

1 **A hidden Markov model to predict early mastitis from**
2 **test-day somatic cell scores.**

3

4

5 J. C. Detilleux

6

7 *Department of Quantitative Genetics, Veterinary Faculty, University of Liège,*

8 *4000 Liège, Belgium.*

9

10

11

12

13

14

15 Corresponding author: Johann C. Detilleux. E-mail: jdetilleux@ulg.ac.be

16

17 Running head: Hidden Markov model in mastitis

18

1 **Abstract**

2

3 The absence of on farm recording systems in most countries precludes the
4 identification of clinical mastitis cases after its occurrence. Therefore, in many
5 countries high somatic cell scores (SCS) in milk are used as indicator for mastitis
6 because they are collected on a routine basis. However, individual test day SCS
7 are not very accurate in identifying infected cows. Mathematical models may
8 improve the accuracy of the biological marker by making better use of the
9 information contained in the available data. Here, a simple hidden Markov model
10 (HMM) was applied on SCS recorded monthly on cows with or without clinical
11 mastitis to evaluate its accuracy in estimating parameters (mean, variance and
12 transition probabilities) under health or disease states. The SCS means were
13 estimated at 1.96 (SD = 0.16) and 4.73 (SD = 0.71) for the hidden healthy and
14 infected states, and the common variance at 0.83 (SD = 0.11). The probabilities to
15 remain uninfected, to recover from infection, to get newly infected and to remain
16 infected between consecutive test-days were estimated at 78.84%, 60.49%,
17 11.70% and 15%, respectively. Three different health related states were
18 compared: clinical stages observed by farmers, subclinical cases defined for
19 somatic cell counts below or above 250,000 cells/mL and infected stages obtained
20 from the HMM. The results showed that HMM identifies infected cows before
21 the apparition of clinical and subclinical signs which may critically improve the
22 power of studies on the genetic determinants of SCS and reduce biases in
23 predicting breeding values for SCS.

1 Key words: **mastitis; hidden Markov model; somatic cell counts.**

2

3

4

5 **Implications**

6

7 In most countries, somatic cell counts (SCC) are routinely used as indicators of
8 mammary infection before milk is exploited for consumption before or after its
9 transformation. However, SCC are not very sensitive in classifying cows as
10 infected or healthy which leads to unnecessary costs and missed profits. Here, a
11 simple hidden Markov model is proposed that improved the diagnostic accuracy
12 of SCC by uncovering the hidden health status of the cows before the apparition
13 of clinical signs and before SCC exceed the threshold of 250,000 cells/mL. This
14 will critically improve the power of genetic studies of mastitis determinants and
15 reduce biases in predicting breeding values.

16

17

1 **Introduction**

2

3 The absence of on farm recording systems in most countries precludes the
4 identification of clinical mastitis cases after its occurrence. Therefore, in many
5 countries high somatic cell counts (SCC) in milk are used as indicator for sub-
6 clinical and clinical mastitis, especially for genetic evaluation to improve
7 resistance to mammary infections that necessitate large amount of data (Shook
8 and Schutz, 1994). However, the problem of identifying infected cows based on
9 their SCC is still not satisfactorily solved as individual SCC are not very sensitive
10 in diagnosing mammary infection, either at the quarter or cow levels (Djabri *et al.*,
11 2002; Sargeant *et al.*, 2001). This has relevant impact in animal selection
12 because imperfect accuracy in the diagnosis of infectious diseases results in a
13 reduction of heritability estimates (Bishop and Woolliams, 2010). It is also a
14 source of misclassification as uninfected animals may have high SCC (and
15 reversely). This may bias prediction of breeding values and decrease the power to
16 detect association between a disease locus and a marker locus (Buyske *et al.*,
17 2009). Selection for very low SCC might even not be the good objective because
18 low initial SCC has been associated with increased susceptibility and severity of
19 subsequent mastitis (Suriyasathaporn *et al.*, 2000).

20 Mathematical models improve the accuracy of SCC measures used to identify
21 infected cows by making better use of the information contained in SCC data. For
22 example, models developed by de Haas *et al.* (2004) lead to the identification of
23 different SCC patterns according to the mammary pathogen: Clinical *E. coli*

1 mastitis is significantly associated with the presence of a short peak in SCC
2 whereas *S. aureus* is associated with long increased SCC. Others have used the
3 finite mixture model (FMM) methodology on SCC to infer the cow's individual
4 probability of being infected (Detilleux and Leroy, 2000; Gianola, 2005). A
5 simple FMM will assign SCC to one of two components hopefully representing
6 SCC from cows with (IMI+) and without (IMI-) intra-mammary infection (IMI),
7 respectively. Then, the identification of animals at risk is computed as the
8 posterior probability of putative IMI, given SCC, rather than on crude SCC.
9 However, after bacteriological examinations of goat milk samples, Boettcher *et al.*
10 (2005) observed their FMM was able to classify correctly only 60% and 48% of
11 the healthy and infected records, respectively. If these results are not
12 encouraging, it should be noted the accuracy for detecting an IMI from
13 bacteriological cultures of single composite or quarter milk samples in
14 subclinically infected cows is known to be low (Lam *et al.*, 1996; Sears *et al.*,
15 1990). This is because pathogens such as *S. aureus* are often shed in an
16 intermittent or cyclical pattern and in numbers too low to be detected by
17 conventional culturing methods (Godden *et al.*, 2002). The *S. aureus* and
18 coagulase negative staphylococci were the most prevalent pathogens in the above
19 mentioned goat study (Moroni *et al.*, 2005).

20 Hidden Markov models (HMM) could be an alternative to FMM. A HMM is
21 defined as a finite set of states, each of which is associated with a probability
22 distribution. Transitions among the states are governed by a set of probabilities
23 called transition probabilities. The joint distribution over all states is a Markov

1 chain. An observation is associated to each state, according to a linked
2 probability distribution, called the emission probability. Only observations are
3 recordable, not the states that are ``hidden"; hence the name (Rabiner, 1989). In
4 the mastitis context, the health status of the mammary gland could be considered
5 as two hidden states and the SCC as the associated observations.

