

Sampling networks of ecological interactions

Pedro Jordano^{*a}

^aIntegrative Ecology Group, Estación Biológica de Doñana, Consejo
Superior de Investigaciones Científicas (EBD-CSIC), Avenida
Americo Vespucio s/n, E-41092 Sevilla, Spain

*Corresponding author: jordano@ebd.csic.es

Summary

1. Sampling ecological interactions presents similar challenges, problems, potential biases, and constraints as sampling individuals and species in biodiversity inventories. Interactions are just pairwise relationships among individuals of two unrelated species, such as those among plants and their seed dispersers in frugivory interactions or those among plants and their pollinators. Sampling interactions is a fundamental step to build robustly estimated interaction networks, yet few analyses have attempted a formal approach to their sampling protocols.

*jordano@ebd.csic.es

2. Robust estimates of the actual number of interactions (links) within diversified ecological networks require adequate sampling effort that needs to be explicitly gauged. Yet we still lack a sampling theory explicitly focusing on ecological interactions.
3. While the complete inventory of interactions is likely impossible, a robust characterization of its main patterns and metrics is probably realistic. We must acknowledge that a sizable fraction of the maximum number of interactions I_{max} among, say, A animal species and P plant species (i.e., $I_{max} = AP$) is impossible to record due to forbidden links, the restrictions imposed by the organisms life-histories. Thus, the number of observed interactions I in robustly sampled networks is typically $I \ll I_{max}$, resulting in extremely sparse interaction matrices with low connectance.
4. Reasons for forbidden links are multiple but mainly stem from spatial and temporal uncoupling of partner species encounters and from intrinsically low probabilities of interspecific encounter for many of the potential pairwise interactions. Adequately assessing the completeness of a network of ecological interactions thus needs a deep knowledge of the natural history details embedded, so that forbidden links can be “discounted” when addressing sampling effort.
5. Here I provide a review and outline a conceptual framework for interaction sampling by building an explicit analogue to individuals and species sampling, thus extending diversity-monitoring approaches to the characterization of complex networks of ecological interactions. This is crucial to assess the fast-paced and devastating effects of defaunation-driven loss of key ecological

interactions and the services they provide.

Keywords

complex networks, food webs, frugivory, mutualism, plant-animal interactions, pollination, seed dispersal

Introduction

Biodiversity sampling is a labour-intensive activity, and sampling is often not sufficient to detect all or even most of the species present in an assemblage. Gotelli & Colwell (2011).

1 Biodiversity assessment aims at sampling individuals in collections and deter-
2 mining the number of species represented. Given that, by definition, samples are
3 incomplete, these collections enumerate a lower number of the species actually
4 present. The ecological literature dealing with robust estimators of species rich-
5 ness and diversity in collections of individuals is immense, and a number of useful
6 approaches have been used to obtain such estimates (Colwell, 2009; ?). Re-
7 cent effort has been also focused at defining essential biodiversity variables (EBV)
8 (?) that can be sampled and measured repeatedly to complement biodiversity
9 estimates. Yet sampling species or taxa-specific EBVs is just probing a single
10 component of biodiversity; interactions among species are another fundamental
11 component, the one that supports the existence of species (?). For example, the
12 extinction of interactions represents a dramatic loss of biodiversity because it en-

13 tails the loss of fundamental ecological functions (?). This missed component of
14 biodiversity loss, the extinction of ecological interactions, very often accompanies,
15 or even precedes, species disappearance. Interactions among species are a key
16 component of biodiversity and here I aim to show that most problems associated
17 to sampling interactions in natural communities have to do with problems associ-
18 ated to sampling species diversity. I consider pairwise interactions among species
19 at the habitat level, in the context of alpha diversity and the estimation of local
20 interaction richness from sampling data (?). In the first part I provide a succinct
21 overview of previous work addressing sampling issues for ecological interaction net-
22 works. In the second part I discuss specific rationales for sampling the biodiversity
23 of ecological interactions.

24 Interactions can be a much better indicator of the richness and diversity of
25 ecosystem functions than a simple list of taxa and their abundances and/or re-
26 lated biodiversity indicator variables (EBVs). Thus, sampling interactions should
27 be a central issue when identifying and diagnosing ecosystem services (e.g., polli-
28 nation, natural seeding by frugivores, etc.). Fortunately, all the whole battery of
29 biodiversity-related tools used by ecologists to sample biodiversity (species, *sensu*
30 *stricto*) can be extended and applied to the sampling of interactions. Analogs
31 are evident between these approaches (?). Monitoring interactions is analogous to
32 any biodiversity sampling (i.e., a species inventory (??)) and is subject to similar
33 methodological shortcomings, especially under-sampling (Coddington *et al.*, 2009;
34 ?; ?; ?). For example, when we study mutualistic networks, our goal is to make an
35 inventory of the distinct pairwise interactions that made up the network. We are
36 interested in having a complete list of all the pairwise interactions among species
37 (e.g., all the distinct, species-species interactions, or links, among the pollinators

38 and flowering plants) that can exist in a given community. Sampling these in-
39 teractions thus entails exactly the same problems, limitations, constraints, and
40 potential biases as sampling individual organisms and species diversity. As Mao
41 & Colwell (?) put it, these are the workings of Preston’s demon, the moving “veil
42 line” between detected and the undetected interactions as sample size increases
43 (?).

