# Leveraging citizen science to assess richness, diversity, and abundance

Tim M. Szewczyk[a], Cleo Bertelsmeier[a], Tanja Schwander[a]

[a]*Department of Ecology and Evolution, University of Lausanne*

**Abstract**

Abstract.

*Keywords:* Keywords

## 1. Introduction

What structures ant species richness, diversity, and community structure at different spatial scales? We know that at a coarse scale, climate is generally important, with some support for phylogenetically conserved temperature preferences. At a local scale, richness typically decreases with canopy cover. In general, abundance seems to be more idiosyncratic and variable, both temporally and spatially. This is particularly true for small species. In ants, abundance can be measured as either the number of colonies in a particular area (i.e., colony density) or as the number of workers (i.e., worker density).

Increasingly, ecologists have access to occurrence data collected in various and haphazard ways, typically in the form of online databases or citizen science projects (e.g., GBIF, BEIN, etc). These data are commonly used for species distribution models (CITE), including for individual species and hierarchical models of species communities (Ovaskainen 2017), which can incorporate trait data, phylogeny, and covariation among species in addition to environmental drivers. However, they are a recent development and so far rely on a single set of occurrence data. However, use of such occurrence data for predicting richness or diversity directly has been somewhat more limited, though it is perhaps conceptually similar to extracting richness from guidebooks, which is fairly common (macroecology examples). There are several reasons for this. First, the data do not generally come from communities or assemblages, but rather an aggregation of detections from many different collectors across a variety of time spans. We like to think of diversity and richness as properties of communities, and these are decidedly not samples of communities,

obscuring the ability to detect or account for interactions among species. Second, the collections for each species may have differing spatial biases, rendering any simple aggregation methods erroneous. Third, there are biases in the species that are more likely to be detected, such that any estimate of richness or diversity will necessarily be of a subset of the community biased toward larger, more active, or more interesting species.

However, that doesn't mean these data can't be useful. Instead, the geographic breadth and rather indiscriminate collection methods can capture occurrences in unexpected locations or detect species that may be missed in alternative, more structured sampling methods. Leveraging the widely available occurrence data could thus clarify patterns of species distributions and diversity, and the drivers that shape them.

Here, we combine species occurrences of ants collected in a citizen science project in western Switzerland with species presence-absence data from a concurrent structured sampling effort. In a hierarchical Bayesian framework, we use the citizen science data to help inform species' responses to regional variables, and the structured samples to inform responses to both regional and local variables, while accounting for bias in geographic and taxonomic sampling effort in the citizen science data. With this model, we detail the patterns of ant colony density and species diversity across the landscapewhile incorporating uncertainty in species compositions, and we evaluate the support for hypothesized drivers across spatial scales. We compare inferences from the combined model with those from each dataset independently, and also assess the differences in observed communities from each sampling method.

## 2. Methods

This is a blurb about Vaud, including some about the climate, topography, as well as current and past human land use. Also a little bit about what is known about the ant fauna here, perhaps.

Ants were collected during the summer of 2019 in two simultaneous collection efforts (Fig. 1). First, a citizen science project involved the distribution of vials of ethanol to interested residents of the canton, who were asked to collect approximately 10 ant workers per colony for each vial. Participants were encouraged to explore under rocks, on bark, inside twigs, and in downed wood, and an online map was updated periodically to highlight data-sparse areas. Collectors returned the vials along with the locality of the sample (latitude and longitude, an address, or an annotated
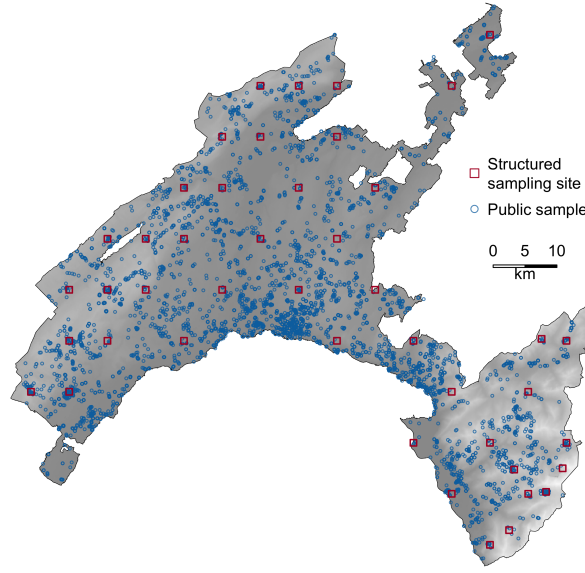
Figure 1: Map of the canton of Vaud, Switzerland. Samples consisted of presence-only data (blue points) and structured presence-absence data collected at 1 $km^2$ long-term biodiversity monitoring sites (red squares).

