

Supplementary Information 1 : Integrating data from different taxonomic resolutions to better estimate community alpha diversity.

Kwaku Peprah Adjei^{1,2}, Claire Carvell³, Nick Isaac³, Francesca Mancini³, Robert B. O'Hara^{1,2}

¹Department of Mathematical Sciences, Norwegian University of Science and Technology, Trondheim Norway

²Center for Biodiversity Dynamics, Norwegian University of Science and Technology, Trondheim Norway

³UK Center of Ecology and Hydrology, Wallingford UK

Exploratory analysis

Introduction

An overlooked step in developing integrated distribution models is exploring a model that fits each dataset. IDM takes advantage of the shared parameters to get better inference, and is not an escape from doing exploratory analysis on each dataset.

This supplementary material performs exploratory analysis on the UK pollinator monitoring scheme (Breeze et al., 2021). As mentioned in the main paper, we have data from three insect groups: bumblebees, solitary bees and hoverflies; and the models are fit independently for each insect group. The exploration is therefore also done for each insect group.

A Poisson and negative binomial generalised linear (mixed) models (GLMM hereafter) were fit for each insect group in the group counts dataset and binomial GLMM to the species occupancy data. The Poisson and logistic regression are the simplest model to fit to the count and occupancy data respectively (Fahrmeir et al., 2022). When there is overdispersion in the count data, the negative binomial regression will fit the count data better (Fahrmeir et al., 2022). The negative binomial regression adds an extra parameter ϕ to model the extra variation in the data. Values of ϕ closer to 0 indicate overdispersion in the count data. A likelihood ratio test of significance for the overdispersion parameter was performed. For each insect group, the log-likelihoods from the Poisson and negative Binomial regression models ($l(\hat{\beta}_{NB})$ and $l(\hat{\beta}_{Poi})$ respectively) were retrieved and test statistic used

Table S1-1: Coefficients from the best models (model equations defined above for BestModelIndex) for the group count and species occupancy data. The BestModelIndex shows which of the model fitted were selected as the best model (where best model is defined as the fitted model with the lowest BIC value); the Intercept and Latitude are fixed effect estimates; VisitVar and SpeciesVar are visit and species random effect variance estimates respectively; Psi is the overdispersion parameter of the negative Binomial and the P.value is the p-value estimate from likelihood ratio test of overdispersion (the difference between the Poisson and Negative Binomial models for each insect group).

Model	Data	BestModel	Intercept	Latitude	VisitVar	SpeciesVar	Psi	P.value
Bumblebees								
Binomial	Occupancy	M3	-4.54	0.08	0.07	1.95	-	-
Negative Binomial	Counts	M1	0.14	-0.38	-	-	1.42	0
Hoverflies								
Binomial	Occupancy	M3	-5.16	0.23	0.02	2.70	-	-
Negative Binomial	Counts	M2	0.96	-0.13	0.03	-	1.36	0
Solitarybees								
Binomial	Occupancy	M3	-5.99	-0.82	0.09	1.49	-	-
Negative Binomial	Counts	M1	-0.99	-0.53	-	-	1.37	0

was:

$$LRT = -2 \times (l(\hat{\beta}_{NB}) - l(\hat{\beta}_{Poi})) \sim \chi_1^2.$$

The negative binomial distribution is chosen for the count data if the p-value is less than 0.05 and Poisson is chosen when the p-value is greater than 0.05.

We performed stepwise model selection technique to obtain the best model for the count and occupancy data. We fit four models which we defined in Equation 1: the first model M_0 was an intercept-only model; the second model M_1 has an intercept and latitudinal gradient as a covariate; the third model M_2 added visits as a random effect to model M_1 ; and the last model M_3 added species as a random effect to model M_2 . Model M_3 was for the species occupancy data only, but the others were for both datasets. We fit the models M_0 and M_1 with the lme4 package (Bates et al., 2009) and M_2 and M_3 with the MASS package (Venables and Ripley, 2002) in R. The models were compared using their BIC values. We chose BIC over AIC values because they have high penalty for model complexity. The best model was the model with the lowest BIC value.