44 Early efforts to recognize and solve sampling problems in analyses of interac-
45 tions stem from researchers interested in food web analyses and in determining
46 the biases of undersampled food web metrics (?Cohen *et al.*, 1993; ?; Bersier,
47 Banasek-Richter & Cattin, 2002; Brose, Martinez & Williams, 2003; Banasek-
48 Richter, Cattin & Bersier, 2004). In addition, the myriad of classic natural history
49 studies documenting animal diets, host-pathogen infection records, plant herbivory
50 records, etc., represent efforts to document interactions occurring in nature. All
51 of them share the problem of sampling incompleteness influencing the patterns
52 and metrics reported. Yet, despite the early recognition that incomplete sampling
53 may seriously bias the analysis of ecological networks (?), only recent studies have
54 explicitly acknowledged it and attempted to determine its influence (?????Chacoff
55 *et al.*, 2012; ?; ?; Bascompte & Jordano, 2014; ?; ?). The sampling approaches
56 have been extended to predict patterns of coextinctions in interaction assemblages
57 (e.g., hosts-parasites) (?). Most empirical studies provide no estimate of sam-
58 pling effort, implicitly assuming that the reported network patterns and metrics
59 are robust. Yet recent evidences point out that number of partner species de-
60 tected, number of actual links, and some aggregate statistics describing network
61 patterns, are prone to sampling bias (???Chacoff *et al.*, 2012; ?; ?; ?). Most of
62 these evidences, however, come from either theoretical, simulation, studies (?) or

63 from relatively species-poor assemblages. Even for species-rich, tropical assem-
64 blages it might be erroneous to conclude that network data routinely come from
65 insufficiently sampled datasets (Chacoff *et al.*, 2012), given the extremely sparse
66 nature of these interaction matrices because of the prevalence of forbidden links
67 (which, by definition, cannot be documented despite extensive sampling effort).
68 However, most certainly, sampling limitations pervade biodiversity inventories in
69 tropical areas (Coddington *et al.*, 2009) and we might rightly expect that frequent
70 interactions may be over-represented and rare interactions may be missed entirely
71 in studies of mega-diverse assemblages (Bascompte & Jordano, 2014); but, to what
72 extent?

73 Sampling interactions: methods

74 When we sample interactions in the field we record the presence of two species
75 that interact in some way. For example, Snow and Snow(?) recorded an inter-
76 action whenever they saw a bird “touching” a fruit on a plant. In a similar way,
77 interactions between pollinators and plants are tallied by recording any visit of a
78 pollinator entering a flower and touching the reproductive parts. We observe and
79 record feeding observations, visitation, occupancy, presence in pollen loads or in
80 fecal samples, etc., of *individual* animals or plants and accumulate pairwise inter-
81 actions, i.e., lists of species partners and the frequencies with which we observe
82 them. Therefore, estimating the sampling completeness of pairwise interactions
83 for a whole network, requires estimating the number (richness) of distinct pairwise
84 interactions accumulated as sampling effort is increased, pooling the data for all
85 partner species. Most, if not all, types of ecological interactions can be illustrated

86 by bipartite graphs, with two or more distinct groups of interacting partners (Bas-
87 compte & Jordano, 2014); for illustration purposes I'll focus more specifically on
88 plant-animal interactions.

89 Sampling interactions requires filling the cells of an interaction matrix with
90 data. The matrix, $\Delta = AP$, is a 2D representation of the interactions among, say,
91 A animal species (rows) and P plant species (columns) (Bascompte & Jordano,
92 2014). An interaction matrix Δ consists of an array of zeroes or ones, or an
93 array of numeric values (including zeroes)- if the data (interaction frequencies) are
94 quantified. The matrix entries illustrate the values of the pairwise interactions
95 visualized in the Δ matrix, and can be 0 or 1, for presence-absence of a given
96 pairwise interaction, or take a quantitative weight w_{ji} to represent the interaction
97 intensity or unidirectional effect of species j on species i (Bascompte & Jordano,
98 2014; ?). Given that the outcomes of most ecological interactions are dependent
99 on frequency of encounters (e.g., visit rate of pollinators, number of records of
100 ant defenders, frequency of seeds in fecal samples), a frequently used proxy for
101 interaction intensities w_{ji} is just how frequent are new interspecific encounters,
102 whether or not appropriately weighted to estimate interaction effectiveness (?).

103 We need to define two basic steps in the sampling of interactions: 1) which type
104 of interactions we sample; and 2) which type of record we get to document the
105 existence of an interaction. In step #1 there are two considerations we need to take
106 into account. First is whether we are sampling the whole community of interactor
107 species (all the animals, all the plants) or we sample just a subset of them, i.e.,
108 a sub matrix $\Delta_{m,n}$ of $m < A$ animal species and $n < P$ plant species of the
109 adjacency matrix Δ_{AP} . Subsets can be: a) all the potential plants interacting with
110 a subset of the animals (Fig. 1a); b) all the potential animal species interacting

111 with a subset of the plant species (Fig. 1b); c) a subset of all the potential
112 animal species interacting with a subset of all the plant species (Fig. 1c). While
113 some discussion has considered how to establish the limits of what represents
114 a network (?) (in analogy to discussion on food-web limits (Cohen, 1978)), it
115 must be noted that situations a-c in Fig. 1 do not represent complete interaction
116 networks. As vividly stated by Cohen et al. (Cohen *et al.*, 1993): “*As more*
117 *comprehensive, more detailed, more explicit webs become available, smaller, highly*
118 *aggregated, incompletely described webs may progressively be dropped from analyses*
119 *of web structure (though such webs may remain useful for other purposes, such as*
120 *pedagogy)*”. Subnet sampling is generalized in studies of biological networks (e.g.,
121 protein interactions, gene regulation), yet it is important to recognize that most
122 properties of subnetworks (even random subsamples) do not represent properties
123 of whole networks (?).

124 0.0.1 Fig. 1 here

125 In step #2 above we face the problem of the type of record we take to sample
126 interactions. This is important because it defines whether we approach the problem
127 of filling up the interaction matrix in a “zoo-centric” way or in a “phyto-centric”
128 way. Zoo-centric studies directly sample animal activity and document the plants
129 ‘touched’ by the animal. For example, analysis of pollen samples recovered from the
130 body of pollinators, analysis of fecal samples of frugivores, radio-tracking data, etc.
131 Phyto-centric studies take samples of focal individual plant species and document
132 which animals ‘arrive’ or ‘touch’ the plants. Examples include focal watches of
133 fruiting or flowering plants to record visitation by animals, raising insect herbivores
134 from seed samples, identifying herbivory marks in samples of leaves, etc.

Most recent analyses of plant-animal interaction networks are phyto-centric; just 3.5% of available plant-pollinator ($N=58$) or 36.6% plant-frugivore ($N=22$) interaction datasets are zoo-centric (see (?)). Moreover, most available datasets on host-parasite or plant-herbivore interactions are “host-centric” or phyto-centric (e.g., (??)). This maybe related to a variety of causes, like preferred methodologies by researchers working with a particular group or system, logistic limitations, or inherent taxonomic focus of the research questions. A likely result of phyto-centric sampling would be adjacency matrices with large $A : P$ ratios. In any case we don’t have a clear view of the potential biases that taxa-focused sampling may generate in observed network patterns, for example by generating consistently asymmetric interaction matrices (?). System symmetry has been suggested to influence estimations of generalization levels in plants and animals when measured as I_A and I_P (?); thus, differences in I_A and I_P between networks may arise from different $A : P$ ratios rather than other ecological factors (?).

