Community Detection: Overlapping Communities

CS224W: Analysis of Networks Jure Leskovec, Stanford University http://cs224w.stanford.edu



Non-overlapping Communities



Overlapping Communities

Non-overlapping vs. overlapping communities



Overlaps of Social Circles

A node can belong to many social "circles"



What if communities overlap?





[Palla et al., '05] Clique Percolation Method (CPM)

- Two nodes belong to the same community if they can be connected through adjacent k-cliques:
 - k-clique:
 - Fully connected graph on k nodes
 - Adjacent k-cliques:
 - overlap in k-1 nodes
- k-clique community
 - Set of nodes that can be reached through a sequence of adjacent *k*-cliques



Two overlapping 3-clique communities

Clique Percolation Method (CPM)

Two nodes belong to the same community if they can be connected through adjacent kcliques:



4-clique



Adjacent 4-cliques



Non-adjacent 4-cliques



Communities for k=4

CPM: Steps

Clique Percolation Method:

Find maximal-cliques

 Def: Clique is maximal if no superset is a clique

Clique overlap super-graph:

- Each clique is a super-node
- Connect two cliques if they overlap in at least *k-1* nodes

Communities:

 Connected components of the clique overlap matrix

How to set k?



Set k so that we get the "richest" (most widely distributed cluster sizes) community structure

CPM method: Example

- Start with graph
- Find maximal cliques
- Create clique overlap matrix A
 - Rows/Cols are maxcliques, entry is number of nodes in common
- Threshold the matrix at value *k-1*
 - If $a_{ij} < k 1$ set 0
- Communities are the connected components of the thresholded matrix



[Palla et al., '07] Example: Phone-Call Network

50300 50007 50010, 50008 50003 🔿 50018 50015 50007 50007 Ó 50007 5001 \050013 **Ø**50006 **Ø**50006

Communities in a "tiny" part of a phone call network of 4 million users [Palla et al., '07]

Example: Website



- No nice way, hard combinatorial problem
- Maximal clique: Clique that can't be extended
 - {a, b, c} is a clique but not maximal clique
 - {a, b, c, d} is maximal clique
- Algorithm: Sketch
 - Start with a seed node
 - Expand the clique around the seed
 - Once the clique cannot be further expanded we found the maximal clique

Note:

This method will generate the same clique multiple times

- Start with a seed vertex a
- Goal: Find the max clique Q that a belongs to
 - Observation:
 - If some x belongs to Q then it is a neighbor of a
 - Why? If $a, x \in Q$ but edge (a, x) does not exist, Q is not a clique!

Recursive algorithm:

- Q ... current clique
- *R* ... candidate vertices to expand the clique to
 Example: Start with *a* and expand around it

Steps of the recursive algorithm

Q=

R=

Jure Leskovec, Stanford CS224W: Analysis of Networks, http://cs224w.stanford.edu

- Start with a seed vertex a
- Goal: Find the max clique Q that a belongs to
 - Observation:
 - If some x belongs to Q then it is a neighbor of a
 - Why? If $a, x \in Q$ but edge (a, x) does not exist, Q is not a clique!

Recursive algorithm:

• Q ... current clique

R ... candidate vertices to expand the clique to **Example:** Start with *a* and expand around it



Steps of the recursive algorithm





$\Gamma(u)$...neighbor set of *u*

- Q ... current clique
- R ... candidate vertices
- Expand(R,Q)
 - while $R \neq \{\}$
 - p = vertex in R
 - $Q_p = Q \cup \{p\}$
 - $R_p = R \cap \Gamma(p)$
 - if R_p ≠ {}: Expand(R_p,Q_p)
 else: output Q_p
 - $R = R \{p\}$



- **Q** ... current clique
- R ... candidate vertices
- Expand(R,Q)
 - while $R \neq \{\}$
 - p = vertex in R
 - $Q_p = Q \cup \{p\}$
 - $R_p = R \cap \Gamma(p)$
 - if R_p ≠ {}: Expand(R_p,Q_p)
 else: output Q_p

 $R = R - \{p\}$

Start: Expand(V, {}) R={a,...f}, Q={} $p = \{b\}$ $Q_{p} = \{b\}$ $R_{p} = \{a, c, d\}$ Expand(R_p , Q): $R = \{a,c,d\}, Q = \{b\}$ p = {a} $Q_{p} = \{b,a\}$ $R_{p} = \{d\}$ Expand(R_p , Q): $R = \{d\}, Q = \{b, a\}$ $p = \{d\}$ $Q_p = \{b,a,d\}$ $R_p = \{\}$: output {b,a,d} $p = \{c\}$ $Q_{p} = \{b,c\}$ $R_{p} = \{d\}$ Expand(R_{p} , Q): $R = \{d\}, Q = \{b, c\}$ $p = \{d\}$

 $Q_{p} = \{b, c, d\}$

R_p = {} : **output {b,c,d**}

- How to prevent maximal cliques from being generated multiple times?
 - Only output cliques that are lexicographically minimum
 - $\{a, b, c\} < \{b, a, c\}$
 - Even better: Only expand to the nodes higher in the lexicographical order



How to Model Networks with Communities?

