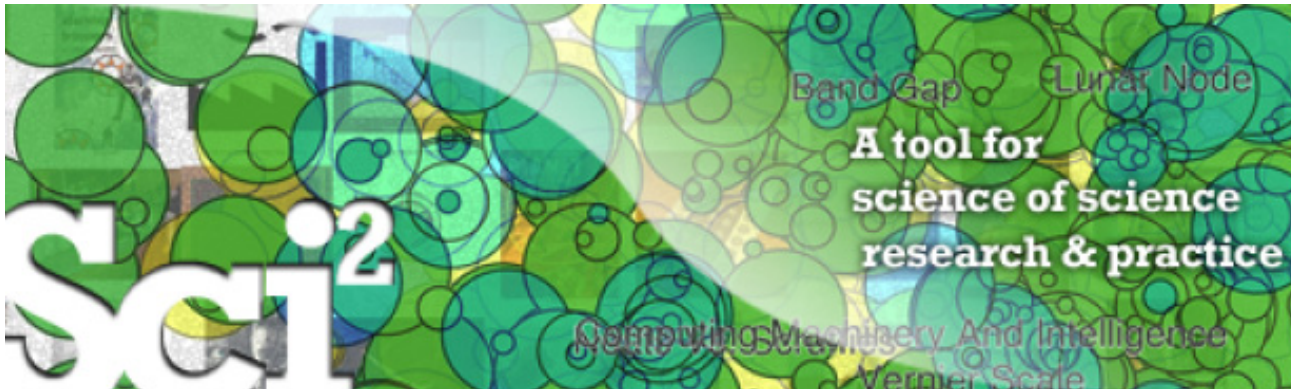# Visualizing Data with the Science of Science (Sci2) Tool



Ted Polley
Research & Editorial Assistant
Cyberinfrastructure for Network Science Center
School of Informatics and Computing
Indiana University Bloomington
http://cns.iu.edu

1

# Presentation Overview

Why should we visualize?
How should we visualize?
Introduction to Sci2
- Introduction
- Macroscopes
- OSGi & Cyberinfrastructure Shell
- Types and levels of analysis
- File formats supported by Sci2
- User Interface
- Supported tools
- Visualizations
- Sci2 Adoption

Brainstorming Session
Break for lunch

Hands-on with Sci2
- Installing Sci2
- Needs Driven Workflow Design
- Introduction to Networks
- Visualizing the Florentine Dataset
- Evolving Co-Authorship Networks
- Visualize Geographic Distribution of Clients
- Visualize Client Brand Networks

Discussion/Questions

# Why should we Visualize?

Humans are surprisingly good at:

- Pattern matching – especially when it comes to identifying trends, gaps, and outliers. We can use this to predict the future.

- Determining placement, orientation, shape, size, color etc.

Ultimately, visualizations provide access to insight much quicker than simply examining the raw data, or even data that has been statistically analyzed. This insight can lead to informed decision making.

Adapted from Noah Iliinsky's Keynote address at the OCLC Symposium, *Four Pillars of Data Visualization*, at ACRL 2013 in Indianapolis, Indiana.

3

# Why should we Visualize?

# How should we Visualize?

Step 1. **Purpose**: identify a clear purpose and focus for the visualization

Step 2. **Content**: obtain the best possible data that suits the purpose

Step 3. **Structure**: identify the best structure for the visualization that suits the purpose

Step 4. **Format**: identity the best way to represent the data so that it is easily digested

Adapted from Noah Iliinsky's Keynote address at the OCLC Symposium, *Four Pillars of Data Visualization*, at ACRL 2013 in Indianapolis, Indiana.

# Introduction to Sci2

The Science of Science (Sci2) Tool is an open-source modular toolset originally designed for the study of science. However it has many uses that support temporal, geospatial, topical, and network analysis and visualization of scholarly datasets.

# Macroscopes

Decision making in science, industry, and politics, as well as in daily life, requires that we make sense of the massive amounts of data that result from complex systems.

Rather than making things larger or smaller, **macroscopes let us observe what is too great, slow, or complex for us to comprehend or sometimes even notice.**

**Microscopes**          **Telescopes**          **Macroscopes**

# Plug-and-Play Macroscopes

While microscopes and telescopes are physical instruments, macroscopes are **continuously changing bundles of software plugins**

Macroscopes make it easy to

- Simply drop plugins into the tool and they appear in the menu, ready to use
- Sharing algorithm components, tools, or novel interfaces becomes as easy as sharing images on Flickr or videos on YouTube

# OSGi & Cyberinfrastructure Shell (CIShell)

- CIShell (http://cishell.org) is an open source software specification for the integration and utilization of datasets, algorithms, and tools
- It extends the Open Services Gateway Initiative (OSGi) (http://osgi.org), a standardized, modularized service platform
- CIShell provides "sockets" into which algorithms, tools, and datasets can be plugged using a wizard-driven process

# OSGi & Cyberinfrastructure Shell (CIShell)

**Developers**

**Users**

Alg

Alg

Alg

Tool

Tool

*CIShell Wizards*

**CIShell**

*Sci2 Tool*

Workflow

Workflow

*NWB Tool*

Workflow

Workflow

# Type of Analysis vs. Level of Analysis

| | Micro/Individual (1-100 records) | Meso/Local (101–10,000 records) | Macro/Global (10,000 < records) |
|---|---|---|---|
| **Statistical Analysis/Profiling** | Individual person and their expertise profiles | Larger labs, centers, universities, research domains, or states | All of NSF, NIH, USDA, all of science |
| **Temporal Analysis (When)** | Funding portfolio of one individual | Mapping topic bursts in 20 years of PNAS | 113 Years of Physics Research |
| **Geospatial Analysis (Where)** | Career trajectory of one individual | Mapping a state's intellectual landscape | PNAS |
| **Topical Analysis (What)** | | flows in research | VxOrd/Topic maps NIH funding |
| **Network Analysis (With Whom?)** | NSF Co-PI network of one investigator | network | NIH's core competency |

# Sci2 Tool – Supported Data Formats

**Input:**
- Network Formats
- GraphML (*.xml or *.graphml)
- XGMML (*.xml)
- Pajek .NET (*.net)
- NWB (*.nwb)
- Scientometric Formats
- ISI (*.isi)
- Bibtex (*.bib)
- Endnote Export Format (*.enw)
- Scopus csv (*.scopus)
- NSF csv (*.nsf)
- Other Formats
- Pajek Matrix (*.mat)
- TreeML (*.xml)
- Edgelist (*.edge)
- CSV (*.csv)

**Output:**
- Network File Formats
- GraphML (*.xml or *.graphml)
- Pajek .MAT (*.mat)
- Pajek .NET (*.net)
- NWB (*.nwb)
- XGMML (*.xml)
- CSV (*.csv)
- JPEG (*.jpg)
- PDF (*.pdf)
- PostScript (*.ps)

# User Interface

# Supported Tools





**Gnuplot**

portable command-line driven interactive data and function plotting utility
http://www.gnuplot.info/.