6 In a simulated data set (Detilleux, unpublished results), the accuracy of estimates
7 obtained with a FMM was increased by incorporating information from previous
8 SCC, as is done in HMM. In another study (Detilleux, 2008), estimates obtained
9 with a mixed HMM were close to the true values unless the prevalence of the
10 disease is low.

11 The objective of this study is to present the mathematical formalism behind the
12 HMM methodology, to apply the model on SCC collected on first parity cows
13 with known clinical status and to compare results on clinical (observed by
14 farmers), subclinical (defined for SCC or above 250,000 cells/mL) and hidden
15 (infected or not as obtained from the HMM) states.

16

17 **Materials and methods**

18

19 *Animals and data collection*

20 Data from the field study of Barkema *et al.* (1998) were used. Briefly,
21 bacteriological samples were collected by the farmers from cows with signs of
22 clinical mastitis. For the present analyses, only the first cases of clinical mastitis
23 per lactation (CM₁) were considered and bacteriological results were consolidated

1 into negative ($B_t = 0$) and positive results ($B_t = 1$), with t representing the month in
2 milk (MIM) at which the case was recorded. Records on clinical case were
3 retrieved between December 1992 and August 1995. Conjointly, a total of
4 526,867 test days with SCC were recorded by the National Milk Recording
5 System (NRS, Arnhem, The Netherlands). The somatic cell scores (SCS) were
6 computed as $\log_2(\text{SCC}/100,000) + 3$ and averaged per MIM. After editing (birth
7 year > 1960, $\text{SCC} < 9,999,000$, $\text{MIM} \leq 10$, calving date \leq test-day date), the data set
8 included 128,748 records on SCC for the first 10 MIM, on 21,829 1st parity cows.
9 A total of 951 mastitis cases were reported of which 774 were bacteriologically
10 positive. Thereafter, clinical cases without positive bacteriological findings were
11 considered as healthy.

12 For each MIM, three (two observed and one hidden) different health states were
13 considered. The records were classified as being from a heifer with (CM+) or
14 without (CM-) a reported clinical case. The SCC may be below (SCM-) or above
15 (SCM+) the threshold of 250,000 cells/ml. This threshold was chosen as an
16 indicator of subclinical mastitis and is the one chosen by de Haas *et al.* (2002) in
17 her previous analyses of the data. The last stage is the hidden infected (IMI+) or
18 uninfected (IMI-) stages that were obtained by the HMM.

19

20 *Statistical analyses*

21 Throughout, k indexes the individual cow, t is the MIM, y_k^t is the SCS observed at
22 t on animal k , and z_k^t is the unknown state with $z_k^t = 0$ if y_k^t is from a hidden IMI-
23 sample and $z_k^t = 1$ if y_k^t is from a hidden IMI+ sample. On each cow, data

1 consists of a series of repeated SCS: $\mathbf{y}_k = \{y_k^1, y_k^2, \dots, y_k^T\}$ and the unobserved
 2 vector is $\mathbf{z}_k = \{z_k^1, z_k^2, \dots, z_k^T\}$, for $t = 1, 2, \dots, T$. For simplicity, T is assumed
 3 constant for all cows.

4
 5 *General formulation of the model.* A simple first-order HMM was assumed with
 6 2 transient states corresponding to the hidden IMI- and IMI+ categories with the
 7 following parameters:

8 - probabilities of transition between hidden states:

$$9 \quad a_k^{00} = \Pr(z_k^{t+1} = 0 \mid z_k^t = 0), a_k^{01} = \Pr(z_k^{t+1} = 1 \mid z_k^t = 0),$$

$$a_k^{10} = \Pr(z_k^{t+1} = 0 \mid z_k^t = 1), a_k^{11} = \Pr(z_k^{t+1} = 1 \mid z_k^t = 1),$$

10 - probability of being IMI- as an initial hidden state $\lambda_k = \Pr(z_k^1 = 0)$, and

11 - probabilities of SCS emission:

$$12 \quad (y_k^t \mid z_k^t = 0) \sim N(\mu_0^t, \sigma^2) \text{ and } (y_k^t \mid z_k^t = 1) \sim N(\mu_1^t, \sigma^2).$$

13 The probabilities of transition represent the probabilities of observing a
 14 particular hidden (unknown) IMI state at time $t + 1$, given the hidden IMI state at
 15 time t . The probabilities of emission represent the probabilities of observing SCS
 16 (at time t) given the hidden IMI state (at time t). It is assumed that correlation
 17 between successive SCS is fully accounted for by the underlying Markov process
 18 structure so that each SCS are independent given the unknown IMI state (output
 19 independence assumption). It is also assumed that state transition probabilities are
 20 independent of the actual time at which the transition takes place and do not
 21 change across time (stationary assumption). Finally, it is assumed that values in

any hidden state are only influenced by the values of the state that directly preceded it (first-order Markov assumption). The suitability of these assumptions for analyzing repeated SCS are discussed afterward.

To obtain the maximum likelihood estimates (MLE) of the parameter set θ_k^t , where $\theta_k^t = (\lambda_k, a_k^{00}, a_k^{01}, a_k^{10}, a_k^{11}, \mu_0^t, \mu_1^t, \sigma^2)$, the likelihood of the data must be maximized over all possible values of θ_k^t and this can be done through the expectation maximization (EM) algorithm.

