Community structure in networks

Santo Fortunato





More links "inside" than "outside"

Graphs are "sparse"

Metabolic







Social

Economical





Outline

- Elements of community detection
- Graph partitioning
- Hierarchical clustering
- The Girvan-Newman algorithm
- New methods
- Testing algorithms
- Conclusions

Null hypothesis

The relations between nodes can be inferred from the topology, i.e.



Questions

- What is a community?
- What is a partition?
- What is a "good" partition?

Communities: definition

- Local criteria
- Global criteria
- Vertex similarity

Local definitions: self-referring

"Look at the subgraph, forget the rest of the network."

Ex. Clique!



n-cliques, k-plexes, etc.

Local definitions: comparative

"Compare links inside and outside of the subgraph"

1) "Strong definition": LS set



2) "Weak definition": internal degree is larger than external degree.

Global definitions

Looking at the community with respect to the whole graph

Null model: graph with no community structure

Ex. Newman-Girvan null model: a random graph with the same degree sequence of the original graph

Definitions based on vertex similarity

Vertices are in the same community if they are "similar" to each other

Similarity measure can be local or global

Ex. Distance between vertices, eigenvector components, structural equivalence, etc.

Warning

Communities are usually implicitly defined by the specific algorithm adopted, without an explicit definition!

The practical definition may depend on the specific system/application

What is a partition?

"A partition is a subdivision of a graph in groups of vertices, such that each vertex is assigned to one group"

Problems:1) Overlapping communities2) Hierarchical structure

Overlapping communities

In real networks, vertices may belong to different modules



G. Palla, I. Derényi, I. Farkas, T. Vicsek, Nature 435, 814, 2005

Hierarchies

Modules may embed smaller modules, yielding different organizational levels



A. Clauset, C. Moore, M.E.J. Newman, LNCS 4503, 1, 2007

What is a "good" partition?





How can we compare different partitions?



Partition \mathcal{P}_1 versus \mathcal{P}_2 : which one is better?

Quality function Q

Is $Q(\mathcal{P}_1) > Q(\mathcal{P}_2)$ or $Q(\mathcal{P}_1) < Q(\mathcal{P}_2)$?

Modularity

$$Q = \frac{1}{L} \sum_{i=1}^{n} \left(l_{i} - \frac{d_{i}^{2}}{4L} \right)$$

 l_i = # links in module i $\frac{d_i^2}{4L}$ = expected # of links in module i



$\frac{d_i}{2L} = \frac{\text{probability that a stub,}}{\text{module i}}$





 $\frac{d_i}{2L} \cdot \frac{d_i}{2L} = \frac{d_i^2}{4L^2} = \frac{d_i^2}{dL^2}$ probability that the link is internal to module i

$$\frac{d_i^2}{4L^2} \cdot L = \frac{d_i^2}{4L} =$$

expected number of links in module i

History

- 1970s: Graph partitioning in computer science
- Hierarchical clustering in social sciences
- 2002: Girvan and Newman, PNAS 99, 7821-7826
- 2002-onward: methods of "new generation", mostly by physicists

Graph partitioning

"Divide a graph in n parts, such that the number of links between them (*cut size*) is minimal"

Problems:

Number of clusters must be specified
 Size of the clusters must be specified

If cluster sizes are not specified, the minimal cut size is zero, for a partition where all nodes stay in a single cluster and the other clusters are "empty"

Bipartition: divide a graph in two clusters of equal size and minimal cut size



Spectral partitioning

Laplacian matrix L

 $L_{ij} = d_i \delta_{ij} - A_{ij}$

Spectral properties of L:

All eigenvalues are non-negative
If the graph is divided in g components, there are g zero eigenvalues
In this case L can be rewritten in a block-diagonal form



If the network is connected, but there are two groups of nodes weakly linked to each other, they can be identified from the eigenvector of the second smallest eigenvalue (*Fiedler vector*)

The Fiedler vector has both positive and negative components, their sum must be 0

If one wants a split into n_1 and $n_2=n-n_1$ nodes, one takes the n_1 largest (smallest) components of the Fiedler vector

Kernighan-Lin algorithm

Start: split in two groups

At each step, a pair of nodes of different groups are swapped so to decrease the cut size

Sometimes swaps are allowed that increase the cut size, to avoid local minima

Hierarchical clustering

Very common in social network analysis

 A criterion is introduced to compare nodes based on their similarity
 A similarity matrix X is constructed: the similarity of nodes i and j is X_{ij}
 Starting from the individual nodes, larger groups are built by joining groups of nodes based on their similarity

Final result: a hierarchy of partitions (dendrogram)



Problems of traditional methods

- Graph partitioning: one needs to specify the number and the size of the clusters
- Hierarchical clustering: many partitions recovered, which one is the best?

