

Michael W. Mahoney

Yahoo Research

(For more info, see:
http://www.cs.yale.edu/homes/mmahoney)

Joint work with: Jure Leskovec, Kevin Lang and Anirban Dasgupta

Lots and lots of large data!

- DNA micro-array data and DNA SNP data
- High energy physics experimental data
- Hyper-spectral medical and astronomical image data
- Term-document data
- Medical literature analysis data
- Collaboration and citation networks
- Internet networks and web graph data
- Advertiser-bidded phrase data
- Static and dynamic social network data

Networks in the wide world

- technological networks
 - AS, power-grid, road networks
- biological networks
 - food-web, protein networks
- social networks
 - collaboration networks, friendships
- language networks
 - semantic networks...
- .





Large Social and Information Networks

Interaction graph model of

networks:

- Nodes represent "entities"
- Edges represent "interaction" between pairs of entities



Large networks at Yahoo:

- social networks
 - -Y! messenger buddylist
 - -interaction on Y! Answers
 - -address book in Y! mail
- information networks
 - -advertiser-query
 - -user-query
 - -query-webpage
 - -webpage-webpage

Sponsored ("paid") Search

Text based ads driven by user specified query

| 🕲 recipe indian food - Yahoo! Search Results - Mozilla Firefox | |
|--|---|
| Eile Edit <u>V</u> iew Hi <u>s</u> tory <u>B</u> ookmarks <u>Y</u> ahoo! <u>T</u> ools <u>H</u> elp | |
| 🔄 🔹 📄 👻 🕑 👔 🎦 http://search.yahoo.com/search?p=recipe+indian+food&fr=yfp-t-501&toggle=1&cop=mss&ei=UTF-8 | 3 🔊 🔹 🕨 🔽 indian food recipes |
| 🕌 Rutgers University Li 🗋 my del.icio.us 🗋 post to del.icio.us | |
| MN • powered by MICOL SEARCH 🖗 🔍 Web Search • 🏂 🔯 • 🖬 • 🖓 34º F • 🔊 News (0) • 🕹 My | / Games 🔻 選 Storage 👻 |
| Y - 🖉 - recipe indian food 🔹 🔶 Search Web - 🖶 - 🔯 Mail - 💽 - 🐼 My Yahoo! 🌾 NCAA Hoops - | 🏆 Fantasy Sports 🔹 📥 Games 🔹 🧦 Music 🔹 🛛 🚿 |
| Yahoo! My Yahoo! Mail Welcome, Guest [Sign In] | Advertiser Sign In Help |
| Web Images Video Local Shopping more Images Video Search recipe indian food Search Search | Answers |
| Search Results 1 - 10 of about 7,260, | 000 for recipe indian food - 0.19 sec. (About this page) |
| Recipe Indian Food www.MonsterMarketplace.com - Browse and compare great deals on recipe indian food. Indian Food sanfrancisco.citysearch.com - Find great Indian restaurants in your area today. Search here. | ILTS SPONSOR RESULTS Indian Food Buy indian food at SHOP.COM. Search our free shipping offers. www.SHOP.com |
| indian food recipe indian food recipe Title: Indian Food Recipe. Yield: 4 Servings. Ingredients. 1 bunch to the echo by: Jonathan Kandell Indian Food Recipes Put recipes.chef2chef.net/recipe-archive/43/231458.shtml - 13k - <u>Cached</u> - <u>More from this site</u> | Recipe India Food Find and Compare prices on recipe india food at Smarter.com. www.smarter.com |
| Recipe Gal: Indian Foods Indian Recipes from Recipe Gal's Archives All Food Posters. Travel Posters. Indian Recipes. Indian Breads Indian Chicken Recipes www.recipegal.com/indian - 10k - <u>Cached</u> - <u>More from this site</u> | Chinese Food Recipe Books on Cataloglink Find chinese food recipe books on CatalogLink. www.CatalogLink.com |
| Indian Recipes, Indian Food Recipe, South Indian Recipes, Indian indian recipes, indian food recipe, south indian Recipes, indian cooking Recipes, Indian Recipes, Indian Food Recipe, South Indian Recipes, Indian Cooking Recipe, www.india4world.com/indian-recipe - 17k - <u>Cached</u> - <u>More from this site</u> Paav Bhaaji - Recipe for Paav Bhaaji - Pao Bhaji | \$19.97 Over 500 Chinese Recipes Cookbook 100% Satisfaction Guaranteed, 543-Page Chinese Cookbook Only \$19.97. ✓ |

Graphs and sponsored search data

 bid, click and impression information for "keyword x advertiser" pair

 mine information at querytime to provide new ads

> - maximize CTR, RPS, advertiser ROI



Sponsored Search Problems

Marketplace depth broadening:

find new advertisers for a particular query/submarket

• Query recommender system:

suggest to advertisers new queries that have high probability of clicks

Contextual query broadening:

broaden the user's query using context information, e.g. past sessions, phrasing, etc.