Likelihood of the data. For one cow, the likelihood of one particular sequence of repeated SCS scores is given by:

$$\text{pr}(\mathbf{y}_k | \theta_k^t) = \alpha_{0,k}^t \beta_{0,k}^t + \alpha_{1,k}^t \beta_{1,k}^t,$$

with $\alpha_{i,k}^t = \text{pr}(y_k^1, y_k^2, \dots, y_k^t, z_k^t = i | \theta_k^t)$ and

$$\beta_{i,k}^t = \text{pr}(y_k^{t+1}, y_k^{t+2}, \dots, y_k^T | z_k^t = i, \theta_k^t),$$

for $i = 0$ and 1 . The $\alpha_{i,k}^t$ represents the probability of a partial sequence and ending up in state i at time t and $\beta_{i,k}^t$ represents the probability of a partial sequence starting from $t + 1$ to T given that the sequence started at state i at time t . This likelihood must be computed, for each cow, over all possible sequences of hidden states (\mathbf{z}_k). To do so, the naive way would be to sum, for each cow, the probabilities over all possible state sequences but their number can be huge ($= 2^T$) and the more efficient forward-backward algorithm is used in practice. This algorithm takes advantages of the sequential nature of the data, going forward ($t =$

1 1, 2, ..., T) and backward (for $t = T, T-1, \dots, 1$) in time, knowing it must end in
 2 some particular state. For a practical description, see Eisner (2002) and its
 3 interactive spreadsheet for teaching the algorithm. After the likelihood is
 4 computed for one cow, the likelihood for all sequences of all cows is computed as
 5 the product of all individual likelihoods (assumption of independence between
 6 cows).

7 *The EM algorithm.* For a detailed derivation of the algorithm for HMM, please
 8 refer to Bilmes (1998) and Rabiner (1989). In short, the EM algorithm consists of
 9 a series of repeated E and M steps. In the E-step, one finds the expected value of
 10 the complete-data log-likelihood with respect to the unknown parameters, given
 11 the observed data (\mathbf{y}_k) and the current parameter estimates ($\theta_k^{(p)}$ at iteration p). To
 12 form the complete data, one assumes both observed (\mathbf{y}_k) and hidden (\mathbf{z}_k) vectors
 13 are known. Then, the expected complete-data log-likelihood is written as:

$$\begin{aligned}
 14 \quad & \sum_{k=1, N} \{ E[z_k^1 = 0 | \mathbf{y}_k, \theta_k^{(p)}] \log(\lambda_k) + E[z_k^1 = 1 | \mathbf{y}_k, \theta_k^{(p)}] \log(1 - \lambda_k) \\
 15 \quad & + \sum_{t=1, (T-1)} \sum_{i,j=0,1} E[z_k^t = i, z_k^{t+1} = j | \mathbf{y}_k, \theta_k^{(p)}] \log(a_k^{ij}) \\
 16 \quad & + \sum_{t=1, (T-1)} \sum_{i,j=0,1} E[z_k^t = i | \mathbf{y}_k, \theta_k^{(p)}] \log p(y_k^t | z_k^t = i) \},
 \end{aligned}$$

17 where the first two terms of the summation involve observations at the start of the
 18 sequence ($t = 1$), the third term counts how many times each i to j transition
 19 occurred in the sequence and the fourth includes all observations generated from
 20 state i .

1 In the M step, one maximizes each term by setting the derivative equal to zero
2 and by using the constraint that $\sum_{i,j=0,1} a_k^{ij} = 1$ to obtain the following MLE ($i = 0, 1$
3 and $j = 0, 1$):

$$4 \quad \hat{\lambda}_k = \gamma_{0,k}^1, \quad \hat{a}_k^{ij} = \frac{\sum_{t=1, (T-1)} \xi_{ij,k}^t}{\sum_{t=1, T} \gamma_{i,k}^t}, \quad \hat{\mu}_i^t = \frac{\sum_{k=1, N} \gamma_{i,k}^t y_k^t}{\sum_{k=1, N} \gamma_{i,k}^t}, \quad \hat{\sigma}^2 = \frac{\sum_{k=1, N} \sum_{t=1, T} \sum_{i=0, 1} \gamma_{i,k}^t (y_k^t - \mu_i^t)^2}{\sum_{k=1, N} \sum_{t=1, T} \sum_{i=0, 1} \gamma_{i,k}^t},$$

$$5$$

$$6 \quad \text{with } \gamma_{i,k}^t = E[z_k^t = i | y_k, \theta_k^{(p)}] = \frac{\alpha_{i,k}^t \beta_{i,k}^t}{\alpha_{0,k}^t \beta_{0,k}^t + \alpha_{1,k}^t \beta_{1,k}^t},$$

$$7 \quad \text{and } \xi_{ij,k}^t = E[z_k^t = i, z_k^{t+1} = j | y_k, \theta_k^{(p)}] = \frac{\alpha_{i,k}^t a_k^{ij} \beta_{j,k}^{t+1} \text{pr}(y_k^{t+1} | z_k^{t+1} = j)}{\alpha_{0,k}^t \beta_{0,k}^t + \alpha_{1,k}^t \beta_{1,k}^t}.$$

8 Note $\gamma_{0,k}^t$ is the individual posterior probability of an IMI- sample, given the
9 whole SCS sequence. Correspondingly, $\xi_{01,k}^t$ is the posterior probability, for the
10 k^{th} cow, that a hidden state sequence that had to generate the SCS sequence went
11 through IMI- at time t and transitioned into IMI+ at time $t + 1$.

12
13 *Evaluation of the MLE.* The HMM described in the preceding sections were
14 used to analyze the SCS records. Missing SCS were restored through a multiple
15 imputation procedure with the MCMC method (proc MI of SAS[®]) in an attempt to
16 avoid loss of statistical power and selection bias associated with loss to follow-up,
17 and to be able to use standard matrix algebra. In this method (refer to Horton and
18 Kleinman, 2007 for a thorough discussion), each missing value is replaced by a

1 set of plausible values that represent the uncertainty about the right value to
 2 impute. After imputation, the set of imputed values were averaged for the
 3 subsequent analyses. Note the SCC were transformed in SCS to ensure normality
 4 which is an assumption of the MCMC method for imputing missing data.

