Complex Networks 2022

# Node Vector Distances: Methods and Applications

Michele Coscia
ITU København
November 7th, 2022
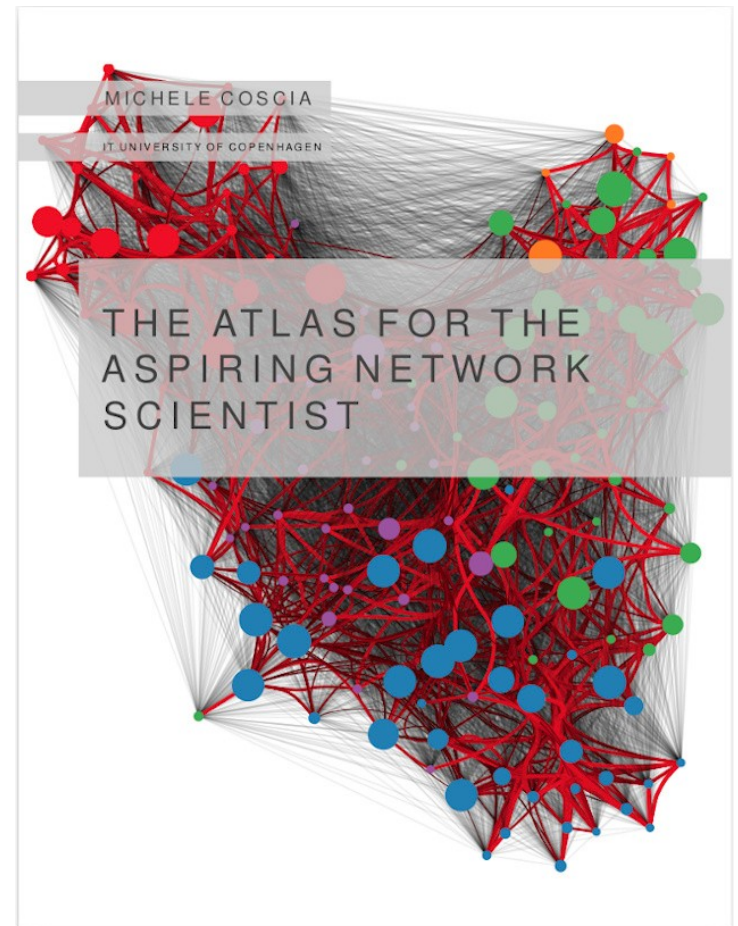
IT UNIVERSITY OF COPENHAGEN

# Download the Material

https://github.com/mikk-c/
complexnetworks22-tutorial

# Readings

- The Atlas for the Aspiring Network Scientist

  – Chapter 40

  – https://www.networkatlas.eu/

- The node vector distance problem in complex networks, ACM Computing Surveys (CSUR) 53 (6), 1-27
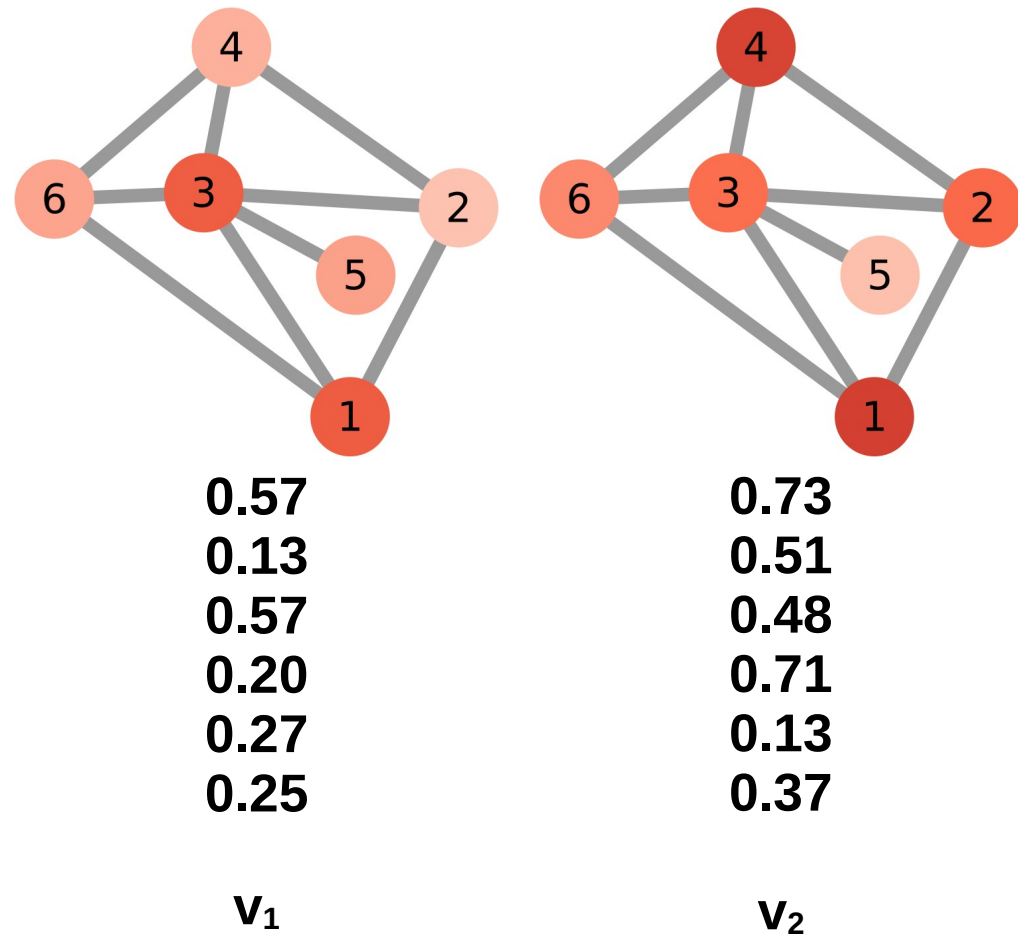
# NVD: What?

- A graph G

# NVD: What?
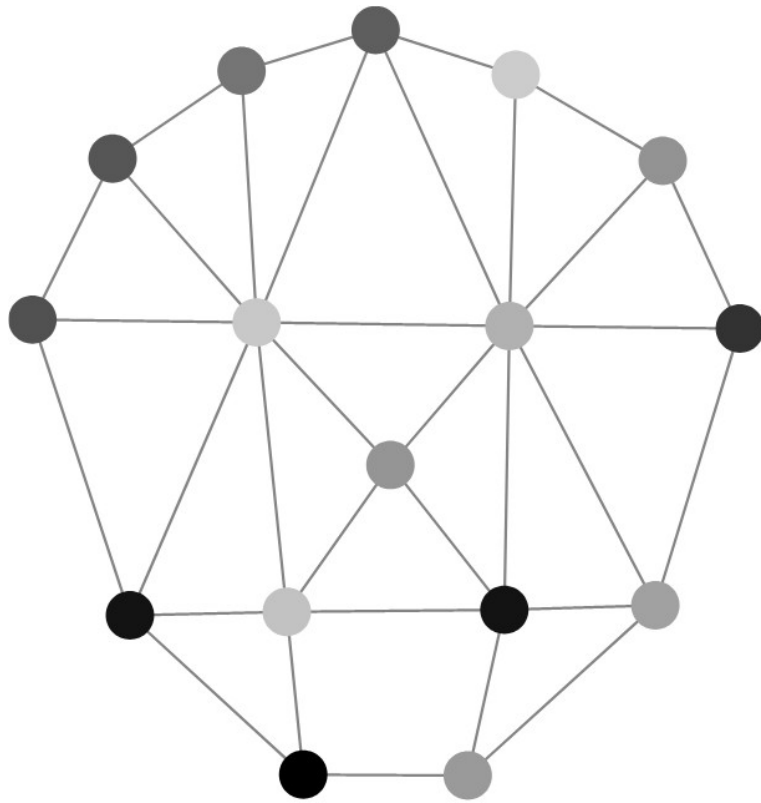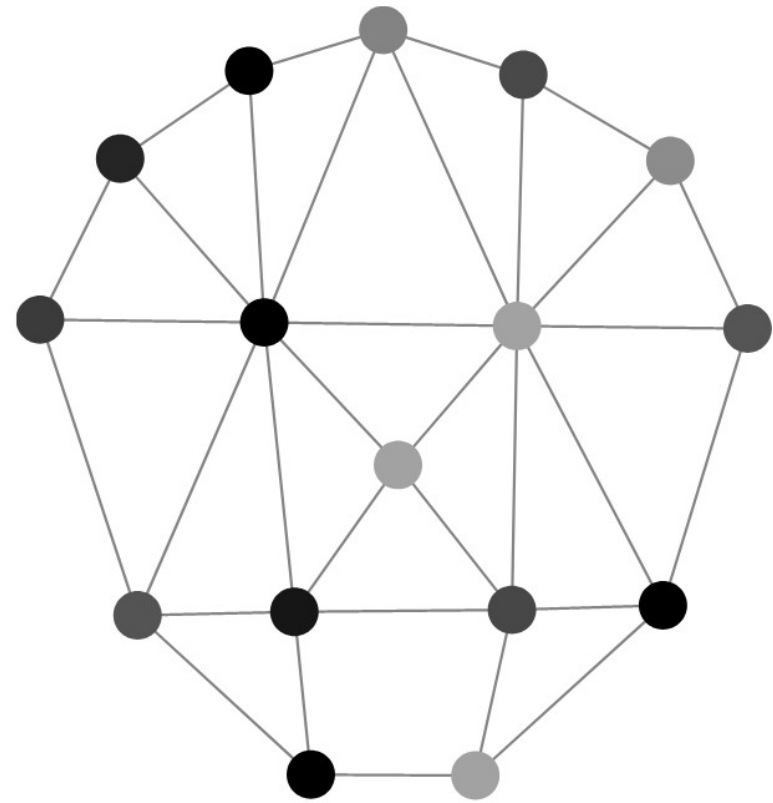
- A graph G
- Two vectors $v_1$ and $v_2$
  - One value per node
- How far is $v_1$ from $v_2$?



