

On the estimation of genetic parameters via variance components

L. DEMPFLÉ, C. HAGGER* and M. SCHNEEBERGER**

Lehrstuhl für Tierzucht der TU München, D-8050 Freising-Weihenstephan, Germany

**Institut of Animal Production, Swiss Federal Institut of Technology, CH-8092 Zurich, Switzerland*

***Herd Book Office for Swiss Braunvieh, CH-6300 Zug, Switzerland*

Summary

Variance components have been estimated by three methods using two different but overlapping data sets from a dairy cattle breeding scheme. The methods were HENDERSON's method III, MINQUE and a new method proposed by HENDERSON in 1980. Two different statistical models of grouping sires were considered. For all methods, the exact variances of the estimators were calculated for given true variance components and assuming normality of the data. As a byproduct, the large sample variances of REML were obtained. A short discussion of the interpretation of the two estimated variance components is given for the two statistical models taking selection into account. A concise description is given of the three estimation methods employed. For a relatively simple model, it is shown that they use different weighting factors for combining means and squares.

The new method proposed by HENDERSON (1980) has two possible disadvantages, namely fewer degrees of freedom for estimating the error variance and one deriving from the relationship with the method of contemporary comparison. From this limited investigation, it is concluded that, in situations where the method might be employed, these disadvantages may not be of great importance. The numerical results of the estimation with the two statistical models lie reasonably well within the expected range. A noteworthy difference in efficiency was found between MINQUE and HENDERSON's method III in favour of MINQUE, given that a reasonable prior estimate of the ratio of the error component to the sire variance component was used in the estimation. As expected, the new method was often inferior to MINQUE but it always retained a surprisingly high efficiency relative to MINQUE for the estimation of the additive genetic variance and the heritability. It is concluded that in situations where MINQUE is very difficult or impossible to compute, the new method appears to be a useful alternative.

Key-words : Efficiency, variance components, genetic parameters, MINQUE, HENDERSON III/IV.

Résumé

Contribution à l'étude de l'estimation des paramètres génétiques par les composantes de la variance

Trois méthodes d'estimations des composantes de la variance ont été testées sur deux échantillons (en partie communs) provenant d'un schéma de sélection de bovins laitiers. La comparaison concernait la méthode III d'HENDERSON, le MINQUE et une nouvelle méthode proposée par HENDERSON en 1980. Deux modèles statistiques de groupage des pères ont été également considérés. Dans tous les cas, on a calculé les variances exactes des estimateurs pour des valeurs données de composantes vraies en supposant la normalité des données. Par extension, on en a déduit les variances du REML pour de grands échantillons. On a discuté également l'interprétation des estimations pour les deux modèles statistiques en prenant en compte des phénomènes de sélection. Les trois méthodes sont décrites brièvement. Partant d'un modèle simple, on montre qu'elles diffèrent par les coefficients de pondération des moyennes et des carrés.

La nouvelle méthode d'HENDERSON présente deux inconvénients possibles, à savoir un moindre nombre de degrés de liberté pour estimer la variance d'erreur et une relation avec la méthode de comparaison aux contemporains. De cette étude limitée, il ressort, toutefois, que ces inconvénients seraient de peu d'importance dans les situations courantes d'application de la méthode. Les résultats numériques relatifs aux deux modèles correspondent assez bien à la gamme de valeurs attendues. Une différence appréciable a été observée en faveur du MINQUE, dans l'efficacité de celui-ci par rapport à celle de la méthode III d'HENDERSON sous réserve d'une valeur satisfaisante de départ du rapport de la variance d'erreur à celle du père. Comme prévu, la nouvelle méthode d'HENDERSON est fréquemment inférieure au MINQUE, mais s'avère étonnamment compétitive en vue de l'estimation de la variance génétique additive et de l'héritabilité. C'est pourquoi, elle doit être considérée comme une alternative intéressante quand le MINQUE devient difficile, voire impossible à calculer.

Mots-clés : Efficacité, composantes de la variance, paramètres génétiques, MINQUE, HENDERSON III/IV.

I. Introduction

This investigation arose from a larger project with the aim of obtaining estimates of genetic parameters for the Swiss Braunvieh population. In this population a heavy amount of crossing with US-Brown-Swiss is practised. Thus, the variance components were estimated separately for three data sets:

- i) offspring of pure Braunvieh sires, born 1971-1972;
- ii) offspring of pure Braunvieh sires, born 1973-1975;
- iii) and offspring of F_1 bulls, born 1972-1975.

The methods used were Maximum Likelihood (ML), Restricted Maximum Likelihood (REML), Minimum Norm Quadratic Unbiased Estimation (MINQUE) and Henderson's method III (H III), (HARTLEY & RAO, 1967; PATTERSON & THOMPSON, 1971; RAO, 1970, 1972; HENDERSON, 1953). For MINQUE and H III the exact variances of the estimators (for given true variance components) were calculated and the large sample variances of REML were obtained as a byproduct. The main results of this study are given elsewhere (HAGGER *et al.*, 1982).

In this paper we concentrate on the smallest data set, dealing only with the F_1 bulls born between 1972 and 1975. With this data set we estimated variance (and covariance) components for milk yield, percent fat (fat %) and percent protein (prot %) using two overlapping data sets, two different statistical models and three estimation procedures, namely MINQUE, H III and a new method proposed by HENDERSON (1980) which in the present paper is called Henderson's method IV (H IV). For all methods used, the estimates as well as their exact variances (for given true variance components and assuming normality) were obtained. Some results on REML were again obtained as a byproduct.

Because the data set is fairly typical for many situations in Central Europe, the main objective was to determine the relative efficiency of the methods, e.g. is it really worthwhile changing from H III to MINQUE? The main criterion for judging this question was the precision achievable (variance of the estimators) by these three unbiased methods. In practice, however, the ease of computing the estimates is also of great importance, whereas the ease of calculating the variances of the estimators is rather unimportant. For practical use a rough estimate of this variance should be sufficient, since we only want to decide whether the estimate should either be ignored (variance very large), or should be used as obtained (variance rather small) or should be combined with other estimates from the literature. In the last case the reciprocals of the variances should be used as weighting factors, but even for this purpose rough estimates should be sufficient.

II. Material and Methods

A. Data set

The data consisted of first lactation records collected from 1978 to 1981. Two overlapping data sets were used. Data set 1 included all daughter records from F_1 bulls having more than 7 daughters whereas data set 2 included all daughter records from F_1 bulls having more than 19 daughters. All bulls were born between 1972 and 1975. Incomplete lactations of 80 to 269 days of cows sold were extended to 305 days by multiplicative factors. Lactation yields were also precorrected multiplicatively for age at calving, days open and additively for alpine pasturing.

TABLE I
Structure of data sets.
Structure des données.