$$\begin{aligned}
M_0 &: \text{intercept} \\
M_1 &: \text{intercept} + \text{latitude} \\
M_2 &: \text{intercept} + \text{latitude} + (1|\text{visit}) \\
M_3 &: \text{intercept} + \text{latitude} + (1|\text{visit}) + (1|\text{species})
\end{aligned} \tag{1}$$

Results and Interpretation

Figure 1, Figure 3 and Figure 5 shows the distribution the bumblebees, hoverflies and solitary bees group count data for each visit to the 74 PoMS sites; Figure 2, Figure 4 and Figure 6 shows the distribution of the proportion of species occupying a given PoMS sites at each survey visit. The

Figure 1 to Figure 6 indicate a significant missing observation across the study regions and visits; with possible overdispersed group counts. The best models for the group count data showed significant overdispersion ($\phi_{BB} = 0.53$, $\phi_{HV} = 0.66$ and $\phi_{SB} = 0.37$ with p-values closer to 0; Table 1).

Furthermore, there were significant intercept and latitudinal gradient effect to both the group counts and species occupancy data (Table 1). Specifically, average abundance decreased as the latitudinal gradient increased for all insect groups ($\beta_{latitude} = -0.47, -0.83, -0.56$ for bumblebees, hoverflies and solitary bees respectively; Table 1); occupancy probability increased as latitudinal gradient increased for bumblebees and hoverflies ($\beta_{latitude} = 0.08, 0.23$ for bumblebees and hoverflies respectively; Table 1), but the opposite for the solitary bees ($\beta_{latitude} = -0.82$; Table 1).

In addition, there was significant species and visit effect on species occupancy probability (since model M_3 was chosen for all the insect groups as the best model; Table 1), but insignificant visit effects on bumblebees and solitary bees abundance (since model M_1 was the best model; Table 1).



Figure S1-1: Distribution of bumblebee counts for the 74 PoMS sites. The counts are faceted by the PoMS site name and colored by the visit number the observations was made. The columns of visits without no bars represent the visit without no group count observation ('NA').



Figure S1-2: Distribution of bumblebees occupancy probability for the 74 PoMS sites. The probabilities for each site are estimated as the proportion of sites with at most one pantrap being occupied by the species. The occupancy probabilities are faceted by the PoMS site name and colored by the visit number the observations was made. The columns of visits without no bars represent the visit without no group count observation ('NA').

Combining both datasets

From the results and discussions above, we model the group count with a negative binomial GLMM with latitudinal gradient as a covariate and visit random effect. The species occupancy probabilities were modeled with a binomial GLMM with latitudinal gradient as a covariate and species and visit random effect. We allow both models to share latitudinal gradient effect as well as species and visit random effects. The details are describe under section 2.2 in the main paper.

