Algorithmic and Statistical Perspectives on Large-Scale Data Analysis

Michael W. Mahoney

Stanford University Feb 2010

(For more info, see: <u>http:// cs.stanford.edu/people/mmahoney/</u> or Google on "Michael Mahoney")



Algorithmic vs. Statistical Perspectives

Lambert (2000)

Computer Scientists

- Data: are a record of everything that happened.
- Goal: process the data to find interesting patterns and associations.
- *Methodology*: Develop approximation algorithms under different models of data access since the goal is typically computationally hard.

Statisticians

- Data: are a particular random instantiation of an underlying process describing unobserved patterns in the world.
- Goal: is to extract information about the world from noisy data.
- *Methodology*: Make inferences (perhaps about unseen events) by positing a model that describes the random variability of the data around the deterministic model.

Perspectives are NOT incompatible

• Statistical/probabilistic ideas are central to recent work on developing improved randomized algorithms for matrix problems.

• Intractable optimization problems on graphs/networks yield to approximation when assumptions are made about network participants.

• In boosting (a statistical technique that fits an additive model by minimizing an objective function with a method such as gradient descent), the computation parameter (i.e., the number of iterations) also serves as a regularization parameter.

Matrices and graphs in data analysis

Graphs:

• *model* information network with graph G = (V,E) -- vertices represent "entities" and edges represent "interactions" between pairs of entities

Matrices:

• *model* data sets by a matrix -- since an m x n matrix A can encode information about m objects, each of which is described by n features

Matrices and graphs represent a nice tradeoff between:

- descriptive flexibility
- algorithmic tractability

But, the issues that arise are very different than in traditional linear algebra or graph theory ...

The gap between NLA and TCS ... (... was the genesis of MMDS!)

Matrix factorizations:

• in NLA and scientific computing - used to express a problem s.t. it can be solved more easily.

• in TCS and statistical data analysis - used to represent structure that may be present in a matrix obtained from object-feature observations.

NLA:

- emphasis on optimal conditioning,
- backward error analysis issues,
- is running time a large or small constant multiplied by n^2 or n^3 .

TCS:

- originally motivated by large data applications
- consider space-constrained or pass-efficient models
- exploiting over-sampling and randomness as computational resources.

Outline

• "Algorithmic" and "statistical" perspectives on data problems

Genetics application

DNA SNP analysis --> choose columns from a matrix

PMJKPGKD, Genome Research '07; PZBCRMD, PLOS Genetics '07; Mahoney and Drineas, PNAS '09; DMM, SIMAX '08; BMD, SODA '09

Internet application

Community finding --> partitioning a graph

LLDM (WWW'08 & TR-Jrnl'08 & WWW'10)

We will focus on what was going on "under the hood" in these two applications --- use statistical properties implicit in worst-case algorithms to make domain-specific claims!

DNA SNPs and human genetics

- Human genome \approx 3 billion base pairs
- 25,000 30,000 genes
- Functionality of 97% of the genome is unknown.
- Individual "polymorphic" variations at ≈ 1 b.p./thousand.

SNPs are known locations at the human genome where two alternate nucleotide bases (alleles) are observed (out of A, C, G, T).

