

# Single and multitrait estimates of breeding values for survival using sire and animal models

T. H. E. Meuwissen<sup>1</sup>, R. F. Veerkamp<sup>1</sup>, B. Engel<sup>1</sup> and S. Brotherstone<sup>2</sup>

<sup>1</sup>*Institute for Animal Science and Health, ID-Lelystad, Box 65, 8200 AB Lelystad, The Netherlands*

<sup>2</sup>*Institute of Cell, Animal and Population Biology, University of Edinburgh, West Mains Road, Edinburgh EH9 3JT, UK*

## Abstract

*Survival data were simulated under the Weibull model in a half-sib family design, and about 50% of the records were censored. The data were analysed using the proportional hazard model (PHM) and, after transformation to survival scores, using a linear and a binary (logit) model (LIN and BIN, respectively), where the survival scores are indicators of survival during time period  $t$  given survival up to period  $t - 1$ . Correlations between estimated and true breeding values of sires (accuracies of selection) were very similar for all three models (differences were smaller than 0.3%). Daughter effects were however less accurately predicted by the LIN model, i.e. taking proper account of the distribution of the survival data yields more accurate predictions of daughter effects. The estimated variance components and regressions of true on estimated breeding values were difficult to compare for the LIN models, because estimated breeding values were expressed as additive effects on survival scores while the simulated true breeding values were expressed on the underlying scale. Also the differences in accuracy of selection between sire and animal model breeding value estimates were small, probably due to the half-sib family design of the data. To estimate breeding values for functional survival, i.e. the component of survival that is genetically independent of production (here milk yield), two methods were compared: (i) breeding values were predicted by a single-trait linear model with a phenotypic regression on milk yield; and (ii) breeding values were predicted by a two-trait linear model for survival and milk yield, and breeding values for survival corrected for milk yield were obtained by a genetic regression on the milk yield breeding value estimates. Both methods yielded very similar accuracies of selection for functional survival, and are expected to be equivalent.*

**Keywords:** *animal models, breeding value, survival.*

## Introduction

Breeding value estimation of longevity suffers from two major difficulties: (i) the records of animals still alive, which include the selection candidates, are censored; (ii) the allocation of cows to contemporary groups is difficult: two cows born in the same herd in 1996 and 1997 which lived for 4 and 3 years, respectively, experienced for most of their lives the same herd management, while a contemporary grouping based on herd-year (-season) of birth would allocate them to different groups. The use of the proportional hazard model (PHM) solves both these difficulties: it takes account of censored records, handles temporary fixed effects, i.e. at different times a cow is allocated to different

contemporary groups, and takes full account of the distribution of longevity data (Ducrocq and Sölkner, 1998a). However, because of computational problems, approximate linear models are often applied to survival data or only a single-trait, sire PHM is used, while (multitrait) animal model estimates of breeding values (EBV) are used for other traits. There are no theoretical limitations to the use of single-trait animal PHM. The multivariate distribution of combined analyses of linear traits and longevity by PHM requires some approximation (Ducrocq, 1999).

Multivariate survival EBV are desirable for several reasons. Culling for production affects survival rates,

such that only genetic improvement of the component of survival that is genetically independent of milk yield is required (assuming a continuous improvement of milk yield). The latter component of survival is often termed functional survival or functional herd life (Ducrocq, 1987) and is obtained in single trait survival analyses by a phenotypic regression of survival or herd life on production, usually expressed as a within-herd deviation. Alternatively, a multitrait analysis of survival and production would yield EBV for functional survival. Whether genetic or phenotypic regression on production is most appropriate is currently debated in the literature (e.g. Visscher *et al.*, 1999). Another benefit of multivariate survival analyses is that accurate survival EBV are only available late in life while selection of young animals is important for genetic gain. We therefore want to include the information from other traits (e.g. functional type traits) to improve the accuracy of survival EBV of young animals. Also, yield is an important cause of culling, which might cause selection bias in the available survival data. A joint (linear and non-linear) evaluation of yield and survival may correct this bias.

Madgwick and Goddard (1989) suggested a multitrait analysis of survival scores, where a survival score in lactation 2 indicates whether the cow survived lactation 2 or not, given that she lived up to lactation 2 (otherwise the survival score of lactation 2 is missing). This approach can also handle some temporary fixed effects, because a cow can be allocated to different groups for each lactation, i.e. each survival score. A drawback of this approach is that survival is modelled by many survival score traits, and the model may be overparameterized. Following the developments of test-day models for daily milk records, Veerkamp *et al.* (2001) suggested the use of random regression models (RRM) for survival scores to avoid this overparameterization. A theoretical derivation showed that RRM for survival scores approximated a PHM analysis when the time-interval for each survival score was small. Here we will use the repeatability model, where a genetic correlation of 1 between survival scores is assumed. The use of a RRM can relax this assumption, and is a potential extension to the work presented here.

The first aim of this simulation study was to compare the accuracy of survival EBV using PHM and using a repeatability model, where the repeatability model either accounted for, or ignored, the binary nature of the survival score data. The latter resulted in a linear model for survival analysis. A second aim was to compare EBVs from sire and animal models using the PHM, binary and linear repeatability models. A

third aim was to compare multitrait prediction of 'genetic' functional survival, i.e. the component of survival that is genetically independent of milk yield, with the commonly used prediction of 'phenotypic' functional survival where a phenotypic regression on milk yield is included in a single trait survival analysis.

