

**Title:**

Formal Modelling of Predator Preferences using Molecular Gut-Content Analysis

**Author Information:**

Edward A. Roualdes<sup>1,2</sup>, Simon Bonner<sup>1</sup>, Thomas D. Whitney<sup>3,\*</sup>, and James D. Harwood<sup>3</sup>

1 Department of Statistics  
University of Kentucky  
Rm. 311 Multidisciplinary Science Building  
725 Rose Street  
Lexington, KY 40536-0082  
USA

2 Email: edward.roualdes@uky.edu  
Phone: 530-570-2674

3 Department of Entomology  
University of Kentucky  
Lexington, KY 40546  
USA

\* Current address:  
Warnell School of Forestry and Natural Resources  
University of Georgia  
Athens, GA 30602  
USA

**Short Title:**

Modelling Predator Preferences

**Word Count:**

This manuscript contains 3709 words including figure captions and headers, but excluding title page and equations.

**Keywords**

electivity; expectation-maximization; predator-prey interactions; generalist predators; food web analysis

## Summary

1. The literature on modelling a predator’s prey selection describes many intuitive indices, few of which have both reasonable statistical justification and tractable asymptotic properties.
2. Here, we provide a simple model that meets both of these criteria, while extending previous work to include an array of data from multiple species and time points.
3. Further, we apply the expectation-maximisation algorithm to compute estimates if exact counts of the number of prey species eaten in a particular time period are not observed.
4. We conduct a simulation study to demonstrate the accuracy of our method, and illustrate the utility of the approach for field analysis of predation using a real dataset, collected on wolf spiders using molecular gut-content analysis.

## 1 Introduction

The indices most commonly used to describe a predator’s food preferences, or selectivity, are relatively old (Ivlev, 1964; Jacobs, 1974; Chesson, 1978; Strauss, 1979; Vanderploeg and Scavia, 1979; Chesson, 1983), and yet many applied papers continue to use them. A quick search of papers published in 2014 returns hundreds of publications that cite these fundamental papers, a few being Clements et al. (2014); Hansen and Beauchamp (2014); Hellström et al. (2014); Lyngdoh et al. (2014); Madduppa et al. (2014). These indices, though intuitive, lack the statistical rigour of a full model, focus on a snapshot in time, and rarely allow more than one prey species to be considered (Lechowicz, 1982). We propose an intuitive statistical model to estimate and statistically test differences in a predators’ prey preferences across an array of time points and between multiple prey species.

A comprehensive overview by Lechowicz (1982), later summarised by Manly et al. (1992), details the benefits and faults of the most popular indices. According to these reviews, a majority of the indices give comparable results, save Strauss’s linear index  $L$ , despite the fact that most of the methods differ by range and linearity of response. While Lechowicz recommends one index,  $E^*$  by Vanderploeg and Scavia (1979) as the “single best” (Lechowicz, 1982), albeit imperfect, index, Manly et al. instead take the approach of excluding the subset of indices which do not “estimate any biologically meaningful value” (Manly et al., 1992). Lechowicz (1982) recommends the index  $E^*$ , an element of the Manly et al. (1992) suggested indices, because the index value 0 denotes random feeding, the index has a range restricted to  $[-1, 1]$  (though  $E^* = 1$  is nigh impossible), and the index is based on the predator’s choice of prey as a function of both the availability of the prey as well as the number of available prey types (assumed known). The downside to this index is its lack of reasonable statistical properties (Lechowicz, 1982), thus making the computation of standard errors, and hypothesis testing difficult. This is, in fact, a common fault amongst most of the indices.

To encourage more formal statistical inference, and simultaneously generalise predators’ selectivity to animal resource selection, Manly et al. (1992) proposed the use of generalised

linear models (GLM). The well established literature on GLMs allows for hypothesis testing to replace the indices, by estimating the proportion of eaten relative to the population of prey species, while using environmental variables as predictors. The model we present here, while restricted to predators' preferences, is a compromise between these two extremes, indices and GLMs. Our model offers formal hypothesis testing and inference similar to the GLMs of Manly et al. (1992), but also provides meaningful single number summaries of the predator's dietary preferences. To do this, we estimate the rate at which a predator consumes the prey of interest instead of estimating the proportion of consumed to available prey. Outcomes of our model, are that we give up the somewhat arbitrary preference for an index to have the range  $[-1, 1]$ , and random feeding, now denoted by 1, is formally testable across time points and across prey species.

Our model enables formal hypothesis testing and statistical inference, while being general enough to perform statistical tests across multiple species and time points. This provides researchers a more detailed analysis of the predator's feeding preferences. Further, because our model is based on underlying Poisson distributions, members of the well studied exponential family, we are able to estimate the parameters of interest even when exact counts of each prey species eaten within any given time period are not observed. Instead, we rely on the researcher being able to detect prey DNA within the predator's gut (Schmidt et al., 2014; Raso et al., 2014; Madduppa et al., 2014) and make a simple binary conclusion: this predator ate some of that prey species during this time period, or did not.

This paper is organised as follows. Section 2 describes our statistical model, for both fully observed count data, and for the non-observed count data for which we use the expectation-maximisation (EM) algorithm, and the statistical tests used to make statements about the population parameters of interest. In section 3, we offer a simulation study that demonstrates the accuracy of our methods. Section 4 provides a real data set, which investigates the eating preferences of wolf spiders (Araneae: Lycosidae), found in the Berea College Forest in Madison County, Kentucky, USA, to demonstrate how those interested in assessing trophic interactions with gut-content analyses could apply our methods. A brief discussion concludes the paper in section 5. Alongside our model, we offer an R (Core Team, 2014) package named **spiders** that implements the methods discussed.

## 2 Methods

We assume data are collected in the following manner. Traps are dispersed, for  $T$  time periods, throughout the habitat of the predator and prey of interest. Predators and prey species are collected in the traps and counted at each time period. We denote the number of predators and the number of prey species caught in each time period  $t \in \{1, \dots, T\}$  by  $J_t$  and  $I_t$ , respectively. Prey species will be indexed by  $s \in \{1, \dots, S\}$ . Let  $X_{jst}$  represent the number of prey species  $s$  that predator  $j$  ate during occurrence  $t$ , where  $j \in \{1, \dots, J_t\}$ . Let  $Y_{ist}$  represent the number of prey species  $s$  found in trap  $i$  during occurrence  $t$ ,  $i \in \{1, \dots, I_t\}$ .

The number of prey species  $s$  that predator  $j$  ate during occurrence  $t$  is assumed to follow a Poisson distribution with rate parameter  $\lambda_{st}$ ,  $X_{jst} \stackrel{\text{iid}}{\sim} \mathcal{P}(\lambda_{st})$ . The parameter  $\lambda_{st}$  represents the rate at which the predator ate prey species  $s$  during time period  $t$ . The number of prey species  $s$  found in trap  $i$  during occurrence  $t$  is assumed to follow a Poisson

distribution with rate parameter  $\gamma_{st}$ . By modelling  $\lambda_{st}$  and  $\gamma_{st}$  we are able to test claims about a predator's eating preferences. Formal statistical statements about the relative magnitudes of the parameters  $\boldsymbol{\lambda} = (\lambda_{11}, \dots, \lambda_{ST})^t$  and  $\boldsymbol{\gamma} = (\gamma_{11}, \dots, \gamma_{ST})^t$  offer insights to the relative rates at which predators eat particular prey species.

