1 **The analysis of disease biomarker data using mixed hidden Markov**

2 **model.**

3

4 Johann C. Detilleux

5

6 Quantitative Genetics Group, Department of Animal Production, Faculty of

7 Veterinary Medicine, University of Liège, Liège, Belgium

8

9 Corresponding author:

10 e-mail: jdetilleux@ulg.ac.be

11 Tel: 32 4 366 4215

12 Fax: 32 4 366 4122

13

14 Short title : Mixed hidden Markov model

15

1

1    **Abstract -**    A mixed hidden Markov model (HMM) is developed for predicting

2    breeding values of a biomarker (here, somatic cell score) and the individual

3    probabilities of health and disease (here, mastitis) based upon measurements of the

4    biomarker.  At a first level, the unobserved disease process (Markov model) is

5    introduced and, at a second level, the measurement process is modelled, making

6    the link between the unobserved disease states and the observed biomarker values.

7    This hierarchical formulation allows joint estimation of the parameters of both

8    processes.  The flexibility of this approach is illustrated on simulated data.  Firstly,

9    lactation curves for the biomarker are generated based upon published parameters

10   (mean, variance and probabilities of infection) for cows with known clinical

11   condition (health or mastitis due to *E. coli* or *S. aureus*).  Next, estimation of the

12   parameters is performed via Gibbs sampling, assuming the health status is

13   unknown.  Results from the simulations and the mathematics showed that the

14   mixed HMM is appropriate to estimate the quantities of interest although the

15   accuracy of the estimates is moderate when the prevalence of the disease is low.

16   The paper ends with some indications for further developments of the

17   methodology.

18

19   **hidden Markov model/ mixed model/ mastitis/ somatic cell score**

20

## 1. INTRODUCTION

Studies have shown variability among cows for natural resistance to intra-mammary infection (IMI). Selection is therefore possible but direct measures of IMI are not readily available. Usually, information on IMI is based upon biomarkers such as somatic cell scores (SCS), electrical conductivity, immunoglobulin or acute phase proteins (reviewed in Detilleux, accepted). One important difficulty in using these biomarkers to find the most resistant animals is that factors known to influence their expression may be different in healthy (IMI-) from infected (IMI+) cows. As these are usually unidentified, breeding values tend to be biased. To reduce this bias and to infer more precisely cows' individual probabilities to be IMI- or IMI+, several authors have used the mixture model methodology on SCS [2], [8], [12], [17]. A generalization of the mixture model is the hidden Markov model (HMM) that presents the advantages to not only estimate individual probabilities of being infected but also to predict individual probabilities of new infection and of recovery, both useful to compute epidemiological measures of IMI spread within a population and to assist mastitis control programs.

The objective of this study is to present the mathematical formalism behind the HMM methodology as it may apply to the analysis of infectious disease biomarkers assumed to be dependent upon the genetic make-up of

the cows. The fit of the HMM will be assessed on data simulated based on parameters obtained in a survey of clinical mastitis cases. Bayesian estimates of the parameters will be obtained using the Gibbs sampler. Finally, limitations and possible extensions of the current approach will be discussed.

## 2. MATERIALS AND METHODS

Throughout, k indexes the individual cow, t (t = 1 to T) is the follow-up time point during the lactation (e.g., month in milk), $y_k^t$ is the value of the biomarker observed at t on animal k, and $z_k^t$ is the corresponding unknown health status (IMI- or IMI+). Let $z_k^t = 0$ if $y_k^t$ is from an unknown IMI- sample and $z_k^t = 1$ if $y_k^t$ is from an unknown IMI+ sample. For simplicity, T is assumed constant for all cows. We use the notation Ødegård *et al.* [17] utilized for their finite mixture model, with slight modifications.

### 2. 1. General formulation of the model

Conditionally on the unknown vector **z**, it was assumed that the vector of observations **y** could be described by the linear model:

$$\mathbf{y} = \mathbf{M_0}\ \mu_0 + \mathbf{M_1}\ \mu_1 + \mathbf{Z}\ \mathbf{a} + \mathbf{e}$$

1    where $\mathbf{y}$ is the (NT X 1) data vector of $y_k^t$, $\boldsymbol{\mu_0}$ and $\boldsymbol{\mu_1}$ are (T X 1) vectors of

2    fixed effects for data on a IMI- or IMI+ cow, respectively, $\mathbf{a}$ is the (Na X 1)

3    vector of random additive genetic effects; $\mathbf{M_0}$ is the (NT X T) matrix with

4    elements = 1 if $z_k^t = 0$ and = 0 otherwise; $\mathbf{M_1}$ is the (NT X T) matrix with

5    elements = 1 if $z_k^t = 1$ and = 0 otherwise; $\mathbf{e}$ is the (NT X 1) vector of

6    residuals; $\mathbf{Z}$ is the (NT X Na) incidence matrix relating $\mathbf{a}$ to $\mathbf{y}$, N is the

7    number of animals with data and Na is the number of animals with pedigree

8    records.

9       The conditional distribution of $\mathbf{y}$, given the vector $\mathbf{z}$, the location and scale

10   parameters, was assumed to be:

11               $(\mathbf{y}|\,\boldsymbol{\mu_0},\,\boldsymbol{\mu_1},\,\sigma_0^2,\sigma_1^2,\,\mathbf{a},\,\mathbf{z}) \sim N\,[(\mathbf{M_0}\,\boldsymbol{\mu_0} + \mathbf{M_1}\,\boldsymbol{\mu_1} + \mathbf{Z}\,\mathbf{a}),\,\mathbf{R}]$

12   with $\mathbf{R} = \mathbf{F_0}\,\sigma_0^2 + \mathbf{F_1}\,\sigma_1^2$, where $\mathbf{F_i}$ is the (NT X NT) diagonal matrix with

13   elements = 1 if $z_k^t = i$ and = 0 otherwise. The parameters $\sigma_0^2$ and $\sigma_1^2$ are the

14   residual variances associated to a record on an IMI- and IMI+ cow,

15   respectively. For the additive effects, it was assumed that $(\mathbf{a}\,|\sigma_a^2) \sim N\,[0,$

16   $\mathbf{A}\,\sigma_a^2\,]$ where $\sigma_a^2$ is the additive genetic variance and $\mathbf{A}$ is the matrix of

17   additive genetic relationship between animals.

18

19   **2. 2. Sampling distribution of the observations given group status**

20

1    The density of the vector $\mathbf{y}$ for the subset of the $N_i$ observations with $z_k{}^t =$

2    i, i.e. $\{\mathbf{z} = i\}$, given the location parameters and the residual variances can be

3    written as:

4    $\mathrm{pr}\,(\mathbf{y}|\,\boldsymbol{\mu_i},\,\sigma_i^2\,,\,\{\mathbf{z}=i\})$

5    $$\propto (\sigma_i^2)^{N_i/2}\,\exp\{(\frac{-1}{2\sigma_i^2})\,(\mathbf{y}-\mathbf{M_i}\,\boldsymbol{\mu_i}-\mathbf{Z\,a})'\,\mathbf{F_i}\,(\mathbf{y}-\mathbf{M_i}\,\boldsymbol{\mu_i}-\mathbf{Z\,a})\}.$$

6

7    **2.3. Prior distributions of parameters and of the unknown status vector**

8

9    For $i = 0$ or $1$, normal prior densities were assumed for the location

10    parameters:

11    $$\mathrm{pr}(\boldsymbol{\mu_i}) \propto (s_i^2)^{-T/2}\,\exp\,\{(-\frac{1}{2s_i^2})\,(\boldsymbol{\mu_i}-\mathbf{1}\,m_i)'\,(\boldsymbol{\mu_i}-\mathbf{1}\,m_i)\},$$

12    where $\mathbf{1}$ is the (TX1) vector of 1.  The prior density for the additive effects,

13    conditionally on the additive variance, was:

14    $$\mathrm{pr}(\mathbf{a}\big|\sigma_a^2) \propto (\sigma_a^2)^{-N/2}\,\exp\{(-\frac{1}{2\sigma_a^2})\,\mathbf{a}'\,\mathbf{A^{-1}\,a}\,\}.$$

15    Under simple mixture models, the individual elements of the classification

16    vector $\mathbf{z}$ are assumed to be independent a priori and to follow the same

17    Bernoulli distribution with the mixing proportion as parameter.  Here, under

18    a equally simple mixed HMM, the variables $z_k{}^t$ do not follow the same

19    distribution.  The first element of the series $(z_k{}^1)$ follows a Bernoulli

1     distribution with $\lambda_k$ as parameter while the other elements follow Bernoulli

2     distributions with state transition probabilities from $z_k^{t-1}$ to $z_k^t$ as parameters.