## Methods

### *The single-trait data sets: DATAS*

Survival data were simulated using a three-step procedure. In step 1, data were simulated on an underlying linear scale ( $x_{ij}$ ), which affects the hazard of an animal by a factor of  $\exp(x_{ij})$ , where the hazard at time  $t$  denotes the instantaneous culling rate at time  $t$  conditional upon survival up to  $t$ . In step 2, the underlying linear data were transformed to survival times using the Weibull model ( $y_{ij}^*$ ); and in step 3, the survival times were censored in case the animal was still alive at the end of our 40-month trial, which resulted in the observed survival times ( $y_{ij}$ ) and an indicator of censoring ( $w_{ij}$ ). The data on the underlying scale,  $x_{ij}$ , were sampled using the model:

$$x_{ij} = h_{(ij)} + s_i + d_{ij}$$

where  $x_{ij}$  = survival of daughter  $j$  of sire  $i$  on the underlying linear scale ( $x_{ij}$  is called linear because all the included effects affect  $x_{ij}$  linearly);  $h_{(ij)}$  = herd effect of the herd of daughter  $j$  of sire  $i$ , and  $s_i$  = sire effect;  $d_{ij}$  = a within-sire deviation that reflected genetic effects (dam plus Mendelian sampling) and environmental effects. There were 20 herd effects sampled from  $N(0, \sigma_h^2)$ , where  $\sigma_h^2$  is the variance of the herd effects, and each daughter was randomly allocated to one of the herds. In order to simplify the analyses, the fixed herd effect was not time-dependent, although time-dependent effects (such as herd  $\times$  month) could have been simulated and analysed by all models that were considered. There were 100 sire effects,  $s_i$  sampled from  $N(0, \sigma_s^2)$  and each sire had 50 or 100 daughters whose within-sire deviation,  $d_{ij}$ , was sampled from  $N(0, \sigma_w^2)$ , i.e. the data set contained either 5000 or 10000 records, and  $\sigma_s^2$  ( $\sigma_w^2$ ) is the variance of the between- (within-) sire effects, with  $\sigma_s^2 = 0.05$  (as estimated by Ducrocq and Sölkner, 1998) and  $\sigma_w^2 = 0.2$ . The simulated daughter effects are partly genetic (variance is  $3 \times \sigma_s^2 = 0.15$ ) and partly due to the environment of the daughter (variance is  $0.2 - 0.15 = 0.05$ ). Details regarding data simulation are given in Table 1.

The Weibull model was used to simulate uncensored survival times,  $y_{ij}^*$ , from the linear predictor,  $x_{ij}$ , by the following algorithm. The survival times,  $t$ , were assumed to be expressed in months after first

**Table 1** Parameters of the simulated data sets DATAS and DATAM

DATAS and DATAM	
Length of experiment	40 months
No. of sires	100
No. of daughters per sire	50 or 100
No. of herds	20
Variance of herd effects ( $\sigma_h^2$ )	0.25
Weibull parameters	$\lambda = 0.04, \rho = 2$
Variations of survival: sire ( $\sigma_s^2$ )	0.05
daughter ( $\sigma_w^2$ )	0.2
DATAM	
Variations of milk yield: sire ( $\sigma_{Ms}^2$ )	0.1
daughter, ( $\sigma_{Mw}^2$ )	0.9
Correlation between survival and yield ( $r$ )	0 or -0.3†
Fraction culled for yield, month of culling	0.2, 11

† Correlation is the same for sire and daughter effects.

calving. For every month,  $t$ , the probability,  $p_{ij,t}$  was calculated that animal  $ij$  was culled during month  $t$  given that it lived up to month  $t$ , which is:

$$p_{ij,t} = \frac{S(t-1; x_{ij}) - S(t; x_{ij})}{S(t-1; x_{ij})} = 1 - \exp[-\lambda^\rho (t^\rho - (t-1)^\rho) e^{x_{ij}}] \quad (1)$$

where  $S(t, x_{ij})$  = the survivor function of the Weibull regression model, i.e.  $S(t, x_{ij}) = \exp[-(\lambda t)^\rho e^{x_{ij}}]$ , with  $\lambda$  and  $\rho$  being the scale and shape parameters of the Weibull distribution, respectively. The derivation of equation (1) is given in the **Appendix**. Whether the animal actually survived month  $t$  or not was simulated by sampling a uniformly distributed variable  $u_{ij,t}$  for every month  $t$ , and the earliest month where  $u_{ij,t} < p_{ij,t}$  resulted in culling of animal  $ij$  such that its survival time was  $y_{ij}^* = t$ . Note that the above algorithm generates discrete survival times, and would not be efficient when very small units of time are used. The shape parameter,  $\rho$ , was set to 2 (as estimated by Ducrocq and Sölkner, 1998), and the scale parameter,  $\lambda$ , was assumed to be 0.04 (based on calculations aiming for about 50% of the records to be censored).

It was assumed that our experiment lasted only 40 months, after which the records from daughters that were still alive became censored. For every daughter a random month of first calving,  $z_{ij}$ , was sampled between months 1 and 40. The daughters that were still alive at the end of the experiment had censored records with survival times,  $y_{ij} = 40 - z_{ij}$ , while daughters that were culled before month 40 had uncensored records with  $y_{ij} = y_{ij}^*$ . The censoring variable  $w_{ij}$  takes the value 1 for uncensored records and 0 otherwise. For example, if a daughter lived for

30 months after first calving,  $y_{ij}^* = 30$ , and she first calved in month 20,  $z_{ij} = 20$ , her actual survival time would be  $y_{ij} = 40 - 20 = 20$ . This record is censored since the daughter is still alive at month 40, so  $w_{ij} = 0$ . The result was that about 50% of the records were censored.

The survival times,  $y_{ij}$ , and censoring indicators,  $w_{ij}$ , were transformed on a monthly basis to binary survival scores. For instance, if a daughter survived for 3 months but was culled in the next month:  $y_{ij} = 3$ ,  $w_{ij} = 1$ , and her survival scores are  $Y_{ij} = [0 \ 0 \ 0 \ 1]$ . If another daughter had a censored record and the experiment ended in her 6th month when she was still alive,  $y_{ij} = 6$ ,  $w_{ij} = 0$  and  $Y_{ij} = [0 \ 0 \ 0 \ 0 \ 0 \ 0]$ .