Interestingly enough, quite complete analyses of interaction networks can be obtained when combining both phyto-centric and zoo-centric sampling. For example, Bosch et al. (Bosch *et al.*, 2009) showed that the addition of pollen load data on top of focal-plant sampling of pollinators unveiled a significant number of interactions, resulting in important network structural changes. Connectance increased 1.43-fold, mean plant connectivity went from 18.5 to 26.4, and mean pollinator connectivity from 2.9 to 4.1; moreover, extreme specialist pollinator species (singletons in the adjacency matrix) decreased 0.6-fold. Zoo-centric sampling has recently been extended with the use of DNA-barcoding, for example with plant-herbivore (?) and plant-frugivore interactions (?). For mutualistic networks we would expect that zoo-centric sampling could help unveiling interactions for

160 rare species or for relatively common species which are difficult to sample by di-
 161 rect observation. Future methodological work may provide significant advances
 162 showing how mixing different sampling strategies strengthens the completeness of
 163 network data. These mixed strategies may combine, for instance, focal analyses,
 164 pollen load or seed contents, camera traps, and DNA barcoding records. We might
 165 expect increased power of these mixed sampling approaches when combining dif-
 166 ferent methods from both phyto- and zoo-centric perspectives (Bosch *et al.*, 2009;
 167 Bluthgen, 2010).

168 **Sampling interactions: rationale**

169 The number of distinct pairwise interactions that we can record in a landscape (an
 170 area of relatively homogeneous vegetation, analogous to the one we would use to
 171 monitor species diversity) is equivalent to the number of distinct classes in which
 172 we can classify the recorded encounters among individuals of two different species.
 173 Yet, individual-based plant-animal interaction networks have been only recently
 174 studied (?). We walk in the forest and see a blackbird Tm picking an ivy Hh fruit
 175 and ingesting it: we have a record for $Tm-Hh$ interaction. We keep advancing and
 176 record again a blackbird feeding on hawthorn Cm fruits so we record a $Tm-Cm$
 177 interaction; as we advance we encounter another ivy plant and record a blackcap
 178 swallowing a fruit so we now have a new $Sa-Hh$ interaction, and so on. At
 179 the end we have a series of classes (e.g., $Sa-Hh$, $Tm-Hh$, $Tm-Cm$, etc.),
 180 along with their observed frequencies. Bunge & Fitzpatrick (Bunge & Fitzpatrick,
 181 1993) review the main aspects and approaches to estimate the number of distinct
 182 classes C in a sample of observations. The sampling of interactions in nature, as

the sampling of species, is a cumulative process. In our analysis, we are not re-sampling individuals, but interactions, so we made interaction-based accumulation curves. If an interaction-based curve points towards a robust sampling, it does mean that no new interactions are likely to be recorded, irrespectively of the species, as it is a whole-network sampling approach (N. Gotelli, pers. com.). We add new, distinct, interactions recorded as we increase sampling effort (Fig. 2). We can obtain an Interaction Accumulation Curve (*IAC*) analogous to a Species cumulating Curve (*SAC*): the observed number of distinct pairwise interactions in a survey or collection as a function of the accumulated number of observations or samples (Colwell, 2009).

0.0.2 Fig. 2 here

Our sampling above would have resulted in a vector $n = [n_1 \dots n_C]'$ where n_i is the number of records in the i^{th} class. As stressed by Bunge & Fitzpatrick (Bunge & Fitzpatrick, 1993), however, the i^{th} class would appear in the sample if and only if $n_i > 0$, and we don't know *a priori* which n_i are zero. So, n is not observable. Rather, what we get is a vector $c = [c_1 \dots c_n]'$ where c_j is the number of classes represented j times in our sampling: c_1 is the number of singletons, c_2 is the number of twin pairs, c_3 the number of triplets, etc. The problem thus turns to be estimating the number of distinct classes C from the vector of c_j values.

Estimating the number of interactions with resulting robust estimates of network parameters is a central issue in the study of ecological interaction networks (?Bascompte & Jordano, 2014). In contrast with traditional species diversity estimates, sampling networks has the paradox that despite the potentially interacting species being present in the sampled assemblage (i.e., included in the A and P

species lists), some of their pairwise interactions are impossible to be recorded. The reason is forbidden links. Independently of whether we sample full communities of subset communities we face a problem: some of the interactions that we can visualize in the empty adjacency matrix Δ with size AP will simply not occur. Thus, independently of the sampling effort we put, we'll never document these pairwise interactions. With a total of AP "potential" interactions, a fraction of them are impossible to record, because they are forbidden (??). Forbidden links are constraints for the establishment of new links, and mainly arise from the biological attributes of the species: no link can be established between a plant and an animal mutualist differing in phenology, i.e. the seeds of a winter-ripening plant cannot be dispersed by a frugivore that is a summer stopover migrant (?). Or, for instance, short-tongued pollinators cannot successfully reach the nectar in long-corolla flowers and pollinate them efficiently (?). Forbidden links are thus represented as structural zeroes in the interaction matrix, i.e., matrix cells that cannot get a non-zero value. So, we need to account for the frequency of these structural zeros in our matrix before proceeding. For example, most measurements of connectance ($C = I/(AP)$) implicitly ignore the fact that by taking the full product AP in the denominator they are underestimating the actual connectance value, i.e., the fraction of actual interactions I relative to the *biologically possible* ones, not to the total maximum $I_{max} = AP$.

