

# Bayesian Stable Isotope Mixing Models

Andrew C. Parnell<sup>\*1</sup>, Donald L. Phillips<sup>†2</sup>, Stuart Bearhop<sup>†3</sup>, Brice X. Semmens<sup>†4</sup>, Eric J. Ward<sup>†5</sup>, Jonathan W. Moore<sup>†6</sup>, Andrew L. Jackson<sup>†7</sup>, and Richard Inger<sup>3</sup>

<sup>1</sup>School of Mathematical Sciences (Statistics), Complex and Adaptive Systems Laboratory, University College Dublin, Ireland

<sup>2</sup>U.S. Environmental Protection Agency, National Health & Environmental Effects Research Laboratory, Oregon, USA

<sup>3</sup>Centre for Ecology and Conservation, School of Biosciences, University of Exeter, UK

<sup>4</sup>Scripps Institution of Oceanography, University of California, San Diego, 9500 Gilman Drive, La Jolla, California, USA

<sup>5</sup>Northwest Fisheries Science Center, National Marine Fisheries Service, National Oceanic and Atmospheric Administration, Seattle, USA

<sup>6</sup>Earth2Ocean Research Group, Simon Fraser University, Burnaby, Canada

<sup>7</sup>School of Natural Sciences, Trinity College Dublin, Ireland

October 1, 2012

## Abstract

In this paper we review recent advances in Stable Isotope Mixing Models (SIMMs) and place them into an over-arching Bayesian statistical framework which allows for several useful extensions. SIMMs are used to quantify the proportional contributions of various sources to a mixture. The most widely used application is quantifying the diet of organisms based on the food sources they have been observed to consume. At the centre of the multivariate statistical model we propose is a compositional mixture of the food sources corrected for various metabolic factors. The compositional component of our model is based on the isometric log ratio (ilr) transform of Egozcue et al. (2003). Through this transform we can apply a range of time series and non-parametric smoothing relationships. We illustrate our models with 3 case studies based on real animal dietary behaviour.

## 1 Introduction

Stable isotope analysis is an increasingly important tool in the study of ecological food webs. The technique utilises the fact that biologically active elements exist in more than one isotopic form. Generally the lighter isotopic form is much more abundant in the environment than the heavier form, although their relative abundance is altered by a range of biological, geochemical and anthropogenic processes. These processes produce isotopic gradients which are reflected in the tissues of plants and animals. Differences in relative abundance of these isotopes within a particular sample can be measured using a mass spectrometer and expressed as the ratio of heavy to light form, which can then be standardised against international reference samples and reported in the delta ( $\delta$ ) notation as parts per thousand or per mil ( $\text{‰}$ ).

---

\*Andrew.Parnell@ucd.ie

†These authors contributed equally to this manuscript and are listed in random order

36 As a consumer's tissues are ultimately derived from the dietary sources they consume, it is possible to use  
37 stable isotope mixing models (SIMMs) to derive the assimilated diet of an individual, or a group of indi-  
38 viduals, given the isotopic ratios of the consumers' tissues and food sources (Phillips, 2012). A number  
39 of recent papers have proposed models to analyse such data, gathering over 1500 citations since their first  
40 introduction. More recently, the models proposed for such data are Bayesian. In this paper we review the  
41 different models proposed and bring them into an over-arching framework. We include three case studies  
42 ranging from the simple to the complex, together with JAGS (Just Another Gibbs Sampler; Plummer, 2003)  
43 code for their implementation<sup>1</sup>.

44  
45 Generally the isotopic ratios of a sample of a consumer's tissues (e.g. blood, feathers, whiskers) are mea-  
46 sured along with a representative sample of potential items from a consumer's diet. The consumer isotopic  
47 values are represented in the model as the convex combination of the source values where the coefficients  
48 in the simplex are the 'dietary proportions'; strictly speaking they are the proportion of the consumers'  
49 dietary proteins obtained from the sources. Estimation of these dietary proportions is the main focus of  
50 our analysis. Most commonly the isotopic observations on the consumers and sources are multivariate with  
51 dimension 2. A thorough description of the uses of stable isotopes can be found in Inger and Bearhop (2008).

52  
53 Once the isotopic data have been collected for both consumer and sources, it is usual to create an *iso-space*  
54 plot which shows the consumer and source values. An example is shown in Figure 1. It is desirable for the  
55 consumer values to lie within the fuzzy convex hull of the sources. However, a further phenomenon is often  
56 observed here, that of *trophic enrichment*, whereby light isotopes are lost during the conversion of source  
57 proteins into consumer tissues. The isotopic values of the consumer (or equivalently the sources) are thus  
58 adjusted by a *trophic enrichment factor* (TEF) which may vary by food source and consumer. These TEF  
59 corrections arise from laboratory studies, and thus contribute another set of (uncertain) data to our analysis.

60  
61  
62 The inference challenge involved in a SIMM is thus to estimate the dietary proportions whilst taking account  
63 of the uncertainty in the source and TEF values. Clearly not all of the consumers will eat exactly the same  
64 diet, so it is common to use a hierarchical model. Furthermore, covariates such as time or age may be  
65 available which are thought to influence the dietary proportions. The model we present in this paper takes  
66 account of these common features in a multivariate hierarchical Bayesian model.

67  
68 The paper is organised as follows. In Section 2 we introduce our general SIMM and show it may be extended  
69 to form part of a larger class of hierarchical compositional generalised linear models. In Section 3 we outline  
70 previous work in this area and show how each of these fits into our general SIMM. Section 4 outlines the  
71 statistical issues concerning the formulation of such models and how they can be fitted. We outline three  
72 case studies in Section 5, showing how the model works in different situations depending on the information  
73 available. We discuss future directions in Section 6.

## 75 2 Model formulation

76 We first provide the notation we use to formulate our model. We suppose that there are  $N$  consumers on  $J$   
77 isotopes, with  $K$  sources. We outline the most important components of our model below:

- 78 •  $Y_{ij}$  represents the isotope measurement on consumer  $i$  for isotope  $j$ . We write  $\mathbf{Y}_i$  as the  $J$ -vector of  
79 isotope values for consumer  $i$

---

<sup>1</sup>See [mathsci.ucd.ie/~parnell\\_a/](http://mathsci.ucd.ie/~parnell_a/)

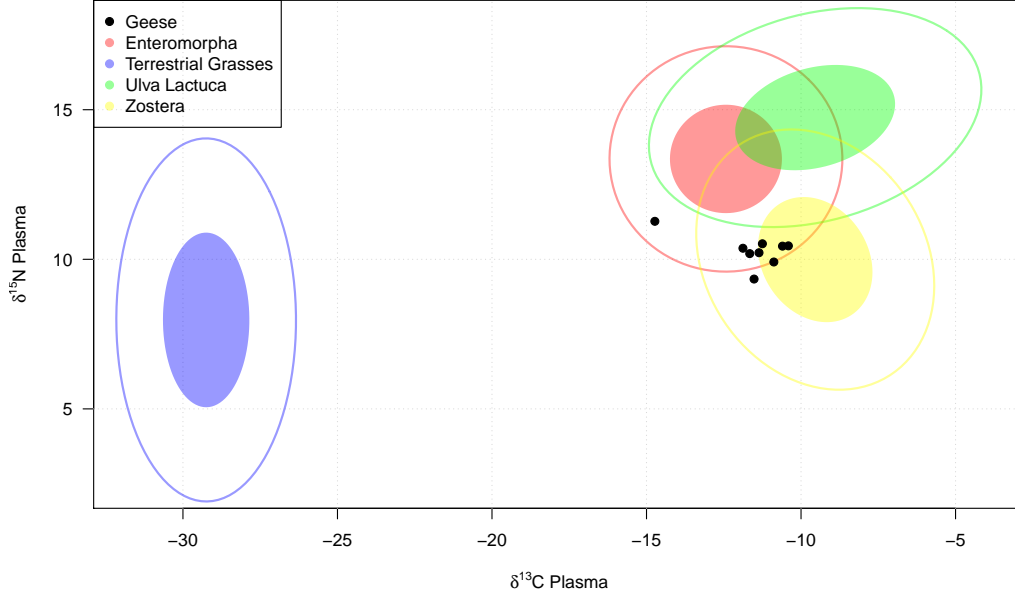


Figure 1: Isospace plot of the geese data of case study 1 (see Section 5.1). The consumer information is shown in the black filled circles, whilst the sources are shown as contours (90% range) and filled ellipses (50% range). The consumers seem to lie close to the *Zostera* source so this is likely to form a substantial part of their diet.