individuals

SNPs

... AG CT GT GG CT CC CC CC AG AG AG AG AG AG AG AA CT AA GG GG CC GG AG CG AC CC AA CC AA GG TT AG CT CG CG CG AT CT CT AG CT AG GG GT GA AG GG TT TT GG TT CC CC CC CC GG AA AG AG AG AA CT AA GG GG CC AG AG GG AA CC AA CC AA GG TT AG CT CG GG GG TT TT CC GG TT GG GG TT GG AA GG TT TT GG TT CC CC CC CC GG AA AG AG AG AA CT AA GG GG CC AG AG CG AC CC AA CC AA GG TT AG CT CG CG CG AT CT CT AG CT AG GG GT GA AG GG TT TT GG TT CC CC CC CC GG AA AG AG AG AA CC GG AA CC GG AG CC AG AC CC AA CC AA GG TT AG CT CG CG CG AT CT CT AG CT AG GG GT GA AG GG TT TT GG TT CC CC CC CC GG AA AG AG AG AA CC GG AA CC CC AG GG CC AC CC AA CG AA GG TT AG CT CG CG CG AT CT CT AG CT AG GT GT GA AG GG TT TT GG TT CC CC CC CC GG AA GG GG GG AA CC AA GG CT GG AA CC AC CG AA CC AA GG TT AG CT CG CG CG AT CT CT AG CT AG GG TT GG AA GG TT TT GG TT CC CC CC CG GA AG GG GG AA CT AA GG GG CT GG AG CC CC CG AA CC AA GT TT AG CT CG CG CG AT CT CT AG CT AG GT TGG AA GG TT TT GG TT CC CC CC CG CA AG AG AG AG AG AA CT AA GG GG CT GG AG CC CC CG AA CC AA GT TT AG CT CG CG CG AT CT CT AG CT AG GG TT GG AA GG TT TT GG TT CC CC CC CG CA AG AG AG AG AG AG AT AA GG GG CT GG AG CC CC CG AA CC AA GT TT AG CT CG CG CG AT CT CT AG CT AG GG TT GG AA GG TT TT GG TT CC CC CC CG CC AG AG AG AG AG AG AG CT AA GG GG CT GG AG CC CC CG AA CC AA GT TT AG CT CG CG CG AT CT CT AG CT AG GG TT GG AA GG TT TT GG TT CC CC CC CC CG GA AA GA AG AG AG AA CT AA GG GG CT GG AG CC CC CG AA CC AA GT TT AG CT CG CG CG AT CT CT AG CT AG GG TT GG AA GG TT TT GG TT CC CC CC CC CG GA AA GA AG AG AG AA CT AA GG GG CT GG AG CC CC CG AA CC AA GT TT AG CT CG CG CG AT CT CT AG CT AG GG TT GG AA GG TT TT GG TT CC CC CC CC CG CA AG AG AG AG AA CT AA GG GG CT GG AG CC CC CG AA CC AA GT TT AG CT CG CG CG AT CT CT AG CT AG GG TT GG AA GG TT TT GG TT CC CC CC CC CG GA AA GA AG AG AG AG AA CT AA GG GG CC AG AG CC CA AG CG AA GT TA AT TG GG GG TT TGG AA ...

SNPs occur quite frequently within the genome and thus are effective genomic markers for the tracking of disease genes and population histories.



DNA SNPs and data analysis

A common modus operandi in applying NLA to data problems:

- Write the gene/SNP data as an $m \times n$ matrix A.
- Do SVD/PCA to get a small number of eigenvectors
- Either: interpret the eigenvectors as meaningful i.t.o. underlying genes/SNPs use a heuristic to get actual genes/SNPs from those eigenvectors

Unfortunately, eigenvectors themselves are meaningless (recall reification in stats):

- "EigenSNPs" (being linear combinations of SNPs) can not be assayed ...
- ... nor can "eigengenes" from micro-arrays be isolated and purified ...
- ... nor do we really care about "eigenpatients" respond to treatment ...

DNA SNPs and low-rank methods

PMJKPGKD, Genome Research '07 (data from K. Kidd, Yale University) PZBCRMD, PLOS Genetics '07 (data from E. Ziv and E. Burchard, UCSF)

<u>Common genetics task</u>: find a small subset of informative actual SNPs

to cluster individuals depending on their ancestry

to determine predisposition to diseases

- Algorithmic question: Can we find the best k actual columns from a matrix?
 Can we find actual SNPs that "capture" information in singular vectors?
 Can we find actual SNPs that are maximally uncorrelated?
- Common formalization of "best" lead to intractable optimization problems.

Column Subset Selection Problem (CSSP)

Input: an m-by-n matrix A and a rank parameter k.