The use of Poisson distributions make the following implicit assumptions: 1) traps independently catch the prey species of interest, 2) predators eat independently of each other.

We consider five variations on the relative magnitude of  $c_{st} = \lambda_{st}/\gamma_{st}$ . These five hypotheses each allow  $c_{st}$  to vary by time, prey species, both, or neither. Because the five hypotheses are nested, a natural testing order is suggested in Figure 1.

1.  $c_{st} = 1$

2.  $c_{st} = c$

3.  $c_{st} = c_s$

4.  $c_{st} = c_t$

5.  $c_{st} = c_{st}$

The first hypothesis states that the relative rate of sampling for the predator and the traps is the same for all species on all occasions. One imagines this is the case if they prey move randomly and the predator simply eats prey which comes within its reach, thus suggesting no selection for a particular prey item. The second states that predators sample prey proportionally across all time periods. The third hypothesis states that predators sample different prey species at different rates, but each rate is steady across time. This implies that the predator expresses preferences for one prey species over another, but is unresponsive to changes due to time. Conversely, the fourth hypothesis implies that each prey species is sampled similarly within each time period, while the rates across time are allowed to change. The fifth hypothesis assumes a predator's selection varies by both time and prey species. This would make sense if environmental and biological variables, such as weather, prey availability, and/or palatability were affecting predators' selection strategies.

[Figure 1 about here.]

## 2.1 Fully Observed Count Data

The likelihood function that allows for estimation of these parameters is as follows. Since we assume  $X_{jst}$  is independent of  $Y_{ist}$  we can simply multiply the respective Poisson probability density functions, and then form products over all  $s, t$  to obtain the likelihood.

$$L(x_{jst}, y_{ist} | \boldsymbol{\lambda}, \boldsymbol{\gamma}) = \prod_{t=1}^T \prod_{s=1}^S \left\{ \prod_{j=1}^{J_t} f_X(x_{jst} | \boldsymbol{\lambda}) \prod_{i=1}^{I_t} f_Y(y_{ist} | \boldsymbol{\gamma}) \right\}. \quad (1)$$

Writing all five hypotheses as  $\lambda_{st} = c_{st}\gamma_{st}$ , we can, in the simplest cases, find analytic solutions for the maximum likelihood estimates (MLEs) of  $c_{st}$ ,  $\gamma_{st}$ , and by invariance  $\lambda_{st}$ , when the data are balanced  $J_t = J$ ,  $I_t = I$ , and  $c_{st} = c$ . Namely, these solutions are

$$\hat{\gamma}_{st} = \frac{X_{\cdot st} + Y_{\cdot st}}{J_t + I_t}, \quad \text{and} \quad \hat{c} = \frac{I \sum_{s,t} X_{\cdot st}}{J \sum_{s,t} Y_{\cdot st}}, \quad \hat{\gamma}_{st} = \frac{X_{\cdot st} + Y_{\cdot st}}{I \left( \frac{\sum_{st} X_{\cdot st}}{\sum_{st} Y_{\cdot st}} + 1 \right)}$$

respectively, where  $X_{\cdot st} = \sum_{j=1}^{J_t} X_{jst}$  and  $Y_{\cdot st} = \sum_{i=1}^{I_t} Y_{ist}$ .

In all other cases, analytic solutions are not readily available and instead we rely on the fact that the log-likelihood  $l(\boldsymbol{\lambda}, \boldsymbol{\gamma}) = \log L$  is concave. To compute MLEs, we maximise the log-likelihood, using coordinate descent (Luo and Tseng, 1992), by iteratively solving partial derivatives of  $l$ , with respect to  $c_{st}$  and  $\gamma_{st}$ , set equal to zero. The partial derivatives of the log-likelihood with respect to the parameters of interest, in the models  $c$ ,  $c_t$ ,  $c_s$ , and  $c_{st}$ , set equal to zero and solved are as follows.

$$\hat{c} = \frac{\sum_{s,t} X_{\cdot st}}{\sum_t J_t \sum_s \gamma_{st}}, \quad \hat{c}_t = \frac{\sum_s X_{\cdot st}}{J_t \sum_s \gamma_{st}}, \quad \hat{c}_s = \frac{\sum_t X_{\cdot st}}{\sum_t J_t \gamma_{st}}, \quad \text{or } \hat{c}_{st} = \frac{X_{\cdot st}}{J_t \gamma_{st}}, \quad \text{and} \quad \hat{\gamma}_{st} = \frac{X_{\cdot st} + Y_{\cdot st}}{J_t c_{st} + I_t}.$$

The equation shown for  $\gamma_{st}$  is the general solution for which  $c_{st}$  could be replaced by  $c$ ,  $c_t$ , or  $c_s$ .

## 2.2 Unobserved Counts

In many applications, such as DNA-based gut-content analysis, it is not possible to count the number of individuals of each prey species that are in a predator's gut. Instead, it is only possible to detect whether or not a predator consumed the prey species during a given time period, based on the rate at which prey DNA decays in the predator gut (Greenstone et al., 2013). In this case we can still make inference about the predators' preferences for the different prey species by using the expectation-maximisation (EM) algorithm to compute MLEs.

We denote the binary random variable indicating whether the  $j^{th}$  predator did in fact eat at least one individual of prey species  $s$  in time period  $t$  by  $Z_{jst} = 1(X_{jst} > 0)$ . Given the Poisson assumptions above, these variables are independent Bernoulli observations with success probability  $p_{st} = P(Z_{jst} = 1) = 1 - \exp\{-\lambda_{st}\}$ . Despite not observing  $X_{jst}$ , we can compute maximum likelihood estimates of the parameters  $\boldsymbol{\lambda}, \boldsymbol{\gamma}$  through the EM algorithm using the complete data log-likelihood

$$l_{comp}(\boldsymbol{\lambda}, \boldsymbol{\gamma}) = \log f_{X,Y,Z}(\mathbf{x}, \mathbf{y}, \mathbf{z} | \boldsymbol{\lambda}, \boldsymbol{\gamma}) = \sum_{s=1}^S \sum_{t=1}^T \left[ \sum_{j=1}^{J_t} \log f_{X,Z}(x_{jst}, z_{jst} | \boldsymbol{\lambda}) + \sum_{i=1}^{I_t} \log f_Y(y_{ist} | \boldsymbol{\gamma}) \right].$$

The density of  $Y_{jst}$  is exactly as in section 2.1 and so we focus on deriving the joint density of  $X_{jst}$  and  $Z_{jst}$ . With the distribution of  $Z_{jst}$  given above, we can compute  $f_{X,Z}(x_{jst}, z_{jst} | \boldsymbol{\lambda})$  by noting that  $X_{jst} = 0$  with probability 1 if  $Z_{jst} = 0$ , and that  $[X_{jst} | Z_{jst} = 0]$  has a truncated Poisson distribution with density

$$f_{X|Y,Z,\boldsymbol{\lambda},\boldsymbol{\gamma}}(x_{jst} | z_{jst}) = \frac{\exp\{-\lambda_{st}\} \lambda_{st}^{x_{jst}}}{(1 - \exp\{-\lambda_{st}\}) x_{jst}!} 1(x_{jst} > 0)$$

