

MIT Open Access Articles

Reconstructing ecological networks with noisy dynamics

The MIT Faculty has made this article openly available. **Please share** how this access benefits you. Your story matters.

Citation: Freilich, Mara A. et al. "Reconstructing ecological networks with noisy dynamics." Proceedings of the Royal Society A 476, 2237 (May 2020): dx.doi.org/10.1098/rspa.2019.0739. © 2020 The Author(s)

Published Version: <http://dx.doi.org/10.1098/rspa.2019.0739>

Publisher: The Royal Society

Permanent Link: <https://hdl.handle.net/1721.1/131016>

Version: Author's final manuscript: final author's manuscript post peer review, without publisher's formatting or copy editing

Terms of use: <http://creativecommons.org/licenses/by-nc-sa/4.0/>



PROCEEDINGS A

rspa.royalsocietypublishing.org



Article submitted to journal

Subject Areas:

complexity, statistics, environmental engineering

Keywords:

food webs, ecological networks, network inference, stochastic model

Author for correspondence:

Mara Freilich

e-mail: maraf@mit.edu

Reconstructing ecological networks with noisy dynamics

Mara A. Freilich¹, Rolando Rebolledo²

Derek Corcoran³ and Pablo A.

Marquet^{3,4,5,6}

¹Massachusetts Institute of Technology-Woods Hole Oceanographic Institution Joint Program, Cambridge, MA, USA

² Instituto de Ingeniería Matemática, Facultad de Ingeniería, Universidad de Valparaíso, Valparaíso, Chile

³ Departamento de Ecología, Facultad de Ciencias Biológicas, Pontificia Universidad Católica de Chile, Santiago, Chile

⁴ Instituto de Ecología y Biodiversidad (IEB), Santiago, Chile

⁵ The Santa Fe Institute, Santa Fe, NM, USA

⁶ Instituto de Sistemas Complejos de Valparaíso (ISCV), Valparaíso, Chile

2 Abstract

3 Ecosystems functioning is based on an intricate web of interactions among living entities. Most
4 of these interactions are difficult to observe, especially when the diversity of interacting entities is
5 large and they are of small size and abundance. To sidestep this limitation, it has become common
6 to infer the network structure of ecosystems from time series of species abundance but it is not
7 clear how well can networks be reconstructed, especially in the presence of stochasticity that
8 propagates through ecological networks. We evaluate the effects of intrinsic noise and network
9 topology on the performance of different methods of inferring network structure from time series
10 data. Analysis of seven different four-species motifs using a stochastic model demonstrates that
11 star shaped motifs are differentially detected by these methods while rings are differentially
12 constructed. The ability to reconstruct the network is unaffected by the magnitude of stochasticity
13 in the population dynamics. Instead, interaction between the stochastic and deterministic parts of
14 the system determines the path that the whole system takes to equilibrium and shapes the species
15 covariance. We highlight the effects of long transients on the path to equilibrium and suggest a
16 path forward for developing more ecologically sound statistical techniques.

17 1. Introduction

18 Species interaction networks offer a quantitative method for understanding the structure and
19 dynamics of complex ecological systems, e.g. [1, 2]. However, in high biodiversity environments,
20 such as microbial communities, it may be difficult to directly observe species interactions. Instead,
21 it may be easier to observe temporal and spatial patterns that result from species interactions.
22 Based on this premise, co-occurrence of species in space and time is increasingly used to infer
23 networks of species interactions [3, 4, 5].

24 While there is a plethora of metrics to construct networks based on observations of species
25 distributions in time and space (see [6] for a review of metrics for time series), very few attempts
26 have been made to verify these metrics e.g. [7, 8] and even fewer to place them on sound ecological
27 *and* mathematical footing. Some of the theoretical issues that are largely unaddressed are how
28 the network architecture affects the ability to detect interactions, what types of interactions and
29 network motifs are best detected, and whether biological systems fit the statistical assumptions
30 of the metrics used [9]. A recent study of a number of different tools for detecting interaction
31 networks found that the tools produced different edges for the same real and simulated data and
32 that the power to detect interactions depends on the distributions of species abundance [9, 10].

33 Covariance and correlation are basic methods to quantify pairwise associations. Assuming
34 that all species have Gaussian distributions, the covariance and the mean value is a complete
35 descriptor of the species distributions. However, the Gaussian assumption is often not satisfied
36 and there is not a direct link between correlation or covariance and interactions [7, 8]. One way
37 to relax the assumption of a Gaussian distribution while ensuring that the statistical technique
38 is not making assumptions about missing data is to use an entropy maximizing method [11, 12].
39 These methods also reduce the effects of indirect interactions or large environmental trends on the
40 inferred associations. There is a similarity between the maximum entropy metric proposed by [11]
41 and the statistical techniques used in press experiments in ecology [13]. An entropy maximizing
42 method does not measure pairwise associations, but instead infers the most probable whole
43 network structure; the ability to detect relationships involving many interconnected members
44 is pertinent to the use of this metric.

45 More than 30 years ago, Peter Yodzis [13] proposed an indeterminacy principle wherein
46 available statistical techniques were insufficient to experimentally determine species interactions
47 in natural communities for a large suite of ecological reasons, including indirect interactions and
48 weak links in species networks. The uncertainty about the ability to detect interactions extends to
49 null models, another commonly used ecological technique [14, 15]. Even with modern statistical
50 techniques, relationships involving more than three members of a community seem to be nearly

impossible to detect reliably [9, 16, 17]. We present an updated indeterminacy principle based on both mathematical and ecological insights. In this study, we use a first principle stochastic species interaction model to investigate the impact of food web structure on both the dynamics of communities and the propagation of stochasticity through them and hence on the ability to detect the network structure of species interactions. We examine the power of three metrics; covariance, Pearson's correlation, and inverse correlation to detect species interactions across a wide parameter space. These metrics are at the heart of most of the network reconstruction techniques. Unlike previous work [10], we focus on small community modules to gain an understanding of both the ecological and statistical reasons for trends in performance.

