

Motifs and bipartite networks

S. Robin

Sorbonne université

Ecole EcoNet, April 2024, Montpellier

Outline

Bipartite networks and motifs

A null model

Motif distribution

Goodness-of-fit and network comparison

Network embedding

Outline

Bipartite networks and motifs

A null model

Motif distribution

Goodness-of-fit and network comparison

Network embedding

Bipartite network

Two types of actors.

- ▶ Mutualistic: plant-pollinator
- ▶ Antagonistic: host-parasite

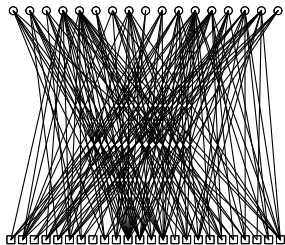
Topological analysis:

understanding the network organisation

Local: node or edge properties (degree, betweenness)

Global: density, connected components, nestedness

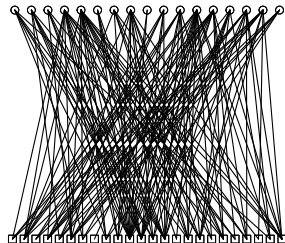
Zackenberg network: [SROB16]



Bipartite network: notations

Species.

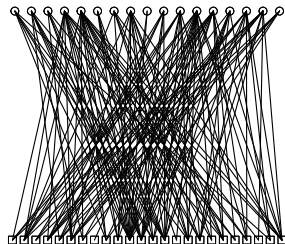
- ▶ $i = 1, \dots, m$ pollinators = rows = bottom nodes
- ▶ $j = 1, \dots, n$ plants = columns = top nodes



Bipartite network: notations

Species.

- ▶ $i = 1, \dots, m$ pollinators = rows = bottom nodes
- ▶ $j = 1, \dots, n$ plants = columns = top nodes



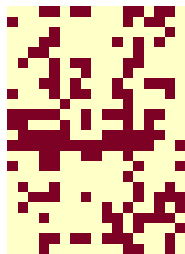
Interactions.

- ▶ $A_{ij} = 1$ if pollinator i interacts with plant j ,
0 otherwise

$$A_{ij} = 1 \Leftrightarrow i \sim j$$

- ▶ adjacency matrix : $m \times n$

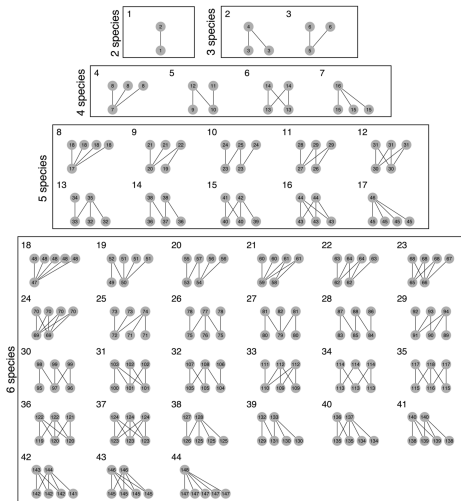
$$A = [A_{ij}]_{1 \leq i \leq m, 1 \leq j \leq n}$$



Bipartite motifs

'Meso-scale' analysis. [SCB⁺19]

- ▶ Motifs = 'building-blocks'
- ▶ between local (several nodes) and global (sub-graph)



Bipartite motifs

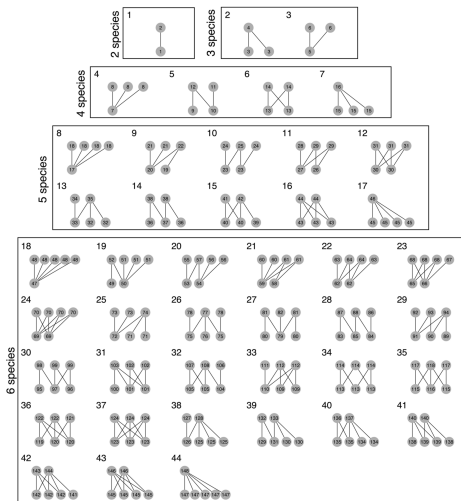
'Meso-scale' analysis. [SCB⁺19]

- ▶ Motifs = 'building-blocks'
- ▶ between local (several nodes) and global (sub-graph)

Interest.

- ▶ Generic description of a network
- ▶ Enables network comparison
- ▶ Even when the nodes are different

(+ 'species-role': out of the scope here)



Bipartite motifs

'Meso-scale' analysis. [SCB⁺19]

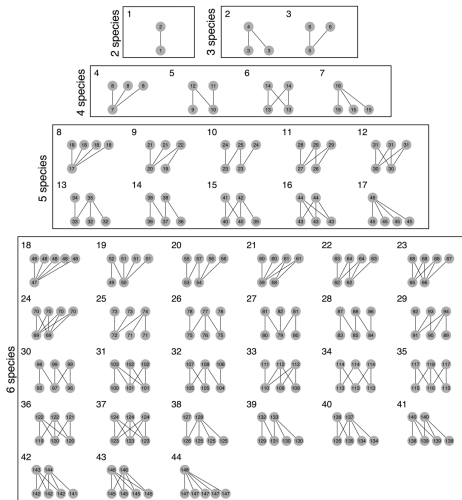
- ▶ Motifs = 'building-blocks'
- ▶ between local (several nodes) and global (sub-graph)

Interest.