- 80 •  $s_{ijk}$  is the source value for consumer  $i$  on isotope  $j$  and source  $k$ . We write  $\mathbf{s}_{ik}$  to be the  $J$ -vector of  
81 isotope source values for consumer  $i$  on source  $k$ , and  $\mathbf{s}_i$  to be the  $J \times K$  matrix of source values for  
82 consumer  $i$ .
- 83 •  $c_{ijk}$  is the TEF value for consumer  $i$  on isotope  $j$  and source  $k$ . We write  $\mathbf{c}_{ik}$  to be the  $J$ -vector of TEF  
84 values for consumer  $i$  on source  $k$ , and  $\mathbf{c}_i$  as the  $J \times K$  matrix of TEF values related to consumer  $i$ .
- 85 •  $p_{ik}$  is the dietary contribution of source  $k$  for consumer  $i$ .  $\mathbf{p}_i$  is the  $K$ -vector of dietary proportions for  
86 consumer  $i$ . Estimation of these dietary proportions is the main focus of our analysis.
- 87 •  $\epsilon_{ijk}$  is a random noise term representing residual variation. We write  $\boldsymbol{\epsilon}_i$  as the  $J$ -vector of residual  
88 terms for consumer  $i$ , and set  $\boldsymbol{\epsilon}_i \sim N(0, \boldsymbol{\Sigma})$  with  $\boldsymbol{\Sigma}$  a covariance matrix.

89 Using the notation above, we write our general model as:

$$\mathbf{Y}_i \sim N(\mathbf{p}_i^T (\mathbf{s}_i + \mathbf{c}_i), \boldsymbol{\Sigma}). \quad (1)$$

90 We assume a hierarchical formulation so that  $\mathbf{s}_{ik} \sim N(\boldsymbol{\mu}_k^s, \boldsymbol{\Sigma}_k^s)$  and  $\mathbf{c}_{ik} \sim N(\boldsymbol{\mu}_k^c, \boldsymbol{\Sigma}_k^c)$  where the means and  
91 covariances of the sources and TEFs are estimated from the source and TEF experimental data.

92  
93 Of particular interest is the modelling structure for the dietary proportions  $\mathbf{p}$ . We use an isometric log-ratio  
94 (ilr) approach as proposed by Egozcue et al. (2003), though other transformations are available (see next  
95 two sections for further discussion). The transformation is written as:

$$\boldsymbol{\phi}_i = \text{ilr}(\mathbf{p}_i) = \mathbf{V}^T \log \left[ \frac{p_{i1}}{g(\mathbf{p}_i)}, \dots, \frac{p_{iK}}{g(\mathbf{p}_i)} \right] \text{ with } g(\mathbf{p}_i) = \left( \prod_{i=1}^K p_{ik} \right)^{1/K}. \quad (2)$$

96 with  $\mathbf{V}$  a  $K - 1 \times K$  matrix of orthonormal basis functions on the simplex. The inverse transformation  
 97  $\mathbf{p}_i = \text{ilr}^{-1}(\phi_i)$  simply involves exponentiating and re-normalising the values. There are two consequences  
 98 of working with the ilr. The first is that we now work in a  $K - 1$  dimensional space. The second is that  
 99 there is no obvious link between the elements of  $\phi_{ik}$  and  $p_{ik}$ , so we lose some degree of interpretability. We  
 100 further parameterise the transformed proportions so that  $\phi_{ik} \sim N(\gamma_{ik}, \kappa_k)$  with  $\kappa_k$  quantifying a random  
 101 effect variance. Were we to use the centred log-ratio (clr) transform (equivalent to setting  $\mathbf{V} = \mathbf{I}$ ; though  
 102 see Section 4 as to why we might not do this) then  $\kappa_k$  would represent the consumer-level variance in diet  
 103 for source  $k$ . Finally we set  $\gamma_{ik}$  to be a mean term restricted in some fashion to allow for estimation. A  
 104 multivariate prior for  $\phi_i$  would also be feasible.

105

106 In situations where covariates  $\mathbf{x}_i$  are available, they are usually linked to the model through the dietary  
 107 proportions. The covariates may take the form of age, sex, time or any other variables upon which diet  
 108 is expected to depend. In certain cases (such as our third case study) both the diet and the sources are  
 109 expected to be functions of a covariate. In the simpler case we apply the covariates by making  $\gamma_{ik}$  functions  
 110 of  $\mathbf{x}_i$ . In the more advanced case we additionally apply them to  $\boldsymbol{\mu}_k^s$  and  $\boldsymbol{\Sigma}_k^s$ .

111

112 We term the TEF data set  $\mathcal{D}_c$  and the source data set  $\mathcal{D}_s$ , each consisting of collections of  $J$ -vectors of  
 113 isotope ratios for each of the  $K$  sources. A full Bayesian posterior distribution gives:

$$\begin{aligned}
 \pi(\mathbf{p}, \phi, \gamma, \boldsymbol{\kappa}, \boldsymbol{\Sigma}, \boldsymbol{\Sigma}_s, \boldsymbol{\Sigma}_c, \boldsymbol{\mu}_s, \boldsymbol{\mu}_c, \mathbf{s}, \mathbf{c} | \mathbf{Y}, \mathbf{x}, \mathcal{D}_c, \mathcal{D}_s) &\propto \left[ \prod_{i=1}^N \pi(\mathbf{Y}_i | \mathbf{p}_i, \mathbf{s}_i, \mathbf{c}_i, \boldsymbol{\Sigma}) \right] \times \left[ \prod_{i=1}^N \pi(\phi_i | \gamma_i, \mathbf{x}_i, \boldsymbol{\kappa}) \right] \\
 &\times \left[ \prod_{i=1}^N \prod_{k=1}^K \pi(\mathbf{s}_{ik} | \boldsymbol{\mu}_k^s, \boldsymbol{\Sigma}_k^s) \right] \times \left[ \prod_{i=1}^N \prod_{k=1}^K \pi(\mathbf{c}_{ik} | \boldsymbol{\mu}_k^c, \boldsymbol{\Sigma}_k^c) \right] \\
 &\times \left[ \prod_{k=1}^K \pi(\boldsymbol{\mu}_k^s, \boldsymbol{\Sigma}_k^s | \mathcal{D}_s) \right] \times \left[ \prod_{k=1}^K \pi(\boldsymbol{\mu}_k^c, \boldsymbol{\Sigma}_k^c | \mathcal{D}_c) \right] \\
 &\times \left[ \prod_{k=1}^K \pi(\boldsymbol{\kappa}_k) \right] \times \pi(\boldsymbol{\Sigma}) \tag{3}
 \end{aligned}$$

114 A Directed Acyclic Graph (DAG) is shown in Figure 2. The sub-models for the sources and TEFs, namely  
 115  $\pi(\boldsymbol{\mu}_k^s, \boldsymbol{\Sigma}_k^s | \mathcal{D}_s)$  and  $\pi(\boldsymbol{\mu}_k^c, \boldsymbol{\Sigma}_k^c | \mathcal{D}_c)$  can be updated as part of the modelling steps, but we prefer a more efficient  
 116 empirical Bayes approach (Carlin and Louis, 2000) whereby we approximate  $\boldsymbol{\mu}_k^s, \boldsymbol{\mu}_k^c, \boldsymbol{\Sigma}_k^s, \boldsymbol{\Sigma}_k^c$  by their sample  
 117 estimates, thus cutting feedback and removing these from the updates and the posterior. Note that this has  
 118 minimal effect on the source and TEF random effects  $s_{ik}$  and  $c_{ik}$  which are still updated as part of the mod-  
 119 elling process. The prior distributions we give to  $\boldsymbol{\Sigma}$  and  $\boldsymbol{\kappa}_k$  are vague Inverse-Wishart and Inverse-Gamma  
 120 respectively.