2. Methods

(a) Stochastic species interaction model

We model food webs (interaction networks) using a stochastic differential equation, which is a generalized Lotka-Volterra model with stochasticity (Eq. 2.1). The generalized Lotka-Volterra model represents the dynamics of biomass in the main system, while the diffusion coefficient represents biomass fluctuations associated with the interactions between the main system and the environment where it is embedded. This includes not only environmental variability, but everything that we do not model explicitly in the main system, such as interaction with species outside the main system and non-trophic interactions with species in the main system [18]. We modeled the biomass dynamics in the main system assuming that each species has resources outside the food web that allow it to grow, which prevents extinction at long times in the stochastic simulation:

$$\left\{ \begin{array}{l} dX_i(t) = X_i(t) \left(r_i + \underbrace{\sum_{j=1}^N D_{ij} X_j}_{\text{generalized Lotka Volterra}} \right) dt + \underbrace{\sqrt{\gamma_i X_i(t)(1 - X_i(t))} dW(t)}_{\text{diffusion}} \\ X_i(0), \text{ with a given probability distribution } \mu \end{array} \right. \quad (2.1)$$

where $X_i(t)$ represents the proportion of biomass in species i (or stochastic abundance of species i) and r_i its per capita growth rate in the absence of inter- and intra-specific interactions. In addition, let a_i denote the strength of interaction. Species interactions are represented by a symmetric N by N matrix $D = (D_{i,j})_{i,j=1}^N$, which, for all $i \neq j$ represents the impact of predation upon the focal species as a biomass loss rate for prey and gain rate for predator. We define the intraspecific interactions as $D_{i,i} = -a_i$, for all $i = 1, \dots, N$. The interspecific interactions can be represented as a network (as in Figure 1). The magnitude of the interspecific interaction rates is given by d . The second term in the right hand part of equation 2.1 represents interactions of the main system with the environment, where $W(t)$ is a Wiener or Brownian motion, γ_i is the intensity of the fluctuation for all $i = 1, \dots, N$. And one assumes that the initial distribution μ is known.

One may write the previous equation in vector and integral notation as follows. Let denote $X(t)$ the column vector with components $X_i(t)$, similarly, call r the vector with coordinates r_i , $\sigma(x)$ the diagonal matrix with components $\sqrt{\gamma_i x_i (1 - x_i)}$, where $x \in [0, 1]^N$, (i.e. each $x_i \in [0, 1]$). Moreover, let denote \bullet the Schur (or component-wise) product of vectors. The equation (2.1) becomes:

$$X(t) = X(0) + \int_0^t X(s) \bullet (r + DX(s)) ds + \int_0^t \sigma(X(s)) dW(s). \quad (2.2)$$

This model provides considerable insight into the way in which species covary and correlate with each other. Past work has used similar generalized Lotka-Volterra models, but with competition rather than predation, to draw conclusions about the influence of niche and neutral processes in ecological communities [19, 20].

92 A proof of the existence and uniqueness of the solution in distribution P_θ to this general
93 equation is given in theorem 1 of [21]. P_θ is a probability defined on the set $C(\mathbb{R}_+, \mathbb{R}^N)$ of
94 continuous functions, depending on the parameter $\theta = (r, D, \gamma)$, where γ is the vector with
95 components γ_i . Notice that here we assume a form of the functional diffusion coefficient $\sigma(x)$
96 depending on γ only, which models the rescaled limit of birth and death processes that appear in
97 a number of ecological and genetics models [18]. This assumption introduces a simplification of
98 the statistical inference problem.

99 (b) Model parameters

100 Predator-prey interactions are modeled using a linear functional response in equation 2.1.
101 The numerical simulation is performed with forward Euler time stepping. Since the predation
102 coefficients are symmetric, the loss to the prey in biomass is equivalent to the gain to the predator
103 in biomass and there is no biomass loss from the system, except due to stochastic fluctuations
104 (an “open system”). We systematically explore the parameter space of the model by varying
105 the amount of stochasticity from the Brownian term and the ratio of interspecific to intraspecific
106 interactions.

107 (i) Interaction motifs

108 We simulate species interactions for all possible directed food webs with four species. These
109 motifs are the building blocks from which larger networks are built. Interaction topologies are
110 shown in figure 1. In the deterministic case, the species grow to an equilibrium in which all
111 species coexist. The parameters a_i and r_i are chosen such that all species grow to the same
112 biomass at equilibrium when there are no interactions. The value of the equilibrium depends on
113 both intraspecific and interspecific interactions associated to the strength of predation, in other
114 words, the amount of energy transferred between prey and predators. In order to systematize
115 the simulations across network geometries, we assign the coefficients $a_i = 0.05$ and $r_i = 0.5$ to the
116 top predator in the simulation (any species that has no predators). We assign the same coefficients
117 as above $a_i = 0.1$ and $r_i = 1$ to basal species and intermediate consumers. As a result, all species
118 have the same carrying capacity in the absence of interactions ($r_i/a_i = 0.1$), but the prey species
119 grow towards this carrying capacity more quickly. If there is more than one species with the same
120 interactions, we assign different interaction rates so that the species dynamics differ between these
121 species. In order to explore the effects of the parameter choices, we vary γ , the coefficient for the
122 stochasticity, between 0 and 0.05. All species have the same stochastic growth coefficient γ , but the
123 noise is independent. We vary the interspecific interaction coefficient between 0.001 and 0.25. We
124 simulated 10000 time series for each parameter combination. We present results from time step
125 150 to time step 350. There are 20 possible undirected networks for the model with four species
126 and three links between species. While the food web models are directed, the networks inferred by
127 correlation and covariance are necessarily undirected. The LIMITS algorithm can detect directed
128 networks.

129 (ii) Niche model trophic network

130 In order to evaluate the implications of the results obtained by studying the dynamics of motifs
131 in detail, we simulate larger networks using the niche model for trophic networks (food webs)
132 [22]. This model is a stochastic model that constructs trophic networks that share many of the
133 properties of observed trophic networks. We specify that the networks have 100 species and are
134 sparse, with a connectivity of 0.05. There is a continuous spectrum of trophic levels in the model so
135 we assign the same intraspecific interaction parameters to all species, $r_i = a_i = 0.09$. We vary the
136 interspecific interaction strength from 0.01 to 1. We assign the same interaction strength to every
137 link present. We also vary the stochasticity γ . We perform 1000 simulations for each parameter
138 combination, each with a different randomly generated trophic network.

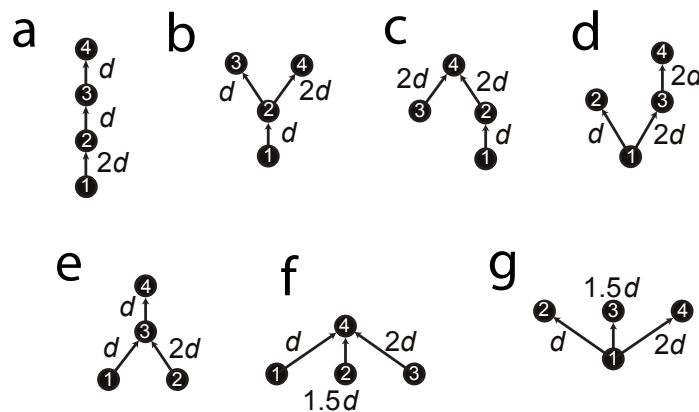


Figure 1. Possible predator-prey networks with four species, when distinguishing between trophic levels. The top predators are the uppermost species in the network. The arrows show movement of biomass from prey to predators. The values near each edge indicate the rate of this interspecific interaction. The parameter d is varied through the simulations.