Graph Mining Recommendations



Micro-markets

find micro-markets by partitioning the "query x advertiser" graph:



advertiser





10 million keywords

What do these networks "look" like?



Questions of interest ...

What are degree distributions, clustering coefficients, diameters, etc.? Heavy-tailed, small-world, expander, geometry+rewiring, local-global decompositions, ... Are there natural clusters, communities, partitions, etc.? Concept-based clusters, link-based clusters, density-based clusters, ... (e.g., isolated micro-markets with sufficient money/clicks with sufficient coherence) How do networks grow, evolve, respond to perturbations, etc.? Preferential attachment, copying, HOT, shrinking diameters, ... How do dynamic processes - search, diffusion, etc. - behave on networks? Decentralized search, undirected diffusion, cascading epidemics, ... How best to do learning, e.g., classification, regression, ranking, etc.? Information retrieval, machine learning, ...

Clustering and Community Finding

• Linear (Low-rank) methods

If Gaussian, then low-rank space is good.

• Kernel (non-linear) methods

If low-dimensional manifold, then kernels are good

Hierarchical methods

Top-down and botton-up -- common in the social sciences

• Graph partitioning methods

Define "edge counting" metric -- conductance, expansion, modularity, etc. -- in interaction graph, then optimize!

"It is a matter of common experience that communities exist in networks ... Although not precisely defined, communities are usually thought of as sets of nodes with better connections amongst its members than with the rest of the world."



Communities, Conductance, and NCPPs

Let A be the adjacency matrix of G=(V,E).

The conductance φ of a set S of nodes is:

$$\phi(S) = \frac{\sum_{i \in S, j \notin S} A_{ij}}{\min\{A(S), A(\overline{S})\}}$$

$$A(S) = \sum_{i \in S} \sum_{j \in V} A_{ij}$$

The Network Community Profile (NCP) Plot of the graph is:

$$\Phi(k) = \min_{S \subset V, |S|=k} \phi(S)$$

Just as conductance captures the "gestalt" notion of cluster/community quality, the NCP plot measures cluster/community quality as a function of size.

Probing Large Networks with Approximation Algorithms

Idea: Use approximation algorithms for NP-hard graph partitioning problems as experimental probes of network structure.

Spectral - (quadratic approx) - confuses "long paths" with "deep cuts" Multi-commodity flow - (log(n) approx) - difficulty with expanders SDP - (sqrt(log(n)) approx) - best in theory Metis - (multi-resolution for mesh-like graphs) - common in practice X+MQI - post-processing step on, e.g., Spectral of Metis

Metis+MQI - best conductance (empirically)

Local Spectral - connected and tighter sets (empirically)

We are not interested in partitions per se, but in probing network structure.

Low-dimensional graphs and expanders



Widely-studied small social networks





Newman's Network Science

Large Social and Information Networks

| • Social nets | Nodes | Edges | Description |
|--|-----------|------------|------------------------|
| LIVEJOURNAL | 4,843,953 | 42,845,684 | Blog friendships [4] |
| Epinions | 75,877 | 405,739 | Who-trusts-whom [35] |
| FLICKR | 404,733 | 2,110,078 | Photo sharing [21] |
| Delicious | 147,567 | 301,921 | Collaborative tagging |
| CA-DBLP | 317,080 | 1,049,866 | Co-authorship (CA) [4] |
| CA-COND-MAT | 21,363 | 91,286 | CA cond-mat [25] |
| • Information networks | | | |
| CIT-HEP-TH | 27,400 | 352,021 | hep-th citations [13] |
| Blog-Posts | 437,305 | 565,072 | Blog post links [28] |
| • Web graphs | | | |
| Web-google | 855,802 | 4,291,352 | Web graph Google |
| Web-wt10g | 1,458,316 | 6,225,033 | TREC WT10G web |
| • Bipartite affiliation (authors-to-papers) networks | | | |
| ATP-DBLP | 615,678 | 944,456 | DBLP [25] |
| ATP-ASTRO-PH | 54,498 | 131,123 | Arxiv astro-ph [25] |
| • Internet networks | | | |
| AS | 6,474 | 12,572 | Autonomous systems |
| GNUTELLA | 62,561 | 147,878 | P2P network [36] |

Table 1: Some of the network datasets we studied.

Large Social and Information Networks



More large networks



10⁵

 10^{4}

10⁶

10⁵

10⁴

"Whiskers" and the "core"

- Whiskers
 - maximal sub-graph detached from network by removing a single edge
 - on average, contains 40% of nodes and 20% of edges
- Core
 - the rest of the graph, i.e., the 2-edge-connected core
 - on average, contains 60% of nodes and 80% of edges
- Global minimum of NCPP is a whisker



Examples of whiskers

Ten largest "whiskers" from CA-cond-mat.