121

### 122 3 Previous work

123 The earliest attempts at applying a mixing model framework to stable isotope data did not use probability  
 124 as the basis for estimation. Initially, SIMMs were restricted to systems involving a single consumer (or  
 125 the mean of multiple consumers), and where the number of isotopes and sources was arranged such that  
 126  $J + 1 = K$ . Such an arrangement yields a linear system with a single solution. Phillips and Gregg (2001)  
 127 provided propagation of error calculations for such a system in their IsoError model to establish confidence  
 128 intervals around the estimates based on the variances of the consumer and source isotopic measurements. The  
 129 method was expanded upon in IsoSource (Phillips and Gregg, 2003) to relax the  $J + 1 = K$  restriction and  
 130 allow for multiple sources but without explicit incorporation of source and consumer variability. IsoSource  
 131 works by simulating values of the dietary proportions  $p$  on a grid to produce multiple valid solutions to the

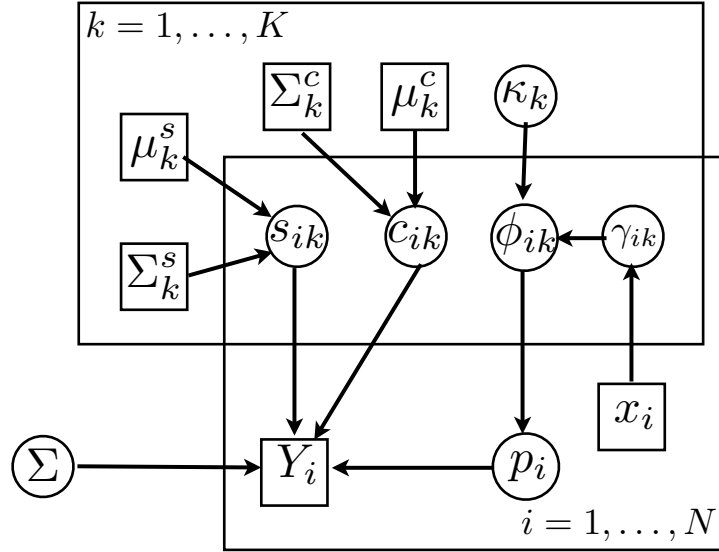


Figure 2: A Directed Acyclic Graph (DAG) of our model. Circles indicate parameters to be estimated whilst squares indicate data. The arrows indicate the direction of information flow.

132 linear system. The valid solutions could be plotted in a histogram-like fashion, though they did not represent  
 133 probability distributions, rather simply the range of values which might be plausible given the geometry of  
 134 the system.

135

136 These initial attempts were formalised in a Bayesian fashion in the models MixSIR (Moore and Semmens,  
 137 2008) and SIAR (Parnell et al., 2008). The MixSIR model can be thought of as a simplification of Equation  
 138 1 without explicit random effects across the dietary proportions, and with  $\Sigma$  set to zero. The sources and  
 139 TEFs are treated as independent across isotopes and are given fixed values for their mean and variance.  
 140 The dietary proportions are given independent Beta-distributed priors or, in a later version, a Dirichlet  
 141 distribution. Since the dietary proportions are the only parameters in the model, it can be fitted extremely  
 142 efficiently using Importance Resampling (e.g. Robert and Casella, 2005) on a grid encompassing the range  
 143 of proportion values. Updated versions of the MixSIR model have included random effects in the dietary  
 144 proportions through the clr transform, and also have allowed for hierarchical models to be fitted, most el-  
 145 egantly in capturing familial relationships affecting the diet of gray wolves in British Columbia (Semmens  
 146 et al., 2009). These latter advanced versions of MixSIR are fitted using the JAGS software (Plummer, 2003).

147

148 The SIAR model (Parnell et al., 2010) is in many ways similar to the basic MixSIR model (and thus still a  
 149 simplification of Equation 1) though includes a residual component which is treated as independent between  
 150 isotopes and given an Inverse Gamma prior. The model also allows for concentration dependencies ( $q_{kj}$ ,  
 151 vectorised as  $\mathbf{q}$ ) which quantify the proportion of the isotope in the given food source (Phillips and Koch,  
 152 2002). They can be added in to our model by replacing  $\mathbf{p}_i$  in Equation 1 with  $\mathcal{C}(\mathbf{p}_i \oplus \mathbf{q})$  where  $\mathcal{C}$  is the  
 153 simplex closure operator and  $\oplus$  is the simplex perturbation (Egozcue et al., 2003). Usually  $\mathbf{q}$  are either fixed  
 154 or given a suitably informative prior distribution. The SIAR model is fitted using standard MCMC with  
 155 Metropolis-Hastings steps to update the dietary proportions.

156

157 More recently, the IsotopeR model has been introduced which extends the SIAR/MixSIR models to a mul-  
158 tivariate setting (both sources and TEFs are multivariate normal) and partitions the residual covariance  $\Sigma$   
159 into a mass spectrometer calibration error and that of residual error. They further allow for the sources,  
160 TEFs and dietary proportions to be random effects with the latter obtained through the clr transform (see  
161 next section for further discussion). The model is fitted in JAGS (Plummer, 2003) but does not allow for  
162 covariate information or for the estimation of the dietary random effect variance.

163

164 Aside from explicit SIMM model development, recent focus has also been on the performance of SIMMs in  
165 non-ideal conditions, most notably with respect to the characterisation of source and TEF values by Bond  
166 and Diamond (2011). Clearly it is absolutely vital that food sources are not excluded from the model as they  
167 will yield biased dietary proportions. Similarly, the estimation of the TEF values must be conducted appro-  
168 priately or there will be some extra uncertainty in the estimated dietary proportions. In related work, Ward  
169 et al. (2011) considered the problem of (dis)aggregating sources and its effect on the resulting estimates.  
170 For example, if a consumer only eats part of a food source it may be hard to obtain isotopic values from  
171 just that part. Similarly if two sources, though different species, lie in the same location in iso-space it may  
172 be impossible for the model to determine the difference in their dietary consumptions. Thus on occasion it  
173 may be pertinent to aggregate sources without any loss of information. Alternatively the aggregation can be  
174 accomplished with fewer assumptions if it is calculated *a posteriori* by combining the negatively correlated  
175 dietary proportions.

176

177 Lastly, it should be noted that there are strong connections between SIMMs and the ‘end member’ analysis  
178 used by geologists to determine the composition of e.g. river sediment. In such cases the sources are usually  
179 known with minimal error and the challenge is to estimate the proportional contributions of different sedi-  
180 ment sources. A number of different methods have been proposed, e.g. Soulsby et al. (2003); Brewer et al.  
181 (2011) (and references therein). Ours most closely resembles that of Palmer and Douglas (2008), though  
182 they use the alr transformed proportions (see below for definition) which are given a spatial prior distribution.

183

## 184 4 Statistical issues in SIMMs

185 The models we fit use ideas from Bayesian hierarchical modelling (e.g. Gelman et al., 2003) and composi-  
186 tional data analysis (Aitchison, 1986). BHM is now part of the standard toolbox of Bayesian statistics and  
187 we do not discuss them further here. Of more interest however is the compositional structure applied to the  
188 dietary proportions as this can strongly affect the behaviour of the posterior distribution. A recent review  
189 of the state of compositional data analysis can be found in Pawlowsky-Glahn and Buccianti (2011). Our  
190 models differ fundamentally from many standard problems as our compositions are not observed directly  
191 but are latent parameters to be estimated, constrained by their geometrical position in iso-space and the  
192 covariates upon which they may depend.

193

194 The starting point for the early Bayesian SIMMs models was that of the Dirichlet distribution, being per-  
195 haps the simplest valid distribution on the simplex. The usual prior distributions used were either flat  
196 where all Dirichlet parameters are set to 1 or the Jeffreys prior where all are set to  $1/K$ . Unfortunately, as  
197 is well known (e.g. Aitchison, 1986), the Dirichlet distribution suffers from a very rigid sub-compositional  
198 independence assumption. This is not necessarily a problem when used as a prior distribution with fixed  
199 components as the posterior distribution may well show interesting sub-compositional properties. However,  
200 if dependence is to be modelled through hyper-parameters of the Dirichlet distribution (e.g. with covariates)  
201 this restriction will remain in the posterior.