139 (c) Network reconstruction

140 We calculated species covariance matrices from the time series generated by the stochastic
 141 simulations as $\text{cov} = \overline{(X_i - \bar{X}_i)(X_j - \bar{X}_j)}$ where $\bar{\cdot}$ represents temporal averaging. We also
 142 calculated species correlation matrices using the Pearson's correlation coefficient $\text{corr} = \frac{\text{cov}}{S_{X_i} S_{X_j}}$
 143 where S_{X_i} is the standard deviation in time of species i . The maximum entropy technique
 144 we used is computed as the inverse Pearson correlation coefficient matrix [11]. This is a more
 145 complete descriptor if the species distributions are Gaussian.

146 For the four species motifs, we list the top three inferred connections for each network
 147 construction technique (covariance, correlation, and inverse correlation [MaxEnt]). The links with
 148 magnitudes that are above a certain threshold are selected when using the covariance and inverse
 149 correlation techniques. The links with the lowest p-values are selected when using correlation.
 150 For the 100 species network we only present the results using correlation to determine the
 151 links. We selected all links with a p-value below 0.01 when using correlation. We then use
 152 a binomial approximation to assign confidence intervals to the probability of detecting each
 153 possible network.

154 LIMITS [23] is one of the most successful network reconstruction techniques for timeseries
 155 data [24]. LIMITS is designed to reconstruct Lotka-Volterra networks, which are the type of
 156 interactions used in this study. LIMITS generates directed networks while the other metrics
 157 generate undirected networks. We implement this metric using Mathematica code provided by
 158 the authors. We use a threshold for inferring an interaction of 0.01. The results from LIMITS
 159 provide a baseline for success of reconstruction using the other metrics.

160 (d) Statistical analysis

161 We used Generalized Linear Models (GLM) [25] to find the relationship between the probability of
 162 finding the true configuration of the network ($prob$) and its predictive variables the strength of the
 163 relationships (d), the level of noise ($gamma$), the network motif ($network$), and the method used to
 164 infer the network structure ($Method$). We performed logistic regression using each of the predictive
 165 variables in isolation and in combination to predict the probability of detecting the true network.
 166 This quantifies the importance of each of the above factors on the network reconstruction.

167 3. The statistical problem

168 Before presenting the results, we outline some expectations about the performance of the
 169 statistical tools from a theoretical perspective. The statistical inference problem consists of the
 170 identification of the probability P_θ that rules the dynamics of the open system described by (2.1).
 171 That probability provides the answer to the query on the network structure as well as a complete
 172 dynamical picture of biomass exchanges between species. This is a hard problem since P_θ is a
 173 probability on an infinite-dimensional space. P_θ represents the state of the whole open system.
 174 In our equation (2.1), the network architecture is carried by the matrix D . That is, one defines a
 175 graph $G = (V, E)$, where $V = \{1, \dots, N\}$ is the set of species and E the set of edges: $\{i, j\} \in E$ if
 176 and only if $D_{i,j} \neq 0$. What statistical techniques can reliably infer $D_{i,j}$?

177 Statistical data have the form of matrices: $(X(\omega_k, t_\ell))$: $k = 1, \dots, K$; $\ell = 1, \dots, L$. Notice that
 178 one needs to observe different trajectories (ω_k) at any arbitrary finite sequence of times (t_ℓ) . This
 179 is a first difficulty for the application of time series methods based on correlation and maximum
 180 entropy to estimate P_θ .

181 The covariance metric is limited as a measure of the network structure because of the influence
 182 of the noise terms. The expectation under P_θ of $X(t)$ from that equation is given by

$$E_\theta(X(t)) = E_\theta(X(0)) + \int_0^t E_\theta [X(s) \bullet (r + DX(s))] ds. \quad (3.1)$$

183 The covariances between components $X_i(t)$ and $X_j(t)$ at a given time t are then

$$\begin{aligned} C(X_i(t), X_j(t)) &= E_\theta [(X_i(t) - E_\theta(X_i(t)))(X_j(t) - E_\theta(X_j(t)))] \\ &= \int_0^t E_\theta(\sigma_i(X(s))\sigma_j(X(s)))ds. \end{aligned} \quad (3.2)$$

184 Therefore, if i and j are not connected nodes, that is $D_{i,j} = 0$ one may have $C(X_i(t), X_j(t)) \neq 0$
 185 if $\sigma_i(X(s))$ and $\sigma_j(X(s))$ are not orthogonal in the space $L^2(\Omega \times [0, t], dP_\theta \otimes ds)$.

186 Using a maximum entropy technique based on the correlations between $X_i(t)$ and $X_j(t)$ to
 187 discover ρ_t may not be possible in most cases. Under P_θ one can find a probability density $\rho_t(x)$
 188 for each fixed time t :

$$P_\theta(X(t) \in A) = \mu_t(A) = \int_A \rho_t(x) dx, \quad (3.3)$$

189 for all measurable subset in \mathbb{R}^N in the space A . Though, as it is well-known, the knowledge of
 190 ρ_t for all $t \geq 0$, does not suffice to identify P_θ . One needs to know an infinite family (μ_{t_1, \dots, t_n})
 191 of measures, where $t_1 < \dots < t_n$ run over all finite sequences of times, and each μ_{t_1, \dots, t_n} is a
 192 measure on the space $(\mathbb{R}^N)^n$, such that

$$\mu_{t_1, \dots, t_{n-1}, t_n}(A_1 \times \dots \times A_{n-1} \times \mathbb{R}^N) = \mu_{t_1, \dots, t_{n-1}}(A_1 \times \dots \times A_{n-1}),$$

193 the so-called Kolmogorov's consistency relation. So, under that hypothesis one could prove the
 194 existence of a probability measure P_θ on the set of trajectories, such that

$$P_\theta(X(t_1) \in A_1, \dots, X(t_n) \in A_n) = \mu_{t_1, \dots, t_n}(A_1 \times \dots \times A_n).$$

195 The maximum entropy at each step t_1, \dots, t_n does not preserve Kolmogorov's consistency
 196 relation and so cannot be applied to each measure μ_{t_1, \dots, t_n} . However, the network underlying a
 197 Gaussian process may be relatively detectable. The probability distribution of Gaussian processes
 198 is entirely determined by the covariance kernel and the mean and it is well known that Gaussian
 199 laws maximize the Shannon entropy among all distribution with second moments. Unfortunately,
 200 as we will demonstrate, our process $X(t)$ here is not a Gaussian one, nor are ecological species
 201 abundance distributions commonly Gaussian. Angulo and coauthors [26] have also shown that
 202 in order for a network to be reconstructed more information beyond a timeseries of abundances,
 203 such as information about the interaction functional forms, must be known.

