

# Connecting the Latent Multinomial

Matthew R. Schofield<sup>1\*</sup>, Simon J. Bonner<sup>2</sup>

<sup>1</sup>Department of Mathematics and Statistics, University of Otago, P.O. Box 56 Dunedin 9054, New Zealand

<sup>2</sup>Department of Statistics, University of Kentucky, Lexington KY 40536, U.S.A.

\**email*: mschofield@maths.otago.ac.nz

SUMMARY: Link et al. (2010) define a general framework for analyzing capture-recapture data with potential misidentifications. In this framework, the observed vector of counts,  $\mathbf{y}$ , is considered as a linear function of a vector of latent counts,  $\mathbf{x}$ , such that  $\mathbf{y} = \mathbf{A}\mathbf{x}$ , with  $\mathbf{x}$  assumed to follow a multinomial distribution conditional on the model parameters,  $\boldsymbol{\theta}$ . Bayesian methods are then applied by sampling from the joint posterior distribution of both  $\mathbf{x}$  and  $\boldsymbol{\theta}$ . In particular, Link et al. (2010) propose a Metropolis-Hastings algorithm to sample from the full conditional distribution of  $\mathbf{x}$ , where new proposals are generated by sequentially adding elements from a basis of the null space (kernel) of  $\mathbf{A}$ . We consider this algorithm and show that using elements from a simple basis for the kernel of  $\mathbf{A}$  may not produce an irreducible Markov chain. Instead, we require a Markov basis, as defined by Diaconis and Sturmfels (1998). We illustrate the importance of Markov bases with three capture-recapture examples. We prove that a specific lattice basis is a Markov basis for a class of models including the original model considered by Link et al. (2010) and confirm that the specific basis used by Link et al. (2010) for their example with two sampling occasions is a Markov basis. The constructive nature of our proof provides an immediate method to obtain a Markov basis for any model in this class.

KEY WORDS: Capture-recapture; Linear constraint; Markov basis; Markov chain Monte Carlo; Misidentification.

## 1. Introduction

The development of capture-recapture methodology has a long history, allowing estimation of demographic parameters of interest for animal populations (see Amstrup et al. 2005, for a review). Similar methods have also been used to study human populations, including intravenous drug users (King et al. 2009) and human rights abuse victims (Lum et al. 2013). In general, a capture-recapture experiment consists of a series of capture occasions on which overlapping subsets of the population are observed. For animal populations the occasions are usually ordered in time while for human populations they may comprise lists obtained from different sources. It is assumed that each individual has a unique identifying mark that is either given or realized when the individual is first captured and this mark can be used to identify the individual on subsequent occasions. In this paper, we are concerned with fitting capture-recapture models to data that provide an incomplete or inaccurate representation of the true encounters of individuals during the experiment. This may occur if the data consist of incomplete summary statistics or if individuals are misidentified on some occasions. Examples of capture-recapture studies that are prone to identification errors include (i) multi-list studies in which individuals may be matched based on personal information such as name, birth date, medical record number (Seber et al. 2000, Lee et al. 2001, Sutherland and Schwarz 2005, Fienberg and Manrique-Vallier 2009), (ii) animal studies in which individual identity is found from non-invasive sampling, e.g. genetic information from scat or hair (Wright et al. 2009, Link et al. 2010, Yoshizaki et al. 2011) or photographic ID of individuals (Yoshizaki et al. 2009, Bonner and Holmberg 2013, McClintock et al. 2013), and (iii) studies in which (at least) two sources of capture-recapture information are available for the same population with little to no information about how the individual IDs in one source corresponds to individual ID from the other sources (Bonner and Holmberg 2013, McClintock et al. 2013). Our focus is on the algorithm for a general class of mark-recapture models allowing for

28 misidentification considered by Link et al. (2010) (hereafter L2010). This class is described by  
 29 the latent multinomial model, in which an observed data vector,  $\mathbf{y}$  can be expressed as a linear  
 30 function of a latent data vector,  $\mathbf{x}$ , modeled by a multinomial distribution with unknown  
 31 parameters  $\boldsymbol{\theta}$ , denoted  $[\mathbf{x}|\boldsymbol{\theta}]$ . The notation  $[x]$  denotes the probability density function  $f_X(x)$   
 32 for a continuous random variable  $X$  or the probability mass function  $\Pr(X = x)$  for a discrete  
 33 random variable  $X$ . The linear function is expressed as

$$34 \quad \mathbf{y} = \mathbf{A}\mathbf{x}, \quad (1)$$

35 where  $\mathbf{A}$  is called the configuration matrix (a matrix of known constants that depends on  
 36 the specific problem) with more columns than rows. We continue to call this modeling setup  
 37 the latent multinomial model, even though the setup is flexible and can accommodate other  
 38 probability mass functions  $[\mathbf{x}|\boldsymbol{\theta}]$ , such as the Poisson model considered by Lee (2002).

39 The goal is to sample from the joint posterior distribution  $[\boldsymbol{\theta}, \mathbf{x}|\mathbf{y}]$  using Markov chain  
 40 Monte Carlo (MCMC) by alternating between sampling from the full conditional distribu-  
 41 tions  $[\boldsymbol{\theta}|\mathbf{x}, \mathbf{y}]$  and  $[\mathbf{x}|\mathbf{y}, \boldsymbol{\theta}]$ . The difficulty with this approach is in specifying an updating  
 42 scheme for  $\mathbf{x}$ . That is, how to efficiently sample from  $[\mathbf{x}|\mathbf{y}, \boldsymbol{\theta}]$  in such a way so that every  
 43  $\mathbf{x}$  vector that satisfies (1) has a positive probability of being reached at some point during  
 44 the updating. We consider three examples demonstrating that the scheme for updating  $\mathbf{x}$   
 45 proposed by L2010 may not produce an irreducible Markov chain for models within the latent  
 46 multinomial framework. We then present theory identifying a class of models for which the  
 47 specific algorithm does produce irreducible Markov chains, and show more generally how  
 48 these methods fit within the framework of algebraic statistics. This allows us to develop an  
 49 extension of the algorithm which can be used to generate valid MCMC samplers for the  
 50 posterior distributions from a broader class of latent multinomial models.

51 The MCMC algorithm we consider throughout this manuscript is presented in Figure 1.  
 52 Starting with an initial state  $\mathbf{x}^0$  satisfying the linear constraint, a proposal is generated on

53 the first iteration by adding or subtracting an element chosen randomly from a subset of  
 54 the kernel (or null space) of  $\mathbf{A}$ ,  $\mathcal{B} = \{\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_m\} \subset \ker(\mathbf{A})$ , with cardinality  $m$ . The  
 55 proposal is then accepted or rejected with probability determined by the Hasting's ratio,  $r$ ,  
 56 and the algorithm continues to the second iteration. This algorithm is a modification of that  
 57 presented by L2010, with three differences: (i) L2010 steps through all  $m$  elements in  $\mathcal{B}$  in  
 58 order instead of selecting an element at random on each iteration, (ii) when stepping through  
 59 every element in  $\mathcal{B}$ , L2010 multiplies element  $\mathbf{a}_i$  by a coefficient  $c \in \{-C_i, \dots, -1, 1, \dots, C_i\}$   
 60 in order to improve convergence, and (iii) L2010 assumes that  $\mathcal{B}$  is a basis for  $\ker(\mathbf{A})$ , while  
 61 we allow  $\mathcal{B}$  to be a more general subset that spans  $\ker(\mathbf{A})$ . The first two differences may  
 62 impact the efficiency of the algorithm but do not change the stationary distribution of the  
 63 resulting Markov chains, and we do not consider these differences further. Our focus is on  
 64 the third difference and the effect that the set  $\mathcal{B}$  can have on the generated Markov chains  
 65 and their stationary distributions.

66 [Figure 1 about here.]

67 To illustrate the problems that may occur if  $\mathcal{B}$  is poorly specified we consider three examples  
 68 of models which fit into the latent multinomial framework. First we consider the same  
 69 closed population mark-recapture model with misidentification considered by L2010. This  
 70 model, called  $M_{t\alpha}$ , assumes that captures occur according to a closed population model with  
 71 time dependent capture probabilities and that errors in identifying an individual are unique  
 72 and create ghost histories with single captures. Second, we consider a multi-list modeling  
 73 problem in which summary statistics are presented in place of the full data set, possibly  
 74 for privacy reasons. Our aim is to sample from possible complete data sets with the given  
 75 sufficient statistics. Finally, we consider a more complicated model of misidentification in  
 76 mark-recapture which allows for one marked individual to be identified as another previously  
 77 marked individual. Full details of these models and the issues regarding the selection of the

78 set  $\mathcal{B}$  to be used in the algorithm in Figure 1 are provided in sections 3, 4, and 5. As  
79 motivation, we consider the output from Markov chains constructed using the algorithm in  
80 Figure 1 for each of the three examples. For each example, we defined  $\mathcal{B}$  to be a basis for  
81  $\ker(\mathbf{A})$  as in L2010 and ran two parallel chains, each of which started from a different initial  
82 value. For both model  $M_{t\alpha}$  and the multi-list model with sufficient statistics, despite strong  
83 evidence that each chain has converged, it is clear that the two chains are not sampling from  
84 the same distribution for a given quantity of interest (Figure 2). This is even more apparent  
85 in the third example where one of the two chains never moves from its initial value.

