

COMMUNITY DETECTION IN LARGE NETWORKS

A Dissertation

Presented to the Faculty of the Graduate School

of Cornell University

in Partial Fulfillment of the Requirements for the Degree of

Doctor of Philosophy

by

June Andrews

May 2012

© 2012 June Andrews
ALL RIGHTS RESERVED

COMMUNITY DETECTION IN LARGE NETWORKS

June Andrews, Ph.D.

Cornell University 2012

Graphs are used to represent various large and complex networks in scientific applications. In order to understand the structure of these graphs, it is useful to treat a set of nodes with similar characteristics as one community and analyze the community's behavior as a whole. Finding all such communities within the graph is the object of community detection. In our research, we compare dozens of existing community detection methods and develop a new class of algorithms for finding communities.

BIOGRAPHICAL SKETCH

June Andrews was born in San Diego, 1985. She attended University of California, Berkeley where she earned her Bachelor's degree in Electrical Engineering and Computer Science, with a minor in Applied Mathematics. She is now completing her doctoral degree in Applied Mathematics at Cornell University.



Figure 1: *Phil Andrews 1955 - 2011*

Here's to you Da.

ACKNOWLEDGEMENTS

It goes without saying, these people have been inspiring forces of nature to work with:

- Mr. Len Kulbacki
- Coach Wilson
- Dr. James Sethian
- Patricia Kovatch
- Dr. John Hopcroft
- Dr. Steve Strogatz
- Dr. Jon Kleinberg

TABLE OF CONTENTS

LIST OF TABLES

LIST OF FIGURES

CHAPTER 1

INTRODUCTION

1.1 Community Detection

Consider an application that studies objects and the interactions between those objects. The application could study anything from people and their friendships, to papers and their citations; a variety of applications fall into this format. If we let nodes represent the objects and edges represent the interactions between those objects, we can store the application's data in a graph. While it can be possible for the application to draw conclusions by looking at every node within the graph, if the graph is large and complex, analyzing every node can be unmanageable and can produce incomprehensible results. We simplify the graph by finding communities of nodes. In particular, we want communities, whose members interact with each other in a particular way and interact with nonmembers of the community in a different way. If such a community is

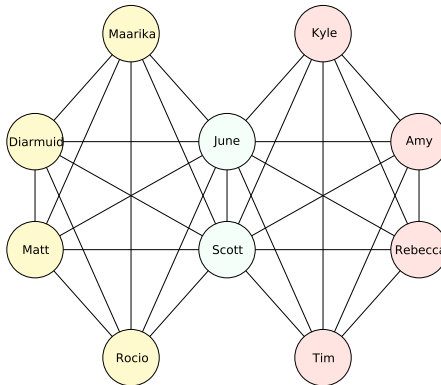


Figure 1.1: A simple graph of people and their friendships. The graph is regular enough to reveal two communities.

found, then two questions arise. How are members of the community related? How does the community interact with the rest of the graph? Given answers to these questions, we can comprehend what is happening in the graph at a local level. Within social networks, we know that communities exist [?], but due to the large, complex nature of social networks, communities can be hard to find. In order to find communities, we must develop the ability to see the forest through the trees. We have to be able to extract the communities of nodes from the interactions of the graph. This is the object of community detection.

Given a graph, there are two prominent questions community detection seeks to answer. First, what is a community and second, what are the communities? Several approaches have been developed to answer these two questions, some with a particular application as motivation. We outline the coupling of a few sciences and one of their preferred detection methods in Table ??.

Prior to 2002, most development of community detection was done within specific fields of applications. Since then, computer scientists have contributed a large volume of advances toward answering these two questions for applications in general. The first goal of this thesis and tie together a portion of these advances into a cohesive understanding of community detection. The second is to use this perspective to create fast and parallel algorithms.

Application	Community Detection Method
Parallel Computation Distribution	k -means clustering [?]
Physics	Belief Propagation [?]
Search Queries	[?]
Sociology	[?]
Storage of Large Matrices	Spectral Analysis [?]
Taxonomy	Neighbor Joining [?]

Table 1.1: A subset of applications and one of their preferred community detection methods.

1.2 Graph Partitioning Methods

For many applications the object is to partition the graph into disjoint components. We call each component a community. There are an exponential number of possible partitions, but not every partition will provide useful information. While much analysis of useful or unuseful information must be left up to the application, there are two characteristics that most applications want in communities. The first is that nodes within a community be well connected. The second is that the community is not well connected to the rest of the graph. The definition of well connected is different for each community detection method.

There are two genres of finding good partitions of the graph, top down approaches that recursively cut the graph and bottom up approaches that united existing partitions.

1.2.1 Top Down Approaches

Top down approaches work by recursively dividing the graph, see Algorithm ?? for their structure. For methods in this category there are two necessary components. The first is the ability to tell if a set of nodes C is a community. The second, if a set of nodes is not a community, then the algorithm finds a way to divide the nodes, without splitting up any communities.

Algorithm 1: RECURSIVE PARTITIONING

Input: $G = (V, E)$

if V is a community **then**

return V

else

divide V into C and $V - C$

return {RECURSIVE PARTITIONING(C), RECURSIVE PARTITIONING($V - C$)}

end if

Conductance

Conductance is a measure of a cut within the graph developed by JTODO [?]. For a given cut, if conductance is low, then there are relatively few edges crossing the cut. Intuitively, this implies that the cut does not divide a community. If further divisions do not improve conductance, then we have found a community.

$$\text{CONDUCTANCE}(C) = \frac{\sum_{u \in C, v \notin C} w(u, v)}{\sum_{u \in C, v \notin C} w(u, v) + \sum_{u, v \in C} w(u, v)} \quad (1.1)$$

While this algorithm is not in heavy use, conductance is used as a measure of whether or not other algorithms that cut the graph have split a community.

Betweenness and Centrality Measures

Betweenness and centrality measures were first presented by Givan and Newman [?]. The intuition is, if an edge lies between two communities, then several

shortest paths between nodes of the two communities will traverse the edge. We remove these edges to divide a network into components. When there are no preferential edges for shortest paths within a component, there are no more edges between communities, and the component is a community.

1.2.2 Bottom Up Approaches

Bottom up approaches work by uniting together subsets of nodes until the subset becomes a community. See Algorithm ?? for their structure. For algorithms in this category there are two components. The first is the determination of which subsets to union. The second is the determination of when a set of nodes is a community. To accomplish these, most bottom up approaches use a metric over the set of subsets. If no two subsets can be united to increase the metric then, every subset is a community.

Algorithm 2: RECURSIVE UNIONING

```

Input:  $S = \{C_1, C_2, \dots\}$ 
  if There exists  $C_i$  and  $C_j$ , such that  $C_i \cup C_j$  is a community then
    return RECURSIVE UNIONING( $\{S - C_i - C_j\} \cup \{C_i \cup C_j\}$ )
  else
    return  $S$ 
  end if

```

Modularity

The overwhelmingly popular metric in this category is modularity. Modularity was first presented by Newman [?]. The metric measures the distance between a provided set of communities and a randomly generated set of communities. Maximizing modularity finds the least random set of communities.

Fast algorithms have been developed for maximizing modularity. In this paper we use the Louvain Algorithm developed by Blondel et. al [?]. The same fast algorithm can be used for the similar metric, modularity ratio [?].

1.3 Overlapping Community Detection

We call two communities overlapping, if there exists a node that is a member of both communities. In practice, these communities are common. For example, think of the community of your colleagues and the community of your family. You are a member of both communities, and while they are different communities, they are overlapping. In fact, for most social networks, we expect there to be many overlapping communities.

1.3.1 Alpha Beta Clustering

In previous sections, communities were the partitions of a graph. Each node was placed in exactly one community. So if it was optimal to place node, n , in community C_1 , then node n would not be placed in community C_2 . Alpha beta clustering makes a change to this step. If adding node n to community C_2 has

a high value, alpha beta clustering adds node n to community C_2 , as well as C_1 . This simple change dramatically restructures community detection. The new structure is a two part process:

1. Create a definition of a community that does not depend on other communities in the graph.
2. Find each community separately.

We now present Mishra's et al [?] approach following these guidelines. Let us say the strength of a connection between a node and a community is the number of edges the node has to members of the community, denoted as $|E(n, C)|$. See Table ?? for a list of all notations. Mishra et al [?] use this notion of strength to define a community satisfying the first guideline. In particular, no node outside of the community is more strongly connect to the community than any of the nodes inside the community. Here is the formal definition of an (α, β) community.

Definition 1 ((α, β)– Community) For community C , let:

$$\alpha(C) = \min_{n \in C} |E(n, C)|$$

$$\beta(C) = \max_{n \notin C} |E(n, C)|$$

If $\alpha(C) > \beta(C)$, then C is an (α, β) community.

Given this definition, Mishra et al[?] are able to find communities quickly and in parallel. In our development of a parallel algorithm we use the same guidelines.

1.4 More Approaches

So far, we have introduced the community detection methods that have provided inspiration for this thesis. There are countless more methods. We briefly outline the most prominent of those methods.

- Kernighan-Lin Algorithm
- k -Clique Percolation
- Belief Propagation
- Hierarchy methods
- Principle Component Analysis

1.5 Desired Improvements

In the field of community detection both algorithms and data sets are increasing in complexity. Hence, a useful theoretical result is the ability to compare and understand complex algorithms. Additionally, a useful experimental result is the ability to compute overlapping communities in parallel on large networks.

We deliver these results:

- A framework for comparing existing community detection methods
- A community definition encouraging overlapping communities
- A parallel algorithm with near perfect scalability to analyze large networks

1.6 Notation

We use the same notation throughout the thesis. A brief description of variables is listed in Table ??.

The assumptions we make are:

- *Self-Loops* We presume there are no self loops in the networks. As a node will always be in the same community as itself, self-loops provide redundant information. Accordingly, $w(u, u) = 0$, for all $u \in V$. We note that this assumption is not held in some of the literature we reference.
- *Edges* We presume that all edges exist and are weighted between 0 and 1. The edge weight function is $w : V \times V \rightarrow \mathbb{R}_{[0,1]}$ Unweighted graphs can easily be adapted into this notation.

We also introduce internal and external edges.

Definition 2 (Internal Edges) *Internal edges are edges between members of the same community C .*

Definition 3 (External Edges) *External edges are edges between a member of community C and a nonmember of C .*

Table 1.2: Notation

Variable Name	Description	Constraints
V	Set of all nodes within the network	$\{u u \in \text{the network}\}$
u and v	Nodes	$u, v \in V$
$w(u, v)$	Edge Weight Function	$w : V \times V \rightarrow \mathbb{R}_{[0,1]}$
G	Network or Graph	$G(V, E)$
C	Community	$C \subset V$
k	Fraction of nodes within C	$k = \frac{ C }{ V }$
$ C $	Size of C	$ C = k V $
S	Set of Communities	$S = \{C_1, C_2, \dots, C_n\}$

Table 1.3: Introduced Functions

Function	Description
$I(C)$	Internal Density of a single Community, C , Definition ??
$E(C)$	External Density of a single Community, C , Definition ??
$I(S)$	Internal Density of a set of Communities, S , Definition ??
$E(S)$	External Density of a set of Communities, S , Definition ??
CONCISENESS(S)	Conciseness of a set of Communities, S , Definition ??

CHAPTER 2

A FRAMEWORK FOR COMPARING METRIC BASED DETECTION

Given the variety of community detection methods, we would like to know the differences and similarities between each. Experimental comparisons from Lancichinetti and Fortunato [?] found that an algorithm's performance depends on the network it is provided. Also from experiments, Leskovec et. al [?] found large communities optimizing metrics diverge from our understanding of a strong community. The process for developing experimental results is to begin with a set of metrics, a set of algorithms to optimize each metric, and a network. Communities found by the algorithms are then compared via their characteristics. Characteristics include: diameter, average path length, degree distribution, size, internal density, etc.

Our approach is the reverse of previous comparisons. We begin with considering the characteristics of a community. The possible values of these characteristics create a multidimensional space. Metrics collapse this multidimensional space onto the real numbers. We can then categorize the multidimensional space according to how the metric evaluates communities in the space. With this method, we can get an understanding of metrics independent of particular networks. We follow through with experiments on four networks to confirm our findings.

2.1 Community Characteristics

A community can be described by a variety of characteristics. Each characteristic provides a dimension in the multidimensional space for describing communities. We now outline the more commonly used characteristics of a community.

- INTERNAL DENSITY is density of edges within the community.
- EXTERNAL DENSITY is the density of edges leaving the community.
- SIZE is the number of nodes within the community.
- DIAMETER is the maximum, shortest path between all pairs of members of the community, using only edges within the community.
- AVERAGE SHORTEST PATH is the average shortest path between any two members of the community.
- OUT DEGREE FRACTION is the fraction of a node's edges leaving the community. Characteristics of the community are then the maximum, minimum, and average of the out degree fraction of all nodes in the community.
- DEGREE DISTRIBUTION is the distribution of the degrees of nodes within the community.

There are more characteristics, but we find that a community can be well described by the above list. The listed characteristics are not independent. A high internal density indicates a small diameter and short average path length. A low external density limits the average out degree fraction. In fact, the value of most characteristics can be bounded by internal and external density values. The size

of a community can not be. Hence, the characteristics of internal density, external density, and community size capture a large amount of information about a community's set of characteristics.