map), and often basic habitat information. A total of X,XXX samples were returned, representing collections between DATE and DATE. For these presence-only data, we discretized the landscape into a 1 $km^2$ grid (3,558 $km^2$) and tallied the number of counts for each species in each cell (X,XXX cells with occurrences).

Second, a structured sampling effort collected local presence-absence data, where samples of ants were collected within 44 sites of 1 $km^2$ each. Thirty-nine of the sites were a part of long-term biodiversity monitoring efforts by the federal government, and as such were arranged on a regular grid with *sim* 5-7 km between adjacent sites. Five additional sites are established monitoring sites by the University of Lausanne. The ants at each site were characterized by 25 plots, distributed among habitat types in approximate proportion to the abundance of each habitat, where each habitat type present within the site was represented by at least one plot. Inaccessible areas (e.g., cliffs, water, property beyond Vaud) were excluded, resulting in several sites with areas less than 1 $km^2$, and the number of plots was reduced proportionally (Appendix). Plots consisted of a 2m radius circle, with soil temperature recorded in the center, and vegetation characterized according to coverage classes within the plot (grass, forb, shrub, litter, bare, moss). Six flags were evenly spaced around the circumference, and within *sim* 20cm of each (total surface area *sim*$0.75m^2$), we

searched for ant colonies within any downed wood, under large rocks, and in *sim* 2L of soil, litter, and small rocks using Hori Hori gardening knives. We haphazardly collected *sim* 10 workers from each colony. All trees within the plot were also inspected for ant workers which were collected regardless of whether or not a colony was identifiable. Additionally, transects were mapped *a priori*, distributed proportionally among habitat types and totalling 2km. Workers were collected from all permanent mounds within $\approx$ 2m of the transect line.

There are many differences between W and Y, but they can be reduced down to three key points:

1. W is collected very broadly spatially, with no constraints within Vaud, while Y was collected in discrete, representative plots within pre-determined squares

2. Effort varies across space for W, but is standardized for Y

3. Species collections are biased in W toward larger, more active species, while Y provides a less biased survey of ground-nesting ants

The model described below accounts for each of these three differences, leveraging the information contained in W to improve estimates for the local diversity estimated by Y.

## 2.1. Model overview

To leverage the information contained in each dataset, we use a community-level hierarchical inhomogenous Poisson point process model (PPM). Inhomogenous Poisson PPMs assume that the distribution of occurrences is dependent on the variation in local intensity, $\lambda(x)$, across space, $x$, which may be observed imperfectly resulting in a thinned point process. One key benefit of PPMs is that the predicted point location intensity can be integrated to arbitrary spatial resolutions. Following the structure of our sampling design, we modelled the expected abundance of each species at two resolutions (coarse: 1 $km^2$, fine: 0.75 $m^2$), representing the area of the sampling sites and the area of the sampling plots respectively.

## 2.2. Model structure

We use two datasets: haphazard citizen science collections (**W**, res: 1 $km^2$), and structured samples (**Y**, 0.75 $m^2$ plots within 1 $km^2$ sites). Thus, for the citizen science dataset, we use only the coarse resolution (1 $km^2$), while the structured samples are incorporated at both resolutions. We assume that the colony density for each species is a function of local and regional environmental

4

conditions, with potential phylogenetic conservatism among species-specific responses. We esti-
mate the effects of regional conditions using both **Y** and **W**, and the effects of local conditions
using **Y**.