156 and expected value

$$\mathbb{E}_{X|Y,Z} X_{jst} = \frac{\lambda_{st} \exp \{\lambda_{st}\}}{\exp \{\lambda_{st}\} - 1}.$$

The joint density of  $X_{jst}, Z_{jst}$  is then

$$f_{X,Z|\lambda}(x_{jst}, z_{jst}) = \begin{cases} \exp \{-\lambda_{st}\}, & x_{jst} = 0 \text{ and } z_{jst} = 0 \\ \frac{\exp \{-\lambda_{st}\} \lambda_{st}^{x_{jst}}}{x_{jst}!}, & x_{jst} > 0 \text{ and } z_{jst} = 1 \\ 0 & \text{otherwise} \end{cases}.$$

157 The EM algorithm works by iterating two steps, the E-step and M-step, until the optimum  
158 is reached (Dempster et al., 1977; McLachlan and Krishnan, 2007). Let  $k$  index the iterations  
159 in the EM algorithm so that  $\boldsymbol{\lambda}^{(k)}$  and  $\boldsymbol{\gamma}^{(k)}$  denote the estimates computed on the  $k^{th}$  M-step.  
160 The E-step consists of computing the expectation of  $l_{comp}$  with respect to the conditional  
161 distribution of  $X$  given the current estimates of the parameters

$$Q^{(k)}(\boldsymbol{\lambda}, \boldsymbol{\gamma}) = \mathbb{E}_{X|Y,Z,\boldsymbol{\lambda}^{(k)}} l_{comp}$$

162 in order to remove the unobserved data. The M-step then involves maximising  $Q^{(k)}(\boldsymbol{\lambda}, \boldsymbol{\gamma})$   
163 with respect to the parameters in the model to obtain updated estimates of the parameters,

$$(\boldsymbol{\lambda}^{(k+1)}, \boldsymbol{\gamma}^{(k+1)}) = \arg \max_{(\boldsymbol{\lambda}, \boldsymbol{\gamma})} Q^{(k)}(\boldsymbol{\lambda}, \boldsymbol{\gamma}).$$

164 These steps are alternated until a convergence criterion monitoring subsequent differences  
165 in the parameter estimates/likelihood is met.

The calculation of  $Q^{(k)}(\boldsymbol{\lambda}, \boldsymbol{\gamma})$  is not difficult and is given by:

$$\begin{aligned} Q^{(k)}(\boldsymbol{\lambda}, \boldsymbol{\gamma}) &= \mathbb{E} \log f_{X,Z|\lambda}(X_{jst}, z_{jst}) + \log f_{Y|\gamma}(y_{ist}) \\ &= \sum_{s=1}^S \sum_{t=1}^T \sum_{j=1}^{J_t} \mathbb{E} \log f_{X,Z|\lambda}(X_{jst}, z_{jst}) + \sum_{s=1}^S \sum_{t=1}^T \sum_{i=1}^{I_t} \log f_{Y|\gamma}(y) \\ &\propto \sum_{s,t,j} (-\lambda_{st} + z_{jst} \log \lambda_{st} \mathbb{E} X_{jst}) + \sum_{s,t} (-I_t \gamma_{st} + Y_{.st} \log I_t \gamma_{st}) \\ &\propto \sum_{s,t} \left( -J_t \lambda_{st} + z_{.st} \log \lambda_{st} \mathbb{E}(X_{jst} | \lambda_{st}^{(k)}, \gamma_{st}^{(k)}) \right) + \sum_{s,t} (-I_t \gamma_{st} + Y_{.st} \log I_t \gamma_{st}). \end{aligned} \tag{2}$$

166 No analytic solution to the M-step exists, however, so we again chose to maximise  $Q$  with  
167 coordinate descent (Luo and Tseng, 1992). In fact, as we only need to find parameters that  
168 increase the value of  $Q$  on each iteration, we forgo fully iterating the coordinate descent algo-  
169 rithm to find the maximum and instead perform just one step uphill within each EM iteration  
170 (Givens and Hoeting, 2012). Since  $Q^{(k)}$  is concave and smooth in the parameters  $\boldsymbol{\lambda}, \boldsymbol{\gamma}$ , we  
171 are able to use the convergence of parameter estimates,  $\|(\boldsymbol{\lambda}^{(k)}, \boldsymbol{\gamma}^{(k)}) - (\boldsymbol{\lambda}^{(k+1)}, \boldsymbol{\gamma}^{(k+1)})\|_{\infty} < \tau$ ,  
172 for some  $\tau > 0$ , as our stopping criterion.

173 As we show in our simulation study, this generalised EM algorithm accurately estimates  
174 the parameters when values of  $\lambda_{st}$  are relatively small, such that zeros are prevalent in the  
175 data  $Z_{jst}$ . In contrast, if the predator consistently eats a given prey species, few to no

zeros will show up in the observed data and  $\mathbb{E}Z_{jst}$  is estimated to be nearly 1. The loss of information is best seen by attempting to solve for  $\lambda_{st}$  in the equation  $1 = \mathbb{E}Z_{jst} = 1 - \exp\{-\lambda_{st}\}$ . As the proportion of ones in the observed data increases, we expect  $\lambda_{st}$  to grow exponentially large. When no zeros are present in the data, so that only ones are observed, the likelihood can be made arbitrarily large by sending the parameter off to infinity.

## 2.3 Testing

The likelihood ratio test statistic is

$$\Lambda(X, Y) := -2 \log \frac{\sup_{\theta_0} L(\theta_0|X, Y)}{\sup_{\theta_1} L(\theta_1|X, Y)},$$

where  $\theta_0, \theta_1$  represent the parameters estimated under the null and alternative hypotheses, respectively. It is well known that the asymptotic distribution of  $\Lambda$  is a  $\chi^2_\rho$  distribution with  $\rho$  degrees of freedom (Wilks, 1938). The degrees of freedom  $\rho$  equal the number of free parameters available in the stated hypotheses under question. If we put the null hypothesis to be  $H_0 : \lambda_t = c_t \gamma_t$ , for all  $t$  and contrast this against  $H_1 : \lambda_{st} = c_{st} \gamma_{st}$  then there are  $\rho = 2(S \cdot T) - S \cdot T - T = S \cdot T - T$  degrees of freedom. When the observations  $X_{jst}$  are not observed, we use  $L_{obs}(Z, Y)$  as the likelihood in the calculation of  $\Lambda$ . The level of significance,  $\alpha$ , is used to reject the null hypothesis in favour of the alternative hypothesis if  $\mathbb{P}(\chi^2_\rho > \Lambda) < \alpha$ .

