# CASOS

# Network Regression
# and
# Visualizing Distributions

## Jeff Reminga

The CASOS Center
COS Program, School of Computer Science, Carnegie Mellon
Summer Institute 2020

**Carnegie Mellon**

**Center for Computational Analysis of
Social and Organizational Systems
http://www.casos.cs.cmu.edu/**

---

**Carnegie Mellon**

# Agenda

- Standard Regression
  - Show how to run a regression using node-level measure and attribute values
  - Node-level means one value per node, a vector of values
  - This is standard, textbook regression
- Network Regression
  - Show how to run a regression using link-level data
  - Link-level means one value per link, a network of values
  - This is network regression
  - QAP Techniques deal with the dependence of the link values (violation of the independence of observation assumption)
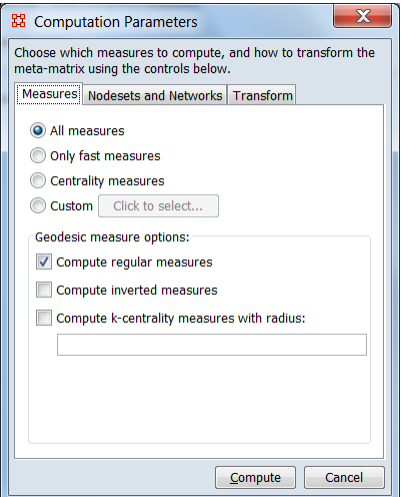
<Your Name>

## Slide 3

# Standard Regression

**Computation Parameters**

Choose which measures to compute, and how to transform the meta-matrix using the controls below.

Measures | Nodesets and Networks | Transform

- ○ All measures
- ○ Only fast measures
- ○ Centrality measures
- ○ Custom    Click to select...

Geodesic measure options:
- ☑ Compute regular measures
- ☐ Compute inverted measures
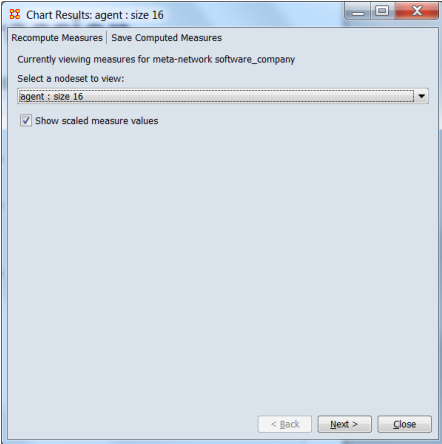- ☐ Compute k-centrality measures with radius:

Compute    Cancel

- Load the "softco.xml" dataset
- Select the softco meta-network in the Meta-Network Manager
- Click the main menu item Visualizations\Measure Charts (or Measure Charts button)
- Select the measures to compute: All, Fast, Centrality, or custom
- Geodesic measures (betweenness, closeness, path lengths) have options:
  - Use Inverted measures when the link values are dissimilarities (such as distance)
  - Use k-centrality for large data: as only considering neighbors within distance k
  - Inverted must be used needed when the links mean "matter of interpretation of link

## Slide 4

# Standard Regression...

**Chart Results: agent : size 16**

Recompute Measures | Save Computed Measures

Currently viewing measures for meta-network software_company

Select a nodeset to view:

agent : size 16

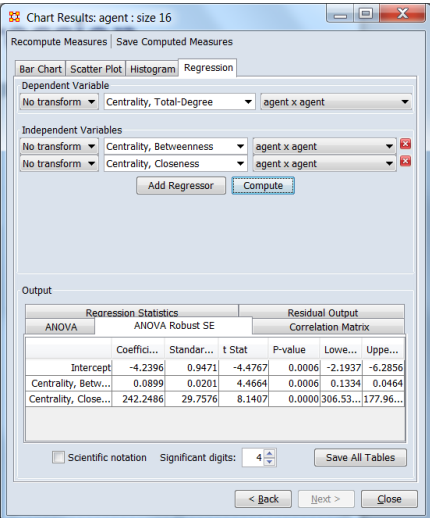☑ Show scaled measure values

< Back    Next >    Close

- After computing measures, you will see this dialog where you select which nodeset values to view
- Click whether to use scaled or unscaled values
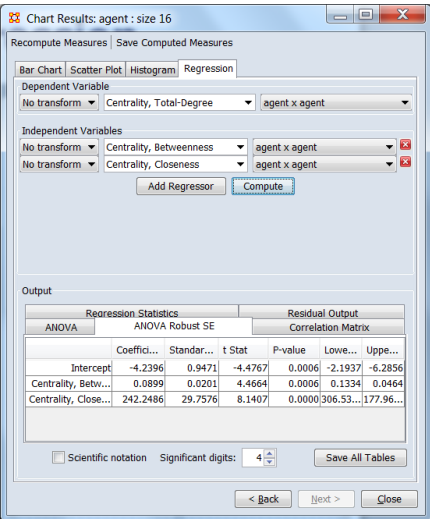- Click next