Network and Communities

- How should we think about large scale organization of clusters in networks?
 - Finding: Community Structure



Network and Communities

- How should we think about large scale organization of clusters in networks?
 - Finding: Core-periphery structure



Nested Core-Periphery

Network and Communities

How do we reconcile these two views? (and still do community detection)



Community structure

Core-periphery

Community Score

- How community-like is a set of nodes?
- A good cluster S has
 - Many edges internally
 - Few edges pointing outside
- What's a good metric:
 Conductance

$$\phi(S) = \frac{|\{(i,j) \in E; i \in S, j \notin S\}|}{\sum_{s \in S} d_s}$$

Small conductance corresponds to good clusters **Note:** We are assuming |S| < |V|/2, d_s degree of node s

Jure Leskovec, Stanford CS224W: Analysis of Networks, http://cs224w.stanford.edu

S

Network Community Profile Plot

(Note |S| < |V|/2)

Define:

Network community profile (NCP) plot

Plot the score of **best** community of size *k*

$$\Phi(k) = \min_{S \subset V, |S|=k} \phi(S)$$





How to (Really) Compute NCP?



Jure Leskovec, Stanford CS224W: Analysis of Networks, http://cs224w.stanford.edu

NCP Plot: Meshes

Meshes, grids, dense random graphs:



NCP plot: Network Science

Collaborations between scientists in networks [Newman, 2005]



Jure Leskovec, Stanford CS224W: Analysis of Networks, http://cs224w.stanford.edu

Natural Hypothesis

Natural hypothesis about NCP:

- NCP of real networks slopes downward
- Slope of the NCP corresponds to the "dimensionality" of the network

What about large networks?



• Social nets	Nodes	Edges	Description
LiveJournal Epinions CA-DBLP	4,843,953 75,877 317,080	$42,845,684 \\ 405,739 \\ 1,049,866$	Blog friendships [5] Trust network [28] Co-authorship [5]
• Information (citation) networks			
Cit-hep-th AmazonProd	$27,400 \\ 524,371$	$352,021 \\ 1,491,793$	Arxiv hep-th [14] Amazon products [8]
• Web graphs			
Web-google Web-wt10g	$855,802 \\ 1,458,316$	$\substack{4,291,352\\6,225,033}$	Google web graph TREC WT10G
• Bipartite affiliation (authors-to-papers) networks			
Atp-DBLP AtM-Imdb	$\begin{array}{c} 615,\!678 \\ 2,\!076,\!978 \end{array}$	$944,456 \\ 5,847,693$	DBLP [21] Actors-to-movies
• Internet networks			
AsSkitter	1.719.037	12.814.089	Autonom. svs.

[Internet Mathematics '09]

Large Networks: Very Different

Typical example: General Relativity collaborations (n=4,158, m=13,422)



[Internet Mathematics '09]

More NCP Plots of Networks



NCP: LiveJournal (n=5m, m=42m)



Jure Leskovec, Stanford CS224W: Analysis of Networks, http://cs224w.stanford.edu

Explanation: The Upward Part

As clusters grow the number of edges inside grows slower that the number crossing



Jure Leskovec, Stanford CS224W: Analysis of Networks, http://cs224w.stanford.edu

Explanation: Downward Part

Empirically we note that best clusters

 (corresponding to green nodes "whiskers") are
 barely connected to the network





NCP plot

\Rightarrow Core-periphery structure

What If We Remove Good Clusters?



Jure Leskovec, Stanford CS224W: Analysis of Networks, http://cs224w.stanford.edu

Suggested Network Structure



Jure Leskovec, Stanford CS224W: Analysis of Networks, http://cs224w.stanford.edu

Part 2: Explanation



How do we reconcile these two views?

Ground-truth Communities

- Basic question: nodes *u*, *v* share *k* communities
- What's the edge probability?
 - Look at networks with ground-truth communities



Communities as Tiles!

Edge density in the overlaps is higher!



"The more different foci (communities) that two individuals share, the more likely it is that they will be tied" - S. Feld, 1981

Communities as "tiles"

Communities as Tiles/Circles



Communities as overlapping tiles

Communities in Networks

What does this mean?