**GUESS**

exploratory data analysis and visualization tool for graphs and networks.

https://nwb.slis.indiana.edu/community/?n=VisualizeData.GUESS.

14

# Supported Tools

Adding more layout algorithms and network visualization interactivity

via Cytoscape http://www.cytoscape.org.

Simply add *org.textrend.visualization.cytoscape_0.0.3.jar* into your /plugin directory.

Restart Sci$^2$ Tool

Cytoscape now shows in the Visualization Menu



Select a network in Data Manager, run Cytoscape and the tool will start with this
network loaded.

# Bridged Tools

R statistical tool



Gephi visualization tool

# Sci2 Visualizations: *General*

Use GnuPlot to visualize the degree distribution of a co-authorship network extracted from ISI data…

# Sci2 Tool Visualizations: *Temporal*

Use Temporal Bar Graph to visualize NSF funding profiles over time…

# Sci2 Tool Visualizations: *Geospatial*

Use the Proportional Symbol Map to size and color symbols proportionally to numeric data, in this case the 20 most populated cities around the world…

**Geospatial Visualization (Proportional Symbol Map)**
Generated from 20 most populous cities
May 02, 2012 | 06:13:38 PM EDT

**Legend**
Interior Color (Linear)
Population density (people per sq. km.)

1,005    5,866    10,727

Exterior Color (Logarithmic)
Area (sq. km.)

1,100    4,435    17,884

Area (Linear)
Population

32,450,000
22,125,000
11,800,000

**How to Read this Map**
This *proportional symbol map* shows 209 countries of the world using the equal-area Eckert IV projection. Each dataset record is represented by a circle centered at its geolocation. The area, interior color, and exterior color of each circle may represent numeric attribute values. Minimum and maximum data values are given in the legend.

CNS (cns.iu.edu)

19

# Sci2 Tool Visualizations: *Geospatial*

Use the Choropleth Map to color regions proportionally to numeric data, in this case the US by state population…



**Geospatial Visualization (Choropleth Map)**
Generated from U.S. state populations
May 02, 2012 | 06:13:42 PM EDT

Alaska (10% actual area)

Puerto Rico

Hawaii (50% actual area)

**Legend**
U.S. State Color (Linear)
Population

568,158    19,130,035  37,691,912

**How to Read this Map**
This *choropleth map* shows 52 U.S. states and other jurisdictions using the Albers equal-area conic projection with Alaska, Puerto Rico, and Hawaii inset. Each U.S. state may be color coded in proportion to a numerical value. Minimum and maximum data values are given in the legend.

*CNS (cns.iu.edu)*

# Sci2 Tool Visualizations: *Geospatial*

Overlay a geospatial network on a base map, in this case Albert-László Barabási and his collaborators…



Geo Map ()
Eckert IV Projection
Apr 19, 2012 | 11:14:48 AM

Created with Sci² Tool | Cyberinfrastructure for Network Science Center (http://cns.iu.edu)

# Sci2 Tool Visualizations: *Topical*

Use the Map of Science via Journals visualization a network drawn the result of mapping a dataset's journals to the underlying sub-discipline(s) those journals contain…



**Topical Visualization**

Generated from 361 Unique ISI Records
90 out of 112 records were mapped to 182 subdisciplines and 13 disciplines.
September 20, 2012 | 11:29 AM EDT

© 2008 The Regents of the University of California and SciTech Strategies.
Map updated by SciTech Strategies, OST, and CNS in 2011.

**Legend**
Circle area: Fractional record count
Unclassified = 22
Minimum = 0
Maximum = 98
Scaling factor = 0.5076673
Color: Discipline
See end of PDF for color legend.

**Area**
98.49
55.15
9.85

**How To Read This Map**
The *UCSD map of science* depicts a network of 554 subdiscipline nodes that are aggregated to 13 main disciplines of science. Each discipline has a distinct color and is labeled. Overlaid are circles, each representing all records per unique subdiscipline. Circle area is proportional to the number of fractionally assigned records. Minimum and maximum data values are given in the legend.

*CNS (cns.iu.edu)*

22

# Sci2 Tool Visualizations: *Networks*

Use GUESS to visualize networks, such as this co-authorship network extracted from ISI data…

# Sci2 Tool Visualizations: *Networks*

Use Circular Hierarchy to visualize networks with community attributes appended…

# Sci2 Tool Visualizations: *Networks*

Use the Bipartite Network visualization to create a network of authors and publication titles…

**Network Visualization**
Generated from Bipartite network from Authors and Title.2
September 20, 2012 | 11:04 AM EDT

**Authors**    **Title**

Xxxx, X    Aaaa

Bbbb

Yyyy, Y    Cccc

Dddd

Zzzz, Z    Eeee

**Legend**
Sorted by
Left side:
Alphabetical
Right side:
Alphabetical

**How To Read This Map**
This *bipartite network* shows two record types and their interconnections. Each record is represented by a labeled circle that is size coded by a numerical attribute value. Records of each type are vertically aligned and sorted, e.g., by node size or alphabetically. Links between records of different type may be weighted as represented by line thickness.

CNS (cns.iu.edu)

25

Topic co-occurrence network of the 2885 cognitive and neuroscience NSF projects funded between 2007 and 2011.

The nodes are labeled based on how the awards were tagged. The nodes are scaled by number of awards (max = 355) with a particular tag and the edges are scaled on number of co-occurrences (max =91) of those tags. The node colors differentiate the different communities of awards, which allows you to identify topic areas.



Cognitive and Neuroscience at the NSF: 2007-2011

This is … an **entirely new way of characterizing and understanding the NSF portfolio**. This is in part because this enables **analysis of the content of the awards/proposals independent of the institutional structure.** One can quickly identify ALL of the Cog/Neuro awards throughout the entire NSF portfolio – so it captures research in all of the unexpected institutional places. This method also allows one to **easily identify areas of parallel or potentially collaborative research being funded by different institutional structures** and … to identify potential areas for advancing science by facilitating collaborations.

*Leah G. Nichols, NSF*

# Questions?

# Needs-Driven Workflow Design

**DEPLOY**

Stakeholders

Validation

Interpretation

**Types and levels of analysis** determine data, algorithms & parameters, and deployment

Data

**READ**

**ANALYZE**

**VISUALIZE**

**Visually encode data**

**Overlay data**

**Select visualization type**

Graphic Variable Types

Modify reference system, add records & links

Visualization Types (reference systems)

29

# Introduction to Networks

## Undirected Networks



**Nodes:**
○

**Edges:**
___

**Node Degree:**
Number of edges connected to nodes

**Isolates:**
Nodes that are not connected to the rest of the network

**Edge Weight:**
Demonstrates relative importance of relationships

## Directed Networks



**Edge Direction:**
Directional relationship is represented by arrows

**In-Degree:**
Number of incoming edges

**Out-Degree:**
Number of outgoing edges

# Visualizing the Florentine Dataset

This example will demonstrate how to visualize data using Sci2. In this workflow we will be working with Padgett's Florentine families dataset which includes 16 different Italian families from the early 15th century. Each family is represented by a node in the network and families are connected by edges that represent either a marriage or business/lending ties. Each node (family) has several attributes: wealth (in thousands of lira), number of priorates (seats on the civic council between 1282-1344), and total ties (total number of business ties and marriages in the dataset).