2.2 Previous Comparisons

All known comparisons have been experimental. The experiments are run by first selecting a set of networks. Then, each detection method finds the communities within each network. Finally, using a set of metrics the communities found by each method are compared. Lancichinetti and Fortunato [?] compared three popular partition algorithms with generated graphs and used normalized mutual information as the comparison metric. Their results conclude that partition algorithms are fast and work well for non-overlapping communities. Leskovec et. al [?] conducted a broader study. They used eight classes of algorithms over 40 networks and compared the results with a series of metrics covered in this chapter.

2.3 Individual Community Based Metrics

Here, we explore metrics that evaluate the strength of a single community. There are three uses of such metrics. The first is for use in a Top Down (Section ??) or Bottom Up (Section ??) style algorithm to find a partitioning of the network. The second is to find a single community within the network. The third is to compare communities found by complex detection techniques. The later use is more common for these metrics, of which conductance is the most popu-

lar and used in [?]. From the multidimensional space describing communities, we will use the subspace of internal density, external density, and community size. We analyze how single community metrics evaluate this three dimensional space.

2.3.1 Internal and External Density

We now provide formal definitions of internal and external density. We briefly restate that internal edges are edges between members of the same community and external edges are edges between members of different communities.

Definition 4 (Internal Density) *Internal density is the total weight of internal edges, compared to the total possible weight. Hence, $I(C) : C \rightarrow \mathbb{R}_{[0,1]}$, where*

$$I(C) = \frac{\sum_{u \in C} \sum_{v \in C} w(u, v)}{|C|(|C| - 1)}. \quad (2.1)$$

For a community C that has no edges between its members, the *internal density* will be minimized with, $I(C) = 0$. For a community C that is a clique, *internal density* will be maximized with $I(C) = 1$. The closer a community, C is to an *internal density* value of 1, the closer it is to being a clique.

Definition 5 (External Density) *External density is the total weight of external edges, compared to the total possible edge weight that could exist leaving the community:*

$$E(C) = \frac{\sum_{u \in C} \sum_{v \notin C} w(u, v)}{|C|(|V| - |C|)}. \quad (2.2)$$

$(I(C), E(C))$	Weak Community Characteristic
$(0, \frac{1}{2})$	Infinite Diameter
$(\frac{1}{2}, 1)$	Large Average Out Degree Fraction

Table 2.1: Examples of internal and external density values and why they represent poor communities.

For a community C that has all possible *external edges*, external density will be maximized at $E(C) = 1$. For a community C , disconnected from the rest of the graph, external density will be minimized at $E(C) = 0$.

There are other representations of $I(C)$ and $E(C)$ that vary how the $|C|$ and $|V|$ terms are used. The analysis and conclusions that follow are not sensitive to such variations.

With our parameterization, all communities can be mapped to a point $(I(C), E(C))$ in the square $\mathbb{R}_{[0,1]} \times \mathbb{R}_{[0,1]}$. Communities with certain values do not correspond to our understanding of a strong community. Such values are listed in Table ???. However, a community mapped to $(\frac{1}{2}, 0)$ has a short average path length, minimal average out degree fraction, and small diameter. This corresponds to a strong community. The higher the internal density and the lower the external density the stronger a community must be in all characteristics. We define the strongest possible community to be ideal.

Definition 6 (Ideal Single Community) *A community, C , is ideal if it is an isolated clique, specifically having the following properties:*

$$I(C) = 1$$

$$E(C) = 0.$$

2.3.2 Study of Relevant Metrics

Given that we can map a community, C , to the point $(I(C), E(C))$, we now analyze how different metric based detection methods operate in the I, E plane. We cover six metrics that evaluate a single community. We use one approximation to simplify the equations, $|C| \approx |C| - 1$. This approximation has a larger impact on smaller communities, but most communities of interest are large enough to allow the approximation. Additionally, we introduce a variable k representing the portion of the nodes within community C such that $|C| = k|V|$

- CONDUCTANCE is the probability that a step in a random walk will leave the community [?].

$$\text{CONDUCTANCE}(C) = \frac{(1 - k)E(C)}{kI(C) + (1 - k)E(C)} \quad (2.3)$$

- CUT RATIO is the fraction of existing edges to possible edges leaving the community [?].

$$\text{CUT RATIO} = E(C) \quad (2.4)$$

- EDGES CUT is the number of edges connecting the community to the rest of the graph [?].

$$\text{EDGES CUT} = k(1 - k)|V|^2 E(C) \quad (2.5)$$

- EXPANSION the average number of edges leaving the community per node [?].

$$\text{EXPANSION} = (1 - k)|V|E(C) \quad (2.6)$$

- INTERNAL DENSITY as a metric, previously existed before our definition of $I(C)$, [?]. However, we stick to our definition of $I(C)$ for intuitive reasoning and note that in previous work internal density represents the mirror

image of our definition.

$$\text{INTERNAL DENSITY} = 1 - I(C) \quad (2.7)$$

- VOLUME is the total degree of nodes within the community [?].

$$\text{VOLUME} = |C|^2 I(C) + k(1 - k)|V|^2 E(C) \quad (2.8)$$

With this parameterization of the metrics, we can already draw some conclusions. All metrics mentioned, besides VOLUME and CONDUCTANCE are a function of either $I(C)$ or $E(C)$, but not both. A metric that considers only $I(C)$ will be optimized by any clique. This is a very restrictive definition of a community and finding all communities in the graph under such a definition is equivalent to finding all the cliques in a graph, a NP-hard problem. A metric that considers only $E(C)$ will be optimized by any disconnected component of the graph, including a community that includes the entire graph. While it is possible to find all disconnected components in linear time, it also provides no useful information about most datasets.

For the metrics that can be parameterized in terms of $I(C)$ and $E(C)$, all mentioned metrics, we can use level sets. Level sets are a way to visually categorize a space. Let us pick the metric conductance. An optimal value of conductance is 0. We can find all points of $(I(C), E(C))$ (without knowing C) that evaluate to $\text{CONDUCTANCE}(C) = 0$. These points form a line in the (I, E) space. Now, we find all the points of $(I(C), E(C))$ that have a conductance value of $\text{CONDUCTANCE}(C) = \delta$. These points will also form a line in the (I, E) space. Because of the continuity of conductance, any community, C , that evaluates to an (I, E) point that lies between these two lines must have a conductance value of $\text{CONDUCTANCE}(C) \in (0, \delta)$. In this way we can visually categorize the space.

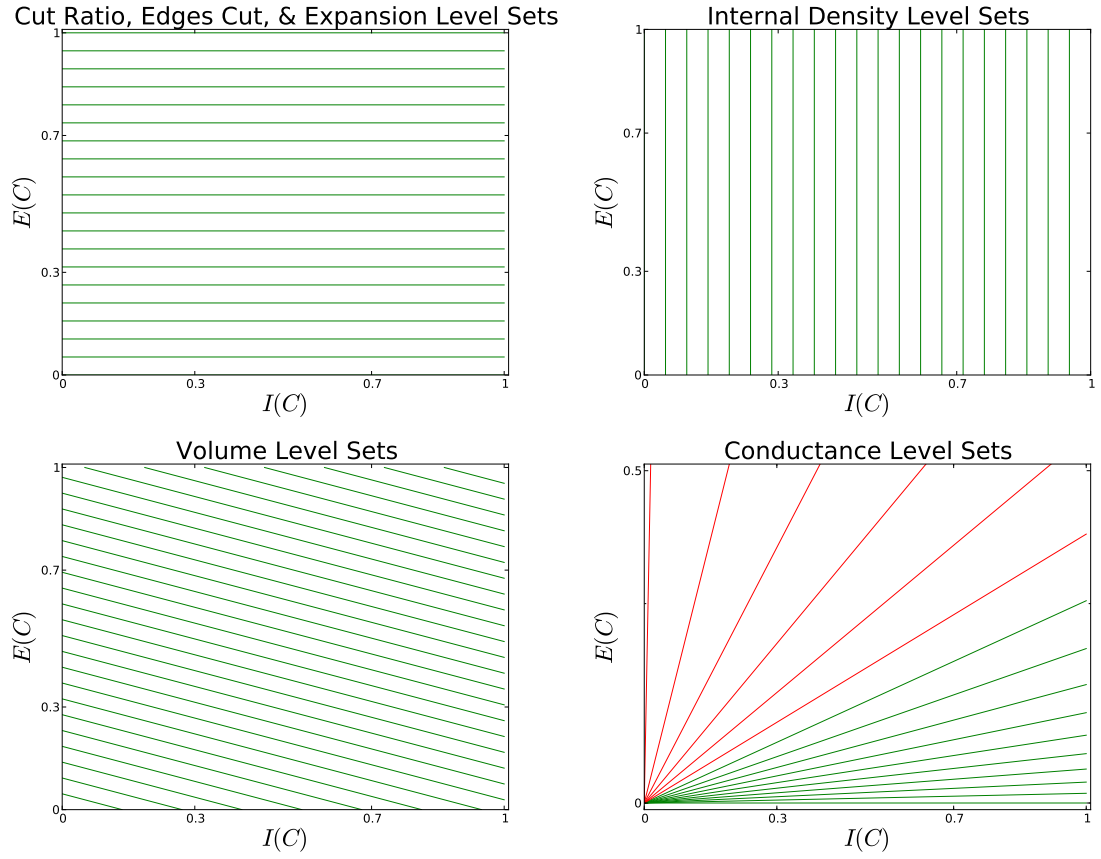


Figure 2.1: Level Sets in the (I, E) plane for different metrics of a single community. There are four ways the (I, E) space is categorized.

In level set figures, any two points in the I, E plane connected by a curve have the same metric value. In our Greedy Algorithm ??, if the algorithm can add a node to the community that crosses a level set to a higher metric valuation, the algorithm will add that node. Visually, the more level sets crossed by a change to the community, corresponds to a higher change in the metric. Traditionally, level sets are used in this manner to show gradient descent to find a local minimum. The minimum that a gradient descent will find, can be found by traveling perpendicular to the level sets. While we find in practice this is a good analogy to understand the behavior of optimizing these metrics, we can

not complete the analogy because the metrics are discrete.

While it is possible to draw conclusions now from the level sets, we proceed with finding communities based on these metrics. In doing so, we confirm and expand experimental results.

Metric	Optimal C	$(I(C), E(C))$
CONDUCTANCE	G	$(x, 0)$
CUT RATIO	G	$(x, 0)$
EDGES CUT	G	$(x, 0)$
EXPANSION	G	$(x, 0)$
INTERNAL DENSITY	any clique	$(1, x)$
VOLUME	G	$(x, 0)$

Table 2.2: Communities that optimize each metric. A value of x , indicates that the optimization is independent of that value.

Greedy Algorithm

The Greedy Single Community Metric Optimization Algorithm ?? takes as input a community and a metric. The algorithm then expands the community, one node at a time, until the metric can not be improved. The resultant community is a local optimum of the metric.

Some metrics require minimization rather than maximization, this algorithm can be adapted accordingly. In the following sections, we use the algorithm by starting with a metric and a subset of two connected nodes. The algorithm produces a series of nested communities, each with an increasing metric score. For each nested community, we compute their $(I(C), E(C))$. This gives us a path through the (I, E) plane. We can use level sets to explain the pattern of node selection that increases the metric.

Algorithm 3: GREEDY SINGLE COMMUNITY METRIC OPTIMIZATION

Input: C , $G = (V, E)$, and METRIC

$inc = 1$

while $inc \geq 0$ and $C \neq V$ **do**

 Let $u \in V$ maximize $METRIC(C \cup u)$.

$inc \leftarrow METRIC(C \cup u) - METRIC(C)$

$C \leftarrow C \cup u$

end while

return C

Expansion, Edges Cut, and Cut Ratio

We now consider metrics that are functions of $E(C)$ and not of $I(C)$: EXPANSION, EDGES CUT, and CUT RATIO. To understand these metrics we plot how they categorize the (I, E) plane with level sets and how iterations of the greedy algorithm choose communities in the (I, E) plane.

For these three metrics, their definitions vary, but their level sets are identical, as shown in Figure ???. The level set of $E(C) = 0$ corresponds to the metric's optimal set of communities. These communities are disconnected from the rest of the graph, and can have an arbitrary internal density. These metrics favor decreases in external density over increases in internal density. In fact, for a community at any position in the (I, E) plane, the node that decreases external density the most will be chosen by the greedy algorithm, rather than a node that improves internal density. The effect of this is visible in the greedy algorithm's path through the (I, E) plane.

