

The Spectre of the Spectrum

An empirical study of the spectra of large networks

David F. Gleich Purdue University

Indiana University Bloomington, IN

Supported by Sandia's John von Neumann postdoctoral fellowship and the DOE Office of Science's Graphs project.

Thanks to Ali Pinar, Jaideep Ray, Tammy Kolda, C. Seshadhri, Rich Lehoucq @ Sandia and Jure Leskovec and Michael Mahoney @ Stanford for helpful discussions.

Spectral density plots



David Gleich (Purdue)

There's information inside the spectra



These figures show the normalized Laplacian. Banerjee and Jost (2009) also noted such shapes in the spectra.

David Gleich (Purdue)

Overview

- Graphs and their matrices
- Data for our experiments
- Issues with computing spectra
- Many examples of graph spectra
- A curious property around the eigenvalue one
- Computing spectra for large networks
- **Ongoing studies**

Why are we interested in the spectra?

Modeling



Network Comparison

Fay et al. 2010 – Weighted Spectral Density

The network is as 19971108 from Jure's snap collect (a few thousand nodes) and we insert random connections from 50 nodes
David Gleich (Purdue) IU Seminar 5/51

Matrices from graphs

Adjacency matrix

 $\mathbf{A} : n \times n, \mathbf{A} = \mathbf{A}^{T}$ $A_{i,j} = 1 \text{ if } (i,j) \in E$ $-d_{\max} \le \lambda(\mathbf{A}) \le d_{\max}$

Laplacian matrix D = diag(Ae) L = D - A $0 \le \lambda(L) \le 2d_{max}$

Normalized Laplacian matrix $\tilde{L} = D^{-1/2}LD^{-1/2} = I - D^{-1/2}AD^{-1/2}$ $0 \le \lambda(\tilde{L}) \le 2$

Random walk matrix $\boldsymbol{P} = \boldsymbol{D}^{-1} \boldsymbol{A}$

Modularity matrix $\mathbf{d} = \mathbf{A}\mathbf{e}$ $\mathbf{M} = \mathbf{A} - 1/(2|E|)\mathbf{d}\mathbf{d}^{T}$

Not covered Signless Laplacian matrix

Incidence matrix (It is incidentally discussed)

Seidel matrix

Heat Kernel

Everything is undirected. Mostly connected components only too.

David Gleich (Purdue)

Erdős–Rényi Semi-circles

Based on Wigner's semi-circle law.

The eigenvalues of the adjacency matrix for n=1000, averaged over 10 trials

Semi-circle with outlier if average degree is large enough.



Observed by Farkas and in the book "Network Alignment" edited by Brandes (Chapter 14)

Previous results

Farkas et al. Significant deviation from the semicircle law for the adjacency matrix

Mihail and Papadimitriou Leading eigenvalues of the adjacency matrix obey a power-law based on the degree-sequence

Chung et al. Normalized Laplacian still obeys a semi-circle law if min-degree large

Banerjee and Jost Study of types of patterns that emerge in evolving graph models – explain many features of the spectra

David Gleich (Purdue)





Encyclopedia of Mathematics and its Applications 66

EIGENSPACES OF GRAPHS

Dragoš Cvetković, Peter Rowlinson, Slobodan Simić

An Introduction to the Theory of Graph Spectra

DRAGOŠ CVETKOVIĆ, PETER ROWLINSON and SLOBODAN SIMIĆ



In comparison to other empiric studies

We use "exact" computation of spectra, instead of approximation.

We study "all" of the standard matrices over a range of large networks.

Our "large" is bigger.

We look at a few random graph models preferential attachment random powerlaw copying model forest fire model

David Gleich (Purdue)

ISSUES WITH COMPUTING SPECTRA

Why you should be very careful with eigenvalues.

Matlab!

Always a great starting point. My desktop has 24GB of RAM (less than \$2500 now!)

24GB/8 bytes (per double) = 3 billion numbers \sim 50,000-by-50,000 matrix

Possibilities

D = eig(A) - needs twice the memory for A,D [V,D] = eig(A) - needs three times the memory for A,D,V

These limit us to \sim 38000 and \sim 31000 respectively.

Bugs – Matlab

eig(A) Returns incorrect eigenvectors

Seems to be the result of a bug in Intel's MKL library. Fixed in R2011a

Bug – ScaLAPACK default

sudo apt-get install scalapack-openmpi Allocate 36000x36000 local matrix

Run on 4 processors

Code crashes

Bug – LAPACK

Scalapack MRRR

- Compare standard lapack/blas to atlas performance
- Result: correct output from atlas
- Result: incorrect output from lapack
- Hypothesis: lapack contains a known bug that's apparently in the default ubuntu lapack

Moral

Always test your software. **Extensively.**



EXAMPLES

Data sources

SNAP	Various	100s-100,000s	
SNAP-p2p	Gnutella Network	5-60k, ~30 inst.	
SNAP-as-733	Autonomous Sys.	~5,000, 733 inst.	
SNAP-caida	Router networks	~20,000, ~125 inst.	
Pajek	Various	100s-100,000s	
Models	Copying Model	1k-100k 9 inst. 324 gs	
	Pref. Attach	1k-100k 9 inst. 164 gs	
	Forest Fire	1k-100k 9 inst. 324 gs	
Mine	Various	2k-500k	
Newman	Various		
Arenas	Various		
Porter	Facebook	100 schools, 5k-60k	
IsoRank, Natalie	Protein-Protein	<10k , 4 graphs	
	1 occur i loccur		

Thanks to all who make data available

Big graphs

Arxiv	86376	1035126	Co-author
Dblp	93156	356290	Co-author
Dictionary(*)	111982	2750576	Word defns.
Internet(*)	124651	414428	Routers
ltdk0304	190914	1215220	Routers
p2p-gnu(*)	62561	295756	Peer-to-peer
Patents(*)	230686	1109898	Citations
Roads	126146	323900	Roads
Wordnet(*)	75606	240036	Word relation
web-nb.edu(*)	325729	2994268	Web

(*) denotes that this is a weakly connected component of a directed graph.