## 2.4 Linear Transformations of $c_{st}$

After determining which model best fits the data, more detail may be extracted through specific hypothesis test of the elements of  $c_{st}$ , or in vector notation as  $\mathbf{c} \in \mathbb{R}^{S \cdot T}$ . Let the elements of  $\hat{\mathbf{c}}$  be the maximum likelihood estimates,  $\hat{c}_{st}$ , as found via the framework above. Since  $\hat{\mathbf{c}}$  is asymptotically normally distributed, any linear combination of the elements is also asymptotically normally distributed. For instance, let  $a$  be a vector of the same dimension of  $\hat{\mathbf{c}}$ . Then  $a^t \hat{\mathbf{c}}$  is asymptotically distributed as  $\mathcal{N}(a^t \mathbf{c}, a^t \Sigma a)$ , where  $\Sigma$  is the covariance matrix of the asymptotic distribution of  $\hat{\mathbf{c}}$ . Tests of the form  $H_0 : a^t \mathbf{c} = \mu$  against any alternative of interest are then approximate  $Z$ -tests. Confidence intervals of any size are similarly, readily obtained. Suppose, for example, that the hypothesis  $c_s$  is determined to best fit the data with  $s$  ranging  $s = 1, 2, 3$ . We can test to see whether or not the first two species are equally preferred under the null hypothesis  $c_1 = c_2$ . This hypothesis is alternatively written in vector notation as  $a^t \mathbf{c} = 0$ , where  $a = (1, -1, 0)^t$ .

## 3 Simulation Study

Our simulations assume two prey species and five time points, throughout. Of the hierarchy of hypotheses, we generate data under three models:  $c, c_s, c_t$ . Sample sizes for both prey species and predator gut count observations are randomly chosen from four overlapping levels: “small” sample sizes are randomly sampled numbers in  $[20, 50]$ , “medium”  $[30, 75]$ , “large”  $[50, 150]$ , and “larger”  $[100, 200]$ . This is repeated for each level of sample size.

We simulate 500 replicate data sets for each of the twelve scenarios above for both types of data, fully observed count data,  $X_{jst}$ , and for non-count data, when we observe only a binary response,  $Z_{jst} = 1(X_{jst} > 0)$ . Each scenario is then fitted with the true model that generated the data. All simulations of non-count data use  $\tau = 10^{-5}$  as the convergence tolerance. A subset of the examples are provided here; the interested reader is referred to the supplementary materials for the complete simulation results. For the simulations we used the R Core Team (2014) package **BatchExperiments** by Bischl et al. (2014).

[Figure 2 about here.]

For all simulated data, the true parameter values for the rate at which prey species are encountered in the wild are fixed to be  $\gamma_{st} = \pi \approx 3.14, \forall s, t$ . The values of  $\lambda_{st}$  are set with respect to each data generating model. For model  $c_{st} = c$ , where predator preferences don't vary by either time or species, we put  $\lambda_{st} = 2\pi, \forall s, t$ . Under model  $c_s$ , the ratio of rates vary by species only, so we put  $\lambda_{1t} = \sqrt{2}$  and  $\lambda_{2t} = \pi$ . Hence,  $c_1 = \sqrt{2}/\pi \approx 0.45$  and  $c_2 = 1$ . For the last model,  $c_t$ , the ratio of rates vary by time  $t$ . Here, we put  $\lambda_{st} = t$  for  $t \in \{1, \dots, 5\}$ .

[Figure 3 about here.]

We consider results when the correct model is fit to the simulated data. Figure 2 shows the density plot of the estimates of  $c_s$  when fitting the true model to the fully observed count data generated under models  $c_s$ , while figure 3 shows the same for the estimates of  $c$  when data is generated under model  $c$ . The plots provide evaluations of parameter estimates under each scenario. For model  $c_s$  in figure 2, the parameters  $c_1 \approx 0.45$  and  $c_2 = 1$  are on average, across all 500 simulations, estimated as  $\hat{c}_1 = 0.45$  and  $\hat{c}_2 = 1.00$ , with sample standard deviations of  $SD(\hat{c}_1) = 0.03$  and  $SD(\hat{c}_2) = 0.06$ . Figure 3 provides results for model  $c_{st} = c$ . Averaging across all 500 simulations, the parameter  $c = 2$  is estimated as  $\hat{c} = 2.00$  with sample standard deviation  $SD(0.06)$ . This is further seen in figure 4, where box plots of the parameter estimates, centred at true parameter values, of the correct model fit to data generated from both  $c_s$  and  $c_t$  show empirically very little bias.

[Figure 4 about here.]

We next generated data with unobserved counts. As noted above under certain circumstances our unobserved counts model accurately estimates the parameters of interest, and at other times can infinitely over-estimate parameters. To investigate this issue further, we consider the same scenarios mentioned above, but reduce all of the count data down to binary observations. For each scenario, we fit the unobserved counts model as if we knew the true underlying model that generated the observed data.

Figures 5 and 6 contain density plots of the estimates of  $c_s, c_t$  for all 500 replications of the data generating models  $c_s, c_t$  with the small and the larger sample sizes, respectively. When data are generated under the model  $c_s$  and the true model is fit to the non-count data, we find even for the small sample size that point estimates are only very slightly biased. When parameter values are of sufficient size to make zeros in the simulated data less common, the estimates from fitting the correct model to the generated data are occasionally over-estimated. This effect is easily seen in figure 6 for the two greatest values of  $c_t$  despite the increased sample size, but is also seen, less dramatically, in the density plot for the  $c_s$  generated data.



[Figure 5 about here.]

The cluster of estimates for  $c_5$  between 3.5 and 4.0 in figure 6 comes from data sets in which  $Z_{js5} = 1$  for all  $j, s$ . For the data shown in figure 6, this happened 73 times out of the 500 replicated data sets. As mentioned above, the estimate of  $c_5$  is infinite in this case. However, the EM algorithm will always provide a finite estimate for all parameters when it terminates. In this case, we set  $\tau = 10^{-5}$  and this caused the algorithm to terminate with  $\hat{c}_5$  between 3.5 and 4.0. To confirm that this is due to the arbitrary choice of  $\tau$ , we repeated the algorithm with smaller values of  $\tau$  for several data sets. As expected,  $\hat{c}_5$  increased without bound as we refit the model with increasingly small values of  $\tau$ .

[Figure 6 about here.]

The over-estimation of parameters, a symptom of the loss of information due to the unobserved counts, can also be seen with box plots of the 500 point estimates centred at their respective true parameter values. Figure 7 contains box plots of the same scenarios in figures 5 and 6. For the 73 cases in which  $Z_{js5} = 1$  for all  $j, s$  under model  $c_t$  with the larger sample size, the bias is infinite since parameter estimates will, theoretically, be infinite. The finite bias shown in these plots is due to the finite estimates provided by the termination of the EM algorithm. Thus, conditional on a mixture of 0s and 1s in the data the corresponding estimators appear to be unbiased, but when no 0s exist in the data the theoretical bias is infinite.

[Figure 7 about here.]

## 4 Application

To illustrate these methods, we analysed a dataset that was collected to investigate the feeding preferences of two species of wolf spider, *Schizocosa ocreata* and *Schizocosa stridulans* (Araneae: Lycosidae). Every 6 – 12 days, 10 to 40 spiders were hand-collected between October 2011 and April 2013 within Berea College Forest in Madison County, Kentucky, USA. Spiders were removed from the litter using an aspirator, placed in separate 1.5 mL microcentrifuge tubes filled with 95% EtOH, and preserved at  $-20^\circ\text{C}$  until DNA extraction. In parallel, we also surveyed availability of forest floor prey using pitfall traps ( $n = 32$ ). For the analysis, both species of *Schizocosa* were pooled and the number of spiders and prey were analysed by month. On average, 69 spiders, 111 Diptera, and 297 Collembola were caught in each time period. The range of the sample sizes across all 18 months was 11 to 181 for caught spiders, 7 to 322 for trapped Diptera, and 101 to 755 for trapped Collembola. Figure 8 plots the total number of each order that was caught during each time period.