- ▶ Generic description of a network
- ▶ Enables network comparison
- ▶ Even when the nodes are different

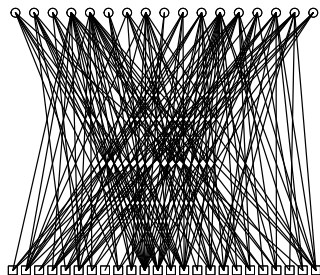
(+ 'species-role': out of the scope here)

Existing tool. `bmotif` package [SSS⁺19]:
counts motif occurrences



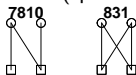
Example

Plant-pollinator network [SROB16]

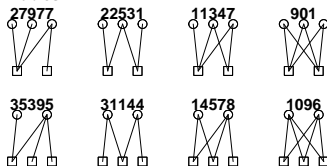


Motif counts.

4 nodes (species)

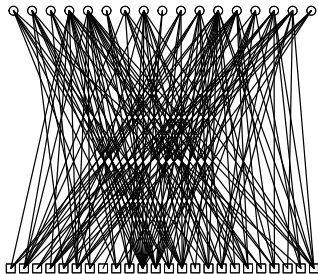


5 nodes



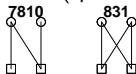
Example

Plant-pollinator network [SROB16]

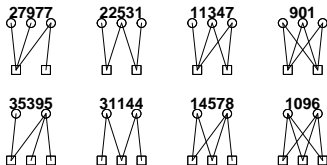


Motif counts.

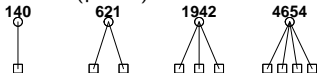
4 nodes (species)



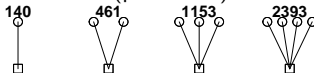
5 nodes



top 'stars' (plants)



bottom 'stars' (pollinators)



Outline

Bipartite networks and motifs

A null model

Motif distribution

Goodness-of-fit and network comparison

Network embedding

Need for a null model

Motif counts obviously depend on

- ▶ the **size** of the network: $n \times m$
- ▶ the **density** of the network
- ▶ the **imbalance** between bottom-node degrees (specialist vs generalist pollinators)
- ▶ the **imbalance** between top-node degrees (specialist vs generalist plants)

Need for a null model

Motif counts obviously depend on

- ▶ the **size** of the network: $n \times m$
- ▶ the **density** of the network
- ▶ the **imbalance** between bottom-node degrees (specialist vs generalist pollinators)
- ▶ the **imbalance** between top-node degrees (specialist vs generalist plants)

Bipartite expected degree distribution (BEDD) model: (in words)

- ▶ Consider m pollinators ($i = 1, \dots, m$):
each plant i has a specific propensity to interact (degree of generalism)
- ▶ Consider n plants ($j = 1, \dots, n$):
each plant j has a specific propensity to interact (idem)
- ▶ The probability for pollinator i and plant j to interact is proportional to the product of their respective propensities.

BEDD model

Bipartite expected degree distribution (BEDD) model: (formaly)

- ▶ ρ = network density
- ▶ g = top node degree imbalance ($\int g = 1$)
- ▶ h = bottom node degree imbalance ($\int h = 1$)

$$\{U_i\}_{i=1,\dots,m} \text{ iid } \sim \mathcal{U}[0, 1] \quad \{V_j\}_{j=1,\dots,n} \text{ iid } \sim \mathcal{U}[0, 1]$$

$$\mathbb{P}\{i \sim j \mid U_i, V_j\} = \rho g(U_i) h(V_j)$$

(Bipartite version of the EDD model [CL02])

BEDD model

Bipartite expected degree distribution (BEDD) model: (formaly)

- ▶ ρ = network density
- ▶ g = top node degree imbalance ($\int g = 1$)
- ▶ h = bottom node degree imbalance ($\int h = 1$)

$$\{U_i\}_{i=1,\dots,m} \text{ iid } \sim \mathcal{U}[0, 1] \quad \{V_j\}_{j=1,\dots,n} \text{ iid } \sim \mathcal{U}[0, 1]$$

$$\mathbb{P}\{i \sim j \mid U_i, V_j\} = \rho g(U_i) h(V_j)$$

(Bipartite version of the EDD model [CL02])

Model parameters:

$$\theta = (\rho, g, h).$$

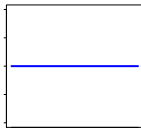
BEDD model

$$\mathbb{P}\{i \sim j \mid U_i, V_j\} = \rho g(U_i) h(V_j)$$

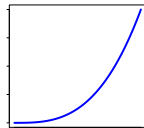
$$\mathbb{E}(D_i \mid U_i) = n \rho g(U_i)$$

$$\mathbb{E}(D_j \mid V_j) = m \rho g(V_j)$$

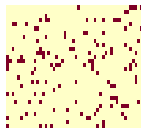
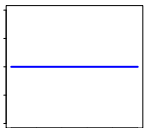
$$h_0(v) =$$



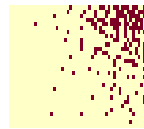
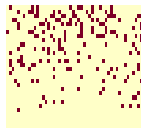
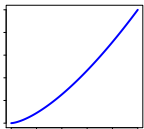
$$h(v) =$$



$$g_0(u) =$$



$$g(u) =$$



Properties of the BEDD model

Assumptions.