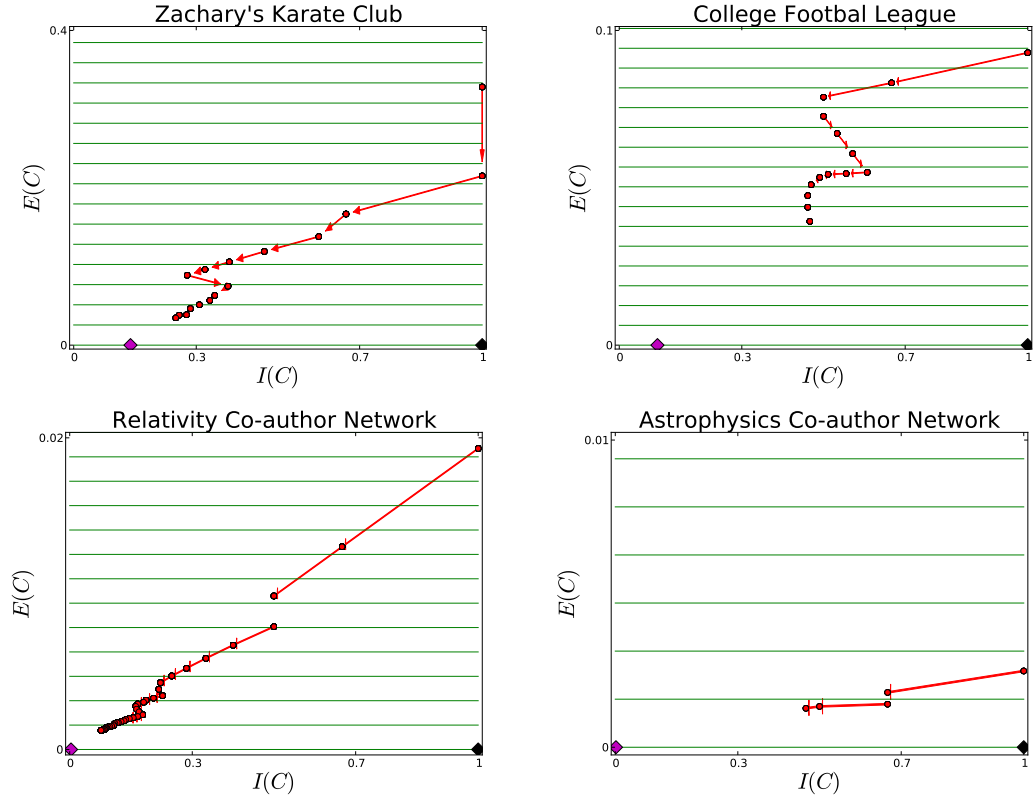


Figure 2.2: External Density based metrics (CUT RATIO, EDGES CUT, and EXPANSION) optimized in different networks. The lower left diamond is the $(I(C), E(C))$ point corresponding to a community of the entire graph. The lower right diamond is the $(I(C), E(C))$ point corresponding to an ideal community. The path corresponds to the intermediary (I, E) values of adding nodes that optimize the metrics using Greedy Algorithm ??

Because all of these metrics only respond to changes in external density, the order of nodes the greedy algorithm adds to the community does not vary between the three metrics. The difference between the three metrics is the point at which they terminate. Termination in this case is determined by the size of the community, $k = \frac{|C|}{|V|}$. Cut ratio is unresponsive to changes in the size of the community, while expansion linearly discounts larger communities. Edges cut heavily favors very large or very small communities. See Fig. ??.

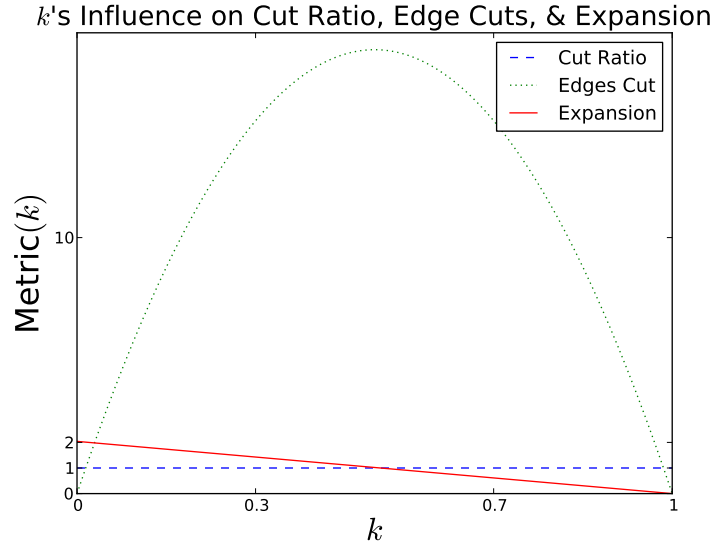


Figure 2.3: Influence of size of community on the values of external density based metrics.

Internal Density as Previously Defined

INTERNAL DENSITY is a function that has been in use before our formalization of $I(C)$. INTERNAL DENSITY is a function of only $I(C)$ and is unresponsive to changes in the external density. Hence, only cliques and subsets of cliques optimize internal density. We do not include indepth analysis, but rather a summary. The level sets of internal density are vertical lines in the I, E plane, as seen in Figure ?? . The greedy algorithm augments our input of two connected nodes to the largest clique it can find (if forced to), as two connected nodes are already a clique.

Volume

A metric that takes both internal and external density into account is volume. The next conclusion is not apparent just from the equation parameterized in

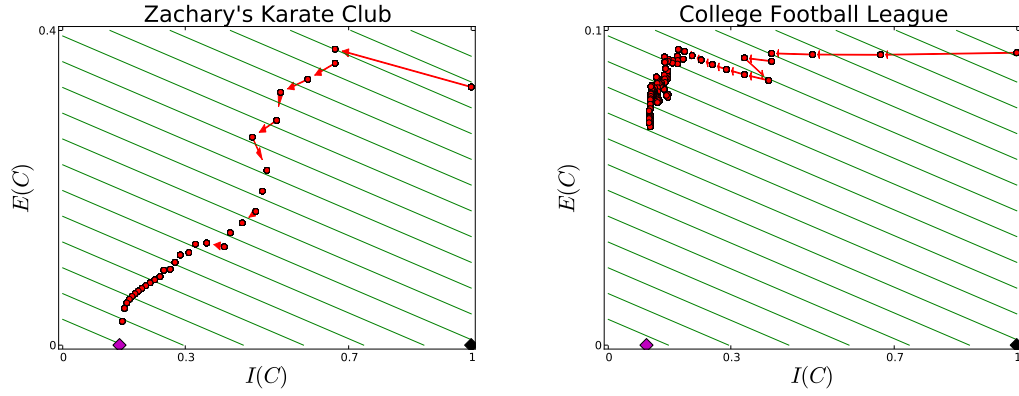


Figure 2.4: Tracing of communities found by volume through the IE plane for a maximum of 100 steps. The lower left diamond is the $(I(C), E(C))$ point corresponding to a community of the entire graph. The lower right diamond is the $(I(C), E(C))$ point corresponding to an ideal community. The path corresponds to the intermediary (I, E) values of adding nodes that optimize the metrics using Greedy Algorithm ??.

terms of internal and external density. However, observing the level sets of volume reveal that the optimal community is at $(I, E) = (0, 0)$ and volume as a metric is optimal for communities with low external density and low internal density. Apart from communities of unconnected nodes, volume can best be optimized by a community encompassing the entire graph. Volume contradicts our intuition that communities should have good internal connectivity.

Conductance

For conductance the level sets are rays radiating from $(I, E) = (0, 0)$, see Fig. ??. As the rays come closer to horizontal, $E(C) = 0$, conductance is closer to optimal. Near $E(C) = 0$, changes to internal density have little effect on the value of conductance. Improvements in conductance come from modifying the community to decrease $E(C)$ as much as possible. If the rays are away from $E(C) = 0$,

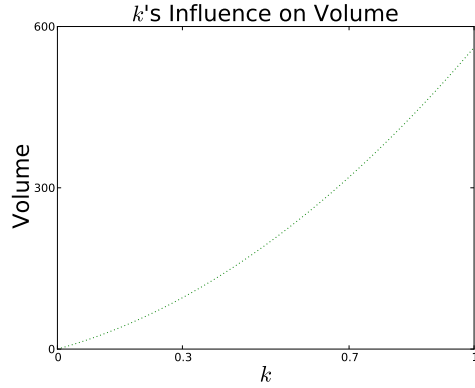


Figure 2.5: The affect, increasing the size of the community has on volume, even for a constant $I(C)$ and $E(C)$.

then improvements to internal density have a larger impact on conductance.

We now analyze the performance of the Greedy Algorithm ?? with conductance and four networks, results are displayed in Figure ?. In the College Football League, the greedy algorithm finds communities in the (I, E) plane where improvements in I and E are balanced. The final community found corresponds to our notion of a good community. For Zachary's Karate Club, the greedy algorithm begins to enter the region where external density determines conductance and returns a community of debateable quality. This effect is more emphasized in the relativity and astrophysics co-author networks. The greedy algorithm initially returns communities in the region of the (I, E) plane with balanced weightings between internal and external density. When external density reaches the region of low external density, the level sets show that small improvements to external density at the cost of lower internal density dramatically improve conductance.

This is the cause of the problem found by Leskovec et. al [?]. As a metric, conductance either incorporates internal density, as in the small College Foot-

ball communities, or does not incorporate internal density, as in the larger Relativity and Astrophysics Co-author communities.

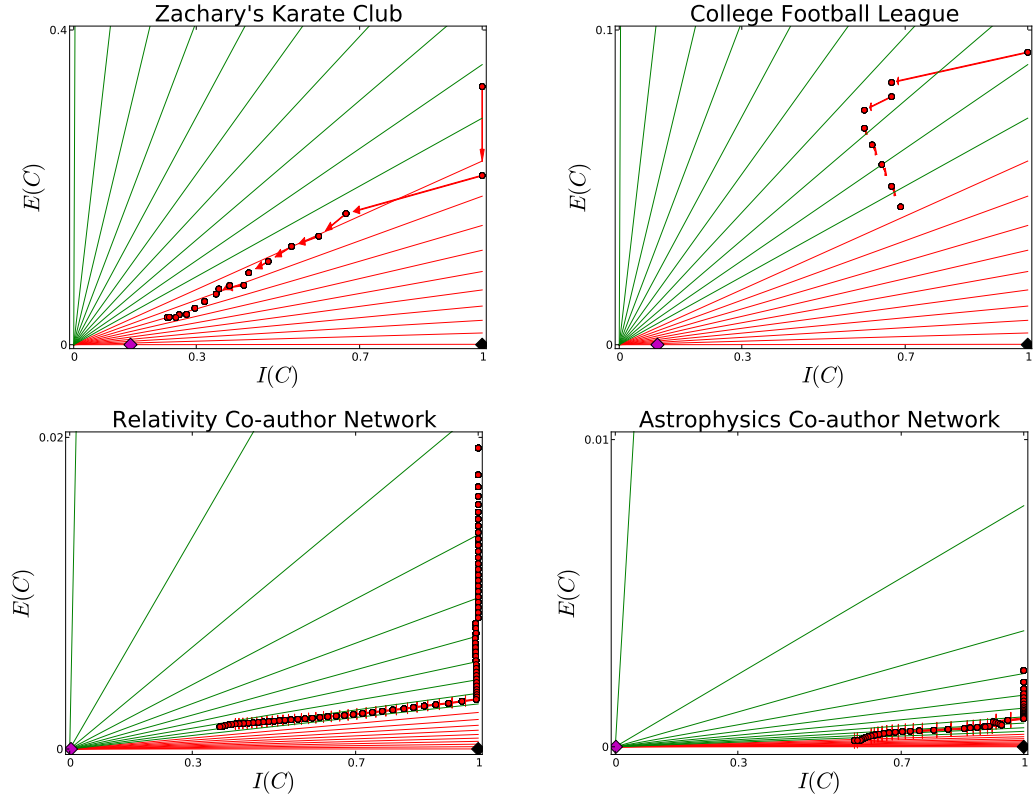


Figure 2.6: The progression of communities that optimize conductance. Note, both the entire graph and the ideal community optimize conductance. In the relativity and astrophysics networks, we stop following the progression of conductance once it becomes clear the entire graph will be engulfed. (In the case of the college football league, a local optimum was reached, but reports an undesirable value of conductance.)

This problem is amplified by the effect a community's size has on conductance. Now we fix the I, E ratio and observe how changes in $|C| = k|V|$ affect conductance in Figure ?? . Conductance always values a larger community more favorably. As long as the community is of small to moderate size and has a large $E(C)$ value, the greedy algorithm will return communities that correspond to our intuition of a good community.

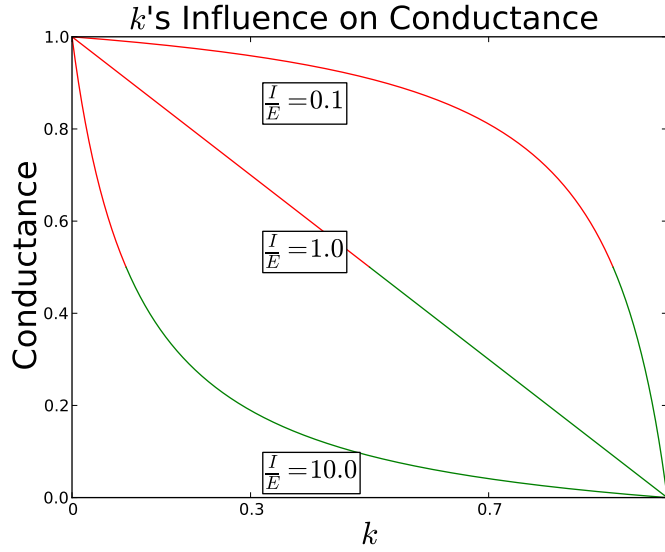


Figure 2.7: Influence of size of community on the value of conductance. The object is to minimize conductance.

2.4 Set of Communities Based Metrics

We now explore metrics that evaluate the strength of a set of communities, $S = \{C_1, C_2, \dots, C_n\}$. Several community detection methods are based on finding a partitioning of the network that optimizes such a metric. The most popular of these metrics is, modularity developed by Newman [?].