204 We expect that the covariance, correlation, and inverse correlation (maximum entropy) metrics
 205 will not perform well in detecting the whole network structure. Simulations allow us to probe the

206 limits of the network inference and expose the diversity of ways in which the network inference
207 depends on the underlying dynamics.

208 4. Numerical results

209 (a) Motifs

210 In order to explore the consequences of these statistical limitations, we perform numerical
211 simulations of the model system (equation 2.1). There are 20 possible networks with four species
212 and three edges (the minimum spanning tree for a four node network), based on combinatorics.
213 Out of these 20 possible networks, there are seven distinct directed motifs in which all four species
214 are connected. The networks are referred to by their labels in figure 1. The power to detect the
215 species interactions is low even for networks that are detected more often than expected by
216 random chance. True species interactions are typically only detected around 6-10% of the time.
217 Since there are 20 unique networks, random detection is 5%. The overall highest probability of
218 detection is for star networks with more basal species when using covariance to detect interactions
219 and with a large species interaction coefficient d (see Table 1 and Figure 3).

220 We summarize the overall performance of the metrics by calculating the sensitivity (or true
221 positive rate), which is the probability of detecting a link when one exists, and specificity (or
222 true negative rate), which is the probability of detecting that there is no link when there is no link
223 (Figure 2). This is a common performance metric that can be compared to other studies of network
224 detection [8, 9, 10]. For these networks with 3 links and 4 species, if edges are chosen at random,
225 specificity and sensitivity will, on average, equal 0.5. We find that both sensitivity and specificity
226 are within one standard deviation of 0.5 for almost all networks and parameter combinations
227 when using covariance and correlation metrics. The LIMITS algorithm [23] can be used as a
228 benchmark for the performance of the other metrics. With low interaction strengths, LIMITS
229 detects few interactions but performs slightly better than the other metrics at high interaction
230 strengths. LIMITS has low specificity, or a high false positive rate for all parameters. The false
231 positive rate is affected by the detection threshold selected by the user.

232 The seven interaction networks with four species can be classified into three different types of
233 systems based on qualitatively similar detection patterns using correlation based metrics (Table
234 1). Logistic regression best predicts the probability of true network detection when using a model
235 that includes interaction strength d , the network type, and reconstruction method as parameters.
236 The model has a Nagelkerke Pseudo R-squared of 0.62. The magnitude of stochasticity γ is not
237 a significant predictor of the probability of detection of the true network. These classifications
238 by network type also align with the undirected network structure, which is either linear or star
239 shaped. The linear networks (networks a,c,d) are best detected by correlation, the star networks
240 with more top predators (networks b,g) are best detected by inverse correlation, and the star
241 networks with more basal species (networks e,f) are best detected by covariance. In all cases,
242 the networks are best detected when interspecific interactions are relatively strong ($d = 0.25$).
243 Detection power is very low for the linear networks using any association metric. The linear-type
244 network d is nearly undetectable with a highest predicted probability less than random.

245 There is a systematic relationship between the most likely network to be inferred and the true
246 network, however the most likely network to be inferred is rarely the true network. The results for
247 three case studies are synthesized in figure 3 the proportion of the 10000 simulations in which each
248 undirected network was detected with each reconstruction metric, with three different interaction
249 strengths (0.001, 0.01, and $0.25 \text{ biomass}^{-1} \text{ day}^{-1}$) and two different stochasticity magnitudes
250 (0.005 and 0.05 day^{-1}). For almost all experiments and parameter combinations, all 20 possible
251 undirected networks were inferred with some non-zero probability (figure 3).

252 When the true network is a star shaped network, the true network is detected across a wider
253 parameter range when there are more basal species than top predators while a ring with three
254 species connected and one species unconnected is inferred for motifs with more top predators
255 (figure 3 top row and electronic supplementary material). For the star shaped networks, the

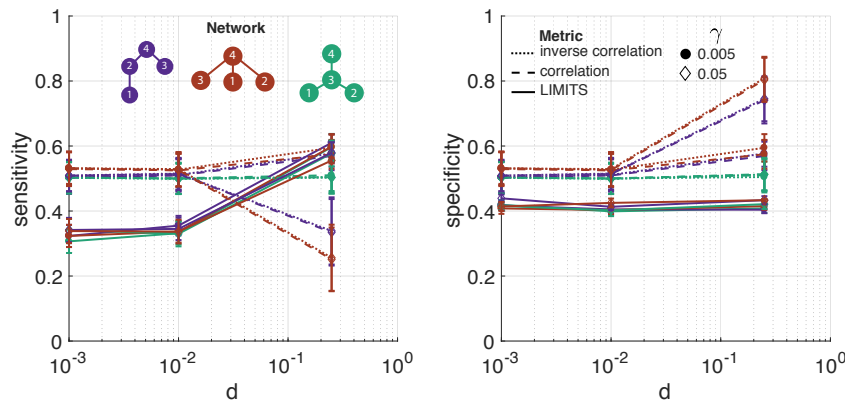


Figure 2. (left) Sensitivity (true positive rate) of reconstruction technique (right) Specificity (true negative rate) of reconstruction technique. The sensitivity and specificity are shown as a function of interaction strength d for a selection of network motifs (colors), reconstruction metrics (line styles), and stochasticity parameters (symbols). The uncertainty ranges are given by the standard deviation.

Table 1. Parameters leading to the highest probability of detecting the network structure arranged by descending probability. The F-value of logistic regression of each parameter as a predictor of detecting the true network is shown below. The stochasticity γ is not a significant predictor with an F-value = 3.2049 and $\Pr(>F) = 0.0740$

	d	network	Method	predicted prob
	0.25	e	Cov	0.0898
	0.25	f	Cov	0.0884
	0.25	c	Corr	0.0646
	0.25	a	Corr	0.0571
	0.25	b	Invcor	0.0568
	0.25	g	Invcor	0.0541
	0.25	d	Corr	0.0474
F-value	33.3005	90.0216	5.0344	
$\Pr(>F)$	0	0	0.0068	

256 complement to the true network is a ring with one species unconnected. The network complement
 257 is purely indirect interactions. Covariance at times has a high probability of detecting the
 258 network complement. For example, with high stochasticity and large interaction coefficients,
 259 the complement of the network with two top predators (network b) is detected in 63% of the
 260 simulations and the complement of the network with three top predators (network g) is detected
 261 in over 99% of the simulations.