<Your Name>



Standard Regression…

- Choose the Regression tab at the top
- The independent variable and the dependent variables are all node-level measures or numeric attributes
- In this regression, the dependent variable is the measure Total Degree Centrality computed on the Agent x Agent network
- The dependent variables are Betweenness and Closeness Centrality also computed on the Agent x Agent network.
- Regression results are print below
- Use the Save All Tables to save values into a CSV file



Standard Regression…

- Regression results are print below
- The results indicate that Total Degree Centrality is determined more by Closeness Centrality than Betweenness Centrality
- This makes sense because an increase in closeness would entail an increase in degree.
- Use the Save All Tables to save values into a CSV file

<Your Name>

**Carnegie Mellon**
ISr institute for SOFTWARE RESEARCH
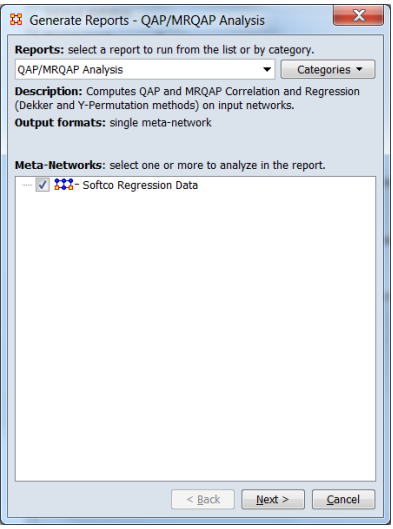
# Network Regression

- In Standard Regression the unit of analysis is a node, and node-level measure and numeric attribute data was used
- In Network Regression the unit of analysis is the dyad (a pair of nodes), and ORA lets you input three types of data:
    - A network directly, such as an Agent x Agent network
    - A vector of node-level numeric attributes (such the Age of each node) repeated by row or column to form a network
    - A vector of node-level measure values (such as Betweeenness Centrality) repeated by row or column to form a network
- Mathematically, the same calculations are used as those in node-level regression
- The link values are "stretched out" into a vector and then input into the standard regression routines
- Special techniques are used to handle the lack of independence of observations in the link-level observations, namely, Multiple Regression Quadratic Assignment Procedure (MRQAP)

CASOS

---

**Carnegie Mellon**
ISr institute for SOFTWARE RESEARCH

# Network Regression...

**Generate Reports - QAP/MRQAP Analysis**

**Reports:** select a report to run from the list or by category.

QAP/MRQAP Analysis    |    Categories ▼

**Description:** Computes QAP and MRQAP Correlation and Regression (Dekker and Y-Permutation methods) on input networks.
**Output formats:** single meta-network

**Meta-Networks:** select one or more to analyze in the report.

☑ Softco Regression Data

< Back    Next >    Cancel

- As an example, load the "softco-regression.xml" dataset
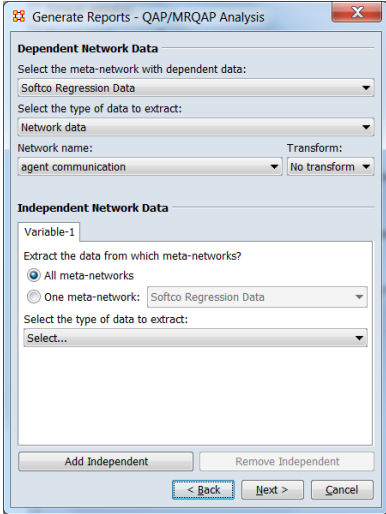- Click Generate Reports
- Select the QAP/MRQAP Analysis report

CASOS

CASOS

‹Your Name›

## Slide 11

**Carnegie Mellon**
isr institute for SOFTWARE RESEARCH

# Network Regression...

**Generate Reports - QAP/MRQAP Analysis**

Select the algorithms to run:
- ☑ Correlation
- ☑ Y-Permutation Regression
- ☑ Double-Dekker Semi-Partialling Regression

Set the algorithm parameters:
- Random seed value: 0
- Number of permutations: 1,000

☐ Return regression network data to the main interface?