2.4.1 Internal Density, External Density, and Conciseness

Our parameterization of internal and external density for single community metrics can not be directly applied to a set of communities, $S = \{C_1, C_2, \dots, C_n\}$. We begin as we did for single communities and consider the characteristics of a good set of communities. A good set of communities is a set of cliques such that every edge is within some community and every community is a maximal

clique. Hence an ideal set of communities has three parameters. Internal density is a representation of how close the set of communities is to being a set of cliques. External density is a representation of how close the set of communities are to covering all edges in the graph. Size of the set of communities is a representation of how concise the set of communities are. With the same methodology for parameterizing and understanding metrics of individual communities we proceed to parameterize metrics for sets of communities with *internal density*, *external density*, and *conciseness*. Formal definitions follow.

Definition 7 (Internal Density of a Set of Communities) For a set of communities, $S = \{C_1, C_2, \dots, C_n\}$, the internal density of the set is the sum of the number of edges that do exist within each community compared to the maximal number of edges that could exist.

$$I(S) = \frac{\sum_{C \in S} (\sum_{u \in C} \sum_{v \in C} w(u, v))}{\sum_{C \in S} |C|(|C| - 1)} \quad (2.9)$$

Definition 8 (External Density of a Set of Communities) In a set of communities, S , the EXT_EDGES is the set of edges not covered by any community. External density is the number of edges in EXT_EDGES compared to the number of edges in the graph.

$$E(S) = \frac{\sum_{(u,v) \in \text{EXT_EDGES}} w(u, v)}{\sum_{u,v \in V} w(u, v)} \quad (2.10)$$

Definition 9 (Conciseness of a Set of Communities) Conciseness is the size of S .

$$\text{CONCISENESS}(S) = |S| \quad (2.11)$$

Our choice of parameter definitions, allows the analysis of any set of communities, including overlapping communities. In particular, our definition of internal density for a set of communities, allows nodes to be placed in multiple communities. External density allows overlapping communities, as well as the conciseness function.

Definition 10 (Ideal Set of Communities) *A set of communities, S , is ideal if it is a set of maximal cliques that cover the graph in very few communities:*

$$I(S) = 1$$

$$E(S) = 0$$

$$|S| = \text{number of connected components of the network.}$$

All three parameters are necessary to ensure a complete description of a set of communities. For any two parameters, there exists a set of communities that can optimize those two parameters. Failure to evaluate the third parameter reveals an undesired characteristic of the set of communities. Figure ??, illustrates the types of communities that can optimize for any two parameters.

2.4.2 Study of Relevant Metrics

Modularity is the most popular of these metrics. It compares the number of internal edges found, to the number of expected edges in a random graph. Modularity was developed by Newman [?] and has found wide spread use due to the fast algorithms for maximizing modularity. In particular, the use of dendrograms in the Louvain Algorithm [?] runs in minutes for large networks.

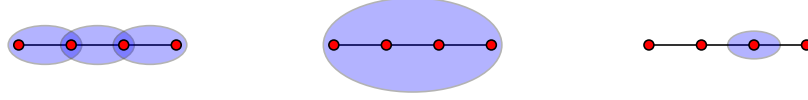


Figure 2.8: The communities that optimize 2 out of 3 parameters. Nodes are in red, lines are edges, and communities are blue ellipses. The left community configuration optimizes $I(S) = 1$ and $E(S) = 0$, but not conciseness at $|S| = 3$. The middle configuration optimizes $E(S) = 0$, $|S| = 1$, but not internal density at $I(S) = \frac{1}{2}$. The right configuration optimizes $I(S) = 1$ and conciseness at $|S| = 1$, but does not optimize external density at $E(S) = 1$.

There is not a closed form parameterization of modularity in terms of our definitions of $I(S)$, $E(S)$, and $|S|$. However, for each module's contribution there is a closed form parameterization in terms of internal and external density for a single community, $I(C)$ and $E(C)$. If we allow, $p = \frac{|C|(|C|-1)}{2L}$ and $q = \frac{|C|(|V|-|C|)}{2L}$, where L is the number of edges in the graph then:

$$\text{MODULARITY}(S) = \sum_{C \in S} pI(C) - (pI(C) + qE(C))^2. \quad (2.12)$$

We first note that if there exists a set of disjoint cliques in the graph, only a partitioning of each clique into a module maximizes modularity. Modularity already aligns more strongly with our understanding of strong communities than previous metrics.

We can not plot the level sets for modularity over a set of communities, but we can plot the level sets for the contribution to modularity from each commu-

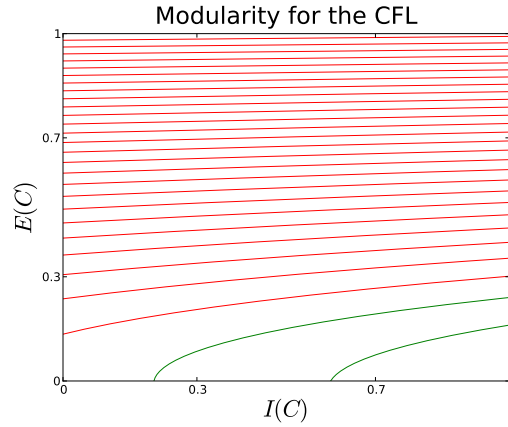


Figure 2.9: The level sets of how MODULARITY treats the $I(C)$, $E(C)$ space for one community of size 9 in the College Football League (CFL). Note the sharp transition from a region that heavily favors improvements in external density to a region that heavily favors improvements in internal density ($E(C) < 0.1$).

nity. In Figure ?? we find that modularity is a two part optimization. When $E(C)$ is large, modularity maximization attempts to decrease $E(C)$ as quickly as possible. Once a threshold of $E(C)$ is crossed, modularity maximization attempts to increase $I(C)$ as quickly as possible. The transition between these two phases of optimization is sudden and revealed by a dramatic turn in the level set curves. The larger the graph the more sudden this transition.

JTODO Modularity has a resolution limit, but that is just because modularity at first tries to optimize external density which is prone to joining communities together. [?]

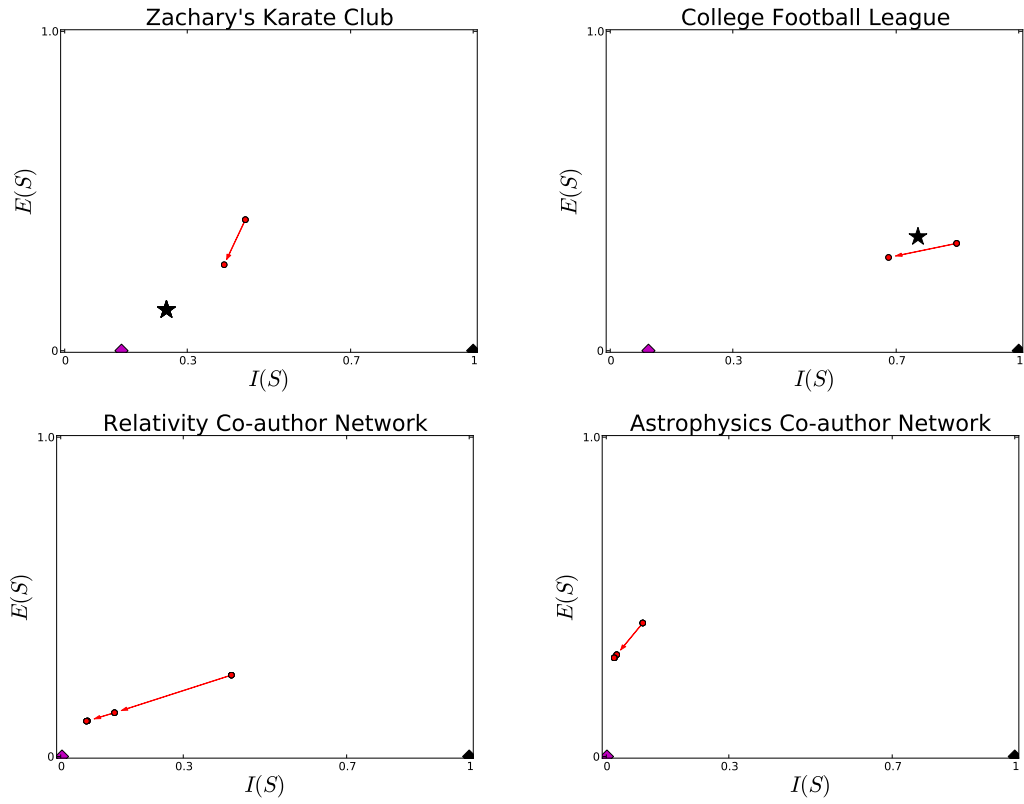


Figure 2.10: Here we run the Louvain Algorithm [?] to maximize modularity. The $(I(S), E(S))$ path is each level of the dendrogram. The $(I(G), E(G))$ value for the entire graph is the diamond in the lower left. In the general relativity and astrophysics co-author networks, modularity does not present much of an improvement over $I(G)$ and has a much higher $E(G)$ value.

CHAPTER 3

A NEW METRIC: LINEARITY

Previously analyzed metrics fell into two categories. In the first, the metrics: edges cut, cut ratio, expansion, and internal density, reflect either internal or external density but not both. These metrics are optimized by sets of nodes that do not provide insight into the structure of the network. The second category, including modularity and conductance, is unpredictable. They have the same values for radically different communities. In being unpredictable, conductance and modularity sometimes produce strong communities and sometimes, especially as the size of the communities increases, return poorly connected communities[?].

In this chapter we will present metrics for single communities and sets of communities that measure both internal and external density and are consistent.

3.1 Single Community Detection

Let us now discuss the criteria of a good metric and find a such a metric. In the previous chapter, we show that internal and external density provide bounds for the characteristics of diameter, average shortest path, etc. While it is possible to design a metric that covers an arbitrary number of characteristics, we argue a metric that reflects both internal and external density provides a good measure of many characteristics. Thus, a good metric should reflect a community's internal and external density. In particular, the metric should be optimized by the ideal community and minimized, or minimal, for communities with poor values of internal and external density. It easy to check the how a metric handles

extreme communities, but we also want an element of predictability for how the metric handles all communities. Here is one definition of predictability. Let communities C and C' have internal and external values: $(I(C), E(C)) = (x, y)$ and $(I(C'), E(C')) = (x + \delta_x, y + \delta_y)$. Then, a metric M as function of internal and external density is predictable if:

$$M(x + \delta_I, y + \delta_E) - M(x, y) = M(\delta_I, \delta_E). \quad (3.1)$$

A linear metric satisfies all mentioned criteria.

Definition 11 (Linearity) *Let C be a community, $\text{LINEARITY}(C)$ is a metric with a linear weighting of internal density, $I(C)$, and external density, $E(C)$, such that $\text{LINEARITY}: C \rightarrow [-1, 1]$.*

$$\text{LINEARITY}(C) = aI(C) - bE(C) \quad (3.2)$$

The constants a and b are restricted to $a, b \in (0, 1]$.

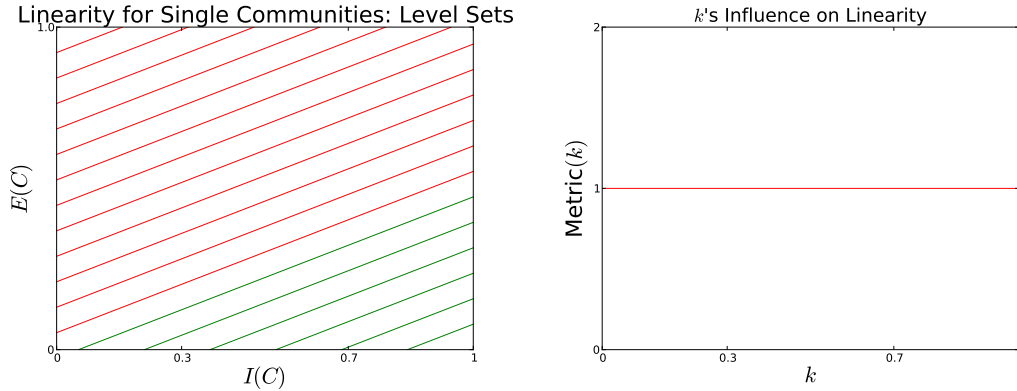


Figure 3.1: The level sets are predictable. The size of a community does not influence LINEARITY.

In some applications, we may want to relax the predictability constraint to find communities of a certain size, internal, or external density. If this is the case, we suggest a polynomial approach to building a metric.

Definition 12 (General Metric) *The general metric for evaluating any single community is a sum of polynomial functions on internal and external density, weighted with a function of the community's size.*

$$\text{GENERAL}(C) = \sum_{i=0} f_i(C)I(C)^i - g_i(C)E(C)^i \quad (3.3)$$

The functions f_i and g_i can be any function of the size of a community.

When using the general equation, the level sets and size of the community's affect should be analyzed. In particular, the local and global maximums should correspond to the desired communities, and the level sets should aid finding desired communities in a manner similar to gradients.

We now analyze LINEARITY in the same way we analyzed other single community metrics. The level sets in Figure ?? reveal a predictable metric that is only optimized by the ideal community. The size of the community does not change the behavior of the metric. We test the LINEARITY metric with the Greedy Algorithm ?? on four networks in Figure ?. The parameter a was set to one, while b required a binary search. The parameter b was set within four steps, such that the greedy algorithm did not return the entire graph or the initial community. There is room for future research on the impact of b . For all possible values of b , we find very few different communities.