#### The multitrait data set: DATAM

Linear survival and milk production were simulated in a multitrait manner. Pairs of records of linear survival and milk yield,

$$\begin{pmatrix} x \\ y_M \end{pmatrix}_{ij}$$

were simulated. First, sire effects

$$\begin{pmatrix} s \\ s_M \end{pmatrix}_i$$

were sampled from

$$MVN \begin{bmatrix} \sigma_s^2 & r\sigma_s\sigma_{Ms} \\ r\sigma_s\sigma_{Ms} & \sigma_{Ms}^2 \end{bmatrix}$$

where  $s_M$  is the sire effect for milk yield,  $\sigma_{Ms}^2$  is the associated sire variance, and  $r$  is the between-sire correlation of the linear survival trait with milk production. Second, residual effects

$$\begin{pmatrix} d \\ d_M \end{pmatrix}_{ij}$$

were sampled from

$$MVN \begin{bmatrix} \sigma_w^2 & r\sigma_w\sigma_{Mw} \\ r\sigma_w\sigma_{Mw} & \sigma_{Mw}^2 \end{bmatrix}$$

where  $d_M$  is the residual effect for yield (contains dam effect and Mendelian sampling effect, and environmental effect of yield),  $\sigma_{Ms}^2$  is its variance, and  $r$  is the residual correlation, i.e. the within- and between-sire correlations were assumed equal. This assumption implies that the genetic and environmental correlations between these traits are also approximately equal. Finally, the herd effect is added to  $x_{ij}$ . Herd effects for milk yield were all set to

zero in the simulation, but were included as fixed effects in the linear mixed model for the analysis of milk yield. In a linear model, estimation of random effects is based on the difference between records and fixed effect estimates, which do not depend on the actual level of the fixed effects, and hence, the true herd levels do not affect the random effect estimates. Putting the above terms together, the linear survival and milk production were simulated by:

$$\begin{pmatrix} x \\ y_M \end{pmatrix}_{ij} = \begin{pmatrix} h_{(ij)} \\ 0 \end{pmatrix} + \begin{pmatrix} s \\ s_M \end{pmatrix}_i + \begin{pmatrix} d \\ d_M \end{pmatrix}_{ij}$$

The transformation from a linear survival trait to actual survival times was the same as in DATAS, except that after 10 months of lactation the milk production records of daughters became available and the daughters with the lowest 20% milk production were culled in month 11 (if they had not been culled already). Parameters for the two data sets are provided in Table 1. The assumption of a zero correlation between milk yield and involuntary culling follows the definition of functional survival, which is the component of survival that is uncorrelated with milk yield. However, true involuntary culling may be expected to show a negative correlation with milk yield because milk yield is negatively correlated with fertility and health traits. To account for this, some data sets were simulated with a correlation of  $r = -0.3$  between  $s_{M_i}$  and  $s_i$ .

True breeding values for functional survival were obtained from:

$$TBV_{FS_i} = a_i - b_g a_{M_i}$$

where  $a_i (= 2s_i)$  and  $a_{M_i} (= 2s_{M_i})$  are the true breeding values of sire  $i$  for survival and milk yield, and  $b_g = r \sigma_s / \sigma_{M_s}$ , i.e.  $b_g = 0$  or  $-0.212$  when  $r = 0$  or  $-0.3$ , respectively (Table 1).

### Analysis of DATAS

DATAS was analysed using PHM, linear and binary models. For the linear and binary models, the survival data were transformed to survival score data, as described above. The binary model (BIN) used the logit-link function,  $\ln(p_{ij,t}/(1 - p_{ij,t}))$ , where  $p_{ij,t}$  is the probability of failure in month  $t$  given survival up to month  $t$ . Hence,  $E(Y_{ij,t}) = p_{ij,t}$ , where the expectation is conditional upon survival up to month  $t$ . If the time classes are sufficiently small, the probability,  $p$ , of failure in any class is small and  $\ln(p/(1 - p)\lambda) \approx \ln(p)$ . The probability of instantaneous failure is  $p_{ij,t} = \lambda \rho(\lambda t)^{\rho-1} \exp(x_{ij}) \delta t$ , where  $t$  is the actual time at the midpoint of time

class  $k$ , and the baseline hazard function is  $\lambda(t) = \lambda \rho(\lambda t)^{\rho-1}$ . The binary model now becomes

$$E(Y_{ij,t}) = p_{ij,t}$$

$$\text{BIN-S: } \text{Logit}(p_{ij,t}) \approx \mu + b \times \ln(t) + h_{(ij)} + s_i$$

$$\text{BIN-AM: } \text{Logit}(p_{ij,t}) \approx \ln(p_{ij,t}) = \mu + b \times \ln(t) + h_{(ij)} + a_i$$

where BIN-S (BIN-AM) denotes the binary sire (animal) model;  $\mu$  = overall mean;  $b = \rho - 1$  = the regression of  $\ln(p)$  on  $\ln(t)$ . The binary model for  $Y_{ij,t}$  therefore contains an overall mean, a regression on log-time, a fixed herd and a random sire effect,  $s_i$  (BIN-S), or a random animal effect  $a_{ij}$  (BIN-AM). In the case of BIN-AM,  $\text{Var}(\mathbf{a}) = \mathbf{A} \sigma_a^2$ , where  $\mathbf{a}$  = vector of  $a_{ij}$  effects,  $\mathbf{A}$  = relationship matrix between the animals, and  $\sigma_a^2$  = additive genetic variance, which is estimated from the data.

