

Using Cooccurrence Networks

Jim Brunner

Los Alamos National Laboratory

Coincidence Network

Constructing a coincidence network. I map the abundances according to

$$a(r_{ji}) = \begin{cases} \lfloor \left(\frac{r_{ji}}{\max_{s_k}(r_{jk})} \right) n \rfloor + 1 & \frac{r_{ji}}{\max_{s_k}(r_{jk})} \geq m \\ 0 & \frac{r_{ji}}{\max_{s_k}(r_{jk})} < m \end{cases}$$

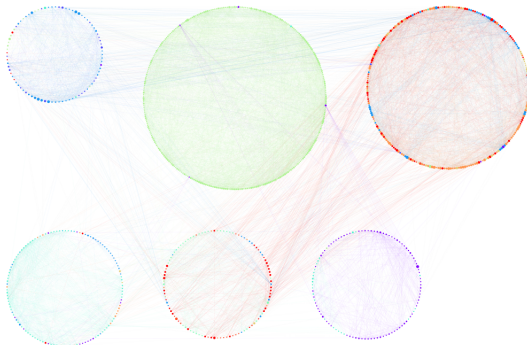
into “bins” relative to the maximum that taxa appears. Then count

$$w_{jk}^1 = \frac{\|\{i : a(r_{ji}) = a(r_{ki}) \neq 0\}\|}{S}$$

how often two organisms appear in the same bin.

Cooccurrence Network

Same idea but now edges weights are compared to a random graph (null model). So we only keep edges that have a higher than “random” weight.



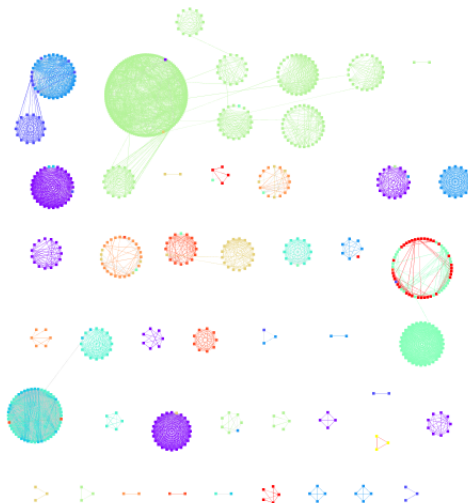
Cooccurrence Network - Pearson Correlation

We can also compute the Pearson correlation coefficient between taxa across the samples.

$$\rho_{xy} = \frac{1}{N} \frac{(\mathbf{x} - \mu_x \mathbf{1}) \cdot (\mathbf{y} - \mu_y \mathbf{1})}{\sigma_x \sigma_y}$$

Then, we keep an edge if $\rho > 0.8$ and $p < 0.05$, where p is the chance of correlation higher than ρ_{xy} in a random model. The random model assigns abundances as a binomial with parameters determined by sample and taxa. The p values are calculated using a Monte Carlo simulation with 1000 trials. The species level network using this method had $\sim 1/2$ as many edges.

Cooccurrence Network - Pearson Correlation



Clustering

We can cluster the network to attempt to determine or evaluate meaningful groups of taxa.

- Community clustering
 - Minimize a function based on “edge betweenness”.
Determines edges that are between clusters.
- Spectral clustering
 - Performs a random walk on the network.

Analyzing a sample

- Filter GOTTCHA results - try to determine probability of seeing groups of organisms - Random Markov Field
- Diffusion on graph

$$\frac{\partial}{\partial t}u(v, t) = -Lu(v, t)$$

where v takes values in the vertex set of the graph. Then, we can encode “known” information in three ways: initial values, boundary values, or a forcing vector.

Analyzing a sample - Diffusion Process

$$\frac{\partial}{\partial t}u(v, t) = -Lu(v, t)$$

where v takes values in the vertex set of the graph. Then, we can encode sample information in three ways: initial values, boundary values, or a forcing vector.

Related to spectral clustering - spectral clustering groups by “distance” in first k eigenmodes of diffusion process. This process gives “distance” in weighted sum of eigenmodes of diffusion process.

Highest ranked nodes are “closest” to the sample data

Initial Values

Initial Value Problem

Let $u_i(t)$ be the solution at node v_i to the discrete diffusion problem

$$\frac{d}{dt}\mathbf{u}(t) = -L\mathbf{u}$$

where L is the graph laplacian with initial conditions determined by sample information.

Then, if \mathbf{K} is the information “known” from the sample and the values of v_k and v_l are unknown,

$$\int_0^\infty u_k(t)dt - \int_0^\infty u_l(t)dt > 0 \Rightarrow P(v_k = 1|\mathbf{K}) > P(v_l = 1|\mathbf{K})$$

Boundary Values

Boundary Value Problem

Let $u_i(t)$ be the solution at node v_i to the discrete diffusion problem

$$\frac{d}{dt}\mathbf{u}(t) = -L\mathbf{u}$$

where L is the graph laplacian with fixed values (which can be regarded as boundary values) $u_i = 1$ if node v_i is known to be “on”, $u_j = 0$ if v_j is known to be “off”.