202

203 A number of extensions to the Dirichlet have been proposed (e.g. Wong, 1998), but we focus here on the  
204 logistic-normal transformations of Aitchison (1986) and Egozcue et al. (2003), through which more flexible  
205 sub-compositional dependence can be obtained. The simplest of these is perhaps the additive log-ratio (alr),

206 where  $g(\mathbf{p}_i)$  in Equation 2 is set to be one of the chosen proportions (which is then removed from the com-  
 207 position) and  $\mathbf{V} = \mathbf{I}$ . However, this can perform poorly when there is no obvious choice of denominator  
 208 and is not permutation-invariant. The centred log-ratio (clr; defined in Equation 2 when  $\mathbf{V}=\mathbf{I}$ ) removes  
 209 the need for a choice of denominator in the log ratio but produces a covariate matrix of rank  $K - 1$ . It is  
 210 also not sub-compositionally coherent (see Pawlowsky-Glahn and Buccianti, 2011, for further discussion of  
 211 these terms). Finally the isometric log-ratio (ilr) uses orthonormal basis functions in the simplex to obtain  
 212 coordinates that are isometric and satisfy the usual compositional requirements of coherence and permuta-  
 213 tion/subcomposition invariance. The choice of basis functions is somewhat subjective; we follow the method  
 214 in Egozcue et al. (2003). From a Bayesian perspective, the latent compositional parameters are more easily  
 215 identifiable when working with the ilr.

216  
 217 Once a suitable transform has been chosen, it is feasible to include covariates or perform any of the traditional  
 218 multivariate analysis techniques. Numerous examples can be found, including a geo-statistical framework  
 219 (Tolosana-Delgado et al., 2011), discrete time series (Barceló-Vidal et al., 2011), or spectral methods (Pardo-  
 220 Igúzquiza and Heredia, 2011). A popular topic is that of zero compositions or zero-inflation (e.g. Butler and  
 221 Glasbey, 2008) which is not a severe issue in SIMMs because we know from experimental observation that  
 222 all dietary sources are consumed.

## 224 5 Case studies

225 We now present three case studies and show how our general model can be used in each of the scenarios.  
 226 In the first case study we analyse the diet of a small sample of geese from data previously studied by Inger  
 227 et al. (2006). This first case study uses the model as proposed in Section 2, and can be seen as a small  
 228 extension to that of Hopkins and Ferguson (2012). The second case study extends the geese model to allow  
 229 for the inclusion of covariates or basis function models. Our final case study includes a compositional time  
 230 series component where the sources and consumers are observed at different time points. In this final case  
 231 study the data are swallows consuming chironomid midges, other freshwater invertebrates and terrestrial  
 232 invertebrates. In all cases we run the models using the JAGS software and check convergence using the  
 233 coda package (Plummer et al., 2006) and the Brooks-Gelman-Rubin diagnostic (Gelman and Rubin., 1992;  
 234 Brooks and Gelman, 1998).

### 236 5.1 Case study 1

237 Our first case study contains 9 ( $\delta^{13}\text{C}$ ,  $\delta^{15}\text{N}$ ) pairs of isotopic values taken from the blood plasma of Brent  
 238 geese sampled on 26th October 2003. The food sources are *Zostera* spp, terrestrial grasses, *Ulva lactuca* and  
 239 *Enteromorpha* spp. Our empirical Bayes estimates for the sources are given in Table 1. The TEFs were taken  
 240 from values in the literature (see references in Inger et al., 2006, for more details) so that  $\boldsymbol{\mu}_k^c = (1.63, 3.54)^T$   
 241 and  $\boldsymbol{\Sigma}_k^c = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$  for all  $k$ . For a simple random effects formulation we set  $\gamma_{ik} = 0$  over all  $i$  and  $k$ . An iso-  
 242 space plot for the data (where the sources have been corrected by the TEFs) is shown in Figure 1. The JAGS  
 243 model was run for 3 chains over 50,000 iterations, removing 10,000 for burn-in and thinning by a factor of 20.

244  
 245 A density plot of the mean dietary proportions is shown in the left panel of Figure 3. They can be seen to  
 246 compare favourably with the simpler SIAR model (see the function `siardemo` in Parnell et al., 2008), though  
 247 in this case we have extra information covering individual dietary estimates, as well as improved estimation  
 248 of the source and TEF random effects. In particular, the flexibility of the hierarchical formulation through  
 249 the ilr transform allows for some multi-modality in the posterior distributions. The right panel of Figure 3  
 250 shows a matrix plot of the joint behaviour of the dietary proportions. These can be useful in determining  
 251 unavoidable model inadequacy, for example when it is impossible to ascertain which food sources are being

| Source       | <i>Enteromorpha</i>                                  | Terr Grasses   | <i>Ulva lactuca</i>  | <i>Zostera</i>   |
|--------------|--|--|--|--|
| $\mu_s^k$    | $(-14.06, 9.82)^T$                                   | $(-30.88, 4.43)^T$                                   | $(-11.17, 11.2)^T$   | $(-11.17, 6.45)^T$   |
| $\Sigma_s^k$ | $\begin{bmatrix} 1.37 & 0 \\ 0 & 1.37 \end{bmatrix}$ | $\begin{bmatrix} 0.41 & 0 \\ 0 & 5.15 \end{bmatrix}$ | $\begin{bmatrix} 3.83 & 0.85 \\ 0.85 & 1.24 \end{bmatrix}$ | $\begin{bmatrix} 1.48 & -0.56 \\ -0.56 & 2.16 \end{bmatrix}$ |

Table 1: Estimates of the source means and covariance matrices for the geese data of case study 1.

252 consumed together. A strong negative correlation indicates that the food sources are indistinguishable. For  
 253 example the strong negative correlation between *Zostera* and *Ulva lactuca* indicates that, whilst it is clear  
 254 they are consuming mainly *Zostera*, the balance between the two cannot be exactly determined.  
 255

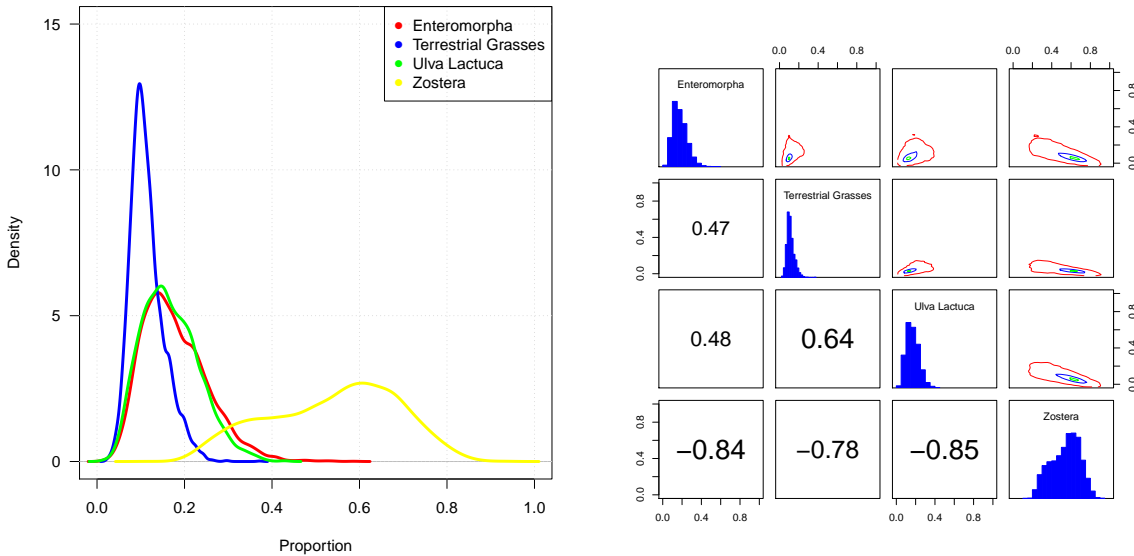


Figure 3: Left panel: Density plot of mean dietary proportions for the geese data of case study 1. Right panel: matrix plot of the posterior dietary proportions obtained from the geese data. The upper-diagonal shows a contour plot, the diagonal a histogram, and the below-diagonal the correlation between the different sources.