The linear models (LIN) are similar to the binary models except that the binary nature of the  $Y_{ij,t}$  data is ignored and effects are introduced on the original scale instead of the log-scale:

$$\text{LIN-S: } E(Y_{ij,t}) = \mu + b \times \ln(t) + h_{(ij)} + s_i$$

$$\text{LIN-AM: } E(Y_{ij,t}) = \mu + b \times \ln(t) + h_{(ij)} + a_i$$

where the regression on  $\ln(t)$  reflects the 'mean survival score curve',  $h_{(ij)}$  = fixed herd effect; and  $s_i$  = random sire effect.

The proportional hazard models were analysed using the SURVIVAL\_KIT package (Ducrocq and Sölkner, 1998b), choosing the Weibull model with fixed herd and random normally distributed sire (PHM-S) or animal (PHM-AM) effects.

### Analysis of DATAM with selection on milk yield

For the prediction of functional survival EBV by phenotypic regression, we performed the following linear single-trait animal model analysis with regression on milk production:

$$\text{LIN-AM1: } E(Y_{ij,t}) = \mu + b \times \ln(t) + b_2 \times y_{M_{ij}} + h_{(ij)} + a_{ij}$$

where  $b_2$  = phenotypic regression coefficient on milk production. The EBV of functional survival are obtained as the estimates of  $a_{ij}$ .

Next we predicted the genetic component of survival that is genetically independent of milk yield,  $a - b_g a_{M_i}$ , with  $a$  ( $a_{M_i}$ ) = additive genetic value of survival (milk yield) and  $b_g$  = genetic regression coefficient. The following bi-variate animal model analysis was performed to predict  $a$  and  $a_{M_i}$ :

$$\text{LIN-AM2: } \begin{pmatrix} Y_{ij,t} \\ y_{Mij} \end{pmatrix} = \begin{pmatrix} \mu \\ \mu_M \end{pmatrix} + b_1 \begin{pmatrix} \ln(t) \\ 0 \end{pmatrix} + \begin{pmatrix} h \\ h_M \end{pmatrix}_{(ij)} + \begin{pmatrix} a \\ a_M \end{pmatrix}_{ij} + \begin{pmatrix} d \\ d_M \end{pmatrix}_{ij} + \begin{pmatrix} \varepsilon_{ij,t} \\ 0 \end{pmatrix}$$

with

$$\text{Var} \begin{pmatrix} \mathbf{a} \\ \mathbf{a}_M \end{pmatrix} = \mathbf{G} \otimes \mathbf{A}$$

and  $\mathbf{G} = 2 \times 2$  genetic correlation matrix of the 2 traits; similarly

$$\text{Var} \begin{pmatrix} \mathbf{d} \\ \mathbf{d}_M \end{pmatrix} = \mathbf{R} \otimes \mathbf{I}$$

and  $\mathbf{R} = 2 \times 2$  environmental correlation matrix of the daughter effects;  $\varepsilon_{ij,t}$  = sampling deviation of individual survival scores  $Y_{ij,t}$  from the average survival of animal  $ij$  ( $\text{Var}(\varepsilon) = \mathbf{I}\sigma_c^2$ ). In this model of analysis, the estimation of the variance components in  $\mathbf{R}$  is problematic because of the structure of the survival score data. In particular element  $R_{11}$ , i.e.  $\text{Var}(\mathbf{d})$ , cannot be separated from  $\text{Var}(\varepsilon)$ . For the same reason,  $\text{Cov}(\mathbf{d}, \mathbf{d}_M)$ , i.e. element  $R_{12}$ , cannot be estimated either. This problem is solved by the following reparameterization. Consider the regression of the permanent environmental component of survival,  $\mathbf{d}$ , on the permanent environmental effect of milk yield,  $\mathbf{d}_M$ :

$$\mathbf{d} = b_r \mathbf{d}_M + \mathbf{d}_e = \mathbf{d}^* + \mathbf{d}_e,$$

where  $b_r$  = regression coefficient of  $\mathbf{d}$  on  $\mathbf{d}_M$ ;  $\mathbf{d}_e$  = residual of this regression;  $\mathbf{d}^* = b_r \mathbf{d}_M$ . It follows that  $\mathbf{d}^*$  is orthogonal to  $\mathbf{d}_e$ , and the correlation between  $\mathbf{d}^*$  and  $\mathbf{d}_M$  is 1. Because of the structure of the survival score data, we cannot distinguish the permanent environmental effect of survival that is independent of milk production,  $\mathbf{d}_e$ , from the sampling deviation,  $\varepsilon$ . In our model of analysis, the combined effect of  $\mathbf{d}_e$  and  $\varepsilon$  will be estimated by  $\varepsilon$ , and, consequently,  $\mathbf{d}^*$  will be estimated by  $\mathbf{d}$ . This reparameterization results in a correlation of 1 between  $\mathbf{d}$  and  $\mathbf{d}_M$ , and setting the correlation between  $\mathbf{d}$  and  $\mathbf{d}_M$  to 1 makes  $R_{11}$  and  $R_{12}$  estimable and results in  $\text{Var}(\mathbf{d})$  estimating the residual variance of survival that can be explained by milk yield. The correlation between  $\mathbf{d}$  and  $\mathbf{d}_M$  was set to 0.99 rather than 1.0, because otherwise  $\mathbf{R}^{-1}$  could not be computed. There are more elegant ways to parameterize LIN-AM2, but in this way it could be programmed in the ASREML package (Gilmour *et al.*, 2000).

The EBV for functional survival of the animals are obtained from:

$$\text{EBV}_{\text{FSi}} = \hat{a}_i - \hat{b}_g \hat{a}_{Mi},$$

where  $\hat{a}_i$  ( $\hat{a}_{Mi}$ ) = estimated breeding value of sire  $i$  for survival (milk production); and  $\hat{b}_g = \text{Cov}_g(a; a_M)_{\text{Est}} / \text{V}_g(a_M)_{\text{Est}}$ , and  $\text{Cov}_g(a; a_M)_{\text{Est}} =$  within replicate estimate of genetic covariance; and  $\text{V}_g(a_M)_{\text{Est}} =$  within replicate estimate of genetic variance of milk production.