- ▶ No preferred or avoided specific connexion
- ▶ **Graph-exchangeable** model: pollinators can be permuted and plants can be permuted

Properties of the BEDD model

Assumptions.

- ▶ No preferred or avoided specific connexion
- ▶ **Graph-exchangeable** model: pollinators can be permuted and plants can be permuted

Properties.

- ▶ Expected degree for pollinator i given U_i : $n \rho g(U_i)$.
- ▶ Expected degree for plant j given V_j : $m \rho h(V_j)$.
- ▶ 'Nested' structure by construction

Properties of the BEDD model

Assumptions.

- ▶ No preferred or avoided specific connexion
- ▶ **Graph-exchangeable** model: pollinators can be permuted and plants can be permuted

Properties.

- ▶ Expected degree for pollinator i given $U_i : n \rho g(U_i)$.
- ▶ Expected degree for plant j given $V_j : m \rho h(V_j)$.
- ▶ 'Nested' structure by construction

Sufficient statistics to fit BEDD:

- ▶ Pollinator degrees + plant degrees
- ▶ or, equivalently, star (single edge, top, bottom) frequencies

Outline

Bipartite networks and motifs

A null model

Motif distribution

Goodness-of-fit and network comparison

Network embedding

Counting motifs

Number of 'positions'.

- ▶ Choose p nodes among m
- ▶ Choose q nodes among n
- ▶ Try all *automorphisms*

$$c_s := \binom{m}{p} \times \binom{n}{q} \times r_s$$

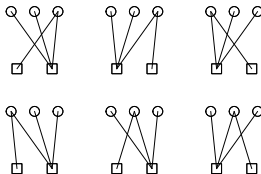
Counting motifs

Number of 'positions'.

- ▶ Choose p nodes among m
- ▶ Choose q nodes among n
- ▶ Try all *automorphisms*

$$c_s := \binom{m}{p} \times \binom{n}{q} \times r_s$$

Automorphisms= non-redundant permutations



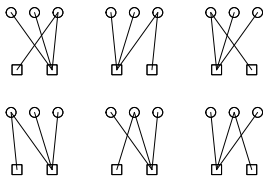
Counting motifs

Number of 'positions'.

- ▶ Choose p nodes among m
- ▶ Choose q nodes among n
- ▶ Try all *automorphisms*

$$c_s := \binom{m}{p} \times \binom{n}{q} \times r_s$$

Automorphisms= non-redundant permutations



Motif count. Try all positions $\alpha = 1, \dots, c_s$, define

$$Y_{s\alpha} = 1 \text{ if match, } \quad 0 \text{ otherwise,}$$

then count the number of matches:

$$N_s = \sum_{\alpha} Y_{s\alpha}$$

→ **Motif frequency:** $F_s := N_s/c_s$

Motif probability

Occurrence probability $\bar{\phi}_s = \mathbb{P}\{Y_{s\alpha} = 1\}$. Under the B-EDD model [OLR22]:

$$\left(\begin{array}{c} \circ \\ \circ \\ \circ \\ \square \\ \square \end{array} \right) = \text{_____}$$

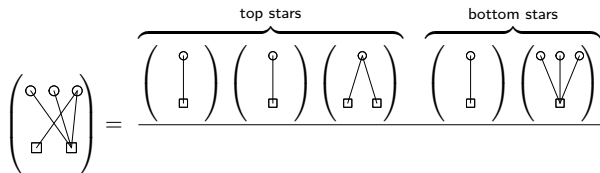
Motif probability

Occurrence probability $\bar{\phi}_s = \mathbb{P}\{Y_{s\alpha} = 1\}$. Under the B-EDD model [OLR22]:

$$\left(\begin{array}{c} \circ \\ \circ \\ \circ \\ \square \\ \square \end{array} \right) = \frac{\overbrace{\left(\begin{array}{c} \circ \\ \square \end{array} \right) \left(\begin{array}{c} \circ \\ \square \end{array} \right) \left(\begin{array}{c} \circ \\ \square \\ \square \end{array} \right)}^{\text{top stars}}}{\hspace{15em}}$$

Motif probability

Occurrence probability $\bar{\phi}_s = \mathbb{P}\{Y_{s\alpha} = 1\}$. Under the B-EDD model [OLR22]:



Motif probability

Occurrence probability $\bar{\phi}_s = \mathbb{P}\{Y_{s\alpha} = 1\}$. Under the B-EDD model [OLR22]:

$$\left(\begin{array}{c} \circ \\ \circ \\ \circ \\ \square \\ \square \end{array} \right) = \frac{\overbrace{\left(\begin{array}{c} \circ \\ \square \end{array} \right) \left(\begin{array}{c} \circ \\ \square \end{array} \right) \left(\begin{array}{c} \circ \\ \square \\ \square \end{array} \right)}^{\text{top stars}} \overbrace{\left(\begin{array}{c} \circ \\ \square \end{array} \right) \left(\begin{array}{c} \circ \\ \circ \\ \square \end{array} \right)}^{\text{bottom stars}}}{\underbrace{\left(\begin{array}{c} \circ \\ \square \end{array} \right)^4}_{\text{edges}}}$$