256

257 There are many other useful statistics we can calculate here quite simply from the posterior distributions.  
 258 In particular, we can focus on individual level variation by calculating, for example, the probability that an  
 259 individual consumes more *Zostera* than another. Similarly with access to the variance parameters  $\kappa_k$  we can  
 260 determine whether there is more variation amongst consumers in some sources rather than others. Lastly, it  
 261 is often desirable to perform model comparison diagnostics to determine whether certain parts of the model  
 262 can be removed without a detrimental effect on prediction. However, we do not perform any further analysis  
 263 on this data set, preferring instead to use these tools when analysing the more sophisticated data below.  
 264

264

## 265 5.2 Case study 2

266 We now extend our Goose model to a larger data set of 248 observations collected over the period of October  
 267 2003 to April 2005 of which the previous case study was just a small part. The sources and TEFs are believed



268 to be stable over the course of the study so the source means and covariances are the same as Table 1. The  
 269 diet of the geese will however vary during the season due to variations in abundance of food sources along  
 270 with social and demographic factors. An iso-space plot for the full data is shown in Figure 4. The iso-space  
 271 plot appears to show that the diet during October to be focussed mainly on *Zostera*, moving on to *Ulva*  
 272 *lactuca* and/or *Enteromorpha* during November/December. In January and February the diet appears to  
 273 be relatively mixed, but focussing almost solely on Terrestrial grasses around April. In addition to the sta-  
 274 ble isotope information we have covariates which state the Goose's sex and whether they are juvenile or adult.  
 275

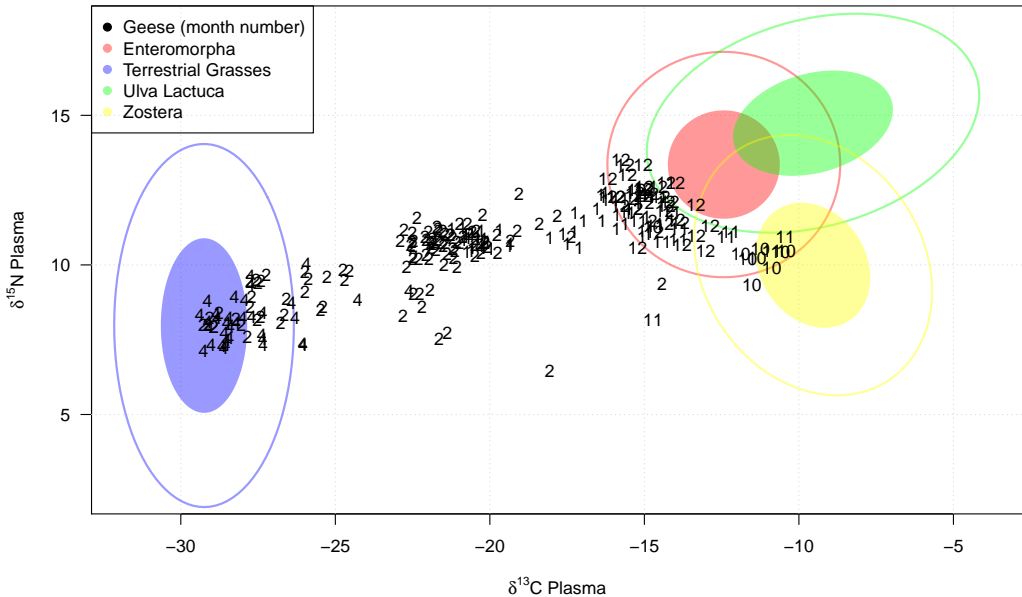


Figure 4: Isospace plot of the full goose data set where the consumers are labelled by month. The data in case study 1 correspond to the data shown for month 10. Note that the sources and the TEFs are unchanged from Figure 1.

276 We consider 6 possible models for the dietary behaviour, accounting for the covariates. In each case we set  
 277  $\gamma_{ik} = \mathbf{X}_i^T \beta_k$  where  $\mathbf{X}_i$  is an  $L$ -vector of covariates or basis functions for consumer  $i$ , and  $\beta_k$  is an  $L$ -vector  
 278 of regression parameters associated with  $\phi_k$ . Note that, when using the ilr transform, the parameters  $\beta_k$   
 279 do not have any association with source  $k$ . With the extra parameters in the model, convergence of the  
 280 MCMC algorithm is a problem because the parameters are often highly correlated across sources, though  
 281 this is somewhat lessened with appropriate choice of  $\mathbf{V}$  in the ilr. We find it helpful to re-parameterise  
 282 with Helmert contrasts across sources so that  $\gamma_{i1} = \mathbf{X}_i^T \beta_1$  and  $\gamma_{ik} = \mathbf{X}_i^T (\beta_1 + \beta_k)$  for  $k \geq 2$ . We use the  
 283 first source (*Enteromorpha*) as  $k = 1$  as this seems to be consumed throughout the season though this is  
 284 obviously not known in advance so does involve some degree of trial and error.

285  
 286 In all cases the parameters  $\beta$  are given vaguely informative  $N(0, 10)$  distributions as a value of  $|\phi|$  in excess  
 287 of 10 is likely to yield dietary proportions near 100% (though again this depends on the correlation structure  
 288 in the data set). We compare the different models using the Deviance Information Criterion (DIC; Spiegel-  
 289 halter et al., 2002) and, for the final chosen model, posterior predictive distributions of the data. We use two  
 290 versions of the DIC, the standard  $p_D$  method of Spiegelhalter et al. (2002), and the  $p_V$  method of Plummer  
 291 (2008) which estimates the optimism as half the variance of the deviance, and thus penalises complex models  
 292 more harshly. Due to the extra complexity in the models, we run them with 3 chains for 200,000 iterations,

293 removing 20,000 for burn-in and thinning by 90.

294

295 The first model we try involves no covariates and is thus the same as that in Section 5.1. The second model  
 296 includes an intercept term and time as a simple linear covariate. The third model replaces linear time with  
 297 a single harmonic component so that  $X_i^T = [1, \cos(\frac{2\pi t_i}{365}), \sin(\frac{2\pi t_i}{365})]$  where  $t_i$  is the Julian day. The fourth,  
 298 fifth and sixth models are expansions of model 3 to include juvenile/adulthood (model 4), sex (model 5) and  
 299 also their interaction (model 6) as covariates. Table 2 shows the different models and the associated DIC  
 300 values. Both versions of DIC seem to prefer models with the harmonic covariate, and also the addition of  
 301 either sex or adulthood.

302

303 Figure 5 shows the relationship between Julian day and dietary proportion for the different sources for model  
 304 4. The switch from *Zostera*, to *Enteromorpha*, to terrestrial grasses is very clearly seen in both juveniles  
 305 and adults, though the uncertainty in juveniles is slightly higher, especially with respect to *Enteromorpha*.  
 306 Figure 6 shows the posterior predictive distribution of the data under model 4. The model seems to pre-  
 307 dict the data well, the three modes corresponding to the main sampling times of October, February and April.  
 308

| Model    | Covariate(s)  | DIC (using $p_V$ ) | DIC (using $p_D$ ) |
|----------|---|--------------------|--------------------|
| 1        | None  | 26907              | 1210.0             |
| 2        | Julian day (linear)                                     | 16957              | 524.3              |
| 3        | Julian day (harmonic)                                   | 16683              | 385.1              |
| <b>4</b> | <b>Julian day (harmonic), Juvenile/Adult</b>            | <b>7551</b>        | 393.6              |
| <b>5</b> | <b>Julian day (harmonic), Sex</b>                       | 8600               | <b>382.8</b>       |
| 6        | Julian day (harmonic), Juvenile/Adult, Sex, Interaction | 10812              | 382.9              |

Table 2: Table of models and model selection criteria. The models with the lowest DIC are shown in bold.