262 Although stochasticity does not significantly influence the probability of detecting the true
 263 network, stochasticity may affect the detection of the true networks through interaction with
 264 the drift term. One way this may happen is through noise-induced large transients away from
 265 equilibrium. Networks with the potential for large transients away from equilibrium can be
 266 identified by calculating the eigenvalues of the symmetric part of the community matrix D . If the
 267 largest eigenvalue of the symmetric part of the community matrix is positive (figure 4), there is the
 268 possibility for large transient growth [27, 28]. The direction of this transient growth is given by the
 269 eigenvector that corresponds to the largest eigenvalues of the symmetric part of the community
 270 matrix. For network b, this vector points most in the direction of species 2, for networks c and f,

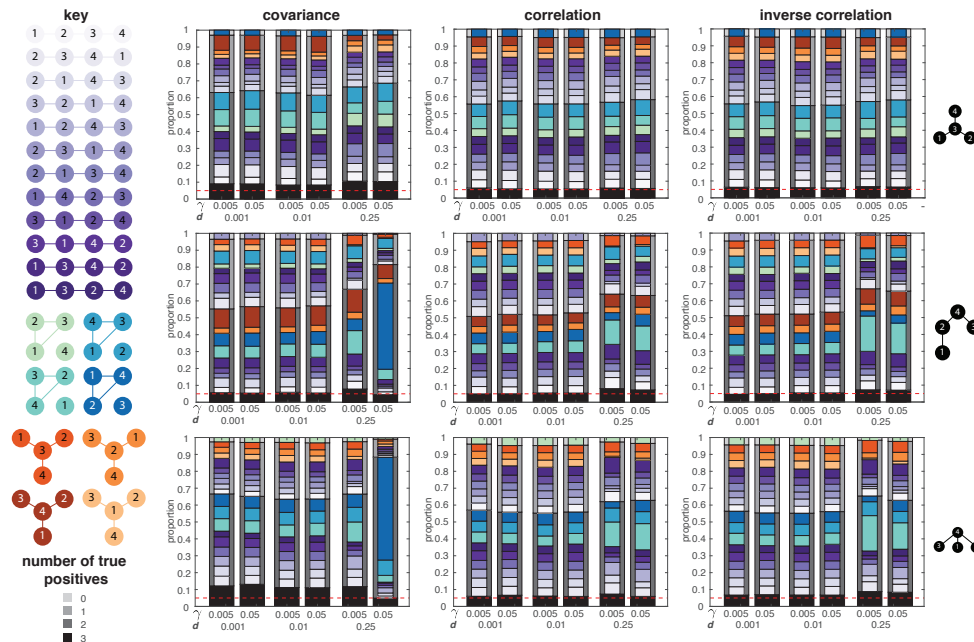


Figure 3. Proportion of trials (out of 10000) for which each network is inferred for each parameter combination (stochasticity parameter γ and strength of predation d) for three different motifs (rows). The color of the bar indicates the network that is inferred by reference to the key at the left. The true network is to the right of each row. The black portion of the bar at the bottom is the true network, consequently the true network is not represented as a colored bar. The gray shaded background bars in the background show the number of links in the true network that are correctly identified as links in the inferred network (true positives). The red line indicates random detection (0.05).

271 this eigenvalue points equally in the direction of species 3 and 4. In all cases with long transients,
 272 one species becomes disconnected and the ring geometry is preferentially constructed. This is
 273 one example of a quantitative evaluation of the way that network structure interactions with the
 274 stochastic drift term may lead to an inferred network that is statistically significant but distinct
 275 from the true network.

276 The abundance distribution for each species in time is non-normal, with most distributions left
 277 skewed, especially for the top consumer (species 4), and increasingly so as the rate of interspecific
 278 interactions and stochasticity increases. An example abundance distribution for two interaction
 279 strengths is given for three of the four-species interaction networks geometries (figure 5). As this
 280 figure shows, and all else being equal, network geometry affect the abundance distribution of
 281 species.

282 (b) Niche model trophic network

283 With the larger, 100 species networks, we find that the relationship between the correlation
 284 network and the interaction networks depends on the interaction strength. With weak
 285 interactions, there is high sensitivity but low specificity due to overly dense correlation networks.
 286 With stronger interactions, by contrast, both the sensitivity and specificity are low, near 0.5 (Figure
 287 6c). In addition, we find that there is a non-random structure to the correlation networks. While
 288 the interaction networks have the highest interaction density away from the diagonals, we find
 289 that the correlation density is highest on the diagonal of the correlation matrix and near the edges
 290 of the matrix (Figure 6d-f). The networks with large interaction strengths are prone to transients
 291 with positive maximum eigenvalues of the symmetric part of the interaction matrix.

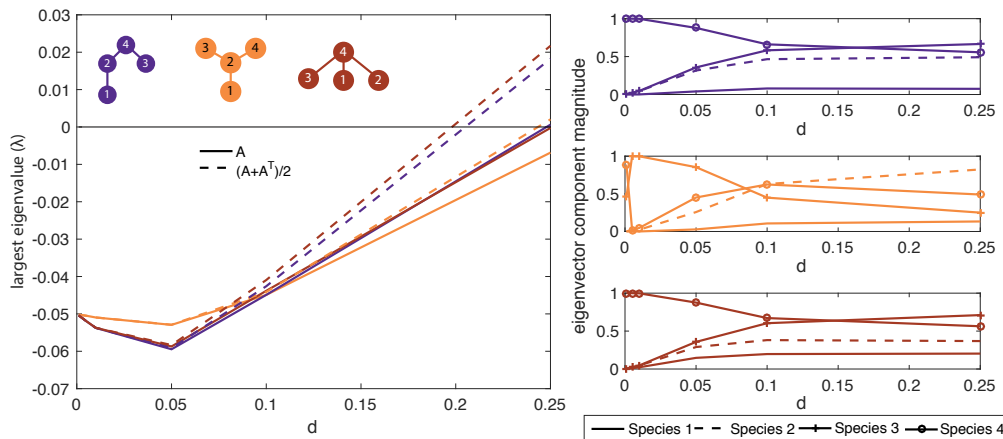


Figure 4. (left) Largest eigenvalues of the species interaction matrix (solid) and the symmetric part of the species interaction matrix (dashed) as a function of species interaction strength d . A positive maximum eigenvalue of the symmetric part of the species interaction matrix indicates non-normal growth. Growth occurs in the direction of the eigenvector that corresponds to the largest eigenvalue. To the right, the direction of this eigenvector in species coordinates.

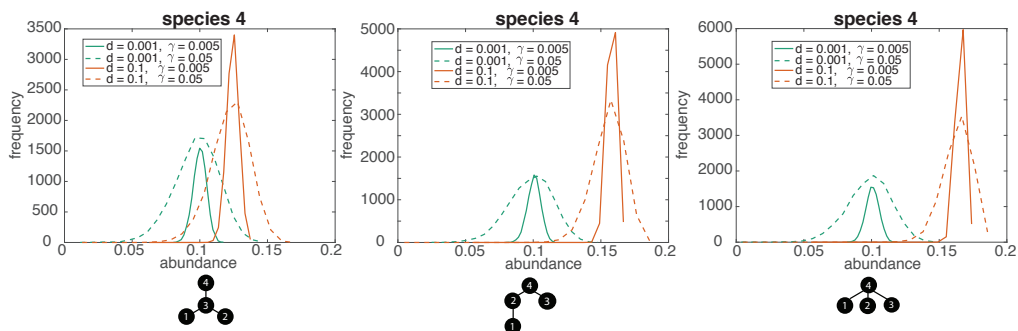


Figure 5. Distribution of each species abundances at the end of 10000 simulations for different parameter combinations.