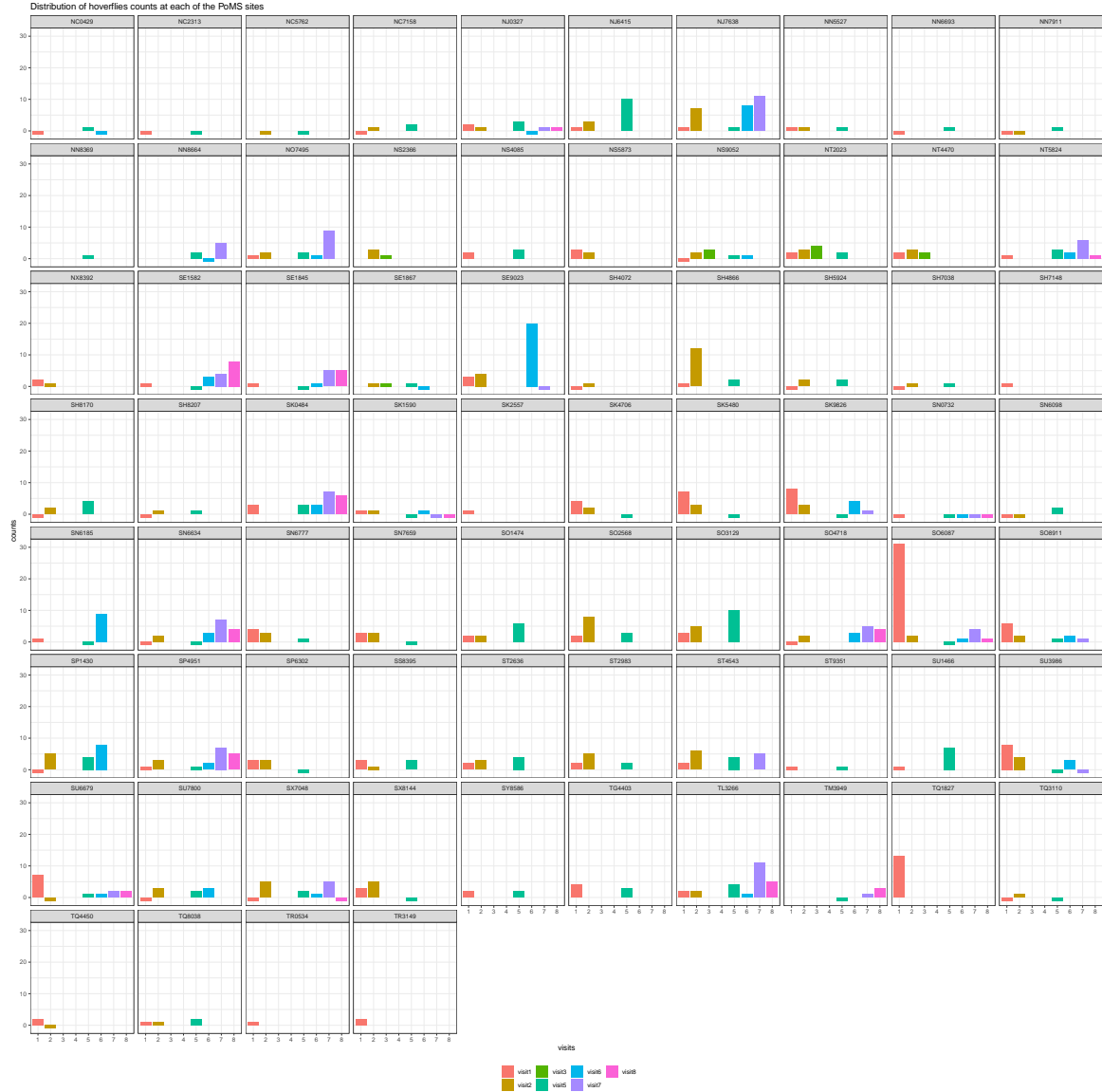


Figure S1-3: Distribution of hoverflies counts for the 74 PoMS sites. The counts are faceted by the PoMS site name and colored by the visit number the observations was made. The columns of visits without no bars represent the visit without no count observation ('NA').