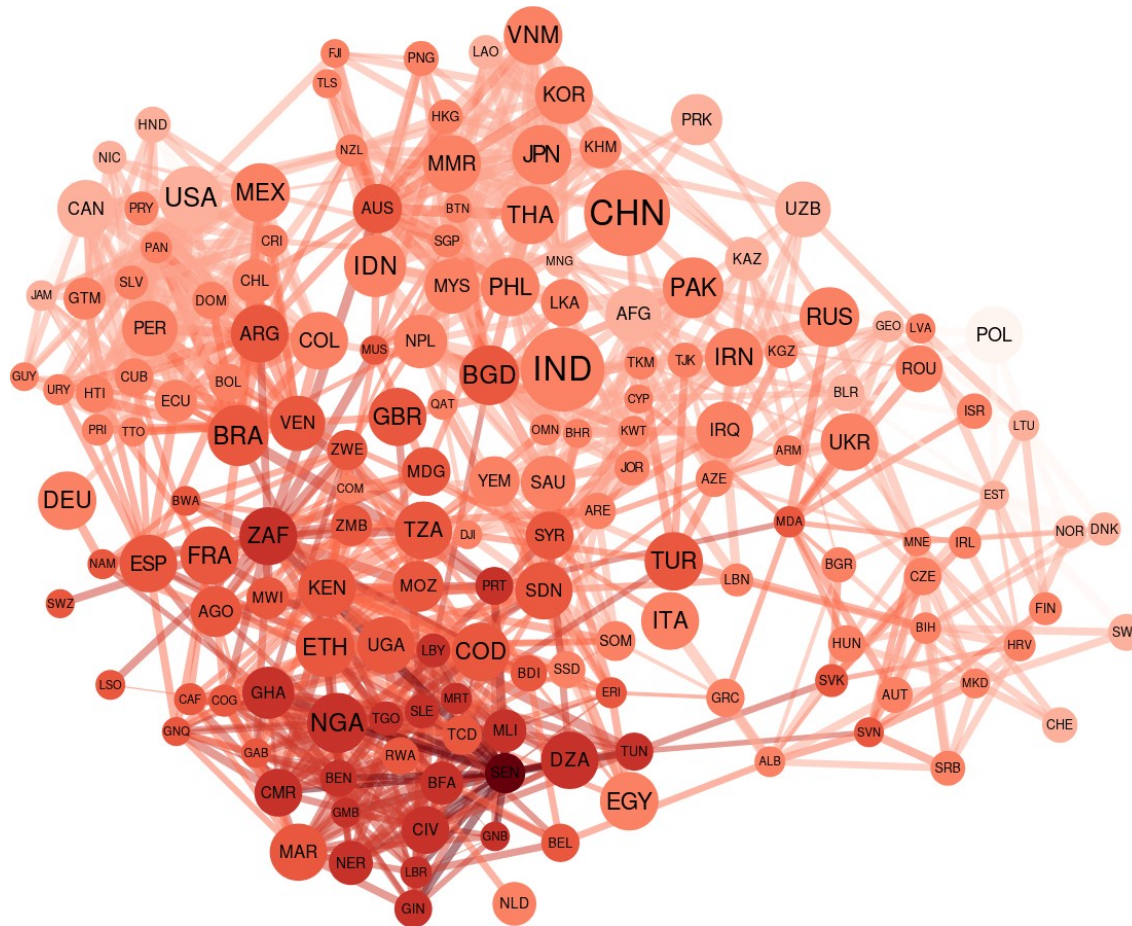| $v_1$ | $v_2$ |
|-------|-------|
| 0.57  | 0.73  |
| 0.13  | 0.51  |
| 0.57  | 0.48  |
| 0.20  | 0.71  |
| 0.27  | 0.13  |
| 0.25  | 0.37  |

# NVD: Why?



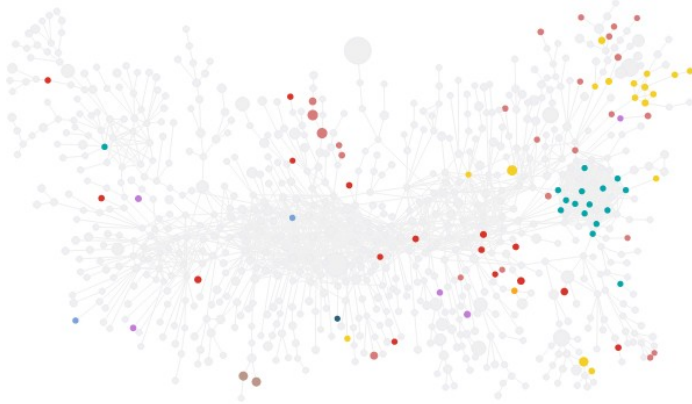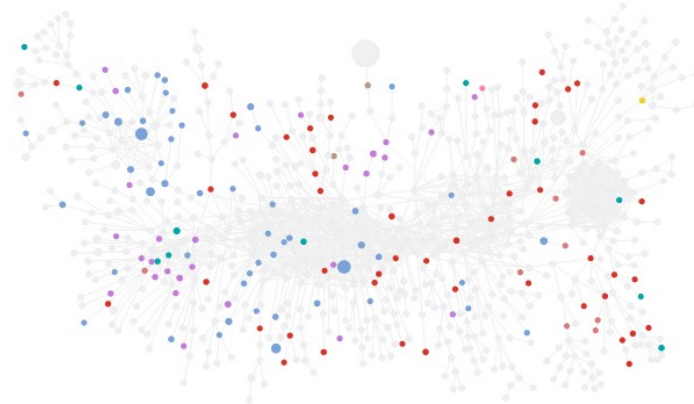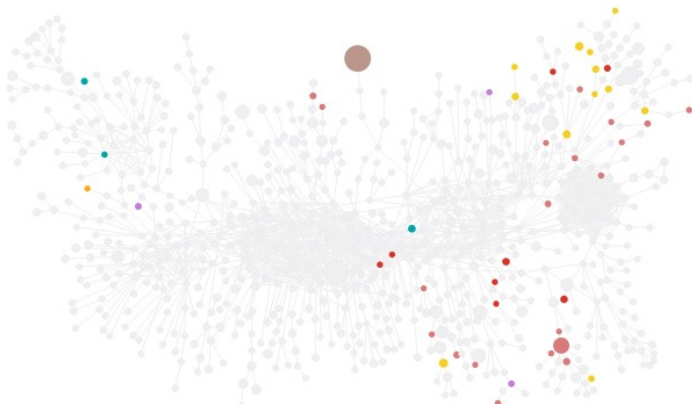(a) Face #1                    (b) Face #2
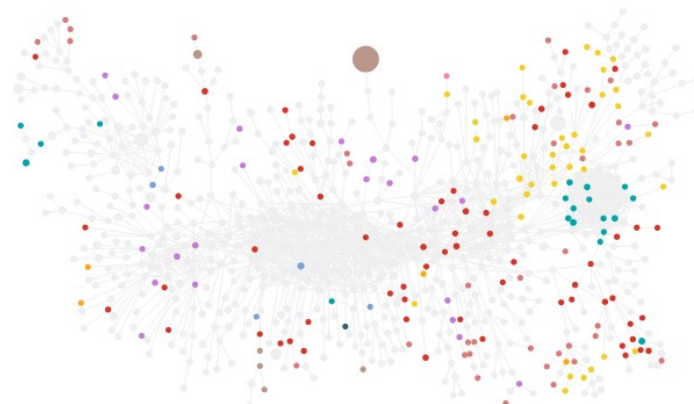
# NVD: Why?