5 Different priors for μ_0 (2 to 5), μ_1 (4 to 8) and σ^2 (1 or 2) were used to start the
 6 EM algorithms. After the MLE of the parameters were computed, the estimated
 7 number of transitions between successive MIM were obtained as

$$8 \quad \hat{n}_{ij,k} = \sum_{t=1,T} \xi_{ij,k}^t$$

9 with $i = j = 0$ if the transition is from IMI- to IMI-, $i = 0$ and $j = 1$ if the transition
 10 is from IMI- to IMI+, $i = 1$ and $j = 0$ if the transition is from IMI+ to IMI-, and $i =$
 11 1 and $j = 1$ if the transition is from IMI+ to IMI+. These numbers were compared
 12 to the observed numbers of transitions between successive MIM with SCM- or
 13 SCM+, and to the observed numbers of transitions between MIM with or without
 14 a clinical case (CM+ and CM-). The comparisons were done for lactations with
 15 or without reported clinical case associated to a positive bacteriological result. In
 16 lactations without reported clinical case, only transitions from CM- to CM- were
 17 achievable. The numbers of transitions from CM+ to CM+ were not computed
 18 because only the first clinical cases were considered.

19

20 **Results**

21

22 Before imputation, 6.01% of the 128,748 monthly records were bacteriologically
 23 positive. The average SCS over all lactations (first parity cows with or without

1 case) was at 2.65 (SD = 1.62) in the first MIM, decreased to a minimum at 2.08
 2 (SD = 1.45) during the second MIM before increasing slowly to 2.70 (SD = 1.37)
 3 at the end of the lactation. A similar pattern was found for lactations without any
 4 case of mastitis (Figure 1), but here, SCS was slightly lower throughout the
 5 lactation. In lactations with mastitis, cases were detected on average on the 128th
 6 DIM and 27.6% of those occurred during the first MIM. The percentage
 7 decreased thereafter, from 12.6% in the second MIM to 5.5% in the last 3 MIM.
 8 The complete sequence (n = 10) of SCS records was available on 10.48% of the
 9 cows and 75.56% of lactations had information on 5 MIM or more. The missing
 10 pattern was not monotone, i.e., a missing SCS was not necessarily followed by
 11 missing SCS. After imputation, the complete data set included 218,290 monthly
 12 SCS records of which 3.23% were considered as CM+. The SCS trend was
 13 similar before and after imputation and the highest difference (about 0.08) was
 14 found when SCS were the smallest, in the second MIM. After imputation, the
 15 SCS averaged 2.64 (SD = 1.37) in the first MIM, decreased to a minimum at 2.00
 16 (SD = 1.27) during the second MIM and increased to reach 2.69 (SD = 1.23) at
 17 the end of the lactation (Figure 1).

18 Figure 1 is about here

19 The estimated means and variance obtained with the HMM were: $\hat{\mu}_0 = 1.96$ (SD
 20 = 0.16), $\hat{\mu}_1 = 4.73$ (SD = 0.71), and $\hat{\sigma}^2 = 0.83$ (SD = 0.11). As comparison, the
 21 observed SCS means for CM- and CM+ lactations were 2.35 (SD = 0.99) and 3.18
 22 (SD = 1.28), respectively. For SCM- and SCM+ lactations, the observed means
 23 were 1.97 (SD = 0.64) and 3.48 (SD = 1.01), respectively.

1 The average number of transitions between hidden (IMI- and IMI+) and
2 observed states (SCM+ and SCM-, CM+ and CM-) states are shown in Figure 2
3 for lactations with or without at least one reported clinical case. The null
4 hypothesis of no differences between these numbers was tested by a t student test
5 ($p < 0.01$): The numbers of transitions from IMI- to IMI- and from SCM- to
6 SCM- were lower than the observed number of transitions from CM- to CM-. For
7 example, when no cases were reported during the entire lactation (Figure 2a),
8 there were 9 transitions from CM- to CM- but the number of transitions from
9 SCM- to SCM- was 8.14 (SD = 2.0) and the number of transitions from IMI- to
10 IMI- was 7.27 (SD = 2.7).

11 Figures 2 is about here

12 The average probabilities of transition between hidden states are given in Figure 3
13 for lactations with or without a reported case of clinical mastitis. Overall, the
14 probability to remain uninfected was $\hat{a}^{00} = 78.84\%$, to recover from infection was
15 $\hat{a}^{10} = 60.49\%$, to get newly infected was $\hat{a}^{01} = 11.70\%$ and to remain infected was
16 $\hat{a}^{11} = 15\%$. No significant differences were found in these probabilities between
17 lactations with or without at least one reported clinical case.

18 Figure 3 is about here

19

20 Discussion

21

22 A naïve HMM is proposed to analyze sequences of monthly SCS as they are
23 collected by the milk recording agencies with the intention of identifying cows

1 with or without mastitis. The data were previously analyzed by de Haas *et al.*
2 (2004) to identify pathogen-specific SCC patterns. The SCS patterns in Figure 1
3 are similar to patterns from the previous study (Figure 1a in Haas *et al.*, 2004),
4 with slight differences mainly due to different editing procedures and considering
5 that SCS were averaged over each MIM.

6 The model provides useful features for genetic and genomic selections. Firstly,
7 results from Figure 2 suggested that analyzing SCS with a HMM lead to the
8 identification of infected cows before the apparition of clinical signs and before
9 SCC gets higher than 250,000 cells/mL. Indeed, among cows for which at least a
10 mastitis case was reported (Figure 2b), the model assigned the state IMI+ on three
11 occurrences while the stage SCM+ was observed on two occasions. In heifers
12 without any reported clinical mastitis case (Figure 2a), there were two IMI+ and
13 one SCM+. The likely sequences for the IMI, SCM and CM stages are shown in
14 Figure 4, considering that most clinical stages were reported in early lactation. It
15 should however be noted that an experimental infection in a well-designed clinical
16 trial is necessary to confirm this findings

17 Figure 4 is about here

18 Although these results should be confirmed in a well-designed clinical trial with
19 experimental infection, this ability will lead to more accurate estimates of
20 breeding values and an earlier and more accurate selection. It will also facilitate
21 the identification of the genetic determinants of mastitis because hidden IMI states
22 may be considered as intermediate phenotypes with stronger genetic determinants
23 than SCM or CM. Secondly, HMM may be used to predict the future health