Motif probability

Occurrence probability $\bar{\phi}_s = \mathbb{P}\{Y_{s\alpha} = 1\}$. Under the B-EDD model [OLR22]:

$$\begin{aligned}
 \binom{\text{graph}}{=} &= \frac{\overbrace{\left(\binom{\text{star}_1}{\square} \right) \left(\binom{\text{star}_2}{\square} \right) \left(\binom{\text{star}_3}{\square \square} \right)}^{\text{top stars}} \overbrace{\left(\binom{\text{star}_4}{\square} \right) \left(\binom{\text{star}_5}{\square} \right)}^{\text{bottom stars}}}{\underbrace{\left(\binom{\text{edge}}{\square} \right)^4}_{\text{edges}}} \\
 \bar{\phi}_s = \mathbb{P}^{BEDD} \binom{\text{graph}}{=} &= \frac{(\phi_1^2 \phi_2) (\phi_1 \phi_4)}{(\phi_1)^4} = \frac{\phi_2 \phi_4}{\phi_1}
 \end{aligned}$$

Motif probability

Occurrence probability $\bar{\phi}_s = \mathbb{P}\{Y_{s\alpha} = 1\}$. Under the B-EDD model [OLR22]:

$$\begin{aligned}
 \left(\begin{array}{c} \circ \\ \circ \\ \circ \\ \square \\ \square \end{array} \right) &= \frac{\overbrace{\left(\begin{array}{c} \circ \\ \square \end{array} \right) \left(\begin{array}{c} \circ \\ \square \end{array} \right) \left(\begin{array}{c} \circ \\ \square \square \end{array} \right)}^{\text{top stars}} \overbrace{\left(\begin{array}{c} \circ \\ \square \end{array} \right) \left(\begin{array}{c} \circ \circ \\ \square \end{array} \right)}^{\text{bottom stars}}}{\underbrace{\left(\begin{array}{c} \circ \\ \square \end{array} \right)^4}_{\text{edges}}} \\
 \bar{\phi}_s = \mathbb{P}^{BEDD} \left(\begin{array}{c} \circ \\ \circ \\ \circ \\ \square \\ \square \end{array} \right) &= \frac{(\phi_1^2 \phi_2) (\phi_1 \phi_4)}{(\phi_1)^4} = \frac{\phi_2 \phi_4}{\phi_1}
 \end{aligned}$$

Estimated probability \bar{F}_s .

$$\bar{\phi}_s := \frac{\phi_2 \phi_4}{\phi_1} \quad \rightarrow \quad \bar{F}_s := \frac{F_2 F_4}{F_1}$$

where F_1, F_2, F_4 = observed frequencies of edges, top stars and bottom stars.

Moments of the count

- ▶ Number of positions: c_s
- ▶ Mean: $\mathbb{E}_{BEDD}(N_s) = c_s \times \bar{\phi}_s$

Moments of the count

► Number of positions: c_s

► Mean: $\mathbb{E}_{BEDD}(N_s) = c_s \times \bar{\phi}_s$

► Variance: Same game, requires to evaluate $\mathbb{E}_{BEDD}(N_s^2) = \mathbb{E}_{BEDD} \left(\sum_{\alpha} Y_{s\alpha} \right)^2$

→ Need to account for overlap between positions (*super-motifs*: [PDK⁺08])



→ Compute the respective expected count in the way as for regular motifs

Moments of the count

► **Number of positions:** c_s

► **Mean:** $\mathbb{E}_{BEDD}(N_s) = c_s \times \bar{\phi}_s$

► **Variance:** Same game, requires to evaluate $\mathbb{E}_{BEDD}(N_s^2) = \mathbb{E}_{BEDD} \left(\sum_{\alpha} Y_{s\alpha} \right)^2$

→ Need to account for overlap between positions (*super-motifs*: [PDK⁺08])



→ Compute the respective expected count in the way as for regular motifs

► **Covariance:** Same game to compute $\text{Cov}(N_s, N_{s'})$

Distribution of the count

Asymptotic normality for non-star motifs. Under BEDD (and sparsity conditions):

$$(F_s - \bar{F}_s) / \sqrt{\widehat{V}(F_s)} \xrightarrow{m, n \rightarrow \infty} \mathcal{N}(0, 1)$$

Proof [OLR22]:

$$\text{decompose } F_s - \bar{F}_s = \underbrace{(F_s - \phi_s)}_{\text{random fluctuations}} + \underbrace{(\phi_s - \bar{\phi}_s)}_{\text{null under BEDD}} + \underbrace{(\bar{\phi}_s - \bar{F}_s)}_{\text{estimation error } \rightarrow 0},$$

+ construct a counting martingale for $F_s - \phi_s$ [GL17].

Distribution of the count

Asymptotic normality for non-star motifs. Under BEDD (and sparsity conditions):

$$(F_s - \bar{F}_s) / \sqrt{\widehat{V}(F_s)} \xrightarrow{m, n \rightarrow \infty} \mathcal{N}(0, 1)$$

Proof [OLR22]:

$$\text{decompose } F_s - \bar{F}_s = \underbrace{(F_s - \phi_s)}_{\text{random fluctuations}} + \underbrace{(\phi_s - \bar{\phi}_s)}_{\text{null under BEDD}} + \underbrace{(\bar{\phi}_s - \bar{F}_s)}_{\text{estimation error } \rightarrow 0},$$

+ construct a counting martingale for $F_s - \phi_s$ [GL17].

Consequence. We know

- ▶ the expected behavior (mean, variance, distribution) of any motif count
- ▶ under the BEDD model (= 'null model').