Goal: choose *exactly k columns* of A s.t. the m-by-k matrix C minimizes the error:

$$\min ||A - P_C A||_{\xi} = \min ||A - CC^+ A||_{\xi} \quad (\xi = 2, F)$$

- Widely-studied problem in numerical linear algebra and optimization.
- Related to unsupervised feature selection.
- Choose the "best" k documents from a document corpus.

Prior work in NLA & TCS

Numerical Linear Algebra:

 Deterministic, typically greedy, approaches.

(E.g., Golub65, Foster86, Chan87, Chan-Hansen90, Bischof-Hansen91, Hong-Pan92, Chandrasekaran-Ipsen94, Gu-Eisenstat96, Bischof-Orti98, Pan-Tang99, Pan00, ...)

- Deep connection with the Rank Revealing QR factorization.
- Good spectral norm bounds.

Theoretical Computer Science:

- Randomized approaches, with some failure probability. (Much work in last 10 years following Frieze, Kannan, and Vempala; Drineas, Kannan, Mahoney, etc..)
 - More than k columns are picked, e.g., O(poly(k)) columns chosen.
 - Very good (1+ε) Frobenius norm bounds. (Drineas, Mahoney, etc. 2005,2006,2008,2009; Deshpande and Vempala 2006)

Subspace sampling probabilities



NOTE: The rows of V_k^{T} are orthonormal, but its columns $(V_k^{T})^{(i)}$ are not.

$$\begin{pmatrix} A_k \\ m \times n \end{pmatrix} = \begin{pmatrix} U_k \\ m \times k \end{pmatrix} \cdot \begin{pmatrix} \Sigma_k \end{pmatrix} \cdot \begin{pmatrix} V_k^T \\ K \end{pmatrix}$$

 V_k : orthogonal matrix containing the top k right singular vectors of A.

 Σ_k : diagonal matrix containing the top k singular values of A.

A hybrid two-stage algorithm

Boutsidis, Mahoney, and Drineas (2007)

Algorithm: Given an m-by-n matrix A and rank parameter k:

* Not so simple ... Actually, run QR on the down-sampled k-by- $O(k \log k)$ version of V_k^{T} .

• (Randomized phase)

Randomly select $c = O(k \log k)$ columns according to "leverage score probabilities".

• (Deterministic phase)

Run a deterministic algorithm on the above columns* to pick exactly k columns of A.

<u>Theorem</u>: Let C be the m-by-k matrix of the selected columns. Our algorithm runs in "O(mmk)" and satisfies, w.p. $\geq 1-10^{-20}$,

$$||A - P_C A||_F \le O\left(k \log^{1/2} k\right) ||A - A_k||_F$$
$$||A - P_C A||_2 \le O\left(k^{3/4} \log^{1/2} (k) (n-k)^{1/2}\right) ||A - A_k||_2$$

Comparison with previous results

Running time is comparable with NLA algorithms.

Spectral norm:

• Spectral norm bound is $k^{1/4}\log^{1/2}k$ worse than previous work.

Frobenius norm:

• An efficient algorithmic result at most (k logk)^{1/2} worse than the previous existential result.

NLA usually interested in columns for the bases they span! Data analysis usually interested in the columns themselves!

Evaluation on term-document data

TechTC (Technion Repository of Text Categorization Datasets)

 lots of diverse test collections from ODP

 ordered by categorization difficulty

• use hierarchical structure of the directory as background knowledge

• Davidov, Gabrilovich, and Markovitch 2004 Fix k=10 and measure Frobenius norm error:



Things to note ...

٠

Different versions of QR (i.e., different pivot rules) perform differently ...

"obviously," but be careful with "off the shelf" implementations.

QR applied directly to V_k^{T} typically does better than QR applied to A ...

- since V_k^{T} defined the relevant non-uniformity structure in A
- since columns "spread out," have fewer problems with pivot rules

"Randomized preprocessing" improves things even more ...