								Data set	
								1	2
Number of observations								19386	17514
Number of region \times herdclass \times year \times seasons								2317	2232
Number of sire groups Model I								4	4
Model II								17	15
Number of sires								293	133
Distribution of daughter numbers									
no of daughters	8-9	10-19	20-29	30-39	40-49	50-59	60-69		
no of bulls	42	118	31	12	6	3	5		
no of daughters	70-79	80-89	90-99	100-149	150-514				
no of bulls	4	2	3	20	47				

B. Statistical models and aspects of selected populations

The following statistical models were used:

$$Y_{ijkl} = h_i + g_j + u_{jk} + e_{ijkl}$$

$$y = X_1 h + X_2 g + Z u + e = X \beta + Z u + e$$

where

- y** is a vector of observations (one trait at a time);
 - h** is a vector of unknown fixed region \times herdclass \times year \times season effects; these effects are used as an equivalent to the more customary herd \times year \times season effects.
 - g** is a vector of unknown fixed sire group effects
 - u** is a vector of random sire effects
 - e** is a vector of random residuals
- X, Z are known design matrices, relating β and u to y .

The difference between the two models lies in the definition of the sire groups.

In *model I* sires born in the same year were assembled in one group, giving 4 groups altogether.

In *model II* groups were formed by grandsires, i.e. paternal half sibs were assembled in one group, giving 17 groups for data set 1 and 15 groups for data set 2.

The following assumptions were made:

$$\begin{aligned} E(\mathbf{u}) &= \mathbf{0} & \text{Var}(\mathbf{u}) &= I\sigma_u^2 \\ E(\mathbf{e}) &= \mathbf{0} & \text{Var}(\mathbf{e}) &= I\sigma_e^2 \\ E(\mathbf{y}) &= \mathbf{X}\boldsymbol{\beta} & \text{Var}(\mathbf{y}) &= \mathbf{Z}\mathbf{Z}'\sigma_u^2 + I\sigma_e^2 = \mathbf{V}. \end{aligned}$$

For calculating the variances of the estimators, it was assumed that \mathbf{e} and \mathbf{u} were independently normally distributed. The vectors of fixed effects are of no interest in our analysis (they are, apart from the definition of sire groups, mere nuisance factors). In the two models the sire effect u_{jk} has different meanings. In model II it is the deviation of the transmitting ability from the true paternal half sib mean, whereas in model I it is the deviation of the transmitting ability from the true average transmitting ability of all bulls born in the same year.

In model II the assumption of independently distributed sire effects $\text{Var}(\mathbf{u}) = I\sigma_u^2$ should be correct (apart from small maternal relationships), whereas with model I certain existing relationships (paternal halfsibs) are ignored. With model I this results in an underestimation of the sire variance. However, in addition to the last mentioned facts, the interpretation of the parameters depends not only on the model but also on the history of the population (BULMER, 1971; DEMPFLER, 1975) as outlined.

If we symbolize the additive genetic variance and the phenotypic variance of the (conceptual) *random mating base population* by σ_A^2 and σ_P^2 ($\sigma_E^2 = \sigma_P^2 - \sigma_A^2$), we have for

$$\begin{array}{rcc} & \frac{\sigma_e^2}{\sigma_E^2 + \frac{3}{4}\sigma_A^2 K} & \frac{\sigma_u^2}{\frac{1}{4}\sigma_A^2 K_I} \\ \text{model I} & & \\ & \frac{\sigma_e^2}{\sigma_E^2 + \frac{3}{4}\sigma_A^2 K} & \frac{3}{16}\sigma_A^2 K_{II} \\ \text{model II} & & \end{array}$$

In the base population we have $K = K_I = K_{II} = 1$. After one generation of truncation selection, where selection is characterized by intensity i , truncation point x and precision ρ , and where the paths are indicated by BB, BC, CB, CC (BC – Bull to Cow, etc.) we get:

$$K = 1 - \frac{1}{3} \rho_{CC}^2 i_{CC} (i_{CC} - x_{CC})$$

$$K_I = 1 - \frac{1}{4} [\rho_{BB}^2 i_{BB} (i_{BB} - x_{BB}) + \rho_{CB}^2 i_{CB} (i_{CB} - x_{CB})]$$

$$K_{II} = 1 - \frac{1}{3} \rho_{CB}^2 i_{CB} (i_{CB} - x_{CB}).$$

After repeated cycles of selection the K -values decrease further and reach an asymptotic value, but even in the extreme case ($\rho^2 i(i-x) \rightarrow 1$) we have $K \geq \frac{2}{3}$; $K_I \geq \frac{1}{2}$;

$$K_{II} \geq \frac{2}{3}.$$

To give an example: a simple well organised selection scheme for milk yield is assumed with $h^2 = 0.25$ in the base population and with selection operating only on first lactation. 70 % of the cows are bred to produce replacement heifers and 0.2 % are bred to produce bulls. The great majority of cows is either sired by selected sires or by test sires. 100 bulls are tested each year on 100 daughter records and the best 5 bulls are then used. For this example Table 2 shows the evolution of K values. These values are only approximate, since it is assumed that even after repeated cycles of selection the breeding values are still normally and independently distributed and that selection is done by truncation and not by the more realistic censoring.

TABLE 2
Evolution of K values.
Évolution des coefficients K .

Generation	Cows sired mainly by					
	proven bulls			test bulls		
	K	K_1	K_{11}	K	K_1	K_{11}
0	1	1	1	1	1	1
1	0.958	0.755	0.923	0.958	0.755	0.923
2	0.900	0.711	0.878	0.950	0.741	0.917
3	0.885	0.700	0.866	0.932	0.730	0.903
4	0.882	0.697	0.863	0.927	0.726	0.900
equilibrium	0.881	0.696	0.862	0.925	0.725	0.898

C. Methods of estimation

Three statistical methods were used, MINQUE, H III and H IV. For MINQUE we have to calculate (notation as given in last section):

$$\begin{aligned}\hat{\sigma}^2 &= S^{-1}q_M & \hat{\sigma}^{2'} &= [\hat{\sigma}_e^2 : \hat{\sigma}_a^2] \\ s_{ij} &= \text{tr}(PV_iPV_j) & i, j &= 0, 1 \\ q_{iM} &= y'PV_iPy \\ P &= \hat{V}^{-1}[I - X(X'\hat{V}^{-1}X)^{-1}X'\hat{V}^{-1}] \\ V_0 &= I; V_1 = ZZ'.\end{aligned}$$

Properties of the estimators are:

$$\begin{aligned}E(\hat{\sigma}^2) &= \sigma^2 \\ \text{Var}(\hat{\sigma}^2) &= S^{-1} \text{Var}(q) S^{-1'} = S^{-1} \Sigma S^{-1'} \\ \Sigma_{ij} &= 2 \text{tr}(PV_iPVPV_jPV).\end{aligned}$$

\hat{V} is proportional to $ZZ' + \bar{\lambda}I$, where $\bar{\lambda}$ is any positive operational value used in the computation. $\bar{\lambda}$ should be as close as possible to the true ratio of σ_e^2/σ_a^2 .

For H III we have to calculate:

$$\begin{aligned}\hat{\sigma}^2 &= S^{-1}q_{H \text{ III}} \\ q_{0H \text{ III}} &= y'[I - W(W'W)^{-1}W']y; \quad W = [X : Z] \\ q_{1H \text{ III}} &= y'[W(W'W)^{-1}W' - X(X'X)^{-1}X']y \\ s_{00} &= n - r(W); \quad s_{01} = 0; \quad s_{10} = r(Z) - 1 \\ s_{11} &= n - \text{tr} Z'X(X'X)^{-1}X'Z.\end{aligned}$$

The formulae for $\text{Var}(\hat{\sigma}^2)$ are similar to the ones given for MINQUE.