0.0.3 Table 1 approx. here

Adjacency matrices are frequently sparse, i.e., they are densely populated with zeroes, with a fraction of them being structural (i.e., unobservable interactions) (Bascompte & Jordano, 2014). It would be thus a serious interpretation error to

231 attribute the sparseness of adjacency matrices for bipartite networks to under-
 232 sampling. The actual typology of link types in ecological interaction networks is
 233 thus more complex than just the two categories of observed and unobserved inter-
 234 actions (Table 1). Unobserved interactions are represented by zeroes and belong
 235 to two categories. Missing interactions may actually exist but require additional
 236 sampling or a variety of methods to be observed. Forbidden links, on the other
 237 hand, arise due to biological constraints limiting interactions and remain unob-
 238 servable in nature, irrespectively of sampling effort (??). Forbidden links are non-
 239 occurrences of pairwise interactions that can be accounted for by biological con-
 240 straints, such as spatio-temporal uncoupling, size or reward mismatching, foraging
 241 constraints (e.g., accessibility), and physiological-biochemical constraints (?). We
 242 still have extremely reduced information about the frequency of forbidden links
 243 in natural communities (???????) (Table 1). Forbidden links FL may actually
 244 account for a relatively large fraction of unobserved interactions UL when sam-
 245 pling taxonomically-restricted subnetworks (e.g., plant-hummingbird pollination
 246 networks) (Table 1). Phenological unmatching is also prevalent in most networks,
 247 and may add up to explain ca. 25–40% of the forbidden links, especially in highly
 248 seasonal habitats, and up to 20% when estimated relative to the total number
 249 of unobserved interactions (Table 2). In any case, we might expect that a frac-
 250 tion of the missing links ML would be eventually explained by further biological
 251 reasons, depending on the knowledge of natural details of the particular systems.
 252 Our goal as naturalists would be to reduce the fraction of UL which remain as
 253 missing links; to this end we might search for additional biological constraints or
 254 added sampling effort. For instance, habitat use patterns by hummingbirds in the
 255 Aroma Valley network (Table 2; (?)) impose a marked pattern of microhabitat

256 mismatches causing up to 44.5% of the forbidden links. There are a myriad of
 257 biological causes beyond those included as *FL* in Table 2 that may contribute
 258 explanations for *UL*: limits of color perception and or partial preferences, pres-
 259 ence of secondary metabolites in fruit pulp and leaves, toxins and combinations of
 260 monosaccharides in nectar, etc. However, it is surprising that just the limited set
 261 of forbidden link types in Table 1 explain between 24.6–77.2% of the unobserved
 262 links. Notably, the Arima Valley, Santa Virgínia, and Hato Ratón networks have
 263 > 60% of the unobserved links explained, which might be related to the fact that
 264 they are subnetworks (Arima Valley, Santa Virgínia) or relatively small networks
 265 (Hato Ratón). All this means that empirical networks may have sizable fractions
 266 of structural zeroes. Ignoring this biological fact may contribute to wrongly infer
 267 undersampling of interactions in real-world assemblages.

268 **0.0.4 Table 2 approx. here**

269 To sum up, two elements of inference are required in the analysis of unobserved
 270 interactions in ecological interaction networks: first, detailed natural history infor-
 271 mation on the participant species that allows the inference of biological constraints
 272 imposing forbidden links, so that structural zeroes can be identified in the adja-
 273 cency matrix; second, a critical analysis of sampling robustness a robust estimate
 274 of the actual fraction of missing links, M , and thus, a robust estimate of I .

275 **Asymptotic diversity estimates**

Let's assume a sampling of the diversity in a specific locality, over relatively ho-
 mogeneous landscape where we aim at determining the number of species present

for a particular group of organisms. To do that we carry out transects or plot samplings across the landscape, adequately replicated so we obtain a number of samples. Briefly, S_{obs} is the total number of species observed in a sample, or in a set of samples. S_{est} is the estimated number of species in the community represented by the sample, or by the set of samples, where *est* indicates an estimator. With abundance data, let S_k be the number of species each represented by exactly k individuals in a single sample. Thus, S_0 is the number of undetected species (species present in the community but not included in the sample), S_1 is the number of singleton species, S_2 is the number of doubleton species, etc. The total number of individuals in the sample would be:

$$n = \sum_{k=1}^{S_{obs}} S_k$$

276

277 A frequently used asymptotic, bias corrected, non-parametric estimator is S_{Chao}
 278 (?Chao, 2005; Colwell, 2013):

$$S_{Chao} = S_{obs} + \frac{S_1(S_1 - 1)}{2(S_2 + 1)}$$

279 Another frequently used alternative is the Chao2 estimator, S_{Chao2} (?), which
 280 has been reported to have a limited bias for small sample sizes (Colwell & Cod-
 281 dington, 1994; Chao, 2005):

$$S_{Chao2} = S_{obs} + \frac{S_1^2}{2S_2}$$

282 A plot of the cumulative number of species recorded, S_n , as a function of some

283 measure of sampling effort (say, n samples taken) yields the species accumulation
 284 curve (SAC) or collector's curve (Colwell & Coddington, 1994). Such a curve
 285 eventually reaches an asymptote converging with S_{est} . In an analogous way, inter-
 286 action accumulation curves (IAC), analogous to SACs, can be used to assess the
 287 robustness of interactions sampling for plant-animal community datasets (???).
 288 For instance, a random accumulator function (e.g., library `vegan` in the R Package
 289 (?)) which finds the mean IAC and its standard deviation from random permuta-
 290 tions of the data, or subsampling without replacement (?) can be used to estimate
 291 the expected number of distinct pairwise interactions included in a given sampling
 292 of records (??). We start with a vectorized interaction matrix representing the
 293 pairwise interactions (rows) recorded during a cumulative number of censuses or
 294 sampling periods (columns) (Table 3) , in a way analogous to a biodiversity sam-
 295 pling matrix with species as rows and sampling units (e.g., quadrats) as columns
 296 (?). In this way we effectively extend sampling theory developed for species di-
 297 versity to the sampling of ecological interactions. Yet future theoretical work will
 298 be needed to formally assess the similarities and differences in the two approaches
 299 and developing biologically meaningful null models of expected interaction richness
 300 with added sampling effort.