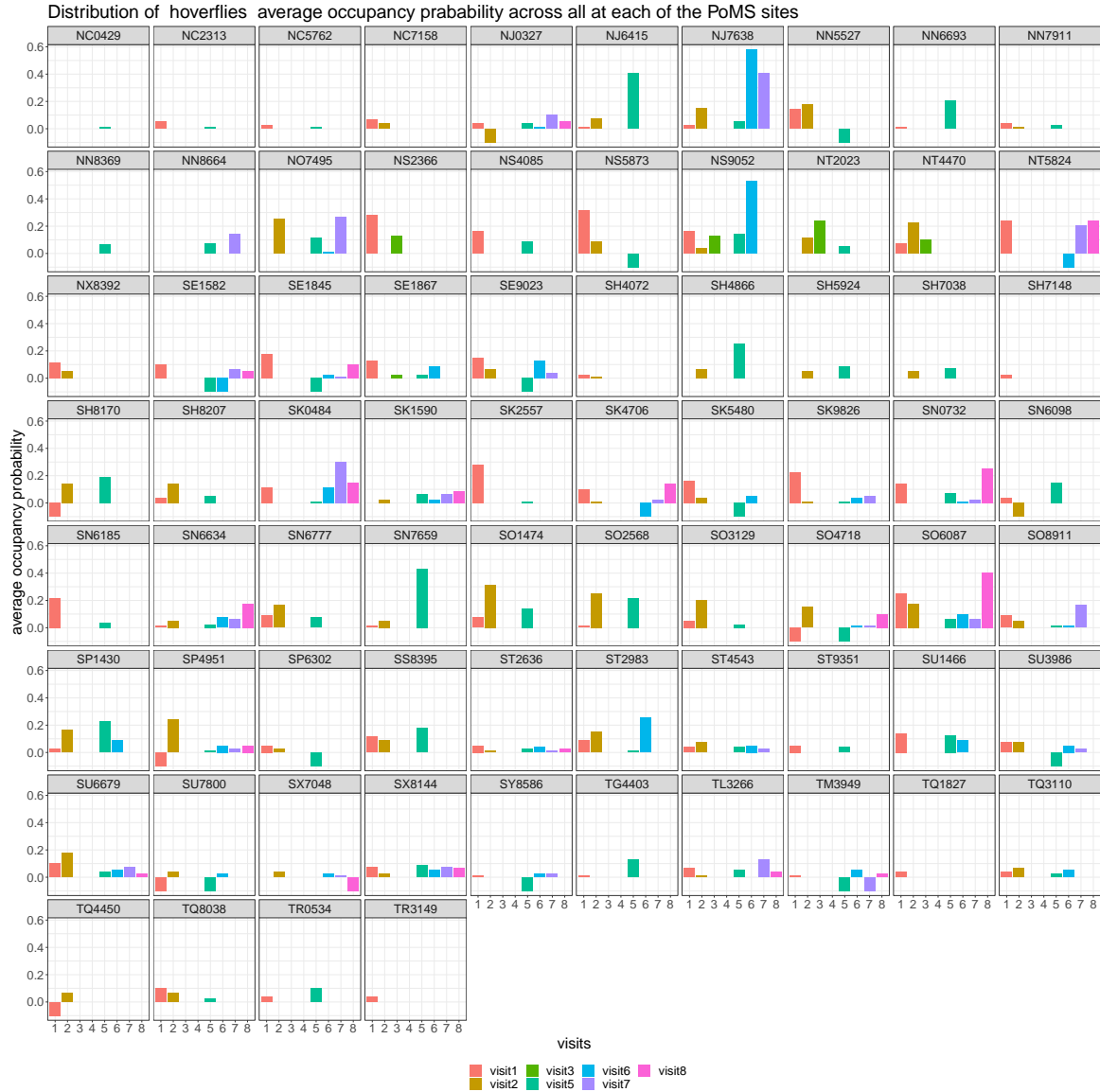


Figure S1-4: Distribution of hoverflies occupancy probability for the 74 PoMS sites. The probabilities for each site are estimated as the proportion of sites with at most one pantrap being occupied by the species. The occupancy probabilities are faceted by the PoMS site name and colored by the visit number the observations was made. The columns of visits without no bars represent the visit without no group count observation ('NA').



Figure S1-5: Distribution of solitarybees counts for the 74 PoMS sites. The counts are faceted by the PoMS site name and colored by the visit number the observations was made. The columns of visits without no bars represent the visit without no group count observation ('NA').



Figure S1-6: Distribution of solitarybees occupancy probability for the 74 PoMS sites. The probabilities for each site are estimated as the proportion of sites with at most one pantrap being occupied by the species. The occupancy probabilities are faceted by the PoMS site name and colored by the visit number the observations was made. The columns of visits without no bars represent the visit without no group count observation ('NA').

62 References

- 63 Douglas Bates, Martin Maechler, Ben Bolker, Steven Walker, Rune Haubo Bojesen Christensen,
64 Henrik Singmann, Bin Dai, Fabian Scheipl, and Gabor Grothendieck. Package ‘lme4’. URL
65 <http://lme4.r-forge.r-project.org>, 2009.
- 66 Tom D Breeze, Alison P Bailey, Kelvin G Balcombe, Tom Brereton, Richard Comont, Mike Edwards,
67 Michael P Garratt, Martin Harvey, Cathy Hawes, Nick Isaac, et al. Pollinator monitoring more
68 than pays for itself. *Journal of Applied Ecology*, 58(1):44–57, 2021.
- 69 Ludwig Fahrmeir, Thomas Kneib, Stefan Lang, and Brian D Marx. Regression models. In *Regression:*
70 *Models, methods and applications*, pages 23–84. Springer, 2022.
- 71 W. N. Venables and B. D. Ripley. *Modern Applied Statistics with S*. Springer, New York, fourth
72 edition, 2002. URL <https://www.stats.ox.ac.uk/pub/MASS4/>. ISBN 0-387-95457-0.