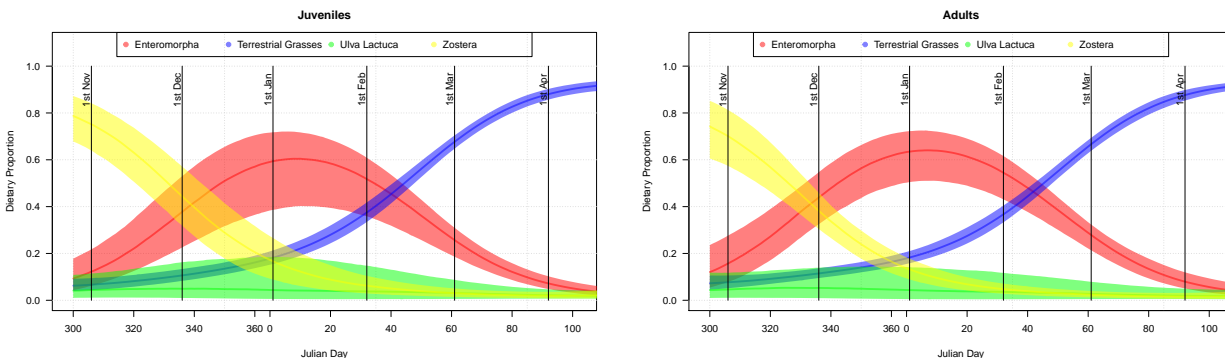


Figure 5: Plot of proportion values against Julian day for model 3. The left panel shows estimates for juveniles whilst the right panel shows adults. The solid lines show median estimates, the outer of the polygons show 90% credibility intervals. The Geese appear to focus on *Zostera* around November before moving on to *Enteromorpha* and then terrestrial grasses. There is slightly more uncertainty in the juveniles than the adults.

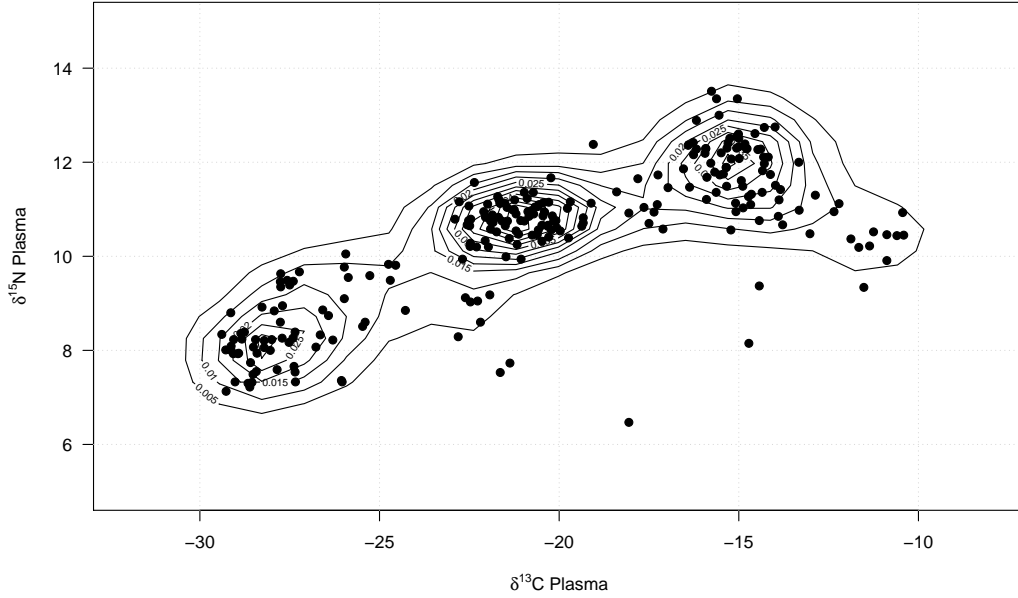


Figure 6: Plot of the predictive distribution of the data from chosen model. The observations are shown as filled circles, whilst the posterior predictive density is represented by the contours.

### 309 5.3 Case study 3

310 Our last case study concerns the dietary behaviour of barn swallows (unpublished data). The data are stable  
311 isotope ratios of blood plasma samples from birds captured between May and August in 2009. Again, we  
312 use Julian day to determine their behaviour over time. However, in this scenario the sources (chironomid  
313 midges, other freshwater invertebrates, and terrestrial invertebrates) are also expected to change over time,  
314 and may have been observed on different days to that of the swallows. We thus expand our model to include  
315 a temporal component on the source vector as well as that of the dietary proportions. In each case, we use  
316 P-splines (Eilers and Marx, 1996) to explain the suitably flexible behaviour. Other time series methods,  
317 such as random walks or Levy processes, may also be appropriate. We write continuous time as  $t_i$  so that  
318 the consumers are now  $\mathbf{Y}(t_i)$ . We set  $\gamma_k(t_i) = \mathbf{X}_i^T \beta_k$  where now  $\mathbf{X}_i$  is an  $L$ -vector of cubic B-spline basis  
319 functions evaluated at time point  $t_i$  and  $\beta_k$  are weights for each basis function on source  $k$ . The P-spline  
320 formulation is completed by giving a random walk prior such that  $\beta_{lk} - \beta_{l-1,k} \sim N(0, \tau_k^{-1})$  where  $\tau_k$  is a  
321 roughness parameter associated with source  $k$  and given a weakly informative  $Ga(2, 1)$  prior.  
322

323 The sources are now described by a multivariate spline model so that source data pairs, denoted  $s'_{jk}(t)$  for the  
324 source experimental data at time  $t$  on isotope  $j$  and source  $k$ , are given  $N\left(\left[\mathbf{X}^T \beta'_1, \dots, \mathbf{X}^T \beta'_J\right]^T, \Sigma'(t)\right)$   
325 independently for each source  $k$ . The number of observations for each source is likely to be different, and  
326 certainly not equal to the number of consumer observations  $N$ . Here, the spline parameters  $\beta'_j$  determine  
327 the mean behaviour of the sources over time on isotope  $j$ . The  $J \times J$  variance matrix  $\Sigma'(t)$  is also allowed to  
328 change over time with diagonal elements given log-splines:  $\log(\Sigma'_{jj}(t)) \sim N(\mathbf{X}^T \beta_{\Sigma}, \kappa_{\Sigma})$ . The cross isotope  
329 covariance  $\Sigma'_{12}$  is parameterised through a single correlation parameter for each source, denoted  $\rho_k$ , and does  
330 not change over time. A spline could also be used here (for example on the arctangent of  $\rho_k$ ) but this was  
331 not found to improve the fit. The source spline model was run for each of the three sources in turn and used  
332 to calculate Maximum a Posteriori (MAP) estimates of  $\mu_k^s(t)$  and  $\Sigma_k^s(t)$  for all times  $t$  at which consumer

333 data were available. An iso-space plot of the swallows data is shown in Figure 7.

334

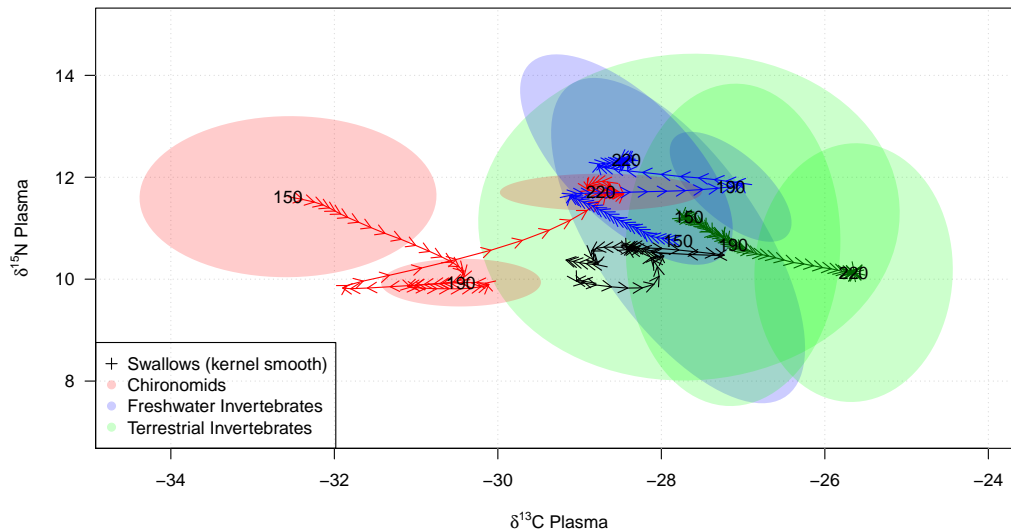


Figure 7: Isospace plot of the swallows data. The source spline model has been run to obtain estimates of the source means and covariances throughout the study period. Arrows indicate the direction of movement over time for the sources and the swallows. The Julian day and the 50% standard ellipses are given for Julian days 150, 190 and 220. Note that the chironomids'  $\delta^{13}\text{C}$  values increase over time from the start of the study period. A similar occurrence can be seen in the terrestrial invertebrates. The data are shown as kernel-smoothed estimates of the swallows' isotope data again over Julian day, starting at 150 and ending on day 220. The swallows can be seen to also increase their  $\delta^{13}\text{C}$  values up until around day 180 whereupon the  $\delta^{13}\text{C}$  returns towards its original value. This plot should be read in conjunction with Figure 8.