A \$8,000 matrix computation



925 nodes and 7400 processors on Redsky for 10 hours normalized Laplacian matrix

David Gleich (Purdue)

Indiana's Facebook Network



Data from Mason Porter. Aka, the start of a \$50,000,000,000 graph.



These are cases where we have multiple instances of the same graph.

Already known?



Already known?



I soon realized I was searching for "spectre" instead of spectrum, oops.

Spikes?

Unit eigenvalue $(\mathbf{I} - \mathbf{D}^{-1}\mathbf{A})\mathbf{x} = \mathbf{x} \Rightarrow \mathbf{A}\mathbf{x} = 0$

Repeated rows

Identical rows grow the null-space.

Banerjee and Jost

Motif doubling and joining small graphs will tend to cause repeated eigenvalues and null vectors.



Banerjee and Jost explained how evolving graphs should produce repeated eigenvalues

David Gleich (Purdue)

Combining Eigenvalues

If A has an eigenvector with a zero component, then





"A + B" (as in the figure) has the same eigenvalue with eigenvector extended with zeros on B.



Bannerjee and Jost observed this for the normalized Laplacian.



David Gleich (Purdue)

Classic spectra



From NASA: http://imagine.gsfc.nasa.gov/docs/science/how II/spectral what.html

Random power law

Random power law 12500 vertices, 500 (2.*) / 400 (1.8) min degree

Generate a power law degree distribution.

Produce a random graph with a prescribed degree distribution using the Bayati-Kim-Saberi procedure.



Preferential Attachment

Start graph with a k-node clique. Add a new node and connect to k random nodes, chosen proportional to degree.



Copying model

Start graph with a k-node clique. Add a new node and pick a parent uniformly at random. Copy edges of parent and make an error with probability α



Obvious follow up here: does a random sample with the same degree distribution show the same thing?

David Gleich (Purdue)

Forest Fire models

Start graph with a k-node clique. Add a new node and pick a parent uniformly at random. Do a random "bfs'/"forest fire" and link to all nodes "burned"



Reality vs. graph models

Real spectra vs model spectra with over 25000 vertices



Where is this going?

We can compute spectra for large networks if needed.

Study relationship with known powerlaws in spectra

Eigenvector localization

Directed Laplacians



Just the degree distribution? No



David Gleich (Purdue)

Facebook is not a copying model





David Gleich (Purdue)

Why is one separated sometimes?



Separation in copying, forest-fire



Note, the axes on these fiures aren't comparable – different plotting scales – but the shapes are.

David Gleich (Purdue)



Strong separation in random powerlaw



Not edge-density alone



Not edge-density alone



David Gleich (Purdue)

COMPUTING SPECTRA OF LARGE NETWORKS



Redsky, Hopper I, Hopper II, and a Cielo testbed. Details if time.

Eigenvalues with ScaLAPACK

Mostly the same approach as in LAPACK

- Reduce to tridiagonal form (most time consuming part)
- 2. Distribute tridiagonals to all processors
- 3. Each processor finds all eigenvalues
- 4. Each processor computes a subset of eigenvectors



ScaLAPACK's 2d block cyclic storage

I'm actually using the **MRRR algorithm**, where steps 3 and 4 are better and faster

MRRR due to Parlett and Dhillon; implemented in ScaLAPACK by Christof Vomel.

David Gleich (Purdue)

Estimating the density directly



A and $\mathbf{F}^{T}\mathbf{AF}$ have the same eigenvalue inertia if \mathbf{F} is nonsingular.

Eigenvalue inertia = (p,n,z)Positive eigenvalues Negative eigenvalues Zero eigenvalues If $\mathbf{F}^T \mathbf{A} \mathbf{F}$ is diagonal, inertia is easy to compute

 $\begin{array}{l} \tilde{\boldsymbol{L}} & \text{has inertia (n-1,0,1)} \\ \tilde{\boldsymbol{L}} - \lambda_e \boldsymbol{I} & \text{has inertia} \\ (\text{sum}(\lambda > \lambda_e), \, \text{sum}(\lambda < \lambda_e), \ldots) \end{array}$

This is an old trick in linear algebra. I know that Fay et al. used it in their weighted spectral density.

David Gleich (Purdue)

Alternatives

Use ARPACK to get extrema

Use ARPACK to get interior around λ_0 via the folded spectrum $((\mathbf{A} - \lambda_0)^2)^k$



Farkas et al. used this approach. Figure from somewhere on the web... sorry!

David Gleich (Purdue)

Adding MPI tasks vs. using threads

Most math libraries have threaded versions (Intel MKL, AMD ACML) Is it better to use threads or MPI tasks?

It depends.

Cray libsci Intel MKL Threads **Ranks** Time Threads Ranks Time-E **Time-T** 1 64 1412.5 36 1271.4 339.0 1 16 4 1881.4 9 1058.1 456.6 4 16 Omitted. 4

Normalized Laplacian for 36k-by-36k co-author graph of CondMat

David Gleich (Purdue)

Weak Parallel Scaling



David Gleich (Purdue)



QUESTIONS



Code will be available eventually. Image from good financial cents.

David Gleich (Purdue)