"Substantively, the data include families who were locked in a struggle for political control of the city of Florence around 1430. Two factions were dominant in this struggle: one revolved around the infamous Medici family, the other around the powerful Strozzis."

More info at http://svitsrv25.epfl.ch/R-doc/library/ergm/html/florentine.html

31

# Visualizing the Florentine Dataset

```
*Nodes
id*int label*string wealth*int totalities*int priorates*int
1 "Acciaiuoli" 10 2 53
2 "Albizzi" 36 3 65
3 "Barbadori" 55 14 0
4 "Bischeri" 44 9 12
5 "Castellani" 20 18 22
6 "Ginori" 32 9 0
7 "Guadagni" 8 14 21
8 "Lamberteschi" 42 14 0
9 "Medici" 103 54 53
10 "Pazzi" 48 7 0
11 "Peruzzi" 49 32 42
12 "Pucci" 3 1 0
13 "Ridolfi" 27 4 38
14 "Salviati" 10 5 35
15 "Strozzi" 146 29 74
16 "Tornabuoni" 48 7 0
*UndirectedEdges
source*int target*int marriage*string business*string
9 1 "T" "F"
6 2 "T" "F"
7 2 "T" "F"
9 2 "T" "F"
5 3 "T" "T"
```

First, load the florentine.nwb by following *File* > *Load* >
**yoursci2directory**/*sampledata/scientometrics/endnote/florentine.nwb.*

Once you have loaded the data in Sci2, it will appear in the Data Manager.

For this workflow we will skip straight to the visualization step, since the network file that we loaded already has the attributes we are interested in visualizing (wealth, priorates, and totalities). For other datasets, you will likely need to extract networks and run some type of analysis to answer the questions in which you are interested.

To visualize this network select the file from the Data Manager and run *Visualization > Networks > GUESS.*

When the network is loaded in GUESS it will be laid out randomly.

The first step in enhancing this network visualization is to apply a different layout. For this visualization we will use the GEM layout *Layout > GEM.* You will notice that the GEM layout is random, you can run it multiple times and the network will appear slightly different each time.

The next step will be to resize the nodes based on the wealth attribute. To do this resize select the *Resize Linear* button and set the parameters to those shown below.

Next we will colorize the nodes based on priorates to add an additional dimension to this visualization.

Next we will color the edges to show the type of relationship between the families. To do this, you will need to select the *Object* edges *based on ->*, set the property to *marriage*, the operator to *==*, and the value to *T*. Next, click the *Color* button and you can select the color of your choice from the pallet that will appear at the bottom of the Graph Modifier pane.



40

You can repeat this process for the *business* property if you want to, or you can leave the edges that represent business ties the default color. In this workflow we will leave them the default color, light gray. The final step is to show all the labels. To do this, you will need to select the "Object" all nodes and the click the *Show Label* button and the labels will appear in the visualization.

Since the GEM layout is random and all the nodes are spaced more or less evenly apart, you do not have to worry about disrupting the layout. However, other layout algorithms may space the nodes according to specific attributes of the network. Manually moving around nodes in this case would disrupt the layout of the network and distort the meaning of the visualization.

The last thing we want to do to our network is color the border of the nodes the same as the nodes themselves. This is not as crucial for networks with only a few nodes, but as the size of your network increases it can become difficult to read with the thick black lines around every node. To color those the same as the node go to the *Interpreter* tab at the bottom of the GUESS window and type in the following commands:

```
for n in g.nodes:
        n.strokecolor = n.color
```

This code basically tells GUESS that for every node (*n*) in this graph of nodes (*g.nodes)* make the border color of the nodes (*n.strokecolor*) equal to the node color (*n.color)*. After you type the first line you will need to hit the "Tab" key before you start typing the second line of code.

# Questions?

# Temporal Analysis: *Evolving Co-Authorship Network*

For this analysis we will be studying the evolution of Alessandro Vespignani's co-authorship network over time. We will see his network of collaborators grow from 1990 to 2006, giving us a sense of how his scholarly output has grown. Each node in the network will represent an author in the data set and the edges that connect them will be weighted based on how many times they have collaborated.

*File > Load > AlessandroVespignani.isi*
Load this file from the sample data folder you copied from the flash drive.

Select *Preprocessing > Temporal > Slice Table by Time* and choose the parameters shown at the right.

# Evolving Co-Authorship Network

Now that the algorithm has been run, you will notice the original dataset has been divided into four tables that cumulatively display the evolution of this data.



49

Select the first table and run *Data Preparation > Extract Co-Author Network*

Repeat this step for each of the tables in the Data Manager

Select the first extracted co-author network and run *Visualization > Networks > GUESS* starting with the network that spans 1990 to 2006 because we will export these node positions and use them to layout the other networks.

The network will be loaded in with a random layout in GUESS
To change the layout select *Layout > Gem*

Resize the nodes based on *number_of_authored_works*
Set the parameters from 5 to 15 and click *Do Resize Linear*

Resize the edges based on *number_of_coauthored_works*
Set the parameters from 1 to 5 and click *Do Resize Linear*

Colorize the nodes based on *times_cited*
Set the parameters from *Gray* to *Black* and click *Do Colorize*

Colorize the edges based on *number_of_coauthored_works*
Set the parameters from *Green* to *Black* and click *Do Colorize*

If you want to remove the borders from the nodes, type the following commands in into the interpreter:

```
for n in g.nodes:
        n.strokecolor=n.color
```

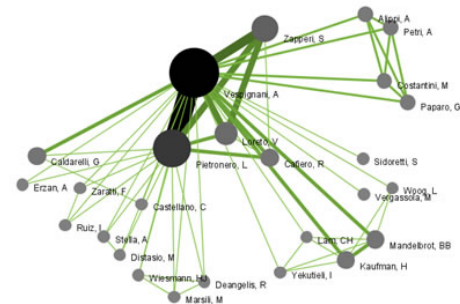Finally add the labels to the nodes by selecting object: *all nodes* and then click Show Label

To save the node positions of the current layout so that the layout is consistent across all time slices select *File > Export Node Positions* and save the file as a CSV file.
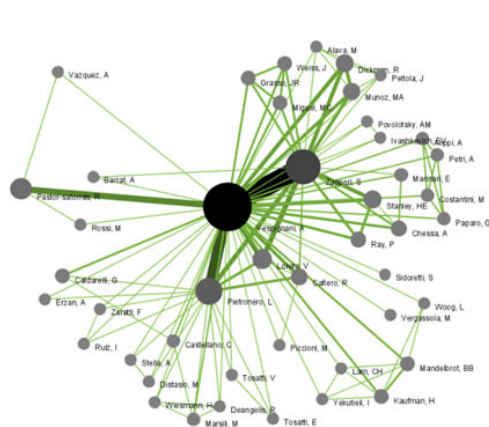
Now when you go to visualize the other three networks you will want to import the node the node positions using *File > Import Node Positions* and the network will be laid out accordingly. When the networks are displayed side-by-side you can see an evolution.
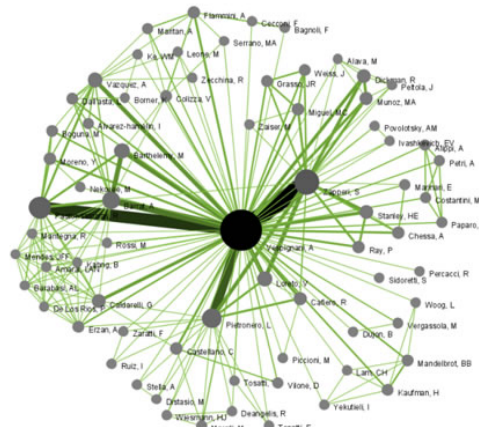