[Table 1 about here.]

To determine whether spiders had consumed dipterans and/or collembolans, we conducted a molecular analysis of their gut-contents. First, DNA from spiders was extracted using Qiagen DNEasy®Tissue Extraction Kit (Qiagen Inc., Chatsworth, California, USA) following

the animal tissue protocol outlined by the manufacturer, with minor modifications. Whole bodies of the spiders were first crushed to release prey DNA from within their alimentary canal for extraction. The 200 $\mu$ L extractions were stored at  $-20^{\circ}\text{C}$  until PCR. Second, order-specific primers from the literature were used to detect the DNA of Collembola and Diptera within the guts of the spiders. Primer pairs designed by Sint et al. (2012), targeting the 18S rDNA gene, were used to detect Collembola predation table 7. A PCR cycling protocol for 12.5 $\mu$ L reactions containing 1 $\times$  Takara buffer (Takara Bio Inc., Shiga, Japan), 0.2 mM dNTPs, 0.2 $\mu$ M of each primer, 0.625 U Takara Ex Taq<sup>TM</sup> and 1.5 $\mu$ L of template DNA, using BioRad PTC-200 and C1000 thermal cyclers (Bio-Rad Laboratories, Hercules, California, USA), was optimised as follows: 95 $^{\circ}\text{C}$  for 1 minute, followed by 35 cycles of 94 $^{\circ}\text{C}$  for 30 seconds, 61.2 $^{\circ}\text{C}$  for 90 seconds, and 72 $^{\circ}\text{C}$  for 60 seconds. Primer pairs designed by Eitzinger et al. (2014), targeting the 18S rDNA gene, were used to detect Diptera predation table 7. PCR cycling protocol for 12.5 $\mu$ L reactions with Takara reagents (as above) and 2 $\mu$ L of template DNA was optimised as follows: 95 $^{\circ}\text{C}$  for 1 minute, followed by 40 cycles of 94 $^{\circ}\text{C}$  for 45 seconds, 60 $^{\circ}\text{C}$  for 45 seconds, and 72 $^{\circ}\text{C}$  for 45 seconds. Both primer pairs were tested for cross-reactivity against a range of prey and predator species from the field site and in all cases, no amplification of DNA was observed, confirming suitable specificity of the primers for this study. Lastly, electrophoresis of 10 $\mu$ L of each PCR product was later conducted to determine success of DNA amplification using 2% Seakem agarose (Lonza, Rockland, Maine, USA) stained with 1 $\times$  GelRed<sup>TM</sup> nucleic acid stain (Biotium, Hayward, California, USA). This procedure allowed us to determine a presence or an absence of Diptera and Collembola DNA within each spider.

[Figure 8 about here.]

[Figure 9 about here.]

These data provide an example of our hierarchy of hypotheses. First, we tested model  $c_{st} = c$  against  $c_{st} = c_s$ , to determine whether or not the wolf spider has different preferences for the two orders Diptera and Collembola. With, one degree of freedom, this likelihood ratio test indicated,  $p - \text{value} < 0.0001$ , that two parameters, one for each order, fits these data better than one parameter for both. Similarly, we tested whether or not there was a significant effect across time by testing model  $c_{st} = c$  against  $c_{st} = c_t$ . Here, the likelihood ratio test, with 17 degrees of freedom, implies that the wolf spiders of the Berea College Forest eat these prey orders at different rates across the months of the year,  $p - \text{value} < 0.0001$ . In fact, we find that the most parameter rich model,  $\lambda_{st} = c_{st}\gamma_{st}$  fits these data better than is expected by chance, when compared to either model  $c_t$  ( $p - \text{value} < 0.0001$ ) or model  $c_s$  ( $p - \text{value} < 0.0001$ ). Model  $c_{st}$  estimates 72 parameters in total; since, in this case, there are two prey of interest and 18 time periods, it takes 36 parameters to estimate each  $c_{st}$  and  $\gamma_{st}$ . Figures 10, 11 plot the point estimates and 95% confidence intervals of  $c_{st}$ , for both prey across all time periods.

[Figure 10 about here.]

[Figure 11 about here.]

With point estimates of  $c_{st}$  under the model  $\lambda_{st} = c_{st}\gamma_{st}$ , we can test any number of linear contrasts. For instance, the hypotheses  $c_{1t} = c_{2t}$ , for  $t \in \{1, \dots, 18\}$  state that wolf spiders equally prefer the orders Diptera and Collembola at each of the 18 time points. Using a level of significance of 0.05, and after making a Bonferroni multiple comparisons adjustment, the data can not say that the two prey are differently preferred in October, November, and December of 2011 and for March and July of 2012.

## 5 Discussion

The model developed here allows for the determination of a predator's eating preferences by testing simultaneously across an array of multiple prey species and time points. This is achieved via a simple, but statistically powerful, likelihood ratio test. Further testing of the ratio of rates for which predators eat to encounter prey species allows researchers to make specific conclusions about predators' preferences. For instance, rates across time can be estimated to make statements about seasonal effects on a predator's eating habits, while relative rates across species groups allows for statements about the relative preferences for different species.

When counts of predators' gut contents are not fully observed, and instead only a binary response indicating the existence of the prey species in the gut is observed, we are able to treat the counts as missing data. By modelling all of the observed data, both the binary responses and the number of prey species caught, and the missing count data, we are able to use the EM algorithm to extract as much information from the data as possible.

Further developments of our model could take into account other environmental variables that might affect a predator's eating habits, such as rain or temperature.

## 6 Acknowledgements

The information reported in this paper (No. 15–08–008) is part of a project of the Kentucky Agricultural Experiment Station and is published with the approval of the Director. Support for this research was provided by the University of Kentucky Agricultural Experiment Station State Project KY008055 and the National Science Foundation Graduate Research Fellowship Program.

## 7 Data Accessibility

An R package, named `spiders`, is available on CRAN at <http://cran.r-project.org/web/packages/spiders/index.html> and fits all the methods discussed above.