Then, if \mathbf{K} is the information “known” and the values of v_k and v_l are unknown, and $\tilde{\mathbf{u}}$ is the equilibrium solution to the diffusion problem,

$$\tilde{u}_k dt > \tilde{u}_l \Leftrightarrow P(v_k = 1|\mathbf{K}) > P(v_l = 1|\mathbf{K})$$

Forcing Function

Forced Problem

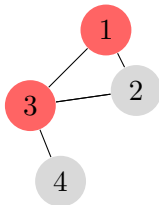
Let $u_i(t)$ be the solution at node v_i to the discrete diffusion problem

$$\frac{d}{dt}\mathbf{u}(t) = -L\mathbf{u} + \mathbf{f}$$

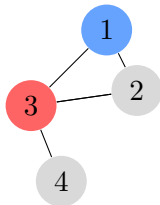
where L is the graph laplacian and \mathbf{f} a forcing vector with $f_i = \alpha_{cc}$ if node v_i is known to be “on”, $f_j = -\beta_{cc}$ if v_j is known to be “off”, where cc denotes a connected component of the graph. We choose α_{cc} and β_{cc} so that on any connected component cc , $\sum \alpha_{cc} = \sum \beta_{cc} = 1$. Then, if \mathbf{K} is the information “known” and the values of v_k and v_l are unknown, and $\tilde{\mathbf{u}}$ is the equilibrium solution to the diffusion problem,

$$\tilde{u}_k > \tilde{u}_l \Leftrightarrow P(v_k = 1 | \mathbf{K}) > P(v_l = 1 | \mathbf{K})$$

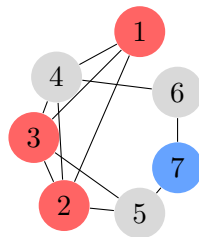
Small Network Examples



(a) Small network,
two on.



(b) small network,
one on one off



(c) Larger network

Figure: Test networks. Gray nodes were “unknown”, blue nodes “off”, and red “on”.

Small Network Examples

Configuration	Method	Ranking	Ties
(a) Small Network, two on.	IVP	2, 4	none
	BVP	4, 2	4, 2
	Forcing	2, 4	none
(b) Small Network, one on one off	IVP	4, 2	none
	BVP	4, 2	none
	Forcing	4, 2	none
(c) Larger network	IVP	4, 5, 6	none
	BVP	4, 5, 6	none
	Forcing	4, 5, 6	none

Network Examples

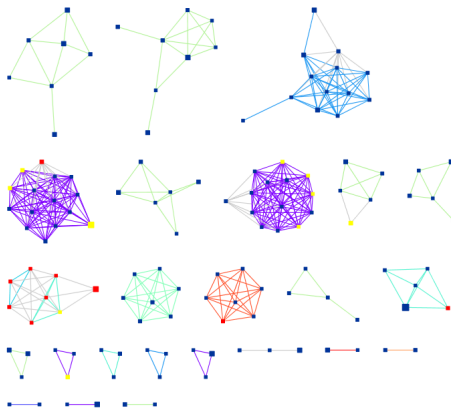


Figure: The “known” set - blue is off and red is on, while yellow is unknown.

Network Examples

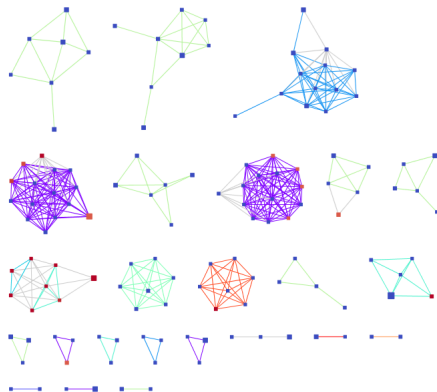


Figure: Result of boundary value problem. Hotter colors indicate higher likelihood, with red indicating assumed “on”.

Assigning a sample

Given a sample, we get ranking from a network N_j . Let $\mathbf{r}^j(\mathbf{s})$ be the ranking given by diffusion on network j . Then,

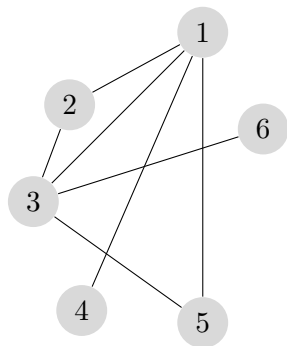
$$F_j(\mathbf{s}) = F(\mathbf{s}, \mathbf{r}^j(\mathbf{s})) = \sum_{i=1}^n c_i^{\mathbf{r}^j} \frac{s_i}{\|\mathbf{s}\|_1}$$

and we can attempt to optimize over the networks we have (if there's only a few that's easy). The conjecture is then

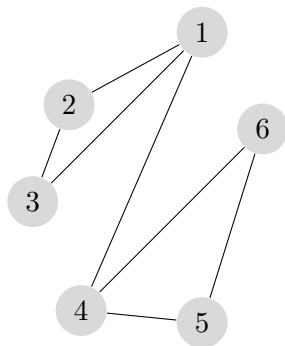
Conjecture

Assume that $F_j(\mathbf{s}) > F_k(\mathbf{s})$. Then $P(\mathbf{s}|N_j) > P(\mathbf{s}|N_k)$.