86 [Figure 2 about here.]

87 The problem in all three examples is that the stationary distribution reached by the Markov  
88 chains produced by the algorithm in Figure 1 may depend on the chosen set,  $\mathcal{B}$  and the initial  
89 value of  $\mathbf{x}$ . Although the values of  $\mathbf{x}$  proposed on each iteration are guaranteed to satisfy the  
90 linear constraint the resulting Markov chains may not reach all points in the sample space  
91 and the stationary distributions may be dependent on the initial values. In the next section  
92 we provide a basic introduction to the field of algebraic statistics and the results of Diaconis  
93 and Sturmfels (1998) and others who have explored approaches for sampling from  $\mathbf{x}$  from  
94 a linear constraint as in (1) in other application areas. We then consider the implications  
95 of this theory to show why the MCMC algorithms failed above (Figure 2), and how valid  
96 MCMC samplers can be constructed for each of the three examples.

## 97 2. Introduction to algebraic statistics

98 Sampling  $\mathbf{x}$  in the presence of the linear constraint in (1) is not unique to capture-recapture  
99 problems. In a seminal paper in algebraic statistics, Diaconis and Sturmfels (1998) considered  
100 a linear constraint of the same form when developing conditional goodness-of-fit tests for  
101 contingency tables (see Karwa and Slavkovic 2013, for a recent review). That is, they

102 considered how to construct an MCMC algorithm to sample different contingency tables  
 103 with common (fixed) row and column sums (such ideas can also be extended to multi-way  
 104 contingency tables).

105 To consider the problem at hand in more detail we will summarize several definitions and  
 106 results from linear algebra in this section (basic definitions regarding kernels and bases are  
 107 provided in the supplementary materials). We will use a  $3 \times 3$  contingency table example to  
 108 illustrate many of the ideas. The table is

$$\begin{array}{cccc|c}
 & x_{11} & x_{12} & x_{13} & x_{1\cdot} \\
 & x_{21} & x_{22} & x_{23} & x_{2\cdot} \\
 & x_{31} & x_{32} & x_{33} & x_{3\cdot} \\
 \hline
 & x_{\cdot 1} & x_{\cdot 2} & x_{\cdot 3} & 
 \end{array}$$

110 where  $x_{ij}$  is the value in the  $i$ th row and  $j$ th column,  $x_{\cdot j}$  refers to the sum of the  $j$ th column  
 111 and  $x_{i\cdot}$  refers to the sum of the  $i$ th row. The column and row sums are vectorized to give  
 112 the vector of summary statistics

$$\mathbf{y} = (x_{\cdot 1}, x_{\cdot 2}, x_{\cdot 3}, x_{1\cdot}, x_{2\cdot})'.$$

114 Note that we need not include the third row sum as this is a derived quantity of the other  
 115 elements of  $\mathbf{y}$ . The individual entries in the table are vectorized to give

$$\mathbf{x} = (x_{11}, x_{21}, x_{31}, x_{12}, x_{22}, x_{32}, x_{13}, x_{23}, x_{33})'.$$

117 The specification is completed with

$$\mathbf{A} = \begin{pmatrix}
 1 & 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\
 0 & 0 & 0 & 1 & 1 & 1 & 0 & 0 & 0 \\
 0 & 0 & 0 & 0 & 0 & 0 & 1 & 1 & 1 \\
 1 & 0 & 0 & 1 & 0 & 0 & 1 & 0 & 0 \\
 0 & 1 & 0 & 0 & 1 & 0 & 0 & 1 & 0
 \end{pmatrix}$$

119 so that the constraints inherent in a contingency table follow (1). If we have column/row

120 sums given by

$$121 \quad \mathbf{y} = (5, 3, 2, 0, 4)'$$

122 then two contingency tables compatible with these constraints have entries

$$123 \quad \mathbf{x}_1 = (0, 2, 3, 0, 1, 2, 0, 1, 1) \quad \text{and} \quad \mathbf{x}_2 = (0, 3, 2, 0, 0, 3, 0, 1, 1). \quad (2)$$

124 Our goal is to specify an MCMC algorithm that samples from the set of vectors  $\mathbf{x}$  that  
125 satisfy (1) for a particular  $\mathbf{y}$ . This is defined as the  $\mathbf{y}$ -fiber (or simply fiber)  $\mathcal{F}_{\mathbf{y}}$ ,

$$126 \quad \mathcal{F}_{\mathbf{y}} = \{\mathbf{x} \in \mathbb{N}^d : \mathbf{y} = \mathbf{A}\mathbf{x}\},$$

127 where  $d$  is the dimension of  $\mathbf{x}$  and  $\mathbb{N} = \{0, 1, \dots\}$ . L2010 refers to  $\mathcal{F}_{\mathbf{y}}$  as the feasible set.

128 To move between elements of the fiber, we make use of the lattice kernel  $\ker_{\mathbb{Z}}(\mathbf{A})$ . The  
129 lattice kernel is the integer valued subset of the kernel,

$$130 \quad \ker_{\mathbb{Z}}(\mathbf{A}) = \ker(\mathbf{A}) \cap \mathbb{Z}^d = \{\mathbf{x} \in \mathbb{Z}^d : \mathbf{A}\mathbf{x} = \mathbf{0}\}.$$

131 In algebraic statistics, a move is defined to be any element of the lattice kernel, such that  
132 the vector  $\mathbf{v}$  is a move if  $\mathbf{v} \in \ker_{\mathbb{Z}}(\mathbf{A})$ . An implication of this is that if  $\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{F}_{\mathbf{y}}$  then  
133  $\mathbf{x}_2 - \mathbf{x}_1$  is a move. The idea is that the elements of the lattice kernel can be added to a vector  
134 that satisfies the linear constraint and the result is guaranteed to still satisfy the constraint.  
135 However, it is not practical to consider all elements of the lattice kernel when updating  $\mathbf{x}$  as  
136  $\ker(\mathbf{A})$  is potentially very large and difficult to compute. Instead we want to find a smaller  
137 set of moves  $\mathcal{B} = \{\mathbf{v}_1, \dots, \mathbf{v}_m\} \subset \ker_{\mathbb{Z}}(\mathbf{A})$  that can be used to update  $\mathbf{x}$ . That is, we require  
138 a smaller set of moves so that it is possible to move between all elements of  $\mathcal{F}_{\mathbf{y}}$  using the  
139 algorithm in Figure 1.

140 The suggestion of L2010 was to use a basis for  $\ker(\mathbf{A})$  for this set of moves. However, we  
141 do not wish to construct a basis for  $\ker(\mathbf{A})$ , but instead a lattice basis for the integer lattice  
142  $\ker_{\mathbb{Z}}(\mathbf{A})$ . A lattice basis is a set of linearly independent vectors where every  $\mathbf{v} \in \ker_{\mathbb{Z}}(\mathbf{A})$   
143 can be found as a linear combination of the lattice basis vectors using integer coefficients.

144 If we insist on using a basis for  $\ker(\mathbf{A})$ , it may not be possible to reach all solutions using  
 145 only integer values of the coefficients,  $c$ , as specified in the algorithm in Figure 1. However,  
 146 even if we choose to use a lattice basis for  $\mathcal{B}$  it may be necessary to pass through one (or  
 147 more) vectors containing negative elements when applying moves one at a time to transition  
 148 between elements in the fiber  $\mathcal{F}_y$ . As vectors  $\mathbf{x}$  containing negative elements can never be  
 149 accepted, the use of a lattice basis for  $\mathcal{B}$  may result in sampling from a subset of the fiber  
 150  $\mathcal{F}_y$  when using the algorithm in Figure 1. This explains the observed results in the three  
 151 examples shown in Section 1: the two chains are exploring different subsets of the fiber.