# NVD: Why?



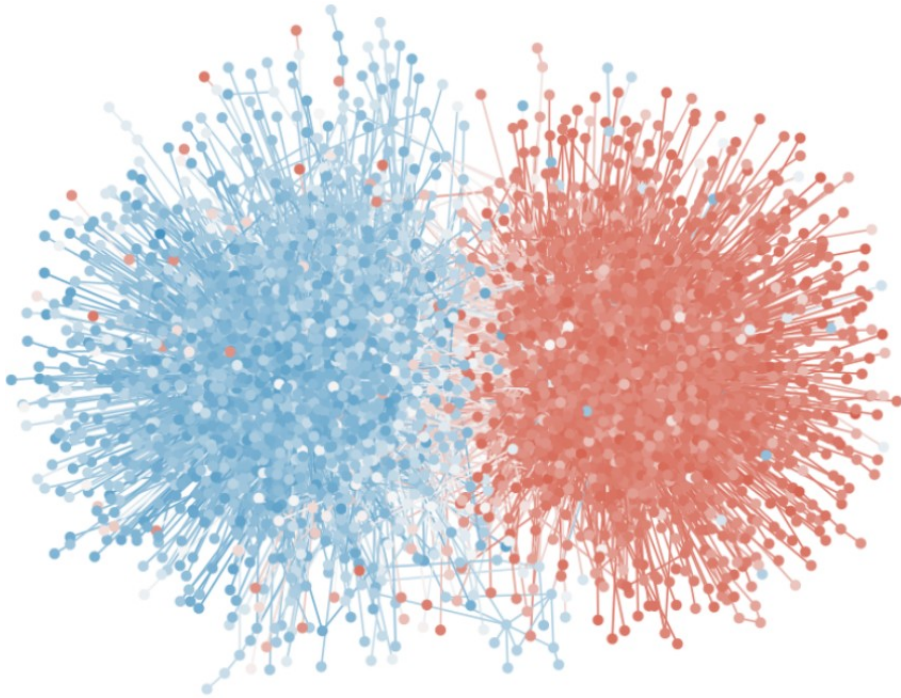(a) Korea 1962

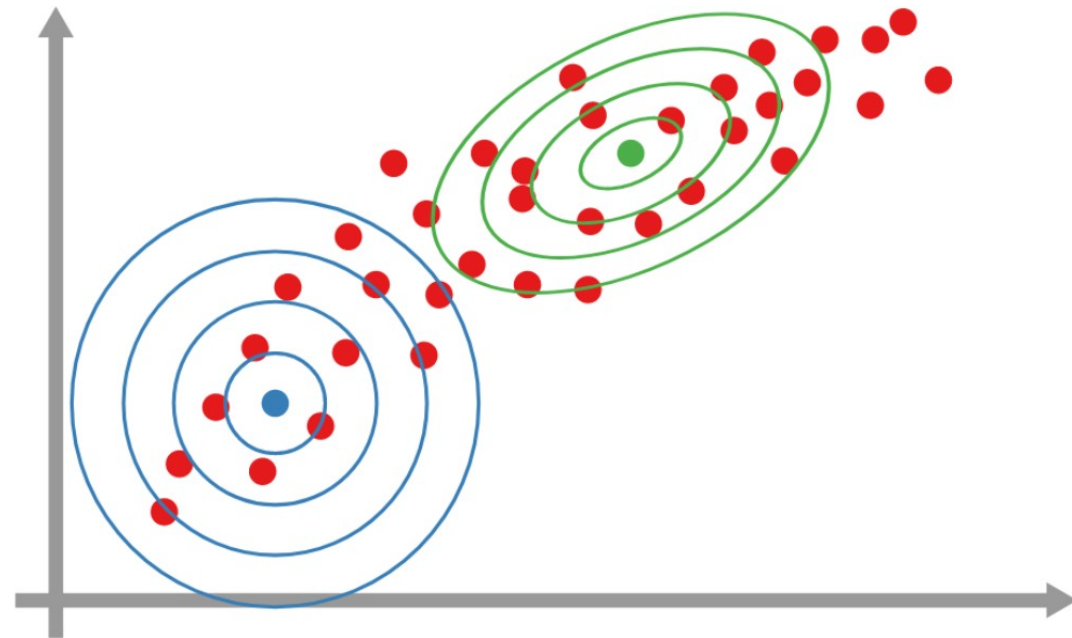(b) Korea 2013

(c) Egypt 1962

(d) Egypt 2013

# NVD: Why?

# Re-Cap on Distances

- Euclidean "straight line"

$$\sqrt{(p-q)^T I (p-q)}$$

- Mahalanobis "bendy" space

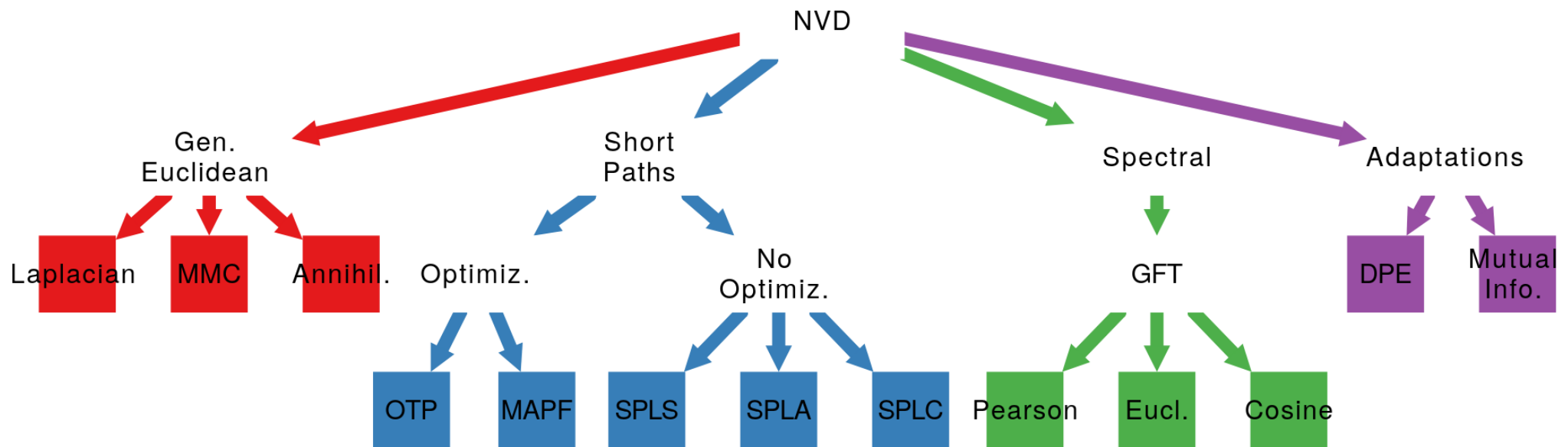$$\sqrt{(p-q)^T \operatorname{cov}(p,q)^{-1}(p-q)}$$
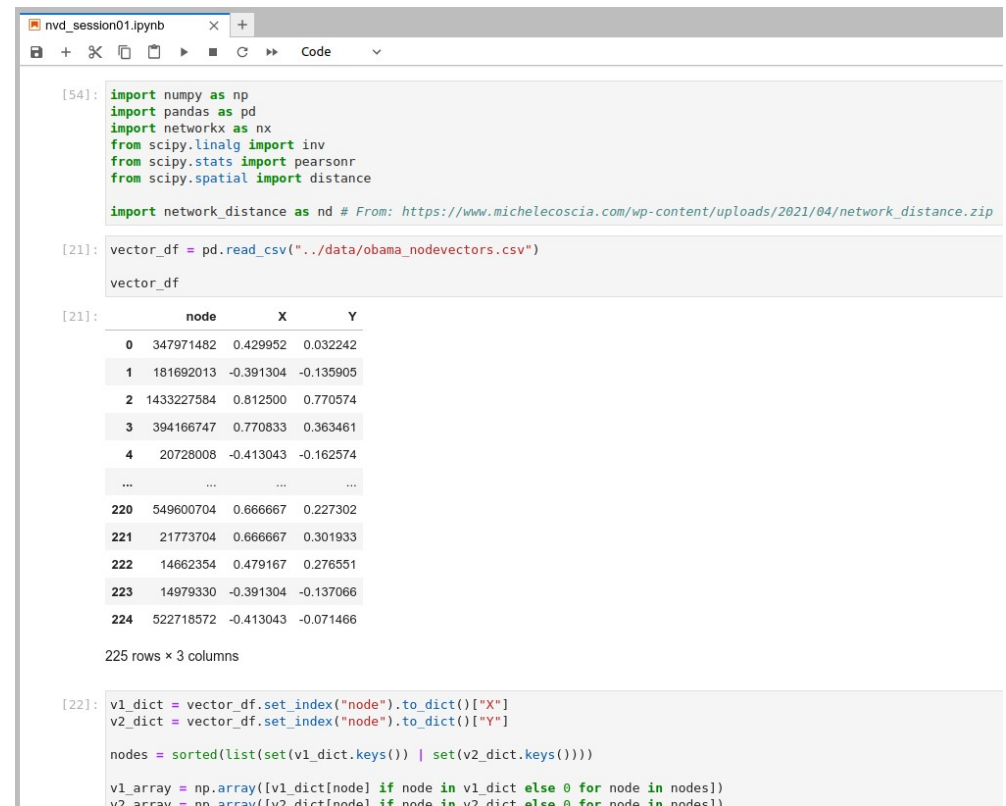
# Euclidean vs Network