3     Formally, the unknown state at time t may be decomposed in:

4 $$pr(z_k^t = i) = p(z_k^t = i | z_k^{t-1} = 0)\, p(z_k^{t-1} = 0) + p(z_k^t = i | z_k^{t-1} = 1)\, p(z_k^{t-1} = 1)$$

5     where $p(z_k^t = i | z_k^{t-1} = j)$ are the state transition probabilities with i, j = 0 or 1.

6     The state transition probabilities are assumed to possess the first-order

7     Markov property namely that, given the present state, the future and past

8     states are independent or that the current value ($z_k^t$) depends solely on the

9     most recent past value ($z_k^{t-1}$). Transition probabilities are also independent

10     of the actual time at which the transition takes place (stationarity

11     assumption).     Then, we have $pr(z_k^t = i | z_k^{t-1} = j) = \pi_k^{ij}$ for all t and

12     $(z_k^t = i | z_k^{t-1} = 0) \sim Ber(\pi_k^{00})$ and $(z_k^t = i | z_k^{t-1} = 1) \sim Ber(\pi_k^{01})$.

13

14     **2. 4. Priors for variance components and probabilities**

15

16     Scale-inverse chi-square distributions with $\nu$ degrees of freedom and scale

17     parameters ($s^2_a$, $s^2_0$, and $s^2_1$) were used for the variance components:

18 $$pr(\sigma_a^2) \propto (\sigma_a^2)^{-(\nu+2)/2} \exp(-\nu s_a^2/2\,\sigma_a^2),$$

19 $$pr(\sigma_0^2) \propto (\sigma_0^2)^{-(\nu+2)/2} \exp(-\nu s_0^2/2\sigma_0^2), \text{ and}$$

1  $$\mathrm{pr}(\sigma_1^2) \propto (\sigma_1^2)^{-(v+2)/2} \exp(-v\,s_1^2/2\,\sigma_1^2).$$

2  Finally, $\lambda_k, \pi_k^{00}$ and $\pi_k^{01}$ were assigned uniform (i.e. Beta(1,1)) prior

3  distributions.

4

5  **2. 5. Joint posterior distributions**

6

7  For all cows, the joint posterior density of all unknown parameters is given

8  by:

9  $$\mathrm{pr}(\mu_0, \mu_1, \sigma_a^2, \sigma_0^2, \sigma_1^2, \mathbf{z}, \mathbf{a}, \boldsymbol{\pi^{00}}, \boldsymbol{\pi^{01}}, \boldsymbol{\lambda}|\mathbf{y})$$

10  $$\propto \quad \mathrm{pr}(\mathbf{y}|\;\mu_0, \mu_1, \sigma_a^2, \sigma_0^2, \sigma_1^2, \mathbf{z}, \mathbf{a}, \boldsymbol{\pi^{00}}, \boldsymbol{\pi^{01}}, \boldsymbol{\lambda})$$

11  $$\mathrm{pr}\,(\mathbf{z}|\,\mu_0, \mu_1, \sigma_a^2, \sigma_0^2, \sigma_1^2, \mathbf{a}, \boldsymbol{\pi^{00}}, \boldsymbol{\pi^{01}}, \boldsymbol{\lambda})$$

12  $$\mathrm{pr}(\mathbf{a}|\,\mu_0, \mu_1, \sigma_a^2, \sigma_0^2, \sigma_1^2, \boldsymbol{\pi^{00}}, \boldsymbol{\pi^{01}}, \boldsymbol{\lambda})$$

13  $$\mathrm{pr}(\mu_0)\,\mathrm{pr}(\mu_1)\,\mathrm{pr}(\sigma_0^2)\,\mathrm{pr}(\sigma_1^2)\,\mathrm{pr}(\sigma_a^2)\;\mathrm{pr}(\boldsymbol{\pi^{00}})\,\mathrm{pr}(\boldsymbol{\pi^{01}})\,\mathrm{pr}(\boldsymbol{\lambda})$$

14  where $\boldsymbol{\lambda} = [\lambda_1, .., \lambda_N]'$, $\boldsymbol{\pi^{00}} = [\pi_1^{00},..., \pi_N^{00}]$ and $\boldsymbol{\pi^{01}} = [\pi_1^{01},..., \pi_N^{01}]'$.

15

16  Explicitly, the joint posterior is:

17  $$(\sigma_0^2)^{-(N_0+v+2)/2} \exp-\frac{1}{2\sigma_0^2}\left\{v\,s_0^2 + (\mathbf{y} - \mathbf{M_0}\,\boldsymbol{\mu_0} - \mathbf{Z}\,\mathbf{a})'\,\mathbf{F_0}\,(\mathbf{y} - \mathbf{M_0}\,\boldsymbol{\mu_0} - \mathbf{Z}\,\mathbf{a})\right\}$$

18  $$(\sigma_1^2)^{-(N_1+v+2)/2} \exp-\frac{1}{2\sigma_1^2}\left\{v\,s_1^2 + (\mathbf{y} - \mathbf{M_1}\,\boldsymbol{\mu_1} - \mathbf{Z}\,\mathbf{a})'\,\mathbf{F_1}\,(\mathbf{y} - \mathbf{M_1}\,\boldsymbol{\mu_1} - \mathbf{Z}\,\mathbf{a})\right\}$$

1      $(s_0^2)^{-T/2} \exp \{(-\dfrac{1}{2s_0^2})(\boldsymbol{\mu_0} - \mathbf{1} \, m_0)' \, (\boldsymbol{\mu_0} - \mathbf{1} \, m_0)\}$

2      $(s_1^2)^{-T/2} \exp \{(-\dfrac{1}{2s_1^2})(\boldsymbol{\mu_1} - \mathbf{1} \, m_1)' \, (\boldsymbol{\mu_1} - \mathbf{1} \, m_1)\}$

3      $(\sigma_a^2)^{-(N+\nu+2)/2} \;\; \exp - \dfrac{1}{2\sigma_a^2} \{\nu \, s_a^2 + \mathbf{a}' \, \mathbf{A^{-1}} \, \mathbf{a}\}$

4      $\displaystyle\prod_{k=1}^{N} (\lambda_k)^{K_k^{0,1}+1} \, (1-\lambda_k)^{K_k^{1,1}+1}$

$\displaystyle\prod_{k=1}^{N} (\pi_k^{00})^{n_k^{00}+1} \, (1 - \pi_k^{00})^{n_k^{10}+1} \;\; \mathbf{x} \;\; \prod_{k=1}^{N} (\pi_k^{01})^{n_k^{01}+1} (1 - \pi_k^{01})^{n_k^{11}+1}$

5      where $K_k^{i,1}$ is an indicator function which takes the value 1 if $z_k^{\,1} = i$ and

6      0 otherwise and $n_k^{ij}$ = number of transitions from $z_k^{\,t} = j$ to $z_k^{\,t+1} = i$ .

7

8      **2. 6. Fully conditional posterior distributions.**

9

10      The conditional posterior distributions of each parameter (or block of

11      parameters) are required for implementing a Gibbs sampler. Conditional on

12      **y** and **z**, these conditional posterior densities are analytical because they only

13      involve one of the possible realizations in the space of all possible

14      sequences of **z**. For the location parameters, we have:

1 $$(\mu_i^t \mid \Theta, \mathbf{y}, \mathbf{z}) \sim N\left( \frac{s_i^2 \sum\limits_{k}^{N}(y_k^t - a_k) K_k^{i,t} + m_i \, \sigma_i^2}{(s_i^2 \sum\limits_{k}^{N} \eta_k^{i,t}) + \sigma_i^2}, \frac{s_i^2 \, \sigma_i^2}{(s_i^2 \sum\limits_{k}^{N} \eta_k^{i,t}) + \sigma_i^2} \right),$$

2 where $\Theta$ refers to values of all parameters that the conditional distributions

3 depend upon (i.e., all parameters except the one under consideration), $\eta_k^{i,t}$ is

4 the number of cows with IMI- $(i = 0)$ or IMI+ $(i = 1)$ unknown state at the $t^{th}$

5 time.