In order to describe H IV, the following observation is of importance: HENDERSON (1972) pointed out that there is a connection between BLUP and MINQUE via the Mixed Model Equations (MME), which is useful for both understanding and computation.

Writing the MME for the model used, we have

$$\begin{bmatrix} X'X & X'Z \\ Z'X & Z'Z + \tilde{\lambda}I \end{bmatrix} \begin{bmatrix} \hat{\beta} \\ \hat{u} \end{bmatrix} = \begin{bmatrix} X'y \\ Z'y \end{bmatrix}. \quad (1)$$

Defining $\hat{e} = y - X\hat{\beta} - Z\hat{u}$ it can be shown that apart from scalars, we have with MINQUE:

$$\begin{aligned} q_{0M} &= \hat{e}'\hat{e} \\ q_{1M} &= \hat{u}'\hat{u}. \end{aligned}$$

In H IV we make use of Eq.(1) and absorb all fixed effects, which leads to :

$$[Z'FZ + \tilde{\lambda}I]\hat{u} = Z'Fy; \quad F = I - X(X'X)^{-1}X'$$

Then the coefficient matrix is replaced by a matrix with diagonal elements identical to those of $Z'FZ + \lambda I$ and with off-diagonal elements equal to zero. This is symbolized by

$$\text{Diag}[Z'FZ + \tilde{\lambda}I]\hat{u} = Z'Fy; \quad D\hat{u} = Z'Fy.$$

The solution for \hat{u} is easy to compute and is used to calculate the following quadratic form:

$$q_{1H IV} = \tilde{u}'\hat{u}.$$

This quadratic form is set equal to its expected value. A second quadratic form for estimating σ_e^2 is needed and it is suggested that «any logical estimator of σ_e^2 , for example the within smallest subclass mean squares» (HENDERSON, 1980) should be utilized. The latter is undoubtedly very easy to compute but there may be other simple estimators which are more efficient.

A solution for \hat{u} can also be obtained directly if Eq.(1) is modified in the following way:

$$\begin{bmatrix} X'X & X'Z \\ Z'X & Z'Z + \tilde{\lambda}I + H \end{bmatrix} \begin{bmatrix} \tilde{\beta} \\ \tilde{u} \end{bmatrix} = \begin{bmatrix} X'y \\ Z'y \end{bmatrix} \quad (2)$$

with

$$H = Z'X(X'X)^{-1}X'Z - \text{Diag}(Z'X(X'X)^{-1}X'Z).$$

D. Computational aspects

For data sets like the one described in Table 1, or larger ones, the computational aspects become very dominant. For all three procedures Eq.(1) was the starting point where, during reading in the sorted data, the region \times herdclass \times year \times season effects were absorbed and other necessary quantities were calculated. Then for MINQUE and H IV an operational $\tilde{\lambda}$ was added to the diagonal elements and \hat{u} was estimated. Using the following notation

$$\begin{bmatrix} X'X & X'Z \\ Z'X & Z'Z + \tilde{\lambda}I \end{bmatrix}^{-1} = \begin{bmatrix} T_{00} & T_{01} \\ T_{10} & T \end{bmatrix}$$

$$Q = I - \tilde{\lambda}T$$

it is well known that T can be calculated from the absorbed set of equations.

For MINQUE the expected values of $\widehat{\mathbf{e}}\widehat{\mathbf{e}}$ and $\widehat{\mathbf{u}}\widehat{\mathbf{u}}$ are calculated and the variances and covariances of $\widehat{\mathbf{e}}\widehat{\mathbf{e}}$ and $\widehat{\mathbf{u}}\widehat{\mathbf{u}}$ are given by:

$$\text{Var}(\widehat{\mathbf{u}}\widehat{\mathbf{u}}) = \sigma_u^4 2 \tilde{\lambda}^0 [\text{tr } \mathbf{Q}^4 + 2\lambda \text{tr } \mathbf{TQ}^3 + \lambda^2 \text{tr } \mathbf{T}^2\mathbf{Q}^2]$$

$$\text{Cov}(\widehat{\mathbf{e}}\widehat{\mathbf{e}}, \widehat{\mathbf{u}}\widehat{\mathbf{u}}) = \sigma_u^4 2 \tilde{\lambda}^2 [\text{tr } \mathbf{TQ}^3 + 2\lambda \text{tr } \mathbf{T}^2\mathbf{Q}^2 + \lambda^2 \text{tr } \mathbf{T}^3\mathbf{Q}]$$

$$\text{Var}(\widehat{\mathbf{e}}\widehat{\mathbf{e}}) = \sigma_u^4 2 \tilde{\lambda}^4 [\text{tr } \mathbf{T}^2\mathbf{Q}^2 + 2\lambda \text{tr } \mathbf{T}^3\mathbf{Q} + \lambda^2 \text{tr } \mathbf{T}^4] + \sigma_u^4 2\lambda^2 [n - r(\text{Eq. (1)})].$$

Having computed $\widehat{\mathbf{e}}\widehat{\mathbf{e}}$ and $\widehat{\mathbf{u}}\widehat{\mathbf{u}}$ with a given operational value of λ , then the true variances can be calculated with these formulae for a range of true λ values. A similar approach was taken for H III and H IV where well known formulae were used.

E. Comparison and discussion of the methods

Before reporting the numerical results, a general discussion of the methods is useful. For discussion the most simple setting is used because otherwise the formulae are too complex to give much insight.

Using the one factor model

$$\mathbf{y} = \mathbf{1}\mu + \mathbf{Z}\mathbf{u} + \mathbf{e}$$

the quadratic forms which are calculated for H III (H III in this case is identical to H I) are :

$$q_{0_{HIII}} = \sum_i \sum_j (y_{ij} - \bar{y}_i.)^2 = \sum_i \sum_j (y_{ij} - \widehat{\mu} + \widehat{u}_i^+)^2$$

$$q_{1_{HIII}} = \sum_i n_i (\bar{y}_i. - \bar{y}..) ^2 = \sum_i w_i^+ (\bar{y}_i. - \widehat{\mu}^+)^2$$

with
$$\widehat{\mu} + \widehat{u}_i^+ = \bar{y}_i \quad \text{and} \quad \widehat{\mu}^+ = \frac{1}{n} \sum_i n_i \bar{y}_i.$$

For MINQUE we calculate:

$$q_{0_M} = \sum_i \sum_j (y_{ij} - \widehat{\mu} + \widehat{u}_i^*)^2$$

$$q_{1_M} = \sum_i \left(\frac{n_i}{n_i + \lambda} \right)^2 (\bar{y}_i. - \widehat{\mu}^*)^2 = \sum_i w_i^* (\bar{y}_i. - \widehat{\mu}^*)^2$$

with

$$\widehat{\mu}^* = \frac{1}{\sum_i \frac{n_i}{n_i + \lambda}} \sum_i \frac{n_i}{n_i + \lambda} \bar{y}_i.$$

$$\widehat{\mu} + \widehat{u}_i^* = \widehat{\mu}^* + \frac{n_i}{n_i + \lambda} (\bar{y}_i. - \widehat{\mu}^*).$$

For H IV use is made of Eq.(2) where we calculate (only q_1 is specified)

$$q_{1_{HIV}} = \sum_i \left(\frac{n_i}{n_i a_i + \lambda} \right)^2 (\bar{y}_i. - \bar{y}..) ^2 = \sum_i w_i^{**} (\bar{y}_i. - \widehat{\mu}^+)^2$$

$$0 < a_i = \frac{n. - n_i}{n.} < 1.$$

Thus, with H III the LS estimate of $\mu + u_i$ regarding u_i as fixed is used for q_0 . For q_1 the LS estimate of μ ignoring u_i is used and the squares are weighted by n_i , the number of observations in group i .