Outline

Bipartite networks and motifs

A null model

Motif distribution

Goodness-of-fit and network comparison

Network embedding

Goodness-of-fit (GOF)

Aim of GOF tests. Test if the observed data arise from a given model.

More cautious: 'if the model fits the data reasonably well'

Goodness-of-fit (GOF)

Aim of GOF tests. Test if the observed data arise from a given model.

More cautious: 'if the model fits the data reasonably well'

Typical approach.

1. Define some statistic (= function of the data) T ,
2. Establish the distribution of T under the model,
3. Compare the observed value of T with its distribution under the model.

Goodness-of-fit (GOF)

Aim of GOF tests. Test if the observed data arise from a given model.

More cautious: 'if the model fits the data reasonably well'

Typical approach.

1. Define some statistic (= function of the data) T ,
2. Establish the distribution of T under the model,
3. Compare the observed value of T with its distribution under the model.

Example.

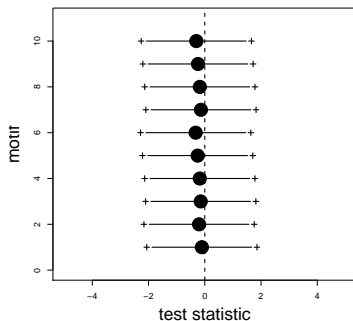
1. Data = observed plant-pollinator network
2. Statistic T = motif count N_5
3. Model = BEDD

Goodness-of-fit (GOF) of the BEDD model

Zackenberg network.

Raw statistic:

$$T_s = \frac{N_s - \widehat{\mathbf{E}}N_s}{\sqrt{\widehat{\mathbf{V}}N_s}}$$



Goodness-of-fit (GOF) of the BEDD model

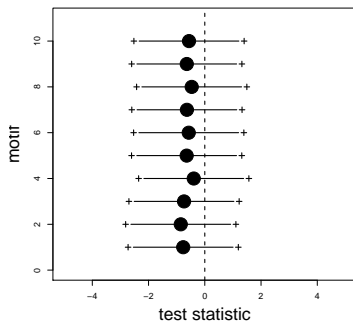
Zackenberg network.

Raw statistic:

$$T_s = \frac{N_s - \hat{\mathbb{E}}N_s}{\sqrt{\hat{\mathbb{V}}N_s}}$$

Corrected stat.: accounts for the estimation error in $\hat{\mathbb{E}}N$

$$T'_s = \frac{N_s - (\hat{\mathbb{E}}N_s - \hat{\mathbb{B}}(\hat{\mathbb{E}}N_s))}{\sqrt{\hat{\mathbb{V}}(N_s - \hat{\mathbb{E}}N_s)}}$$



Testing degree imbalance

Question. Is there some degree imbalance between plants?

Testing degree imbalance

Question. Is there some degree imbalance between plants?

Statistical test.

- ▶ Assume $A \sim BEDD(\rho, g, h)$,

$$H_0 = \{h = 1\}$$

- ▶ For motif s , evaluate $\widehat{\mathbb{E}}_0(N_s)$ and $\widehat{\mathbb{V}}_0(N_s)$ and compare

$$W_s = (N_s - \widehat{\mathbb{E}}_0(N_s)) / \sqrt{\widehat{\mathbb{V}}_0(N_s)}$$

with $\mathcal{N}(0, 1)$

Testing degree imbalance

Question. Is there some degree imbalance between plants?

Statistical test.

- ▶ Assume $A \sim BEDD(\rho, g, h)$,

$$H_0 = \{h = 1\}$$

- ▶ For motif s , evaluate $\widehat{\mathbb{E}}_0(N_s)$ and $\widehat{\mathbb{V}}_0(N_s)$ and compare

$$W_s = (N_s - \widehat{\mathbb{E}}_0(N_s)) / \sqrt{\widehat{\mathbb{V}}_0(N_s)}$$

with $\mathcal{N}(0, 1)$

Example. (only one significant difference)

plant-pollinator					
s	5	6	10	15	16
W_s	$-6.45 \cdot 10^{-2}$	$9.96 \cdot 10^{-1}$	$-6.63 \cdot 10^{-2}$	$7.52 \cdot 10^{-1}$	2.43
seed dispersal					
s	5	6	10	15	16
W_s	$-2.14 \cdot 10^{-1}$	$-2.14 \cdot 10^{-1}$	$-2.93 \cdot 10^{-1}$	$-2.95 \cdot 10^{-1}$	$-3.56 \cdot 10^{-1}$

Comparing network imbalances

Question. Do network A and B share the same imbalance for pollinators?

Comparing network imbalances

Question. Do network A and B share the same imbalance for pollinators?

Statistical test.

- ▶ Assume $A \sim BEDD(\rho^A, g^A, h^A)$ and $B \sim BEDD(\rho^B, g^B, h^B)$

$$H_0 = \{g^A = g^B\}$$

- ▶ For motif s , evaluate $\widehat{\mathbb{E}}_{\widehat{\rho}^A, \widehat{g}^B, \widehat{g}^A}(N_s^A)$ and $\widehat{\mathbb{E}}_{\widehat{\rho}^B, \widehat{g}^A, \widehat{g}^B}(N_s^B)$ and compare

$$W_s = \frac{(N_s^A - \widehat{\mathbb{E}}_0(N_s^A)) - (N_s^B - \widehat{\mathbb{E}}_0(N_s^B))}{\sqrt{\widehat{\mathbb{V}}_0(N_s^A) + \widehat{\mathbb{V}}_0(N_s^B)}}$$

with $\mathcal{N}(0, 1)$

Comparing network imbalances

Question. Do network A and B share the same imbalance for pollinators?