6     Let $\mathbf{W} = [\mathbf{Z}\ \mathbf{M_0}\ \mathbf{M_1}]$ and the vector of parameters $\boldsymbol{\theta} = [\mathbf{a}\ \mu_0\ \mu_1]$'. Hence,

7 one can write the model as: $\mathbf{y} = \mathbf{Z}\,\mathbf{a} + \mathbf{M_0}\,\mu_0 + \mathbf{M_1}\,\mu_1 + \mathbf{e} = \mathbf{W}\,\boldsymbol{\theta} + \mathbf{e}$. By

8 partitioning the parameter vector $\boldsymbol{\theta}$ as $\boldsymbol{\theta_1} = \mathbf{a}$ and $\boldsymbol{\theta_2} = [\mu_0\ \mu_1]$', we can

9 compute the conditional posterior distribution of the vector of additive

10 genetic values as $(\mathbf{a} \mid \Theta, \mathbf{y}, \mathbf{z}) \sim N(\hat{\mathbf{a}}, \mathbf{C}_{11}^{-1})$ with $\hat{\mathbf{a}} = \mathbf{C}_{11}^{-1}\left[\mathbf{r_1} - \mathbf{C}_{12}\boldsymbol{\theta}_2\right]$ and $\mathbf{r_1}$,

11 $\mathbf{C_{11}}, \mathbf{C_{12}}$ = the corresponding partition of $\mathbf{C} = [\mathbf{W'}\,\mathbf{R^{-1}}\,\mathbf{W} + \mathbf{A^{-1}}/\sigma_a^2]$ and $\mathbf{r} =$

12 $\mathbf{W'}\,\mathbf{R^{-1}}\,\mathbf{y}$.

13     The fully conditional posterior density of the genetic variance is

14 $$\mathrm{pr}\,(\sigma_a^2 \mid \Theta, \mathbf{y}, \mathbf{z}) \propto (\sigma_a^2)^{-(N+\nu+2)/2} \ \exp -\frac{1}{2\sigma_a^2}\{\nu\, s_a^2 + \mathbf{a'}\,\mathbf{A^{-1}}\,\mathbf{a}\}$$

15 which is in the form of a scale-inverse chi-square density, with $[N + \nu]$

16 degrees of freedom and scale parameter $[\mathbf{a'}\ \mathbf{A^{-1}}\ \mathbf{a} + \nu\ s^2_a]$. Likewise, the

1   fully conditional densities of the residual variances for IMI- and IMI+

2   observations are:

3   $pr\,(\sigma_i^2 \mid \Theta,\,\mathbf{y},\,\mathbf{z})$

4   $\propto (\sigma_i^2)^{-(N_i+\nu+2)/2}\,\exp-\dfrac{1}{2\sigma_i^2}\{\nu\,s_i^2 + (\mathbf{y}-\mathbf{M_i}\,\boldsymbol{\mu_i}-\mathbf{Z}\,\mathbf{a})'\,\mathbf{F_i}\,(\mathbf{y}-\mathbf{M_i}\,\boldsymbol{\mu_i}-\mathbf{Z}\,\mathbf{a})\}$

5   which are in the form of scale-inverse chi-square densities, with $[N_i + \nu]$

6   degrees of freedom, and with scale parameter $= \{\nu\,s_i^2 + (\mathbf{y}-\mathbf{M_i}\,\boldsymbol{\mu_i}-\mathbf{Z}\,\mathbf{a})'\,\mathbf{F_i}$

7   $(\mathbf{y}-\mathbf{M_i}\,\boldsymbol{\mu_i}-\mathbf{Z}\,\mathbf{a})\}$ for i = 0 and 1.

8   For the $k^{th}$ cow, the fully conditional posterior densities of the parameters

9   $\lambda_k,\ \pi_k^{00}$ and $\pi_k^{01}$ are:

10                   $pr\ (\lambda_k \mid \Theta,\,\mathbf{y},\,\mathbf{z})$

11       $\propto \lambda^{K_k^{0,1}+1}(1-\lambda)^{K_k^{1,1}+1},\,pr\,(\pi_k^{00} \mid \Theta) \propto (\pi_k^{00})^{n_k^{00}+1}(1-\pi_k^{00})^{n_k^{10}+1},$

12   and $pr\,(\pi_k^{01} \mid \Theta,\,\mathbf{y},\,\mathbf{z}) \propto (\pi_k^{01})^{n_k^{01}+1}(1-\pi_k^{01})^{n_k^{11}+1}$ which are in the form of beta

13   distributions.

14     Finally, one must compute the fully conditional distribution for each

15   individual $z_k^t$. These may be obtained either from the $pr(\mathbf{z}\mid\Theta,\,\mathbf{y})$ or by

16   considering $pr(z_k^t \mid \mathbf{z}(-z_k^t),\,\Theta,\,\mathbf{y})$ where $\mathbf{z}(-z_k^t)$ represent the hidden vector $\mathbf{z}$

1      without $z_k^t$, as suggested by one referee. Under the first alternative, $pr(\mathbf{z}|\Theta)$

2      can be decomposed as :

3
$$pr(\mathbf{z}|\Theta,\mathbf{y}) = pr(z_k^1|\Theta,\mathbf{y}) \prod_{t=2}^{T} pr(z_k^t|z_k^{t-1},\Theta,\mathbf{y}),$$

4      which leads to a stochastic version of the forward-backward algorithm in

5      which $z_k^1$ is sampled from a Bernoulli distribution with parameter

6      $pr(z_k^1 = 0 \cap \mathbf{y})$ and each $z_k^t$ is sampled successively (for t = 2 to T) from

7      Bernoulli distributions with parameter $\xi_k^{ij,t} = pr(z_k^t = i \,|\, z_k^{t-1} = j, \mathbf{y})$. The

8      computations are reduced as components of $\xi_k^{ij,t} = \dfrac{\alpha_k^{j,t-1}\, \pi_k^{ij}\, b_k^{i,t}\, \beta_k^{i,t}}{\alpha_k^{j,t-1}\, \beta_k^{j,t-1}}$ may be

9      stored gradually as t increases from 1 to T:

10
$$\alpha_k^{j,t} = pr([y_k^1, y_k^2, \ldots y_k^t] \cap z_k^t = j),$$
$$\beta_k^{i,t} = pr([y_k^{t+1}, \ldots, y_k^T] \,|\, z_k^t = i),$$
$$\pi_k^{ij} = pr(z_k^t = i \,|\, z_k^{t-1} = j) \text{ and } b_k^{i,t} = pr(y_k^t \,|\, z_k^t = i).$$

11      The forward and backward probabilities can be efficiently calculated by the

12      following recursion formulae [10]:

13
$$\alpha_k^{j,t} = [\alpha_k^{0,t-1}\, \pi_k^{j0} + \alpha_k^{1,t-1}\, \pi_k^{j1}]\; b_k^{j,t}$$
$$\beta_k^{i,t} = \left[\beta_k^{0,t+1}\, \pi_k^{0i}\, b_k^{0,t+1}\right] + \left[\beta_k^{1,t+1}\, \pi_k^{1i}\, b_k^{1,t+1}\right]$$

1  with initial conditions given by: $\alpha_k^{0,1} = \lambda_k \ b_k^{0,1}$, $\alpha_k^{1,1} = (1 - \lambda_k) \ b_k^{1,1}$ and

2  $\beta_k^{i,T} = 1$ for i = 0 and 1.

3  In the second alternative, $\text{pr}(z_k^t \mid \mathbf{z}(-z_k^t), \Theta, \mathbf{y})$ is reduced to $\text{pr}(z_k^t \mid z_k^{t-1}, z_k^{t+1},$

4  $\Theta, \mathbf{y})$ because of the first-order Markov property on $\mathbf{z}$. Then, $\text{pr}(z_k^t =$

5  $i \mid z_k^{t-1} = j, z_k^{t+1} = r, \Theta, \mathbf{y}) \propto \text{pr}(y_k^1 \mid z_k^1 = i) \, \text{pr}(z_k^1 = i)$ if t = 1. It is proportional

6  to $\text{pr}(z_k^t = i \mid z_k^{t-1} = j) \, \text{pr}(y_k^t \mid z_k^t = i, \Theta) \, \text{pr}(z_k^{t+1} = r \mid z_k^t = i)$ for t = 2 to T-1 and to

7  $\text{pr}(y_k^T \mid z_k^T = i) \, \text{pr}(z_k^T = i \mid z_k^{T-1} = j)$ if t = T. Note this alternative uses T

8  different components while the first alternative generates a realization of $\mathbf{z}$

9  directly from its conditional $p(\mathbf{z} \mid \mathbf{y}, \Theta)$. It presents also a more complicated

10  correlation structure (since each $z_k^t$ depends on both $z_k^{t-1}$ and $z_k^{t+1}$) than the

11  first alternative which may lead to a slower mixing chain.