1 status of a cow, based on its previous sequence of SCS. Mathematically, this
 2 prediction is given by combining forward (α) and backward (β) probabilities used
 3 in the algorithm as:

$$4 \quad \text{pr}(z_k^t = i \mid y_k^1, y_k^2, \dots, y_k^{t-1}, \theta_k) = \frac{\alpha_{0,k}^{t-1} \hat{a}_k^{0i} + \alpha_{1,k}^{t-1} \hat{a}_k^{1i}}{\alpha_{0,k}^{t-1} + \alpha_{1,k}^{t-1}}$$

5 In these probabilities, the uncertainty about the time of exposure to infection, if it
 6 has occurred, is reduced because data on the entire available sequence of SCS is
 7 exploited. Therefore, it may lower the biases due to incomplete exposure on
 8 estimable heritabilities (Bishop and Wooliams, 2010). Thirdly, the model
 9 provides estimates of the probability of recovery (IMI+ to IMI- = a^{10}) and of new
 10 infection (IMI- to IMI+ = a^{01}) for each animal. These parameters are directly
 11 related to well-established selection objectives for better udder health and
 12 epidemiological concepts. For example, the force of infection (ω = the rate at
 13 which susceptible individuals become infected) and the recovery rate (δ = the rate
 14 at which infected individuals recover) may be obtained from a^{10} and a^{01} as:

$$15 \quad a^{01} = \frac{\omega}{\omega + \delta} (1 - e^{-(\omega + \delta)t}) \text{ and } a^{10} = \frac{\delta}{\omega + \delta} (1 - e^{-(\omega + \delta)t}),$$

16 assuming a SI model (Anderson and May, 1992; Detilleux *et al.*, 2006). Then,
 17 data from genetic and epidemiological studies could be combined to analyze the
 18 impact of selecting for a better ability to recover from disease on the spread of the
 19 disease at the population level. Finally, the model can be extended by adding
 20 genetic random effects to obtain breeding values for SCS (Detilleux, 2008) or
 21 even for the hidden IMI variable (Altman, 2007), considering the total genetic

1 effects on SCS is a combination of the effects of genes responsible for presence or
2 not of infection and for the magnitude of the SCS response after infection.

3 The model is very flexible and allows the inclusion of prior knowledge (e.g.,
4 clinical or laboratory records) to the SCS information. The effects of covariates
5 (e.g., treatment or culling, breed, parity) on the progression of the IMI could also
6 be studied by comparing transition rates.

7 The HMM methodology presents also some limitations. The HMM, as proposed
8 here, necessitated that the sequence of SCS was complete. One possibility was to
9 discard lactations with incomplete information but this would have decreased the
10 amount of available data. Missing data were instead imputed and a multiple
11 imputation procedure was chosen as it increases robustness to departures from the
12 true imputation model considerably compared to single imputation approaches
13 that do not reflect uncertainty about the imputed values. The MCMC method was
14 chosen because SCS were distributed normally and because the missing pattern
15 was not monotone. After imputation, the SCS curves were slightly lower than
16 before imputation (Figure 1). This may be explained by the fact that, in the
17 MCMC method, missing SCS were replaced by randomly selecting a value (at
18 any MIM) and that SCS at different MIM are correlated with the SCS being
19 imputed. Another drawback was the assumption that probability of staying in a
20 given state was independent of the duration of the state. It could have been
21 modeled explicitly as $a_{11}^{d-1} (1 - a_{11})$ which is the probability of staying d times in
22 state IMI+. The transition probabilities were assumed constant across time
23 although it is known that susceptibility to IMI vary across lactation stages (Paape

1 *et al.*, 2002). This stationary assumption is very strong but it could be relaxed by
2 parameterizing the mean of the IMI+ distribution to account for various trend or
3 seasonality in the data (Le Strat and Carrat, 1999). Another assumption of the
4 HMM, the independence between successive SCS, could be released in its
5 autoregressive form by allowing previous SCS to assist in predicting the current
6 SCS (Lavery *et al.*, 2002; Ephraim and Roberts, 2005). Finally, the assumption
7 of homoscedasticity can be relaxed by modeling different variances for the IMI+
8 and IMI- samples (Dettelleux, 2008)

9 The maximum likelihood estimation via the EM algorithm has also some
10 disadvantages. For example, it does not provide an estimated covariance matrix
11 for the parameters. Bootstrap methods can be used but they are computationally
12 intensive for this type of model. Other alternatives are to estimate parameters via
13 the Gibbs sampler or Bayesian variational methods (Jaakkola and Jordan, 2000).
14 Collinearity between parameter estimates can lead to identifiability problems
15 (Brookhart *et al.*, 2002) and the EM may converge toward singular estimates at
16 the boundary of the parameter space. It may also fail to converge. The problem
17 becomes particularly severe when time series are short and data sparse (Cooper
18 and Lipstich, 2004).

19

20 **Conclusions**

21

22 A simple hidden Markov model (HMM) was applied on SCS recorded monthly on
23 cows with or without clinical mastitis to evaluate its accuracy in estimating

1 parameters under health or disease states. The SCS means were estimated at 1.96
2 (SD = 0.16) and 4.73 (SD = 0.71) for the hidden healthy and infected states, and
3 the common variance at 0.83 (SD = 0.11). The probabilities to remain uninfected,
4 to recover from infection, to get newly infected and to remain infected between
5 consecutive test-days were estimated at 78.84%, 60.49%, 11.70% and 15%,
6 respectively. Three different health related states were compared: clinical stages
7 observed by farmers, subclinical cases defined for somatic cell counts below or
8 above 250,000 cells/mL and infected stages obtained from the HMM. The results
9 showed that HMM identifies infected cows before the apparition of clinical and
10 subclinical signs which may critically improve the power of studies on the genetic
11 determinants of SCS and reduce biases in predicting breeding values for SCS.
12 The HMM provides also epidemiological parameters that describe the spread of
13 mastitis at the population level.

14

15 **Acknowledgements**

16

17 My special thanks go to CRV BV and Prof. H. Barkema for providing the
18 data and to EADGENE (European Animal Disease Genomics Network of
19 Excellence for Animal Health and Food Safety) for financial support.