With MINQUE we use the BLUP estimate of $\mu + u_i$ for q_0 and the BLUE estimate (GLS estimate regarding u_i as random) of μ for q_1 and $(n_i/(n_i + \bar{\lambda}))^2$ as weighting factor. If $\bar{\lambda}$ is zero (implying no variation within sires) the square of each sire is equally weighted, regardless of n_i , which is completely in agreement with intuition. If $\bar{\lambda}$ is very large, each square has a weight proportional to the square of n_i . Thus, depending on $\bar{\lambda}$ the weights of the squares can vary from being proportional to 1 up to n_i^2 . For a given distribution of n_i there should be a $\bar{\lambda}$ where the weights of MINQUE are in similar proportion but not identical to n_i , the weights used in H III. For the same model a discussion of the weightings of the squares (using always $\hat{\mu}^+$) being in agreement with the above mentioned results, but using the F-value of the Analysis of Variance instead of $\bar{\lambda}$, was presented by ROBERTSON (1962).

It should be further noted that, if μ were known, then the weights used in MINQUE for q_1 are proportional to the reciprocals of the variance of the squares, and therefore well known weighting factors are used to combine these squares.

With H IV the LS estimate of μ is used (as in H III), whereas the weights are similar but not identical to those of MINQUE.

With regard to H IV several comments can be made:

i) Methods that have a high efficiency relative to MINQUE and that are easier to compute are very desirable and urgently needed.

ii) Using the obvious estimator for σ_c^2 (the within smallest subclass mean squares) quite a lot of available information may not be utilized. Consider the simple model in sire evaluation

$$y_{ij} = \mu + h_i + u_j + e_{ijk}.$$

If there is a total of n daughter records from n_u sires which are distributed over n_h herds, then, with H III $n - n_u - n_h + 1$ degrees of freedom (df) are used to estimate σ_c^2 . A similar number of df is used by MINQUE. For the obvious estimator only $n - c$ df are used ($c \sim$ number of filled subclasses). In the extreme case of a completely balanced block design we have $(n_h - 1)(n_u - 1)$ df for H III and zero for the obvious estimator, since there is only one observation in each smallest subclass. In a typical dairy sire evaluation scheme there may be few half-sibs in a herd \times year \times season, which would lead to a drastic reduction in df. Even in our example using region \times herdclass \times year \times season we had 16777 df (15150 df) in data set 1 (data set 2) for H III and only 7395 df (6808 df) for the obvious estimator, resulting in the error-variance of $\hat{\sigma}_c^2$ being more than 2.2 times larger than with H III. As already mentioned, other estimators for σ_c^2 than the «obvious» one could be used, like the H III estimator or the MINQUE estimator (e.g. with $\bar{\lambda} \rightarrow \infty$). However, as can be seen from fig. 1, the MINQUE estimator for $\bar{\lambda} \rightarrow \infty$ (sometimes referred to as MINQUE(0)) can be very inefficient; whereas the H III estimator always has a high efficiency. Choosing a different estimator than the obvious one, it should still be easy to compute, since this is the only justification for changing from MINQUE to H IV.

iii) In a progeny testing situation, where β contains only fixed herd effects (herd \times year \times season) and u the transmitting abilities, the solutions of u are the Contemporary Comparison (CC) estimates as was pointed out by POWELL & FREEMAN (1974). In sire evaluation there were good reasons to move away from CC and use more sophisticated methods. The question is whether the disadvantages of the CC

method are carried over to H IV. One major disadvantage of the CC method lies in the fact that the competition, a sire has in a certain herd is not taken into account. It is implicitly assumed that the mean of competing sires is the same in all herds. However, if we have several subpopulations the effects of the subpopulations (the group effects) are accounted for in H IV. In the context of estimating variance components we must always have a random sample of sires and the daughters of these sires should be distributed randomly over the herds. In this case we would expect that the disadvantages of the CC method would not be of great importance in the estimation of variance components. In order to investigate if there could be more bias with H IV than with MINQUE or H III, the following example was considered: there is a number of herds available, which are considered as fixed, thus no further assumptions about them need to be made. A random sample of sires is drawn out of a well defined population. Given that bulls were mated randomly over herds, without any assortative mating and without any preferential treatment of the daughters, we would have good conditions for estimating variance components unbiased. However, what happens if after drawing a random sample of bulls, we get some information on them and order these bulls according to this information (consider the trait type score at the age of one year, where we could have a random sample of male calves, conduct a performance test and then use all bulls in a progeny testing scheme for the same trait, allowing farmers the choice of bulls). If we relabel the bulls according to the ordering (1 labelling the bull with the highest order) we no longer have $E(\mathbf{u}) = \mathbf{0}$ and $\text{Var}(\mathbf{u}) = I\sigma_u^2$ but we have instead $E(\mathbf{u}) = \rho\mu_0\sigma_u$ and $\text{Var}(\mathbf{u}) = (1 - \rho^2)I\sigma_u^2 + \rho^2V_0\sigma_u^2$ where ρ is the correlation between the true sire value and the information on which the ordering is based. μ_0 is the vector of expected values for order-statistics from the unit normal distribution and V_0 is likewise the variance-covariance matrix of the vector of order-statistics. The values for μ_0 and V_0 are given e.g. by SARHAN & GREENBERG (1962, p. 193) and the formulae for $E(\mathbf{u})$ and $\text{Var}(\mathbf{u})$ are standard results for associate variables (DAVID, 1970, p. 41). Now in the dairy industry, it is not unlikely that some farmers use only the «very best testbulls» whereas others use average or even below average bulls. This may even apply to a trait like milk yield.

With all three methods considered, we compute quadratic forms, and in the standard case set these equal to the expected values derived under the assumption of $E(\mathbf{u}) = \mathbf{0}$, $\text{Var}(\mathbf{u}) = I\sigma_u^2$. In the example it is possible to derive the expectation under the condition of ordering and nonrandom use of the sires and thus the bias can be calculated. Some results are given in Table 3. From the few cases investigated out of the large number of conceivable ones it seems that with larger daughter number the bias of H IV is somewhat larger than with MINQUE and that H III is more robust against this departure from the usual assumptions. It is well known (SEARLE, 1968) that H III gives unbiased estimates of the variance components if there are nonzero covariances between the factors of the model. However, the case investigated here, is different, because there is essentially a correlation between the sires of the same herd. Knowing the value of one sire utilised in a herd enables one to make informative predictions about the other sires used in the same herd. In the standard application of H III the expectation is taken under the assumption of $\text{Var}(\mathbf{u}) = I\sigma_u^2$ which does not apply for this example. However, from this limited inference, these results cannot be used as a strong argument against H IV in comparison to MINQUE.