- due to *implicit* regularization
- (if you are careful with various parameter choices)
- and it improves worse QR implementations more than better code



FIG. 6







• Most NLA codes don't even run on this 90 x 2M matrix.

• Informativeness is a state of the art supervised technique in genetics.

Selecting PCA-correlated SNPs for individual assignment to four continents (Africa, Europe, Asia, America)



SNPs by chromosomal order

Paschou et al (2007) PLoS Genetics

An Aside on: Least Squares (LS) Approximation



Ubiquitous in applications & central to theory:

Statistical interpretation: best linear unbiased estimator.

Geometric interpretation: orthogonally project b onto span(A).

Algorithmic and Statistical Perspectives

$$\begin{aligned} \mathcal{Z}_2 &= \min_{x \in R^d} ||b - Ax||_2 \\ &= ||b - A\hat{x}||_2 \end{aligned}$$

Algorithmic Question: How long does it take to solve this LS problem?
Answer: O(nd²) time, with Cholesky, QR, or SVD*
Statistical Question: When is solving this LS problem the right thing to do?

Answer: When the data are "nice," as quantified by the leverage scores.

*BTW, we used statistical leverage score ideas to get the first (1+ε)-approximation worst-caseanalysis algorithm for the general LS problem that runs in o(nd²) time for *any* input matrix. Theory: DM06,DMM06,S06,DMMS07

Numerical implementation: by Tygert, Rokhlin, etc. (2008)

Statistical Issues and Regression Diagnostics

Statistical Model: $b = Ax + \varepsilon$

 ε = "nice" error process

b' = $A \times_{opt} = A(A^T A)^{-1}A^T b$ = prediction

 $H = A(A^{T}A)^{-1}A^{T}$ is the "hat" matrix, i.e. projection onto span(A)

Statistical Interpretation:

 H_{ij} -- measures the leverage or influence exerted on b'_i by b_j, H_{ii} -- leverage/influence score of the i-th constraint Note: $H_{ii} = |U^{(i)}|_2^2 = row$ "lengths" of spanning orthogonal matrix

(Note: these are the sampling probabilities we used for our worst-case algorithms!)

Hat Matrix and Regression Diagnostics

See: "The Hat Matrix in Regression and ANOVA," Hoaglin and Welsch (1978)



Figure A. The Two Carriers for the Wood Beam Data (Plotting symbol is beam number.).

	j									
i	l	2	З	4	5	6	7	8	9	10
1234567	.418	-,002 .242	.079 .292 .417	274 .136 019 .604	048 .243 .273 .197 .252	.181 .128 .187 038 .111 .148	.128 041 126 .168 030 .042 .262	.222 .033 .044 022 .019 .117 .145	.050 035 153 .275 010 .012 .277	.242 .004 028 010 .111 .174

Things to note:

- Point 4 is a bivariate outlier and H_{44} is largest, just exceeds 2p/n=6/10.
- Points 1 and 3 have relatively high leverage extremes in the scatter of points.
- \cdot H_{1.4} is moderately negative opposite sides of the data band.
- \cdot H₁₈ and H₁₁₀ moderately positive those points mutually reinforce.
- \cdot H₆₆ is fairly low point 6 is in central position.

Leverage Scores of "Real" Data Matrices



Leverage scores of Zachary karate network edge-incidence matrix.



Cumulative leverage score for the Enron email data matrix.

Leverage Scores and Information Gain



Similar strong correlation between (unsupervised) Leverage Scores and (supervised) Informativeness elsewhere!

A few general thoughts

- Q1: Why does a statistical concept like leverage help with worst-case analysis of traditional NLA problems?
- A1: If a data point has high leverage and is *not* an error, as worst-case analysis *implicitly* assumes, it is very important!
- Q2: Why are statistical leverage scores so non-uniform in many modern large-scale data analysis applications?
- A2: Statistical models are often *implicitly* assumed for computational and not statistical reasons---many data points "stick out" relative to obviously inappropriate models!