152 These ideas are formalized using the concept of connectivity. Elements  $\mathbf{x}_j, \mathbf{x}_k \in \mathcal{F}_y$  are  
 153 connected using the set  $\mathbf{V} = (\mathbf{v}_1, \dots, \mathbf{v}_m)$  if there are moves  $\mathbf{v}_i \in \mathbf{V}$ ,  $i \in \{1, \dots, M\}$  so  
 154 that we can start from  $\mathbf{x}_j$  and add or subtract these moves one at a time to reach  $\mathbf{x}_k$   
 155 without any element in any of the partial sums ever being negative (note that the elements  
 156  $\mathbf{v}_i$ ,  $i = 1, \dots, M$  need not be distinct and some elements may be repeated multiple times).  
 157 That is, there exist  $\epsilon_1, \dots, \epsilon_M \in \{-1, 1\}$  such that

$$158 \quad \mathbf{x}_k = \mathbf{x}_j + \sum_{j=1}^M \epsilon_j \mathbf{v}_j \quad \text{and} \quad \mathbf{x}_1 + \sum_{k=1}^L \epsilon_k \mathbf{v}_k \in \mathcal{F}_y, \quad L = 1, \dots, M - 1.$$

159 We then say that the fiber  $\mathcal{F}_y$  is connected by  $\mathbf{V}$  if every pair of elements in the fiber are  
 160 connected.

161 We can apply the algorithm in Figure 1 to the  $3 \times 3$  contingency table example using the  
 162 elements of a lattice basis. A lattice basis can be found using the Hermite normal form (Aoki  
 163 et al. 2012, pg. 53). Unless otherwise stated, all lattice bases provided in this manuscript  
 164 are found using this approach. We note that the lattice basis obtained is not unique and a  
 165 different basis is often found if one reorders the columns of  $\mathbf{A}$  (and corresponding entries of



166  $\mathbf{x}$ ). For the contingency table, a lattice basis is given by elements LB1 – LB4 in (3)

	$x_{11}$	$x_{21}$	$x_{31}$	$x_{12}$	$x_{22}$	$x_{32}$	$x_{13}$	$x_{23}$	$x_{33}$
LB1	1	-1	0	-1	1	0	0	0	0
LB2	-1	0	1	1	0	-1	0	0	0
LB3	1	-1	0	0	0	0	-1	1	0
LB4	0	0	0	1	0	-1	-1	0	1

(3)

168 If we attempt to apply any of the elements LB1 — LB4 to either  $\mathbf{x}_1$  or  $\mathbf{x}_2$  in (2) we  
 169 immediately find a problem. Either adding or subtracting any of LB1 – LB4 results in  
 170 at least one negative count in the proposal and will lead to it being automatically rejected.  
 171 That means there is no way to use the elements LB1 – LB4 as moves in the algorithm in  
 172 Figure 1 and successfully transition between the two solutions in (2). In fact, we are unable  
 173 to move between any two valid solutions. As a result, the lattice basis in (3) does not connect  
 174 the fiber for this example. One solution is to change the algorithm in Figure 1 to use elements  
 175 of a lattice basis in a linear combination instead of one-at-a-time. While attractively simple,  
 176 Diaconis and Sturmfels (1998) implemented this for several examples and found that it was  
 177 inefficient and did not work well in practice. We do not consider this further.

178 To overcome the shortcomings of constructing moves via integer multiples of an element  
 179 from a lattice basis, we take a Markov basis for the set  $\mathcal{B}$  (Diaconis and Sturmfels 1998). A  
 180 Markov basis is a larger set of elements in  $\ker_{\mathbb{Z}}(\mathbf{A})$  that connects all fibers  $\mathcal{F}_{\mathbf{y}}$  irrespective  
 181 of the given values in  $\mathbf{y}$ . A finite set  $\mathcal{M} \subset \ker_{\mathbb{Z}}(\mathbf{A})$  is a Markov basis if, for any  $\mathbf{y}$  such that  
 182  $\mathcal{F}_{\mathbf{y}} \neq \emptyset$  and for all elements  $\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{F}_{\mathbf{y}}$ ,  $\mathbf{x}_1 \neq \mathbf{x}_2$ , there exist  $M > 0$ ,  $\mathbf{v}_1, \dots, \mathbf{v}_M \in \mathcal{M}$  and  
 183  $\epsilon_1, \dots, \epsilon_M \in \{-1, 1\}$  such that

$$184 \quad \mathbf{x}_2 = \mathbf{x}_1 + \sum_{j=1}^M \epsilon_j \mathbf{v}_j \quad \text{and} \quad \mathbf{x}_1 + \sum_{k=1}^L \epsilon_k \mathbf{v}_k \in \mathcal{F}_{\mathbf{y}}, \quad L = 1, \dots, M - 1.$$

185 The first condition says that we can use moves from a Markov basis as in the algorithm in  
 186 Figure 1 to move between any two elements of our fiber. The second condition says that  
 187 when moving between any two elements in the fiber, we always remain in the fiber (i.e. we  
 188 never encounter a negative count).

189 Although Markov bases are relatively easy to describe there is no simple algorithm for their  
 190 computation. Diaconis and Sturmfels (1998) show how a Markov basis can be computed  
 191 using techniques from commutative algebra. The theory is based on what is now known as  
 192 the Fundamental Theorem of Markov Bases which describes how finding a Markov basis is  
 193 equivalent to finding a set of generators of a toric ideal in a polynomial ring associated with  
 194 the matrix  $\mathbf{A}$ . We refer the interested reader to Cox et al. (2007) for details on commutative  
 195 algebra and to Diaconis and Sturmfels (1998), Drton et al. (2009), Aoki et al. (2012) and the  
 196 references therein for additional information on the generation of Markov bases in algebraic  
 197 statistics. Unless otherwise stated, we use the freely available software `4ti2` (Hemmecke et al.  
 198 2013) to compute the Markov bases for the examples in this manuscript.

199 For the  $3 \times 3$  contingency table, a Markov basis consists of the nine elements in (4)

	$x_{11}$	$x_{21}$	$x_{31}$	$x_{12}$	$x_{22}$	$x_{32}$	$x_{13}$	$x_{23}$	$x_{33}$
MB1	0	0	0	0	1	-1	0	-1	1
MB2	0	0	0	1	-1	0	-1	1	0
MB3	0	0	0	1	0	-1	-1	0	1
MB4	0	1	-1	0	-1	1	0	0	0
MB5	0	1	-1	0	0	0	0	-1	1
MB6	1	-1	0	-1	1	0	0	0	0
MB7	1	-1	0	0	0	0	-1	1	0
MB8	1	0	-1	-1	0	1	0	0	0
MB9	1	0	-1	0	0	0	-1	0	1

(4)

201 It is a straightforward exercise to confirm that we can transition between the two solutions in  
 202 (2) by adding or subtracting moves from (4) one-at-a-time without encountering a negative  
 203 count. More importantly, the moves in (4) can be used to connect any two solutions in the  
 204 same fiber, no matter what value of  $\mathbf{y}$  is observed.

205 There is often a need to analytically find a Markov basis for a given problem. Even though  
 206 tools like `4ti2` are freely available, computation of Markov bases remains challenging. As we  
 207 discuss later, for many of the capture-recapture examples we have explored, `4ti2` can fail to  
 208 compute Markov bases for studies with a moderate to large number of sampling occasions.

209 As we know of no simple test to confirm whether a specified set of moves  $\mathcal{B}$  is a Markov  
 210 basis, we often need to rely on theoretically derived Markov bases to confirm that our MCMC  
 211 algorithms are valid. In the following section we find such a theoretical result for a class of  
 212 capture-recapture models including  $M_{t\alpha}$ .

### 213 3. Model $M_{t\alpha}$ and Simple Corruptions

214 Here, we examine model  $M_{t\alpha}$ , the specific model of misidentification considered by L2010.  
 215 We fit this model into a larger class of models in which any identification error results in  
 216 what we refer to as a simple corruption. We then show that for any model in this class,  
 217 we can construct a lattice basis that is guaranteed to connect every element of the fiber,  
 218 irrespective of  $\mathbf{y}$ , i.e. it is also a Markov basis.

219 Model  $M_{t\alpha}$  builds on the standard closed population model with time-dependent capture  
 220 probabilities, model  $M_t$  of Otis et al. (1978), by allowing for individuals to be misidentified  
 221 when captured. The model assumes that all errors are unique meaning that an individual  
 222 cannot be identified as another individual and the same error cannot occur multiple times.  
 223 The result is that an error on the  $j^{th}$  capture occasion leads to a ghost observed history  
 224 containing a single observation on the  $j^{th}$  occasion.

225 For this model, the vector of summary statistics,  $\mathbf{y}$ , contains the counts of the  $2^K - 1$   
 226 observable capture histories. The vector of latent variables contains the counts of the possible  
 227 true histories constructed from the events:

- 228 • 0 – the individual was not captured,
- 229 • 1 – the individual was captured and correctly identified,
- 230 • 2 – the individual was captured and incorrectly identified.