1990-1991

1990-1996

1990-2001
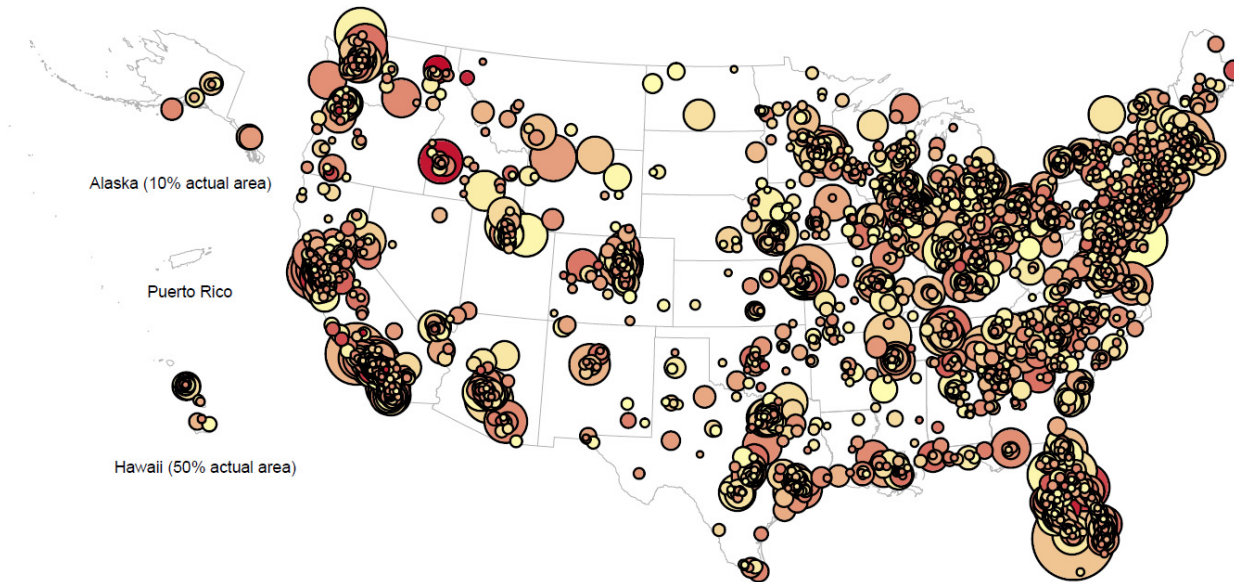
1990-2006

# Questions?

# Geospatial Analysis:
# *Proportional Symbol Map of Clients*

For this analysis we will be exploring the geographic distribution of clients and encoding various attributes of the data, such as equity and tenure as a TD Ameritrade client. The sample data has been provided to you prior to the workshop.
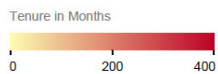
# How much is too much?



**Geospatial Visualization (Proportional Symbol Map)**

Geographic Distribution of Clients
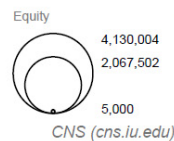Sep 20, 2013 | 11:01:54 AM EDT

Alaska (10% actual area)

Puerto Rico

Hawaii (50% actual area)

**Legend**
Interior Color (Linear)

Tenure in Months

0       200       400
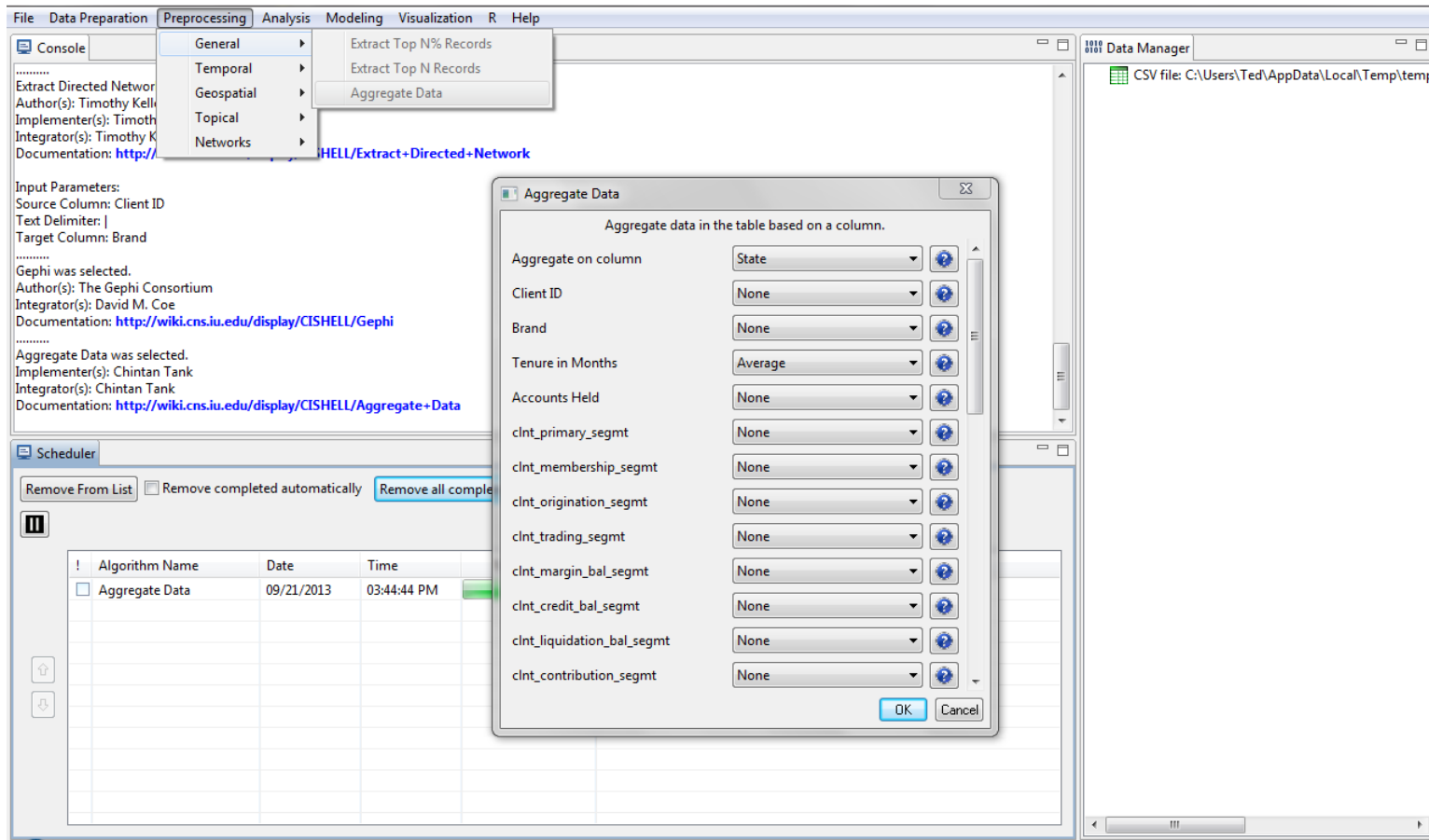
**Area (Linear)**
Equity