In each of the 50 replicated DATAS and 10 replicated DATAM data sets, variance components were estimated using ASREML (Gilmour *et al.*, 2000) for the LIN and BIN models and the SURVIVAL\_KIT (Ducrocq and Sölkner, 1998b) for the PHM models. The latter also explicitly estimated the Weibull shape parameter  $\rho$ . Within each replicate, accuracies of selection were obtained by calculating the correlation between true (simulated) and estimated sire effects (Yadzi *et al.*, 2000).

## Results

Table 2 shows the accuracies of selection of the 100 sires when breeding values were predicted by LIN, BIN or PHM. When used as sire or animal models, LIN, BIN and PHM gave very similar accuracies of selection. The animal model EBV seemed only marginally more accurate than the sire model EBV. Also, the linear model EBV were slightly less accurate than binary- and PHM-EBV, which were equally accurate. These differences were however very small.

Table 3 gives estimates of the sire and animal variances from the sire and animal models, respectively. The sire effects of LIN are expressed on the scale of the survival scores instead of on the

**Table 2** Correlation between true and estimated breeding values of the sires for functional survival in the single-trait data set DATAS†

Analysis‡	Sire model	Animal model
50 daughters per sire		
LIN	0.687	0.687
BIN	0.689	0.690
PHM	0.689	0.690
100 daughters per sire		
LIN	0.804	0.805
BIN	0.805	0.806
PHM	0.805	0.807

† Results averaged over 50 replicated DATAS data sets; standard errors varied between 0.005 and 0.007.

‡ LIN = linear model; BIN = binary model; PHM = proportional hazard model.

**Table 3** Estimates of sire and animal variances obtained from analysis of the single-trait data set DATAS with 100 daughters per sire (the simulated sire variance was 0.05, which results in an animal variance of 0.2†)

Analysis‡	Sire model	Animal model
LIN	$2.9 \times 10^{-5}$	$3.6 \times 10^{-5}$
BIN	0.048	0.274
PHM	0.045	0.221

† Results averaged over 50 replicated DATAS data sets; standard errors varied between  $1.2 \times 10^{-6}$  and  $1.8 \times 10^{-6}$  for LIN, and between 0.0014 and 0.0073 for BIN and PHM.

‡ LIN = linear model; BIN = binary model; PHM = proportional hazards model.

underlying linear scale, which results in much smaller sire variances. The **Appendix** shows that the average culling probability of the daughters of a sire (conditional on survival up to the previous month) may be approximated by:

$$p_i \approx m_i (1 + s_i),$$

with  $m_i$  = the average 'baseline' culling probability conditional on survival up to the previous month, where averaging is over the months that the daughters of sire  $i$  actually lived, and the term 'baseline' denotes that these conditional culling probabilities are due to the baseline hazard, i.e. they are calculated under the assumption that  $x_{ij} = 0$ . Note that the values of  $m_i$  differ between sires, because the months that a daughter actually lived (and thus was a candidate for culling) depends on the longevity of the daughters of the sire, and thus on  $s_i$ . The sire variances of the LIN model equal the variances of the conditional culling probabilities of the sires,  $p_i$ :

$$\text{Var}(p_i) \approx [E(m_i)]^2 \sigma_s^2 + 2\text{Cov}(m_i; s_i)E(m_i) + \text{Var}(m_i) \quad (2)$$

which reduces to  $m^2 \sigma_s^2$  when all survival scores,  $m_i$ , are equal to  $m$ , i.e. when the baseline hazard is a straight line ( $\rho = 1$ ); and the approximation is from a first order Taylor series expansion of  $m_i$  and  $(1 + s_i)$  around their means. The sire variances of the LIN model were indeed accurately predicted by  $m^2 \sigma_s^2$  in simulated data sets with  $\rho = 1$  (result not shown). Note that the term  $\text{Cov}(m_i; s_i)$  is usually negative: sires with low  $s_i$  (longer living daughters) obtain a higher  $m_i$ , because their daughters reach ages where the baseline hazard (the average culling) is large and *vice versa*. When  $m_i$  was calculated by averaging the conditional culling probabilities,  $p_{ij,t}$ , of equation (1), with  $x_{ij}$  set to 0, we obtained  $E(m_i) = 0.0312$ ;  $\text{Var}(m_i) = 7.39 \times 10^{-6}$ ; and  $\text{Cov}(m_i; s_i) = -0.000302$ . Substitution of these values in equation (2) results in

$\text{Var}(p_i) = 3.71 \times 10^{-5}$ , which is somewhat larger than the LIN sire variance of  $2.9 \times 10^{-5}$  in Table 3. A possible explanation for this difference is that the average of within-sire means yields an overestimate of the sire variance.

For normally distributed traits, we would expect the variance of animal effects to be four times larger than the variance of sire effects. There are two reasons why, for these analyses, this is not the case: (i) the sire model treats all survival scores of sire  $i$  as independent records, while the animal model accounts for the fact that survival scores of daughter  $j$  of sire  $i$  are correlated, and that the better survivors among the daughters of sire  $i$  have more survival scores; (ii) the animal model averages the culling probabilities first within daughters, which gives equal weight to the average culling probabilities of short and long lifetimes (in the sire model, longer lifetimes result in more survival scores and thus more weight). When the baseline hazards were first averaged per daughter before being averaged per sire, the terms of equation (2) were  $E(m_i) = 0.023$ ,  $\text{Var}(m_i) = 5.04 \times 10^{-6}$ , and  $\text{Cov}(m_i; s_i) = -0.000239$ , which gives a sire variance of  $2.1 \times 10^{-5}$ . The animal variance is  $4 \times 2.1 \times 10^{-5} = 8.4 \times 10^{-5}$ , which should be compared with the estimate of  $3.6 \times 10^{-5}$  in Table 3. Apparently the estimates of the daughter effects have also affected the estimate of the animal variance in the LIN animal model analysis, and biased this estimate downwards.