301 **Assessing sampling effort when recording interac-** 302 **tions**

303 The basic method we can propose to estimate sampling effort and explicitly show
 304 the analogues with rarefaction analysis in biodiversity research is to vectorize the

305 interaction matrix AP so that we get a vector of all the potential pairwise inter-
 306 actions (I_{max} , Table 1) that can occur in a community of A animal species and
 307 P plant species. The new “species” we aim to sample are the pairwise interac-
 308 tions (Table 3). So, if we have in our community *Turdus merula* (Tm) and *Rosa*
 309 *canina* (Rc) and *Prunus mahaleb* (Pm), our problem will be to sample 2 new
 310 “species”: $Tm - Rc$ and $Tm - Pm$. In general, if we have $A = 1...i$, animal
 311 species and $P = 1...j$ plant species, we’ll have a vector of “new” species to sample:
 312 $A_1P_1, A_1P_2, ...A_2P_1, A_2P_2, ...A_iP_j$. We can represent the successive samples where
 313 we can potentially get records of these interactions in a matrix with the vectorized
 314 interaction matrix and columns representing the successive samples we take (Table
 315 3). This is simply a vectorized version of the interaction matrix.

316 **0.0.5 Table 3 approx. here**

317 For example, mixture models incorporating detectabilities have been proposed to
 318 effectively account for rare species (?). In an analogous line, mixture models could
 319 be extended to samples of pairwise interactions, also with specific detectability
 320 values. These detection rate/odds could be variable among groups of interactions,
 321 depending on their specific detectability. For example, detectability of flower-
 322 pollinator interactions involving bumblebees could have a higher detectability than
 323 flower-pollinator pairwise interactions involving, say, nitidulid beetles. These more
 324 homogeneous groupings of pairwise interactions within a network define modules
 325 (Bascompte & Jordano, 2014), so we might expect that interactions of a given mod-
 326 ule (e.g., plants and their hummingbird pollinators; Fig. 1a) may share similar
 327 detectability values, in an analogous way to species groups receiving homogeneous
 328 detectability values in mixture models (?). Such sampling, in its simplest form,

329 would result in a sample with multiple pairwise interactions detected, in which
 330 the number of interaction events recorded for each distinct interaction found in
 331 the sample is recorded (i.e., a column vector in Table 3, corresponding to, say, a
 332 sampling day). The number of interactions recorded for the i_{th} pairwise interaction
 333 (i.e., $A_i P_j$ in Table 3), Y_i could be treated as a Poisson random variable with a
 334 mean parameter λ_i , its detection rate. Mixture models (?) include estimates for
 335 abundance-based data (their analogous in interaction sampling would be weighted
 336 data), where Y_i is a Poisson random variable with detection rate λ_i .
 337 This is combined with the incidence-based model, where Y_i is a binomial ran-
 338 dom variable (their analogous in interaction sampling would be presence/absence
 339 records of interactions) with detection odds λ_i . Let T be the number of samples
 340 in an incidence-based data set. A Poisson/binomial density can be written as (?):

$$g(y; \lambda) = \begin{cases} \frac{\lambda^y}{y! e^\lambda} & [1] \\ \binom{T}{y} \frac{\lambda^y}{(1+\lambda)^T} & [2] \end{cases}$$

341 where [1] corresponds to a weighted network, and [2] to a qualitative network.
 342 The detection rates λ_i depend on the relative abundances ϕ_i of the interactions,
 343 the probability of a pairwise interaction being detected when it is present, and the
 344 sample size (the number of interactions recorded), which, in turn, is a function
 345 of the sampling effort. Unfortunately, no specific sampling model has been devel-
 346 oped along these lines for species interactions and their characteristic features. For
 347 example, a complication factor might be that interaction abundances, ϕ_i , in real
 348 assemblages are a function of the abundances of interacting species, that determine
 349 interspecific encounter rates; yet they also depend on biological factors that ulti-

350 mately determine if the interaction occurs when the partner species are present. It
351 its simplest form, ϕ_i could be estimated from just the product of partner species
352 abundances, an approach recently used as a null model to assess the role of biolog-
353 ical constraints in generating forbidden links and explaining interaction patterns
354 (?). Yet more complex models should incorporate not only interspecific encounter
355 probabilities, but also phenotypic matching and incidence of forbidden links.

356 Rarefaction analysis and diversity-accumulation analysis (??) come up imme-
357 diately with this type of dataset. This procedure plots the accumulation curve
358 for the expected number of distinct pairwise interactions recorded with increasing
359 sampling effort (??). Asymptotic estimates of interaction richness and its associ-
360 ated standard errors and confidence intervals can thus be obtained (?). It should
361 be noted that the asymptotic estimate of interaction richness implicitly ignores
362 the fact that, due to forbidden links, a number of pairwise interactions among the
363 I_{max} number specified in the adjacency matrix Δ cannot be recorded, irrespective
364 of sampling effort. Therefore, the asymptotic value most likely is an overestimate
365 of the actual maximum number of links that can be present in an assemblage. If
366 forbidden links are taken into account, the asymptotic estimate should be lower.
367 Yet, to the best of my knowledge, there is no theory developed to estimate this
368 “biologically real” asymptotic value. Not unexpectedly, most recent analyses of
369 sampling effort in ecological network studies found evidences of undersampling
370 (Chacoff *et al.*, 2012). This needs not to be true, especially when interaction sub-
371 webs are studied (??), and once the issue of structural zeroes in the interaction
372 matrices is effectively incorporated in the estimates.

373 The *real* missing links

374 Given that a fraction of unobserved interactions can be accounted for by for-
375 bidden links, what about the remaining missing interactions? We have already
376 discussed that some of these could still be related to unaccounted constraints, and
377 still others would be certainly attributable to insufficient sampling. Would this
378 always be the case? Multispecific assemblages of distinct taxonomic relatedness,
379 whose interactions can be represented as bipartite networks (e.g., host-parasite,
380 plant-animal mutualisms, plant-herbivore interactions- with two distinct sets of
381 unrelated higher taxa), are shaped by interspecific encounters among individuals
382 of the partners (Fig. 2). A crucial ecological aspect limiting these interactions is
383 the probability of interspecific encounter, i.e., the probability that two individuals
384 of the partner species actually encounter each other in nature.