`< Back`  `Next >`  `Cancel`

- After clicking Next, we have some options to select
- Correlation will compute the correlation between the dependent and independent variables
- Y-Permutation Regression will compute regression using a QAP technique involving permuting the independent network
- Double-Dekker Semi-Partialling Regression will compute the regression using another technique to deal with the inherent dependence of link data
- Correlation and Regression will use a bootstrapping technique that inputs a random seed and number of iterations
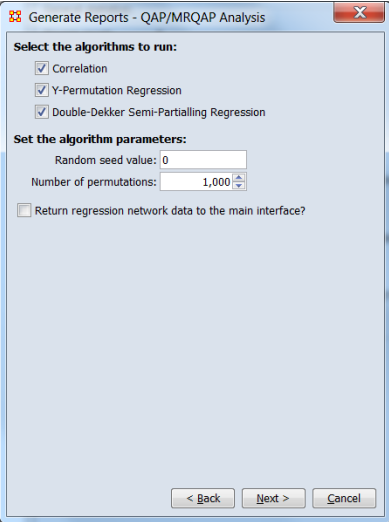
June 2020     © 2020 CASOS, Director Kathleen M. Carley     11

## Slide 12

**Carnegie Mellon**
isr institute for SOFTWARE RESEARCH

# Network Regression...

**Correlation (Dependent to Independent)**

This shows the correlation and related statistics between the dependent network variable and each independent network variable.

At least one input network has non-binary link values, and therefore the Euclidean distance was computed.

| Variable Name | Variable Meta-Network | Variable Description | Correlation | Significance | Euclidean Distance |
|---|---|---|---|---|---|
| X1 | Softco Regression Data | Network: agent shared knowledge | 0.153 | 0.077 | 28.671 |
| X2 | Softco Regression Data | Network: agent shared tasks | 0.221 | 0.013 | 18.762 |

- Click Next, choose an output file, and click Finish
- This is the Correlation output from the report
- Note that Agent Shared Tasks is more highly correlated than shared Knowledge, but with less significance
- Euclidean distance is used because the independent networks are both valued

June 2020     © 2020 CASOS, Director Kathleen M. Carley     12

<Your Name>

## Slide 13

**Carnegie Mellon**
ISR institute for SOFTWARE RESEARCH

# Network Regression...

**Regression Results**

Reports the results from the regression. There are three computations for standard errors: the classical formula is reported in column Std.Errors; heteroskedasticity robust standard errors are reported in column Robust Std.Errors; finally, bootstrapped standard errors are reported in column Bootstrapped Std.Errors.
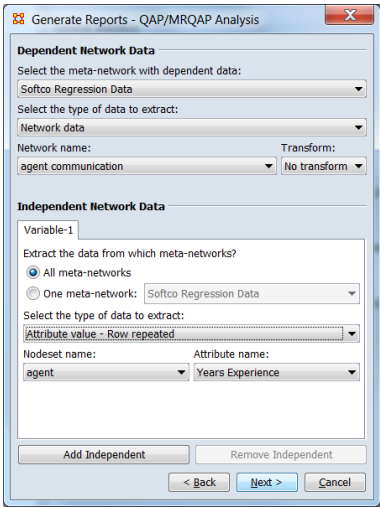
| R-Squared (R2) | 0.052 |
| Residual Sum Of Squares | 48.531 |
| Total Sum Of Squares | 51.183 |
| Standard Error | 0.453 |

| Variable Name | Variable Meta-network | Variable Description | Coef | Std.Coef | Std.Errors | Robust Std.Errors | Bootstrapped Std.Errors | Sig.Y-Perm | Sig.Dekker |
|---|---|---|---|---|---|---|---|---|---|
| Intercept | | | 0.204 | | 0.072 | | | | |
| X1 | Softco Regression Data | Network: agent shared knowledge | 0.024 | 0.060 | 0.029 | 0.030 | 0.045 | 0.288 | 0.303 |
| X2 | Softco Regression Data | Network: agent shared tasks | 0.079 | 0.192 | 0.030 | 0.032 | 0.042 | 0.036 | 0.042 |

- The Regression results show as well that Shared Tasks has a greater effect on communication than Shared Knowledge
- Again, however, the significance of Shared Tasks using both using Y-Permutation and Double Dekker Semi-Partialling is less

**CASOS**

## Slide 14

**Carnegie Mellon**
ISR institute for SOFTWARE RESEARCH

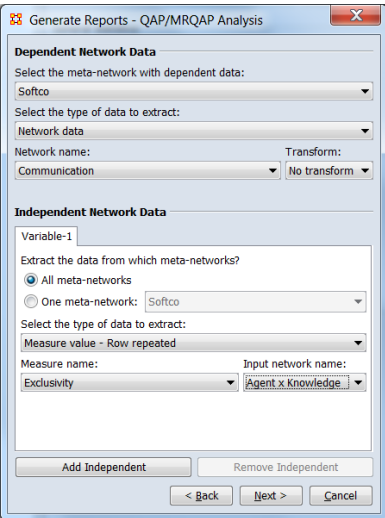# Network Regression: Attributes



- Another example is to assess the degree to which a communication between two agents is explained by the years experience of the actors.