12

13  **2. 7. Implementation of a Gibbs Sampler**

14

15  The following steps describe how a Gibbs sampling can be implemented

16  for our model, using the stochastic version of the forward-backward

17  algorithm to sample $\mathbf{z}$:

18  1.  Set initial values for parameters as needed.

2.  Select the block ($\mathbf{\theta}_1$) of the vector $\mathbf{\theta}$, compute $\widetilde{\theta}_1 = C_{11}^{-1}\left[r_1 - C_{12}\theta_2\right]$ and replace $\mathbf{a}$ with $[\widetilde{\theta}_1 + C_{11}^{-0.5}\,\text{rannor}(0)]$ where rannor(0) is a random draw from a standard normal distribution.

3.  Replace $\mu_i$ (i = 0 and 1) with

$$
\left[\frac{s_i^2 \sum\limits_{k}^{N}(y_k^t - a_k)\,K_k^{1,t} + m_i\,\sigma_i^2}{(s_i^2 \sum\limits_{k}^{N} n_{i,k}) + \sigma_i^2}\right] + \left[\left(\frac{s_i^2\,\sigma_i^2}{(s_i^2 \sum\limits_{k}^{N} n_{i,k}) + \sigma_i^2}\right)^{0.5} \text{rannor}(0)\right].
$$

4.  Replace $\sigma_a^2$ with $(\mathbf{a'}\,\mathbf{A^{-1}}\,\mathbf{a} + v\,s_a^2)/\chi_{N+v}^2$, where $\chi_{N+v}^2$ is a random draw from a central chi-square distribution with $[v + N]$ degrees of freedom.

5.  Replace $\sigma_i^2$ with $\{v\,s_i^2 + (\mathbf{y} - \mathbf{M_i}\,\mathbf{\mu_i} - \mathbf{Z}\,\mathbf{a})'\,\mathbf{F_i}\,(\mathbf{y} - \mathbf{M_i}\,\mathbf{\mu_i} - \mathbf{Za})\}\,/\chi_{N_i+v}^2$ for i = 0 or 1, where $\chi_{N_i+v}^2$ is a random draw from a central $\chi$-square distribution with $[N_i+v]$ degrees of freedom.

6.  Compute $\zeta_k^{0,1} = \alpha_k^{0,1}\,\beta_k^{0,1} = \text{pr}\,(z_k^1 = 0 \cap \mathbf{y})$ and sample $z_k^1$ from Ber($\zeta_k^{0,1}$).

7.  Compute and store $\zeta_k^{0j,t}$ for t = 2, ... , T and j = 0 or 1.  Then, sample $z_k^t$ from Ber($\zeta_k^{0j,t}$) if $z_k^{t-1} = j$ for t = 2, .., T.

8. Sample $\lambda_k$ and $\pi_k^{ij}$, from their corresponding beta distributions with parameters $K_k^{i,1} + 1$ and $n_k^{ij} + 1$, for $i,j = 0$ and 1, respectively.

9. Repeat (2 through 8) $\rho$ times for burn-in as needed. Then, sample all parameters $\delta$ times. The total number of cycles is $\rho + \delta$.

In this study, values for the hyperparameters are: $s^2_0 = 0.5$, $s^2_1 = 1$, $m_0 =$ overall average computed from the data, $m_1 = m_0 + 3$, $\nu = 2$, $s^2_a = h^2 s^2_p$ ($s^2_p$ = variance computed from the data) and $h^2 = 0.1$.

**2. 8. Simulations**

The model was evaluated using simulated values for the biomarker (here, SCS) with genetic effects considered as having the same distributions for cows with IMI+ and IMI- samples. Each simulation was replicated 10 times. Simulated rather than real data were used because a negative diagnosis, even based on absence of bacteria in cell culture, is not a guarantee of health and the reverse has also been observed [22].

**2. 8. 1. Simulated data**

Results from the field study of de Haas *et al.* [6], [7] on pathogen-specific SCC curves among multiparous cows were used to simulate the means of monthly samples from IMI- and IMI+ cows. In the Figure 3b of de Haas's paper [6], it is shown that, in cows clinically infected with *E. coli*, SCC increase rapidly after infection occurring around the 2$^{nd}$ month-in-milk, peak at 2,000 cells per μL above pre-infection values and return to pre-infection levels one month later. On the other hand, presence of a long increased SCC, without recovery within 4 consecutive months, was common in lactations with clinical *Staph. aureus* mastitis. In the cows without clinical mastitis, SCC followed an approximate inverse lactation curve. The SCC values were log$_2$-transformed in SCS and used to simulate the SCS means, as explained below. In the simulations, it was also considered that cows may be classified as high and moderate responders on the basis of the extent of their immune response to a particular infection [14]. Therefore, SCS were considered at higher values and of longer duration in high than moderate responders (Figure 1).

In the simulations, 3 discrete generations were considered with 400 cows per generation. No selection was applied, sires were selected from 30 different bulls, each cow was replaced by a daughter, and mating was at random. Breeding values for base animals were sampled from a normal distribution with null mean and additive variance of 0.15 or 0.25. The values for the additive variance was found in the literature [4]. Breeding

1     values for non-base animals were sampled from a normal distribution with

2     the mid-parent value as mean and variance = 0.15/2 or 0.25/2.  Inbreeding

3     was ignored.

4     Somatic cell scores under healthy ($\textbf{SCS}_0$) and infected ($\textbf{SCS}_1$) states were

5     simulated as follows:

6     $$\textbf{SCS}_0 = \textbf{M}_0\ \boldsymbol{\mu}_0 +\ \textbf{a} + \textbf{e}_0 \ \text{ and } \textbf{SCS}_1 = \textbf{M}_1\ \boldsymbol{\mu}_1 +\ \textbf{a} + \textbf{e}_1 \ ,$$

7     where $\boldsymbol{\mu}_0$ and $\boldsymbol{\mu}_1$ are the (T X 1) vector means of both distributions, $\textbf{a}$ is  the

8     (N X 1) vector of breeding values (computed as above), and $\textbf{M}_0$ and $\textbf{M}_1$ are

9     the incidence matrices relating $\boldsymbol{\mu}_0$ and $\boldsymbol{\mu}_1$ to $\textbf{SCS}_0$ and $\textbf{SCS}_1$, respectively.

10     The number of observations per cow was set at T = 10 or 20 observations.

11     The vectors $\textbf{e}_0$ and $\textbf{e}_1$ were sampled from 2 normal distributions with null

12     means and residual variances set at 1.0 and 1.4.  The values for the residual

13     variances were found in the literature [13].  Each element of $\boldsymbol{\mu}_0$ and $\boldsymbol{\mu}_1$ was

14     taken from the curves observed in cows without and with mastitis, and for

15     high and low responders (Figure 1).  The cows were assigned to a group

16     (IMI+ or IMI-) at random using appropriate membership probabilities:  The

17     proportion of cows with at least one IMI+ sample was set at $P_{cow}$ = 20% and

18     50% and, among IMI+ cows, the proportion infected with *E. coli* was set at

19     $P_{coli}$ = 0%, 50% and 100% (the other IMI+ cows were considered infected

20     with *S. aureus*).  If a cow was assigned to the IMI+ group, the time at which

21     the  clinical  episode  starts  (= t*)  was  sampled  from  an  exponential

1    distribution with scale parameter 3 which is in agreement with the reported

2    median time of first occurrence of mastitis, i.e., 2 to 3 months [6].