Estimated sire variances from the BIN and PHM sire models were close to the simulated sire variance of 0.05 (Table 3), because here both the simulated and modelled sire effects were expressed on the same scale, i.e. the underlying scale. The animal model estimates were also reasonably close to their expectation of 0.2 ( $= 4 \times 0.05$ ), although the BIN-AM estimate, in particular, was somewhat too large, possibly due to difficulties with variance component estimation of binary animal models (Engel and Buist, 1998).

Table 4 shows the regression of true on estimated sire effects, i.e.  $\text{Cov}(s_i; \hat{s}_i) / \text{Var}(\hat{s}_i)$ , where  $\hat{s}_i$  = estimate of the sire effect. In normal linear model analyses this regression equals 1. In the LIN model, the estimates of the sire effects,  $s_i$ , are however expressed on the scale of the culling probabilities while the true sire effects,  $s_i$ , are expressed on the underlying scale, such that their regression deviates from 1. For the BIN and PHM models, the regressions of true on estimated sire effects are close to their expectation of 1, because here both the true and estimated sire effects are expressed on the same (underlying) scale.



**Table 4** Regression of true on estimated sire effects in the single-trait data set DATAS with 100 daughters per sire†

Analysis‡	Sire model	Animal model§
LIN	44.9	92.8
BIN	1.02	0.90
PHM	1.06	0.97

† Results averaged over 50 replicated DATAS data sets; standard errors varied between 1.0 and 4.2 for LIN, and between 0.012 and 0.016 for BIN and PHM.

‡ LIN = linear model; BIN = binary model; PHM = proportional hazard model.

§ Animal model breeding value estimates were divided by 2 in order to obtain estimates of sire effects.

Table 5 shows the accuracies of selection for functional survival of the 100 sires when breeding values were predicted by a single-trait model using phenotypic regression (LIN-AM1) or a bi-variate model using genetic regression (LIN-AM2). Note that in the data set with zero genetic and environmental correlations between milk production,  $y_{Mij}$ , and the linear survival trait,  $x_{ij}$ , a relationship between actual survival,  $y_{ij}$ , and milk production is introduced by the culling for production. The accuracies of selection are very similar for phenotypic and genetic regression.

## Discussion

Single and multitrait, sire and animal model breeding value estimation of survival were compared. Only marginal differences in accuracies were observed between the linear and non-linear animal and sire models, which is expected because the LIN model is an approximation of BIN, which is an approximation of PHM (Veerkamp *et al.*, 2001). The high accuracies of LIN models are encouraging, because the LIN models are much more easily extended to large scale multitrait animal model evaluations, which are needed to predict total profit EBVs of animals, than are non-linear models. Meuwissen *et al.* (2000) used unrealistically high sire and daughter variances, which exaggerated the non-linearities of the model, and found that the linear model was about 4% less accurate than PHM and BIN. The latter two models were equally accurate. However, Table 2 shows that these differences are much smaller when more realistic (smaller) variances are used. Note, though, that the simulated data sets were perfectly designed (equal numbers of progeny per sire, no dam relationships, no overlapping generations, limited environmental effects), which favours accurate predictions for all models and limits any differences in accuracy between models. Further research is needed to compare the models under

**Table 5** The correlation between true and estimated breeding values of the sires for functional survival in the two-trait data set DATAM with 100 daughters per sire†

Analysis‡	$r = 0§$	$r = -0.3$
LIN-AM1 using phenotypic regression	0.720	0.708
LIN-AM2 using genetic regression	0.720	0.707

† Results averaged over 10 replicated DATAM data sets.

‡ LIN-AM1 = single-trait linear animal model with regression on milk production; LIN-AM2 = bi-variate linear animal model for milk production and survival and genetic regression is used to predict functional survival.

§  $r$  = the correlation between survival and milk yield for the between- and within-sire effects.

more challenging circumstances, and also to compare the models for possible estimation biases, and their ability to estimate genetic trends (e.g. due to selection for milk) and to deal with various environmental factors.

Because the sum of genetic and permanent environmental effects of animals,  $s_i + d_{ij}$ , were simulated, the accuracy of daughter EBVs could not be obtained by calculating the correlation between daughter genetic effects and their EBVs. However, the correlation between  $s_i + d_{ij}$  and  $EBV_{ij}$  was 0.113 and 0.159 for the linear and binary model, respectively. Note that this correlation would be  $0.8 = (4 \times 0.05) / (0.05 + 0.2)$ , when the EBVs had an accuracy of 1. The figures of 0.113 and 0.159 suggest rather low accuracies of selection for females, as expected by Ducrocq (1999), although they were also low because  $\pm 50\%$  of the daughters had only a censored record. In contrast to the accuracies of selection of sires, the accuracies of selection of females are substantially increased when taking proper account of the distribution of the survival data, which may be partly due to LIN-EBVs and  $s_i + d_{ij}$  being expressed on different scales, i.e. the non-linear scale transformation reduces the correlation. Although the accuracies of survival EBVs of females are rather low, their prediction is still needed in order that a total merit index for each cow can be calculated.