From Networks to Communities



Community-Affiliation Graph Model (AGM)



- Generative model: How is a network generated from community affiliations?
 Model parameters:
 - Nedee V Communities C Manch
 - Nodes V, Communities C, Memberships M
 - Each community c has a single probability p_c

AGM: Generative Process



Given parameters (V, C, M, {p_c})

- Nodes in community c connect to each other by flipping a coin with probability p_c
- Nodes that belong to multiple communities have multiple coin flips: Dense community overlaps

If they "miss" the first time, they get another chance through the next community"

 $p(u, v) = 1 - \prod_{c \in M_u \cap M_v} (1 - p_c)$

Note: If two nodes u and v have no communities in common, then p(u,v)=0. We resolve this by having a "background" community that every node is a member of.

AGM: Dense Overlaps



[icdm '12]

AGM: Flexibility

 AGM can express a variety of community structures: Non-overlapping, Overlapping, Nested

11/30/17



Jure Leskovec, Stanford CS224W: Analysis of Networks, http://cs224w.stanford.edu

Detecting Communities

Detecting communities with AGM:



Given a Graph, find the Model

Affiliation graph *M* Number of communities C Parameters *p_c*

[wsdm '13]

"Relaxing" AGM

"Relax" the AGM: Memberships have strengths



• F_{uA} : The membership strength of node uto community A ($F_{uA} = 0$: no membership)

BigCLAM Model

 Prob. of nodes linking is proportional to the strengths of shared memberships: P(u, v) = 1 - exp(-F_u · F_v^T)

 Now, given a network, we estimate F

$$l(F) = \sum_{(u,v)\in E} \log(1 - \exp(-F_u F_v^T)) - \sum_{(u,v)\notin E} F_u F_v^T$$

Non-negative matrix factorization:

- Update F_{uC} for node u while fixing the memberships of all other nodes
- Updating takes linear time in the degree of $oldsymbol{u}$

[wsdm '13]

BigCLAM Model

Apply block coordinate gradient ascent

$$\nabla l(F_u) = \sum_{v \in \mathcal{N}(u)} F_v \frac{\exp(-F_u F_v^T)}{1 - \exp(-F_u F_v^T)} - \sum_{v \notin \mathcal{N}(u)} F_v$$

- Step size: backtracking line search
- Project F_u back to a non-negative vector

Pure gradient ascent is slow! However:

$$\sum_{v \notin \mathcal{N}(u)} F_v = \left(\sum_v F_v - F_u - \sum_{v \in \mathcal{N}(u)} F_v\right)$$

By caching F_v the gradient step takes linear time in the degree of u

Experimental Setup



- How well do inferred communities correspond to ground-truth?
 - **F1** score, Ω-index, Mutual Information

 We can rank algorithms based on their ability to detect ground-truth

Experiments: Ground-truth



- BigClam improves:
 - 79% over Link clustering
 - 48% over CPM
 - 15% over MMSB

(while being orders of magnitude faster)

Methods

- L Link Clustering
- C Clique Percolation
- M Mixed-Membership Stochastic Block Model
- A AGM
- Measures
- Number of Communities
- Normalized Mutual Information
- F1-score

Experiments: PPI Nets



Protein-Protein interaction networks:
 Gene Ontology based quality of detected communities

Communities: Issues and Questions

Communities: Issues and Questions

Some issues with community detection:

- Many different formalizations of clustering objective functions
- Objectives are NP-hard to optimize exactly
- Methods can find clusters that are systematically "biased"

Questions:

- How well do algorithms optimize objectives?
- What clusters do different methods find?

Many Different Objective Functions

Single-criterion:

- Modularity: *m*-*E(m)*
- Edges cut: cMulti-criterion:
 - Conductance: c/(2m+c)
 - Expansion: c/n
 - Density: 1-m/n²
 - CutRatio: c/n(N-n)
 - Normalized Cut: c/(2m+c) + c/2(M-m)+c
 - Flake-ODF: frac. of nodes with more than ½ edges pointing outside S



n: nodes in Sm: edges in Sc: edges pointing outside S

Many Classes of Algorithms

Many algorithms that implicitly or explicitly optimize objectives and extract communities: Heuristics:

- Girvan-Newman, Modularity optimization: popular heuristics
- Metis: multi-resolution heuristic [Karypis-Kumar '98]

Theoretical approximation algorithms:

Spectral partitioning

NCP: Live Journal



Jure Leskovec, Stanford CS224W: Analysis of Networks, http://cs224w.stanford.edu

Properties of Clusters (1)

500 node communities from Spectral:





500 node communities from Metis:





[WWW `09]

Properties of Clusters (2)



- Metis gives sets with better conductance
- Spectral gives tighter and more well-rounded sets





[WWW `09]

Multi-criterion Objectives



Single-criterion Objectives



Observations:

- All measures are monotonic
- Modularity
 - prefers large clusters
 - Ignores small clusters

Edges cut