficient matrix of (40). This produces an EM-type algorithm ([9]).

Throughout the algorithm, in order to avoid numerical problems due to null extreme categories, null probabilities  $\Pi_{ij}$  were set to a minimum value (0.01 here) like in Misztal *et al.* [18].

Programmes are written in fortran 77 using the NAG library [19] and are available on request.