4,130,004
2,067,502

5,000

CNS (cns.iu.edu)

**How to Read this Map**
This *proportional symbol map* shows 52 U.S. states and other jurisdictions using the Albers equal-area conic projection with Alaska, Puerto Rico, and Hawaii inset. Each dataset record is represented by a circle centered at its geolocation. The area, interior color, and exterior color of each circle may represent numeric attribute values. Minimum and maximum data values are given in the legend.

# Aggregate clients by state

*Preprocessing > General > Aggregate Data* and set the Aggregate on Column to State, find the Average Tenure in Months and Equity for the clients in these states.
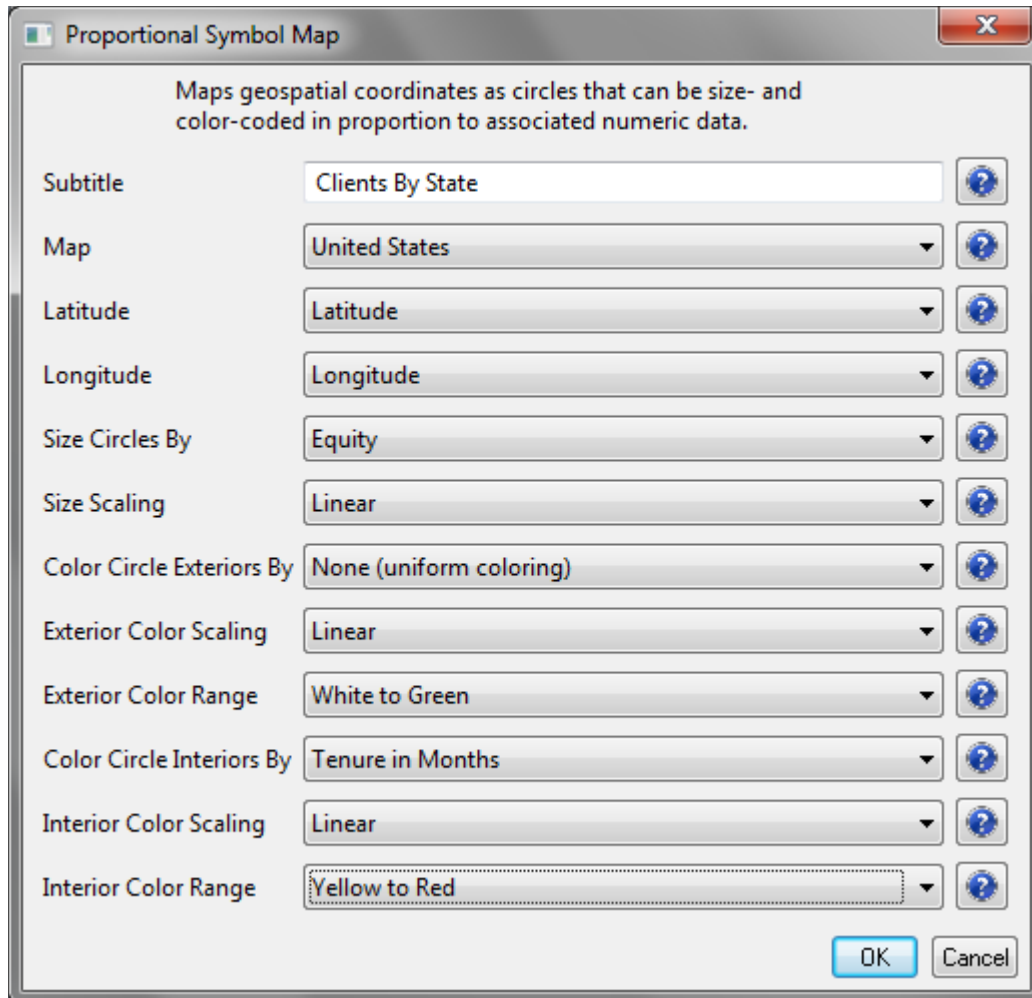
The result is a simplified dataset with the each state present. The *Count* column gives the total number of clients in each state for this dataset. The *Tenure in Months* column gives the average tenure in months for those clients, and the *Equity* column gives the average equity for clients in those states.

| | A | B | C | D |
|---|---|---|---|---|
| 1 | Tenure in Months | Equity | State | Count |
| 2 | 106 | 165412 | HI | 27 |
| 3 | 98 | 145186 | FL | 308 |
| 4 | 95 | 128056 | TX | 322 |
| 5 | 104 | 111266 | LA | 41 |
| 6 | 105 | 81466 | AL | 55 |
| 7 | 59 | 70413 | MS | 9 |
| 8 | 99 | 112742 | GA | 114 |
| 9 | 96 | 158025 | AZ | 117 |
| 10 | 92 | 150764 | SC | 56 |
| 11 | 110 | 137150 | CA | 891 |
| 12 | 93 | 116514 | NM | 24 |
| 13 | 97 | 194911 | NC | 116 |
| 14 | 111 | 153475 | OK | 35 |
| 15 | 77 | 62068 | AR | 19 |
| 16 | 99 | 191937 | TN | 54 |
| 17 | 112 | 141320 | NV | 32 |