The animal models were as accurate as the sire models in predicting sire EBVs (Table 2). This result is expected when traits are normally distributed and bulls have unrelated mates as assumed here. However, in the case of survival data the sire model is inconsistent in assuming a normally distributed sire effect but a Weibull distributed residual, while part of this residual will be due to genetic effects of dams and Mendelian sampling that are also

normally distributed (given normally distributed sire effects) (Korsgaard *et al.*, 2000). Such extra variation, e.g. due to dam or Mendelian sampling effects, is referred to as 'overdispersion' (Louis, 1991; Tempelman and Gianola, 1996). The latter model was also used for the simulation of the data, except that here a normally distributed environmental effect was also assumed. Meuwissen *et al.* (2000) did find improved accuracies of about 3% when using animal models compared with sire models, but they used unrealistically high variance components in their simulations. In the case of more realistic parameter estimates, these additional normally distributed within family effects are small relative to the error due to Weibull sampling (their variances are 0.2 and  $\pi^2/6 = 1.645$ , respectively). Treating the combination of these residual effects as one residual due to Weibull sampling, as in the sire model, has apparently only a marginal effect on the accuracy of predicting EBV (Table 2). Although the animal model analyses more closely resembled the simulated data, i.e. they did contain normally distributed within family effects, the accuracies of EBV were very similar to those of the sire models.

Although the accuracy of LIN-EBVs was competitive with those of BIN and PHM models, they are expressed on different scales, i.e. on the scale of culling probabilities *versus* the underlying scale. The transformation between the scales is complicated by the fact that the additive effect of a sire on culling probabilities depends on the baseline hazard, which depends on the part of hazard function that the daughters experience on average, which in turn depends on the longevity of the daughters and thus on the sire effect. In the case of the LIN animal model, the estimate of the animal variance was smaller than expected. Further research is needed into why this happens, and whether the estimate of the animal variance increases when maternal relationships between the daughters are included in the analysis.

The environmental effect of daughters on survival scores was not fitted in the single trait models, because it cannot be separated from the error due to the Weibull sampling. In the analysis of survival scores, the variance due to a permanent environmental effect on these scores involves estimation of the covariance between successive scores. However, this is not possible because the second score is only known when the first score is 0 (the animal survived), i.e. in pairs of first and second score observations, the first score is always 0 such that this covariance cannot be estimated (Visscher *et al.*, 1999). Following Visscher *et al.* we therefore assumed that the covariance between successive

survival scores was zero, i.e. a permanent environmental effect was not fitted in the analyses. If this permanent environmental effect was fitted, the resulting accuracies of selection were very variable (result not shown).

In the multitrait analysis LIN-AM2, a permanent environmental effect of survival scores was fitted, because it could be estimated through the environmental correlation with milk production by setting this correlation to 1.0. Thus, the residual of the survival scores is split into two effects: (1) the permanent environmental effect, which is predicted from the environmental effect of milk yield; and (2) an error which is orthogonal to this permanent environmental effect. This is similar to writing  $B = b \times A + e$ , where A and B are two traits,  $b$  is a regression coefficient, and  $(b \times A)$  has a correlation of 1 with A. Note that splitting B into a term with a correlation of 1 with A  $(b \times A)$  and an orthogonal residual is possible, irrespective of the correlation between A and B, and enables the common environmental variance to be estimated.

Phenotypic and genetic regression on milk production gave almost identical selection accuracies (Table 5). This is in agreement with selection index results for two linear traits (Kennedy *et al.*, 1993). The optimal linear selection index weights for improvement of genetic functional survival, i.e.  $a_f - b_g a_{Mi}$ , are 1 and  $-b_p$  for survival ( $y_i$ ) and milk production ( $y_{Mi}$ ), respectively, where  $b_p$  is the phenotypic regression coefficient of survival on milk. Note that the breeding goal is orthogonal to milk yield, i.e.  $\text{Cov}(y_{Mi}, a_f - b_g a_{Mi}) = \text{Cov}(a_{Mi}, a_f - b_g a_{Mi}) = \text{Cov}(a_{Mi}, a_i) - b_g \text{Var}(a_{Mi}) = 0$ , i.e. the breeding goal is predicted by the component of survival that is orthogonal to  $y_{Mi}$ , i.e.  $y_f - b_p y_{Mi}$ . This explains the result that genetic functional survival is equally accurately predicted from a single trait analysis of functional herd-life with phenotypic regression on production as from a bi-variate prediction of survival and production, followed by a genetic regression of the bi-variate EBV. Although, this result holds always for linear traits as shown by Kennedy *et al.* (1993), a sensitivity analysis is needed to determine whether the results of Table 5 also hold with different parameters, e.g. different correlations between the genetic and between the environmental effects.

We attempted multitrait analyses of production and survival, where the model for survival contained a phenotypic regression on production, but this gave inconsistent results. The inconsistencies occurred because one of the  $y$ -variables, i.e.  $y_{Mi}$  is simultaneously used as an explanatory variable, i.e. as a  $x$ -variate. These inconsistent results will



disappear if selection and regression was for another trait, e.g. predicted profitability of the cow. A phenotypic regression on within-herd deviations of milk yield will not remove these problems because the model for survival also contains a herd effect, which would fit any across-herd differences in  $b_p Y_{Mf}$ . Since herd effects for milk yield were zero in the simulations, the actual milk yields also represented within-herd deviations.

Despite the equivalence of the bi-variate model for survival and production, and the single-trait survival model with a phenotypic regression on survival, multitrait models still have an advantage in that they can use the information from early predictors of survival, such as functional type traits. Hopefully the latter will substantially improve the accuracy of the survival breeding values of young animals that are eligible for selection.

## References

**Ducrocq, V.** 1987. An analysis of length of productive life in dairy cattle. *Ph.D. thesis, Cornell University.*

**Ducrocq, V.** 1999. Topics that deserve further attention in survival analysis applied to dairy cattle breeding — some suggestions. *INTERBULL Bulletin 21*: 181-189.

**Ducrocq, V. and Sölkner, J.** 1998a. Implementation of a routine breeding value evaluation of dairy cows using survival analysis techniques. *Proceedings of the sixth world congress on genetics applied to livestock production, Armidale, vol. 23* pp. 359-362.