One would like a method that can predict the number and the size of the partition and indicate a subset of "good" partitions

Girvan-Newman algorithm

M. Girvan & M.E.J Newman, PNAS 99, 7821-7826 (2002)

Divisive method: one removes the links that connect the clusters, until the latter are isolated

How to identify intercommunity links? Betweenness

Link-betweenness: number of shortest paths crossing a link



Steps

- 1. Calculate the betweenness of all links
- 2. Remove the one with highest betweenness
- **3. Recalculate the betweenness of the remaining edges**
- 4. Repeat from 2

The process delivers a hierarchy of partitions: which one is the best?

The best partition is the one corresponding to the highest modularity Q

M.E.J. Newman & M. Girvan, Phys. Rev. E 69, 026113 (2004)

The algorithm runs in a time $O(n^3)$ on a sparse graph (i.e. when m ~ n)

New methods

- Divisive algorithms
- Modularity optimization
- Spectral methods
- Dynamics methods
- Clique percolation

Divisive algorithms

Based on link removal (like GN)

Ex. Algorithm by Radicchi et al. (PNAS 101, 2658-2663, 2004)

Edge clustering:



Main idea: inter-community links have low edge clustering coefficient



Steps

- 1. Calculate the edge clustering of all links
- 2. Remove the one with lowest edge clustering
- **3. Recalculate the measure for the remaining edges**
- 4. Repeat from 2

Advantage over GN: fast! The CPU time scales as O(n²) on a sparse graph

Modularity optimization

4L

Goal: find the maximum of Q over all possible network partitions

Problem: NP-complete!

Greedy algorithms
 Simulated annealing
 Extremal optimization

Greedy algorithm

M.E.J. Newman, Phys. Rev. E 69, 066133, 2004

- Start: partition with one node in each community
- Merge groups of nodes so to obtain the highest increase of Q
- Continue until all nodes are in the same community
- Pick the partition with largest modularity

CPU time O(n²)

Resolution limit of modularity



S.F. & M. Barthélemy, PNAS 104, 36 (2007)







Spectral methods

Finding communities from spectral properties of graph matrices: A, L, etc. Ex. Algorithm by Donetti & Muñoz (JSTAT, P10012, 2004)

The first few eigenvectors of the Laplacian are computed

Eigenvector components act like coordinates to represent nodes in space



Nodes are then grouped with hierarchical clustering

Dynamic algorithms

- Potts model
- Synchronization
- Random walks

Clique percolation

G. Palla, I. Derényi, I. Farkas, T. Vicsek, Nature 435, 814, 2005



Communities in weighted networks



Most methods based on topology fail

Some methods can be easily adapted

Sparsity may not be crucial!

Weighted modularity

$$Q_w = \frac{1}{W} \sum_{i=1}^n \left(s_i^{in} - \frac{s_i^2}{4W} \right)$$

 S_i^{in} = sum of weights of links in module i (strength of module i)



expected strength of module i

Markov Cluster Algorithm

S. Van Dongen, PhD thesis (2000)

Basic idea: diffusion flow on a network

 $W_{ij} \rightarrow S_{ij}$

S_{ij} is the stochastic matrix:

S_{ij}=W_{ij}/s_i

Three parameters: p, , k

Steps:

- 1. (Diffusion) Raise the stochastic matrix to the power p (e.g. p=2)
- 2. (Inflation) Raise each resulting matrix element to the power α
- 3. Normalize the elements of the resulting matrix (by row)
- 4. Keep only the k largest elements per column
- 5. Repeat from 1.

After a sufficient number of iterations the matrix converges to a matrix with 0s and 1s, with disconnected components!