Next, select *Analysis >  Geospatial > Generic Geocoder*

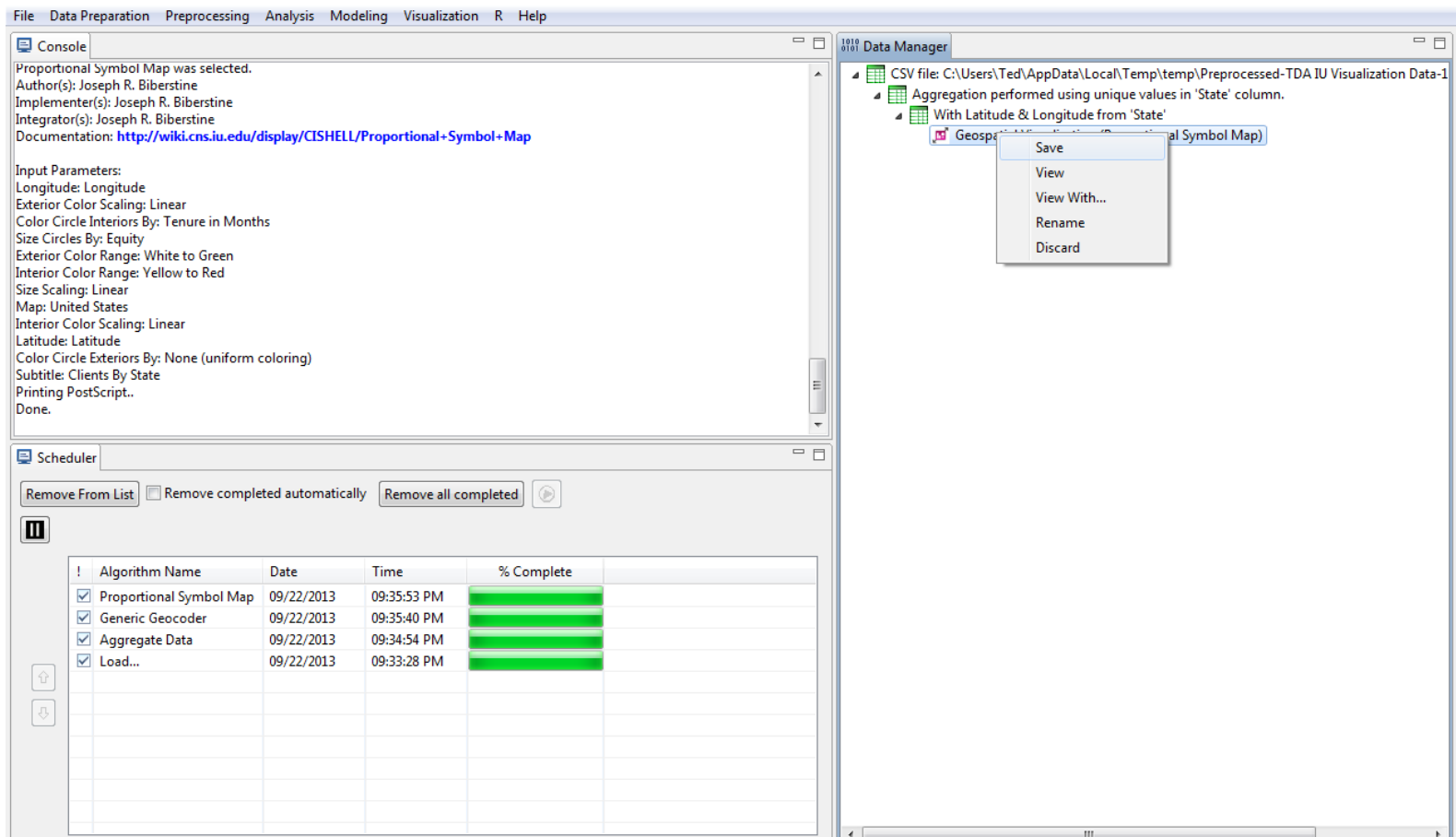Next, select *Visualization > Geospatial > Proportional Symbol Map*

The result is a PostScript file in the Data Manager. Right-click on the file and select *Save*.
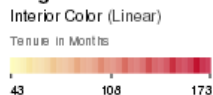
You will have to convert the PostScript file to a PDF to view it. You can use Adobe Distiller to convert, or an online service, such as http://ps2pdf.com.



**Geospatial Visualization (Proportional Symbol Map)**
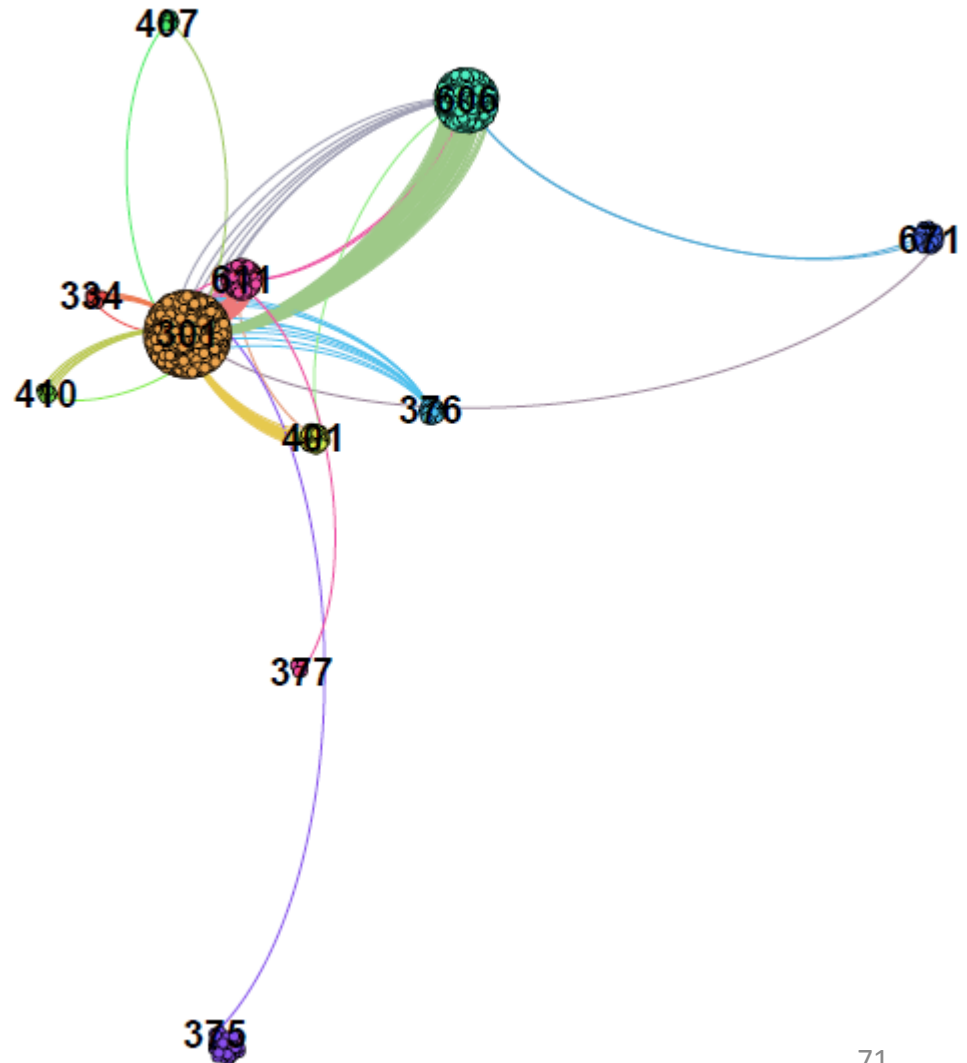Clients By State
Sep 22, 2013 | 09:40:24 PM EDT

Alaska (10% actual area)

Puerto Rico

Hawaii (50% actual area)