Same Euclidean distances, different network distances
The network "bends" the space, like in Mahalanobis

# Current Methods

# Tutorial Part #1

- Objectives:
  - Refresher on vector distances
  - Understand I/O of library
  - Set up the data

# Generalized Euclideans

# Heat Diffusion

- Imagine each node as a thermometer
- Connected to other thermometers to pass heat
- If we know the heat h at time t
- What would be the heat at t + 1?

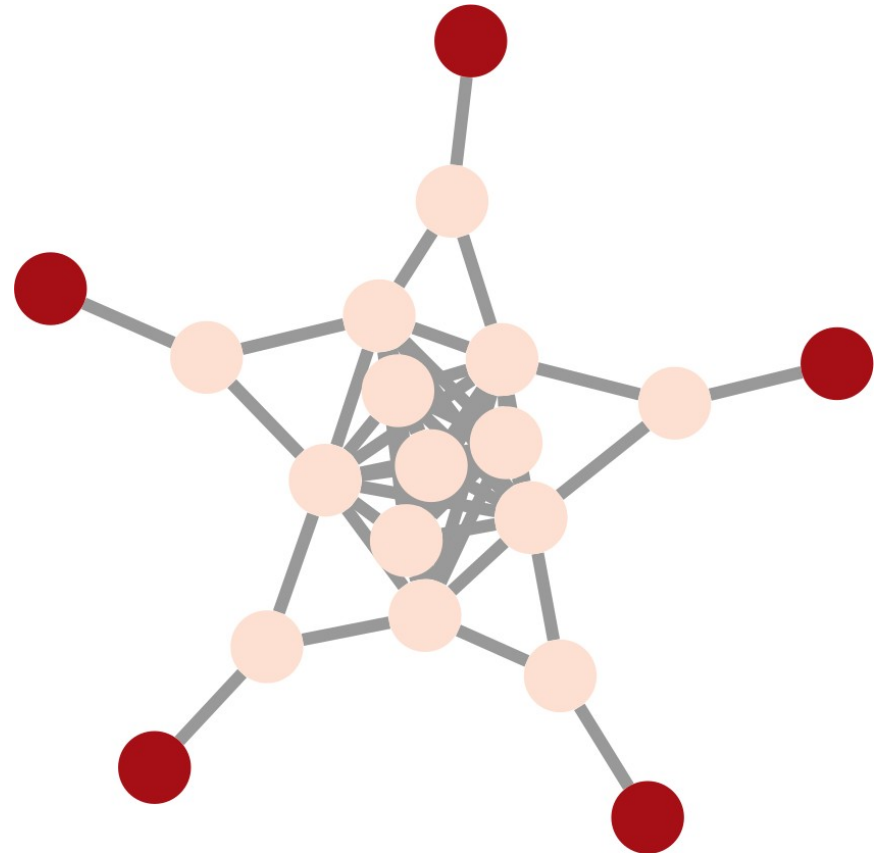$$\frac{\partial h}{\partial t} = -Lh$$

**Laplacian**

# Generalized Euclidean: Laplacian

- L tells us how easy it is to pass heat between nodes

- Which is via random walks

- Thus its inverse gives a sense of distance

$$\sqrt{(p-q)^T L^{-1} (p-q)}$$

# Network Variance

- Not only about distances!

- Vector v: how spread out is on the network?

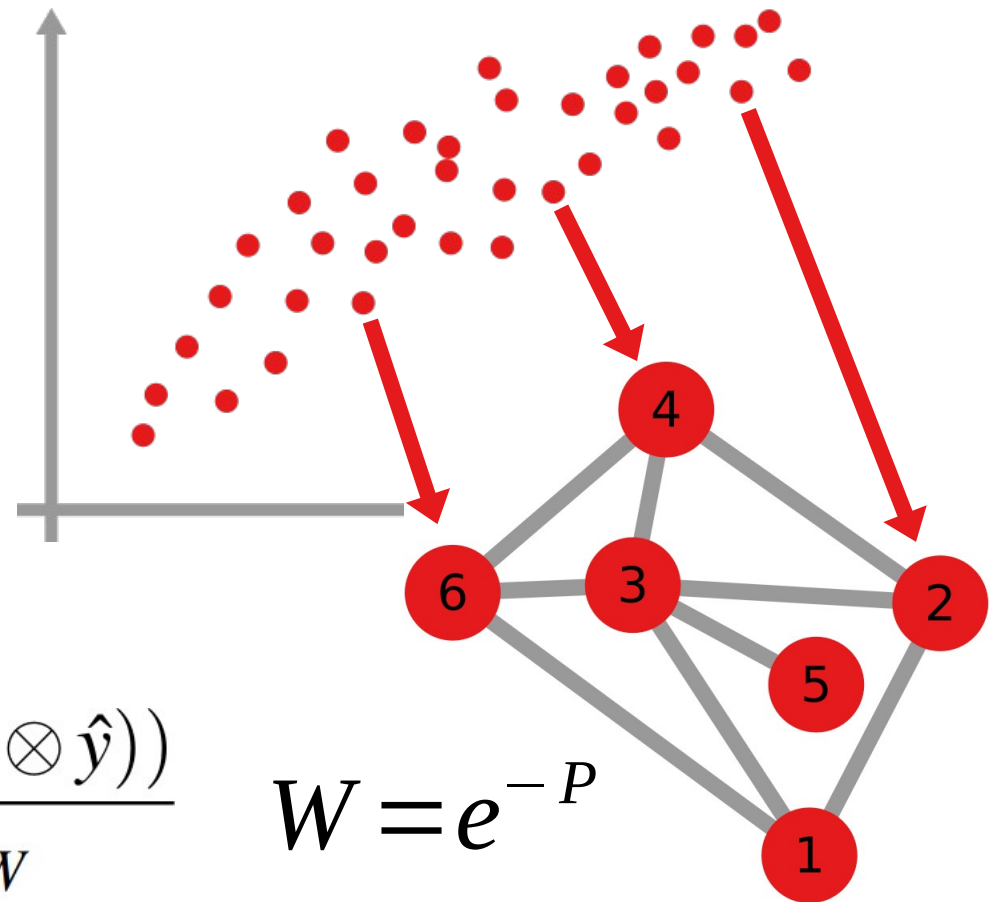- Same as variance, now network variance

# Network Correlation

- How are two variables related?