385 Given log-normally distributed abundances of the two species groups, the ex-
386 pected “neutral” probabilities of interspecific encounter (PIE) would be simply the
387 product of the two lognormal distributions. Thus, we might expect that for low
388 PIE values, pairwise interactions would be either extremely difficult to sample, or
389 just simply non-occurring in nature. Consider the Nava de las Correhuelas inter-
390 action web (NCH, Table 2), with $A = nnn$, $P = nnn$, $I = nnn$, and almost half
391 of the unobserved interactions not accounted for by forbidden links missing links,
392 $M = 53.1\%$. Given the robust sampling of this network (?), a sizable fraction of
393 these possible but missing links would be simply not occurring in nature, most
394 likely by extremely low PIE , in fact asymptotically zero. Given the vectorized
395 list of pairwise interactions for NCH, I computed the PIE values for each one by
396 multiplying element wise the two species abundance distributions. The $PIE_{max} =$

0.0597, being a neutral estimate, based on the assumption that interactions occur in proportion to the species-specific local abundances. With $PIE_{median} 1.410^{-4}$ we may safely expect (note the quantile estimate $Q_{75\%} = 3.2710^{-4}$) that a sizable fraction of these missing interactions may simply not occur according to this neutral expectation (?) (?) (neutral forbidden links, *sensu* (Canard *et al.*, 2012)). Which is the expected frequency for pairwise interactions? and, which is the expected probability for unobserved interactions? More specifically, which is the probability of missing interactions, M (i.e., the unobserved ones that cannot be accounted for as forbidden links)?

When we consider the vectorized interaction matrix, enumerating all pairwise interactions for the AP combinations, the expected probabilities of finding a given interaction can be estimated with a Good-Turing approximation (?). The technique, developed by Alan Turing and I.J. Good with applications to linguistics and word analysis (?) has been recently applied in ecology (Chao *et al.*, 2015), estimates the probability of recording an interaction of a hitherto unseen pair of partners, given a set of past records of interactions between other species pairs. Let a sample of N interactions so that n_r distinct pairwise interactions have exactly r records. All Good-Turing estimators obtain the underlying frequencies of events as:

$$P(X) = \frac{(N_X + 1)}{T} \left(1 - \frac{E(1)}{T}\right) \quad (1)$$

where X is the pairwise interaction, N_X is the number of times interaction X is recorded, T is the sample size (number of distinct interactions recorded) and $E(1)$ is an estimate of how many different interactions were recorded exactly once.

419 Strictly speaking Equation (1) gives the probability that the next interaction type
420 recorded will be X , after sampling a given assemblage of interacting species. In
421 other words, we scale down the maximum-likelihood estimator $\frac{n}{T}$ by a factor of
422 $\frac{1-E(1)}{T}$. This reduces all the probabilities for interactions we have recorded, and
423 makes room for interactions we haven't seen. If we sum over the interactions we
424 have seen, then the sum of $P(X)$ is $1 - \frac{1-E(1)}{T}$. Because probabilities sum to one,
425 we have the left-over probability of $P_{new} = \frac{E(1)}{T}$ of seeing something new, where
426 new means that we sample a new pairwise interaction.

427 Note, however, that Good-Turing estimators, as the traditional asymptotic
428 estimators, do not account in our case for the forbidden interactions. To account for
429 these FL I re-scaled the asymptotic estimates, so that a more meaningful estimate
430 could be obtained (Table 4). The scaling was calculated as $Chao1 * (I + ML) / AP$,
431 just correcting for the FL frequency, given that $I + ML$ represent the total *feasible*
432 interactions when discounting the forbidden links (Table 1). After scaling, observed
433 I values (Table 2) are within the $Chao1$ and ACE asymptotic estimates but below
434 the ACE estimates for Hato Ratón and Zackenberg (Table 4). Thus, even after
435 re-scaling for FL , it is likely that adequate characterization of most interaction
436 networks will require intensive sampling effort.

437 Discussion

438 Recent work has inferred that most data available for interaction networks are
439 incomplete due to undersampling, resulting in a variety of biased parameters and
440 network patterns (Chacoff *et al.*, 2012). It is important to note, however, that
441 in practice, many surveyed networks to date have been subnets of much larger

442 networks. This is true for protein interaction, gene regulation, and metabolic net-
443 works, where only a subset of the molecular entities in a cell have been sampled
444 (?). Despite recent attempts to document whole ecosystem meta-networks (?), it
445 is likely that most ecological interaction networks will illustrate just major ecosys-
446 tem compartments. Due to their high generalization, high temporal and spatial
447 turnover, and high complexity of association patterns, adequate sampling of ecolog-
448 ical interaction networks requires extremely large sampling effort. Undersampling
449 of ecological networks may originate from the analysis of assemblage subsets (e.g.,
450 taxonomically or functionally defined), and/or from logistically-limited sampling
451 effort. It is extremely hard to robustly sample the set of biotic interactions even
452 for relatively simple, species-poor assemblages; yet, concluding that all ecological
453 network datasets are undersampled would be unrealistic. The reason stems from
454 a biological fact: a sizeable fraction of the maximum, potential links that can be
455 recorded among two distinct sets of species is simply unobservable, irrespective of
456 sampling effort (?).

457 Missing links are a characteristic feature of all plant-animal interaction net-
458 works, and likely pervade other ecological interactions. Important natural history
459 details explain a fraction of them, resulting in unobservable interactions (i.e., for-
460 bidden interactions) that define structural zeroes in the interaction matrices and
461 contribute to their extreme sparseness. Sampling interactions is a way to monitor
462 biodiversity beyond the simple enumeration of component species and to develop
463 efficient and robust inventories of functional interactions. Yet no sampling theory
464 for interactions is available. Some key components of this sampling are analo-
465 gous to species sampling and traditional biodiversity inventories; however, there
466 are important differences. Focusing just on the realized interactions or treating

467 missing interactions as the expected unique result of sampling bias would miss
468 important components to understand how mutualisms coevolve within complex
469 webs of interdependence among species.

470 Contrary to species inventories, a sizable fraction of non-observed pairwise
471 interactions cannot be sampled, due to biological constraints that forbid their oc-
472 currence. A re-scaling of traditional asymptotic estimates for interaction richness
473 can be applied whenever the knowledge of natural history details about the study
474 system is sufficient to estimate at least the main causes of forbidden links. More-
475 over, recent implementations of inference methods for unobserved species (Chao
476 *et al.*, 2015) can be combined with the forbidden link approach, yet they do not
477 account either for the existence of these ecological constraints.