## References

- Bernd Bischl, Michel Lang, and Olaf Mersmann. *BatchExperiments: Statistical experiments on batch computing clusters.*, 2014. URL <http://CRAN.R-project.org/package=BatchExperiments>. R package version 1.3.
- Jean Chesson. Measuring preference in selective predation. *Ecology*, 59(2):211–215, 1978.
- Jean Chesson. The estimation and analysis of preference and its relationship to foraging models. *Ecology*, 64(5):1297–1304, 1983.
- Hayley S Clements, Craig J Tambling, Matt W Hayward, and Graham IH Kerley. An objective approach to determining the weight ranges of prey preferred by and accessible to the five large african carnivores. *PloS one*, 9(7):e101054, 2014.
- R Core Team. *R: A language and environment for statistical computing*. R Foundation for Statistical Computing, Vienna, Austria, 2014. URL <http://www.R-project.org/>.
- Arthur P Dempster, Nan M Laird, and Donald B Rubin. Maximum likelihood from incomplete data via the em algorithm. *Journal of the Royal Statistical Society. Series B (Methodological)*, 39(1):1–38, 1977.
- Bernhard Eitzinger, E Michael Unger, Michael Traugott, and Stefan Scheu. Effects of prey quality and predator body size on prey dna detection success in a centipede predator. *Molecular ecology*, 23(15), 2014.
- Geof H Givens and Jennifer A Hoeting. *Computational statistics*, volume 708. John Wiley & Sons, 2012.
- Matthew H Greenstone, Mark E Payton, Donald C Weber, and Alvin M Simmons. The detectability half-life in arthropod predator–prey research: what it is, why we need it, how to measure it, and how to use it. *Molecular ecology*, 23(15), 2013.
- Adam G Hansen and David A Beauchamp. Effects of prey abundance, distribution, visual contrast and morphology on selection by a pelagic piscivore. *Freshwater Biology*, 59(11): 2328–2341, 2014.
- Peter Hellström, Jesper Nyström, and Anders Angerbjörn. Functional responses of the rough-legged buzzard in a multi-prey system. *Oecologia*, 174(4):1241–1254, 2014.
- Viktor Sergeevich Ivlev. *Experimental ecology of the feeding of fishes*. London, 1964.
- Jürgen Jacobs. Quantitative measurement of food selection. *Oecologia*, 14(4):413–417, 1974.
- Martin J Lechowicz. The sampling characteristics of electivity indices. *Oecologia*, 52(1): 22–30, 1982.
- Zhi-Quan Luo and Paul Tseng. On the convergence of the coordinate descent method for convex differentiable minimization. *Journal of Optimization Theory and Applications*, 72(1):7–35, 1992.

396 Salvador Lyngdoh, Shivam Shrotriya, Surendra P Goyal, Hayley Clements, Matthew W Hay-  
397 ward, and Bilal Habib. Prey preferences of the snow leopard (*Panthera uncia*): Regional  
398 diet specificity holds global significance for conservation. *PloS one*, 9(2):e88349, 2014.

399 Hawis H Madduppa, Neviaty P Zamani, Beginer Subhan, Unggul Aktani, and Sebastian CA  
400 Ferse. Feeding behavior and diet of the eight-banded butterflyfish *Chaetodon octofasciatus*  
401 in the thousand islands, indonesia. *Environmental Biology of Fishes*, pages 1–13, 2014.

402 Bryan FJ Manly, Lyman L McDonald, Dana L Thomas, Trent L McDonald, and Wallace P  
403 Erickson. *Resource selection by animals*. Springer, 1992.

404 Geoffrey McLachlan and Thriyambakam Krishnan. *The EM algorithm and extensions*, vol-  
405 ume 382. John Wiley & Sons, 2007.

406 Lorna Raso, Daniela Sint, Rebecca Mayer, Simon Plangg, Thomas Recheis, Silvia Brun-  
407 ner, Rüdiger Kaufmann, and Michael Traugott. Intraguild predation in pioneer predator  
408 communities of alpine glacier forelands. *Molecular ecology*, 23(15), 2014.

409 Jason M Schmidt, Sarah K Barney, Mark A Williams, Ricardo T Bessin, Timothy W Coo-  
410 long, and James D Harwood. Predator–prey trophic relationships in response to organic  
411 management practices. *Molecular ecology*, 23(15), 2014.

412 Daniela Sint, Lorna Raso, and Michael Traugott. Advances in multiplex PCR: balancing  
413 primer efficiencies and improving detection success. *Methods in Ecology and Evolution*, 3  
414 (5):898–905, 2012.

415 Richard E Strauss. Reliability estimates for Ivlev’s electivity index, the forage ratio, and a  
416 proposed linear index of food selection. *Transactions of the American Fisheries Society*,  
417 108(4):344–352, 1979.

418 HA Vanderploeg and D Scavia. Two electivity indices for feeding with special reference to  
419 zooplankton grazing. *Journal of the Fisheries Board of Canada*, 36(4):362–365, 1979.

420 Samuel S Wilks. The large-sample distribution of the likelihood ratio for testing composite  
421 hypotheses. *The Annals of Mathematical Statistics*, 9(1):60–62, 1938.

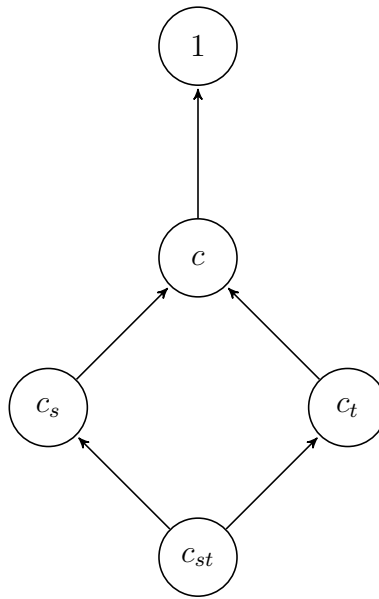


Figure 1: This hierarchy of hypotheses suggests the order in which the discussed models should be tested. One begins with the most complex models at the bottom and sequentially, following the arrows, tests simpler hypotheses using the formal test described in section 2.3 until a final model is established.

### Small Sample Size

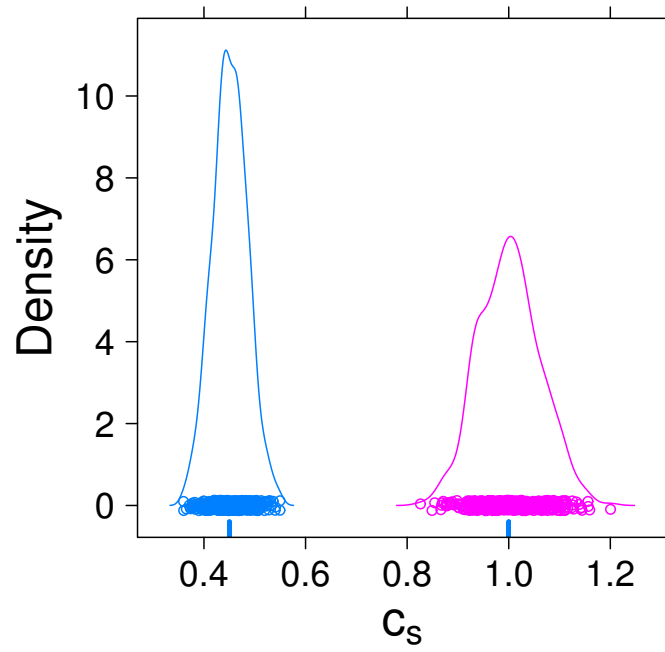


Figure 2: The density plot of all 500 estimates of fitting the true model to the data generated from models  $c_s$  with the small sample size is shown. Each element of  $c_s$  is colour coded for clarity, and ticks on the x-axis show the true parameter values.