TABLE 3
Bias of the three methods (H III, MINQUE, H IV) resulting from non-random use of bulls given as multiples of σ_u^2 ($\lambda = 15$).
Biais des trois méthodes (H III, MINQUE, H IV) dû à une utilisation non aléatoire des taureaux ($\lambda = 15$) exprimé en multiples de σ_u^2 .

Design - N*	p	bias in σ_u^2 H III	MINQUE	H IV	bias in σ_u^2 H III	MINQUE	H IV
1, [a a a a a]	any	0	0	0	0	0	0
2, [2 2 2 0 0] [0 0 2 2 2]	0 0.2 0.5 0.8 1.0	0 0 0 0 0	0 0.009 0.055 0.141 0.220	0 0 0 0 0	0 -0.020 -0.126 -0.323 -0.504	0 0 -0.031 -0.194 -0.498 -0.778	0 0 -0.028 -0.175 -0.448 -0.700
3, [8 8 8 0 0] [0 0 8 8 8]	0 0.2 0.5 0.8 1.0	0 0 0 0 0	0 0.004 0.025 0.064 0.099	0 0 0 0 0	0 -0.020 -0.126 -0.323 -0.504	0 -0.025 -0.157 -0.401 -0.627	0 -0.027 -0.169 -0.433 -0.676
4, [16 16 16 0 0] [0 0 16 16 16]	0 0.2 0.5 0.8 1.0	0 0 0 0 0	0 0.003 0.017 0.042 0.066	0 0 0 0 0	0 -0.020 -0.126 -0.323 -0.504	0 -0.022 -0.137 -0.350 -0.548	0 -0.026 -0.165 -0.422 -0.660
5, [16 16 16 16 16 16 16 16]	0 0.5 1.0	0 0 0	0 0.013 0.051	0 0 0	0 -0.176 -0.704	0 -0.192 -0.768	0 -0.207 -0.830

*In N rows represent the herds (\times year \times season) columns represent the ordered sires the entries give the number of daughters.

a any value greater 1.

p correlation between true sire value and information, on which the ordering is base.

III. Results and Discussion

A. Influence of the models on heritability estimates

Whereas with H III only one result is obtained, with MINQUE and H IV a multitude of results are obtained depending on the values of $\bar{\lambda}$ used. The heritability estimates for $\bar{\lambda} = 15$ for milk yield and $\bar{\lambda} = 9$ for fat % and protein % are reported (Table 4). The variance of these estimators from Model II is indicated in the last section in connection with the figures 7 and 8. The variance from Model I is somewhat smaller. The \hat{h}^2 were estimated under the assumption that $K = K_I = K_{II} = 1$. The resulting estimates for σ_A^2 (milk yield, MINQUE, data set 2) are 117751 kg² for model I and 138232 kg² for model II leading to an estimate of K_I/K_{II} of 0.85 which is well within the expected range.

TABLE 4

Estimates of heritabilities.
Estimations de l'héritabilité.

Data set 1	Milk kg	Fat %	Protein %
Model I			
H III	0.32	0.37	0.48
MINQUE	0.34	0.40	0.55
H IV	0.35	0.41	0.57
Model II			
H III	0.39	0.45	0.53
MINQUE	0.40	0.45	0.65
H IV	0.41	0.46	0.69
Data set 2			
Model I			
H III	0.30	0.36	0.46
MINQUE	0.30	0.35	0.51
H IV	0.31	0.36	0.52
Model II			
H III	0.38	0.44	0.49
MINQUE	0.35	0.41	0.59
H IV	0.36	0.42	0.65

Now the question is which \hat{h}^2 to use in practical situations e.g. for estimating sires. This depends again on the model used. If we have a model like model I (sires grouped by year, no relationship matrix) then from a bayesian point of view the applicable \hat{h}^2 is that from model I, since it parameterizes best the *a priori* distribution of the transmitting ability of test bulls. If, on the contrary, we use the full numerator relationship matrix relative to the base population, the parameters of the base population should be used and thus, the estimates from model II are more appropriate. However, in theory they still underestimate the parameters of the base population since K and K_{II} not being unity is not accounted for in the estimation. In practice, however, it may be very difficult to determine those coefficients with any reasonable precision.

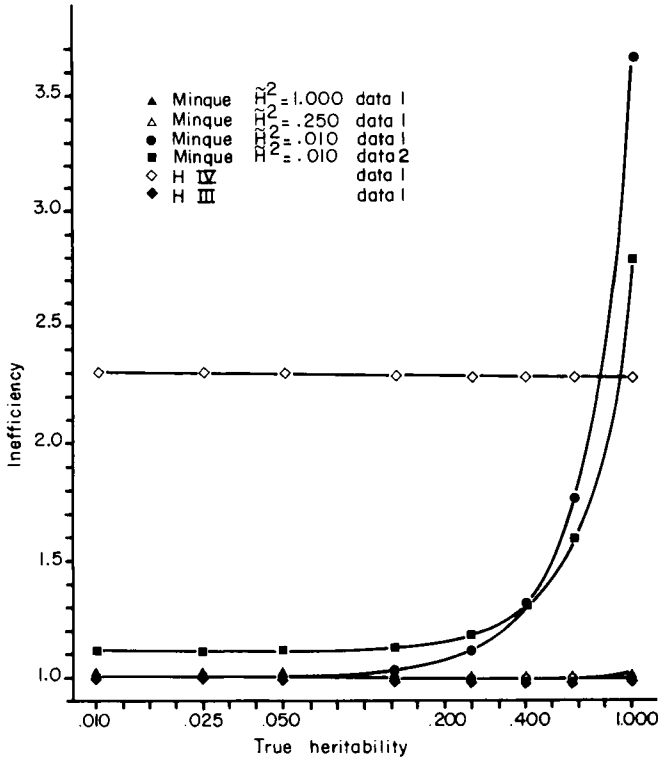


FIG. 1

Inefficiency (increase of variance) of procedures for estimating σ_e^2 compared to MINQUE with $h^2 = \bar{h}^2$. True heritabilities are on a log scale.

Inefficacité (accroissement de variance) des méthodes dans l'estimation de σ_e^2 par rapport au MINQUE au point $h^2 = \bar{h}^2$. Les valeurs vraies de l'héritabilité sont données selon une échelle logarithmique.

B. Efficiency of the methods

The comparison of the efficiency of the estimators is shown in the figures 1-9. There the following attitude is taken: each version of MINQUE or H IV with a given operational value of λ (symbolised as $\hat{\lambda}$) has to be regarded as a procedure in itself, since in practice only one such procedure will be utilized, where of course the true state of nature, that means the true λ , is unknown. Quite often, however, we can put reasonable lower and upper bounds on it. For milk yield e.g., we are rather sure that under our condition the following is true: $0.1 < h^2 < .4$. In addition, with paternal half sibs we have the relation $\lambda = (4 - h^2)/h^2$. Instead of λ and $\hat{\lambda}$, we can therefore use h^2 and \hat{h}^2 , a parameter more familiar to geneticists.