231 For example, for a study with  $K = 5$  capture occasions the true history 01221 would generate  
 232 three observed histories: 01001, 00100, and 00010. Including the null history  $0 \dots 0$ , the vector

233 of true counts has length  $3^K$ . The configuration matrix,  $\mathbf{A}$ , has dimension  $(2^K - 1) \times 3^K$   
 234 and  $A_{ij} = 1$  if the  $j^{\text{th}}$  true history generates the  $i^{\text{th}}$  observed history and is equal to zero  
 235 otherwise. For example, the column corresponding to the history 01221 would contain three  
 236 non-zero entries in the rows associated with the observable histories 01001, 00100, and 00010.  
 237 A description of the model along with the vectors  $\mathbf{x}$  and  $\mathbf{y}$  and matrix  $\mathbf{A}$  for  $K = 2$  are  
 238 given in the supplementary materials, with more details in L2010.

239 A feature of  $M_{t_\alpha}$  is that whenever an error in identification occurs, it involves only one  
 240 individual and results in one or more observed histories. We define such an error as a simple  
 241 corruption. For example, the errors in true history 01221 above affect no other true history  
 242 and lead to three observed histories. Another example of simple corruptions are the errors  
 243 that occur when multiple marks cannot be matched, as described in Bonner and Holmberg  
 244 (2013) and McClintock et al. (2013). Suppose that a study uses photographs to identify  
 245 individuals and that photographs taken from the left or right side cannot be matched without  
 246 further information. In this case, any individual that is photographed from both the left and  
 247 right sides on different occasions will contribute two histories to the observed data set. Using  
 248 the events  $L$  and  $R$  to denote photographs from the left and right, the true history  $0LRRL$   
 249 would generate observed histories  $0L00L$  and  $00RR0$ . In this case, each true history will  
 250 contribute one or two histories to the observed data set.

251 For a model that contains only simple corruptions, we have the following theorem:

252 **THEOREM 1:** *Suppose that: (i)  $\mathbf{A}$  contains only the values 0 and 1 and (ii) the columns*  
 253 *of  $\mathbf{A}$  contain all of the columns of the identity matrix. Then there exists a lattice basis*  
 254  $\mathcal{L} = \{\mathbf{v}_1, \dots, \mathbf{v}_m\}$ , *which is also a Markov basis.*

255 The first condition (values of 0 and 1) occurs under the assumption of simple corruption, while  
 256 the second condition (columns of the identity matrix) occurs when every observable history  
 257 is also a true history in which there is no misidentification. Provided these assumptions hold,

258 then we can use the algorithm in Figure 1 with a suitable lattice basis  $\mathcal{L}$  and connect the  
 259 fiber. The proof of this theorem is provided in the supplementary materials, along with a  
 260 description of how to construct the lattice (Markov) basis  $\mathcal{L}$ .

261 The conditions of Theorem 1 are satisfied for model  $M_{t\alpha}$ , so that for  $K = 2$  we obtain the  
 262 Markov basis in (5)

	$x_{00}$	$x_{01}$	$x_{02}$	$x_{10}$	$x_{11}$	$x_{12}$	$x_{20}$	$x_{21}$	$x_{22}$
MB1	1	0	0	0	0	0	0	0	0
MB2	0	-1	1	0	0	0	0	0	0
MB3	0	-1	0	-1	0	1	0	0	0
MB4	0	0	0	-1	0	0	1	0	0
MB5	0	-1	0	-1	0	0	0	1	0
MB6	0	-1	0	-1	0	0	0	0	1

(5)

264 The basis in (5) is identical to that presented by L2010 for model  $M_{t\alpha}$  when  $K = 2$ .

265 The approach of L2010 to finding a basis involves choosing pivotal (or constraining)  
 266 variables when solving the set of equations  $\mathbf{A}\mathbf{x} = \mathbf{0}$  (a full description is available either in  
 267 L2010, pg 180–181, or in the supplementary materials). L2010 chose specific pivotal variables  
 268  $(x_{01}, x_{10}$  and  $x_{11})$  when finding the basis for model  $M_{t\alpha}$  when  $K = 2$ . However, it was implied  
 269 that this choice was arbitrary and no guidance was given as to how to select pivotal variables  
 270 when  $K > 2$ . It turns out that changing the pivotal variables can lead to different sets of  
 271 basis vectors which may not be Markov bases. We show in the supplementary materials that  
 272 for  $K = 2$  and a different set of pivotal variables,  $x_{22}$ ,  $x_{20}$  and  $x_{11}$ , the resulting basis differs  
 273 from that in (5). We also show that when the conditions of Theorem 1 are satisfied, there is  
 274 a specific choice of pivotal variables guaranteed to return the Markov basis  $\mathcal{L}$ . In particular,  
 275 if we order  $\mathbf{x}$  as in L2010 for model  $M_{t\alpha}$  and take the variable corresponding to the leading  
 276 non-zero entry in each row of  $\mathbf{A}$  as pivotal (as was done by L2010 for  $K = 2$ ), the basis  
 277 found will be the Markov basis  $\mathcal{L}$ .

278 Theorem 1 ensures that there is at least one lattice basis which is also a Markov basis for  
 279 model  $M_{t\alpha}$ . However, it does not imply that every lattice basis is a Markov basis. For model

280  $M_{t\alpha}$  and  $K = 2$  another lattice basis (found by hand) is given in (6)

	$x_{00}$	$x_{01}$	$x_{02}$	$x_{10}$	$x_{11}$	$x_{12}$	$x_{20}$	$x_{21}$	$x_{22}$
LB1	1	0	0	0	0	0	0	0	0
LB2	0	-1	1	0	0	0	0	1	-1
LB3	0	0	1	0	0	0	1	0	-1
LB4	0	0	0	1	0	0	-1	0	0
LB5	0	0	-1	0	0	1	-1	1	-1
LB6	0	0	0	0	0	0	0	1	-1

(6)

282 Suppose the observed data are  $\mathbf{y} = (363, 22, 174)$  (as in L2010), then two elements in the fiber  
 283 are  $\mathbf{x}_1 = (0, 363, 0, 22, 174, 0, 0, 0, 0, 0)'$  and  $\mathbf{x}_2 = (0, 361, 2, 22, 174, 0, 0, 0, 0, 0)'$ . We are unable to  
 284 move between these two using LB1 – LB6 in (6) as moves in the algorithm in Figure 1. In  
 285 particular, if we start at (the observed history)  $\mathbf{x}_1$  the moves LB2, LB3, LB5 and LB6 will  
 286 lead to automatic rejections because they will always propose a negative value. This means  
 287 that  $\mathbf{x}_1$  and  $\mathbf{x}_2$  are not connected and thus the fiber is not connected.

288 We repeated the analysis of L2010 using both the Markov basis in (5) and the lattice  
 289 basis in (6) using the same prior distributions as in L2010 (we used only one of the priors  
 290 L2010 considered for  $\alpha$ ; a beta distribution with parameters 19 and 1). In both cases we  
 291 implemented the algorithm in Figure 1 using  $\mathbf{x}_1$  as the starting value with interest in the  
 292 abundance  $N$ . We checked convergence via trace plots and plotted the resulting distribution  
 293 for  $N|y$  in both cases (Figure 3). The lattice basis in (6) leads to a distribution for  $N$  that  
 294 is substantially different from the true posterior distribution and could lead to incorrect  
 295 decision making.

296 [Figure 3 about here.]

297 We note that efficiency gains can be made if there are observable histories with zero count.  
 298 In particular, we can delete the entries in  $\mathbf{y}$  and the rows of  $\mathbf{A}$  corresponding to the zero  
 299 counts before deleting any columns of  $\mathbf{A}$  and corresponding entries of  $\mathbf{x}$  that are known to  
 300 have zero count. Provided the assumptions of Theorem 1 are still satisfied by the resulting

301 configuration matrix then we can still find a set of moves guaranteed to connect all elements  
 302 in the fiber. The resulting set of moves is no longer a Markov basis but a Markov subbasis  
 303 (Chen et al. 2006) as it is only valid for the observed  $\mathbf{y}$ . This corresponds to the approach  
 304 taken by both Bonner and Holmberg (2013) and McClintock et al. (2013) for data with  
 305 multiple marks that could not be matched.

306 This section shows that we must take care even with simple corruptions to ensure that the  
 307 lattice basis we are using is also a Markov basis. The following two sections give examples  
 308 where we do not have simple corruptions (in one of these it does not even make sense to  
 309 think of corruptions in the sense of model  $M_{t\alpha}$ ) and a Markov basis has greater cardinality  
 310 than a lattice basis.