478 Ecological interactions provide the wireframe supporting the lives of species,
479 and they also embed crucial ecosystem functions which are fundamental for sup-
480 porting the Earth system. Yet we still have a limited knowledge of the biodiversity
481 of ecological interactions, but they are being lost (extinct) at a very fast pace, fre-
482 quently preceding species extinctions (?). We urgently need robust techniques to
483 assess the completeness of ecological interactions networks because this knowledge
484 will allow the identification of the minimal components of ecological complexity
485 that need to be restored after perturbations to rebuild functional ecosystems.

486 Acknowledgements

487 I am indebted to Jens M. Olesen, Alfredo Valido, Jordi Bascompte, Thomas
488 Lewinshon, John N. Thompson, Nick Gotelli, Carsten Dormann, and Paulo R.
489 Guimaraes Jr. for useful and thoughtful comments and discussion at different

490 stages of this manuscript. The study was supported by a Junta de Andalucía
491 Excellence Grant (RNM-5731), as well as a Severo Ochoa Excellence Award from
492 the Ministerio de Economía y Competitividad (SEV-2012-0262). The Agencia
493 de Medio Ambiente, Junta de Andalucía, provided generous facilities that made
494 possible my long-term field work in different natural parks.

495 Data accessibility

496 Please state where you have deposited the raw data underlying your analyses. It
497 will need to include the name of the repository (e.g. Dryad, figshare, GenBank
498 etc.) and location of the data (i.e DOI). For authors archiving at Dryad, we can
499 facilitate the process when your paper is accepted.

500

501 References

- 502 Banasek-Richter, C., Cattin, M. & Bersier, L. (2004) Sampling effects and the ro-
503 bustness of quantitative and qualitative food-web descriptors. *Journal of Theo-*
504 *retical Biology*, **226**, 23–32.
- 505 Bascompte, J. & Jordano, P. (2014) *Mutualistic networks*. Monographs in Popu-
506 lation Biology, No. 53. Princeton University Press, Princeton, NJ.
- 507 Bersier, L., Banasek-Richter, C. & Cattin, M. (2002) Quantitative descriptors of
508 food-web matrices. *Ecology*, **83**, 2394–2407.
- 509 Bluthgen, N. (2010) Why network analysis is often disconnected from community

- ecology: A critique and an ecologist's guide. *Basic And Applied Ecology*, **11**,
185–195.
- Bosch, J., Martín González, A.M., Rodrigo, A. & Navarro, D. (2009) Plant-pollinator networks: adding the pollinator's perspective. *Ecology Letters*, **12**,
409–419.
- Brose, U., Martinez, N. & Williams, R. (2003) Estimating species richness: Sensitivity to sample coverage and insensitivity to spatial patterns. *Ecology*, **84**,
2364–2377.
- Bunge, J. & Fitzpatrick, M. (1993) Estimating the number of species: a review. *Journal of the American Statistical Association*, **88**, 364–373.
- Canard, E., Mouquet, N., Marescot, L., Gaston, K.J., Gravel, D. & Mouillot, D. (2012) Emergence of structural patterns in neutral trophic networks. *PLoS ONE*, **7**, e38295.
- Chacoff, N.P., Vazquez, D.P., Lomascolo, S.B., Stevani, E.L., Dorado, J. & Padrón, B. (2012) Evaluating sampling completeness in a desert plant-pollinator network. *Journal of Animal Ecology*, **81**, 190–200.
- Chao, A. (2005) Species richness estimation. *Encyclopedia of Statistical Sciences*, pp. 7909–7916. Oxford University Press, New York, USA.
- Chao, A., Hsieh, T.C., Chazdon, R.L., Colwell, R.K. & Gotelli, N.J. (2015) Unveiling the species-rank abundance distribution by generalizing the Good-Turing sample coverage theory. *Ecology*, **96**, 1189–1201.

- 531 Coddington, J.A., Agnarsson, I., Miller, J.A., Kuntner, M. & Hormiga, G. (2009)
532 Undersampling bias: the null hypothesis for singleton species in tropical arthro-
533 pod surveys. *Journal of Animal Ecology*, **78**, 573–584.
- 534 Cohen, J.E. (1978) *Food webs and niche space*. Princeton University Press, Prince-
535 ton, New Jersey, US.
- 536 Cohen, J.E., Beaver, R.A., Cousins, S.H., DeAngelis, D.L., Goldwasser, L., Heong,
537 K.L., Holt, R.D., Kohn, A.J., Lawton, J.H., Martinez, N., O'Malley, R., Page,
538 L.M., Patten, B.C., Pimm, S.L., Polis, G., Rejmanek, M., Schoener, T.W.,
539 Schenly, K., Sprules, W.G., Teal, J.M., Ulanowicz, R., Warren, P.H., Wilbur,
540 H.M. & Yodis, P. (1993) Improving food webs. *Ecology*, **74**, 252–258.
- 541 Colwell, R. & Coddington, J. (1994) Estimating terrestrial biodiversity through ex-
542 trapolation. *Philosophical Transactions Of The Royal Society Of London Series*
543 *B-Biological Sciences*, **345**, 101–118.
- 544 Colwell, R.K. (2009) Biodiversity: concepts, patterns, and measurement. *The*
545 *Princeton Guide to Ecology* (ed. S.A. Levin), pp. 257–263. Princeton University
546 Press, Princeton.
- 547 Colwell, R.K. (2013) EstimateS: Biodiversity Estimation.

Figure captions

Figure 1. Sampling ecological interaction networks (e.g., plant-animal interactions) usually focus on different types of subsampling the full network, yielding submatrices $\Delta[m, n]$ of the full interaction matrix Δ with A and P animal and plant species. a) all the potential plants interacting with a subset of the animals (e.g., studying just the hummingbird-pollinated flower species in a community); b) all the potential animal species interacting with a subset of the plant species (e.g., studying the frugivore species feeding on figs *Ficus* in a community); and c) sampling a subset of all the potential animal species interacting with a subset of all the plant species (e.g., studying the plant-frugivore interactions of the rainforest understory).