- Use *Attribute value - Row Repeated* to compare the link value (a,b) with the attribute value of node b

- Use *Attribute value – Column Repeated* to compare the link value (a,b) with the attribute value of a

- Select the Independent variable as shown to use the Years Experience attribute

**CASOS**

**CASOS**

<Your Name>

## Slide 15

# Network Regression: Measures

**Generate Reports - QAP/MRQAP Analysis**

**Dependent Network Data**

Select the meta-network with dependent data:
Softco

Select the type of data to extract:
Network data

Network name:      Transform:
Communication      No transform

**Independent Network Data**

Variable-1

Extract the data from which meta-networks?
◉ All meta-networks
◯ One meta-network:   Softco

Select the type of data to extract:
Measure value - Row repeated

Measure name:      Input network name:
Exclusivity      Agent x Knowledge

Add Independent      Remove Independent

< Back   Next >   Cancel

- Another example, using the "softco.xml" dataset, is to assess the degree to which a communication from agent A → B is explained by exclusive access of B to knowledge.

- That is, to what extent there is a communication from A → B so that A can access the knowledge of B.

- Use *Measure value - Row Repeated* to compare the link value (a,b) with the attribute value of node b

- Select the Independent variable as shown to use the Exclusivity measure computed on Agent x Knowledge

## Slide 16

# Visualizing Distributions

- Nodesets can have node attributes
- The distribution of node attribute values can be shown for categorical number, categorical text, or numerical data types
- Networks have link values whose distribution can be shown

<Your Name>

## Slide 19

**Carnegie Mellon**
ISR institute for SOFTWARE RESEARCH

# Node Attribute Distribution –
# Categorical Text values

Now click the down arrow for the US Resident column



The values are shown as a frequency table because US Resident is text category variable.

The data type of an attribute can be changed using the **Edit tab**
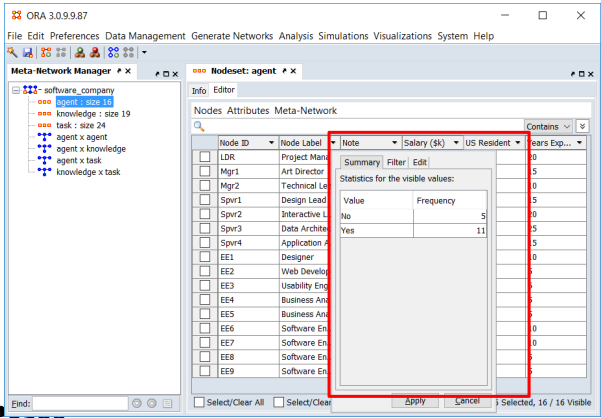(Summary | Filter | Edit)

## Slide 20

**Carnegie Mellon**
ISR institute for SOFTWARE RESEARCH

# Node Attribute Distribution -
# Charts



- The Nodeset Editor also contains a charting feature to show the distribution of values
- Go to the nodeset editor's main menu: Attributes and select "Charts" (at bottom)
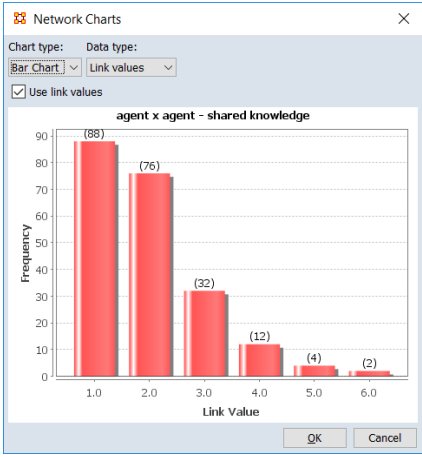- Shown is the bar chart for Years Experience

**Visualizing Network Value Distributions**

- The network editor has a similar chart tool to display the distribution of row sum values, column sum values, and links values

- This is located under the Network Editor's main menu: Nodes \ Charts

- Shown is the distribution of link values for the shared knowledge (folded agent x knowledge) network