3

4    **2.8.2. Evaluation of the accuracy of the estimates**

5

6    The estimates $(\hat{\mu}_i^t, \hat{\sigma}_0^2, \hat{\sigma}_1^2, \hat{\sigma}_a^2, \hat{\mathbf{a}})$ of the parameters $(\mu_i^t, \sigma_0^2, \sigma_1^2, \sigma_a^2, \mathbf{a})$ were

7    computed, after burn-in, as the means of the posterior distributions. Their

8    accuracies were assessed over the range of parameter values (sensitivity

9    analysis) as follows.   For the predicted breeding values, the Spearman

10   correlation coefficient ($corr_{BV}$) with the true breeding values was computed

11   for each replicate and averaged over the 10 replicates.   For residual and

12   additive variances, the differences ($bias_{\sigma 0}$, $bias_{\sigma 1}$, and $bias_{\sigma a}$) between

13   estimates and simulated values were computed for each replicate and

14   averaged over the 10 replicates.   For the location parameters, the biases

15   ($bias_{\mu 0}$ and $bias_{\mu 1}$) were calculated between the estimates and $\overline{y}_i^t$ , where

16   $$\overline{y}_i^t = \frac{\sum\limits_{k=1,n_t^i} (y_k^t \mid z_k^t = i)}{n_t^i}$$ computed with known values for $z_k^t$.   Finally,

17   sensitivity (SE), specificity (SP), and probability of correct classification

18   (PCC), were computed at each iterative step as:

1 $$SE = \sum_{k=1,N} \sum_{t=1,T} p\left(\hat{z}_k^t = 1 \middle| z_k^t = 1\right), SP = \sum_{k=1,N} \sum_{t=1,T} pr\left(\hat{z}_k^t = 0 \middle| z_k^t = 0\right), \text{ and}$$

2 $$PCC = \sum_{k=1,N} \sum_{t=1,T} pr\left[(z_k^t = 1 \cap \hat{z}_k^t = 1) \cup (z_k^t = 0 \cap \hat{z}_k^t = 0)\right].$$

3  After burn-in, these were averaged over the $\delta$ Gibbs rounds and the 10

4  replicates.

5

6  **3. RESULTS AND DISCUSSION**

7

8  From visual inspection of the algorithmic convergence, it was found that a

9  total of 1,000 cycles and a burn-in ($\rho$) of 200 runs were sufficient to remove

10  the influence of the prior values and obtain stable estimates.  All results

11  presented thus correspond to the last ($\delta = 800$) runs of the Gibbs algorithm.

12  This may seem very few cycles but results were checked for 3 simulated

13  data sets over a higher number of cycles of the Gibbs sampler.  Convergence

14  rates were also checked with an EM algorithm and the Gibb sampler on

15  models similar to those used in the simulation of this study but without

16  genetic covariance structure ($\mathbf{SCS_i} = \mathbf{M_i}\ \boldsymbol{\mu_i} + \mathbf{e_i}$ ).  Explanations may be

17  linked to the simplicity of the pedigree structure, small number of cows and

18  the fact that value for $m_0$ and $s^2_p$ were obtained from the data.

19

**3. 1. Overall accuracy of the estimates**

Overall, the sensitivity was high (SE ~ 90%) but the specificity low (SP ~ 60%).  Because of this high sensitivity, we can be confident that a cow with $\hat{z}_k^t = 0$ is healthy and spare the costs of further testing (e.g., bacteriological cultures) or useless treatment. On the other end, the low specificity indicates that cows with $\hat{z}_k^t = 1$ should be further tested to confirm the clinical suspicion.  These observations may suggest some economic interest in HMM.

Before any testing, the probability for a cow to be IMI+ can only be estimated from the prevalence of the disease in the population, while, after testing, this probability is estimated from the posterior probability of being IMI+ given a positive test (also called the positive predictive value).  With SE = 90% and SP = 60%, the difference between prior and posterior probabilities is maximum at disease frequencies between 20% and 50%, with posterior probabilities 20% higher than the prior probabilities.  These frequencies are within the range of prevalence typically reported for mastitis, as illustrated in the following few studies.  In Finland, Pitkälä *et al.* [18] reported 31% of cows with SCC>300,000/mL (mastitis) in 2001.  In Switzerland, Roesch *et al.* [19] reported 40% of cows showing at least 1 positive California Mastitis Test in at least one quarter at 31 d and 102 d

1    post partum. In a survey of clinical and subclinical mastitis in England and

2    Wales, the mean incidence of clinical mastitis recorded by the farmer was 47

3    cases per 100 cows per year [3]. In Canada, Sargeant *et al.* [21] observed

4    19.8% of cows experienced one or more cases of clinical mastitis during a 2-

5    year observational study. Therefore, HMM may also be of interest in field

6    studies, when it is necessary to precisely identify infected cows.

7    Breeding values from the HMM seemed accurate in predicting the true

8    additive genetic merit of the cows. Indeed, the correlation (corr$_{BV}$) between

9    simulated and estimated breeding values varied from 65% to 79% over the

10   whole data sets. This is close to the correlations of 70 to 75% computed as

11   the square root of the coefficient of determination (CD), where

12   $CD = 1 - \dfrac{PEV}{V}$, PEV = prediction error variance = $[\mathbf{W'R^{-1}\,W} + \mathbf{A}^{-1}/\sigma_a{}^2]^{-1}$

13   and V = true additive variance = $\mathbf{A}\,\sigma_a{}^2$ [11]. The PEV were computed with

14   the values of the parameters used in the simulation and weighted by the true

15   proportion of IMI- and IMI+ per cow.

16   On the other hand, the HMM was less efficient in estimating the

17   parameters for the IMI+ group. Indeed, $\hat{\sigma}_1^2$ had a tendency to underestimate

18   and $\hat{\mu}_1^t$ to overestimate the values used in the simulation. The biases varied

19   from -1.33 to -0.13 (mean = -0.59) for $\hat{\sigma}_1^2$ and from -0.02 to 3.26 (mean =

20   1.14) for $\hat{\mu}_1^t$. The magnitude of the biases decreased when the amount of

information available on IMI+ cows increased, as discussed in the sensitivity

analyses below.


**3. 2. Sensitivity analyses**


The robustness of the HMM approach was assessed by computing the

biases in the estimates over a wide range of values for the simulated

parameters.

Over all, estimates of means and variances were rather insensitive to the

values of the corresponding simulated values but they were sensitive to the

proportion of cows with at least one IMI+ sample ($P_{cow}$) and to the

proportion of *E. coli* among infected cows ($P_{coli}$). This suggests HMM

estimates are sensitive to the amount of data available to compute them. For

example, biases in the estimation of both location parameters ($\hat{\mu}_0^t, \hat{\mu}_1^t$) were

highest when $P_{cow}$ was lowest (Figure 2), suggesting that it is necessary to

have a sufficient number of observations per cow when the disease

prevalence is low. Similarly, SE, SP and PCC decreased as the proportion

of *E. coli* infection ($P_{coli}$) increased (Figure 3). This was not surprising

because, in cows infected with *E. coli*, only few simulated SCS were higher

than SCS for IMI- samples, as is observed in naturally occurring *E. coli*

infections usually of short duration.

1    Level of response to infection influenced estimates of transition

2    probabilities, in opposition to estimates of both location parameters and of

3    breeding values.  For example, SE and PCC were higher among high (SE =

4    92%; PCC = 64%) than moderate (SE = 80%; PCC = 60%) responders

5    suggesting that HMM is more accurate when IMI- and IMI+ distributions

6    are farther apart.  Conversely, accuracy of $\hat{\sigma}_1^2$ worsened when the distance

7    between IMI- and IMI+ distributions increased with $bias_{\sigma 1}$ = -0.51 for

8    moderate and $bias_{\sigma 1}$ = -0.80 for high responders.

9    Note that SE and SP were insensitive to change in disease frequency

10   ($P_{cow}$), as they should by definition, conversely to PCC that is, by definition,

11   a function of the disease frequency:  PCC = [SE*pr(IMI+)] + [SP*pr(IMI-)].

12   Finally, note SE and SP reported here are different from SE and SP in

13   Ødegård *et al.* [17] in which

14   $$SE = \frac{\sum\limits_{i=1,n} t_i PPM_i}{\sum\limits_{i=1,n} t_i} \text{ and } SPE = \frac{\sum\limits_{k=1,n}(1 - t_i)(1 - PPM_i)}{n - \sum\limits_{i=1,n} t_i}$$

15   where $PPM_i$ is the posterior mean of the estimates of $z_i$ averaged over Gibbs

16   samples (after burn-in), $t_i$ = 0 if IMI-, $t_i$ = 1 if IMI+, and i = 1 to n cows.