- Pearson works in a flat space

- What if we have a network?

- E.g. correlation between age and sharing activity in a social network

$$\rho_{x,y,G} = \frac{\text{sum}(W \times (\hat{x} \otimes \hat{y}))}{\sigma_{x,W}\,\sigma_{y,W}}$$
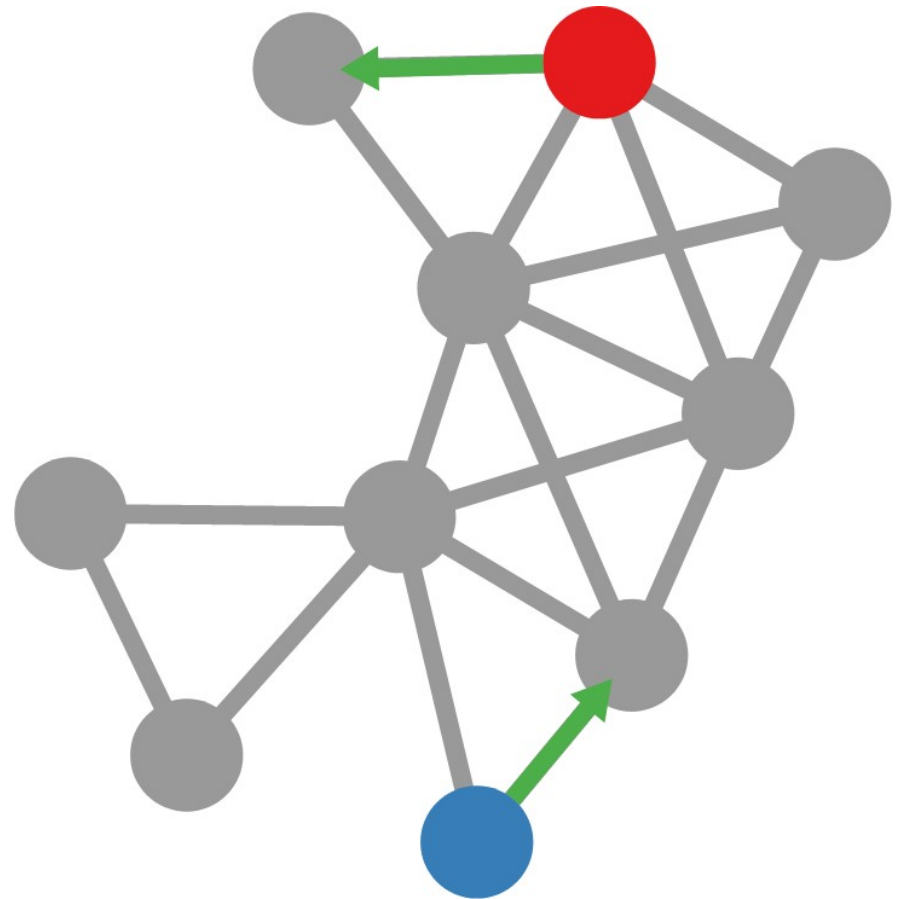
$$W = e^{-P}$$

# Random Walks

- If we move randomly, we expect bump into each other

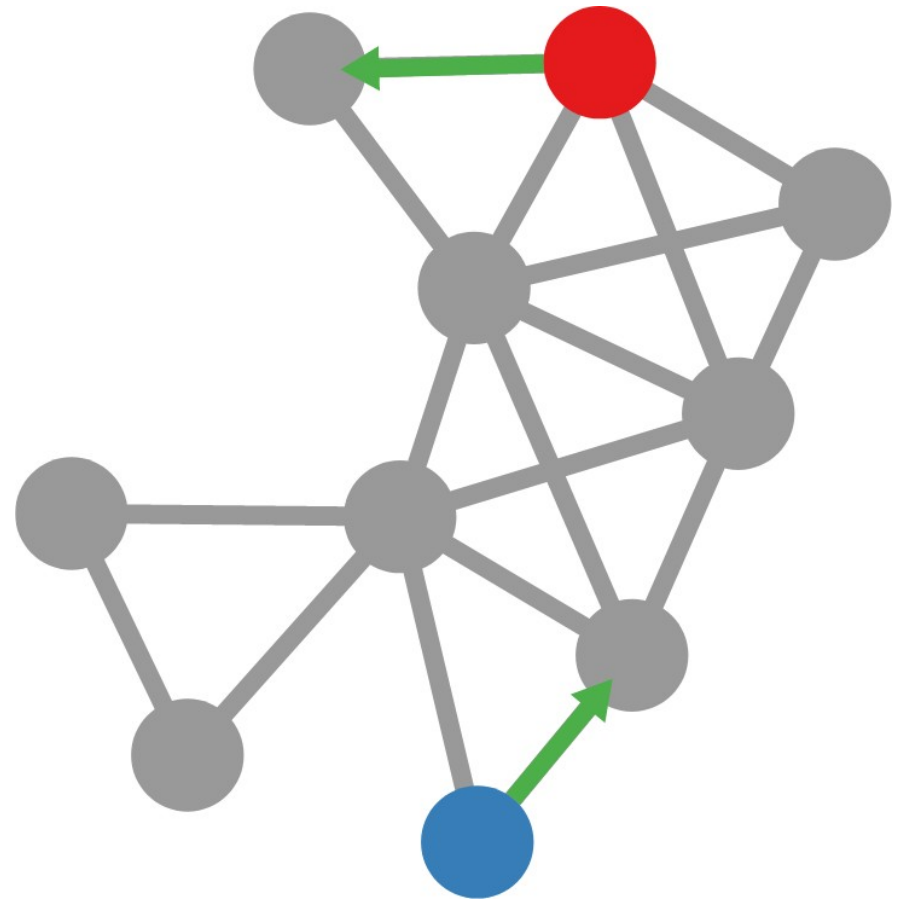- The later we do so from what we expected, the farther we were

- Z-scores!

$$\sqrt{(p-q)^T Z^{-1} (p-q)}$$

# Random Walks #2

- $A^k$: p of going from $v_1$ to $v_2$ in k steps (RW)

- For which k would our vectors completely overlap each other?

$$\sqrt{(p-q)^T \sum_{k=0}^{\infty} A^k (p-q)}$$

# Tutorial Part #2

- Objectives:
  - Build an intuition of the measure with a simple toy experiment
  - Calculate Euclidean, variance, and correlation on a network
  - Calculate alternative Euclideans
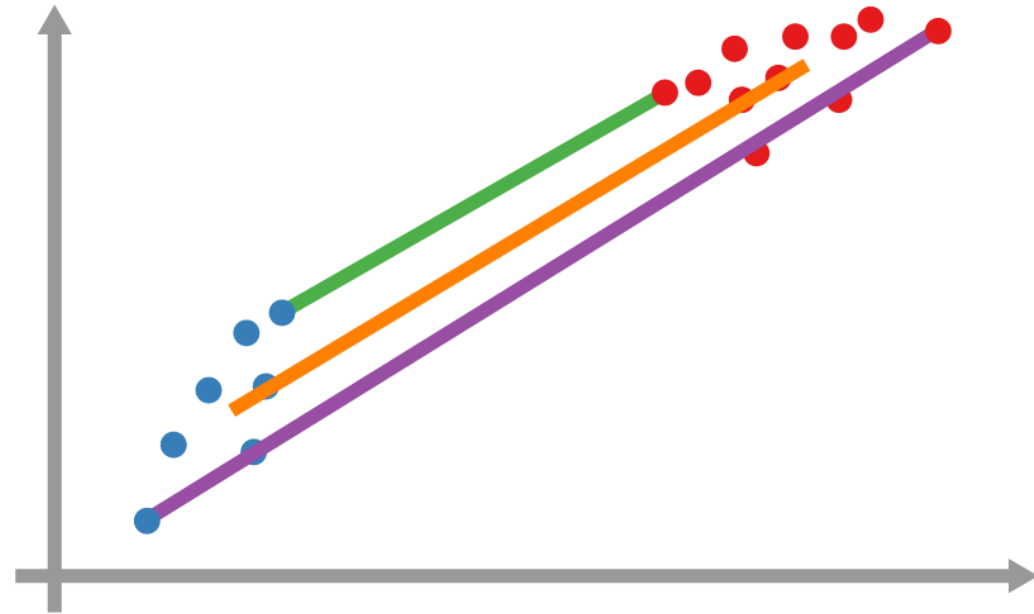
# Shortest Path Distances