**Legend**
Interior Color (Linear)
Tenure in Months
43    108    173

Area (Linear)
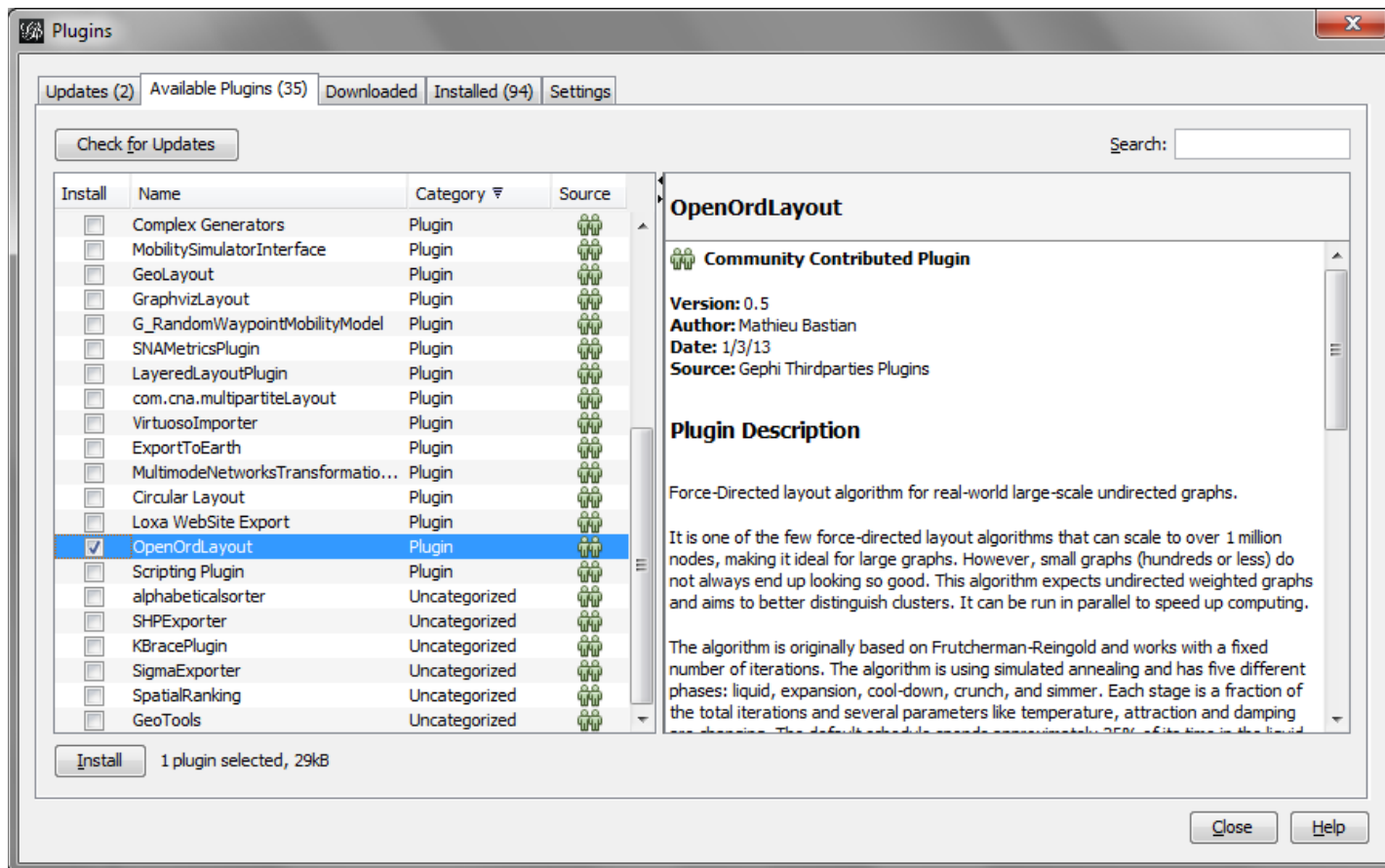Equity
336,988
191,133
45,278

**How to Read this Map**
This *proportional symbol map* shows 52 U.S. states and other jurisdictions using the Albers equal-area conic projection with Alaska, Puerto Rico, and Hawaii inset. Each dataset record is represented by a circle centered at its geolocation. The area, interior color, and exterior color of each circle may represent numeric attribute values. Minimum and maximum data values are given in the legend.
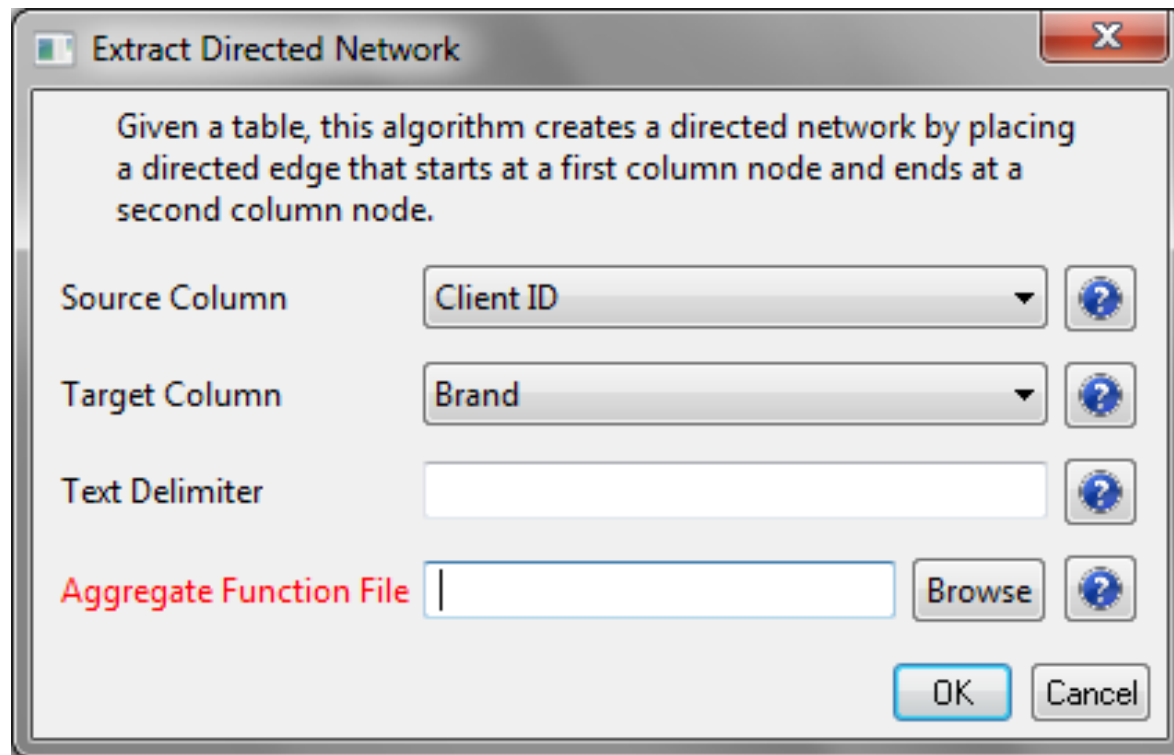
CNS (cns.iu.edu)

70

Extract a directed network showing the connections between clients and brands.

This network visualization requires the OpenOrd Network Layout in Gephi. It is available as a plugin. Simply open Gephi and select, *Tools > Plugins.* Then select the Available Plugins tab and find the OpenOrdLayout plugin. Install this plugin and restart Gephi.