17

18   **3. 3. General discussion**

19

1    The main advance of this paper is the presentation of a HMM in which

2    genetic random effects were added to the conditional model for the observed

3    data.  In the subject-area literature, HMMs with random effects have been

4    used in a very limited way.  Only recently, Altman [1] introduced a mixed

5    HMM to study lesion counts in multiple sclerosis patients.  In her model,

6    parameters for the observed and hidden data are allowed to vary randomly

7    among patients, although they are assumed independent from each other (no

8    genetic relationship).  This suggests a natural extension of the present

9    HMM, i.e., to allow the parameters of the hidden Markov chain to vary

10   randomly among cows.  However, interpretation of the results of such

11   extended model will be delicate because sets of identical genes may be

12   associated to both IMI and SCS (confounding effects).  Stated otherwise, the

13   total genetic effects on SCS would be a combination of the effects of genes

14   responsible for presence or not of IMI (resistance to infection) and for the

15   magnitude of the SCS response after IMI (tolerance after infection).

16   Structural equation modeling is a technique to evaluate models with

17   different hypothesized relationships among variables.  In this context, it

18   would be interesting to evaluate the different models proposed in Figure 4 to

19   determine the amount of relationships between genes insuring tolerance or

20   resistance to infection.  In the model proposed here, biomarker value at one

21   specific time is independently influenced by the IMI status and by some

22   genes.  But, both the IMI status and the biomarker values could also be

1 under the influence of this same set of genes (model b of Figure 4). The

2 relationship between genes, biomarker and IMI status can become even

3 more complicated with different sets of correlated genes influencing the

4 expression of both traits (model e). This is important for the long-term

5 because some epidemiological models predict that selection for resistant

6 cows (no infection) may not be as durable as selection for tolerant (infection

7 but no disease) cows [16], [20]. Increased resistance would reduce disease

8 transmission, reducing the fitness advantage of carrying the resistant genes,

9 and possibly impose pressure upon the pathogen to evade the control

10 strategy. By contrast, as genes conferring disease tolerance spread within a

11 population, the disease incidence rises, increasing the evolutionary

12 advantage of carrying the tolerance genes, without leading to genetic

13 changes in the parasite population.

14 Other extensions of the HMM are possible. Trends and seasonality in

15 SCS can be readily accommodated to relax the assumption of time-

16 independence between transition probabilities [15]. Prior information on the

17 parameters can be included to increase accuracy and speed up convergence.

18 Location parameters can be made more realistic by considering effects

19 affecting SCS values, such as age, herd, or season, as a few examples.

20 Elements of the M matrices could take different values than zeroes or ones

21 to reflect the different effects on SCS for different parts of the lactation.

1  The genetic variance could also be different for IMI- and IMI+ samples and

2  would allow for genetic difference in the response in SCS to IMI.

3  The first-order Markov assumption is also a limiting feature of the HMM

4  and mechanisms of transmission of the IMI between cows could also be

5  considered more precisely in deriving the transition probabilities.  Indeed,

6  transmission of infection is a complex process that involves the mixing

7  structure of the population (as it determines the probability of contact

8  between animals), the infectiousness of the contagious animal (or infective

9  dose) and the susceptibility of a healthy cow (i.e., its probability of getting

10  infected after contact with a contagious animal).  To solve these issues,

11  Cooper and Lipsitch [5] proposed to model the transition probabilities of the

12  hidden Markov chain in terms of the parameters of epidemiologic models

13  used to describe the transmission of an infectious disease at the population

14  level.

15

16  **3. 4. Conclusions**

17

18  In summary, it is shown that the mixed HMM provides a good fit to the

19  data sets simulated in this study.  The advantages of the HMM over other

20  approaches are the prediction of health or disease status, the reduction of

21  confirmatory diagnosis costs, and the increased accuracy in breeding values.

22  However, future work needs to be done to extend the HMM proposed here,

1 the most important piece of which is to quantify the level of resistance and

2 tolerance to infection while considering the mechanisms of transmission

3 between healthy and sick cows.

4

5 **4. ACKNOWLEDGMENTS**

6

9

10 **5. REFERENCES**

11

12 [1] Altman R. M., Mixed hidden Markov model:  An extension of the

13 hidden Markov model to the longitudinal data setting, J. Am. Stat. Assoc.

14 102 (2007) 201-210.

15

16 [2] Boettcher P.J., Moroni P., Pisoni G., Gianola D., Application of finite

17 mixture model to somatic cell scores of Italian goats, J. Dairy Sci. 88 (2005)

18 2209-2216.

19

20 [3] Bradley A. J., Leach K. A., Breen J. E., Green L. E., Green M. J., Survey

21 of the incidence and aetiology of mastitis on dairy farms in England and

22 Wales,  Vet. Rec. 160 (2007) 253-257.

1

[4] Carlén E., Strandberg E., Roth A., Genetic parameters for clinical mastitis, somatic cell score, and production in the first three lactations of Swedish Holstein cows, J. Dairy Sci. 87 (2004) 3062-3070.

[5] Cooper B., Lipsitch M., The analysis of hospital infection data using hidden Markov model, Biostat. 5 (2004) 223-237.

[6] de Haas Y., Barkema H. W., Veerkamp R. F., The effect of pathogen-specific clinical mastitis on the lactation curve for somatic cell count, J. Dairy Sci. 85 (2002) 1314-1323.

[7] de Haas Y., Veerkamp R. F., Barkema H. W., Gröhn Y. T., Schukken Y. H., Associations between pathogen-specific cases of clinical mastitis and somatic cell count patterns, J. Dairy Sci. 87 (2004) 95 – 105.

[8] Detilleux J.C., Leroy P., Application of a mixed normal mixture model for the estimation of mastitis-related parameters, J. Dairy Sci. 83 (2000) 2341–2349.

[9] Detilleux J., Genetic factors affecting susceptibility to udder pathogens, Vet. Microbiol. (accepted)

1

[10] Eisner J., An interactive spreadsheet for teaching the Forward-Backward algorithm, in: Proceedings of the ACL workshop on effective tools and methodologies for teaching NLP and CL, July 2002, Philadelphia, pp.10-18.

[11] Fouilloux M-N., Laloë D., A sampling method for estimating the accuracy of predicted breeding values in genetic evaluation, Genet. Sel. Evol. 33 (2001) 473-486.

[12] Gianola D., Prediction of random effects in finite mixture models with Gaussian components, J. Anim. Breed. 122 (2005) 145-159.

[13] Heringstad B., Gianola D., Chang Y. M., Ødegård J., Klemetsdal G., Genetic associations between clinical mastitis and somatic cell score in early first-lactation cows, J. Dairy Sci., 89 (2006) 2236 – 2244.

[14] Hernández A., Karrow N., Mallard B. A., Evaluation of immune responses of cattle as a means to identify high and low responders and use of a human microarray to differentiate gene expression, Genet. Sel. Evol. 35 (2003) 67-81.

1  [15] Le Strat Y., Carrat F.,  Monitoring epidemiologic surveillance data

2  using hidden Markov models,  Stat. Medicine 18 (1999) 3463-3478.

3

4  [16] Miller M. R., White A., Boots M., The evolution of host resistance:

5  Tolerance and control as distinct strategies,  J. Theor. Biol. 236 (2005) 198-

6  207.

7

8  [17] Ødegård J., Jensen J., Madsen P., Gianola D., Klemetsdal G.,

9  Heringstad B., Detection of mastitis in dairy cattle by use of mixture models

10  for repeated somatic cell scores: A Bayesian approach via Gibbs sampling,

11  J. Dairy Sci. 86 (2003) 3694–3703.

12

13  [18] Pitkälä A., Haveri M., Pyörälä S., Myllys V., Honkanen-Buzalski T.,

14  Bovine mastitis in Finland 2001 − Prevalence, Distribution of bacteria, and

15  antimicrobial resistance,  J. Dairy Sci., 87 (2004) 2433-2441.

16

17  [19] Roesch M., Doherr M. G., Schären W., Schällibaum M., Blum J. W.,

18  Subclinical mastitis in dairy cows in Swiss organic and conventional

19  production systems,  J. Dairy Res. 74 (2007) 86-92.

20

21  [20] Roy B. A., Kirchner J. W., Evolutionary dynamics of pathogen

22  resistance and tolerance,  Evolution 54 (2000) 51-63.

1

2    [21] Sargeant J.M., Scott H. M., Leslie K. E., Ireland M. J., Bashiri A.,

3    Clinical mastitis in dairy cattle in Ontario: frequency of occurrence and

4    bacteriological isolates, Can. Vet J. 39 (1998) 33-38.