3.2 Multiple Community Detection

As we constructed a linear metric for a single community, we now construct a linear metric for sets of communities, S . The characteristics of a set of communities can be summarized by internal density (Definition ??), external density

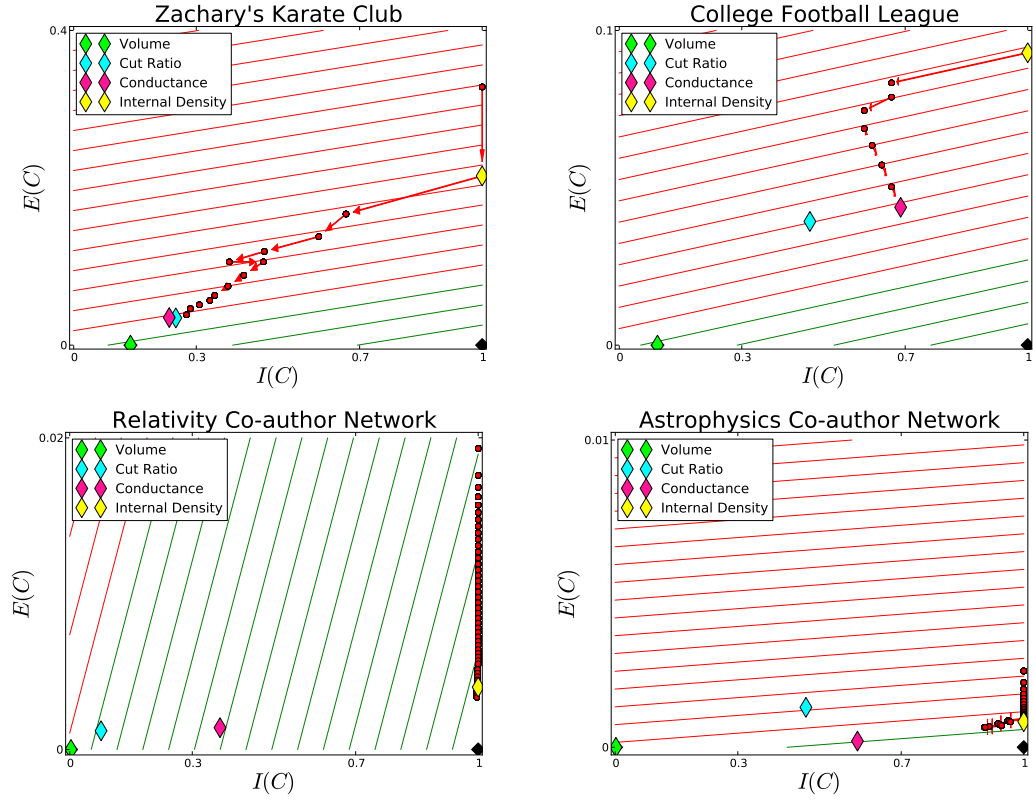


Figure 3.2: Single Communities produced by Linearity, in red. The colored diamonds are the $(I(C), E(C))$ values produced by previously tested single community metrics. In the Karate Club and CFL network Linearity returns a community close to conductance. In the relativity and astrophysics network Linearity returns a community closer to internal density. The black diamond is the ideal community, but does not exist in these networks.

(Definition ??) and the number of communities in the set (Definition ??). A good metric should reflect a set of community's $E(S)$, $I(S)$, and $|S|$ values. In particular, the metric should be optimized by the ideal set of communities and minimal for sets of communities with poor values of $E(S)$, $I(S)$, or $|S|$. As with single community metrics, we want a predictable metric for sets of communities. The most predictable metric is linear.

Definition 13 (Linearity) Given a set of communities, $S = \{C_1, C_2, \dots\}$

$\text{LINEARITY}(S)$ is a metric mapping S to $[-2, 1]$.

$$\text{LINEARITY}(S) = aI(S) - bE(S) - c|S|, \quad (3.4)$$

where $a, b, c \in (0, 1]$.

Depending on the application, communities with particular values of internal density, external density, or size may be desired. In these cases we recommend a polynomial expression for the metric.

Definition 14 (General Metric) *Our metric for single communities in its greatest generality:*

$$\text{GENERAL}(S) = \sum_{i=0} f_i(S)I(S)^i - g_i(S)E(S)^i - h_i|S|^i \quad (3.5)$$

Whenever creating a metric of this form it is recommended to check the level sets for elements of unpredictability.

To maximize our linear algorithm for sets of communities, we will create a greedy algorithm with two stages. The first is to use an adapted Louvain algorithm [?] to find a partition maximizing linearity. The final stage will be to expand each partition to include individual nodes. This algorithm is a heuristic to maximize linearity, but runs in complexity equivalent to the Louvain Algorithm $O(JTODD)$.

We first state a conjecture about greedy algorithms.

Conjecture 15 (Maintaining Internal Density) *Let community C have internal density $I(C)$ and external density $E(C)$. If an expansion of C to include node v_1 results in a decrease in internal density, ie $I(C \cup v_1) < I(C)$, then expansion will only*

create a community with internal density $I(C \cup v_1 \cup v_2 \cup \dots \cup v_i) = I(C)$ by including a large clique, v_1, v_2, \dots, v_i .

We have stated the conjecture corresponding to single communities, and a similar conjecture exists for sets of communities. The conjecture comes from our experience that once internal density is decreased it can rarely be increased by a greedy algorithm. When internal density is decreased and then increased by a greedy algorithm, a clique is involved. Improvements to internal density are hard, improvements to external density are easy. External density can be lowered by incorporating more nodes into the community and minimized by including the entire connected component containing the community. This leads to the development of a greedy algorithm that maintains or improves internal density, until only improvements in external density can be made. This order of greedy algorithm is opposite the order of modularity maximization that first minimizes external density and then tries to maximize internal density, see Section ??.

JTODO include pseudo code for the Louvain Algorithm

To adapt the Louvain Algorithm, we must show that the following property holds:

Property 16 (Louvain Criteria) *Let M be any metric, $S = \{C_1, C_2, \dots\}$, and communities C_i and C_j have no edges between them. Let S' be the set of communities S , with communities C_i and C_j replaced by their union, ie $S' = S - C_i - C_j + C_i \cup C_j$, then:*

$$M(S) \geq M(S') \tag{3.6}$$

The contributions of the characteristics $I(S)$ and $E(S)$ decrease linearity by joining unconnected sets of nodes. The third characteristic of our linearity metric $|S|$ can increase linearity by joining unconnected sets of nodes. However, we will limit ourselves to a , b , and c values such that overall linearity is not increased and use the Louvain Algorithm.

To adapt the Louvain Algorithm, we could exchange the modularity metric for the linearity metric and get a good partition. We take it one step further and use the Conjecture ???. So far we have not set the parameters a , b , and c in linearity. From the conjecture, the algorithm should optimize internal density first and then external density. We begin with the parameters set to $a = 1$, $b = 0$, and $c = \frac{1}{|V|}$. Maximizing the linearity metric with these parameter values results in a partitioning of maximal cliques. We now relax the parameter $b = \delta_b$ and complete the Louvain Algorithm. This will result in a partitioning of near cliques. The process is continued, gradually increasing b and completing the Louvain Algorithm on the new parameters. The question is when to stop increasing b . In practice we increase b until the partitioning of the graph is the entire graph and then retract b by one increment.

JTODO include demonstration of increasing b .

The final step of our algorithm for maximizing linearity uses the advantage of overlapping communities. Given the partitions produced by the previous step, we augment each partition by nodes that increase linearity. Note, because the partition produced in the previous step was a local maximal partitioning, no partition will be augmented to include another partition.

We now run our algorithm on four data sets and compare to known results and modularity results. The object is to provide a preliminary analysis of the algorithm. In depth results are provided in Chapter ??.

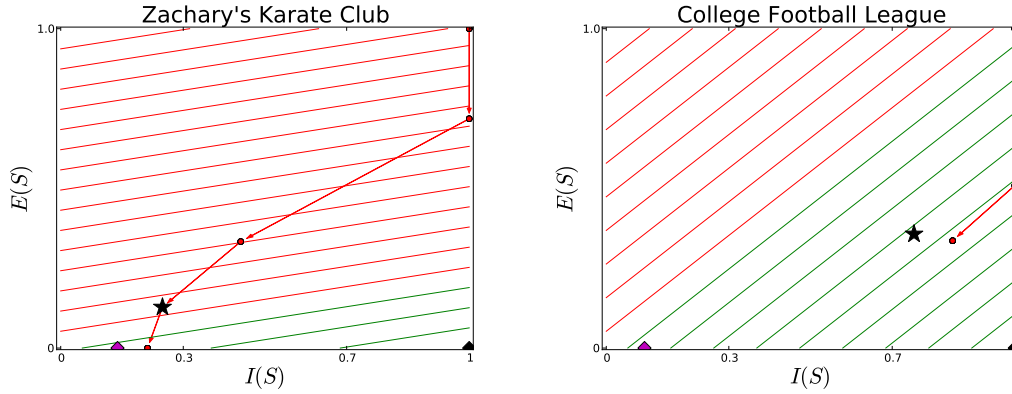


Figure 3.3: Linearity for sets of communities compared to Known Solutions (the black stars). Linearity produces sets of communities with better values of internal and external density. In depth analysis is provided in Figure: ??.

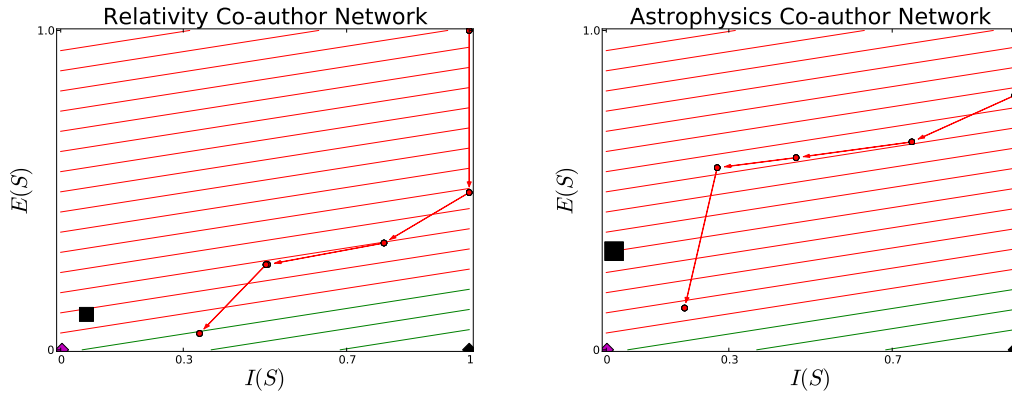


Figure 3.4: The $I(S)$, $E(S)$ values produced by modularity are provided by black squares. The linearity path is traced in red. The first segment corresponds to finding maximal cliques, the middle segments correspond to improvements in the partition due to the Louvain Algorithm. The last segment is from expanding the partitions to produce overlapping communities. In both cases linearity produces sets of communities with better internal and external density values.

CHAPTER 4

PARALLEL COMMUNITY DETECTION

The metric based detection methods can analyze networks up to $\approx 50k$ nodes. At that point, the computational time to complete these methods becomes prohibitive, but there is no limit of a network's size. The Amazon network has 500 thousand nodes, the Twitter network has 17 million nodes, and the memetracker network has 96 million nodes. Analysis of these large networks will require a parallel algorithm. In this chapter we adapt and modify our understanding of communities to develop an embarrassingly parallel algorithm for community detection.

4.1 Parallel Algorithms

Here we briefly introduce the two important aspects of good parallel algorithms.

Parallel algorithms take an application and divide the computational cost into units. Each unit of computation is then assigned to a processor. The wall clock time of an application is the time it takes from beginning the computation to receiving an answer. The computational time is the number of processors times the wall clock time, in other words the man hours of cpu's. For most applications, we are concerned with lowering the wall clock time. For parallel algorithms, this can be achieved by increasing the number of processors used. Ideally, if we double the number of processors, we halve the wall clock time. This is known as perfect scalability. How close an algorithm comes to perfect scalability is the scalability of the algorithm. The primary bottleneck prevent-

ing perfect scalability is the communication cost between processors. In order for one processor to finish computations, it may rely on results produced by another processor. For a good parallel algorithm it is essential to minimize communication between processors.

The second aspect of good parallel algorithms, is that they redefine how we perceive an application. The intuitions and techniques used for small datasets do not scale to large datasets. And often times, the intuitions and techniques needed to process large datasets do not work on small datasets.

We will design a parallel algorithm with these two aspects in mind. The perspective of algorithms for small networks, metric based detection methods, has been, *is this community a good community to include in the set of communities*. Modularity and linearity have a parameter that depends on the entire set of communities. This perspective requires knowledge of all other communities within the network. Parallelizing this perspective leads to a high communication cost between processors. For a good parallel algorithm we need a new perspective. The perspective we use is *given a node and a set of nodes, do they belong to the same community?* There is no communication necessary for this perspective to be parallelized.

4.1.1 Previous Parallel Algorithms

JTODO, there have been a few parallel algorithms:

- John and Laura with Ising like approach
- (α, β) communities for finding near clique communities

- Microsoft for finding near clique communities.

Each approach defines a local definition of a community and searches for communities independently. We will use the same outline.