Thus the choice of \hat{h}^2 is often not difficult and the procedure has also to be judged only in this range. All results are given relative to the best possible procedure (in the sense of minimum variance) having the properties of unbiasedness and translation-invariance and utilizing all data. For each true h^2 there exists an optimal procedure, but it is unknown to the user. The minimum variance utilised in the comparison is identical to the large sample variance of REML.

For the comparison shown in the figures the inefficiency is defined as follows:

$$\text{Inefficiency} = \frac{\text{Var}(\hat{\sigma}^2 | \text{for given data set and procedure})}{\text{Var}(\hat{\sigma}^2 | \text{for all data and best procedure})}$$

If the variance of procedure A is x times as large as the variance of the best procedure it can be roughly interpreted as follows: in order to reach the same precision with procedure A as with the best procedure the design (with the given unbalancedness and average daughter number) has to be x times as large. Sometimes, however, the higher precision may not be very crucial e.g. for the estimate of σ_e^2 , since with any procedure (e.g. H IV) we may get a reasonably good estimate.

C. Efficiency for estimating σ_e^2

In the figure 1 the inefficiencies of the procedures with respect to the best procedure are shown. As expected the efficiency of the estimator used for H IV is low since it utilises much fewer df. The H III estimator is only slightly inferior to the best estimator whereas the MINQUE estimator with \hat{h}^2 much smaller than h^2 is very inefficient. There it can even occur ($h^2 = 1$, $\hat{h}^2 = 0.01$), that using the reduced data set the estimate is more precise than using the full data set.

D. Efficiency for estimating σ_u^2

In the figures 2, 3 and 4 the inefficiencies of the procedures for estimating σ_u^2 with respect to the best procedure are shown. The main conclusions from these figures are:

i) By a good choice of \hat{h}^2 a large superiority of the MINQUE estimator over the H III estimator is often achieved.

ii) By using an appropriate value of \hat{h}^2 (such that $|h^2 - \hat{h}^2|$ is small) the H IV estimator is, as expected, inferior to the MINQUE estimator. However, it always retains a high efficiency. This efficiency is highest for very small h^2 , since with respect to the quadratic form for q_1 , H IV and MINQUE converge for $\hat{h}^2 \rightarrow 0$, but they are different for q_0 , where a form is used for H IV which is less efficient. In our data set, the inefficiency of H IV is 1.013 for $h^2 = \hat{h}^2 = .01$ and 1.151 for $h^2 = \hat{h}^2 = 1$.

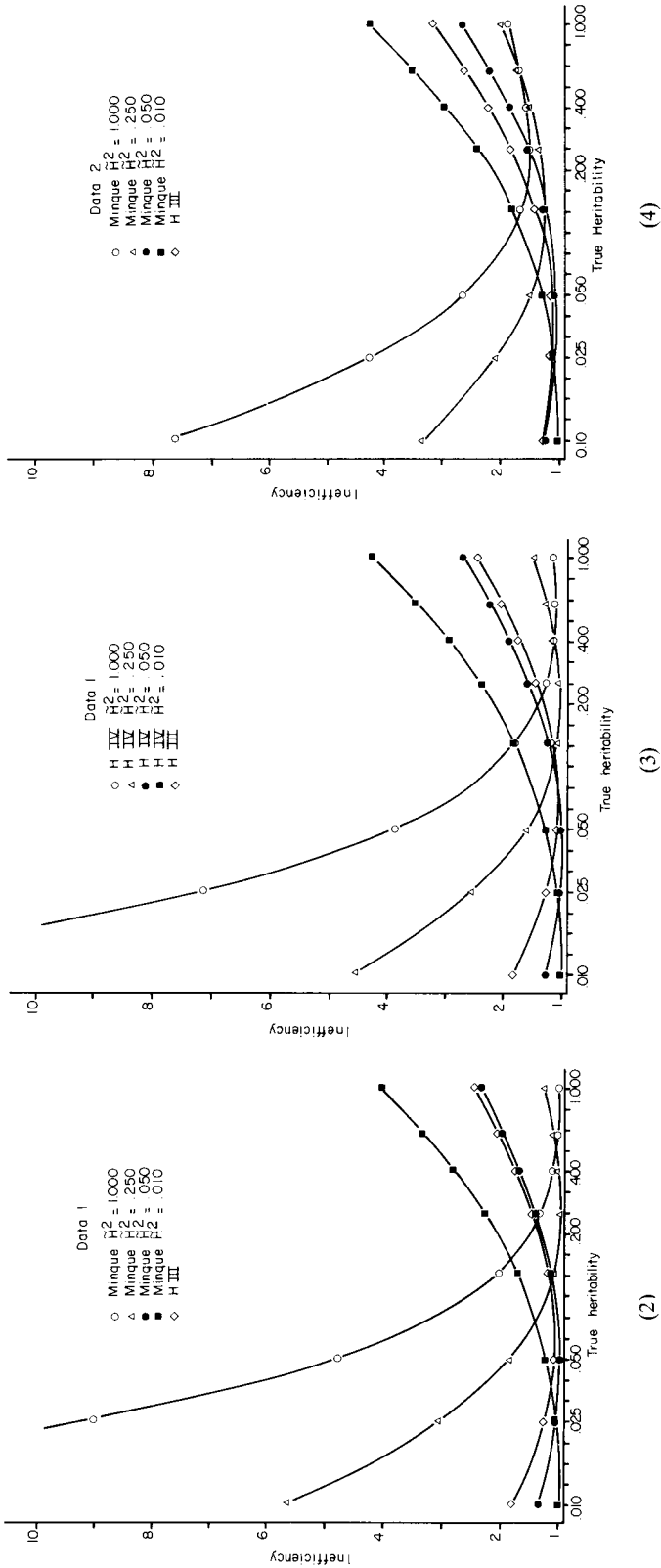


FIG. 2, 3, 4
 Inefficiency (increase of variance) of procedures for estimating σ_u^2 compared to MINQUE
 with $h^2 = h^2$.
 Inefficacité (accroissement de variance) des méthodes dans l'estimation de σ_u^2 par rapport au
 MINQUE au point $h^2 = h^2$.

iii) By using a \hat{h}^2 which is far off the true value of h^2 both MINQUE and H IV are very inefficient. For MINQUE with $\hat{h}^2=0$ (MINQUE (0)) this was also shown by QUASS & BOLGIANO (1979). If h^2 is large but a small value of h^2 is used, H IV decreases somewhat faster in efficiency than MINQUE and if h^2 is small and h^2 large, the efficiency of MINQUE decreases faster. The reason for this behaviour is not obvious to us and it is unclear if this is just peculiar to the present design.

iv) Comparing the figure 2 and the figure 4 for the optimal method, it can be seen that reducing the data set has quite different effects depending on h^2 . If h^2 is very low e.g. $h^2=0.01$, the inefficiency is small (1.05) whereas with $h^2=1.0$ the inefficiency is large (1.88).

v) If a procedure other than the optimal one is used, reducing the data set can improve the estimate. This is true for all three methods considered.