**Ducrocq, V. and Sölkner, J.** 1998b. The Survival Kit — a Fortran package for the analysis of survival data. *Proceedings of the sixth world congress on genetics applied to livestock production, Armidale, vol. 27*, pp. 447-448.

**Engel, B. and Buist, W.** 1998. Bias reduction of approximate maximum likelihood estimates for heritability in threshold models. *Biometrics 54*: 1155-1164.

**Gilmour, A. R., Cullis, B. R., Welham, S. J. and Thompson, R.** 2000. *ASREML reference manual version 2000*. New South Wales Agriculture, Orange, Australia.

**Kennedy, B. W., Werf, J. H. J. van der and Meuwissen, T. H. E.** 1993. Genetic and statistical properties of residual feed intake. *Journal of Animal Science 71*: 3239-3250.

**Korsgaard, I. R., Andersen, A. H. and Jensen, J.** 2000. On different models, on heritability, reliability and related quantities of survival analysis. *Proceedings of the European Association for Animal Production, The Hague, vol. 51*, p. 80.

**Louis, T.** 1991. Assessing, accommodating and interpreting the influences of heterogeneity. *Environment and Health Perspectives 80*: 215-222.

**Madgwick, P. A. and Goddard, M. E.** 1989. Genetic and phenotypic parameters of longevity in Australian dairy cattle. *Journal of Dairy Science 72*: 2624-2632.

**Meuwissen, T. H. E., Engel, B., Veerkamp, R. F. and Brotherstone, S.** 2000. A linear approximation to proportional hazard models for the analysis of survival

data. *Proceedings of the European Association for Animal Production, The Hague, vol. 51*, p. 81.

**Tempelman, R. J. and Gianola, D.** 1996. A mixed model for overdispersed count data in animal breeding. *Biometrics 52*: 265-279.

**Veerkamp, R. F., Brotherstone, S., Engel, B. and Meuwissen, T. H. E.** 2001. Analysis of censored survival data using random regression models. *Animal Science 72*: 1-10.

**Visscher, P. M., Thompson, R., Yadzi, H., Hill, W. G. and Brotherstone, S.** 1999. Genetic analysis of longevity in the UK: present practice and considerations for the future. *INTERBULL Bulletin 21*: 16-22.

**Yadzi, M. H., Thompson, R., Ducrocq, V. and Visscher, P. M.** 2000. Genetic parameters and response to selection in proportional hazard models. *Proceedings of the European Association for Animal Production, The Hague, vol. 51*, p. 81.

(Received 10 September 2001—Accepted 21 March 2002)

## Appendix

*The culling probability in time period  $t-1$  to  $t$  conditional on survival up to  $t-1$*

The culling probability during a finite time period  $t-1$  to  $t$  given that the animal survived up to time  $t-1$ ,  $p_{ij,t}$  equals the fraction of the animals that are culled in this time period divided by the fraction of the animals that live up to this time period. The latter is given by the survivor function:

$$S(t-1, x_{ij}) = \exp[-(\lambda(t-1))^\rho e^{x_{ij}}],$$

with  $\lambda$  and  $\sigma$  being the scale and shape parameters of the Weibull model, respectively, and  $x_{ij}$  = the underlying linear effect of the  $j$ th survival record of the  $i$ th sire. Because the fraction of the animals that is culled in the period  $t-1$  to  $t$  is  $S(t-1, x_{ij}) - S(t, x_{ij})$ , it follows that:

$$\begin{aligned} p_{ij,t} &= \frac{S(t-1; x_{ij}) - S(t; x_{ij})}{S(t-1; x_{ij})} \\ &= 1 - \frac{S(t; x_{ij})}{S(t-1; x_{ij})} \\ &= 1 - \frac{\exp[-(\lambda t)^\rho e^{x_{ij}}]}{\exp[-(\lambda(t-1))^\rho e^{x_{ij}}]} \\ &= 1 - \exp[-\lambda^\rho (t^\rho - (t-1)^\rho) e^{x_{ij}}], \end{aligned}$$

which equals equation (1) in the main text.

If the time period  $t-1$  to  $t$  is sufficiently small, the conditional culling probability simplifies to:

$$\begin{aligned} p_{ij,t} &\approx 1 - [1 - \lambda^\rho \{t^\rho - (t-1)^\rho\} e^{x_{ij}}] \\ &\approx \lambda^\rho \{t^\rho - (t-1)^\rho\} e^{x_{ij}}. \end{aligned}$$

This equation may be used to approximate the average additive effect of a sire. Let us assume that the number of daughters is sufficiently large such that the average effect of daughters and herds is approximately zero, i.e. the effect on the linear scale simplifies from  $x_{ij}$  to  $s_i$  when averaging over all daughters of sire  $i$ . Hence, the conditional culling probability of sire  $i$  is:

$$\begin{aligned} p_{i,t} &\approx \lambda^\rho [t^\rho - (t-1)^\rho] e^{s_i} \\ &\approx m(t) e^{s_i} \\ &\approx m(t) (1 + s_i) \end{aligned}$$

where  $m(t) = \lambda^\rho [t^\rho - (t-1)^\rho]$ . Next we average over the time periods that the daughters of sire  $i$  experience these culling probabilities, i.e. the average conditional culling probability of the daughters of sire  $i$  is:

$$p_i = E_{t,i}(p_{i,t}) = m_i (1 + s_i),$$

where  $m_i = E_{t,i}(m(t))$  with  $E_{t,i}()$  denoting the expectation over the time periods  $t$  that the daughters of sire  $i$  actually lived. Note that  $m_i$  differs between sires, because for a sire with longer living daughters the averaging  $E_{t,i}(m(t))$  is over a different set of time periods compared with when the daughters have shorter lifetimes. If  $\text{Var}(m_i)$  is sufficiently small, the additive effect of sire  $i$  on the conditional culling probabilities may be approximated by  $E(m_i)s_i$ .