292 5. Discussion

293 There is a low probability of detecting the true interaction network, but we find that it is possible
 294 to detect four-species interaction modules for a wide range of predation rates and especially for
 295 star shaped interaction geometries. Detection of the true network is robust to stochasticity for all
 296 metrics and, counter intuitively, aided by stochasticity in some cases. We can explain some of the
 297 covariance behavior through understanding the transient dynamics of the model system.

298 (a) Ability to detect interaction modules

299 For most parameter combinations where the true interaction matrix can be detected, it is detected
 300 in at most only 10-11% of the trials. This level of detection of whole modules yields specificity
 301 and sensitivity of about 0.5, which is consistent with the levels of sensitivity and specificity
 302 obtained by other simulation studies [8, 9, 10]. These values of sensitivity and specificity would
 303 seem to suggest that edges are selected randomly, however certain motifs are more likely to
 304 be detected than others, namely star shaped geometries, and certain motifs are consistently
 305 constructed, namely motifs with three species completely connected. For all metrics there is a

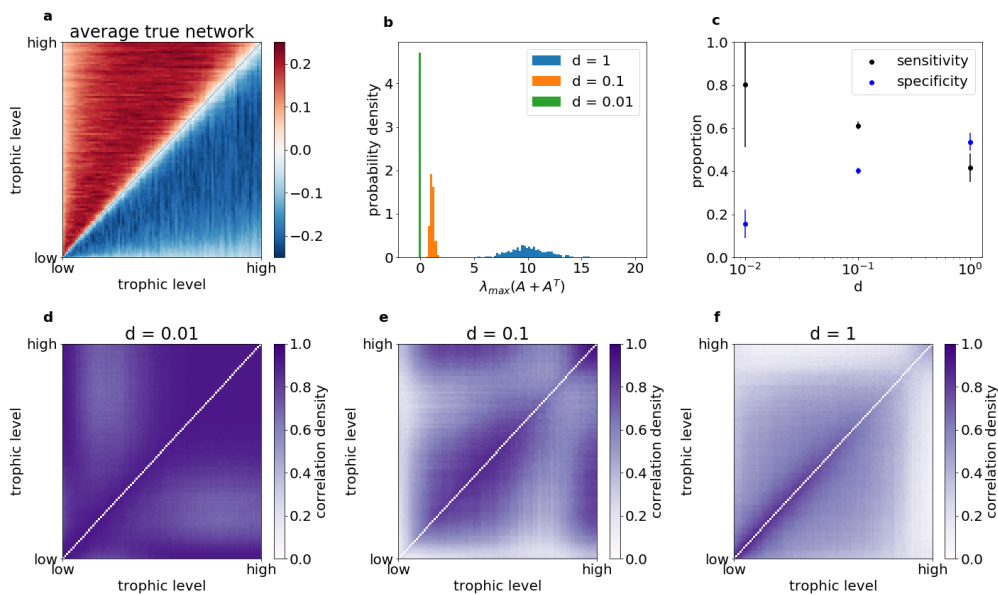


Figure 6. (a) Average interaction strength in the trophic networks used to simulate the many species communities. Darker shades indicate that the interaction is more frequent, lighter shades indicate less frequent interactions. All 100 species are shown on the axis arranged from low to high trophic levels. (b) Distribution of the maximum eigenvalue of the symmetric part of the species interaction matrix. Positive values indicate that the interaction matrix is susceptible to non-normal growth of perturbations. (c) Sensitivity and specificity for different values of interaction strength. The mean value for all 1000 simulations is shown with the confidence interval indicating the variance. (d-f). Constructed species interaction matrices using correlation. The darker shades indicate that the particular interaction is detected more frequently by correlation.

306 bias against detecting linear networks, which are almost never detected. It is important that we
 307 found that there is some ability to detect the whole network structure with the inverse correlation
 308 technique because the maximum entropy method relies on a reconstruction of the most probable
 309 configuration for the whole network, whereas the correlation and covariance metrics are pairwise
 310 metrics.

311 (b) Network geometry

312 The motifs used in this study are prevalent in true ecological and genetic networks. A chain of
 313 length three and a motif with four nodes all connected in a loop was found to occur more often
 314 than expected by chance in ecological networks while an undirected chain of length four, which
 315 is nearly undetectable in our analysis, occurs less often and the star shaped motif, which is more
 316 easily detectable, occurs more often than expected by chance in a protein interaction network
 317 [29]. The networks used in this study differ from true networks in that true networks are thought
 318 to be more sparse, meaning they have a much lower connectance [30]. This difference might be
 319 expected to affect the total number of edges and hence the specificity and sensitivity.

320 The network geometry affects which networks are inferred, even if the inferred networks
 321 are not the true network. Researchers may conclude that species networks inferred statistically
 322 (“association networks” [8]) are interesting and useful objects to study that have some
 323 relationship to the underlying interaction network structure even if there is not direct
 324 correspondence between the observed associations and true interactions. This result of non-
 325 random generation of a network holds for large interaction networks as well.

326 (c) Transience and other effects of stochasticity

327 There is an implicit assumption in the literature that increased noise decreases the ability to detect
328 interactions within networks [11, 31]. Noise or stochasticity in either abundance or occurrence,
329 however, is an essential assumption behind the use of correlation based analysis, since in
330 deterministic environments there is no correlation. Therefore, noise can be exploited to make
331 inferences about network structure [32]. We find that the magnitude of the stochasticity is not
332 a predictor of the success in inferring network structure.

333 Propagation of stochasticity and of biomass through the web affects the inferred structures.
334 A few of the networks used (networks b, c, and f) have the potential for non-normal growth
335 of perturbations (as indicted by positive eigenvalues of the symmetric part of the community
336 matrix), especially with strong interactions between species. Although all species have identical
337 and independent stochastic diffusion terms, the interaction network directs perturbations
338 through the web such that certain species may be perturbed away from equilibrium more than
339 others. This is particularly true of networks that have the potential for non-normal growth.
340 Stochasticity can have an apparent organizing effect because perturbations are amplified in a
341 specific direction and recovery to equilibrium is constrained to occur along certain pathways.
342 In biological terms, we find that the effects of top predators seem to propagate through food
343 webs, affecting all species present in the system. For example, the network with species 1, 2,
344 and 4 connected in a ring is commonly constructed for a motif with two intermediate consumers
345 (networks c and d). In these networks, species 1, 2, and 4 are connected in a chain in the true
346 network. Connection of species 1 and species 4 is meaningful in that it represents a trophic
347 cascade. Even more strikingly, the network in which all species are connected to species 4 (the top
348 predator) is constructed by covariance across a wide range of parameters for network d. Species
349 4 is only connected to species 2 in the true food web. This strong influence of top predators could
350 be one reason that linear chains are almost never detected - the top predator is likely to covary
351 with more species than just its immediate prey.