335 We use the source model predictions in our standard SIMM with the spline formulation on the proportions as  
336 outlined above. An alternative model where all parameters are estimated simultaneously would be possible,  
337 if rather slow. Instead, we retain the cutting feedback assumption so that the source and dietary proportions  
338 are run separately. For both the source and the dietary proportion models we run for 200,000 iterations,  
339 removing 20,000 for burn-in and thinning by 90. For both models, we use 25 knots, which seems to cover  
340 the flexibility in the data adequately.

341

342 Figure 8 shows the posterior dietary proportion estimates over Julian day for the swallows, predicted from  
343 the resulting spline parameter estimates. The results clearly indicate that they are feeding on mainly fresh  
344 water invertebrates during the early part of the data before concentrating on chironomids around the start of  
345 August. Figure 9 shows the predictive distribution of the data under this model. The fit appears satisfactory.

346

## 347 6 Discussion and future directions

348 The SIMM formulation outlined in Section 2 allows for a rich framework upon which to include a variety  
349 of other statistical structures. The basis of such models is a mixture compositional structure applied to the

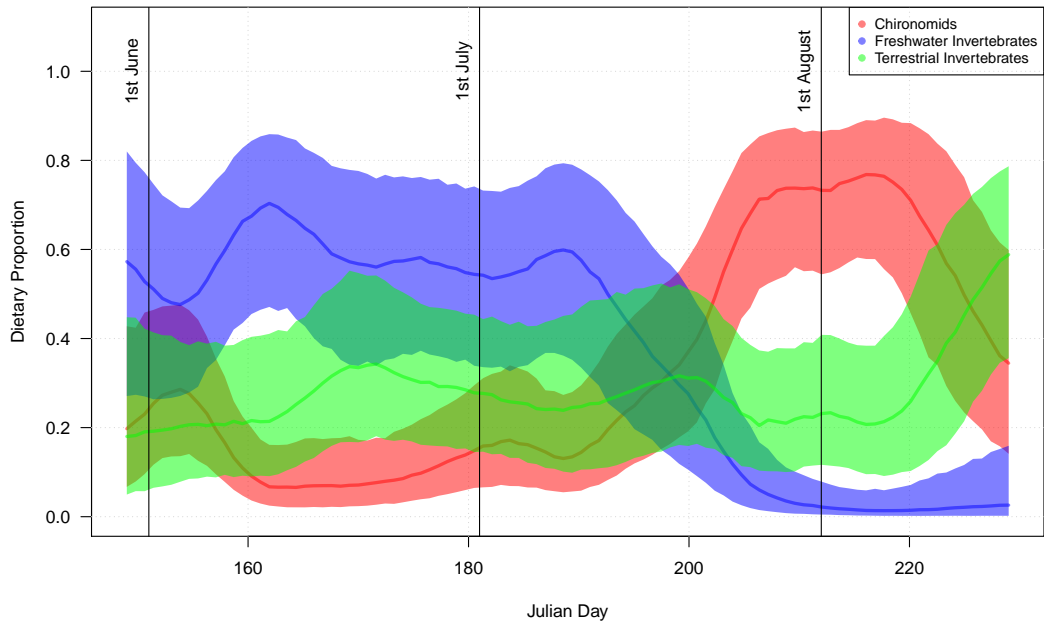


Figure 8: Plot of dietary proportion values against Julian day for the swallows data of Case Study 3. The solid lines show median estimates, whilst the filled polygons show 90% credibility intervals for each source.

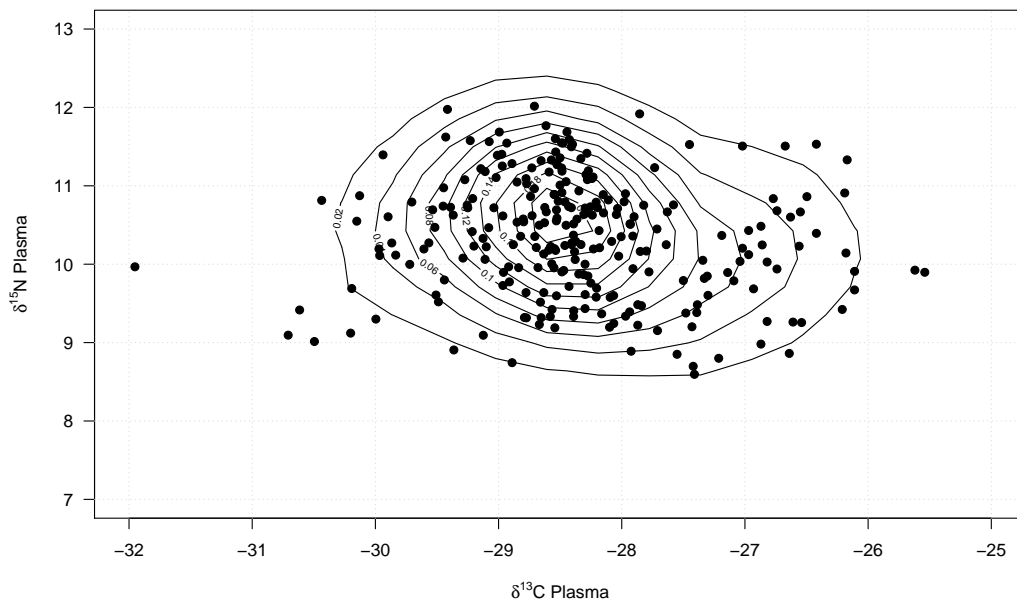


Figure 9: Predictive distribution of the data for the swallows example. The observations are shown as filled circles, whilst the posterior predictive density is represented by the contours.

350 dietary sources. In the case studies above we have illustrated how some simple regression and smoothing

351 models might be included. The results seem useful and allow for many interesting findings which would not  
352 have otherwise been possible without, e.g. the time series or spline components.

353

354 The main challenges in such models are that the source and TEF values are fully and correctly characterised.  
355 It is a simple geometrical exercise to verify that if sources are missing from the data set, or that the TEF  
356 means and errors are poorly estimated, the dietary proportion estimates will be biased. This bias is com-  
357 pounded by the compositional nature of the problem where, if one dietary proportion is poorly estimated,  
358 then the others may well be too. There is no statistical test for missing sources; we rely on the ability of the  
359 ecologist to observe the system adequately. This may be particularly important if multiple organisms and  
360 sources are to be analysed simultaneously in a dietary network. A single missing source may cause biases  
361 across the dietary proportion estimates. Such models are not to be encouraged unless there is very strong  
362 evidence that no sources are missing and that TEFs are estimated correctly.

363

364 It is reasonably straightforward to expand our approach for other values of  $J$  where differing numbers of  
365 isotopes may be available for analysis. It should be noted, however, that if  $J \geq 3$  it becomes harder to  
366 determine model fit, especially as no obvious iso-space plot can be created. The solutions to this problem  
367 may lie in closer scrutiny of predictive distributions, and some subjective judgement over the size of the  
368 residual covariance matrix  $\Sigma$ . In such scenarios extra care is required in reporting the output of the SIMM.

369

370 There are further opportunities for expansion of the models. Some of these may include:

- 371 • The simultaneous inclusion of multiple tissue varieties. It is often the case that different tissues are  
372 sampled on the same consumer as they are captured. These differentially replenished tissues will  
373 represent the dietary proportions consumed over different periods. For example, whilst blood plasma  
374 might represent the immediately sampled food sources, feathers might represent the diet consumed over  
375 the previous few months. A long-term data set where multiple tissues are analysed simultaneously may  
376 allow for increased precision in the dietary proportions. Alternatively, it may be possible to estimate  
377 the time scale over which the tissues are responding.
- 378 • Clustering/Mixture models to determine groupings. If there are hidden groupings amongst the organ-  
379 isms it may be possible to discern them using a model-based clustering approach (sensu Fraley and  
380 Raftery, 2002). Even without such groupings, increased flexibility can be obtained by using mixtures  
381 of Gaussian distributions to model non-parametric behaviour.
- 382 • Long-tailed multivariate distributions to account for outliers or small sample sizes. Clearly where  
383 sample sizes are small the multivariate Gaussian assumption (especially for sources) may be invalid,  
384 and thus heavier-tailed distributions may be required. One such which seems to be most easily fitted  
385 is the Multivariate Normal-Inverse Gaussian (MNIG Barndorff-Nielsen, 1997) which has been used  
386 previously in clustering and financial settings.