Problem: the final configuration depends on the parameters p and (mostly!) α

Complexity: O(nk²)

http://www.micans.org/mcl/

Matrix ordering

Goal: to put a matrix in block-diagonal form!



Cost function optimization

$$C = \frac{1}{N} \sum_{i,j=1}^{N} W_{ij} |i - j|$$

Quadratic assignment problem, NP complete!

Standard optimization techniques, e.g. simulated annealing

M. Sales-Pardo, R. Guimerà, A.A. Moreira, L.A.N. Amaral, PNAS 104, 15224 (2007)

K-clustering

Set of data points, distance d(x,y) for each pair of points x,y

Goal: dividing the points in k groups such to maximize/minimize a given measure

Problem: number of clusters given as input!

Minimum k-clustering : minimizing the "diameter" of a cluster, i.e. the largest distance between points of the cluster

k-clustering sum : minimizing the average intracluster distance

k-center : minimizing the maximum distance of cluster points from a "centroid"

k-median : minimizing the average distance of cluster points from a "centroid"

k-means : minimizing the average squared distance of cluster points from an arbitrary "centroid" point

Example: k-means clustering

Points embedded in metric space

- 1. k points are randomly chosen, as far as possible from each other ("centroids")
- 2. Each data point is assigned to the nearest centroid
- 3. Recalculate positions of centroids by determining the centers of mass of the k clusters
- 4. Repeat from 2.

Problem: result sensitive to initial conditions!

Testing algorithm

- Artificial networks
- Real networks with known community structure

Benchmark of Girvan & Newman



Problems

- All nodes have the same degree
- All communities have equal size

In real networks the distributions of degree and community size is highly heterogeneous!

New benchmark (A. Lancichinetti, S. F., F. Radicchi, PRE 78, 046110, 2008)

- Power law distribution of degree
- Power law distribution of community size
- A mixing parameter μ sets the ratio between the external and the total degree of each node







Real networks



Outlook

A long way to go ... more questions than answers from clustering

- Overlapping communities
- Hierarchies
- Testing
- Computational complexity
- Clustering in dense correlation matrices (i.e. neither sparse nor complete)

More rigorous definition of the problem!

Community Structure in Graphs

Santo Fortunato^a, Claudio Castellano^b

 a Complex Networks Lagrange Laboratory (CNLL), ISI Foundation, Torino, Italy

^b SMC, INFM-CNR and Dipartimento di Fisica, "Sapienza" Università di Roma, P. le A. Moro 2, 00185 Roma, Italy

Abstract

Graph vertices are often organized into groups that seem to live fairly independently of the rest of the graph, with which they share but a few edges, whereas the relationships between group members are stronger, as shown by the large number of mutual connections. Such groups of vertices, or communities, can be considered as independent compartments of a graph. Detecting communities is of great importance in sociology, biology and computer science, disciplines where systems are often represented as graphs. The task is very hard, though, both conceptually, due to the ambiguity in the definition of community and in the discrimination of different partitions and practically, because algorithms must find "good" partitions among an exponentially large number of them. Other complications are represented by the possible occurrence of hierarchies, i.e. communities which are nested inside larger communities, and by the existence of overlaps between communities, due to the presence of nodes belonging to more groups. All these aspects are dealt with in some detail and many methods are described, from traditional approaches used in computer science and sociology to recent techniques developed mostly within statistical physics.

1 Introduction

The origin of graph theory dates back to Euler's solution [1] of the puzzle of Königsberg's bridges in 1736. Since then a lot has been learned about graphs and their mathematical properties [2]. In the 20th century they have also become extremely useful as representation of a wide variety of systems in different areas. Biological, social, technological, and information networks can be studied as graphs, and graph analysis has become crucial to understand the features of these systems. For instance, social network analysis started in the 1930's and has become one of the most important topics in sociology [3, 4]. In recent times, the computer revolution has provided scholars with a huge amount of data and computational resources to process and analyse these data. The size of real networks one can potentially handle has also grown considerably, reaching

1

S. F., C. Castellano, arXiv:0712.2716

arXiv:0712.2716v1 [physics.soc-ph] 17 Dec 2007