# Prediction and Comparison of Two or More Networks: Hamming Distance, Correlation, QAP, MRQAP

## CASOS

Ramon Villa-Cox

rvillaco@andrew.cmu.edu

School of Computer Science, Carnegie Mellon

Summer Institute 2020

**Carnegie Mellon**

institute for SOFTWARE RESEARCH

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

---

**Carnegie Mellon**
institute for SOFTWARE RESEARCH

# Motivation

- How can we compare 2 different networks?
  - Famous work by Bernard and Killworth

- Fraternity Dataset
  - 58 Nodes (Frat Members)
  - 2 Different Networks
  - Number of interactions between students
    - Seen by unobtrusive observer
    - BKFRAB in ORA
  - Rank of perceived interaction
    - Surveyed from participants
    - BKFRAC in ORA

<Your Name>

**Carnegie Mellon**
**isr** institute for SOFTWARE RESEARCH

# Motivation

## How similar is the cognitive network to the behavioral network?

Lets load the data and check in ORA

**CASOS**

9 June 2020                                                                3

---

**Carnegie Mellon**
**isr** institute for SOFTWARE RESEARCH

# First Attempt

- Visualize the networks
  - They look different
  - Doesn't tell us much more than we already knew

- Cut links less than the mean
  - They look more different
  - Still hard to tell

- Lesson: visual tools help, but actual differences are hard to define from visuals

**CASOS**

9 June 2020                                                                4

<Your Name>

**Carnegie Mellon**
isr institute for SOFTWARE RESEARCH

# How do we compare networks?

- That is, given two networks, what should we do to understand their similarities and differences?
- "Tools"
  - Visual analysis, Metrics, Statistics
- "Approaches"
  - Node level metrics, network level metrics, motifs, network structure

CASOS
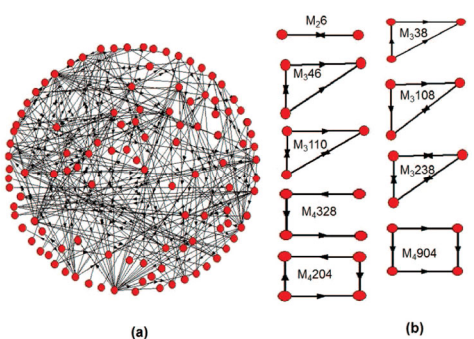
9 June 2020      5

---

**Carnegie Mellon**
isr institute for SOFTWARE RESEARCH

# What is a motif?

- Partial subgraph
  - Introduced by Uri Alon

- Also called local patterns

- Compare how frequently they occur to occurrence in random network
  - Over representation shows that it is an important characteristic of the network



(a)      (b)

Image From "Identification of Important Nodes in Directed Biological Networks: A Network Motif Approach" Wang, Lu, and Yu

CASOS

9 June 2020      6

CASOS

<Your Name>

**Carnegie Mellon**
institute for
SOFTWARE
RESEARCH

# Motifs in ORA

- Measure Charts

- All Measures

- Clique Count

- Doesn't work for fully-connected weighted graph!
  - Have to binarize first

CASOS

9 June 2020                                                                                              7

---

**Carnegie Mellon**
institute for
SOFTWARE
RESEARCH

# Motifs in ORA



CASOS

9 June 2020                                                                                              8

CASOS

<Your Name>

**Carnegie Mellon**
ISR institute for SOFTWARE RESEARCH

# Comparing Network Structures

- We can compare networks more generally by looking at its structure

- Specifically, we look at the structure of its adjacency matrix

- Compute distance metrics between adjacency matrices
  - Hamming Distance
  - Euclidean Distance

- Use Correlations

**CASOS**

**Carnegie Mellon**
ISR institute for SOFTWARE RESEARCH
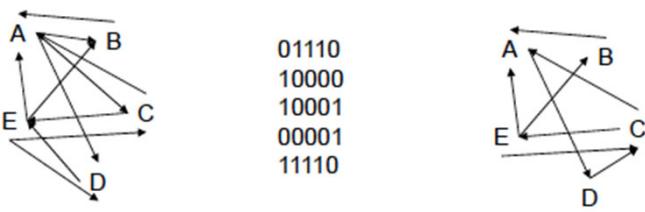
# Hamming Distance

- Data assumed to be binary string (list of 0's and 1's)

- How many digits need to be flipped in A to obtain B?
  - Or vice versa
  - Formally: $d_h = \sum_i |A_i - B_i|$
  - Could also apply the above to weighted data

- Normalization bounds distance from 0 to 1
  - Number of non-diagonal spaces in an adjacency matrix: N*(N-1)
    - N = number of nodes

- Normalized formula: $\widehat{d_h} = \frac{1}{N*(N-1)} \sum_i |A_i - B_i|$

**CASOS**

**CASOS**

**Carnegie Mellon**
ISr institute for SOFTWARE RESEARCH

# Example



01110
10000
10001
00001
11110

00010
10000
10001
00100
11100

3) String

4) Calculate

Distance = 5
5/20, .25, 25%

0111010000100010000111110
0001010000100010010011100

**CASOS**

9 June 2020                                                                11

---

**Carnegie Mellon**
ISr institute for SOFTWARE RESEARCH

# Euclidean Distance

- The distance metric most people are familiar with

- Assumes Euclidean space
  - Normal space (straight dimensions with orthogonal axis)
  - Not necessarily true for networks

