# Modeling and Mining of Networked Information Spaces

## http://www.mathstat.dal.ca/~mominis
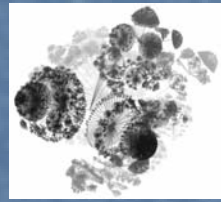


A collaboration graph

## The Opportunity

*Show me your friends and I will tell you what you are.*

- Google: the more incoming links, the more worthy a web page is (Page Rank).
- Citation Graph (network)
- Email Graph (network)
- Phone Call Graph (network)
- Collaboration graph (network)
- **All analysis is based on the network links, not the content !!!**

## Citation indexing of the scientific literature

- Used to be done manually, updated periodically
- Citeseer automated the process (in CSCI)
- Focused crawler collects papers off the Web
- Intelligent document processing extracts title, authors, abstract, references
- Builds citation graph
  - nodes are papers
  - directed links are references/citations
- Analyses citation graph

## Citation graph

- In-degree follows power law
  - Fraction of web pages having $k$ incoming links is proportional to $1/k^2$
- Tightly connected
  - **even** after removing high hub and authority articles,
  - Bridges between subareas offer insight

## Dynamic graphs

- Reflect evolution in social networks over time
  - Email graph
  - Phone call graph
  - Citation graph
  - Coauthorship graph
- Abnormal patterns may signify
  - Unusual event
  - Fraud
  - Terrorist activity

## K-cores

- *k-core* is the subgraph generated by recursively removing all nodes of degree less than *k*.
- Here the 1-core is the full graph.
- The 2-core is composed of the red and green nodes.
- The 3 and 4-cores consist of only the red nodes.



## Citation graphs



## Power Law Distribution



## Connectivity of the Citation graph



## Identifying Interesting Regions

- Real world graphs are often too big to visualize.
- We need to extract features from these graphs to describe and compare them.
- Highly connected regions of a graph can viewed as forming a graph structure.
- We use this structure to characterize a large graph.

## How to find "hot spots" in dynamic graphs?

- Open question
- Things to try:
  - Features of node neighbourhoods (as a time series)
  - Feature evolution
    - Over time
      - As neighbourhood size changes
  - Features used to classify nodes
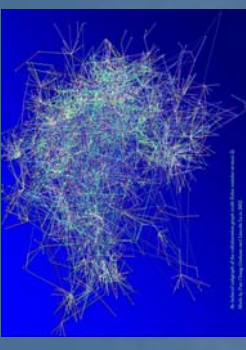
## Email graph of Dal Computer Science

- Classify nodes into
  - faculty,
  - students,
  - staff, or
  - mailing lists
- Detect events
  - exams,
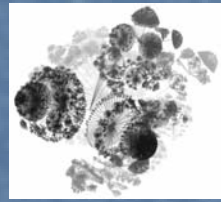  - Study breaks,
  - open house,
  - local conference