Statistical test.

- Assume $A \sim BEDD(\rho^A, g^A, h^A)$ and $B \sim BEDD(\rho^B, g^B, h^B)$

$$H_0 = \{g^A = g^B\}$$

- For motif s , evaluate $\hat{\mathbb{E}}_{\hat{\rho}^A, \hat{g}^B, \hat{g}^A}(N_s^A)$ and $\hat{\mathbb{E}}_{\hat{\rho}^B, \hat{g}^A, \hat{g}^B}(N_s^B)$ and compare

$$W_s = \frac{(N_s^A - \hat{\mathbb{E}}_0(N_s^A)) - (N_s^B - \hat{\mathbb{E}}_0(N_s^B))}{\sqrt{\hat{\mathbb{V}}_0(N_s^A) + \hat{\mathbb{V}}_0(N_s^B)}}$$

with $\mathcal{N}(0, 1)$

Example. (no significant difference)

s	5	6	10	15	16
F_s^A	$9.21 \cdot 10^{-5}$	$1.00 \cdot 10^{-5}$	$8.12 \cdot 10^{-6}$	$3.32 \cdot 10^{-7}$	$4.47 \cdot 10^{-8}$
$\hat{\mathbb{E}}_0 F_s^A$	$1.96 \cdot 10^{-4}$	$3.75 \cdot 10^{-5}$	$1.74 \cdot 10^{-5}$	$4.25 \cdot 10^{-6}$	$1.33 \cdot 10^{-6}$
F_s^B	$5.13 \cdot 10^{-4}$	$1.15 \cdot 10^{-4}$	$5.07 \cdot 10^{-5}$	$1.79 \cdot 10^{-5}$	$5.96 \cdot 10^{-6}$
$\hat{\mathbb{E}}_0 F_s^B$	$2.66 \cdot 10^{-4}$	$2.92 \cdot 10^{-5}$	$2.85 \cdot 10^{-5}$	$1.50 \cdot 10^{-6}$	$1.69 \cdot 10^{-7}$
W_s	-1.56	-1.56	-0.97	-1.28	-0.96

Outline

Bipartite networks and motifs

A null model

Motif distribution

Goodness-of-fit and network comparison

Network embedding

Network embedding: Multivariate analysis

Analysing multiple networks. Principle

- ▶ 'Embed' each network into a convenient space (e.g. \mathbb{R}^d)
- ▶ Use standard multivariate analysis (clustering, PCA, MDS, ...)

Network embedding: Multivariate analysis

Analysing multiple networks. Principle

- ▶ 'Embed' each network into a convenient space (e.g. \mathbb{R}^d)
- ▶ Use standard multivariate analysis (clustering, PCA, MDS, ...)

Using motifs. K networks

$$(\text{Network})_k \rightarrow (N_1^k, \dots, N_S^k) \in \mathbb{R}^S$$

but need to correct for: network sizes, correlation between motif frequencies, etc...

Network embedding: Multivariate analysis

Analysing multiple networks. Principle

- ▶ 'Embed' each network into a convenient space (e.g. \mathbb{R}^d)
- ▶ Use standard multivariate analysis (clustering, PCA, MDS, ...)

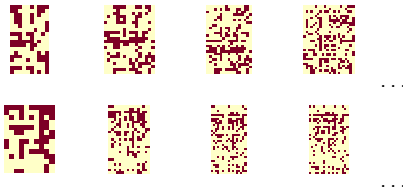
Using motifs. K networks

$$(\text{Network})_k \rightarrow (N_1^k, \dots, N_S^k) \in \mathbb{R}^S$$

but need to correct for: network sizes, correlation between motif frequencies, etc...

Zackenberg dataset. $K = 46$ networks

- ▶ 2 years
- ▶ 1 network observed every few days

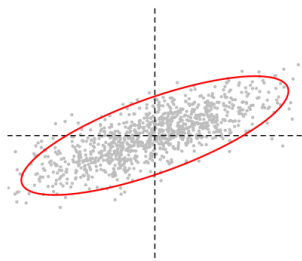


Choleski transform

Aim: 'Remove' correlation and variance heterogeneity

Covariance matrix of (X_1, X_2) :

$$\Sigma_{X_1, X_2} = \begin{bmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{12} & \sigma_2^2 \end{bmatrix}$$



Diagonalization: $\Sigma = P\Lambda P^{-1}$

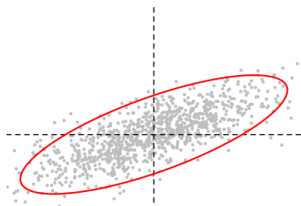
Choleski matrix: $\Sigma^{-1/2} = P\Lambda^{-1/2}P^{-1}$

Choleski transform

Aim: 'Remove' correlation and variance heterogeneity

Covariance matrix of (X_1, X_2) :

$$\Sigma_{X_1, X_2} = \begin{bmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{12} & \sigma_2^2 \end{bmatrix}$$

