# Moving from Data to Latent Spaces and Networks

Captain Iain Cruickshank
icruicks@Andrew.cmu.edu
Summer Institute 2020
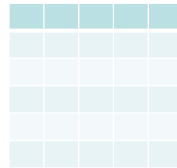
**Carnegie Mellon**

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

---

**Carnegie Mellon**

# Non-Network Sociometric Data

- What happens when we get data about entities, but not a network?
  - Often its easy to get attributes, but difficult or impossible to get relations between entities
- Also, how do we deal with complex data types, like categorical variables?
  - Categorical variables common for describing persons (i.e. 'is a smoker', 'hair type', etc.)
- We still want to analyze that data and have a flexible, accurate model of the data



**?**

June 2020

# The Workflow

| Obtain Data about entities of interest | → | Project the data into a latent space | → | Learn a graph on the data in the latent space | → | Analyze the graph to answer questions |

*The overall idea is that given some data, which may be categorical, high-dimensional, or combination thereof is to *model* that data as something which *preserves relationships* and can be *easily analyzed* (i.e. a network)

CASOS
June 2020

---

# Putting Data into a Latent Space

| Obtain Data about entities of interest | → | Project the data into a latent space | → | Learn a graph on the data in the latent space | → | Analyze the graph to answer questions |

- After collecting data, we place the data into a latent space
- We will cover Socio-Cultural Cognitive Mapping (SCM) to place data into a latent space

CASOS
June 2020

CASOS

## Overview of SCM

- Take a set of node attributes or network data and use the information to place nodes in space.
  - User defines the geometry of the space
  - User provides data

- Nodes that are highly similar will be near each other, while nodes that are quite different will be far apart.

- Overall goodness-of-fit is evaluated with a Chi-Squared Test

June 2020

## SCM Process



June 2020

**Carnegie Mellon**
isr institute for SOFTWARE RESEARCH

# SCM Model

$$F(i, j) = R_i \times C_j \times 2^{-d_{ij}^a}$$

    – Where $i$ and $j$ are entities, R and C are row and column multipliers, and the final term is an interaction term

$$d_{ij} = \left( \sum_k |x_{ik} - x_{jk}|^M \right)^{\frac{1}{M}}$$

    – d is the Minkowski distance between the entities $i$ and $j$ in the data matrix of $X$.

**CASOS**

June 2020

---

**Carnegie Mellon**
isr institute for SOFTWARE RESEARCH

# Creating a Model of the Latent Space Data

Obtain Data about entities of interest → Project the data into a latent space → **Learn a graph on the data in the latent space** → Analyze the graph to answer questions

- Now that the data has been placed into a latent space, we want to have a model of the data
- Graphs (networks) make good models of data
  - Have emergent structures
  - Interpretable
  - Allow for local heterogeneity in the data

**CASOS**

June 2020

**CASOS**

## Overview of Unsupervised Graph Learning

Carnegie Mellon
ISR institute for SOFTWARE RESEARCH

- The fundamental idea of graph learning is to find the best graph representation of some data
  - It could be considered as a way of approximating the manifold of the data
  - A recent survey of the field is available in Qiao et al. *Data-driven graph construction and graph learning: A review* and Brugere et al. *Network Structure Inference, A Survey: Motivations, Methods, and Applications*
- Used in everything from subspace learning, clustering, dimensionality reduction, manifold learning, metric learning, etc.

CASOS

June 2020

---

## k-NN Network Modularity

Carnegie Mellon
ISR institute for SOFTWARE RESEARCH

- Procedure that takes an *affinity matrix*, constructs a graph where each entity receives a connection to their $k$ nearest neighbors, and then finds subgroups via modularity maximization
- Try for several values of $k$ and pick that one which has the best modularity