- Definition: $d_E = \sqrt{\sum_i (A_i - B_i)^2}$

- Note: in the binary case: $d_E = \sqrt{d_h}$

- Not bounded

**CASOS**

9 June 2020                                                                12

## Correlation

- Correlation measures the strength of relationship between two things
  - In our case: links occurring / not occurring in different networks

- Definition: $r = \dfrac{\sum_i (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_i (x_i - \bar{x})^2} * \sqrt{\sum_i (y_i - \bar{y})^2}}$

- Bounded from -1, 1

- Values far from 0 indicate strong relationship

- Negative values indicate inverse relationship

**CASOS**

## Regression

- These concepts are very closely related to regression

- Regression assumes that one variable (dependent) is a function of another variable (independent)

- The function is then found by estimating the conditional expectation

- For networks: is one network a function of another network?
  - Is the perceived friendship network a function of the actual contact network?

**CASOS**

<Your Name>

Carnegie Mellon
isr institute for SOFTWARE RESEARCH

# Thinking about distances

- Original motivation: how similar are these networks?

- Now we can put a number on it
  - Allows us to say which networks are more/less similar

- But how do we know these numbers matter?

- Use statistics!
  - Could use a bootstrapped t-test, for example

- **What makes this hard for networks?**

CASOS

9 June 2020                                                                 15

---

Carnegie Mellon
isr institute for SOFTWARE RESEARCH

# The problem with regression/correlation

- Regression
  - Y: friendship network
  - X: knowledge homophily network



- Naïve approach
  - Write networks as vectors
  - Run OLS on vectors

CASOS

9 June 2020                                                                 16

**Carnegie Mellon**
ISr institute for SOFTWARE RESEARCH

# The problem with regression/correlation

- Regression
  - Y: friendship network
  - X: knowledge homophily network

| Friendshi... | | | | | ...mophily | | | |
|---|---|---|---|---|---|---|---|---|
| | | | | | | .8 | .7 | .6 |
| | | | | | .8 | | .8 | 0 |
| | | | | | .7 | .8 | | 0 |
| | | | | | .6 | 0 | 0 | |



**Wrong!**
**Networks are fundamentally correlated and violate i.i.d. assumption of classical statistics**

| .9 | .8 | 0 | .9 | .7 | 0 | .8 | ... | .7 | .8 | .0 | .6 | 0 | 0 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

- Naïve approach
  - Write networks as vectors
  - Run OLS on vectors

**CASOS**

9 June 2020                                                                 17

---

**Carnegie Mellon**
ISr institute for SOFTWARE RESEARCH

# Another way of looking at this



- What is the correlation?
  - Krackhardt, 1987

- If represented as vectors, these would look very different
  - Graph isomorphism

**CASOS**

9 June 2020                                                                 18

**CASOS**

<Your Name>

## QAP: Quadratic Assignment Procedure

Carnegie Mellon
ISf institute for SOFTWARE RESEARCH

- How do we account for re-namings? QAP!

- The procedure:
  - Compute your statistic (distance, correlation, etc.)
  - Repeat for all possible namings:
    - Shuffle the node names in one of the networks
    - Re-compute your statistic
  - These recomputed samples makeup the null distribution
  - Compare your statistic to the null model
    - Can get a p-value, etc.

- Similar approach to bootstrapping

CASOS

9 June 2020                                                        19

---

## Statistical comparison – an example

Carnegie Mellon
ISf institute for SOFTWARE RESEARCH

- Let's just look at correlation between our network and a "random" network
- Process:
  - Create a new network
  - Fill it with random data
- Run the QAP/MRQAP report
  - What would you expect to see?
  - What do you see?

CASOS

9 June 2020                                                        20

CASOS

<Your Name>

**Now Lets Compare our Networks**

9 June 2020

21



**Running QAP in ORA**

Generate Reports - QAP/MRQAP Analysis

**Select Report**
Filter Data
Negative Links
Transform Data
Remove Nodes

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

QAP/MRQAP Analysis     Categories

Description  Input Requirements Output Formats

Computes QAP and MRQAP Correlation and Regression (Dekker and Y-Permutation methods) on input networks.