# NVD Seems Easy

- Just calculate shortest paths!
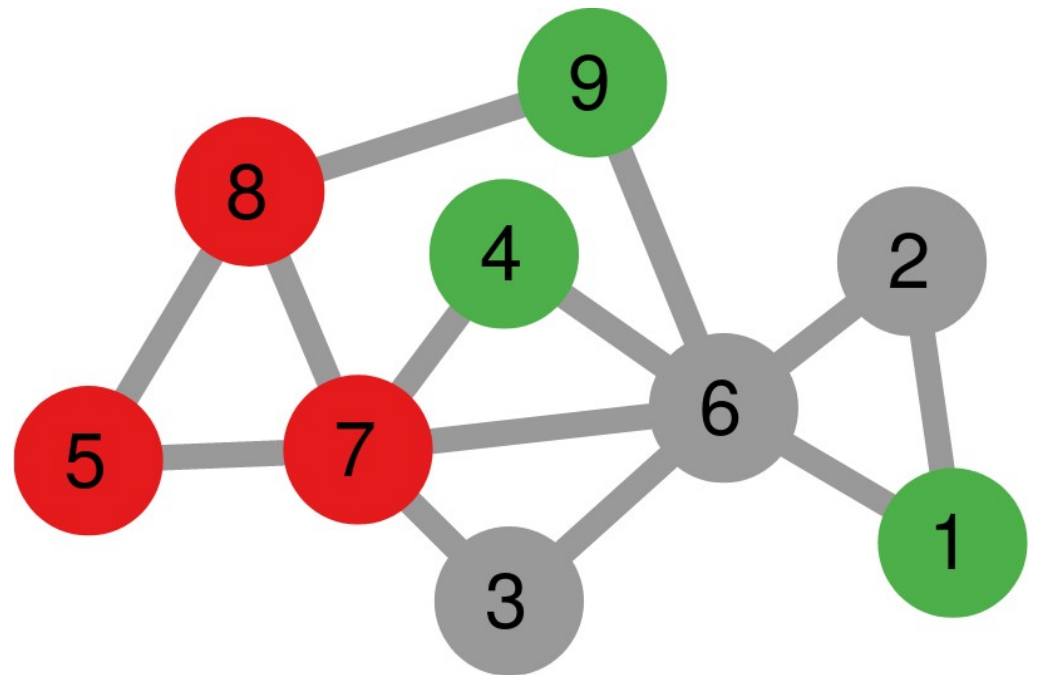
- But then: which ones do we choose?

# Non-Optimized

- **Single** linkage
  - Always pick the shortest available

- **Complete** Linkage
  - Always pick the longest available

- **Average** Linkage
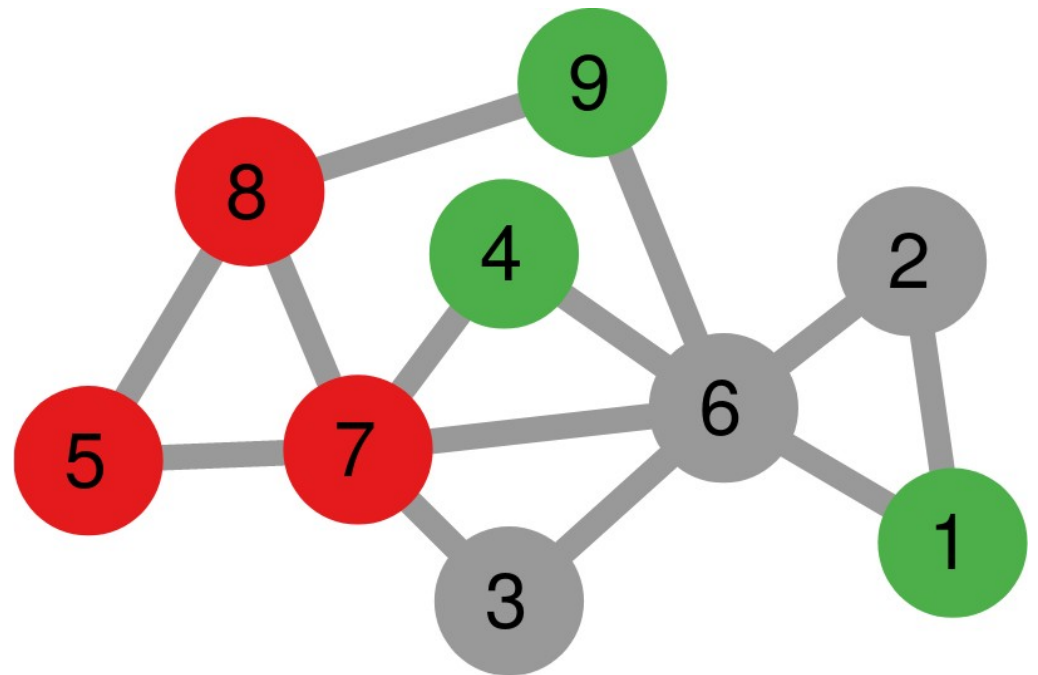  - Pick them all and weigh the alternatives

# Non-Optimized

- ## Single Linkage
  - 1 + 1 + 3
- ## Complete Linkage
  - 3 + 2 + 2
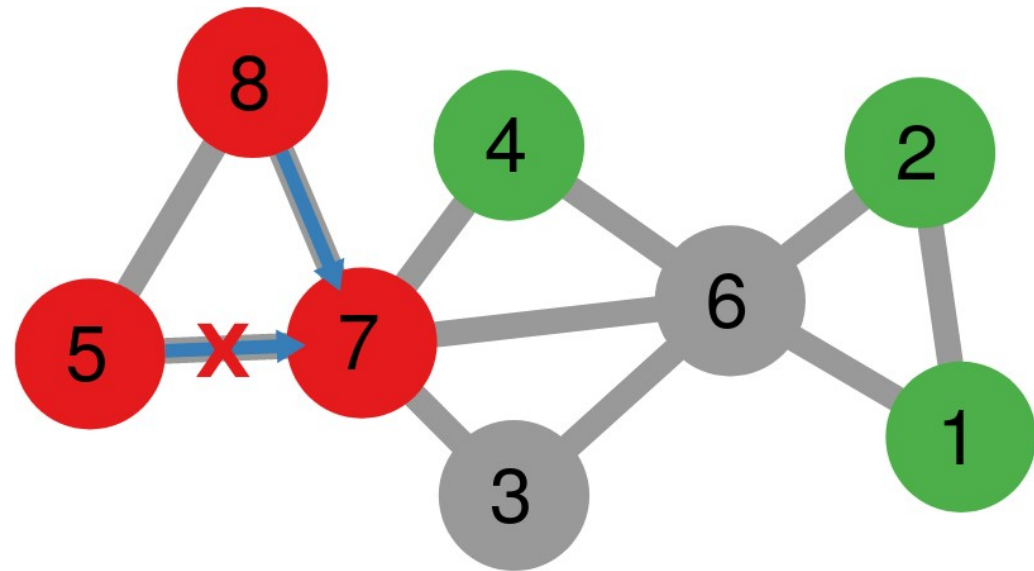- ## Average Linkage
  - 18 / 3