### Medium Sample Size

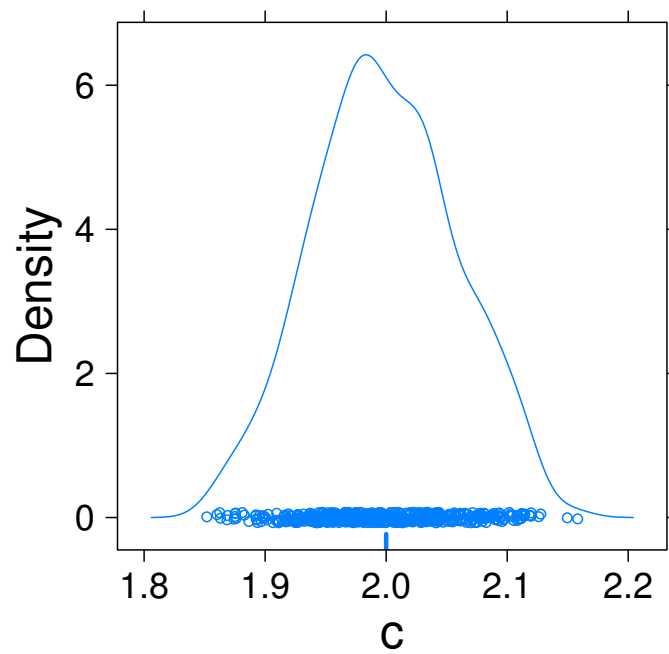


Figure 3: The density plot of all 500 estimates of fitting the true model to the data generated from model  $c$  with the medium sample size is shown. A tick on the x-axis shows the true parameter value.



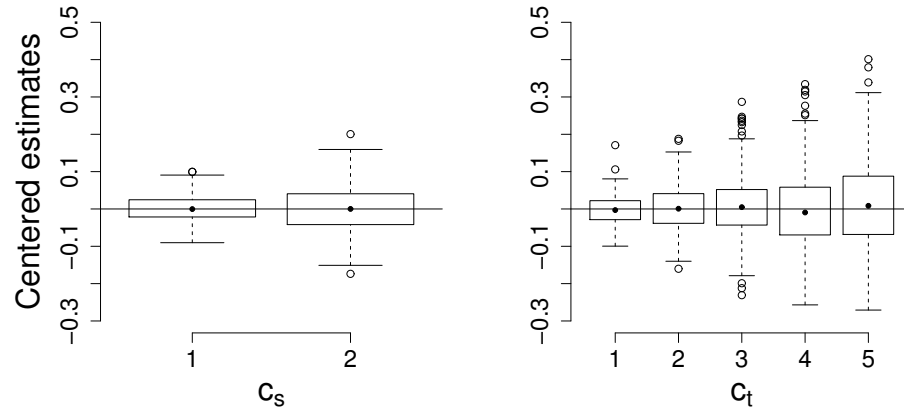


Figure 4: Shown are all 500 estimates, centred at the true parameter values, from fitting the true model to the data generated from models  $c_s, c_t$  with sample sizes small and medium, respectively.

### Small Sample Size, Non-Count Data

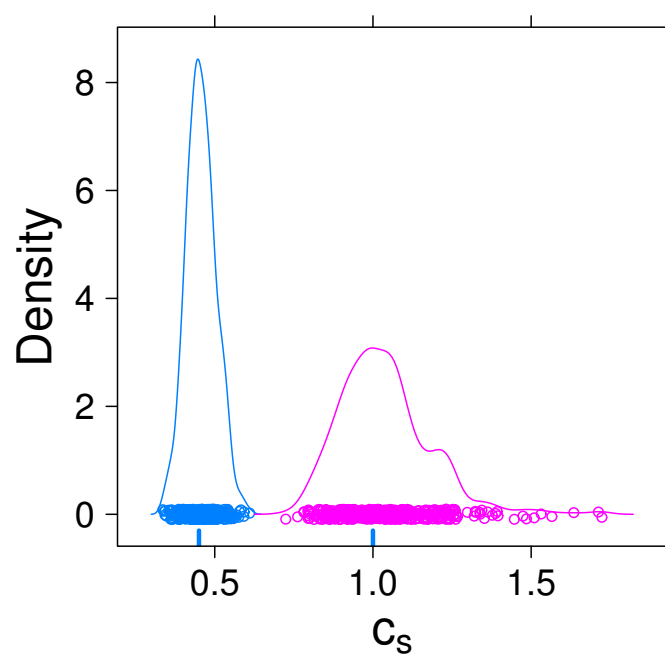


Figure 5: The density plot of all 500 estimates of fitting the true model to the data generated from model  $c_s$ , when counts are not observed, is shown with the small sample size. Each element of  $c_s$  is colour coded for clarity, and ticks on the x-axis show the true parameter values.

### Larger Sample Size, Non-Count Data

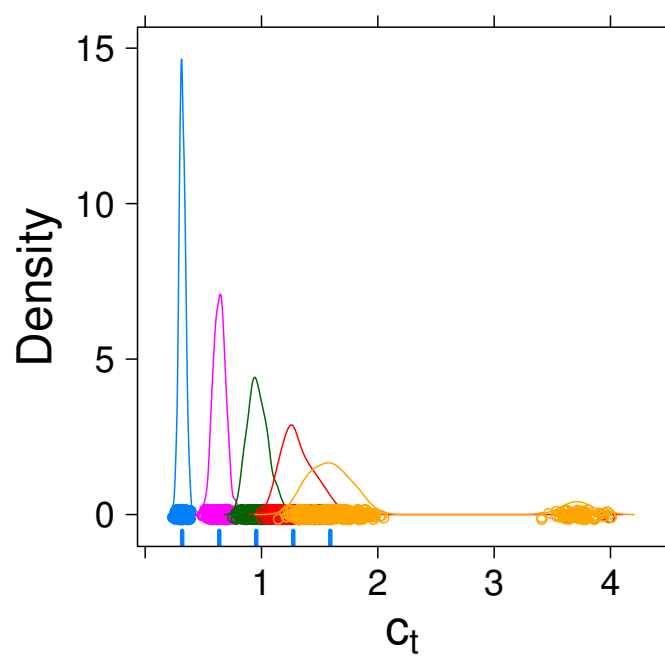


Figure 6: The density plot of all 500 estimates of fitting the true model to the data generated from model  $c_t$ , when counts are not observed, is shown with the larger sample size. Each element of  $c_t$  is colour coded for clarity, and ticks on the x-axis show the true parameter values.

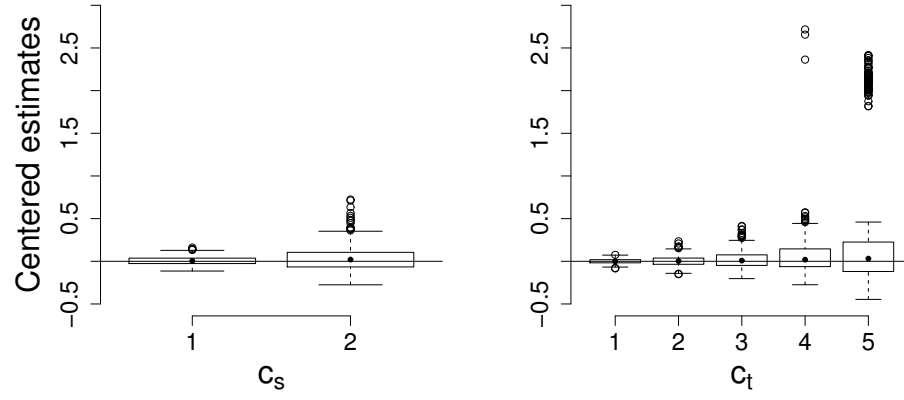


Figure 7: Shown are all 500 estimates, centred at the true parameter values, from fitting the true model to the data generated from models  $c_s, c_t$ , when counts are not observed, with sample sizes small and larger, respectively.