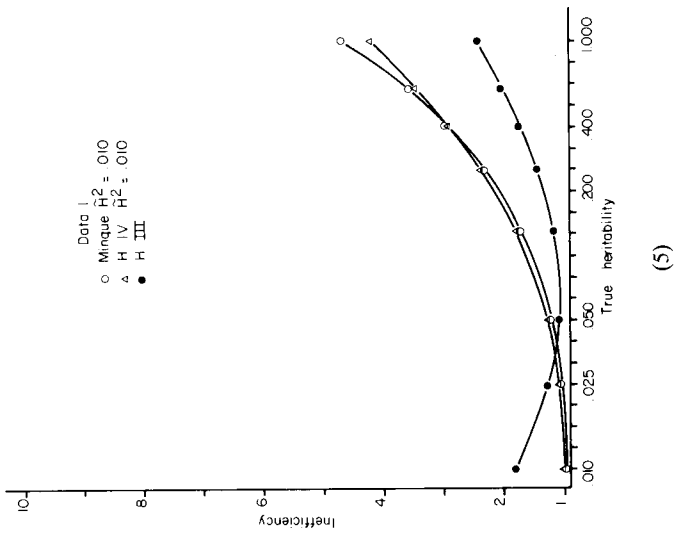
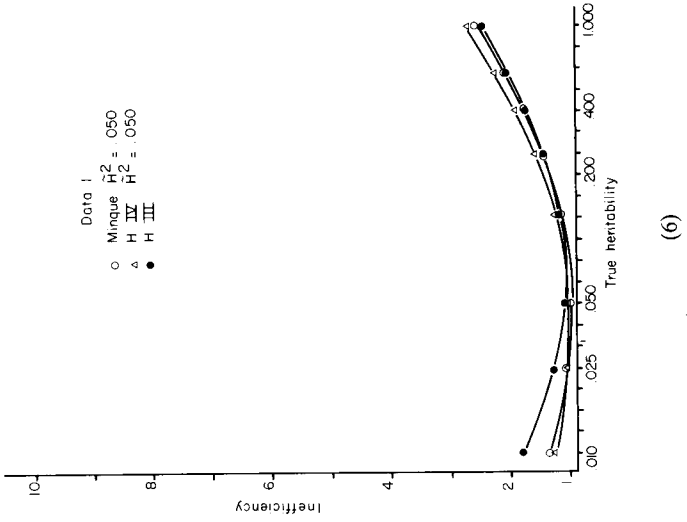
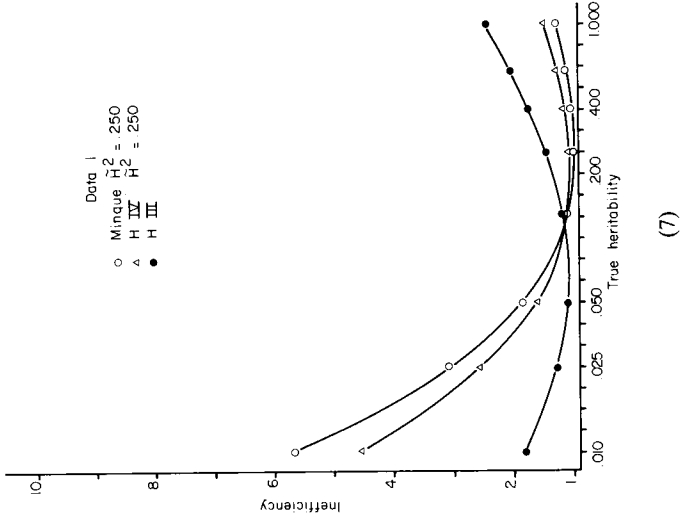
It is at first sight surprising that an estimate can be improved by ignoring data i.e. ignoring information. For the Analysis of Variance method in the one way classification (then identical to H III) this was also pointed out by ROBERTSON (1962) and by SWIGER *et al.* (1964). A look at the formulae in section II.E. explains that paradox. The h^2 are applied to calculate the weights used to combine the means and to combine the squares. If the weights are far off the optimal values then it can easily happen that the estimator combining all squares is less precise than the estimator combining only a subset of the squares. (If we have two estimates of μ , $\hat{\mu}_1$ and $\hat{\mu}_2$ with

$$\text{Var}(\hat{\mu}_1)=1, \quad \text{Var}(\hat{\mu}_2)=5 \quad \text{then} \quad \hat{\mu}_c = \frac{1}{2}(\hat{\mu}_1 + \hat{\mu}_2)$$

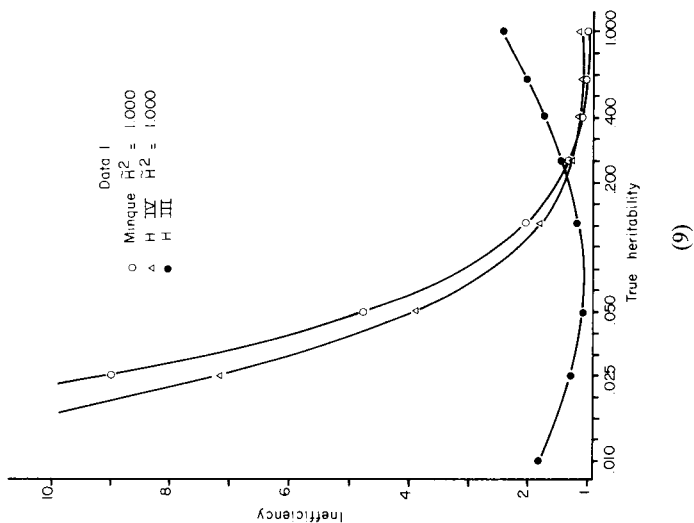
is less precise than $\hat{\mu}_1$). With optimal weights that will never happen. With H III the weights are completely given by the method and they are in no case optimal (except all n_i are equal) but in the present data they are never very extreme. It should be observed, however, that in this data set MINQUE with $\hat{h}^2=.05$ is always better than H III (strictly speaking the superiority was determined for $h^2=.01, .025, .05, .129, .15, .20, .25, .40, .60, 1.0$), and the MINQUE with $\hat{h}^2=.25$ is inferior only with very small h^2 but is considerably better than H III over the remaining range.

A look at the formulae in section II.E. also explains the observation noted under iv). With a low h^2 , bulls having few daughters do not contribute much information. In the optimal method they are weighted not very heavily, whereas with $\lambda=0$ each bull, regardless of daughter number gets equal weight (for q_1). With progeny testing in a random mating population it is always true that $\lambda \geq 3$, ($h^2 \leq 1$) thus for the breeding scheme considered, the weights would differ, but not much for $h^2 \rightarrow 1$. In this case reducing the data set implies ignoring a lot of valuable information.