5

6    [22] Wenz J. R., Barrington G. M., Garry F. B., McSweeney K. D.,

7    Dinsmore P., Goodell G., Callan R. J., Bacteremia associated with naturally

8    occurring coliform mastitis in dairy cows, J. Am. Vet. Med. Assoc., 219

9    (2001) 976 – 981.

10

11

**FIGURE CAPTIONS**

**Figure 1.** Means of SCS for lactations without clinical mastitis (plain line) and lactations with clinical mastitis associated with *Staphylococcus aureus* (square) or *Escherichia coli* (triangle) occurring on the median MIM for multiparous cows (Adapted from de Haas *et al.*, [6]) .

**Figure 2.** Differences between simulated and estimated values for the means of the distributions for healthy (dot bar) and infected (slash bar) cows as a function of the proportion of infected cows.

**Figure 3.** Sensitivity (plain bar), specificity (open bar) and probability of correct classification (circled bar) as a function of the proportion of *E. coli* among infected cows.

**Figure 4.** Five different hypothetical models of the relationship between genetic background (G), intra-mammary infection (IMI) and biomarker (Bio). The first model (a) is the model of this study (The dependent variables are the targets of one-headed arrows).
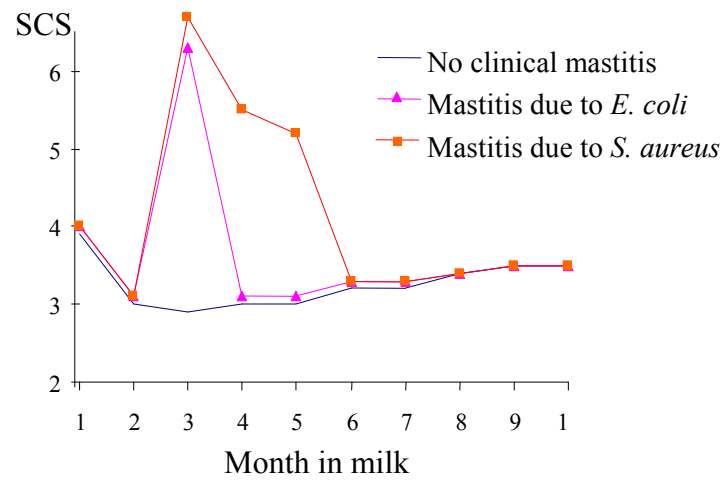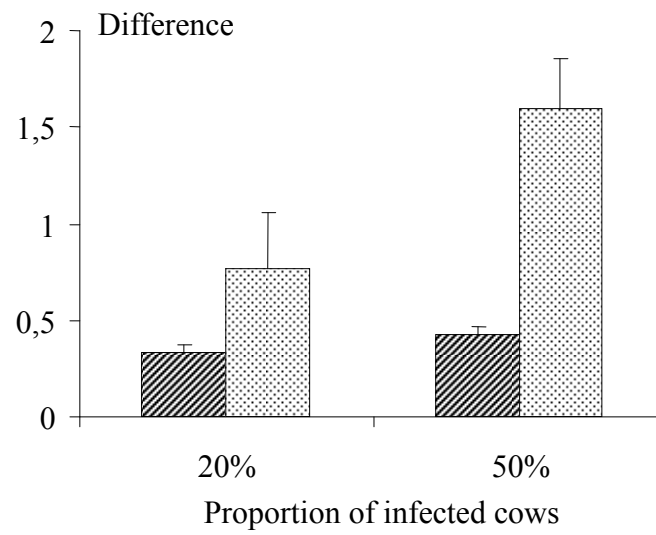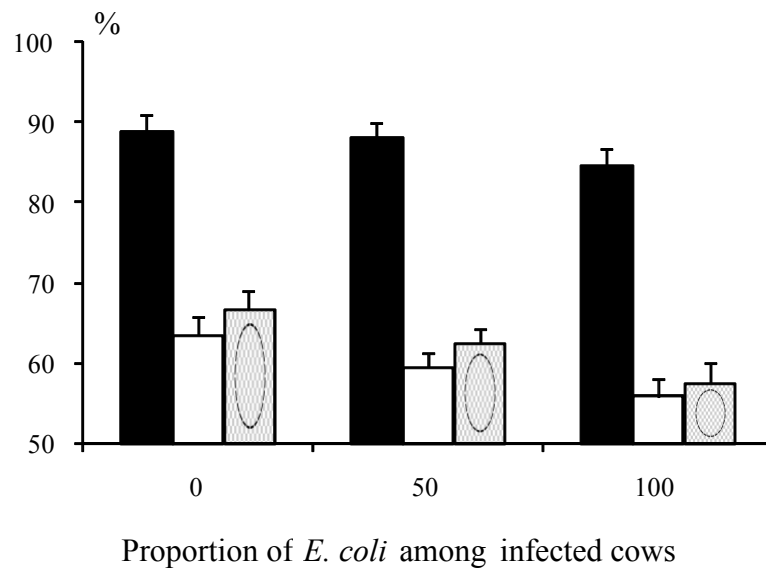
1

2    Figure 1

3
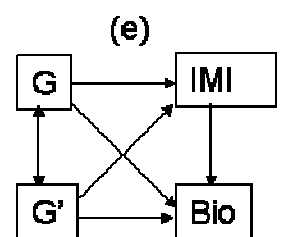
4

1

2    Figure 2.



3

4

1    Figure 3.

2



Proportion of *E. coli* among infected cows

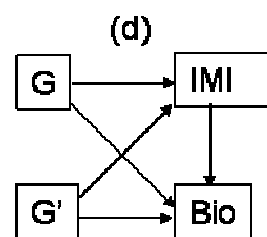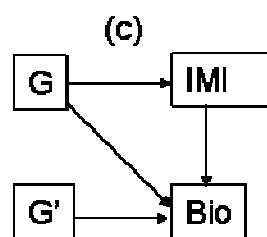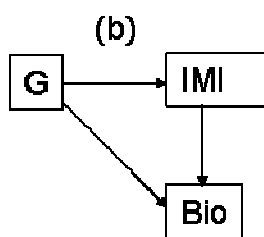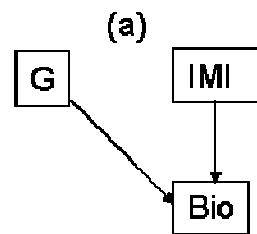3

4

5

6

1    Figure 4.



2

3

4

2

**Table I**.  Sensitivity (SE), specificity (SP), and probability of correct

classification (PCC) as a function of the level of response to infection, high

(H) or moderate (M) responders, number of samples per cow (T), percentage

of cows with at least one IMI+ sample ($P_{cow}$), percentage infected with *E.*

*coli* ($P_{coli}$), and residual and additive genetic variances ($\sigma_0^2, \sigma_1^2, \sigma_a^2$ ).  Data

sorted by SE.

9

| SE | SP | PCC | T | $P_{cow}$ | $P_{coli}$ | $\sigma_0^2$ | $\sigma_1^2$ | $\sigma_a^2$ |
|---|---|---|---|---|---|---|---|---|
| High responders | | | | | | | | |
| 95.03 | 59.65 | 63.70 | 10 | 50 | 50 | 1.0 | 1.0 | 0.15 |
| 94.50 | 58.19 | 60.64 | 10 | 20 | 0 | 1.4 | 1.4 | 0.15 |
| 94.25 | 49.59 | 56.73 | 10 | 20 | 50 | 1.4 | 1.4 | 0.15 |
| 94.03 | 58.05 | 59.90 | 20 | 20 | 50 | 1.0 | 1.0 | 0.25 |
| 93.92 | 62.71 | 65.98 | 20 | 50 | 0 | 1.0 | 1.0 | 0.25 |
| 93.79 | 58.88 | 60.63 | 20 | 20 | 50 | 1.4 | 1.4 | 0.25 |
| 93.20 | 57.51 | 59.31 | 20 | 20 | 50 | 1.4 | 1.4 | 0.25 |
| 93.08 | 55.15 | 56.95 | 10 | 20 | 50 | 1.4 | 1.4 | 0.25 |
| 92.64 | 58.23 | 62.16 | 10 | 50 | 50 | 1.4 | 1.4 | 0.15 |
| 92.64 | 65.99 | 68.16 | 20 | 20 | 0 | 1.4 | 1.4 | 0.25 |
| 92.63 | 57.49 | 58.34 | 20 | 20 | 50 | 1.4 | 1.4 | 0.25 |
| 92.03 | 59.91 | 61.49 | 20 | 20 | 50 | 1.4 | 1.4 | 0.25 |
| 90.41 | 50.89 | 51.65 | 10 | 20 | 100 | 1.4 | 1.4 | 0.15 |
| 89.58 | 50.60 | 51.34 | 10 | 20 | 100 | 1.4 | 1.4 | 0.15 |
| 89.05 | 69.75 | 73.53 | 20 | 50 | 0 | 1.0 | 1.0 | 0.15 |
| 88.81 | 68.09 | 72.19 | 20 | 50 | 0 | 1.4 | 1.4 | 0.25 |
| 88.19 | 66.02 | 70.42 | 20 | 50 | 0 | 1.4 | 1.4 | 0.25 |
| 88.14 | 68.43 | 72.38 | 20 | 50 | 0 | 1.0 | 1.4 | 0.15 |
| 85.06 | 68.53 | 71.84 | 20 | 50 | 0 | 1.0 | 1.4 | 0.25 |