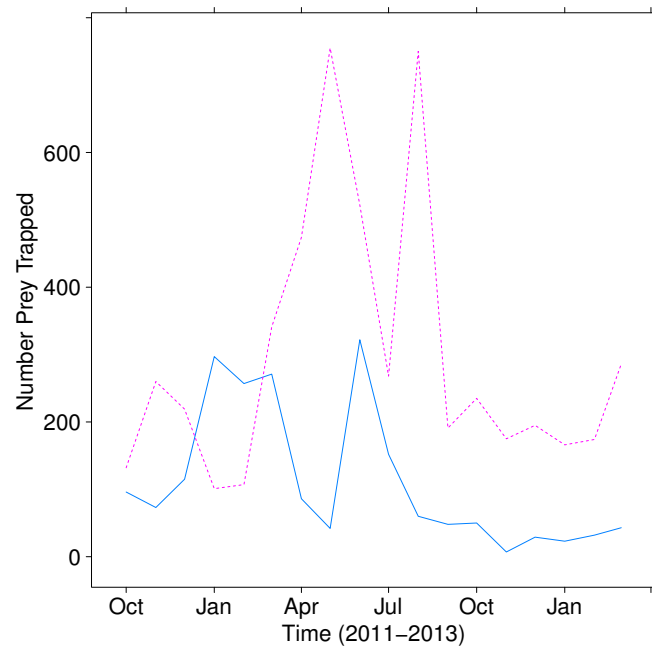


Figure 8: For both Collembola (pink/dashed) and Diptera (blue/solid), the plot shows the number of the prey trapped in each time period.

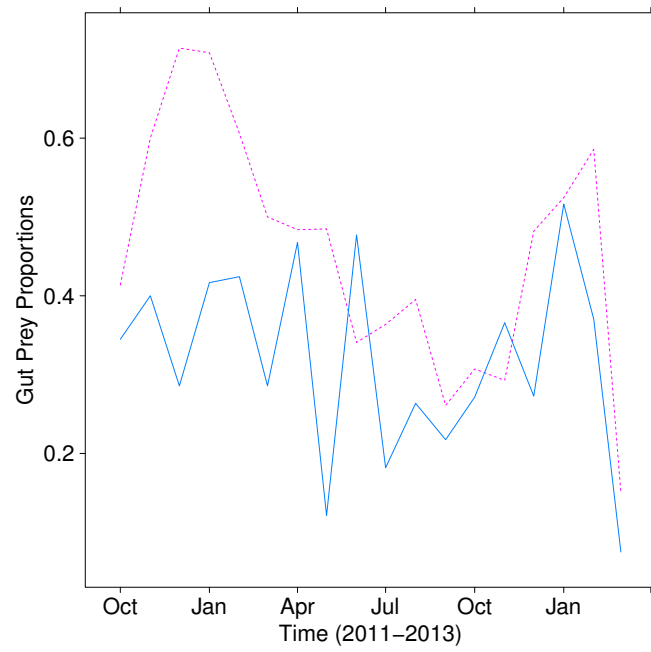


Figure 9: For both Collembola (pink/dashed) and Diptera (blue/solid), the plot shows the prey proportions in the sampled wolf spiders' guts in each time period.

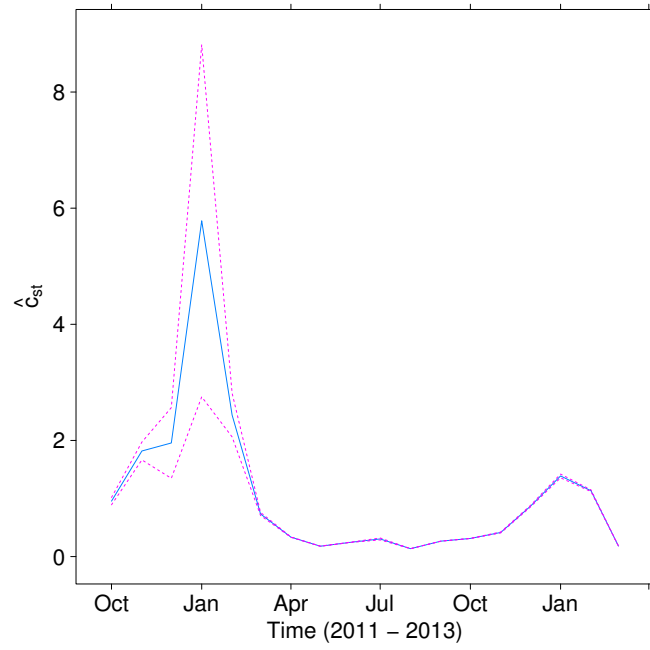


Figure 10: Point estimates (blue/solid) and 95% confidence intervals (pink/dashed) as estimated from the model  $c_{st}$  for Collembola.

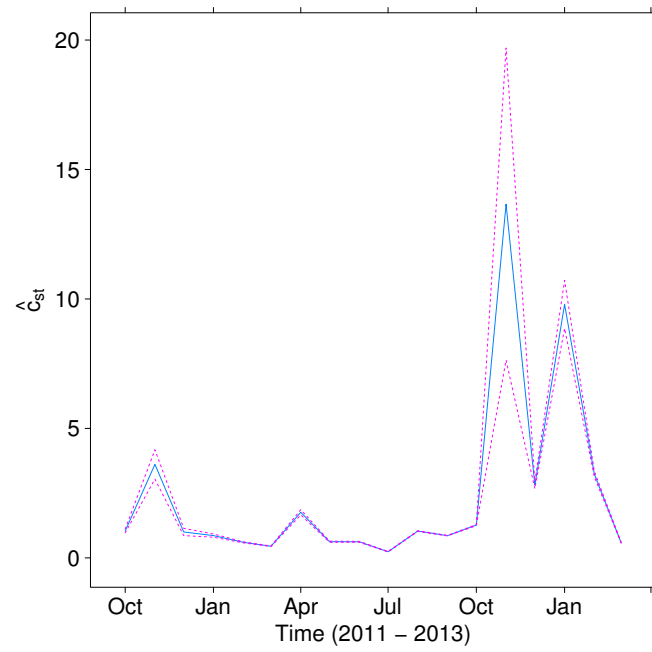


Figure 11: Point estimates (blue/solid) and 95% confidence intervals (pink/dashed) as estimated from the model  $c_{st}$  for Diptera.



Target group	Primer names and sequences 5' – 3'	Size (bp)	Source
Collembola	Col3F: GGACGATYTTRTTRGTTTCGT Col-gen-A246: TTTCACCTCTAACGTCGCAG	228	Sint et al. (2012)
Diptera	DIPS16: CACTTGCTTCTTAAATrGACAAATT DIPA17: TTyATGTGAACAGTTTCAGTyCA	198	Eitzinger et al. (2014)

Table 1: Targeted prey orders, primer names and sequences, size of amplicon, and source of design for the detection of prey taxa within the guts of Schizocosa spiders. Both primer sets were used in singleplex PCR assays.