20

21 **R eferences**

22

1 Anderson RM and May RM, Infectious diseases of Humans. Oxford Science
2 Publications, Oxford, 1992.
3

4 Altman RM 2007. Mixed hidden Markov model: An extension of the hidden
5 Markov model to the longitudinal data setting. Journal Animal Science
6 Association 102: 201-210.
7

8 Barkema, HW, Schukken YH, Lam TJGM, Beiboer ML, Wilmink H, Benedictus
9 G, and Brand A 1998. Incidence of clinical mastitis in dairy herds grouped in
10 three categories by bulk milk somatic cell counts. Journal Dairy Science 81:411-
11 419.
12

13 Bilmes JA 1998. A gentle tutorial of the EM algorithm and its application to
14 parameter estimation for Gaussian mixture and Hidden Markov models.
15 Technical report, University of Berkeley.
16

17 Bishop SC, and Wooliams JA 2010. On the genetic interpretation of disease data.
18 PLoS ONE 5(1): e8940. doi:10.1371/journal.pone.0008940
19

20 Boettcher PJ, Moroni P, Pisoni G, and Gianola D 2005. Application of finite
21 mixture model to somatic cell scores of Italian goats. Journal Dairy Science 88,
22 2209-2216.
23

1 Brookhart MA, Hubbard AE, van der Laan MJ, Colford JM Jr, Eisenberg JN
2 2002. Statistical estimation of parameters in a disease transmission model:
3 analysis of a *Cryptosporidium* outbreak. *Statistics in Medicine* 21, 3627–3638.
4

5 Buyske S, Yang G, Matise TC, Gordon D 2009. When a case is not a case:
6 Effects of phenotype misclassification on power and sample size requirements for
7 the transmission disequilibrium test with affected child trios. *Human Heredity* 67,
8 287-292.
9

10 Cooper B and Lipstich M 2004. The analysis of hospital infection data using
11 hidden Markov models. *Biostatistics* 5, 223–237.
12

13 de Haas Y, Barkema HW, and Veerkamp RF 2002. The effect of pathogen-
14 specific clinical mastitis on the lactation curve for somatic cell count. *Journal*
15 *Dairy Science* 85, 1314-1323.
16

17 de Haas Y, Veerkamp RF, Barkema HW, Gröhn YT, and Schukken YH 2004.
18 Associations between pathogen-specific cases of clinical mastitis and somatic cell
19 count patterns. *Journal Dairy Science* 87, 95–105.
20

21 Detilleux JC, and Leroy P 2000. Application of a mixed normal mixture model
22 for the estimation of mastitis-related parameters. *Journal Dairy Science* 83, 2341–
23 2349.

1
2 Detilleux JC, Vangroenweghe F, and Burvenich C 2006. Mathematical model
3 of the acute inflammatory response to *Escherichia coli* intramammary challenge.
4 Journal Dairy Science 89,3455-65.
5
6 Detilleux JC 2008. The analysis of disease biomarker data using a mixed hidden
7 Markov model. Genetics Selection Evolution 40, 491-509.
8
9 Djabri B, Bareille N, Beaudeau F, and Seegers H 2002. Quarter milk somatic
10 cell count in infected dairy cows: a meta-analysis. Veterinary Research 33, 335-
11 357.
12
13 Ephraim Y, and Roberts W 2005. Revisiting autoregressive hidden Markov
14 modeling of speech signals. IEEE Signal Processing Letters 12, 166-169.
15
16 Eisner J 2002. An interactive spreadsheet for teaching the forward-backward
17 algorithm. Conference at the ACL-02 Workshop on effective tools and
18 methodologies for teaching natural language processing and computational
19 linguistics, Philadelphia, Pennsylvania, pp 10-18.
20
21 Gianola D 2005. Prediction of random effects in finite mixture models with
22 Gaussian components. Journal of Animal Breeding and Genetics 122, 145-159.
23

1 Godden SM, Jansen JT, Leslie KE, Smart NL, and Kelton DF 2002. The effect
2 of sampling time and sample handling on the detection of *Staphylococcus aureus*
3 in milk from quarters with subclinical mastitis. Canadian Veterinary Journal 43,
4 38-42.
5
6 Horton NJ and Kleinman KP 2007. Much ado about nothing: A comparison of
7 missing data methods and software to fit incomplete data regression models. The
8 American Statistician, 61, 79-90.
9
10 Jaakkola T, and Jordan MI 2000. Bayesian parameter estimation via variational
11 methods. Statistics and Computing, 10: 25-37.
12
13 Lam T, van Wuijckhuise LA, Franken P, Morselt ML, Hartman EG, and
14 Schukken YH 1996. Use of composite milk samples for diagnosis of
15 *Staphylococcus aureus* mastitis in dairy cattle. Journal American Veterinary
16 Medical Association 208, 1705–1708.
17
18 Lavery WH, Milet MJ, and Kelly IW 2002. Simulation of hidden Markov
19 models with EXCEL. The statistician 51, 31-40.
20
21 Le Strat Y, and Carrat F 1999. Monitoring epidemiologic surveillance data
22 using hidden Markov models. Statistics in Medicine 18, 3463-3478.
23

1 Moroni P, Pisoni G, Vimercati C, Rinaldi M, Castiglioni B, Cremonesi P, and
2 Boettcher P 2005. Characterization of *Staphylococcus aureus* isolated from
3 chronically infected dairy goats. Journal Dairy Science 88, 3500 – 3509.
4
5 Paape M, Mehrzad J, Zhao X, Detilleux J, and Burvenich C 2002. Defense of the
6 bovine mammary gland by polymorphonuclear neutrophil leukocytes. Journal
7 Mammary Gland Biology Neoplasia 7, 109–121.
8
9 Rabiner LR 1989. A tutorial on hidden Markov models and selected
10 applications in speech recognition. Proc. IEEE 77:257-268, 1989.
11
12 Sargeant JM, Leslie KE, Shirley JE, Pulkrabek BJ, and Lim GH 2001.
13 Sensitivity and specificity of somatic cell count and California mastitis test for
14 identifying intramammary infection in early lactation. Journal Dairy Science 84,
15 2018-2024.
16
17 SAS OnlineDoc™:Version 8. Statistical Analysis System.
18
19 Sears PM, Smith BS, English PB, Herer PS, and Gonzales RN 1990. Shedding
20 pattern of *Staphylococcus aureus* from bovine intramammary infections. Journal
21 Dairy Science 73, 2785–2789.
22