387 These are just three of the active areas of research to which SIMMs are being applied by our group. Many  
388 others are likely to appear as the field advances. We hope to report soon on these exciting new developments.

389

## 390 References

- 391 Aitchison, J. (1986). *The statistical analysis of compositional data*. London, UK: Chapman & Hall, Ltd.
- 392 Barceló-Vidal, C., L. Aguilar, and J. A. Martín-Fernández (2011). Compositional VARIMA Time Series.  
393 In Pawlowsky-Glahn, V and A. Buccianti (Eds.), *Compositional Data Analysis: Theory and Applications*.  
394 Chichester: John Wiley & Sons.
- 395 Barndorff-Nielsen, O. E. (1997). Normal Inverse Gaussian Distributions and Stochastic Volatility Modelling.  
396 *Scandinavian Journal of Statistics* 24(1), 1–13.

- 397 Bond, A. L. and A. W. Diamond (2011). Recent Bayesian stable-isotope mixing models are highly sensitive  
398 to variation in discrimination factors. *Ecological Applications* 21(4), 1017–1023.
- 399 Brewer, M. J., D. Tetzlaff, I. A. Malcolm, and C. Soulsby (2011). Source distribution modelling for end-  
400 member mixing in hydrology. *Environmetrics* 22, 921–932.
- 401 Brooks, S. P. and A. Gelman (1998). General Methods for Monitoring Convergence of Iterative Simulations.  
402 *Journal of computational and graphical statistics* 7, 434–455.
- 403 Butler, A. and C. Glasbey (2008). A latent Gaussian model for compositional data with zeros. *Journal of*  
404 *the Royal Statistical Society: Series C (Applied Statistics)* 57(5), 505–520.
- 405 Carlin, B. P. and T. A. Louis (2000). *Bayes and Empirical Bayes Methods for Data Analysis*, Volume 7.  
406 CRC Press.
- 407 Egozcue, J., V. Pawlowsky-Glahn, G. Mateu-Figueras, and C. Barceló-Vidal (2003). Isometric Logratio  
408 Transformations for Compositional Data Analysis. *Mathematical Geology* 35(3), 279–300.
- 409 Eilers, P. H. C. and B. D. Marx (1996). Flexible smoothing with B-splines and penalties. *Statistical*  
410 *Science* 11(2), 89–121.
- 411 Fraley, C. and A. E. Raftery (2002). Model-Based Clustering, Discriminant Analysis, and Density Estimation.  
412 *Journal of the American Statistical Association* 97(458), 611–631.
- 413 Gelman, A., J. B. Carlin, H. S. Stern, and D. B. Rubin (2003). *Bayesian Data Analysis, Second Edition*  
414 (*Chapman & Hall/CRC Texts in Statistical Science*) (2 ed.). Chapman and Hall/CRC.
- 415 Gelman, A. and D. B. Rubin. (1992). Inference from iterative simulation using multiple sequences. *Statistical*  
416 *Science* 2, 457–472.
- 417 Hopkins, J. B. and J. M. Ferguson (2012). Estimating the Diets of Animals Using Stable Isotopes and a  
418 Comprehensive Bayesian Mixing Model. *PLoS ONE* 7(1), e28478.
- 419 Inger, R. and S. Bearhop (2008). Applications of stable isotope analyses to avian ecology. *Ibis* 150, 447–461.
- 420 Inger, R., G. D. Ruxton, J. Newton, K. Colhoun, J. A. Robinson, A. L. Jackson, and S. Bearhop (2006).  
421 Temporal and intrapopulation variation in prey choice of wintering geese determined by stable isotope  
422 analysis. *Journal of Animal Ecology* 75, 1190–1200.
- 423 Moore, J. W. and B. X. Semmens (2008). Incorporating uncertainty and prior information into stable isotope  
424 mixing models. *Ecology Letters* 11(5), 470–480.
- 425 Palmer, M. J. and G. B. Douglas (2008). A Bayesian statistical model for end member analysis of sedi-  
426 ment geochemistry, incorporating spatial dependences. *Journal of the Royal Statistical Society: Series C*  
427 (*Applied Statistics*) 57(3), 313–327.
- 428 Pardo-Igúzquiza, E. and J. Heredia (2011). Spectral Analysis of Compositional Data in Cyclostratigraphy.  
429 In V. Pawlowsky-Glahn and A. Buccianti (Eds.), *Compositional Data Analysis: Theory and Applications*.  
430 Chichester: John Wiley & Sons.
- 431 Parnell, A. C., R. Inger, S. Bearhop, and A. L. Jackson (2008). SIAR: Stable isotope analysis in R:  
432 <http://cran.r-project.org/web/packages/siar/index.html>.
- 433 Parnell, A. C., R. Inger, S. Bearhop, and A. L. Jackson (2010). Source partitioning using stable isotopes:  
434 coping with too much variation. *PLoS ONE* 5(3), e9672.
- 435 Pawlowsky-Glahn, V. and A. Buccianti (2011). *Compositional Data Analysis: Theory and Applications*.  
436 Wiley-Blackwell.

- 437 Phillips, D. L. (2012). Converting isotope values to diet composition: the use of mixing models. *Journal of*  
438 *Mammalogy* 93(2), 342–352.
- 439 Phillips, D. L. and J. W. Gregg (2001). Uncertainty in source partitioning using stable isotopes. *Oecolo-*  
440 *gia* 127, 171–179.
- 441 Phillips, D. L. and J. W. Gregg (2003). Source partitioning using stable isotopes: coping with too many  
442 sources. *Oecologia* 136(2), 261–9.
- 443 Phillips, D. L. and P. L. Koch (2002). Incorporating concentration dependence in stable isotope mixing  
444 models. *Oecologia* 130(1), 114–125.
- 445 Plummer, M. (2003). JAGS: A program for analysis of Bayesian graphical models using Gibbs sampling.
- 446 Plummer, M. (2008). Penalized loss functions for Bayesian model comparison. *Biostatistics (Oxford, Eng-*  
447 *land)* 9(3), 523–39.
- 448 Plummer, M., N. Best, K. Cowles, and K. Vines (2006). CODA: convergence diagnosis and output analysis  
449 for MCMC.
- 450 Robert, C. and G. Casella (2005). *Monte Carlo Statistical Methods (Springer Texts in Statistics)*. Springer.
- 451 Semmens, B. X., E. J. Ward, J. W. Moore, and C. T. Darimont (2009). Quantifying Inter- and Intra-  
452 Population Niche Variability Using Hierarchical Bayesian Stable Isotope Mixing Models. *PLoS ONE* 4(7),  
453 9.
- 454 Soulsby, C., J. Petry, M. Brewer, S. Dunn, B. Ott, and I. Malcolm (2003). Identifying and assessing  
455 uncertainty in hydrological pathways: a novel approach to end member mixing in a Scottish agricultural  
456 catchment. *Journal of Hydrology* 274(1-4), 109–128.
- 457 Spiegelhalter, D. J., N. G. Best, B. P. Carlin, and A. van der Linde (2002). Bayesian measures of model  
458 complexity and fit. *Journal of the Royal Statistical Society (Series B)* 64, 583–639.
- 459 Tolosana-Delgado, R., K. G. van den Boogaart, and V. Pawlowsky-Glahn (2011). Geostatistics for Com-  
460 positions. In V. Pawlowsky-Glahn and A. Buccianti (Eds.), *Compositional Data Analysis: Theory and*  
461 *Applications*, Chapter 6. Chichester, UK: John Wiley & Sons.
- 462 Ward, E. J., B. X. Semmens, D. L. Phillips, J. W. Moore, and N. Bouwes (2011). A quantitative approach  
463 to combine sources in stable isotope mixing models. *Ecosphere* 2(2), art19.
- 464 Wong, T. T. (1998). The generalized Dirichlet distribution in Bayesian analysis. *Applied Mathematics and*  
465 *Computation* 97, 165–181.