Figure 2. Sampling species interactions in natural communities. Suppose an assemblage with $A = 3$ animal species (red, species 1–3 with three, two, and 1 individuals, respectively) and $P = 3$ plant species (green, species a-c with three individuals each) (colored balls), sampled with increasing effort in steps 1 to 6 (panels). In Step 1 we record animal species 1 and plant species 1 and 2 with a total of three interactions (black lines) represented as two distinct interactions: $1 - a$ and $1 - b$. As we advance our sampling (panels 1 to 6, illustrating e.g., additional sampling days) we record new distinct interactions. Note that we actually sample and record interactions among individuals, yet we pool the data across species to get a species by species interaction matrix. Few network analyses have been carried out on individual data(?).

Figure 1:

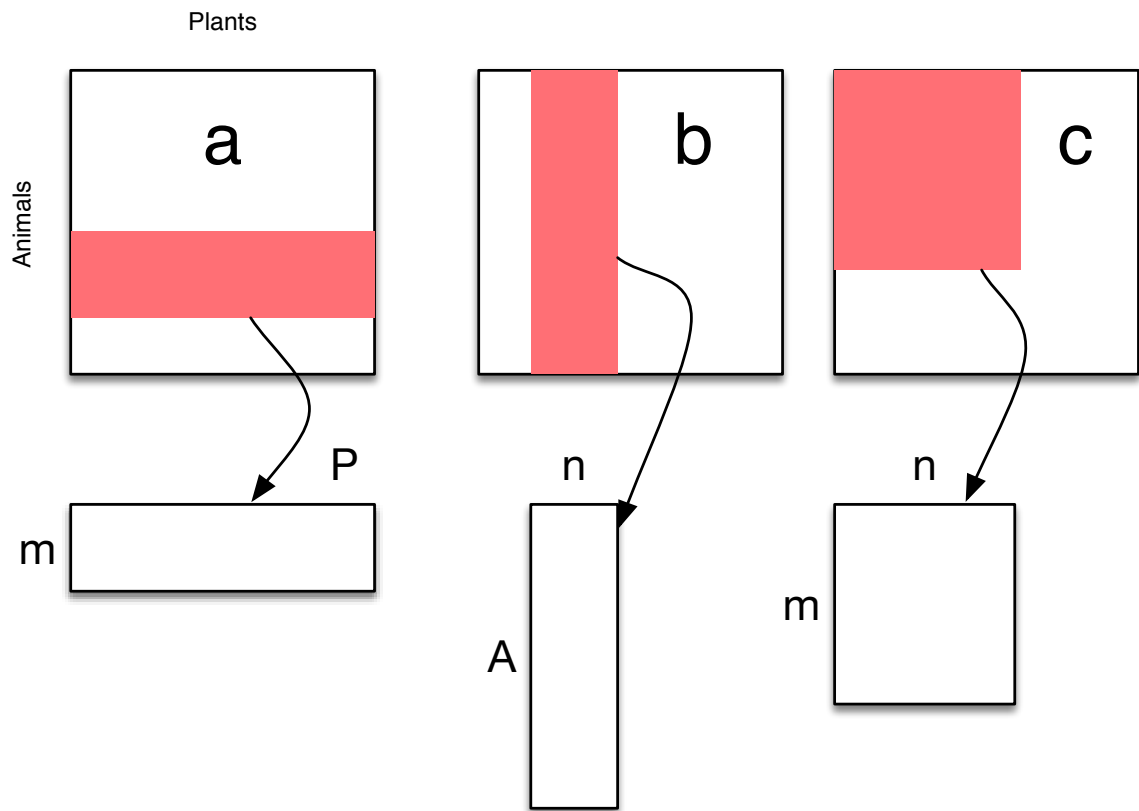
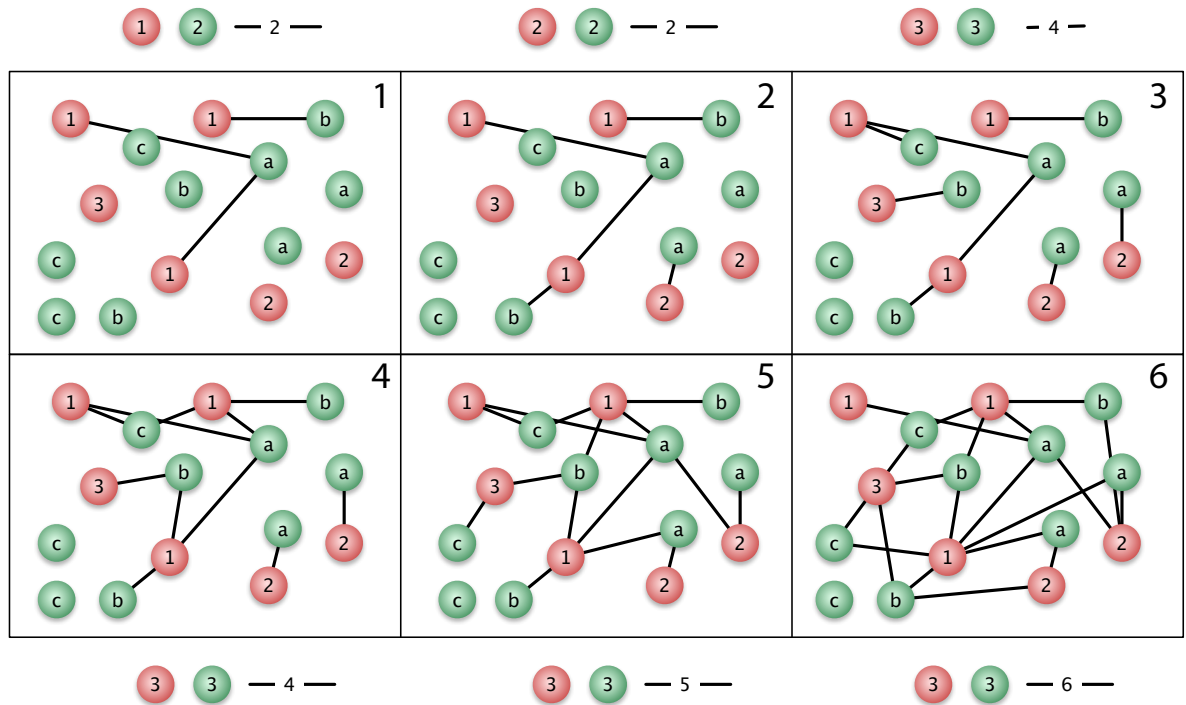


Figure 2:



Jordano – Figure 1