352 Stochasticity only affects detection if there is interaction between the stochasticity and the drift
353 terms. In the numerical experiments, we find that γ , the strength of the stochastic noise, is not a
354 significant predictor of the ability to detect the true network (Supplementary Table 1).

355 The Brownian motion used in this study is on a particular (fast) timescale. The functional
356 form used results in larger magnitude stochastic jumps at intermediate abundances, which could
357 be another organizing force on the communities because movement towards equilibrium can
358 generate correlation. Species should on average be at equilibrium and move about it randomly,
359 but they are knocked farther away due to stochasticity when they are near equilibrium. While
360 we use uncorrelated Brownian motion in these models, it is possible that the stochasticity alone
361 could generate covariance if the diffusion terms are not orthogonal.

362 (d) Comparison of metrics

363 Covariance is in many cases the most successful metric at detecting species interactions. However,
364 we find that it is in general the metric that is most likely to detect non-random structure, including
365 false networks. Consequently, correlation and inverse correlation, which are normalized by the
366 single species variance at times outperform covariance in detecting the true network, with the
367 caveat that these metrics infer each of the possible networks with more uniform probability. The
368 LIMITS algorithm had lower sensitivity than the others methods when interactions were weak,
369 but increased as interaction strength increases, but within similar ranges as for the other methods.
370 Specificity, on the other hand remained low, due to the low threshold necessary for detecting
371 the interspecific interactions when they are weaker than or the same magnitude as intraspecific
372 interactions. It is important to note that covariance and correlation only detect symmetric matrices
373 while LIMITS can detect asymmetric matrices. This may affect the performance of the metrics and
374 affects interpretation of the results. Evaluating a large network simulated using the niche model
375 [22], which generates networks that are similar in their properties to observed food webs, and

376 the correlation method, shows that sensitivity decreases and specificity increases as interactions
377 strength increases. The large values of associated eigenvalues implies a larger role for transients
378 in affecting the correlation structure.

379 Inverse correlation removes indirect interactions by minimizing the large scale trends. In
380 the simulations presented here, in which there are no deterministic external forcing or “hubs”
381 controlling species interactions, correlation and inverse correlation perform similarly. The inverse
382 correlation method might perform better with a more complex network of species interactions.
383 The covariance metric is able to make use of species relative abundances while the normalization
384 used by correlation removes this information. As we show in Figure 5, species relative abundances
385 do provide meaningful information on network geometry to the extent that the species abundance
386 distribution is affected by the network structure. As shown in this figure, the same species under
387 the same parameters changes in abundance distribution due to network geometry. It is important
388 to note that the interaction coefficients are symmetric and that adding additional complexity to
389 this model by working with asymmetric community matrices may alter the reported simulation
390 results. However, the mathematical exposition is agnostic to the structure of the interaction
391 matrix.

392 This study suggests a few paths forward for the development of improved metrics of species
393 interactions. The metrics currently in use perform best when variables are normally distributed.
394 While use of correlation does not require that sample values are normally distributed, it is only an
395 exhaustive measure of association if the joint distribution of the samples is a multivariate normal.
396 The maximum entropy method used here similarly makes the approximation that the samples
397 values are drawn from a normal distribution, however the maximum entropy method could be
398 generalized to be appropriate for the observed abundance distributions. In our case the samples
399 (and hence the joint distribution) are not normally distributed, which is true for most communities
400 [18, 33]. Metrics based on the abundance distribution of biological species across different trophic
401 levels might be better suited for network detection. For example, a linkage disequilibrium metric
402 in genetics is specialized for beta distributed observations [34]. In addition, this system is in a non-
403 equilibrium steady state and the absence of detailed balance means that in addition to maximizing
404 the entropy, understanding entropy production could help to detect the true network. A distinct
405 approach would be to use parameter estimation or constrained optimization to fit a model to
406 observed time series.

407 6. Conclusions

408 Whole static networks of species interactions are detectable using existing methods for network
409 inference including a maximum entropy method, but with low probability. Our finding that
410 specificity and sensitivity do not differ significantly from random while there is non-random
411 selection of network motifs demonstrates that it is important to consider not only the success in
412 detecting pairwise interactions but also the way in which correlation metrics may systematically
413 select certain sets of edges. Counter intuitively, increased stochasticity does not necessarily make
414 detection of interactions between species less likely. Instead, the path that the system takes to
415 equilibrium once perturbed is determined by the links between species, however this path does
416 not necessarily facilitate detection by the existing metrics. Existing metrics have systematic biases.
417 Indirect interactions may be more likely to be detected than direct interactions, even using inverse
418 correlation. This is particularly true for systems prone to long transients.

419 While we focus on the problem of ecological network inference in this paper, network
420 inference is an important tool in other domains including genetic networks [35, 36]. There are
421 many network inference techniques for researchers to choose including those based on machine
422 learning [37]. We recommend that the mathematical and scientific basis of each of these techniques
423 be evaluated carefully before their application in new domains.

424 **Data Accessibility.** The simulation codes are available as supplementary material

425 **Authors' Contributions.** MAF, RR, and PAM designed the study, MAF and RR designed the stochastic
426 model and performed the theoretical derivations, MAF ran the numerical experiments, MAF and DC

performed the data analysis, MAF, RR, and PAM analyzed the results, MAF, RR, and PAM wrote the manuscript

Funding. This work was supported by a Fulbright Student Grant awarded to MAF. PAM acknowledges support from project FONDECYT 1161023, FONDECYT 1200925, and AFB 170008 from CONICYT, RR acknowledges partial support of CONICYT-Proyecto Redes 180018 and CIMFAV-CIDI.

Acknowledgements. We thank the Fulbright Chile program for research support and the Santa Fe Institute for hosting a research visit. The authors also appreciate insightful conversations with Sergio Rojas. We acknowledge constructive comments from two anonymous reviewers.