| 84.27 | 55.36 | 55.94 | 20 | 20 | 100 | 1.4 | 1.4 | 0.25 |
|---|---|---|---|---|---|---|---|---|
| Moderate responders | | | | | | | | |
| 94.24 | 57.41 | 59.28 | 20 | 20 | 50 | 1.0 | 1.0 | 0.25 |
| 79.74 | 52.41 | 52.95 | 20 | 20 | 50 | 1.0 | 1.0 | 0.25 |
| 79.09 | 54.89 | 56.74 | 20 | 20 | 0 | 1.4 | 1.4 | 0.25 |
| 77.95 | 53.64 | 54.81 | 20 | 20 | 50 | 1.4 | 1.4 | 0.25 |
| 77.67 | 64.32 | 67.03 | 20 | 50 | 0 | 1.0 | 1.4 | 0.15 |
| 77.06 | 63.14 | 65.90 | 20 | 50 | 0 | 1.0 | 1.4 | 0.25 |
| 75.77 | 51.78 | 52.24 | 20 | 20 | 100 | 1.4 | 1.4 | 0.25 |
| 73.04 | 58.81 | 61.60 | 20 | 50 | 0 | 1.0 | 1.4 | 0.25 |

1

2

**Table II.** Accuracy of the estimates of the mixed hidden Markov model as a function of the level of response to infection, high (H) or moderate (M) , number of samples per cow (T), percentage of cows with at least one IMI+ sample ($P_{cow}$), percentage infected with *E. coli* ($P_{coli}$), and residual and additive genetic variances ($\sigma_0^2, \sigma_1^2, \sigma_a^2$ ). The accuracy is determined as the differences between values used in the simulations and estimates of means (bias$_{\mu0}$, bias$_{\mu1}$) and residual variances (bias$_{\sigma0}$, bias$_{\sigma1}$) in IMI- and IMI+ cows, respectively; the differences between values used in the simulations and estimates of additive genetic variance (bias$_{\sigma a}$); and the correlation between predicted and simulated breeding values (corr$_{BV}$). Data sorted by corr$_{BV}$.

| corr$_{BV}$ | bias$_{\sigma0}$ | bias$_{\sigma1}$ | bias$_{\sigma a}$ | bias$_{\mu0}$ | bias$_{\mu1}$ | T | $P_{cow}$ | $P_{coli}$ | $\sigma_0^2$ | $\sigma_1^2$ | $\sigma_a^2$ |
|---|---|---|---|---|---|---|---|---|---|---|---|
| High responders | | | | | | | | | | | |
| 0.79 | 0.00 | -0.66 | -0.08 | 0.24 | 0.47 | 20 | 50 | 0 | 1.0 | 1.4 | 0.15 |
| 0.79 | 0.02 | -0.65 | -0.02 | 0.21 | 0.28 | 20 | 50 | 0 | 1.0 | 1.0 | 0.15 |
| 0.78 | -0.02 | -0.78 | 0.00 | 0.22 | 0.43 | 20 | 50 | 0 | 1.0 | 1.4 | 0.25 |
| 0.77 | 0.01 | -0.70 | 0.01 | 0.28 | 0.51 | 20 | 50 | 0 | 1.4 | 1.4 | 0.25 |
| 0.77 | 0.02 | -0.63 | 0.04 | 0.23 | 0.52 | 20 | 50 | 0 | 1.4 | 1.4 | 0.25 |
| 0.74 | -0.01 | -0.29 | 0.05 | 0.41 | 2.16 | 20 | 20 | 100 | 1.4 | 1.4 | 0.25 |
| 0.74 | 0.06 | -0.46 | -0.01 | 0.50 | 2.93 | 10 | 20 | 100 | 1.4 | 1.4 | 0.15 |
| 0.73 | 0.04 | -0.57 | 0.02 | 0.31 | 0.80 | 20 | 20 | 0 | 1.4 | 1.4 | 0.25 |
| 0.73 | 0.09 | -0.48 | -0.03 | 0.55 | 3.26 | 10 | 20 | 100 | 1.4 | 1.4 | 0.15 |
| 0.72 | 0.03 | -0.42 | 0.04 | 0.52 | 1.26 | 20 | 20 | 50 | 1.4 | 1.4 | 0.25 |
| 0.71 | 0.02 | -0.46 | 0.04 | 0.42 | 1.22 | 20 | 20 | 50 | 1.4 | 1.4 | 0.25 |
| 0.71 | 0.03 | -0.48 | 0.05 | 0.40 | 1.13 | 20 | 20 | 50 | 1.4 | 1.4 | 0.25 |
| 0.71 | 0.09 | -0.65 | -0.02 | 0.44 | 1.86 | 10 | 20 | 50 | 1.4 | 1.4 | 0.15 |
| 0.70 | 0.02 | -0.44 | 0.04 | 0.38 | 1.17 | 20 | 20 | 50 | 1.4 | 1.4 | 0.25 |
| 0.70 | 0.09 | -0.60 | 0.06 | 0.51 | 1.73 | 10 | 20 | 50 | 1.4 | 1.4 | 0.25 |
| 0.69 | 0.03 | -0.57 | 0.04 | 0.36 | 0.87 | 20 | 50 | 0 | 1.0 | 1.0 | 0.25 |

| | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 0.69 | 0.11 | -0.74 | -0.03 | 0.40 | 1.69 | 10 | 20 | 0 | 1.4 | 1.4 | 0.15 |
| 0.68 | 0.08 | -1.25 | -0.02 | 0.38 | 1.48 | 10 | 50 | 50 | 1.0 | 1.0 | 0.15 |
| 0.67 | 0.03 | -0.44 | 0.06 | 0.43 | 1.06 | 20 | 20 | 50 | 1.0 | 1.0 | 0.25 |
| 0.67 | 0.07 | -1.21 | -0.03 | 0.39 | 1.46 | 10 | 50 | 50 | 1.4 | 1.4 | 0.15 |
| Moderate responders | | | | | | | | | | | |
| 0.76 | -0.02 | -0.46 | -0.02 | 0.24 | 0.00 | 20 | 50 | 0 | 1.0 | 1.4 | 0.15 |
| 0.75 | -0.01 | -0.13 | 0.05 | 0.48 | 1.61 | 20 | 20 | 100 | 1.4 | 1.4 | 0.25 |
| 0.75 | -0.01 | -0.14 | 0.07 | 0.47 | 1.30 | 20 | 20 | 50 | 1.0 | 1.0 | 0.25 |
| 0.75 | -0.03 | -0.21 | 0.04 | 0.32 | 0.70 | 20 | 20 | 0 | 1.4 | 1.4 | 0.25 |
| 0.74 | -0.02 | -0.18 | 0.06 | 0.32 | 0.82 | 20 | 20 | 50 | 1.4 | 1.4 | 0.25 |
| 0.73 | -0.03 | -0.46 | 0.04 | 0.32 | 0.19 | 20 | 50 | 0 | 1.0 | 1.4 | 0.25 |
| 0.72 | -0.04 | -0.36 | 0.05 | 0.39 | -0.02 | 20 | 50 | 0 | 1.0 | 1.4 | 0.25 |
| 0.66 | 0.03 | -0.45 | 0.06 | 0.44 | 1.22 | 20 | 20 | 50 | 1.0 | 1.0 | 0.25 |

1