Analysis of Complex Systems

Lecture 4: Local network organisation: Clusters, motifs, Jaccard index, betweenness centrality

> Marcus Kaiser m.kaiser@ncl.ac.uk

Objectives

- Clusters
- Motifs
- Jaccard index (matching index)
- Edge and node betweenness

Clusters

Clusters

Clusters: nodes within a cluster tend to connect to nodes in the same cluster but are less likely to connect to nodes in other clusters

Quantitative measure: modularity Q (Newman & Girvan, Physical Review E, 2004)

important terms:

hierarchical (cluster, sub-cluster, ...)
overlapping or non-overlapping

(one node can only be member of one cluster)

predefined number of clusters

(e.g. k-means algorithm)

Potential time problem for large networks, O(k^N)

Hundreds of algorithms for cluster detection!





Cluster detection – algorithm 1

Non-hierarchical, overlapping

Genetic algorithm



- Have as few as possible connections *between* them
- Have as few as possible absent connections *within* them



- Random starting configurations
- Evolution:
 - Mutation
 - Evaluation
 - Selection
- Validation

- : Area relocation
- : Cost function
- : Threshold

Hilgetag et al. (2000) Phil. Trans. Roy. Soc. Lond. B.



Hilgetag et al. (2000) Phil Trans R Soc 355: 91

Cluster detection – algorithm 2

hierarchical, non-overlapping

Monte-Carlo approach

Each internal node *r* of the *dendrogram* is associated with a probability p_r that a pair of vertices in the left and right subtrees of that node are connected. (The shades of the internal nodes in the figure represent the probabilities.)





Clauset et al. (2008). Nature 453, 98-101

Cluster detection – algorithm 2

Example result for a food web





Clauset et al. (2008). Nature 453, 98-101

Motifs



Idea: determine building blocks of networks.

Hope: structural building blocks correspond to functional units.

Pattern: possible connection configuration for a k-node subgraph (see list of all 3-node configurations)

Motif: pattern that occurs significantly more often than for rewired benchmark networks (same number of nodes and edges and same degree distribution)



* Milo et al. (2002) Science; http://www.weizmann.ac.il/mcb/UriAlon/groupNetworkMotifSW.html

Motif detection – algorithm



Motif detection – results



Network	Nodes	Edges	N _{real}	$N_{\rm rand} \pm {\rm SD}$	Z score	N _{real}	$N_{\text{rand}} \pm \text{SD}$	Z score	N _{real}	$N_{\rm rand} \pm {\rm SD}$	Z score
Gene regulation (transcription)			$ \begin{vmatrix} \mathbf{x} \\ \mathbf{\psi} \\ \mathbf{y} \\ \mathbf{z} \end{vmatrix} $		Feed- forward loop			Bi-fan			
E. coli S. comprisiant	424	519	40 70	7 ± 3	10	203	47 ± 12 300 ± 40	13			
Neurons		$ \begin{array}{c c} & & \\ & $		Feed- forward loop	X Y Z W		Bi-fan	$ \begin{array}{c} \swarrow^{X} & \aleph \\ $		Bi- parallel	
C. elegans†	252	509	125	90 ± 10	3.7	127	55 ± 13	5.3	227	35 ± 10	20

Milo et al. Science, 2002

Motif detection – problems

Advantages:

- Identify special network patterns which *might* represent functional modules
- Disadvantages:
- Slow for large networks and unfeasible for large (e.g. 5-node) motifs (#patterns: 3-node – 13; 4-node – 199; 5-node: 9364; 6-node - 1,530,843)
- Rewired benchmark networks do not retain *clusters*; most patterns become insignificant for clustered benchmark networks*

Jaccard index (matching index)

Jaccard index

Jaccard index = similarity of incoming and outgoing connections of two nodes

$$J(A,B) = \frac{|A \cap B|}{|A \cup B|}$$

|M| : number of elements in the set M A \cap B: common elements in sets A and B (intersection) A U B: all elements in sets A and B (union)

A=0 1 0 0 0 0 1 0 1 B=0 1 0 1 0 0 1 0 0 $|A \cap B| = 2$ $|A \cup B| = 4$ J(A,B) = 2 / 4 = 0.5

Similarity and compensation

- Use of non-metric multi-dimensional scaling (NMDS)
- Idea: similar connections -> similar function



Centrality measures

Node betweenness

Node betweenness: number of shortest paths that go through one node



Edge betweenness

Edge betweenness: number of shortest paths that go through one edge



High edge betweenness

Low edge betweenness

Centrality measure example



Node 8 has the highest node betweenness Edge 8-9 has the highest edge betweenness

Summary

- What are clusters? Which kinds of clustering algorithms exist?
- What are motifs? Which features of the original network are retained for the benchmark networks and which are lost?
- What is the Jaccard index? What could similar connectivity indicate?
- What are betweenness measures?



1. A cluster algorithm is hierarchical, non-overlapping, and has a predermined number of clusters. Which of these features is the most devastating one for biological network analysis?

2. The Jaccard index for two proteins is very high. Would you expect the nodes to be in the same or in different protein clusters?

3. You have identified the edge with the highest edge betweenness. How would network properties change if that edge were removed? Which measure would be most affected?