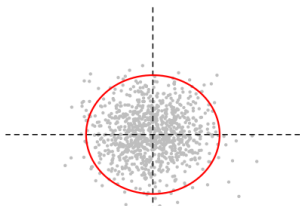


Diagonalization: $\Sigma = P\Lambda P^{-1}$

Choleski matrix: $\Sigma^{-1/2} = P\Lambda^{-1/2}P^{-1}$

Choleski transform:

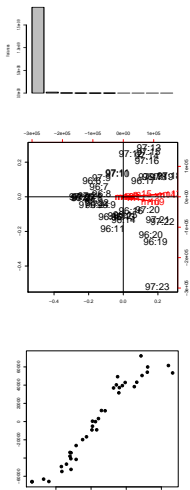
$$\begin{bmatrix} X'_1 \\ X'_2 \end{bmatrix} = \Sigma^{-1/2} \begin{bmatrix} X_1 \\ X_2 \end{bmatrix}$$



$$\Sigma_{X'_1, X'_2} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

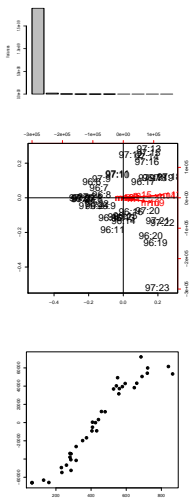
Network embedding: Zackenber's data [SROB16]

Raw counts

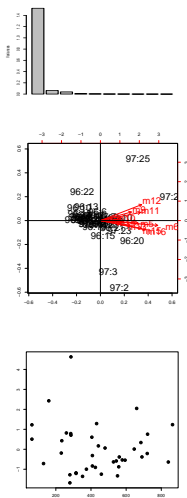


Network embedding: Zuckerberg's data [SROB16]

Raw counts

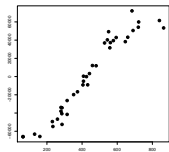
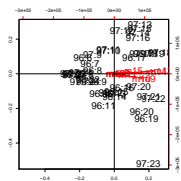


Corrected stat.

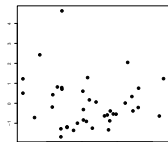
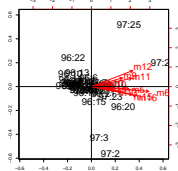
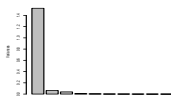


Network embedding: Zuckerberg's data [SROB16]

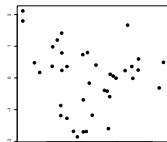
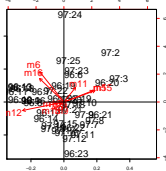
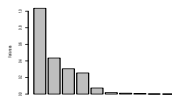
Raw counts



Corrected stat.

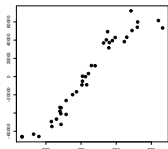
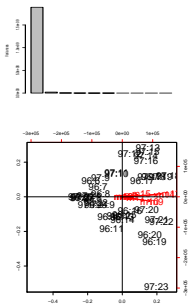


Choleski

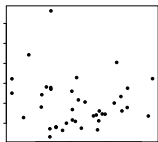
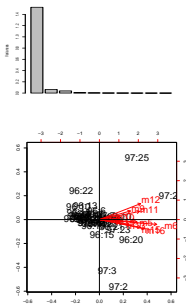


Network embedding: Zuckerberg's data [SROB16]

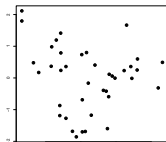
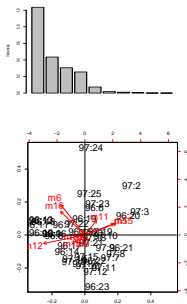
Raw counts



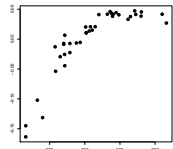
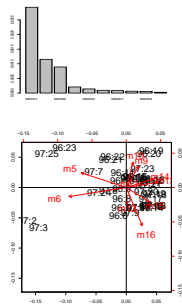
Corrected stat.



Choleski



Bray-Curtis MDS

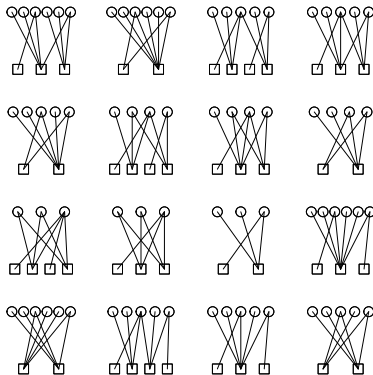


Super-motifs

Motif:



Some super-motifs:



... 396 super-motifs

Variance:

$$\begin{aligned}
 N_s^2 &= \left(\sum_{\alpha} Y_{s\alpha} \right)^2 \\
 &= \sum_{\alpha, \beta: \alpha \cap \beta = \emptyset} Y_{s\alpha} Y_{s\beta} \\
 &\quad + \underbrace{\sum_{\alpha, \beta: \alpha \cap \beta \neq \emptyset} Y_{s\alpha} Y_{s\beta}}_{\text{occurrence of a super-motif}}
 \end{aligned}$$

Covariance: same game, for $Y_{s\alpha} Y_{s'\beta}$ with $s \neq s'$

Asymptotic distribution of the count

Estimated probability.

$$\bar{\phi}_s := \phi_2 \phi_4 / \phi_1 \quad \rightarrow \quad \bar{F}_s := F_2 F_4 / F_1$$

where F_1, F_2, F_4 = observed frequencies of top stars, bottom stars and edges.