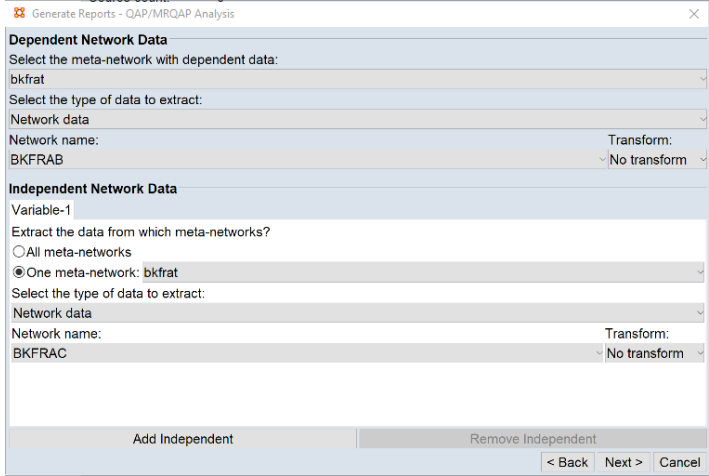**Meta-Networks:** select one or more to analyze in the report.

bkfrat

< Back    Next >    Cancel

9 June 2020

22

<Your Name>

# Running QAP in ORA

**Carnegie Mellon**
isr institute for SOFTWARE RESEARCH

**Correlation (Dependent to Independent)**

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

Significance for Pearson Correlation is the fraction of trial bootstrap values that are higher than the actual.

Significance for Hamming and Euclidean Distance is the fraction of trial bootstrap values that are lower than the actual.

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

Think of these like p-values, Similarities are significant!

| Variable | Variable Meta-Network | Variable Description | Correlation | Significance | Euclidean Distance | Significance |
|---|---|---|---|---|---|---|
| X1 | bkfrat | Network: BKFRAC | 0.370 | 0 | 191.520 | 0 |

The table below has information about how the above significance values were computed. The observed (i.e. actual) values are computed on the input data and then a number of trials are run in which the input data is permuted and the values recalculated. This creates a sequence of trial values. The statistics of these trial values are reported in the table below, and the significance is either the proportion higher or lower than the observed.

Number of trials: 1000

| Variable | Method | Trial Values | | | | | Proportion ≥ Observed | Proportion ≤ Observed |
|---|---|---|---|---|---|---|---|---|
| | | Observed | Min | Max | Average | Std.dev | | |
| X1 | Correlation | 0.370 | -0.113 | 0.110 | -8.533e-04 | 0.039 | 0 | 1 |
| X1 | Euclidean Distance | 191.520 | 208.878 | 222.675 | 215.817 | 2.422 | 1 | 0 |

CASOS

---

# Running QAP in ORA

**Carnegie Mellon**
isr institute for SOFTWARE RESEARCH

**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.

The input data has been centered and therefore the constant term in the regression will always have value zero and is not reported below.

| Model Fit | |
|---|---|
| R-Squared (R2) | 0.137 |
| Residual Sum Of Squares | 33,159.454 |
| Total Sum Of Squares | 38,421.992 |
| Standard Error | 3.168 |

| Variable | Variable Meta-network | Variable Description | Coef | Std. Coef | Std. Errors | Robust Std.Errors | Bootstrapped Std.Errors | Sig.Y-Perm |
|---|---|---|---|---|---|---|---|---|
| X0 | bkfrat | Network: BKFRAC | 1.065 | 0.370 | 0.047 | 0.065 | 0.111 | 0 |

The table below has information about how the above significance values were computed. The observed (i.e. actual) values are computed on the input data and then a number of trials are run in which the input data is permuted and the values recalculated. This creates a sequence of trial values. The statistics of these trial values are reported in the table below, and the significance is either the proportion higher or lower than the observed.

Number of trials: 1000

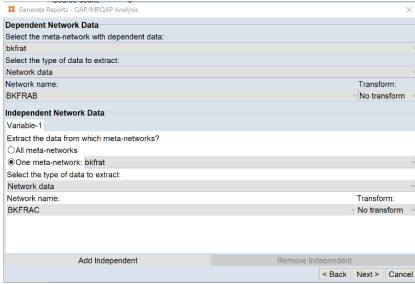| Variable | Method | Trial Values | | | | | Proportion ≥ Observed | Proportion ≤ Observed |
|---|---|---|---|---|---|---|---|---|
| | | Observed | Min | Max | Average | Std.dev | | |
| X0 | Y-Permutation | 1.065 | -0.350 | 0.330 | -0.002 | 0.109 | 0 | 1 |

CASOS

CASOS

<Your Name>

## Slide 27

**Carnegie Mellon**
ISR institute for SOFTWARE RESEARCH

# MR-QAP

- What if we want to model multiple relationships?

- Regression -> Multiple Regression

- QAP -> MR-QAP

- In ORA: "add independent" allows you to add more variables



**CASOS**

9 June 2020                                                                 27

## Slide 28

**Carnegie Mellon**
ISR institute for SOFTWARE RESEARCH

# Recap

- Networks can be compared in a variety of ways

- Motifs allow you to see/compare "building blocks" of a network

- Distances/Correlation allow you to quantitatively find differences in network structure

- To analyze distances/correlation QAP must be used
  - Due to graph isomorphism and i.i.d. samples

- Multiple regression can also be performed using MRQAP

**CASOS** Be careful with binary outcome variables!
  - Since the model is linear regression

9 June 2020                                                                 28

**CASOS**