#### 311 4. Example: Sufficient Statistics

312 Next we consider the problem of modeling data from a closed population when sufficient  
 313 statistics from one or more models are provided in place of the raw data. The raw data  
 314 may not be available for a variety of reasons, e.g. privacy concerns. Here we assume that  
 315 the population is closed and that we have the sufficient statistics associated with three  
 316 commonly used models  $M_t$ ,  $M_b$  and  $M_h$  (Otis et al. 1978). From model  $M_h$  we have the  
 317 statistics  $f_1, \dots, f_K$ , where  $f_j$  is the number of individuals who were caught  $j$  times from a  
 318 total of  $K$  sampling occasions; from model  $M_t$  we have the statistics  $n_1, \dots, n_K$ , where  $n_j$   
 319 is the number of individuals captured in the  $j$ th sample; and from model  $M_b$  we have the  
 320 statistic  $M. = \sum_{j=1}^t M_j$ , with  $M_j$  the number of marked individuals in the population in  
 321 sample  $j$ . Note that we do not include the other sufficient statistics for model  $M_t$  and  $M_b$   
 322 noted by Otis et al. (1978) as they are deterministic functions of  $f_1, \dots, f_K$ .

323 All of these statistics are linear functions of the data which means that this problem can  
 324 be expressed using the linear constraint in (1). In this example,  $\mathbf{x}$  represents the vector of  
 325 counts for the  $2^K - 1$  true histories;  $\mathbf{y}$  represents the vector of counts for the  $2K + 1$  sufficient

326 statistics; and the configuration matrix,  $\mathbf{A}$ , is a  $(2K + 1) \times (2^K - 1)$  matrix. Details of how  
 327 to find  $\mathbf{A}$  along with an example for a study with  $K = 4$  occasions are provided in the  
 328 supplementary materials.

329 Here we explore this scenario using multi-list data from a South Auckland, New Zealand,  
 330 diabetes study from the Ph.D. research of Huakau (2001) and included in the Ph.D. research  
 331 of Sutherland (2003). We ignore the potential errors in matching individuals between lists and  
 332 assume that each individual is correctly matched (see Lee (2002) for how such errors could  
 333 also be accounted for using the linear constraint (1)). There are  $K = 4$  lists: general prac-  
 334 titioners records (G), pharmacy records (P), outpatient records (O) and inpatient discharge  
 335 records (D) that we assume are ordered as written. We use the data for males and reduce the  
 336 full data (which is available in Sutherland 2003) to the statistics:  $\mathbf{n} = (n_G, n_P, n_O, n_D)' =$   
 337  $(629, 622, 6279, 1623)'$ ,  $\mathbf{f} = (f_1, f_2, f_3, f_4)' = (6030, 1312, 161, 4)'$  and  $M = 8680$  to give

$$338 \quad \mathbf{y} = (6030, 1312, 161, 4, 629, 622, 6279, 1623, 8680)'$$

339 As well as  $\mathbf{y}$  being sufficient for models  $M_t$ ,  $M_h$  and  $M_b$ , it is also sufficient for the two-factor  
 340 quasi-symmetric version of model  $M_{th}$  that is induced by a Rasch model (see Agresti 1994,  
 341 for details of this model).

342 The vector  $\mathbf{x}$  is indexed by  $\boldsymbol{\omega} = (\omega_G, \omega_P, \omega_O, \omega_D)$ , where  $\omega_j = 1$  denotes inclusion on list  
 343  $j$  with  $\omega_j = 0$  otherwise, so that  $x_{1101}$  is the number of individuals on lists G, P and D and  
 344 not on list O. Our focus here is to attempt to make inference about  $x_{1000}$ , the number of  
 345 individuals who appear only in list  $G$ . We may also wish to fit a model to  $\mathbf{x}$  for which  $\mathbf{y}$   
 346 are not sufficient statistics. By definition, the resulting model would be nonidentifiable, but  
 347 this does not necessarily mean that there is no information about parameters of this model,  
 348 including the abundance  $N$ . The latent multinomial model can be used in either of these  
 349 situations.



A lattice basis found using the Hermite normal form is

	$x_{0001}$	$x_{0010}$	$x_{0011}$	$x_{0100}$	$x_{0101}$	$x_{0110}$	$x_{0111}$	$x_{1000}$	$x_{1001}$	$x_{1010}$	$x_{1011}$	$x_{1100}$	$x_{1101}$	$x_{1110}$	$x_{1111}$
LB1	0	0	0	0	0	0	0	0	-1	0	1	1	0	-1	0
LB2	0	0	0	0	0	0	0	0	-1	1	0	0	1	-1	0
LB3	0	-1	1	0	0	0	0	1	-1	0	0	0	0	0	0
LB4	1	-2	0	1	0	0	0	0	0	0	1	0	-2	1	0
LB5	1	-2	0	0	1	0	0	1	-1	0	1	0	-2	1	0
LB6	1	-2	0	0	0	1	0	1	-1	0	1	0	-1	0	0
LB7	1	-2	0	0	0	0	1	1	0	0	0	0	-2	1	0

Using the seven moves LB1 – LB7 in the algorithm in Figure 1 it is impossible to move between the two solutions  $\mathbf{x}_1$  and  $\mathbf{x}_2$

$$\mathbf{x}_1 = (652, 4865, 794, 253, 18, 234, 62, 260, 26, 221, 67, 19, 0, 32, 4)'$$

$$\mathbf{x}_2 = (684, 4901, 694, 253, 31, 154, 161, 192, 49, 365, 0, 19, 0, 0, 4)'$$

If we are currently at  $\mathbf{x}_2$ , it is clear that all moves (except LB3) will lead to at least one negative cell count and will be automatically rejected. The vector LB3 can be used to update  $\mathbf{x}_2$ , but we are unable to get to  $\mathbf{x}_1$  using LB3 alone. Again, we have at least two sets of elements in the fiber that we can move within, but are unable to move between.

A Markov basis for this problem can be constructed in `4ti2` and is made up of the 16 elements given in the supplementary materials. Since (i) `4ti2` finds a minimal Markov basis, and (ii) the cardinality of the Markov basis is larger than that of a lattice basis, we can be certain that a lattice basis can never be a Markov basis for this problem. Even though it is likely possible to construct another lattice basis that can move between  $\mathbf{x}_1$  and  $\mathbf{x}_2$  there will be either (i) another two elements in the fiber that are not connected, or (ii) another two elements in the fiber for a different  $\mathbf{y}$  that we cannot move between with such a lattice basis.

Here we fit model  $M_t$  and run the algorithm in Figure 1 with both the Markov basis given in the supplementary materials and the lattice basis specified above (details of the model are given in the supplementary materials). We make use of the factorization theorem (e.g, see Casella and Berger 2002, pg. 276) that states that a model  $f(\mathbf{x}|\boldsymbol{\theta})$  with sufficient statistics

372  $\mathbf{y}$  can be expressed as

$$373 \quad f(\mathbf{x}|\boldsymbol{\theta}) = g(\mathbf{x}|\mathbf{y})h(\mathbf{y}|\boldsymbol{\theta}).$$

374 A practical implication is that only  $g(\mathbf{x}|\mathbf{y})$  is required if interest is in a function of  $\mathbf{x}$  such as  
 375  $x_{1000}$ , and the parameters  $\boldsymbol{\theta} = (N, p_1, \dots, p_K)$  need not be specified. A related implication is  
 376 that if we do choose to update  $\boldsymbol{\theta}$  the resulting chains will converge to the correct posterior  
 377  $[\boldsymbol{\theta}|\mathbf{y}]$  even if we (i) do not update  $\mathbf{x}$ , or (ii) update  $\mathbf{x}$  using a set of moves that is unable to  
 378 connect the fiber, such as the lattice basis above; provided we specify an appropriate MCMC  
 379 sampler for  $\boldsymbol{\theta}$ .

380 Using the lattice basis and starting at  $\mathbf{x}_2$  the resulting distributions for  $x_{1000}$  are qualita-  
 381 tively different from the posterior distribution found using the Markov basis even though the  
 382 individual chains appear to have converged to the stationary distribution (Figure 4). The  
 383 true value of  $x_{1000} = 260$  has some posterior mass when using a Markov basis (despite being  
 384 in the tail). If we were to believe the results when using the lattice basis  $x_{1000} = 260$  is so  
 385 far in the tail, we would conclude it has negligible posterior mass.

386 [Figure 4 about here.]