Outline

- "Algorithmic" and "statistical" perspectives on data problems
- Genetics application

DNA SNP analysis --> choose columns from a matrix

Internet application

Community finding --> partitioning a graph

In many large-scale data applications, "algorithmic" and "statistical" perspectives interact in fruitful ways --- we use statistical properties implicit in worst-case algorithms to make domain-specific claims!

Networks and networked data

Lots of "networked" data!!

- technological networks
 - AS, power-grid, road networks
- biological networks
 - food-web, protein networks
- social networks
 - collaboration networks, friendships
- information networks

- co-citation, blog cross-postings, advertiser-bidded phrase graphs...

language networks

• ...

- semantic networks...

Interaction graph model of networks:

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



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.

Motivation: Sponsored ("paid") Search

Text based ads driven by user specified query

The process:

- Advertisers bids on query phrases.
- Users enter query phrase.
- Auction occurs.
- Ads selected, ranked, displayed.
- When user clicks, advertiser pays!

barcelona chair	Search	Options -	YAHOO!
•	1-10 of 4,220,000 fo	r barcelona chair (<u>About</u>) - 0.09 s	ec SPONSOR RE
	15 BT 155 D		Barcelona Chair Direct from
Also try: barcelona style chair, l	Importer		
		SPONSOR RESTILT	Barcelona Sofa, Barcelona Chai
Barcelona Chair: Sale Weekend		of orthogic relation.	designs.
www.PGMod.com/Barcelona-Chair - Cu Free S&H.	stomer Appreciation Sale! Sa	we 5% on Barcelona Chair	+ www.WickedElements.com
Press land Ohele Free Ohlening			Barcelona Chairs
Barcelona Chair - Free Shipping	ean imitations Our Barcelor	a Chair offers denuine	Chairs & Seats from 152+ Shops
quality	icap mitations: our Darocio r	a onan onoio genane	www.Calibex.com
Barcelona Chairs			
BizRate.com - We Offer 2,500+ Chair Ch	Barcelona Chair - \$659.99		
Olassia Bassalana Ohair Os Oala	* ***		Free Shipping
Classic Barcelona Chair On Sale \$899 funkceata com. Al colore available. The Barcelona Chair is a classic nices that			Loveseat, daybed, ottoman. Free shipping Up to 60% off
Tankysola.com - Ar colors avalable. The		o ploce main	www.modabode.com
	Later Fore		
Yahoo!s: Report bad results or ads. But	cket test: F655		Buy Barcelona Chairs
			Barcelona Chairs on Sale.
Barcelona Chair - Volo Leather	hair and Stool (1929), origin:	ully created to furnish his	www.NexTag.com/sofas
German Pavilion at the International Exhi	bition in Barcelona , have cor	ne	
www.dwr.com/productdetail.cfm?id=7200) - 17k		Barcelona Chair
Description of the Million of the second			The Right Style For Your Space.
Darcelona chair - wikipedia, the fi	ree encyclopedia	o for Baraalana Chair	Shopzilla.com/chairs

The Barcelona chair and ottoman was designed by Mies van der Rohe for ... Barcelona Chair, inspired by its predecessors, the campaign and folding chairs ...

Bidding and Spending Graphs



A "social network" with "term-document" aspects.

Uses of Bidding and Spending graphs:

- "deep" micro-market identification.
- improved query expansion.

More generally, user segmentation for behavioral targeting.

What do these networks "look" like?



Micro-markets in sponsored search

Goal: Find *isolated* markets/clusters with *sufficient money/clicks* with *sufficient coherence*. Ques: Is this even possible?



10 million keywords

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 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:

 $\Phi(k) = \min_{S \subset V, |S| = k} \phi(S)$ **A** "size-resolved" community-quality measure!

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

• Use approximation algorithms to approximate it (even better than exactly)

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.

Approximation algorithms as experimental probes?

