**CPT Iain J. Cruickshank**
icruicks@Andrew.cmu.edu

**Dr. Kathleen M. Carley**
kathleen.carley@cs.cmu.edu

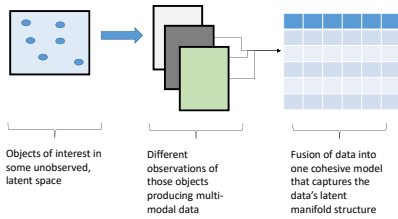**Carnegie Mellon University**
www.casos.cs.cmu.edu

# Multi-modal Graph Learning and Misinformation

## Misinformation is a Problem

- "Fake news is now viewed as one of the greatest threats to democracy, justice, public trust, freedom of expression, journalism and economy" (Shu et al. 2019).
- Credited with significant measurable impacts from effecting elections to impacting the stock market
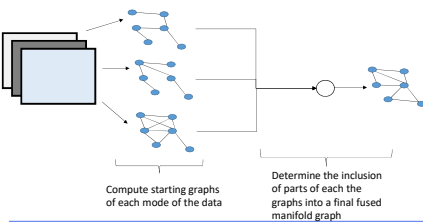- Misinformation is also difficult to characterize (Shu et al. 2019, 2017)

## Misinformation is better characterized by Multi-modal Data

- There are many ways to characterize misinformation, from the use of false information, to writing style, to how it propagates
- Recent studies in supervised learning (Shu et al. 2019) have found better performance in identifying misinformation by supervised learning by *incorporating all of these characterizations*
- Thus, if we treat the characterizations of misinformation as multi-modal data, then we should have better unsupervised detection of misinformation
  - Multi-modal data: different types or sources of data that describe the same event or persons



Objects of interest in some unobserved, latent space

Different observations of those objects producing multi-modal data

Fusion of data into one cohesive model that captures the data's latent manifold structure
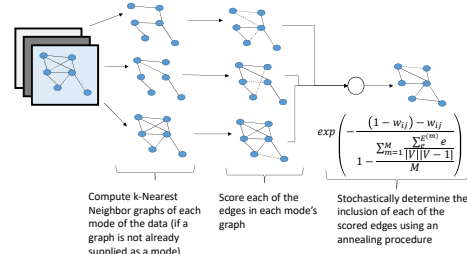
## Multi-modal Graph Learning

- Graphs make ideal structures for representing data
  - Have emergent structures, like clusters
  - Allow for local heterogeneity
  - Interpretable by man and machine
- The fundamental idea of graph learning is to find the best graph representation of some data
  - Can be used with multi-modal data



Compute starting graphs of each mode of the data

Determine the inclusion of parts of each the graphs into a final fused manifold graph

## Graph Annealing
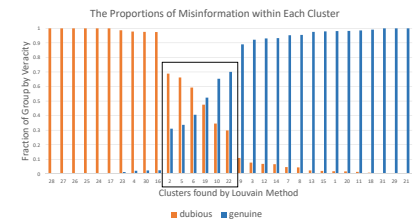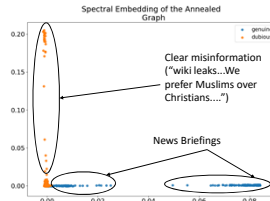
- Our proposed solution for creating a suitable data model for unsupervised detection of fake news from multi-modal data is Graph Annealing



Compute k-Nearest Neighbor graphs of each mode of the data (if a graph is not already supplied as a mode)

Score each of the edges in each mode's graph

$$exp\left(-\frac{(1-w_{ij})-w_{ij}}{1-\frac{\sum_{m=1}^{M}\frac{\sum \tilde{e}^{(m)} e}{|V||V-1|}}{M}}\right)$$

Stochastically determine the inclusion of each of the scored edges using an annealing procedure

- The final step of Edge selection uses a metropolis criterion balancing the probability of an edge existing with the sparsity of the graphs, like the decision function in simulated annealing

## Annealed Graphs for Analyzing Misinformation

- To test the data fusion process of graph annealing we used a simple misinformation dataset from Kaggle.
  - ~18,000 articles with ~40/60 'misinformation'/genuine
  - Only preprocessing was to remove articles that did not have text and those with foreign text
- We used three characterizations of the data
  - The actual words used (Frequency Matrix)
  - The Parts of Speech used (document-by-PoS)
  - Re-publication network (classic network)
- Visualizing the manifold of the annealed graph by Laplacian Eigenmaps displayed distinct and overlapping regions of misinformation and genuine information



Spectral Embedding of the Annealed Graph

Clear misinformation ("wiki leaks...We prefer Muslims over Christians....")

News Briefings



The Proportions of Misinformation within Each Cluster

- Clustering the annealed graph by the Louvain method resulted in 31 clusters
- Clusters present in the annealed graph reflect both topic groups and veracity of the particular article
  - Group 5 contains a high amount 2016 election news, most of which is dubious
  - Group 6 contains mostly conflict related news (e.g. airstrikes in Africa, casualties in Afghanistan), also with a high amount of fake news.
  - Group 22 contains a lot of articles on gun control and recent public shooting events
- Overall, using an annealed graph model made for ease of identifying not only areas of likely misinformation, but those themes and topics which are being subjected to heavy amounts of misinformation, in an unsupervised setting

institute for SOFTWARE RESEARCH