## 387 5. Example: Band Misreading in Mark-Resight

388 As a final example we consider a mark-resight model which allows for the possibility that  
 389 individuals are misidentified when resighted in the field. Imagine that there are  $K_1$  distinct  
 390 occasions, on which researchers capture a number of unmarked individuals, mark them, and  
 391 release them back into the population. Along with that are a series of  $K_2$  resighting occasions,  
 392 on which the researchers conduct visual surveys to identify previously marked individuals.  
 393 Data from the experiment consist of the observed resighting histories for each individual. If  
 394 there were no errors then standard mark-resight models could be used to estimate survival  
 395 or movement rates (e.g. Hestbeck et al. 1991); or abundance (e.g. McClintock et al. 2006).

396 Suppose now that individuals may be misidentified when they are resighted. In direct con-  
397 trast to model  $M_{t\alpha}$ , which assumes that errors are unique and never match other individuals,  
398 we assume that errors may be repeated and always match the identity of previously marked  
399 individuals. The justification for this assumption is that the available set of marks is known on  
400 each occasion when individuals are identified by man-made marks instead of natural markers  
401 (e.g., genotypes or photo-id). Erroneous sightings of marks which have not been released can  
402 then be identified and removed from the data prior to the analysis. The only time an error  
403 cannot be detected and discarded is when one previously marked individual is misidentified  
404 as another previously marked individual. We note that removal of erroneous sightings is only  
405 justified when estimating survival. Removing erroneous sightings when including unmarked  
406 individuals would lead to biased estimators of abundance (McClintock et al. 2014).

407 For the remainder of the section, we assume that the capture and resighting occasions  
408 occur simultaneously so that  $K = K_1 = K_2$ . The true capture histories for each individual  
409 can now be constructed in terms of four possible events. On each occasion, individual  $i$  may  
410 be:

- 411 • not captured or resighted (event 0),
- 412 • captured or resighted and correctly identified (event 1), or
- 413 • resighted and incorrectly identified (event 2).

414 Further to this, another individual may be resighted and incorrectly identified as individual  $i$   
415 (event 3). Events 2 and 3 represent false negative and false positive resightings. For example,  
416 the history 123 for individual  $i$  would indicate that  $i$  was captured and marked on the first  
417 occasion, was resighted and misidentified on the second occasion, and that another individual  
418 was resighted and identified as  $i$  on the third occasion of a study with  $K = 3$  occasions. To  
419 simplify the example, we assume that individuals cannot be misidentified when they are first  
420 captured and that multiple events involving the same individual cannot occur on a single

421 occasion (e.g., it is not possible to resight  $i$  and incorrectly identify another individual as  
 422  $i$  on the same occasion). This assumption may be unrealistic in some situations and was  
 423 made to make the approach tractable. Developing methodology to relax this assumption is  
 424 ongoing research.

425 For an experiment with  $K$  occasions, the model has  $(4^K - 1)/3$  possible true histories and  
 426 the usual  $2^K - 1$  observable histories. Further to this, there are  $K - 1$  extra constraints that  
 427 equate the number of false negatives and false positives (2s and 3s) on occasions 2 through  
 428  $K$ . As a result,  $\mathbf{A}$  has dimension  $(2^K + K - 2) \times (4^K - 1)/3$  and a basis for  $\ker_{\mathbb{Z}}(\mathbf{A})$  has  
 429  $(4^K - 1)/3 - (2^K + K - 2)$  elements.

430 To make this more concrete, we consider the specific case of an experiment comprising  
 431  $K = 3$  occasions. In this case, there are  $(4^3 - 1)/3 = 21$  possible true histories,  $2^3 - 1 = 7$   
 432 observable histories, and  $3 - 1 = 2$  extra constraints on the number of false positive and  
 433 negative resightings (2s and 3s) on occasions 2 and 3. Details of how to construct  $\mathbf{A}$  along  
 434 with  $\mathbf{x}$  and  $\mathbf{y}$  for a study with  $K = 3$  capture occasions are provided in the supplementary  
 435 materials. In this case, a basis for  $\ker_{\mathbb{Z}}(\mathbf{A})$  has 12 elements and the specific lattice basis  
 436 obtained using the Hermite normal form is provided in the supplementary materials, along  
 437 with the Markov basis, computed using `4ti2`, that has 63 elements.

To illustrate the problems that can occur with this model we first consider the analysis of  
 a single (fake) data set. Suppose that each observable history is recorded one time so that

$$\mathbf{y} = (1, 1, 1, 1, 1, 1, 1).$$

438 An exhaustive search confirms that the fiber defined by  $\mathbf{y}$  contains exactly 120 unique  
 439 elements. However, the lattice basis given in the supplementary materials does not connect  
 440 all of the elements in the fiber. Instead, the lattice basis divides the fiber into two distinct  
 441 pieces including a large set of 87 connected elements; and a further set of 33 isolated elements  
 442 which connect to nothing else. As a result, the distribution of the sample generated by the

443 algorithm in Figure 1 using the elements of the lattice basis in the supplementary materials  
 444 as moves will depend on the starting point.

445 To show this, we have investigated the output from the algorithm in Figure 1 when using  
 446 a lattice basis as our set of moves. We have chosen a starting point that lies in the largest  
 447 part of the fiber and connects with 86 other elements:

$$448 \quad \mathbf{x}_1 = (1, 0, 1, 1, 0, 1, 1, 0, 0, 1, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0)'$$

449 Assuming a multinomial distribution for  $[\mathbf{x}|\boldsymbol{\theta}]$  is not appropriate to account for the band  
 450 misreading process and specification of a more complex  $[\mathbf{x}|\boldsymbol{\theta}]$  is ongoing research. As our  
 451 goal is to show that a lattice basis is unable to connect the fiber, we simplify the model by  
 452 setting  $[\mathbf{x}|\boldsymbol{\theta}] \propto 1$ . A valid sampler should then sample uniformly from the 120 elements in  
 453 the fiber. For comparison, we have also run a chain using the full Markov basis starting at  
 454  $\mathbf{x}_1$ . As expected, the first chain visits 87 unique solutions and the second visits all 120. To  
 455 visualize the impact this can have on inference, Figure 5 compares the distributions of the  
 456 number of errors in the solutions identified by each chain. Using the lattice basis, the first  
 457 chain oversamples the solutions with too few errors, placing too much mass on solutions with  
 458 one or two errors and not enough on solutions with three, four, or five errors. In comparison,  
 459 the distribution generated using the full Markov basis matches the true distribution of the  
 460 number of errors in the 120 elements almost exactly.

461 [Figure 5 about here.]

## 462 6. Discussion

463 Here we have presented examples of capture-recapture models that show the importance  
 464 of using a Markov basis when sampling from a linearly constrained vector of counts. In  
 465 particular, we have demonstrated the danger of using elements of a lattice basis as one-at-a-  
 466 time moves in an algorithm as in Figure 1. In many situations a set referred to as a Markov

467 basis is needed to ensure we can move between various elements of the fiber without passing  
468 through invalid (negative) counts. Even when a Markov basis is a lattice basis, we must take  
469 care because not every lattice basis is a Markov basis.

470 For a given matrix  $\mathbf{A}$  the need for a Markov basis over a simpler lattice basis depends on  
471 the lattice basis chosen, as well as the data observed. If we consider the lattice basis for the  
472  $3 \times 3$  contingency table in section 2, difficulties arose because our data had a row sum of 0.  
473 A related issue is that even when a lattice basis is unable to connect the fiber, it may still  
474 be able to connect nearly all elements in the fiber. In such a case, using a lattice basis may  
475 lead to a distribution that is an acceptable approximation of the true posterior distribution.  
476 This is especially the case if the elements of the fiber that are not connected to the initial  
477 value are in areas of low probability in the model  $[\mathbf{x}|\boldsymbol{\theta}]$ . This can be seen in the example from  
478 Section 4: using the lattice basis and starting at the second starting value (Figure 2; right  
479 panel) results in an estimated posterior density that is practically indistinguishable from the  
480 true posterior distribution (Figure 4). However, there is no guarantee that any given lattice  
481 basis will provide a good approximation to the fiber. It is possible that even with multiple  
482 starting values we may choose values that only connect a small proportion of the fiber.

483 One important aspect that we have only briefly mentioned is the difficulty in constructing  
484 Markov bases. For the purposes of this manuscript we have overcome this difficulty through  
485 (i) analytical results, or (ii) the use of the software package `4ti2` (Hemmecke et al. 2013).  
486 While the latter is possible for the examples we explored, it is unable to evaluate a Markov  
487 basis for some capture-recapture examples with a moderate to large number of sampling  
488 occasions. For example, `4ti2` was unable to compute a Markov basis (on the lead authors  
489 work machine) for the band read error model in section 5 for  $K > 4$ . If we were to use `4ti2`  
490 for model  $M_{t\alpha}$  (ignoring the theorem presented in section 3), `4ti2` was unable to compute  
491 a Markov basis for  $K > 5$ . The implication of this is that for an algorithm in the spirit of

492 Figure 1 to be implemented for problems not involving simple corruptions, methodological  
493 work is likely to be necessary to ensure a potential set of moves is a Markov basis.