# Earth Mover Distance

- Find the best set of paths

- Minimize the cost

- Similar to single linkage

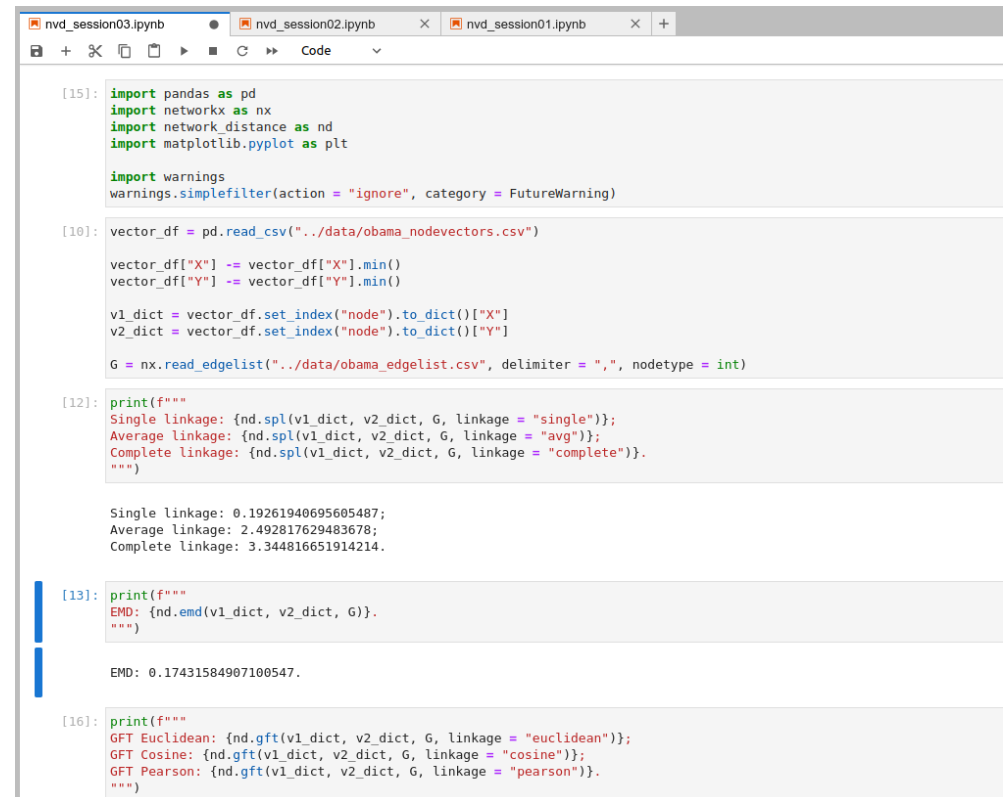# Multi Agent Path Finding

- Similar to EMD

- With constraints

- Nodes & edges have capacity limits

- E.g. cannot use the same link at the same time

# Tutorial Part #3

- Objectives:
  - Experiment with different linkage criteria
  - Use optimal EMD solution
  - (Skipping ahead) Test spectral distance