Asymptotic distribution of the count

Estimated probability.

$$\bar{\phi}_s := \phi_2 \phi_4 / \phi_1 \quad \rightarrow \quad \bar{F}_s := F_2 F_4 / F_1$$

where F_1, F_2, F_4 = observed frequencies of top stars, bottom stars and edges.

Asymptotic normality for non-star motifs. Under BEDD (and sparsity conditions):

$$(F_s - \bar{F}_s) / \sqrt{\widehat{V}(F_s)} \xrightarrow[m, n \rightarrow \infty]{} \mathcal{N}(0, 1)$$

Proof:

- ▶ decompose

$$F_s - \bar{F}_s = \underbrace{(F_s - \phi_s)}_{\text{random fluctuations}} + \underbrace{(\phi_s - \bar{\phi}_s)}_{\text{null under BEDD}} + \underbrace{(\bar{\phi}_s - \bar{F}_s)}_{\text{estimation error } \rightarrow 0},$$

- ▶ construct a counting martingale [GL17] for $F_s - \phi_s$

Asymptotic distribution of the count

Estimated probability.

$$\bar{\phi}_s := \phi_2 \phi_4 / \phi_1 \quad \rightarrow \quad \bar{F}_s := F_2 F_4 / F_1$$

where F_1, F_2, F_4 = observed frequencies of top stars, bottom stars and edges.

Asymptotic normality for non-star motifs. Under BEDD (and sparsity conditions):

$$(F_s - \bar{F}_s) / \sqrt{\widehat{\mathbb{V}}(F_s)} \xrightarrow[m, n \rightarrow \infty]{} \mathcal{N}(0, 1)$$

Proof:

► decompose

$$F_s - \bar{F}_s = \underbrace{(F_s - \phi_s)}_{\text{random fluctuations}} + \underbrace{(\phi_s - \bar{\phi}_s)}_{\text{null under BEDD}} + \underbrace{(\bar{\phi}_s - \bar{F}_s)}_{\text{estimation error } \rightarrow 0},$$

► construct a counting martingale [GL17] for $F_s - \phi_s$

Test statistic. Under BEDD:

$$N_s \approx \mathcal{N}(\widehat{\mathbb{E}}(N_s), \widehat{\mathbb{V}}(N_s)) \quad \Leftrightarrow \quad (N_s - \widehat{\mathbb{E}}(N_s)) / \sqrt{\widehat{\mathbb{V}}(N_s)} \approx \mathcal{N}(0, 1)$$

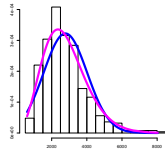
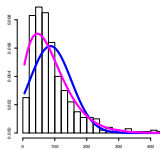
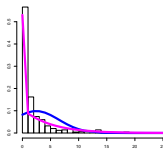
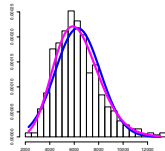
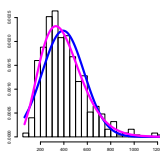
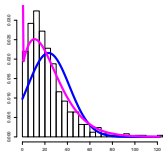
In practice: Asymptotic normality

$(n = 2m/3)$

$m = 50$

$m = 100$

$m = 200$



Normal distribution, Poisson-geometric distribution with same mean and variance [Sta01,PDK⁺08]

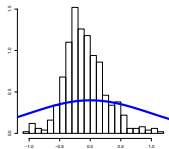
In practice: Test statistic

Need to account for the estimation error of $\widehat{\mathbb{E}N}$

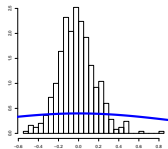
Regular stat.:

$$\frac{N - \widehat{\mathbb{E}N}}{\sqrt{\widehat{\mathbb{V}N}}}$$

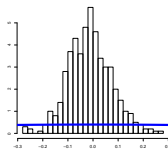
$m = 50$



$m = 100$



$m = 200$



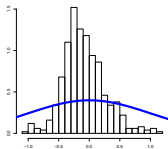
In practice: Test statistic

Need to account for the estimation error of $\hat{\mathbb{E}}N$

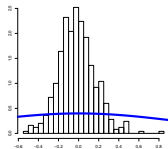
Regular stat.:

$$\frac{N - \hat{\mathbb{E}}N}{\sqrt{\hat{\mathbb{V}}N}}$$

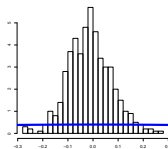
$m = 50$



$m = 100$

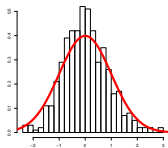
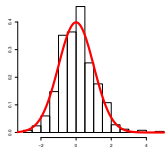
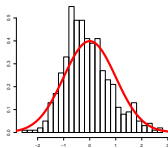


$m = 200$



Correction:

$$\frac{N - (\hat{\mathbb{E}}N - \hat{\mathbb{B}}(\hat{\mathbb{E}}N))}{\sqrt{\hat{\mathbb{V}}(N - \hat{\mathbb{E}}N)}}$$



- Need to evaluate $\mathbb{V}(N - \hat{\mathbb{E}}(N))$ and $\mathbb{B}(\hat{\mathbb{E}}N)$: resort to Taylor expansion (Δ -method)