494 Several alternative algorithms and methods have been proposed for sampling from the fiber  
495 that avoid the calculation of a full Markov basis. We anticipate that such approaches may  
496 be useful for a range of capture-recapture examples. These include independent sampling  
497 of elements of the fiber (e.g., see Chen et al. 2005), extending the algorithm to allow  
498 limited travel through vectors  $\boldsymbol{x}$  that contain negative values while using a set of moves  
499 that is not guaranteed to connect the fiber (e.g., see Bunea and Besag 2000) and approaches  
500 that dynamically find a Markov basis as the algorithm runs (e.g., see Dobra 2012). While  
501 promising, we expect these approaches will require adapting to the particular challenges  
502 faced in problems involving misidentification in capture-recapture data.

#### 503 ACKNOWLEDGEMENTS

504 We thank three anonymous referees for comments on a previous version of this manuscript  
505 as well as Ruriko Yoshida for participating in valuable discussion on algebraic statistics.

#### 506 SUPPLEMENTARY MATERIALS

507 Web Appendices referenced in Sections 2, 3, 4 and 5 are available with this paper at the  
508 Biometrics website on Wiley Online Library.

#### 509 REFERENCES

- 510 Agresti, A. (1994). Simple capture-recapture models permitting unequal catchability and  
511 variable sampling effort. Biometrics **50**, 494–500.
- 512 Amstrup, S. C., McDonald, T. L., and Manly, B. F. J. (2005). Handbook of  
513 Capture-Recapture Analysis. Princeton University Press.
- 514 Aoki, S., Hara, H., and Takemura, A. (2012). Markov Bases in Algebraic Statistics. Springer.

- 515 Bonner, S. J. and Holmberg, J. (2013). Mark-recapture with multiple non-invasive marks.  
516 Biometrics **69**, 766–775.
- 517 Bunea, F. and Besag, J. (2000). MCMC in  $i \times j \times k$  contingency tables. Fields Institute  
518 Communications **26**, 23–36.
- 519 Casella, G. and Berger, R. L. (2002). Statistical Inference. Duxbury Pacific Grove, CA.
- 520 Chen, Y., Diaconis, P., Holmes, S. P., and Liu, J. S. (2005). Sequential Monte Carlo methods  
521 for statistical analysis of tables. Journal of the American Statistical Association **100**,  
522 109–120.
- 523 Chen, Y., Dinwoodie, I., and Sullivant, S. (2006). Sequential importance sampling for  
524 multiway tables. The Annals of Statistics **34**, 523–545.
- 525 Cox, D., Little, J., and O’Shea, D. (2007). Ideals, Varieties, and Algorithms: An Introduction  
526 to Computational Algebraic Geometry and Commutative Algebra. Springer.
- 527 Diaconis, P. and Sturmfels, B. (1998). Algebraic algorithms for sampling from conditional  
528 distributions. The Annals of Statistics **26**, 363–397.
- 529 Dobra, A. (2012). Dynamic Markov bases. Journal of Computational and Graphical Statistics  
530 **21**, 496–517.
- 531 Drton, M., Sturmfels, B., and Sullivant, S. (2009). Lectures on Algebraic Statistics. Springer.
- 532 Fienberg, S. E. and Manrique-Vallier, D. (2009). Integrated methodology for multiple systems  
533 estimation and record linkage using a missing data formulation. AStA Advances in  
534 Statistical Analysis **93**, 49–60.
- 535 Hemmecke, R., Hemmecke, R., Koeppe, M., Malkin, P., and Walter, M. (2013). User’s guide  
536 for 4ti2 version 1.6.
- 537 Hestbeck, J. B., Nichols, J. D., and Malecki, R. A. (1991). Estimates of movement and site  
538 fidelity using mark-resight data of wintering Canada geese. Ecology **72**, 523–533.
- 539 Huakau, J. T. (2001). New methods for analysis of epidemiological data using



- 540 capture-recapture methods. PhD thesis, The University of Auckland.
- 541 Karwa, V. and Slavkovic, A. (2013). Conditional inference given partial information in  
542 contingency tables using Markov bases. Wiley Interdisciplinary Reviews: Computational  
543 Statistics **5**, 207–218.
- 544 King, R., Bird, S. M., Hay, G., and Hutchinson, S. J. (2009). Estimating current injectors  
545 in Scotland and their drug-related death rate by sex, region and age-group via Bayesian  
546 capture-recapture methods. Statistical Methods in Medical Research **18**, 341–359.
- 547 Lee, A. (2002). Effect of list errors on the estimation of population size. Biometrics **58**,  
548 185–191.
- 549 Lee, A. J., Seber, G. A. F., Holden, J. K., and Huakau, J. T. (2001). Capture–recapture,  
550 epidemiology, and list mismatches: several lists. Biometrics **57**, 707–713.
- 551 Link, W. A., Yoshizaki, J., Bailey, L. L., and Pollock, K. H. (2010). Uncovering a latent  
552 multinomial: analysis of mark–recapture data with misidentification. Biometrics **66**,  
553 178–185.
- 554 Lum, K., Price, M. E., and Banks, D. (2013). Applications of multiple systems estimation  
555 in human rights research. The American Statistician **67**, 191–200.
- 556 McClintock, B. T., Conn, P., Alonso, R., and Crooks, K. R. (2013). Integrated modeling of  
557 bilateral photo-identification data in mark-recapture analyses. Ecology **94**, 1464–1471.
- 558 McClintock, B. T., Hill, J. M., Fritz, L., Chumbley, K., Luxa, K., and Diefenbach, D. R.  
559 (2014). Mark-resight abundance estimation under incomplete identification of marked  
560 individuals. Methods in Ecology and Evolution **5**, 1294 – 1304.
- 561 McClintock, B. T., White, G. C., and Burnham, K. P. (2006). A robust design mark-  
562 resight abundance estimator allowing heterogeneity in resighting probabilities. Journal  
563 of Agricultural, Biological, and Environmental Statistics **11**, 231–248.
- 564 Otis, D. L., Burnham, K. P., White, G. C., and Anderson, D. R. (1978). Statistical inference

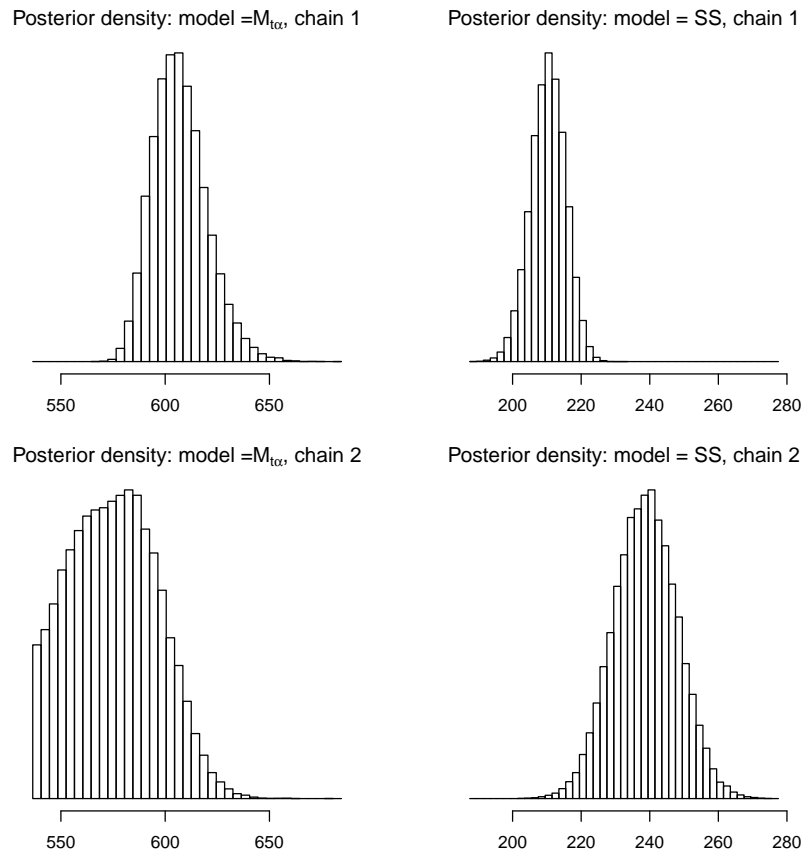
- 565 from capture data on closed animal populations. Wildlife Monographs **62**, 1–135.
- 566 Seber, G. A., Huakau, J. T., and Simmons, D. (2000). Capture-recapture, epidemiology, and  
567 list mismatches: two lists. Biometrics **56**, 1227–1232.
- 568 Sutherland, J. and Schwarz, C. (2005). Multi-list methods using incomplete lists in closed  
569 populations. Biometrics **61**, 134–140.
- 570 Sutherland, J. M. (2003). Multi-list methods in closed populations with stratified or  
571 incomplete information. PhD thesis, Simon Fraser University.
- 572 Wright, J. A., Barker, R. J., Schofield, M. R., Frantz, A. C., Byrom, A. E., and Gleeson,  
573 D. M. (2009). Incorporating genotype uncertainty into mark-recapture-type models for  
574 estimating abundance using DNA samples. Biometrics **65**, 833–840.
- 575 Yoshizaki, J., Brownie, C., Pollock, K. H., and Link, W. A. (2011). Modeling misidentification  
576 errors that result from use of genetic tags in capture–recapture studies. Environmental  
577 and Ecological Statistics **18**, 27–55.
- 578 Yoshizaki, J., Pollock, K. H., Brownie, C., and Webster, R. A. (2009). Modeling misidenti-  
579 fication errors in capture-recapture studies using photographic identification of evolving  
580 marks. Ecology **90**, 3–9.