# Spectral Methods
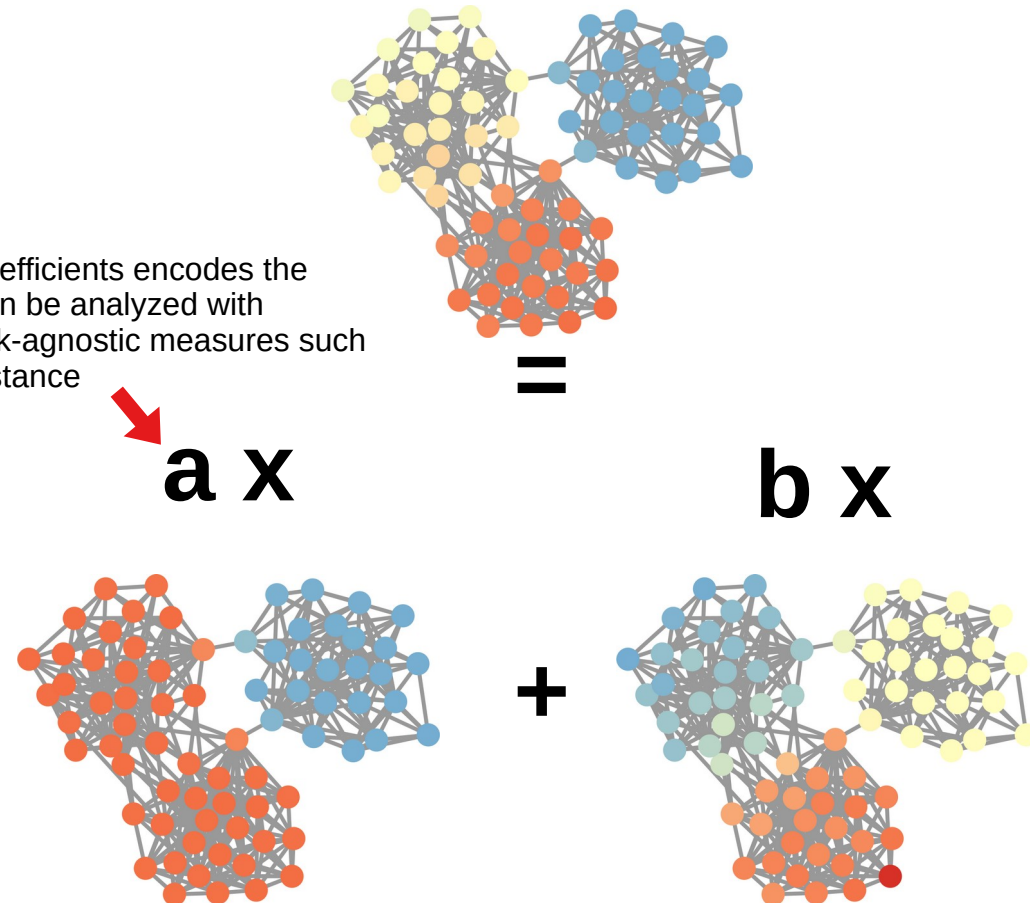
# Graph Fourier Transform

- FT: decompose any signal into its component frequencies

=

a x
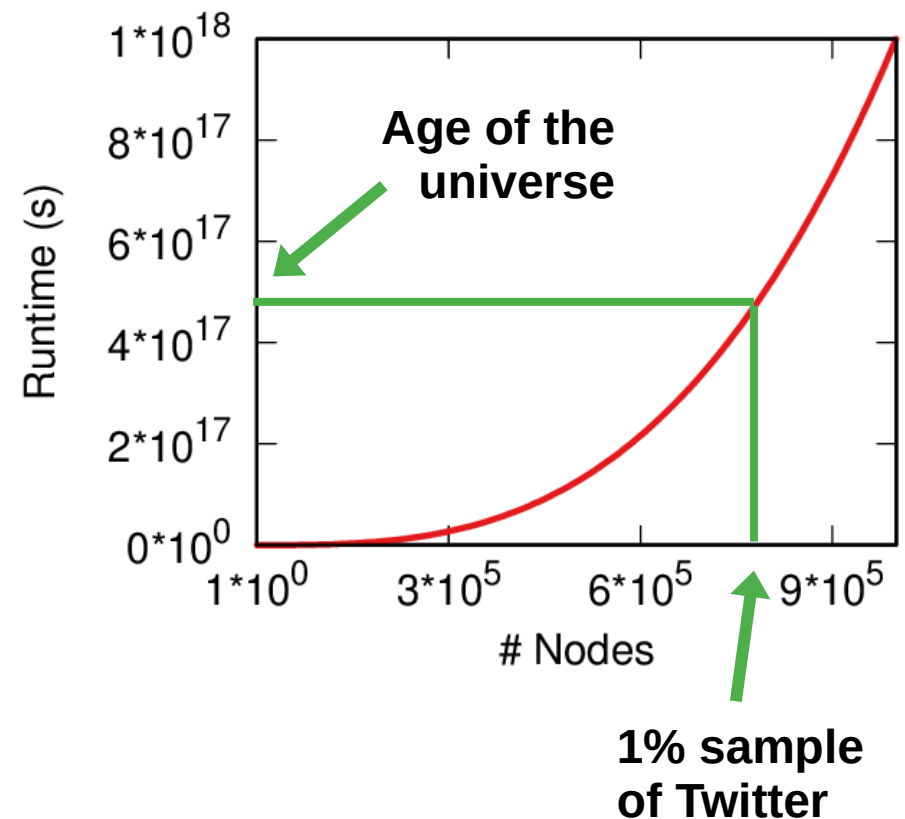
+

b x

# Graph Fourier Transform

- FT: decompose any signal into its component frequencies

- GFT: same, but the components are the eigenvectors of the graph

- Remember: eigenvectors encode structure

The vector of coefficients encodes the structure and can be analyzed with classical network-agnostic measures such as Euclidean distance
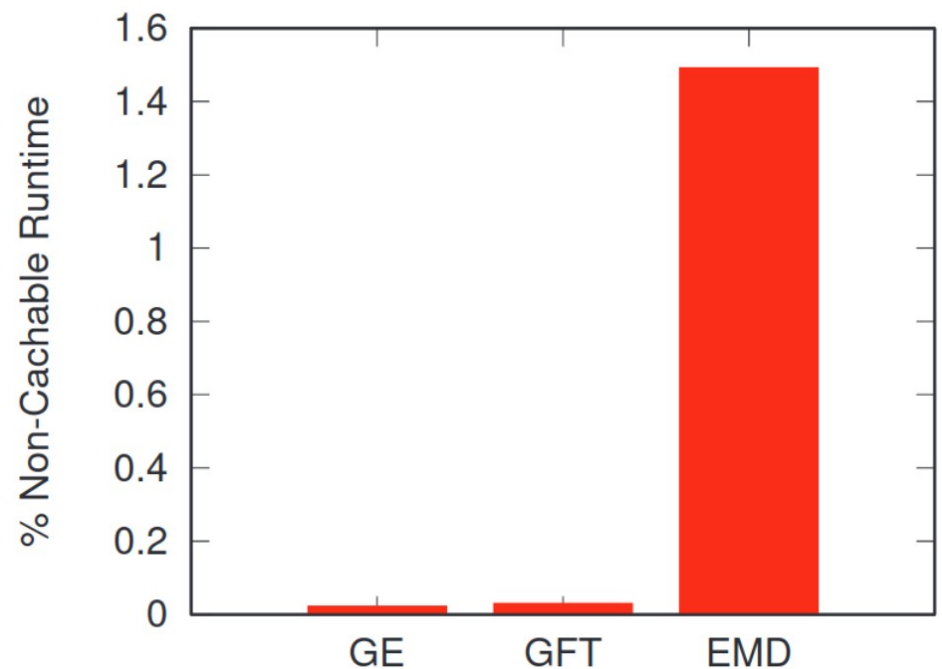
$=$

**a x**

**b x**

**+**

# Computational Efficiency

- Expensive steps can be $O(|V|^3)$…
  - Pseudoinverse for GE
  - Shortest paths
  - Eigenvectors of Laplacian

# Some Tricks

- Expensive part needs to be done once per network

- You can re-use it for all node vector pairs you have
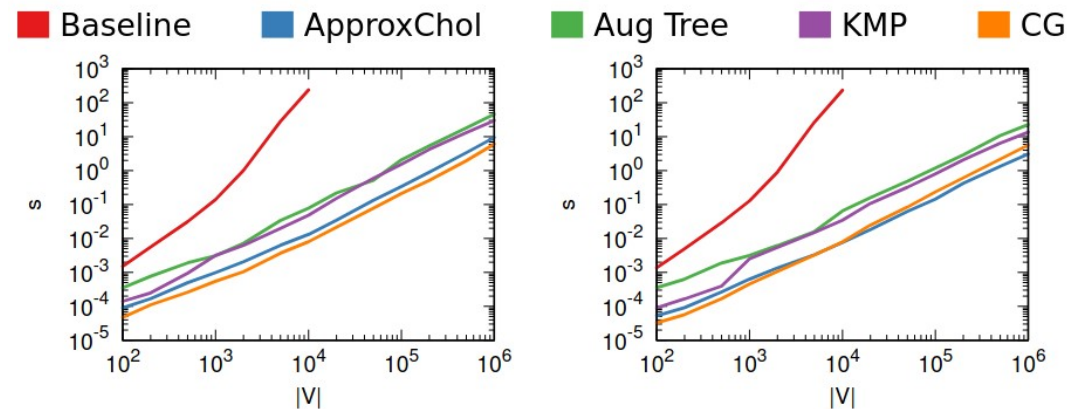
- E.g. disease spreading on an unchanging social network

# Laplacian Solvers

$$\sqrt{(p-q)^T \boxed{L^{-1}(p-q)}}$$
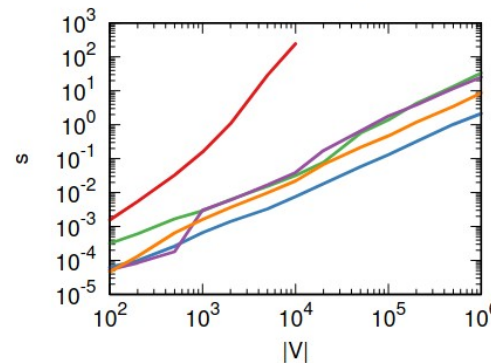
This part can be
calculated efficiently

**Without** pseudo-inverting
the Laplacian!

With Laplacian solvers (in
Julia)



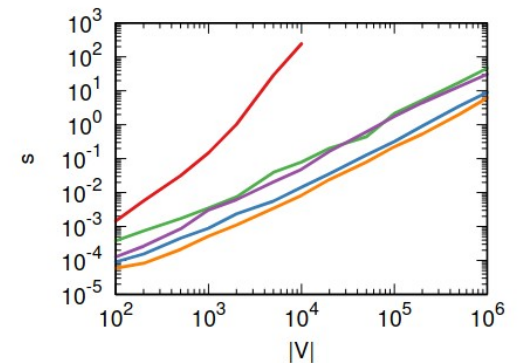Legend: Baseline, ApproxChol, Aug Tree, KMP, CG
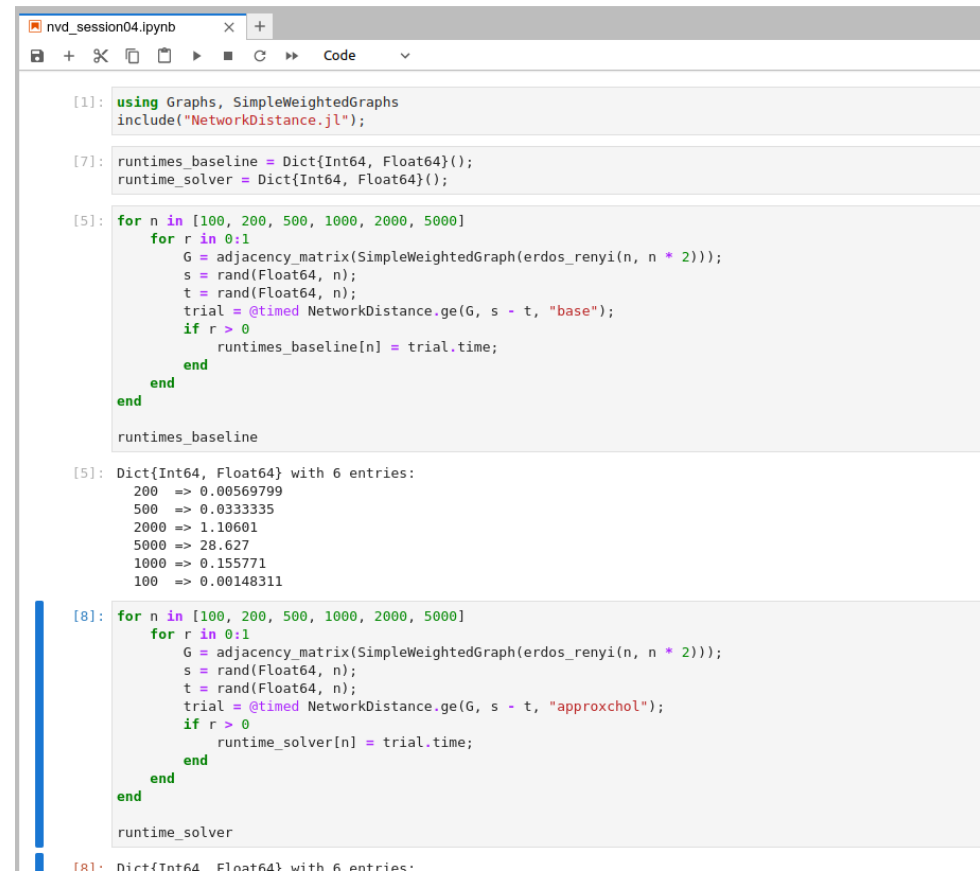
(a) Erdos-Renyi

(b) Barabasi-Albert

(c) Watts-Strogatz

(d) Stochastic Blockmodel

# Tutorial Parts #4 & #5

- Objectives:
  - Testing runtimes
  - Experimenting with Laplacian solvers
  - Trying out some application scenarios