4.2 Characteristics and Statistical Significance

Previously, all characteristics have been of a community to determine the strength of a community. Now, we consider the characteristics between a node and a set of nodes to determine the likelihood that the node and set of nodes belong to the same community. In a social network a characteristic between a person, n , and a set of people, C is the number of friends n has in C . In a citation network a characteristic between a paper, n , and a set of papers, C , is number of papers n cites in C . We name the characteristic representing the number of connections between node, n , and set of nodes, C , to be χ_e .

Definition 17 (Characteristic χ_e) *The characteristic χ_e is the number of edges between a node, n , and a set, C :*

$$\chi_e(n, C) = |E(n, C)|. \quad (4.1)$$

Another characteristic relies on the percentage of edges from n to other members of the graph are within C . In a social network, this characteristic is the percentage of friends node n has that are within C . We name this characteristic χ_p .

Definition 18 (Characteristic χ_p) *The characteristic χ_p is the percentage of edges n*

has that lead to members of C :

$$\chi_p(n, C) = \frac{|E(n, C)|}{\text{degree}(n)}. \quad (4.2)$$

Given more details about nodes and their connections to a set of nodes more characteristics can be defined. In a social graph with dates of when friendships begin, a duration characteristic reflects average time of friendship between a node and a set of nodes. Depending on the network more characteristics may exist. For an undirected, unweighted graph, the only two characteristics between a node, n , and a set of nodes, C , are χ_e and χ_p .

Given a node, n , and a set of nodes C , we would like to find the probability that n and C belong to some larger community, C' . In particular, the two probabilities we would like to compute are:

$$P(n \cup C \subset C' | \chi_e(n, C))$$

$$P(n \cup C \subset C' | \chi_p(n, C))$$

Direct calculation of these quantities is only possible for networks where internal and external edges are created with known probability distributions. If either of these probabilities is particularly high, we say the probability that n and C belong to a larger community is statistically significant. We assign the bounds at which $\chi_e(n, C)$ and $\chi_p(n, C)$ are statistically significant to be b_e and b_p . These thresholds b_e and b_p are set by each application. We can now define the set of nodes with a statistically significant χ_e value to be $\Phi_e(C) = \{n | \chi_e(n, C) \geq b_e\}$. Similarly, $\Phi_p(C) = \{n | \chi_p(n, C) \geq b_p\}$.

We can use the $\Phi_e(C)$ and $\Phi_p(C)$ to define some characteristics of communities.

Definition 19 (Closed Community) *A community C is closed under two conditions. The first is for every node n in C , either χ_e or χ_p is statistically significant. The second is that no statistically significant node in either $\Phi_e(C)$ or $\Phi_p(C)$ is not included in C .*

$$\Phi_e(C) \cup \Phi_p(C) = C \quad (4.3)$$

An application will set the thresholds b_e and b_p to determine statistically significant, we use these to describe the strength of a community.

Definition 20 (Strength of a Community) *The strength of a community C is a set of bounds, $\{b_e, b_p\}$.*

4.3 Algorithm

We would like to find all of the closed communities of a network. Given the definition of χ_e and χ_p a greedy algorithm presents itself. Let us say a seed is a set of nodes with a very high probability of belonging to the same community. This seed can be a large clique. How to find seeds is covered in Section ??.

Given a seed, S , it may already be a closed community. If it is not closed, there is some node, n , such that $n \cup S$ has the highest probability of being contained in a community. We can augment the set of nodes to include n and continue in a similar fashion until we have found a set of closed nodes. The pseudocode is in Algorithm ?? and the algorithm for GET SEEDS is in Section ?. The rest of this section analyzes the greedy algorithm.

Algorithm 4: FIND ALL COMMUNITIES

Input: $G = (V, E)$

Seeds = GET SEEDS(G)

Communities = {}

for $S \in \text{Seeds}$ **do**

while S is not closed, ie $S \neq \Phi_e(S) \cup \Phi_p(S)$ **do**

Find node $n \in \Phi_e(S) \cup \Phi_p(S)$ with highest probability $n \cup S$ is a community.

$S \leftarrow S \cup n$

end while

Communities \leftarrow Communities $\cup \{S\}$

end for

return Communities

4.3.1 Correctness of Expansion

The step of augmenting the seed to include another node is the expansion step. The algorithm will expand to include the node with either the maximum $\chi_e(n, S)$ or $\chi_p(n, S)$ value. We would like to know the probability that the algorithm expands to include a node such that n and S do belong to the same community. Let the probability of an external edge existing be drawn from the distribution P_E . Similarly, let the probability of an internal edge existing follow the distribution P_I .

Let us first consider the node, n , that has the maximum $\chi_e(n, S)$ value. We can calculate the probability that n and S belong to the some community, C :

$$P(n \cup S \subset C | n \text{ maximizes } \chi_e(n, S)) = 1 - \left(1 - \sum_{x=0}^{|S|} P_I(X = x) P_E(X < x)^{|V-C|} \right)^{|C-S|}.$$

This equation can be difficult to calculate. We now simplify the equation. If we want the probability to grow as:

$$P(n \cup S \subset C | n \text{ maximizes } \chi_e(n, S)) \geq 1 - \left(\frac{1}{2} + a\right)^{|C-S|}.$$

The more nodes needed to find the closed community C , the higher the probability of selecting a node that belongs to C . The factor a controls how quickly the probability of selecting such a node converges to 1. To satisfy this bound, we need:

$$\frac{1}{2} - a \leq \sum_{x=0}^{|S|} P_I(X = x) P_E(X < x)^{|V-C|}$$

For certain P_I and P_E , this may be hard to calculate a . A trick is to use a step function to bound $P_E^{|V-C|}$ from below. The function $P_E^{|V-C|}$ is monotonically increasing from 0 to 1 and at $y \in (0, |S|)$ crosses $\frac{1}{2}$. We bound $P_E^{|V-C|}$ with the step function. If $x < y$, then $P_E(X < x)^{|V-C|}$ is bounded below by 0. If $x > y$, then:

$$P_E(X < x)^{|V-C|} \geq \frac{1}{2}.$$

This lower bounding function can be used to find a stricter and simpler requirement:

$$1 - 2a \leq \sum_{x=y}^{|S|} P_I(X = x)$$

We calculate these values for the Binomial and Power Law Distributions in the following sections.

Binomial Distribution

We now presume the graph is a random binomial graph. Edges between members of the same community exist with probability p and edges between members of different communities exist with probability q . The number of edges

a node has into a set of nodes, S , is either represented by the random variable $B(|S|, p)$ or $B(|S|, q)$, depending on whether or not it belongs to a community containing S . The binomial distributions have defined probability and cumulative density distributions that we can use in Equation ???. The equation allows us to calculate what p and q must be for the algorithm to choose a node n such that n and S belong to some community C .

We want to find a y and $p, q \in (0, 1)$ that satisfy the previously found Equation ???:

$$\sum_{x=y}^{|S|} P_I(X = x) \geq 1 - 2a.$$

The equation determining y is:

$$P_E(X < y) = \frac{1}{2^{\frac{1}{|V-C|}}}.$$

If $|S| > 20$, we can approximate $B(|S|, q)$ with a normally distributed random variable and use standard deviations to calculate y . For graphs of the size 10^z the equation is roughly $1 - 10^{-z}$. For a million node network, y must be at least 5.5 standard deviations from the mean $q|S|$. For any network with less than a million nodes we approximate $P_E(X < x)^{|V-C|}$ by the step function $\frac{1}{2}I_{x>7q|S|}$. We now apply y :

$$1 - 2a \leq \sum_{7q|S|}^{|S|} P_I(X = x).$$

The function $P_I(X = x)$ is determined by $B(|S|, p)$. If we set a , we can find what p must be. If we set a to $\frac{1}{4}$, then $p = 7q$. The probability of choosing an n such that $n \cup S$ is a subset of a community C is then:

$$P(n \cup S \subset C | n \text{ maximizes } \chi_e(n, S)) \geq 1 - \left(\frac{3}{4}\right)^{|C-S|}.$$

If the network contains a million nodes, 10 of which we want to find from $C - S$, then with probability 94% the algorithm selects a desired node.

JTODO: provide calculations for $P(n \cup S \subset C | n \text{ maximizes } \chi_p(n, S))$.

For small graphs the bound that the probability of edges between members of different communities is less than $\frac{1}{7}$ does not always hold. For all available large networks, the probability of an edge between members of different communities is close to 0. In these large networks requiring $p \geq 7q$ is reasonable.

Pareto Distribution

We now perform a similar calculation for the Pareto Distribution, a power law distribution. The distribution has the parameters: x_m , the minimum connectivity of n to S and α the rate of decay. For nodes not in C , the probability and cumulative density of $\chi_e(n, S)$ are:

$$\begin{aligned} P_E(X = x) &= \frac{\alpha x_m^\alpha}{x^{\alpha+1}} \\ P_E(X < x) &= 1 - \left(\frac{x_m}{x}\right)^\alpha. \end{aligned}$$

For nodes in C , we assume P_I is P_E shifted by δ . The probability and cumulative densities for $\chi_e(n, S)$, where n is in S :

$$\begin{aligned} P_I(X = x) &= \frac{\alpha(x_m + \delta)^\alpha}{x^{\alpha+1}} \\ P_I(X < x) &= 1 - \left(\frac{x_m + \delta}{x}\right)^\alpha \end{aligned}$$

We limit the graph to all nodes that are connected to S . This sets x_m to one. A node that is in C is connected to S by some minimal amount, $1 + \delta$.

Given P_I and P_E , the equation for $P(n \cup S \subset C | n \text{ maximizes } \chi_e(n, S))$ can then

be calculated:

$$P(n \cup S \subset C | n \text{ maximizes } \chi_e(n, S)) = 1 - \left(1 - \frac{(1 + \delta)^\alpha}{|V - C| + 1} \left(1 - \left(1 - \frac{1}{(1 + \delta)^\alpha} \right)^{|V - C| + 1} \right) \right)^{|C - S|}.$$

We provide sample values in Table ?? . In the table, let set A be the set of nodes connected to S , but that do not belong to a community containing S . Let set B be the set of nodes connected to S and belong to a community containing S . The probability $P(n \cup S \subset C | n \text{ maximizes } \chi_e(n, S))$ is the probability the expansion step picks a node from set B .

$ A $	$ B $	δ	$P(n \cup S \subset C n \text{ maximizes } \chi_e(n, S))$	description
200	20	5	78%	S is half of C
200	30	5	90%	S is $\frac{1}{4}$ of C
400	30	5	67%	number of nodes doubled
200	30	2.5	63%	P_I shift halved

Table 4.1: The probability that the expansion step recovers a desired node from set B .

JTODO include calculation of $P(n \cup S \subset C | n \text{ maximizes } \chi_p(n, S))$.

4.3.2 Seeds

Our algorithm and probability of correctness rely on having a set S such that S is a subset of some community C . From our analysis, the larger S is, the higher the probability the algorithm recovers C . We would like to find maximally sized seeds in a local manner. We introduce the function $B_r(S)$ as the set of all nodes within r hops from a node in the set S . This is called the ball of radius r around S . We now prove a short theorem.

Theorem 4.3.1 (Size of S with diameter at most 3) *Within a community, C , there exists a subset S with diameter at most 3, and of size:*

$$|S| \geq \max_{L \subset C} \{|B_1(L) \cap C|\} > \max_{L \subset C} |L|, \quad (4.4)$$

where L is a clique within C .

Proof 21 *The proof is straight forward. Let L be a clique such that $L \subset C$. Consider the ball $B_1(L) \cap C$, ie all nodes in C that are connected to a member of L . For all $u, v \in B_1(L) \cap C$, let u be connected to the node $n_u \in L$ and v be connected to the node $n_v \in L$. Then a candidate shortest path between u and v is $(u, n_u), (n_u, n_v), (n_v, v)$, a path of length 3. Thus, $B_1(L) \cap C$ has a diameter of at most 3. Since this is true for all such balls centered around cliques within the community, it must be true for the largest such ball.*

The theorem allows us to create a local algorithm for finding seeds. We presume that the nodes of $B_1(L) \cap C$ form a community in the subgraph $B_1(L)$. In practice, we find these seeds are easy to recover from the network.

Theorem 4.3.2 (The size of a community containing two seeds) *Let seeds, S_1 and S_2 be found by Algorithm ???. Then the minimum size of a community, C , containing both S_1 and S_2 is $|C| \geq 4 \max(|S_1|, |S_2|)$.*

Proof 22 (JTODO complete proof) *Since both seeds were found by Algorithm ??, the shortest path between any pair of nodes, $n \in S_1$ and $m \in S_2$ is at least 4. If there is a community containing both S_1 and S_2 , then we can say there is a set of nodes B between seeds S_1 and S_2 . If the entire community is ...*

Algorithm 5: GET_SEEDS

Input: G

Seeds = {}

while There exists a clique of size ≥ 5 in G **do**

 Let L be a large clique in G .

 Let S be a community in the subgraph of $B_1(L)$.

 Seeds \leftarrow Seeds $\cup \{S\}$

 Remove all nodes in $B_1(L)$ from G

end while

return Seeds

4.4 Scalability

The algorithm for finding all communities first finds all seeds and then expands all seeds. Each step can be performed in an embarrassingly parallel manner.

Finding seeds can be run in parallel by assigning a randomly selected node to each processor. Each processor then finds a seed containing the node. The nodes considered by one round of seed selection are then removed from the graph and the processes is continued until all seeds are found.