*Received MMM 2015. Revised MMM 2015.*

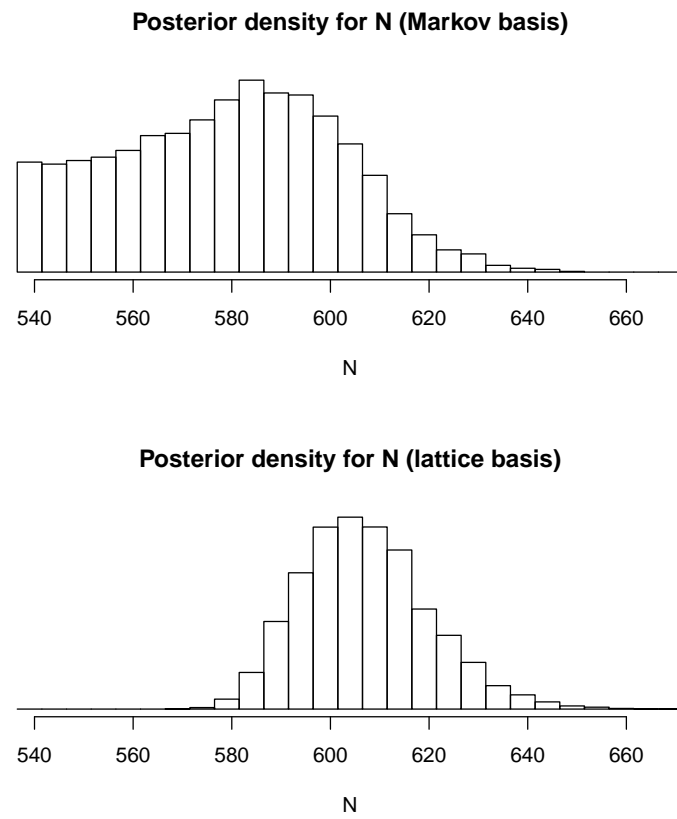
*Accepted MMM 2015.*

- 1: Initialize  $\mathbf{x}^0$  so that  $\mathbf{y} = \mathbf{A}\mathbf{x}^0$
- 2: **for**  $i = 1 : n$  **do**
- 3:     Sample  $k \in \{1, 2, \dots, m\}$  with equal probability
- 4:     Sample  $c \in \{-1, 1\}$  with equal probability
- 5:     Set  $\mathbf{x}_{\text{cand}} = \mathbf{x}^{i-1} + c\mathbf{a}_k$
- 6:     Calculate the metropolis acceptance probability:  $r = \min\left(1, \frac{[\mathbf{x}_{\text{cand}}]_{|\theta|}}{[\mathbf{x}^{i-1}]_{|\theta|}}\right)$
- 7:     Accept  $\mathbf{x}_{\text{cand}}$  with probability  $r$  (if accepted  $\mathbf{x}^i = \mathbf{x}_{\text{cand}}$ ; otherwise  $\mathbf{x}^i = \mathbf{x}^{i-1}$ )
- 8: **end for**

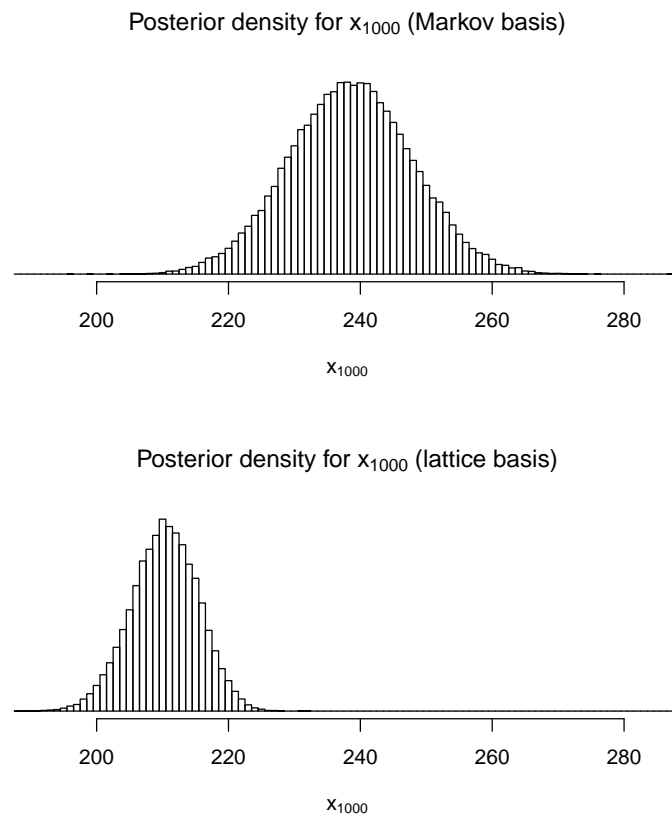
**Figure 1.** Algorithm for updating the latent counts  $\mathbf{x}$ . The value  $n$  is the number of iterations in the algorithm and the vectors  $\mathcal{B} = \{\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_m\}$  are a subset of the kernel of  $\mathbf{A}$ .



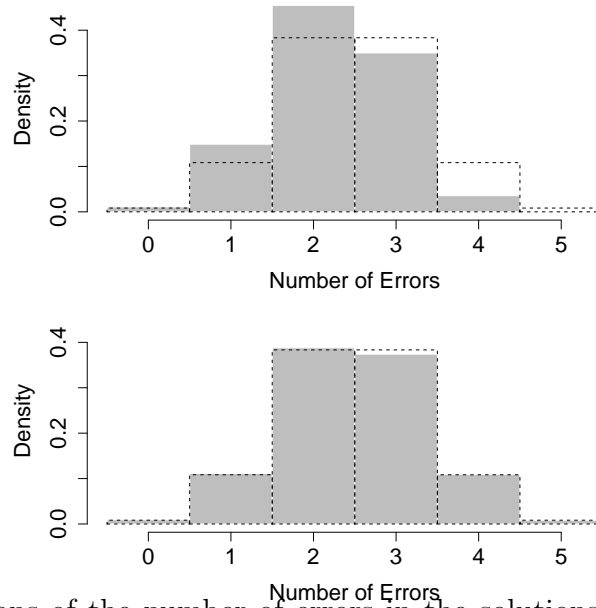
**Figure 2.** Estimated posterior densities of a quantity of interest for model  $M_{t\alpha}$  (left panel) and a multi-list model where summary statistics are presented in place of full data (SS; right panel). Within each model, the resulting density estimates are plotted separately from the output of two parallel MCMC algorithms (for each model) with different starting values.



**Figure 3.** Histograms of the estimated posterior density of  $N|y$  when using the Markov basis from (5) (top) and the lattice basis from (6) (bottom) when starting from  $\mathbf{x}_1$ .



**Figure 4.** Posterior densities of  $x_{1000}$  when using the Markov basis from the supplementary materials (top) and the lattice basis specified in section 4 (bottom) when starting at  $\mathbf{x}_2$ .



**Figure 5.** Distributions of the number of errors in the solutions sampled given the data  $\mathbf{y}$ . The top histogram illustrates the distribution generated using the lattice basis with the starting value  $\mathbf{x}_1$ . The bottom plot illustrates the distribution obtained using the full Markov basis with the same starting value. In each plot, the gray bars represent the distribution of the number of errors while the dashed bars represent the true distribution over all 120 unique solutions.