Assigning a sample



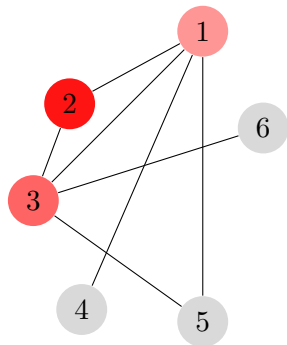
(a) Network A_1



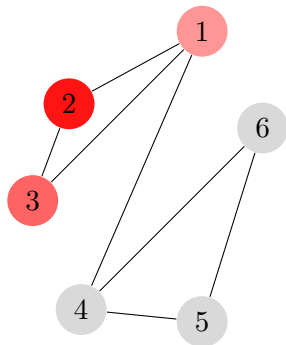
(b) Network A_2

Figure: Test Networks

Assigning a sample



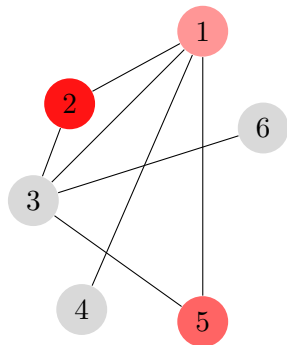
(a) Network A_1



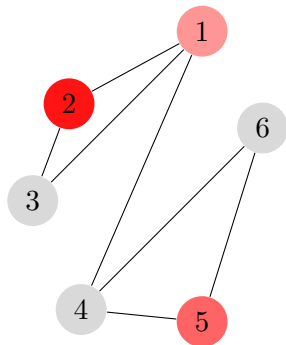
(b) Network A_2

Figure: “Sample” $u_1 = (1/6, 1/2, 1/3, 0, 0, 0)$

Assigning a sample



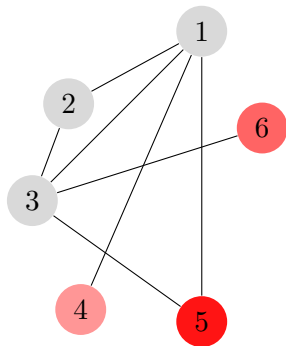
(a) Network A_1



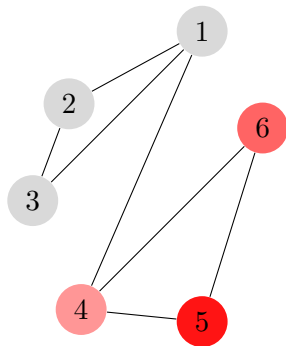
(b) Network A_2

Figure: “Sample” $u_2 = (1/6, 1/2, 0, 0, 1/3, 0)$

Assigning a sample



(a) Network A_1

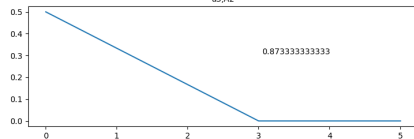
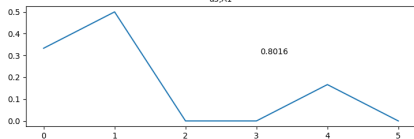
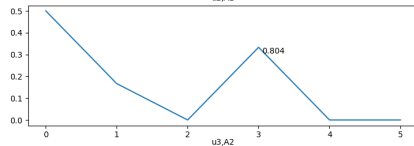
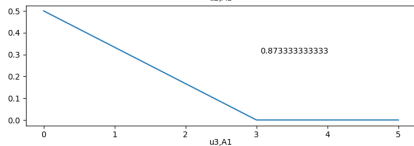
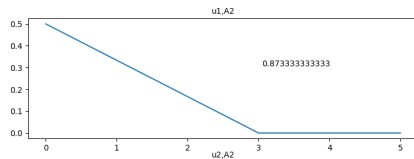
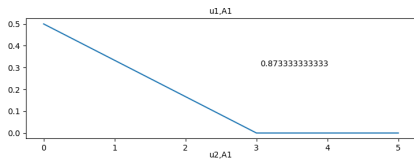


(b) Network A_2

Figure: “Sample” $u_3 = (0, 0, 0, 1/6, 1/2, 1/3)$

Assigning a sample

Abundance v Rank



Assigning a sample