CASOS

June 2020

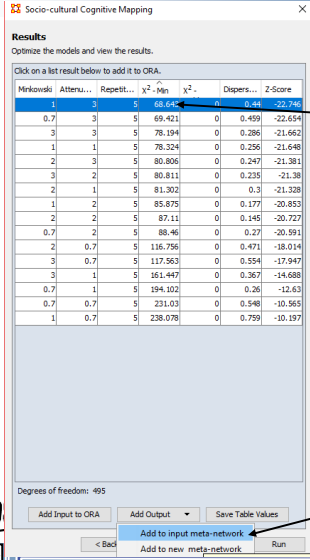**Time for an Example!**

---

Obtain Data about entities of interest → Project the data into a latent space → Learn a graph on the data in the latent space → Analyze the graph to answer questions
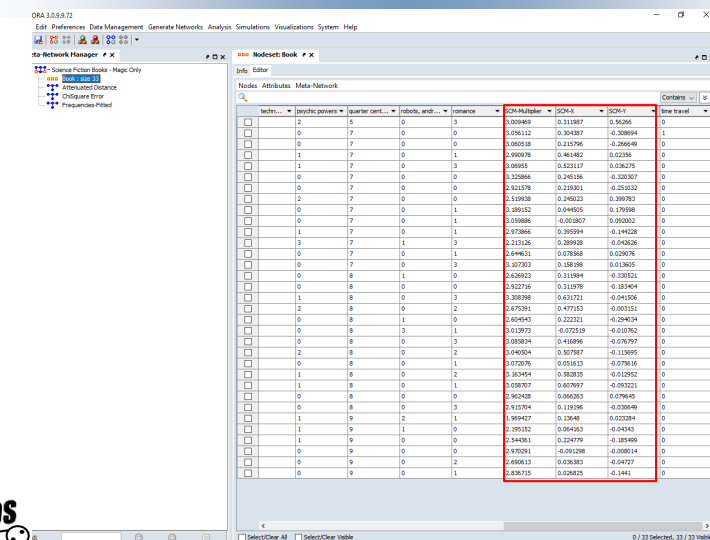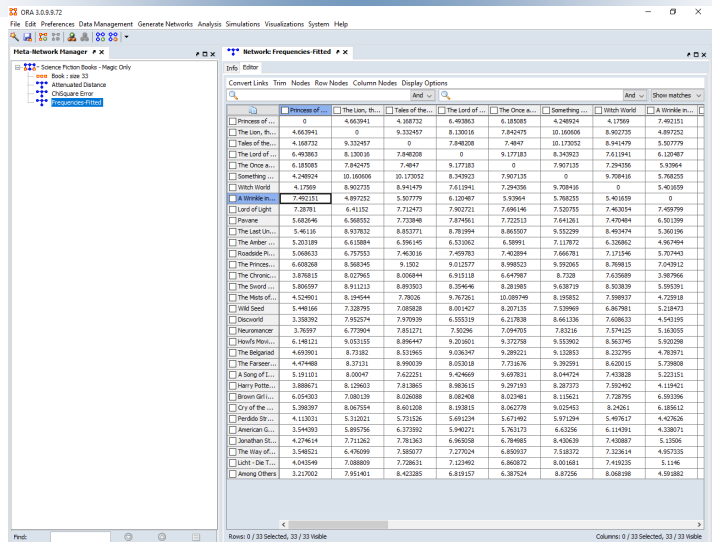
# Step 1: Find the Data

- Read in "Science Fiction Books – Magic Only.xml"
- 33 Books
- Set of Attributes for each Book

**Step 2: Start the SCM**



**Select the Frequency Data**



Can select a network or attributes of a node, which creates a frequency network

Can select different attributes and different levels of attributes

Check for mutually exclusive and redundant attributes. Generally you always want to do this to improve performance

CASOS

**Select from SCM Results**

Generally speaking, you will want to use the output which places points that generates the smallest Chi Squared Value

Finally, add your selected result to ORA (note: you can also add the frequency network input and the actual table of results, too).

June 2020



**SCM Results Meta-Network**

June 2020

SCM Results Meta-Network



Visualizing SCMs

CASOS



Go to "Multi-Dimensional Layout"



Configure the Layout

We will visualize in 2-d, since we found spatial points in 2-d

Select 'SCM-X'

Select 'SCM-Y'

Select 'Run Layout'

June 2020

See Layout!



Explore the Layout with Node Coloring: Gender

**Step 3: Learn a Graph**

Go to 'Generate Reports', 'Locate Groups', and navigate to the specific algorithm.

Go over to the 'General Options' tab

Make sure to select 'Add located groups network to the input network' (that's how we get back the best fit graph!)

June 2020



**Select the Latent Space Attributes**

Go to 'Generate Reports', 'Locate Groups', and navigate to the specific algorithm

Only select our new latent space positions, 'SCM-X' and 'SCM-Y'

Finally, run the analysis

June 2020

Step 4: Analyze the Results

Now, we have learned the best fit k-NN graph for our data, using modularity as the means of determining the goodness of fit.



Step 4: Analyze the Results

Node coloring by sub group.
Node size by degree centrality

Hatfields and McCoys, based on historical documentation



8th Ukrainian Parliament, based on votes

Network of Sakula virus samples, based on binary attributes of the code

## Recap

- In research we often get data that may be complex and have uncertain relationships
- We can deal with the data by creating an analyzable, flexible and interpretable model of that data through the presented procedure
  - Place the data in a latent space
  - Learn a graph on the data
  - Analyze the graph
- Graph-based models of data can be used for many, many different types of data