The usual modus operandi for approximation algorithms:

- define an objective, the numerical value of which is intractable to compute
- develop approximation algorithm that returns approximation to that number

• graph achieving the approximation may be unrelated to the graph achieving the exact optimum.

But, for randomized approximation algorithms with a geometric flavor (e.g. matrix algorithms, regression algorithms, eigenvector algorithms; duality algorithms, etc):

- often can approximate the vector achieving the exact solution
- randomized algorithms compute an ensemble of answers -- the details of which depend on choices made by the algorithm
- maybe compare different approximation algorithms for the same problem.

Analogy: What does a protein look like?



Three possible representations (all-atom; backbone; and solvent-accessible surface) of the three-dimensional structure of the protein triose phosphate isomerase.

Experimental Procedure:



- Generate a bunch of output data by using the unseen object to filter a known input signal.
- Reconstruct the unseen object given the output signal and what we know about the artifactual properties of the input signal.

Low-dimensional and small social networks



d-dimensional meshes









Newman's Network Science



RoadNet-CA

NCP for common generative models



What do large networks look like?

Downward sloping NCPP

small social networks (validation)

"low-dimensional" networks (intuition)

hierarchical networks (model building)

existing generative models (incl. community models)

Natural interpretation in terms of isoperimetry

implicit in modeling with low-dimensional spaces, manifolds, k-means, etc.

Large social/information networks are very very different

We examined more than 70 large social and information networks We developed principled methods to interrogate large networks Previous community work: on small social networks (hundreds, thousands)





Focus on the red curves (local spectral algorithm) - blue (Metis+Flow), green (Bag of whiskers), and black (randomly rewired network) for consistency and cross-validation.

"Whiskers" and the "core"

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





If remove whiskers, then the lowest conductance sets (the "best" communities) are "2-whiskers":



How do we know this plot it "correct"?

Lower Bound Result

Spectral and SDP lower bounds for large partitions

Modeling Result

Very sparse Erdos-Renyi (or PLRG wth $\beta \epsilon$ (2,3)) gets imbalanced deep cuts

Structural Result

Small barely-connected "whiskers" responsible for minimum

• Algorithmic Result

Ensemble of sets returned by different algorithms are very different

Spectral vs. flow vs. bag-of-whiskers heuristic

Spectral method implicitly regularizes, gets more meaningful communities

Random graphs and forest fires

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. $p_{ij} = w_i w_j / \sum_k w_k.$

A "power law random graph" model (Chung-Lu)





A "forest fire" model (LKF05)



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:





A few general thoughts

Regularization is typically *implemented* by adding a norm constraint

• makes the problem harder (think L1-regularized L2-regression).

Approximation algorithms for intractable graph problems *implicitly* regularize

- relative to combinatorial optimum
- incorporate empirical signatures of bias-variance tradeoff.

Use statistical properties *implicit* in worst-case algorithms to provide insights into informatics graphs

• good since networks are large, sparse, and noisy.

A "claimer" and a "disclaimer":

- Today, mostly took a "10,000 foot" view:
 - But, "drilled down" on two specific examples that illustrate "algorithmicstatistical" interplay in a novel way
- Mostly avoided* "rubber-hits-the-road" issues:
 - Multi-core and multi-processor issues
 - Map-Reduce and distributed computing
 - Other large-scale implementation issues





*But, these issues are very much a motivation and "behind-the-scenes" and important looking forward!

Conclusion

- "Algorithmic" and "statistical" perspectives on data problems
- Genetics application
 - DNA SNP analysis --> choose columns from a matrix
- Internet application

Community finding --> partitioning a graph

In many large-scale data applications, "algorithmic" and "statistical" perspectives interact in fruitful ways.

MMDS Workshop on "Algorithms for Modern Massive Data Sets"

(http://mmds.stanford.edu)

at Stanford University, June 15–18, 2010

Objectives:

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

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

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

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

Registration will be available soon!