At a local scale, $Y_{is}$ represents the number of colonies of species $s$ detected at plot $i$ of site
$j$, which is an estimate of the true number of colonies $\lambda_{is}$, which in turn depends on the local
environment and the overall abundance of species $s$ at site $j$:

$$Y_{is} \sim Poisson(\lambda_{is}) \tag{1}$$

$$log(\lambda_{is}) = a_s V_i + log(h\Lambda_{js}) \tag{2}$$

where $V_i$ is a vector of local environmental covariates, $a_s$ is a vector of species-specific responses,
$h$ is a constant representing the proportion of site $j$ sampled by plot $i$ ($0.75\ m^2 / 1\ km^2 = 7.5e\text{-}7$), and
$\Lambda_{js}$ is the total expected abundance of species $s$ at site $j$. Thus, $log(h\Lambda_{js})$ functions as an intercept,
determining the baseline expected abundance at each plot within each site, which are then further
refined according to the effects of the local environment.

At the site scale, the total abundance within a $1km^2$ cell is observed indirectly by $Y_{is}$ as describe
above, and directly by $W_{ks}$. The species counts from the unstructured survey in cell $k$ are thinned
such that:

$$W_{ks} \sim Multinomial(\Lambda_{ks} D_s) \tag{3}$$

where $D_s$ is the detection bias for each species, which accounts for bias in the community com-
position based on species that are more readily observed. Thus, for each species $s$ in cell $k$, the
observed count $W_{ks}$ for species $s$ is expected to be higher if the species is relatively abundant in cell
$k$ ($\Lambda_{ks}$ is larger), more samples were collected in $k$ , or if species $s$ is likely to be over-represented
in haphazard collections by non-experts ($D_s$ is larger).

Thus, both $\Lambda_{js}$ and $\Lambda_{ks}$ represent the expected abundance of species $s$ per $km^2$, but the exact
cells are not required to overlap. The ecological processes driving $\Lambda_{(jk)s}$ are assumed to be the
same, however, such that:

$$log(\Lambda_{(j,k)s}) = b_s X_{(j,k)} \tag{4}$$

where $b_s$ is a vector of species-specific slopes, and **X** is a matrix of environmental covariates.

5

The slopes $a_s$ and $b_s$ are species-specific responses at $0.75m^2$ and $1km^2$ resolutions, respectively, and are distributed about genus-level means $A_g$ and $B_g$ with standard deviations $\sigma^a$ and $\sigma^b$. The genus-level means are in turn distributed about aggregate means, $\alpha$ and $\beta$ with covariance matrices $\Sigma^A$ and $\Sigma^B$, which reflect the overall responses of the ant community to environmental variables at each resolution while accounting for phylogenetic relatedness at the genus level.

Finally, we calculated several quantities at the site-scale and plot-scale. Following Hefley & Hooten (2016), we calculated the probability of presence as $\psi_i s = 1 - e^{-\lambda_{is}}$ and $\Psi_j s = 1 - e^{-\Lambda_{js}}$, as well as Shannon's H at each plot, $\beta$ diversity among plots within each site, and among sites, and overall $\gamma$ diversity at each site.