Expanding each seed is an independent process. Each seed is assigned to a processor and expanded until a closed community is found. The hierarchical structure of a graph can be found by finding a larger community that is also closed.

Finding seeds and communities has near perfect scalability.

Analysis of the communities is not a parallel process. Future work could involve how to process the communities found and find the hierarchical structure in a parallel manner. For now we find this structure in a brute force manner.

4.5 Results

Our algorithm runs quickly and efficiently. We compare the communities found with this algorithm to others in the next chapter. The algorithm is named PARALLEL.

JTODO include a bound on the number of communities that exist.

CHAPTER 5

CASE STUDIES OF NETWORKS

So far we have been developing algorithms. Here, we show how those algorithms perform on a variety of networks.

5.1 Known Community Comparisons

We cover in depth community detection performance on two networks with known community structure.

5.1.1 Karate Club Network

The Karate Club Network represents a set of students belonging to a karate club. Zachary studied the students in [?] and found that students interacted with each other outside of the club's practice times. In our representation the students are the nodes and their interactions are the edges. In the course of Zachary's observations, the club split into two groups that wanted to practice separately. We consider the two groups the club split into to be the known communities. We now compare the communities found by different detection methods with the known communities in Figure ??.

Each of the detection methods produces different sets of communities, each with its own merit. Compared to the known communities, the communities produced by linearity are the most similar. Linearity produces two communities with an overlap. Within the network there is a set of centrally connected nodes, these are the ones in the overlap. Maximizing modularity produces four

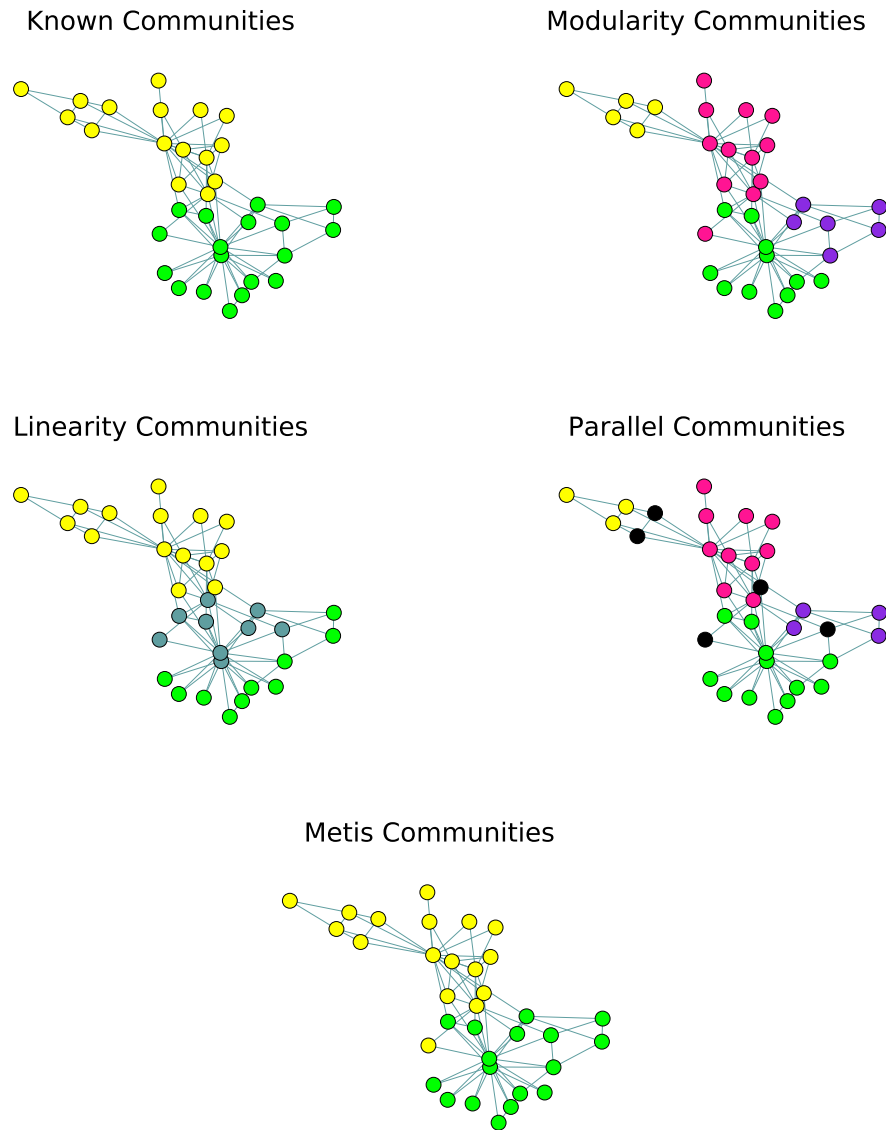


Figure 5.1: Communities produced by the different community detection methods, communities are marked by coloration. Linearity produces two communities with the overlap colored in grey. Parallel produces four communities with four nodes in black not belonging to any community.

communities. Two of the communities are large portions of the two known communities. The additional two communities are more independent and near clique like. Parallel produces four communities that are subsets of the commu-

Communities From	$I(S)$	$E(S)$	$ S $	Average Diameter	NCP
Known	0.25	0.13	2		
Modularity	0.39	0.27	4		
Linearity	0.21	0.0	2		
Parallel	0.43	0.35	4		
Metis					

Table 5.1: Internal Density, External Density, and number of communities for the set of communities returned by each detection method on the Karate Club network. NCP is the network community profile.

nities produced by maximizing modularity. Parallel does not classify five nodes, marked in black. These nodes have exactly two edges, each going to a different community. These nodes do not to have a statistically strong connection to any community.

5.1.2 College Football Network

The college football network represents the 115 collegiate football teams and their games. The nodes are the teams and the edges represent pairs of teams that played a game. The collegiate teams are split into divisional conferences, we consider these to be the known communities. For the games, a team must play nearly every member of its conference. Additionally, each team plays teams from other conferences. If team A belongs to a conference with only a few teams, then A must play more teams from other conferences than a team within a larger conference. This makes large conferences easy to detect and smaller conferences harder to detect.

The communities returned by different detection methods are in Figure ??.

There are twelve conferences and we plot each at an hour on a clock face. Each

Communities From	$I(S)$	$E(S)$	$ S $	Average Diameter	NCP
Known	0.75	0.36	12.0		
Modularity	0.68	0.29	10.0		
Linearity	0.85	0.34	13.0		
Parallel	0.87	0.35	13.0		
Metis					

Table 5.2: Internal Density, External Density, and number of communities for the set of communities returned by each detection method on the CFL network.

color represents a different community. For a majority of the nodes, all detection methods produce the same communities. The smallest conferences are located at five o'clock and ten. Each detection method breaks these conferences up and handles their nodes in a different way. Modularity places each node in another conference. The nodes from the five o'clock conference are incorporated into the six and seven o'clock conferences. Linearity and parallel create two communities for the split five o'clock conference. Parallel does not classify four of the nodes, in black, in the ten o'clock conference. These nodes do not have a strong connection to another community. The last difference is that linearity and parallel break the nine o'clock conference into two communities. Overall, modularity deviated from the known communities by creating fewer communities to cover more of the edges. Linearity and parallel deviated from the known communities by creating more and denser communities.

5.2 Collaboration Networks

The set of collaboration networks consist of authors as nodes and edges representing a collaboration. We modify these networks slightly. In collaboration networks there are several graduate students that collaborate with only one pro-

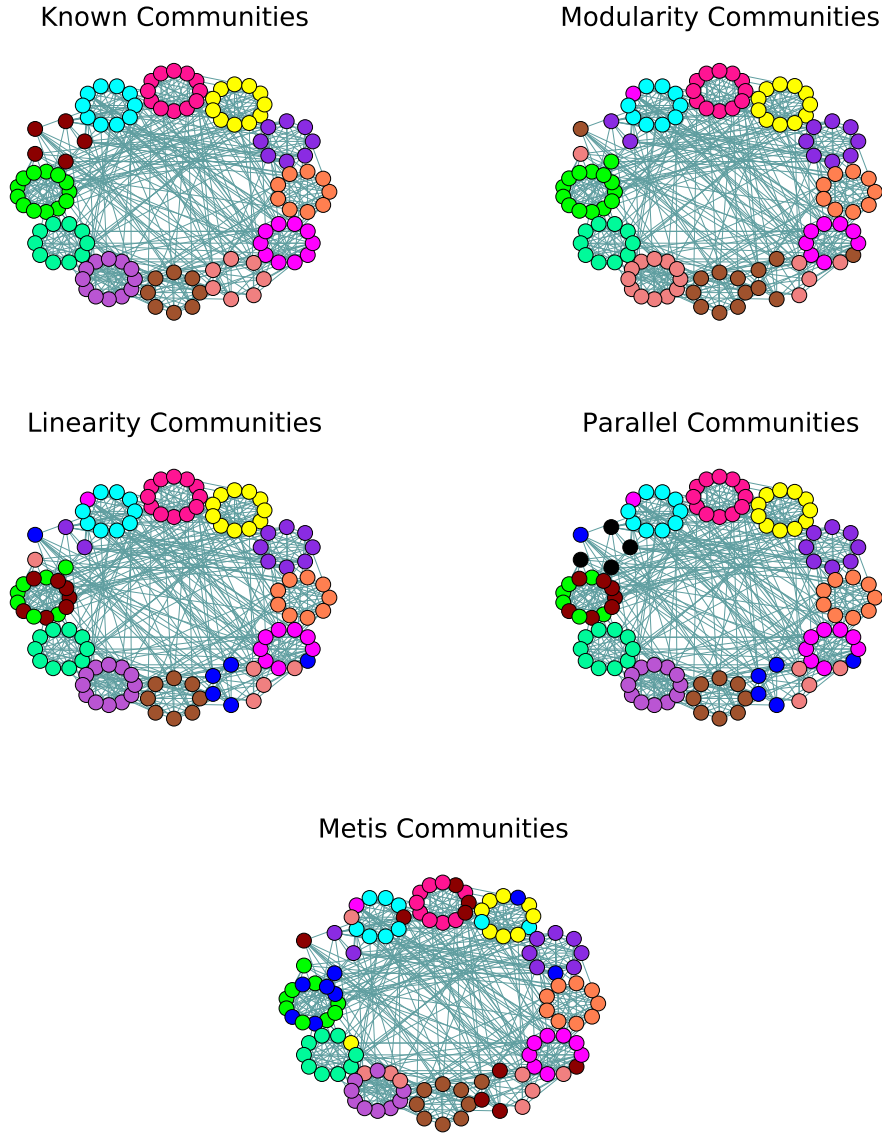


Figure 5.2: Different solutions produced by the different community detection methods, communities are marked by coloration.

fessor. We eliminate these nodes in a way that maintains the structure of the graph.

Definition 23 (Bridge) *A bridge is an edge that replaces a low degree node connecting two high degree nodes. If the graph has edges (i, j) and (j, k) , but not (i, k) and i, k are*

nodes with degree greater than d and j is a node with degree less than or equal d , we replace node j with a bridge connecting i and k .

JTODO include diagram

There are no known solutions to the collaboration networks, but we can compare the communities returned by each method. Our comparison reveals the resolution limit of modularity. Modularity returns fewer and larger communities than our linearity and parallel methods. In particular, we can see how modularity returns communities that deviate from our understanding of a strong community. Linearity and Parallel illuminate different community structure, much closer to our understanding of strong communities.

5.2.1 General Relativity

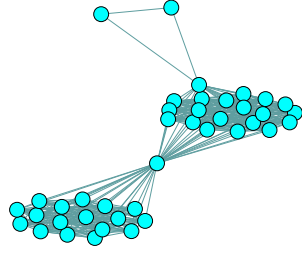
Modularity produces 34 communities in the general relativity graph. It is the smallest of the coauthor networks, but we begin to see the resolution limit of modularity.

Communities From	$I(S)$	$E(S)$	$ S $	Avg C Diameter	Avg NCP
Linearity	0.33	0.05	235	3.99	0.03
Parallel	0.18	0.07	408	4.94	0.05
Modularity	0.06	0.11	34	7.15	0.09
Metis	0.19	0.50	117	4.57	0.00006

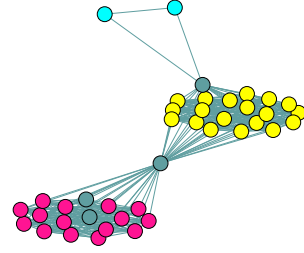
Table 5.3: Internal Density, External Density, and number of communities for the set of communities returned by each detection method on the astrophysics coauthor network.

We now provide a sample of communities found. Figure ?? has a striking example of

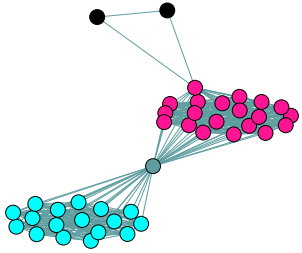
Modularity Communities



Linearity Communities



Parallel Communities



Metis Communities

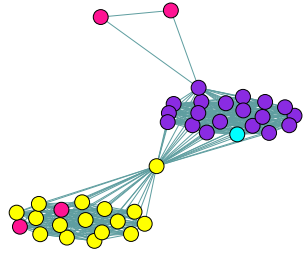


Figure 5.3: One of the smaller communities produced by Modularity with 40 nodes. Grey nodes are nodes shared by communities and black nodes are nodes not within a community. Note the resolution limit of modularity.