Select, *Data Preparation > Extract Directed Network*

# Visualize the network with Gephi, *Visualization> Networks > Gephi*

# Gephi will layout the nodes randomly

Select the OpenOrd layout. This layout is extremely efficient for large networks. Set the Edge Cut parameter to 0.0. This will result in a slightly less clustered result.
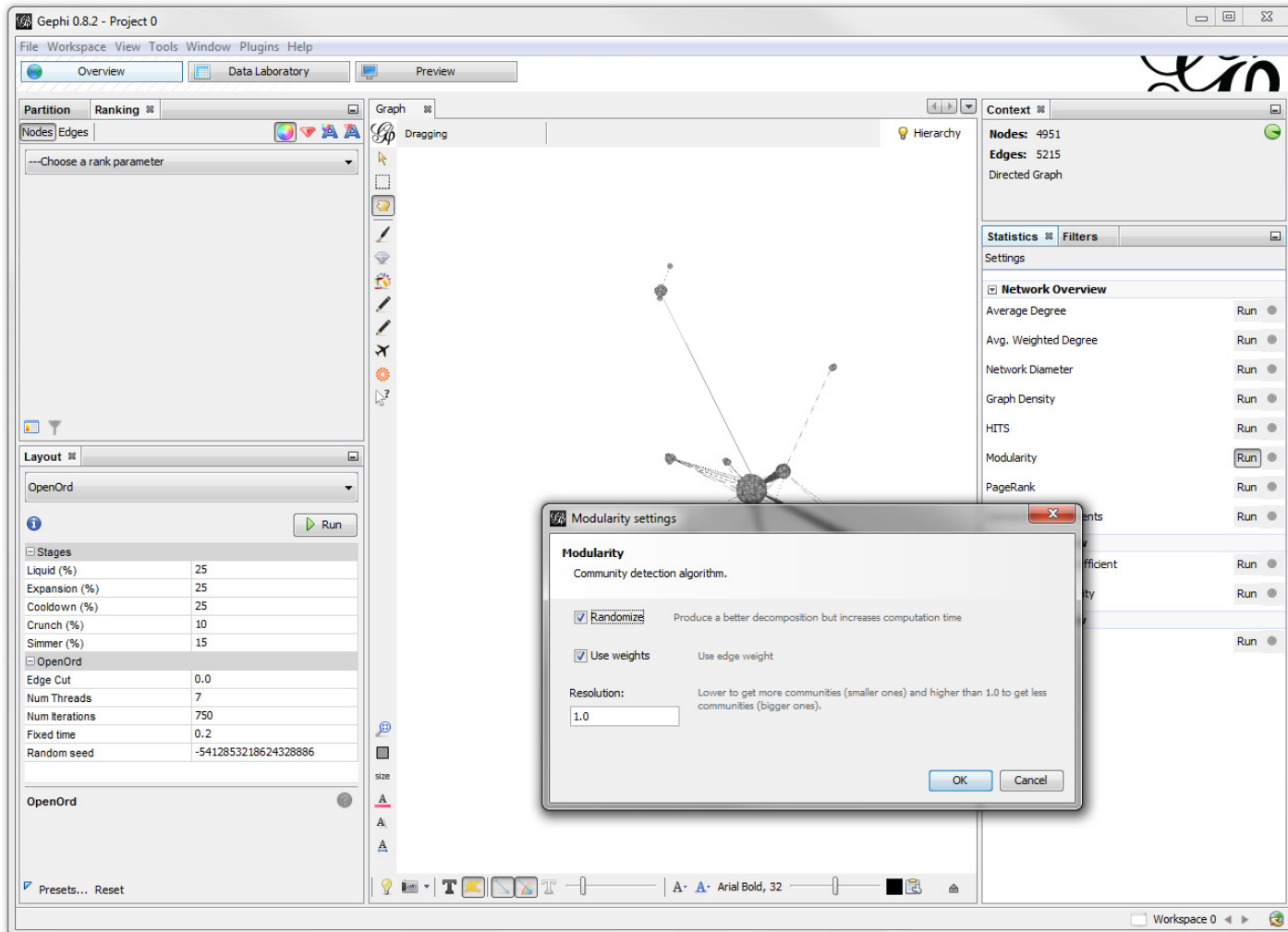
In the statistics tab, run Modularity. This identifies communities of clients based on the brands with which they are associated.

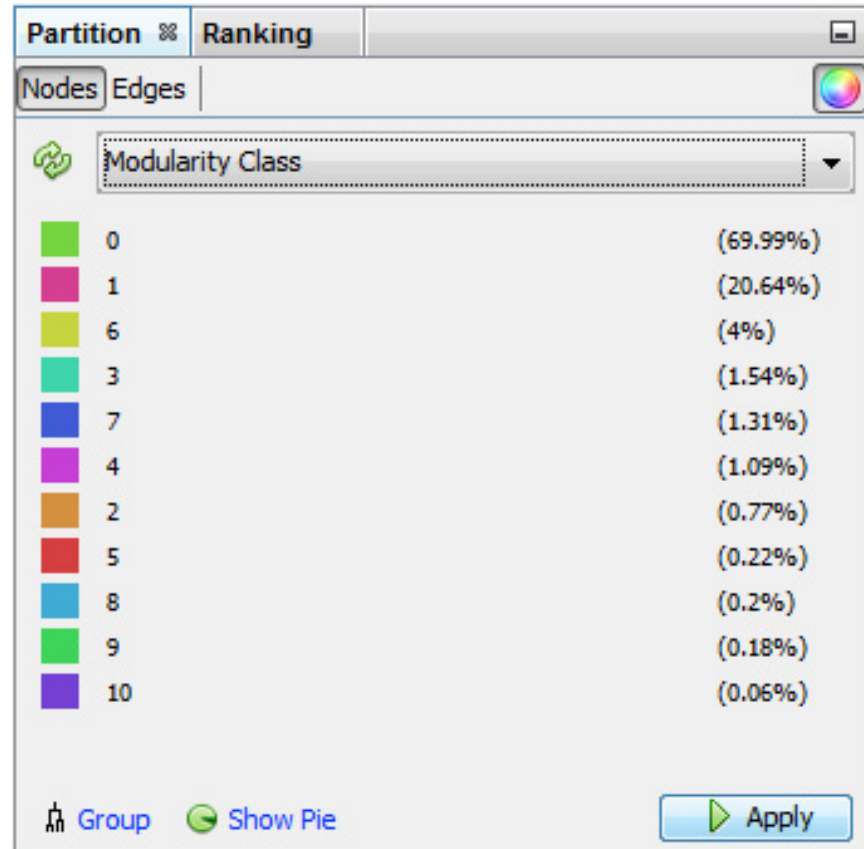Use the partition tab in the upper right-hand corner. Select Modularity Class. You may have to click refresh. Select Apply and the colors shown will be applied to each modularity class.

This presents a network where the nodes are distinct clusters based on brand. You can determine brand popularity based on cluster size. Each brand cluster will have a distinct color. There are also a lot of connections between brand clusters.

# Brainstorming/Questions?