| Parameter | Description | Type |
|---|---|---|
| $i$ | structured sampling plots ($0.75\ m^2$) | index |
| $j$ | structured sampling cells ($1\ km^2$) | index |
| $k$ | citizen science cells ($1\ km^2$) | index |
| $s$ | species | index |
| $g$ | genus | index |
| $l$ | local covariates ($0.75\ m^2$) | index |
| $r$ | regional covariates ($1\ km^2$) | index |
| $\mathbf{Y}_{is}$ | structured sampling counts ($0.75\ m^2$) | data |
| $\mathbf{W}_{ks}$ | citizen science counts ($1\ km^2$) | data |
| $\mathbf{V}_{il}$ | local covariates ($0.75\ m^2$) | data |
| $\mathbf{X}_{(jk)r}$ | regional covariates ($1\ km^2$) | data |
| $h$ | structured sampling proportional effort | data |
| $\lambda_{is}$ | colony intensity ($0.75\ m^2$) | latent |
| $\mathbf{\Lambda}_{(jk)s}$ | colony intensity ($1\ km^2$) | latent |
| $\alpha_l$ | aggregate ant responses ($0.75\ m^2$) | slopes |
| $\mathbf{A}_{lg}$ | genus-level ant responses ($0.75\ m^2$) | slopes |
| $\mathbf{a}_{ls}$ | species-level ant responses ($0.75\ m^2$) | slopes |
| $\sigma_l^a$ | response sd among congeners | sd |
| $\mathbf{\Sigma^A}_{gg}$ | genus-level covariance matrix | cov mx |
| $\beta_r$ | aggregate ant responses ($1\ km^2$) | slopes |
| $\mathbf{B}_{rg}$ | genus-level ant responses ($1\ km^2$) | slopes |
| $\mathbf{b}_{rs}$ | species-level ant responses ($1\ km^2$) | slopes |
| $\sigma_r^b$ | response sd among congeners | sd |
| $\mathbf{\Sigma^B}_{gg}$ | genus-level covariance matrix | cov mx |
| $D_s$ | citizen science species bias (proportional) | random effect |

Table 1: Parameters in the model. Could be moved to an appendix, or shortened since many of these don't need to be highlighted.
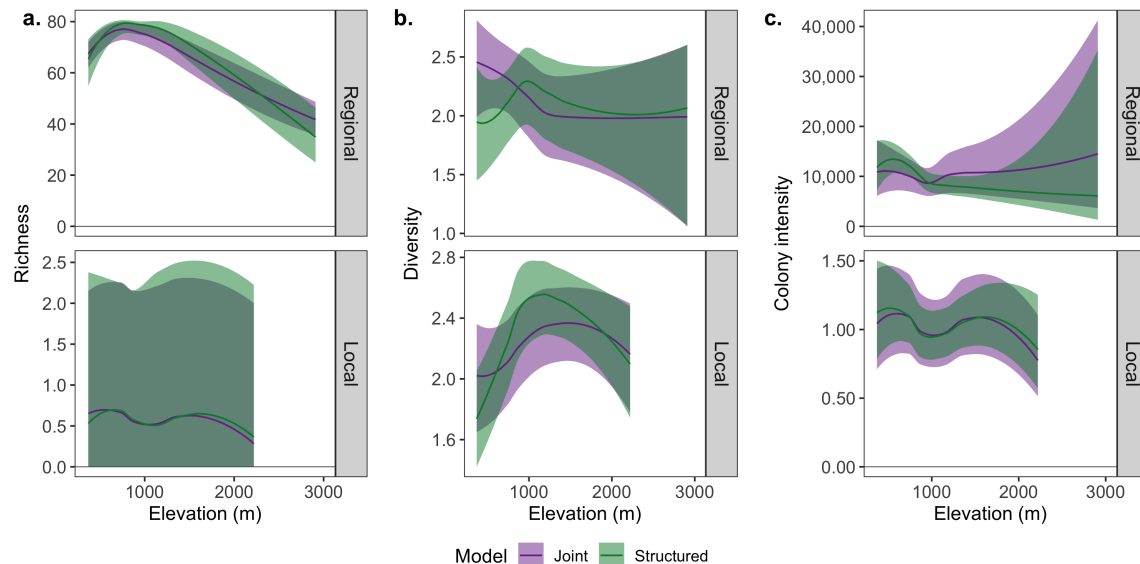
Figure 2: Elevational patterns of posterior distributions at regional and local scales for (a) predicted species richness, (b) predicted Shannon H diversity, and (c) colony intensity. Lines and ribbons are loess lines using the posterior medians and 95% Highest Density Interval, respectively, for the model using both datasets (Joint: purple) and the model using only the presence-absence data (Structured: green). All metrics were calculated in each sample from the posterior distribution using the predicted species-level colony intensities.

*2.3. Validation and model selection*

To evaluate the effect of including **W** and to perform variable selection, we compared a version of the model with and without **W**. We randomly divided the 44 sites with data for **Y** into training (75%) and testing (25%) subsets (**Appendix 1**). For each combination of covariates, we parameterized the models with the training subset and evaluated the ability to predict $\lambda$ with the testing subset using the log predictive density, a metric that uses the full posterior predictive distribution.