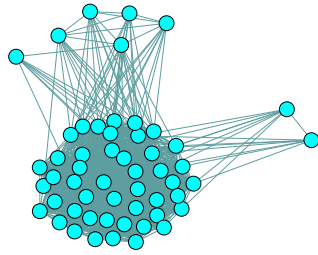
5.2.2 Condensed Matter

Condensed Matter is the largest of the coauthor networks

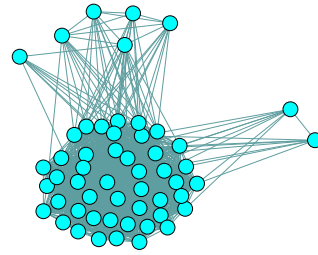
Communities From	$I(S)$	$E(S)$	$ S $
Linearity	0.18	0.11	626
Parallel	0.18	0.39	825
Modularity	0.01	0.25	52

Table 5.4: Internal Density, External Density, and number of communities for the set of communities returned by each detection method on the astrophysics coauthor network.

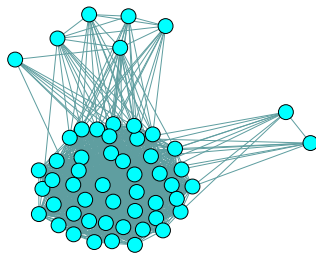
Modularity Communities



Linearity Communities



Parallel Communities



Metis Communities

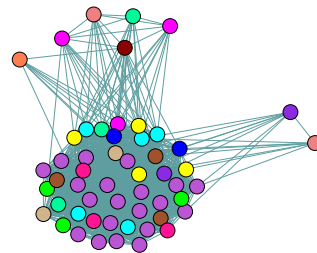


Figure 5.4: The largest community produced by Linearity at 54 nodes.

5.3 Physics ArchiveX

5.4 Enron Email Network

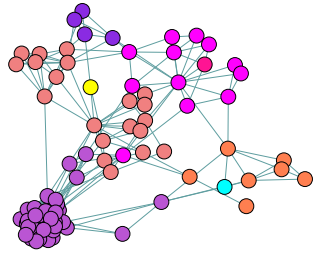
5.5 Epinions Social Network

5.6 Gnutella P2P Network

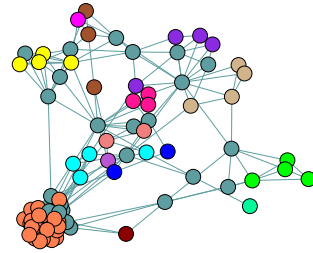
5.7 Web Graphs

5.7.1 Berkeley Webpage

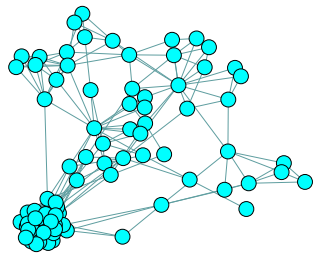
Modularity Communities



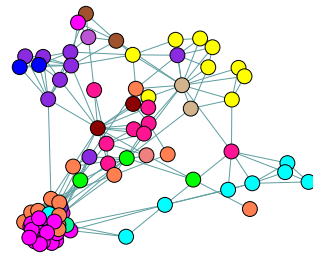
Linearity Communities



Parallel Communities



Metis Communities



Parallel Subcommunities

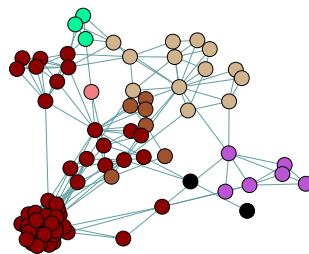
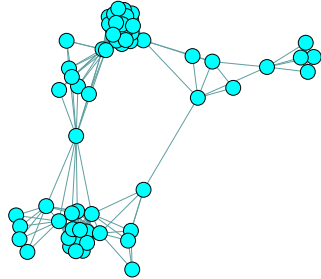
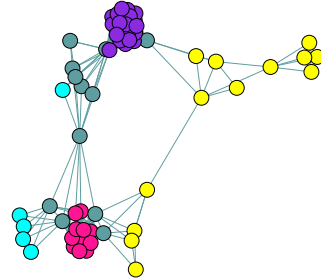


Figure 5.5: The largest community produced by Parallel within 10 forces at 94 nodes. Five subcommunities were created by parallel.

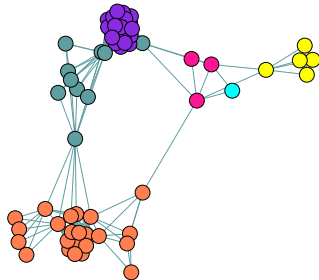
Modularity Communities



Linearity Communities



Parallel Communities



Metis Communities

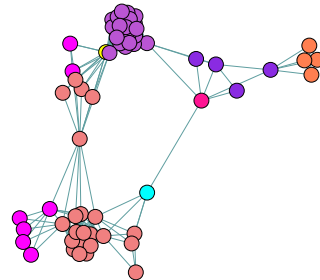
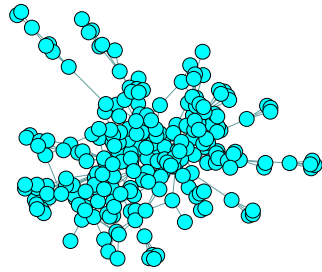
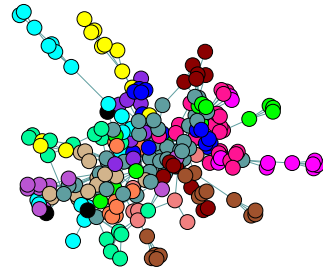


Figure 5.6: The median sized community for Modularity with 63 nodes. Grey nodes are nodes shared by communities and black nodes are nodes not within a community. Note the resolution limit of modularity.

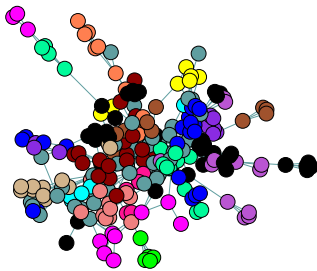
Modularity Communities



Linearity Communities



Parallel Communities



Metis Communities

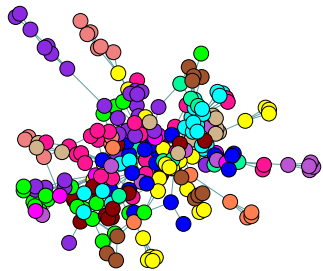


Figure 5.7: The largest community produced by Modularity with 234 nodes. The most discernable difference is that the whiskers connected by two or fewer edges are broken off.

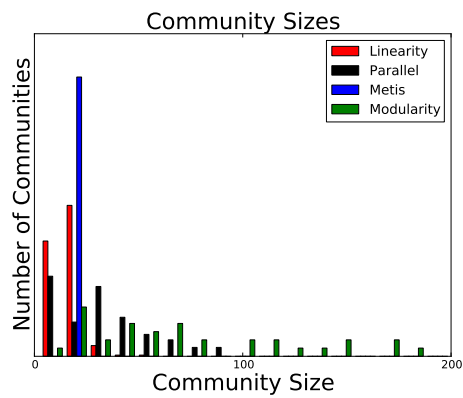


Figure 5.8: The distribution of community sizes.

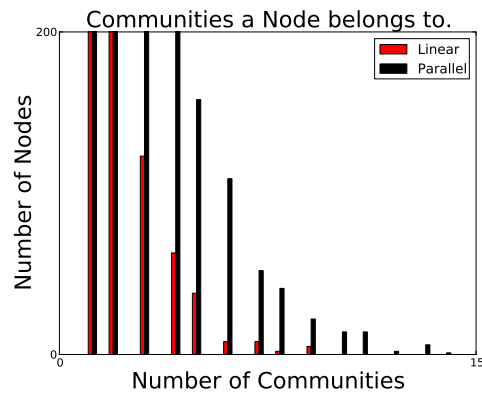
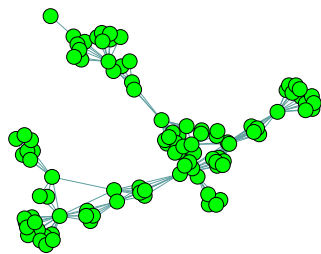
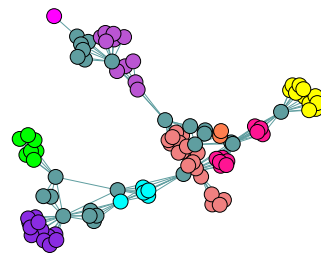


Figure 5.9: The number of communities a node belongs to. Follows a power law distribution. Comes from the degree distribution

Modularity Communities



Linearity Communities



Parallel Communities

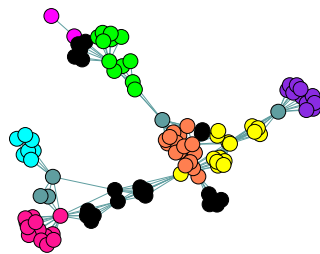


Figure 5.10: One of the smaller communities produced by Modularity with 103 nodes. Grey nodes are nodes shared by communities and black nodes are nodes not within a community. Note the resolution limit of modularity.

CHAPTER 6

CONCLUSIONS

This thesis makes three contributions. The first is a framework for comparing metric based community detection. The other two contributions are fast community detection algorithms. The advances to community detection provided by each are:

- LINEARITY
 - Provides a reliable detection method to recover large and small communities, Chapter ??
- PARALLEL
 - Has near perfect scalability to analyze enormous networks, Chapter ??
 - Recovers sets of communities with complex overlapping patterns, Chapter ??

When these algorithms are applied to applications, in Chapter ??, we find another set of conclusions:

- A power law distribution represents the number of communities a node belongs to.
- Communities in citation networks are created by the union of previous topics.
- If everyone votes, the results of elections change [JTODO].

BIBLIOGRAPHY

- [1] R. Andersen, F. Chung, and K. Lang. Local graph partitioning using pagerank vectors. 2006.
- [2] V. Blondel, J. Guillaume, R. Lambiotte, and E. Lefebvre. Fast unfolding of communities in large networks. *J. Stat. Mech.*, JTOD0, 2008.
- [3] A. Cappocci, V.D.P. Servedio, G. Calarelli, and F. Colaiori. Detecting communities in large networks. *Physica A*, 352:669–676, 2005.
- [4] D. Chen, Y. Fu, and M. Shang. A fast and efficient heuristic algorithm for detecting community structures in complex networks. *Physica A*, 388:2741–2749, 2009.
- [5] L. Danon, A. Diaz-Guilera, J. Duch, and A. Arenas. Comparing community structure identification. *Journal of Statistical Mechanics: Theory and Experiment*, JTOD0, 2005.
- [6] J. Duch and A. Arenas. Community detection in complex networks using extremal optimization. *Physical Review E*, 72(027104 JTOD0), 2005.
- [7] S. Fortunato. Community detection in graphs. *Physics Reports*, 486:75–174, November 2010.
- [8] M. Girvan and M. E. J. Newman. Community structure in social and biological networks. *PNAS*, 99(12):7821–7826, June 2002.
- [9] M. B. Hastings. Community detection as an inference problem. *Archive JTOD0*, 2006.
- [10] P. Hui, E. Yoneki, S. Chan, and J. Crowcroft. Distributed community detection in delay tolerant networks. Kyoto, Japan, August 2007. MobiArch.
- [11] A. Jain. Data clustering: 50 years beyond k-means. 2008.
- [12] R. Kannan, S. Vempala, and A. Vetta. On clusterings: Good, bad, and spectral. *Journal of the ACM*, 51(3):497–515, May 2004.
- [13] A. Lancichinetti and S. Fortunato. Community detection algorithms: A comparative analysis. *Physical Review E*, 80(056117 JTOD0), 2009.

- [14] A. Lancichinetti, S. Fortunato, and F. Radicchi. Benchmark graphs for testing community detection algorithms. *Physical Review E*, 78(046110 JTODO), 2008.
- [15] A. Lancichinetti, M. Kivela, J. Saramaki, and S. Fortunato. Characterizing the community structure of complex networks. *PloS ONE*, 5, August 2010.
- [16] J. Leskovec, K. Lang, A. Dasgupta, and M. Mahoney. Statistical properties of community structure in large social and informatino networks. Beijing, China, April 2008. WWW 2008.
- [17] J. Leskovec, K. Lang, and M. Mahoney. Empirical comparison of algorithms for network community detection. Raleigh, NC, April 2010. WWW 2010.
- [18] A. Maiya and T. Berger-Wolf. Sampling community structure. Raleigh, NC, April 2010. WWW 2010.
- [19] N. Mishra, R. Schreiber, I. Stanton, and R. Tarjan. Finding strongly-knit clusters in social networks. *Internet Mathematics*, 5:155–174, 2009.
- [20] M. Newman. Modularity and community structure in networks. *PNAS*, 103(23):8577–8582, 2006.
- [21] M. Porter, J. Onnela, and P. Mucha. Communities in networks. *Notices of the AMS*, 56(9):1082–1097, October 2009.
- [22] W. W. Zachary. An information flow model for conflict and fission in small groups. *Journal of Anthropological Research*, 33:452–473, 1977.
- [23] S. Zhang, R. Wang, and X. Zhang. Identification of overlapping community structure in complex networks using fuzzy c-means clustering. *Physica A*, 374:483–490, 2007.