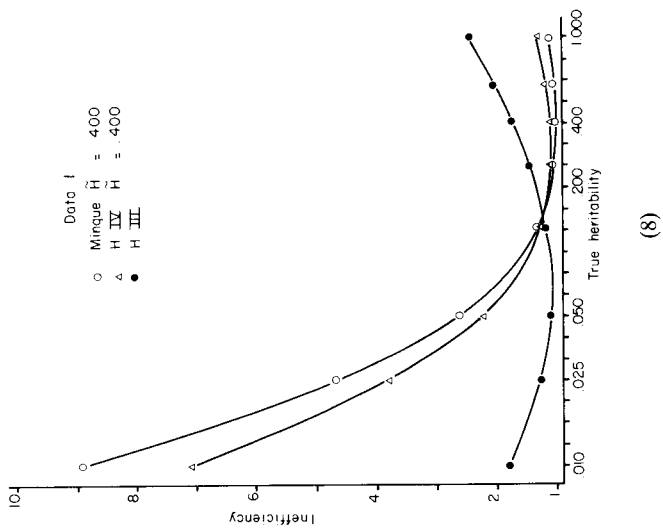
Another observation, which is given in Table 5, indicates that with H IV the smallest variance is not achieved if $h^2 = \hat{h}^2$. For example if $h^2 = .25$ then $\hat{h}^2 = .40$ gives a slightly smaller variance than $\hat{h}^2 = .25$. From our calculations it is not possible to give empirically the best value of \hat{h}^2 for our data set. This observation agrees with one made by HENDERSON (1980).



I.N.R.A. - C.N.R.Z.
 Département de Génétique Animale
 BIBLIOTHÈQUE
 F 78350 JOUY EN JOSAS



(9)



(8)

FIG. 5 à 9

Inefficiency (increase of variance) of procedures for estimating heritabilities compared to MINQUE with $h^2 = h^2$.

Inefficacité (accroissement de variance) des méthodes dans l'estimation de l'héritabilité par rapport au MINQUE au point $h^2 = h^2$.

TABLE 5

*Inefficiency (increase of variance) of H IV for estimating σ_u^2 for some combinations of h^2 and \bar{h}^2 .
Inefficacit  (avec l'accroissement de la variance) de la m thode H IV dans l'estimation de σ_u^2
pour certaines combinaisons de h^2 et \bar{h}^2 .*

h^2	\bar{h}^2		
	.60	.40	.25
.60	1.112	1.160	1.293
.40	1.087	1.090	1.168
.25	1.137	1.074	1.075
.129	1.453	1.245	1.095

E. Efficiency for estimating h^2

In the figures 5 to 9 the variance of \bar{h}^2 is shown. These variances were computed using the usual Taylor Series approximation (KENDALL & STUART 1969, p. 232). The main conclusion from these figures is the relatively high efficiency of H IV compared to MINQUE in spite of the low efficiency of the estimator used for σ_e^2 . In the case investigated this does not have a large effect, since the variance of \bar{h}^2 is dominated by the variance of $\hat{\sigma}_u^2$. For the data set given, the lowest possible s.e. for \bar{h}^2 are 0.006, 0.012, 0.033, 0.045 and 0.077 for $h^2 = 0.01, 0.05, 0.25, 0.40$ and 1.0 respectively.

A further observation can be made by comparing the figure 2 and the figure 6. Though MINQUE with $\bar{h}^2 = 0.05$ was always superior to H III for estimating σ_u^2 , this is not true for estimating h^2 . The reason is found from the figure 1, where it can be seen that for estimating σ_e^2 H III has always a very high efficiency, whereas MINQUE can be quite inefficient for a large value of $|h^2 - \bar{h}^2|$. Since for estimating h^2 both $\hat{\sigma}_e^2$ and $\hat{\sigma}_u^2$ are needed, the lower variance of $\hat{\sigma}_u^2$ from MINQUE is more than compensated for by the larger variance for $\hat{\sigma}_e^2$ in case of $\bar{h}^2 = 0.05$ and $h^2 \rightarrow 1$.

IV. Conclusion

From the results presented and from the more theoretical considerations we conclude that in data sets and models like the ones investigated (which we believe are very common) the judicious use of MINQUE can improve the estimates of genetic parameters quite considerably compared to the H III estimates. The H IV estimators are, as expected, not as good as the MINQUE estimators, but they showed nevertheless a very high efficiency for estimating σ_A^2 and h^2 . One suspected weakness of the H IV estimator against violation of the model assumptions which it inherited from the CC method does not seem to be of great importance according to our limited study. Thus if MINQUE is impossible or very difficult to compute, H IV seems to be a useful alternative.

Received October 29, 1982.

Accepted April 29, 1983.

References

- BULMER M.G., 1971. The effect of selection on genetic variability. *Am., Nat.*, **105**, 201-211.
- DAVID H.A., 1970. *Order statistics*, Wiley, New York.
- DEMPFLE L., 1975. A note on increasing the limit of selection through selection within families. *Genet. Res. Camb.*, **24**, 127-135.
- HAGGER C., SCHNEEBERGER M., DEMPFLE L., 1982. ML, REML, MINQUE and Henderson 3 estimates of variance and covariance components for milk yield, fat and protein content of Braunvieh and Brown Swiss x Braunvieh sires. *Proc., 2nd World Congr., Genet., Appl., Livest., Prod., Madrid 4-8 oct., 1982*.
- HARTLEY H.O., RAO J.N.K., 1967. Maximum likelihood estimation for the mixed analysis of variance model. *Biometrika*, **54**, 93-108.
- HENDERSON C.R., 1953. Estimation of variance and covariance components, *Biometrics*, **9**, 226-252.
- HENDERSON C.R., 1972. Sire evaluation and genetic trends. *Proc., Anim., Breed. Genet., Symp. in honor of Dr J.L. LUSH, ASAS, ADSA, Champaign, Illinois*, 10-41.
- HENDERSON C.R., 1980. A simple method for unbiased estimation of variance components in the mixed model. Mimeo, Cornell University.
- KENDALL M.G., STUART A., 1969. *The Advanced Theory of Statistics*. Vol. 1, Griffin, London.
- PATTERSON H.D., THOMPSON R., 1971. Recovery of interblock information when block sizes are unequal. *Biometrika*, **58**, 545-554.
- POWELL R.L., FREEMAN A.E., 1974. Estimators of sire merit. *J. Dairy Sci.*, **57**, 1228-1233.
- QUAAS R.L., BOLGIANO D.C., 1979. Sampling variances of the MINQUE and Method 3 estimators of the sire component of variance. *Proc., Conf., Var., Comp., Anim., Breed., Cornell Univ., Ithaca, New York, July 16-17, 1979*, 99-106.
- RAO C.R., 1970. Estimation of heteroscedastic variances in linear models. *J. Am., Stat., Assoc.*, **65**, 161-172.
- RAO C.R., 1972. Estimation of variance and covariance components in linear models. *J. Am., Stat., Assoc.*, **67**, 112-115.
- ROBERTSON A., 1962. Weighting in the estimation of variance components in the unbalanced single classification. *Biometrics*, **18**, 413-417.
- SARHAN A.E., GREENBERG B.G., 1962. *Contribution to Order Statistics*. Wiley, New York.
- SEARLE S.R., 1968. Another look at Henderson's methods of estimating variance components. *Biometrics*, **24**, 749-787.
- SWIGER L.A., HARVEY W.R., EVERSON D.O., GREGORY K.E., 1964. The variance of intraclass correlation involving groups with one observation. *Biometrics*, **20**, 818-826.