## 3. Results

We used simulated data to validate the model, testing versions of the model using only **W**, only **Y**, and both **W** and **Y**. The model using both datasets performs best at predicting both site-level and plot-level abundances for testing subsets. I don't have any statistics yet, and will need to do full runs to be more confident. It is, of course, unsurprising that the site-level abundance $\Lambda$ is predicted more accurately, as the use of both datasets effectively increases the sample size for estimating the site-level slopes. The more accurate prediction of $\lambda$ is more exciting, since **W** doesn't directly include any information about it, but rather helps indirectly by improving predictions of $\Lambda$ at each site, which helps to constrain the intercept in the plot-level equation.
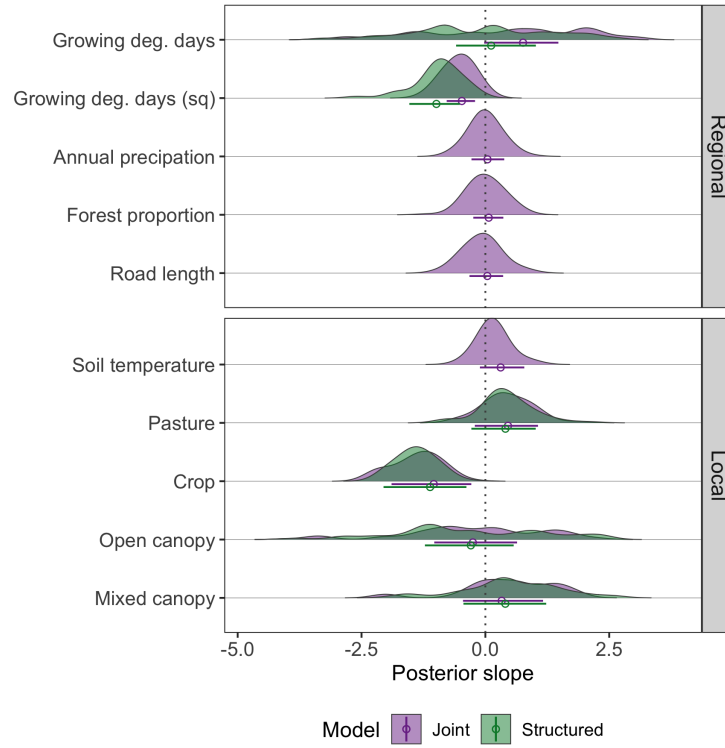
7

Figure 3: Distribution of species-level responses in optimal models. Density curves represent the distribution of species-level posterior medians for the responses to local and regional variables for the model using both datasets (Joint: purple) and the model using only the presence-absence data (Structured: green). Points and lines show the posterior median and 95% Highest Density Interval for the aggregate ($\beta$) responses.
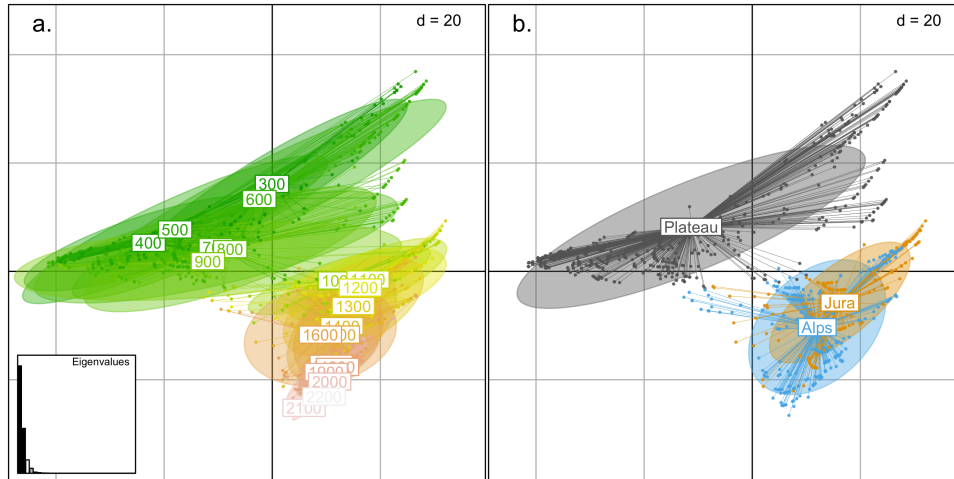