- 1 Shook GE and Schutz MM 1994. Selection on somatic cell score to improve
2 resistance to mastitis in the United States. *Journal of Dairy Science* 77, 648-658.
3
- 4 Suriyasathaporn W, Schukken YT, Nielen M, and Brand A. 2000. Low somatic
5 cell count: a risk factor for subsequent clinical mastitis in a dairy herd. *Journal*
6 *Dairy Science* 83, 1248-1255.
7

1 Figure 1 Monthly average somatic cell scores for all lactations (square) and
2 for lactations without clinical mastitis (cross), before (straight line) and after
3 (broken lines) imputation.
4

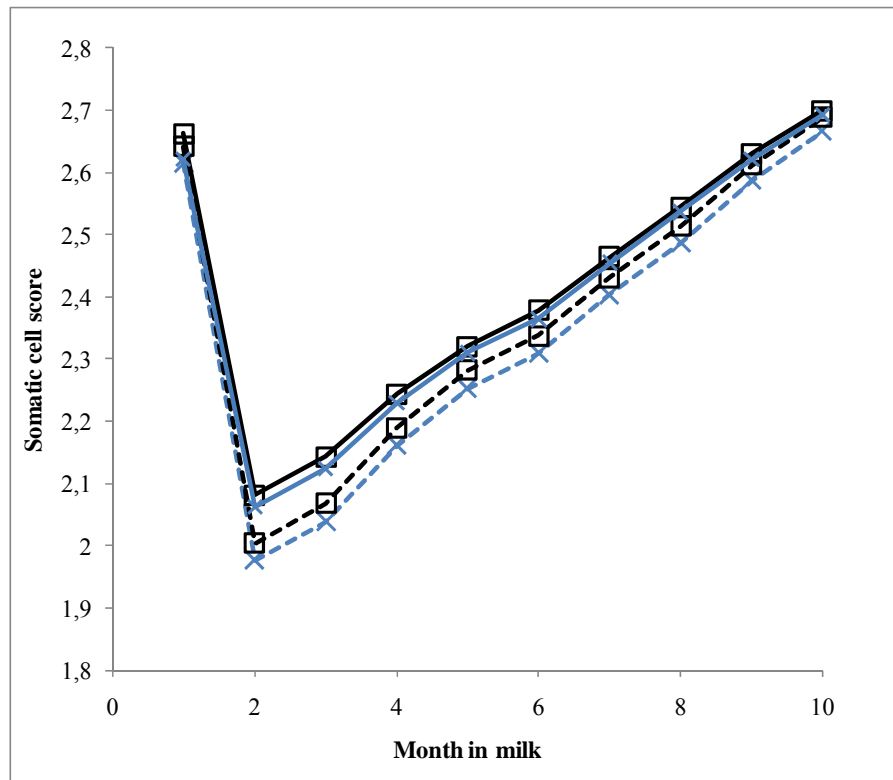
5 Figure 2 Mean number of transitions from one mastitis state to another
6 across the 10 test-days for cows without (Figure 2a) and with (Figure 2b) at least
7 one reported clinical case of mastitis. States are designed as CM+ or CM- (plain
8 bar) when a clinical case is reported or not, as SCM+ or SCM- (spotted bar) when
9 SCC are above or below 250,000 cells/ml, and IMI+ or IMI- (stripped bar) when
10 the sample is classified as infected or not by the model, respectively.
11

12 Figure 3 Average probabilities of transition between IMI states for lactations
13 with (stripped bar) and without (plain bar) a reported clinical case. The hidden
14 states are IMI+ and IMI- when the sample is classified as infected or not by the
15 model, respectively.
16

17 Figure 4 Examples of sequences for the IMI, SCM and CM stages across the
18 10 test-days based on the results shown in Figure 2, for cows without (Figure 4a)
19 and with (Figure 4b) at least one reported clinical case of mastitis. The sign is +
20 when the sample is positive for the stage at the test day. The sign is - when the
21 sample is negative for the stage at the test day.
22