"Regularization" and spectral methods

• regularization properties: spectral embeddings stretch along directions in which the random-walk mixes slowly

-Resulting hyperplane cuts have "good" conductance cuts, but may not yield the optimal cuts







spectral embedding

notional flow based embedding

Regularized and non-regularized communities (1 of 2)



- Metis+MQI (red) gives sets with better conductance.
- Local Spectral (blue) gives tighter and more well-rounded sets.



Regularized and non-regularized communities (2 of 2)

Two ca. 500 node communities from Local Spectral Algorithm:



Two ca. 500 node communities from Metis+MQI:





Lower Bounds ...

- ... can be computed from:
- Spectral embedding

(independent of balance)

SDP-based methods

(for volume-balanced partitions)





Lots of Generative Models

• Preferential attachment - add edges to high-degree nodes

(Albert and Barabasi 99, etc.)

• Copying model - add edges to neighbors of a seed node

(Kumar et al. 00, etc.)

- Hierarchical methods add edges based on distance in hierarchy (Ravasz and Barabasi 02, etc.)
- Geometric PA and Small worlds add edges to geometric scaffolding (Flaxman et al. 04; Watts and Strogatz 98; etc.)
- Random/configuration models add edges randomly

(Molloy and Reed 98; Chung and Lu 06; etc.)

NCPP for common generative models



A simple theorem: random graphs

Let $\mathbf{w} = (w_1, \dots, w_n)$, where $w_i = ci^{-1/(\beta-1)}, \quad \beta \in (2,3).$ Connect nodes *i* and *j* w.p. $p_{ij} = w_i w_j / \sum_k w_k.$



NCP Plot for G(w) model, with power-law degrees and $\beta \epsilon$ (2,3).



Structure of the G(w) model, with $\beta \epsilon$ (2,3).

Note: Sparsity is the issue, not heavy-tails per se. (Power laws with $\beta \epsilon$ (2,3) give us the appropriate sparsity.)

A "forest fire" model

Leskovec, Kleinberg, and Faloutsos 2005

At each time step, a new node vi:

- links to a "seed" node w, chosen uniformly at random.
- selects x "outlinks" and y "inlinks" of w at random.

• forms "outlinks," i.e., burns, to those selected nodes and then proceeds to burn recursively.

Notes:

• Preferential attachment flavor - second neighbor is not uniform at random.

- Copying flavor since burn seed's neighbors.
- Hierarchical flavor seed is parent.

• "Local" flavor - burn "near" -- in a diffusion sense -the seed vertex.



NCPP Of the FF Model

Two different parameter values:



Note: for these parameters, this model also reproduces "densification" and "shrinking diameters" of real graphs (Leskovec et al. 05).

Comparison with "Ground truth" (1 of 2)

Networks with "ground truth" communities:

- LiveJournal12:
 - users create and explicitly join on-line groups
- CA-DBLP:
 - publication venues can be viewed as communities
- AmazonAllProd:
 - each item belongs to one or more hierarchically organized categories, as defined by Amazon
- AtM-IMDB:
 - countries of production and languages may be viewed as communities (thus every movie belongs to exactly one community and actors belongs to all communities to which movies in which they appeared belong)

Comparison with "Ground truth" (2 of 2)





Miscellaneous thoughts ...

Sociological work on community size (Dunbar and Allen)

- 150 individuals is maximum community size
- On-line communities have 60 members and break down at 80
 - Military companies, divisions of corporations, etc. close to the Dunbar's 150

Common bond vs. common identity theory

- Common bond (people are attached to individual community members) are smaller and more cohesive
- Common identity (people are attached to the group as a whole) focused around common interest and tend to be larger and more interpersonally diverse

What edges "mean" and community identification

social networks - reasons an individual adds a link to a friend can vary enormously citation networks or web graphs - links are more "expensive" and are more semantically uniform.

Conclusions

• about networks and data:

Can use approximation algorithms as experimental probes "Best" communities get less and less "community-like"

"Octopus" or "Jellyfish" model - with "whiskers" and "core"

• about modeling these networks:

Common generative models don't capture community phenomenon Graph locality - important for realistic network generation Local regularization - important due to sparsity

Workshop on "Algorithms for Modern Massive Data Sets" (http://mmds.stanford.edu)

Stanford University and Yahoo! Research, June 25-28, 2008

Objectives:

- Address algorithmic, mathematical, and statistical challenges in modern statistical data analysis.

- Explore novel techniques for modeling and analyzing massive, high-dimensional, and nonlinearstructured data.

- Bring together computer scientists, mathematicians, statisticians, and data analysis practitioners to promote cross-fertilization of ideas.

Organizers: M. W. Mahoney, L-H. Lim, P. Drineas, and G. Carlsson.

Sponsors: NSF, Yahoo! Research, PIMS, DARPA.