Figure 4: Double principle coordinate analysis (DPCoA) of predicted local communities, with plots colored by (a) elevational bin, and (b) region. The central plateau includes hills from 300m to 1000m, with the Jura and the Alps rising steeply in the east and west, respectively.
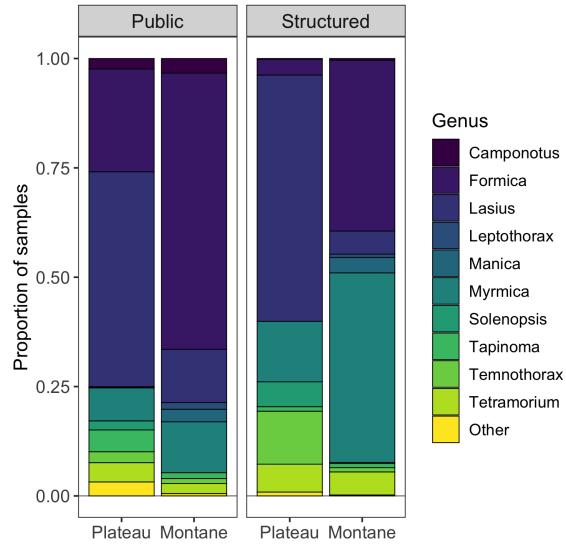
Figure 5: Relative abundance of genera in the presence-only (public) and presence-absence (structured) datasets across plateau and montane environments. Only genera that constitute ≥ 1% of at least one subset are shown, with all others indicated as 'Other'.
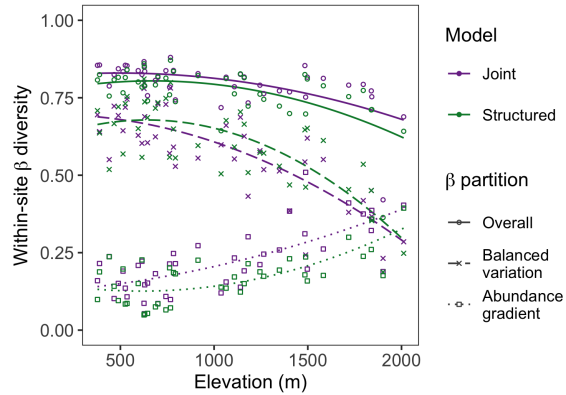


Figure 6: Multi-site $\beta$-diversity and its components across elevation. Within each $1km^2$ site, the species-level colony intensity posterior medians were used to calculate a multi-site $\beta$-diversity index and its components, representing the variation among local plots. 'Balanced variation' represents changes in relative abundance among species, while the 'abundance gradient' partition denotes changes in total abundance.

9

## 4. Discussion

Discussion goes here.

Conclusion

## 5. Acknowledgments

Thanks to everyone.

## 6. Appendixes

**Appendix 1.** Supplementary methods.

**Appendix 2.** Model code.

**Full equations**

$$Y_{is} \sim Poisson(\lambda_{is}) \tag{5}$$

$$log(\lambda_{is}) = a_s V_i + log(h\Lambda_{js}) \tag{6}$$

$$W_{ks} \sim Poisson(\Lambda_{ks}E_k D_s) \tag{7}$$

$$log(\Lambda_{(j,k)s}) = b_s X_{(j,k)} \tag{8}$$

$$logit(E_k) = \eta U_k \tag{9}$$

$$\begin{aligned} a_s &\sim Norm(A_g, \sigma^a) \\ A_g &\sim Norm_G(\alpha, \Sigma^A) \end{aligned} \tag{10}$$

$$\begin{aligned} b_s &\sim Norm(B_g, \sigma^b) \\ B_g &\sim Norm_G(\beta, \Sigma^B) \end{aligned} \tag{11}$$

## 7. Bibliography

## References