1 Figure 1.

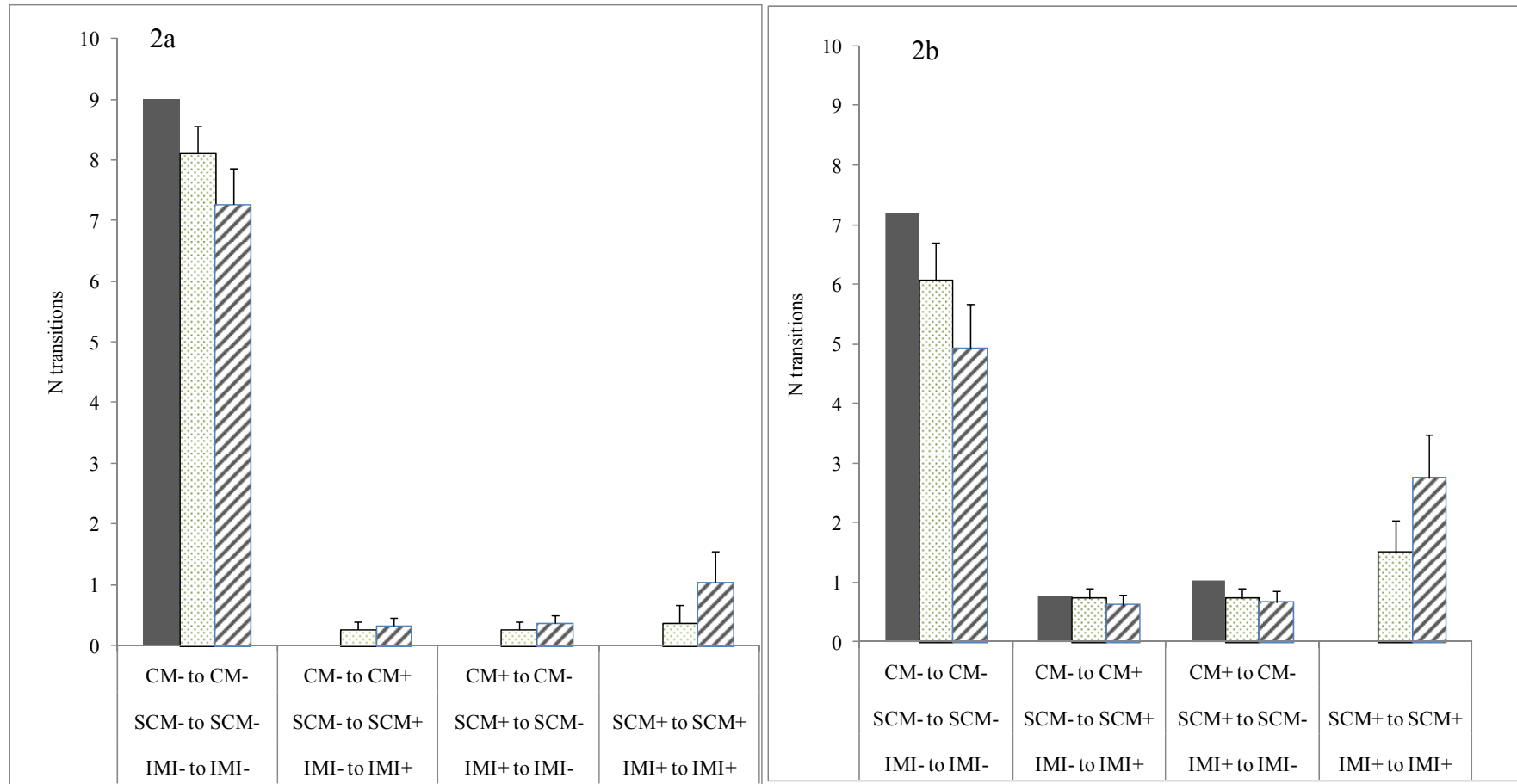
2



3

4

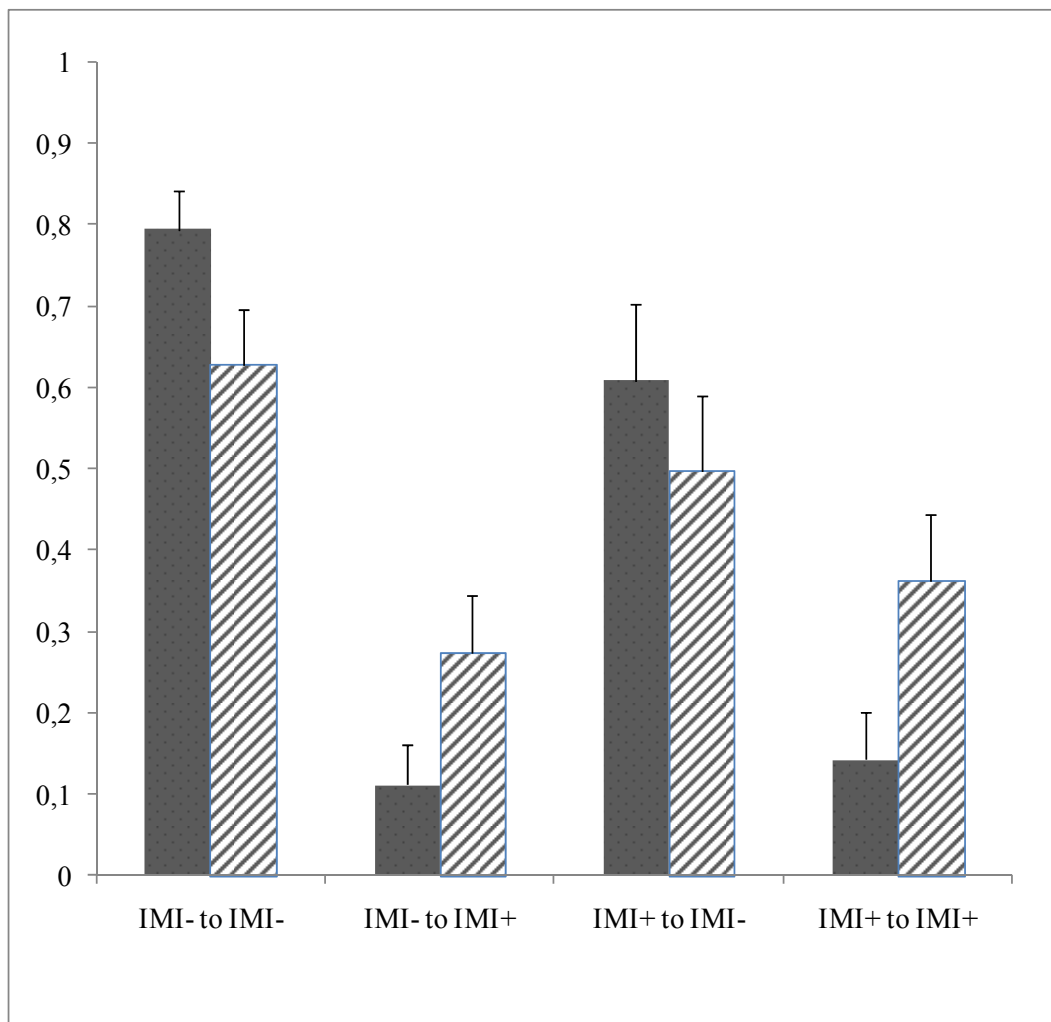
1 Figure 2.



2

1 Figure 3

2



3

4

1 Figure 4.

2

3

4

5