References

- 1 Dunne JA. The network structure of food webs. *Ecological networks: linking structure to dynamics in food webs.* 2006;p. 27–86.
- 2 Bascompte J, Jordano P. *Mutualistic networks.* Princeton University Press; 2013.
- 3 Stephens CR, Heau JG, González C, Ibarra-Cerdeña CN, Sánchez-Cordero V, González-Salazar C. Using biotic interaction networks for prediction in biodiversity and emerging diseases. *PLoS One.* 2009;4(5):e5725.
- 4 Steele JA, Countway PD, Xia L, Vigil PD, Beman JM, Kim DY, et al. Marine bacterial, archaeal and protistan association networks reveal ecological linkages. *The ISME journal.* 2011;5(9):1414–1425.
- 5 Faust K, Raes J. Microbial interactions: from networks to models. *Nature Reviews Microbiology.* 2012;10(8):538–550.
- 6 Faust K, Lahti L, Gonze D, de Vos WM, Raes J. Metagenomics meets time series analysis: unraveling microbial community dynamics. *Current opinion in microbiology.* 2015;25:56–66.
- 7 Barner AK, Coblenz KE, Hacker SD, Menge BA. Fundamental contradictions among observational and experimental estimates of non-trophic species interactions. *Ecology.* 2018;99(3):557–566.
- 8 Freilich MA, Wieters E, Broitman BR, Marquet PA, Navarrete SA. Species co-occurrence networks: Can they reveal trophic and non-trophic interactions in ecological communities? *Ecology.* 2018;99(3):690–699.
- 9 Weiss S, Van Treuren W, Lozupone C, Faust K, Friedman J, Deng Y, et al. Correlation detection strategies in microbial data sets vary widely in sensitivity and precision. *The ISME journal.* 2016;.
- 10 Berry D, Widder S. Deciphering microbial interactions and detecting keystone species with co-occurrence networks. *Frontiers in microbiology.* 2014;5:219.
- 11 Lezon TR, Banavar JR, Cieplak M, Maritan A, Fedoroff NV. Using the principle of entropy maximization to infer genetic interaction networks from gene expression patterns. *Proceedings of the National Academy of Sciences.* 2006;103(50):19033–19038.
- 12 Azaele S, Muneeppeerakul R, Rinaldo A, Rodriguez-Iturbe I. Inferring plant ecosystem organization from species occurrences. *Journal of theoretical biology.* 2010;262(2):323–329.
- 13 Yodzis P. The indeterminacy of ecological interactions as perceived through perturbation experiments. *Ecology.* 1988;69(2):508–515.
- 14 Gotelli NJ, Graves GR. *Null models in ecology;* 1996.
- 15 Freilich MA, Connolly SR. Phylogenetic community structure when competition and environmental filtering determine abundances. *Global ecology and biogeography.* 2015;24(12):1390–1400.
- 16 Gilarranz LJ, Hastings A, Bascompte J. Inferring topology from dynamics in spatial networks. *Theoretical ecology.* 2015;8(1):15–21.
- 17 Coenen AR, Weitz JS. Limitations of Correlation-Based Inference in Complex Virus-Microbe Communities. *mSystems.* 2018;3(4):e00084–18.
- 18 Marquet PA, Espinoza G, Abades SR, Ganz A, Rebolledo R. On the proportional abundance of species: Integrating population genetics and community ecology. *Scientific Reports (Nature Publisher Group).* 2017;7:1–10.

- 478 19 Fisher CK, Mehta P. The transition between the niche and neutral regimes in ecology.
479 Proceedings of the National Academy of Sciences. 2014;111(36):13111–13116.
- 480 20 Haegeman B, Loreau M. A mathematical synthesis of niche and neutral theories in community
481 ecology. *Journal of theoretical biology*. 2011;269(1):150–165.
- 482 21 Rebolledo R, Navarrete SA, Kéfi S, Rojas S, Marquet PA. An Open-System Approach to
483 Complex Biological Networks. *SIAM Journal on Applied Mathematics*. 2019;79(2):619–640.
- 484 22 Williams RJ, Martinez ND. Simple rules yield complex food webs. *Nature*. 2000;404(6774):180–
485 183.
- 486 23 Fisher CK, Mehta P. Identifying keystone species in the human gut microbiome from
487 metagenomic timeseries using sparse linear regression. *PloS one*. 2014;9(7).
- 488 24 Röttgers L, Faust K. From hairballs to hypotheses—biological insights from microbial networks.
489 *FEMS microbiology reviews*. 2018;42(6):761–780.
- 490 25 McCullagh P. Generalized linear models. *European Journal of Operational Research*.
491 1984;16(3):285–292.
- 492 26 Angulo MT, Moreno JA, Lippner G, Barabási AL, Liu YY. Fundamental limitations
493 of network reconstruction from temporal data. *Journal of the Royal Society Interface*.
494 2017;14(127):20160966.
- 495 27 Farrell BF, Ioannou PJ. Generalized stability theory. Part I: Autonomous operators. *Journal of*
496 *the atmospheric sciences*. 1996;53(14):2025–2040.
- 497 28 Neubert MG, Caswell H. Alternatives to resilience for measuring the responses of ecological
498 systems to perturbations. *Ecology*. 1997;78(3):653–665.
- 499 29 Milo R, Shen-Orr S, Itzkovitz S, Kashtan N, Chklovskii D, Alon U. Network motifs: simple
500 building blocks of complex networks. *Science*. 2002;298(5594):824–827.
- 501 30 Proulx SR, Promislow DE, Phillips PC. Network thinking in ecology and evolution. *Trends in*
502 *Ecology & Evolution*. 2005;20(6):345–353.
- 503 31 Deng Y, Jiang YH, Yang Y, He Z, Luo F, Zhou J. Molecular ecological network analyses. *BMC*
504 *bioinformatics*. 2012;13(1):113.
- 505 32 Lipinski-Kruszka J, Stewart-Ornstein J, Chevalier MW, El-Samad H. Using dynamic noise
506 propagation to infer causal regulatory relationships in biochemical networks. *ACS synthetic*
507 *biology*. 2015;4(3):258–264.
- 508 33 Harte J, Kinzig A, Green J. Self-similarity in the distribution and abundance of species. *Science*.
509 1999;284(5412):334–336.
- 510 34 Gianola D, Manfredi E, Simianer H. On measures of association among genetic variables.
511 *Animal genetics*. 2012;43(s1):19–35.
- 512 35 Feist AM, Herrgård MJ, Thiele I, Reed JL, Palsson BØ. Reconstruction of biochemical networks
513 in microorganisms. *Nature Reviews Microbiology*. 2009;7(2):129.
- 514 36 Dondelinger F, Mukherjee S. Statistical network inference for time-varying molecular data
515 with dynamic bayesian networks. In: *Gene Regulatory Networks*. Springer; 2019. p. 25–48.
- 516 37 Camacho DM, Collins KM, Powers RK, Costello JC, Collins JJ. Next-generation machine
517 learning for biological networks. *Cell*. 2018;173(7):1581–1592.