# What is network science?
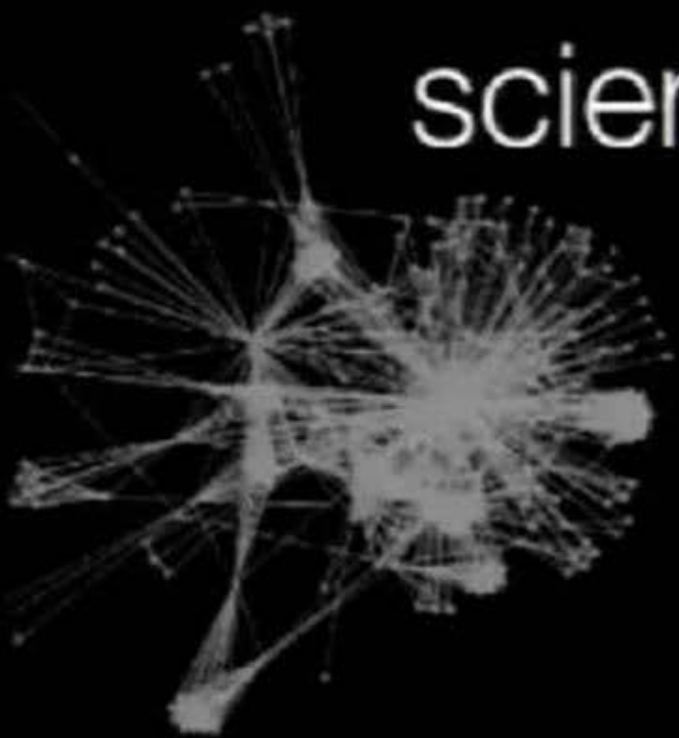
David Gleich · Purdue

# Network Science

the study of **network representations** of physical, biological, and social phenomena leading to **predictive models** of these phenomena
National Research Council (via Wikipedia)

- Models

- Algorithms

- Data

Network Science is
*CS&E applied to graphs*
Me

# Dimensionality of social networks using motifs and eigenvalues

Facebook network    LinkedIn network

HAND-ESU method
for graphless

Estimate mGEO-P
parameters $\alpha$, $\beta$

Spectral histogram of
normalized Laplacian

prediction          training

Support vector
machine classifier

choices          target

Minimum
KL divergence

prediction

best match

Estimated dimension

Estimated dimension

mGEO-P networks of
various dimension

$G_1$        $G_2$        $G_3$        $G_4$        $G_5$        $G_6$        $G_7$        $G_8$

# Dimensionality of social networks using motifs and eigenvalues

# Local methods in
# network science

David F. Gleich

Purdue University

David Gleich · Purdue

5

# Local methods identify small, meaningful regions in massive networks



Co-purchasing in Amazon
330k Vertices, 1M edges

Ground-truth set has
~10 vertices

Two local methods
- pprgrow
- hkgrow

find most of the true set starting from a seed inside the set.

David Gleich · Purdue

# Local methods help characterize the graph!



## Vertex neighborhood or egonet
The induced subgraph of a vertex its neighbors

Egonets of social networks should show "structural holes" [Burt95,Kleinberg08].

Used for anomaly detection [Akoglu10],
community seeds [Huang11.Schaeffer11],
overlapping communities [Schaeffer07.Rees10].

# Characterizing local anomalies in graphs using egonets



OREGON

Verizon
Sprint
AT&T

Number of edges

$1.6636x + (-0.23992) = y$
$1.0737x + (-0.15151) = y$
$2.0737x + (-0.45254) = y$

Number of nodes

From OddBall: Spotting Anomalies in Weighted Graphs
Akoglu et al. PAKDD2010

From Knowledge Sharing and Yahoo Answers
Adamic et al. WWW2008

David Gleich · Purdue

## Our perspective

Local diffusions are some of the best community detection algorithms available!

# Local Community Detection

Given seed(s) *S* in *G*, find a community that contains *S*.

seed

## Community
A set of vertices with high internal and low external connectivity

# Low-conductance sets are communities

$$\text{conductance}(T) = \frac{\text{\# edges leaving } T}{\text{\# edge endpoints in } T}$$

$$= \text{" chance a random step exits } T \text{ "}$$



conductance(comm) =
39/381 = .102

## How to find these ?

# Graph diffusions find low-conductance sets

A diffusion propagates "rank" from a seed across a graph.



seed

● = high
● = low
} diffusion value

↳ = local community /
low-conductance set

But don't diffusions go everywhere in the graph?

# The most used diffusions stay localized even in massive graphs

$$(I - \beta P)x = (1 - \beta)s$$

plot(**x**)                    nnz(**x**) $\approx$ 800*k*



Crawl of flickr from 2006 ~800k nodes, 6M edges, beta=1/2

# The most used diffusions stay localized even in massive graphs

$$(I - \beta P)x = (1 - \beta)s$$



plot(x)    nnz(x) ≈ 800k

Crawl of flickr from 2006 ~800k nodes, 6M edges, beta=1/2

David Gleich · Purdue

# Our mission

Find the solution with work
roughly proportional to the
*localization*, not the matrix.

# Our Point

Coordinate relaxation methods yield localized algorithms for diffusions
in a pleasingly wide variety of settings.

# Our Results

New empirical and theoretical insights into *why* and *how* a specific form is so effective.

# Localized methods for diffusions use the push coordinate relaxation strategy

**The push method**

Coordinate relaxation
for $A\,x = b$

Update $\quad \mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \rho_j \mathbf{e}_j$

such that $[\mathbf{A}\mathbf{x}^{(k+1)}]_j = [\mathbf{b}]_j$

or $\qquad [\mathbf{A}\mathbf{x}^{(k+1)}]_j = [\mathbf{b}]_j + \varepsilon_j$

Used in coordinate descent, Gauss-Seidel,
Gauss-Southwell, and many other methods.

David Gleich · Purdue

# Localized methods for diffusions use the push coordinate relaxation strategy

## Push on a graph-based linear system
has a super-duper awesome property

The
Push
Method
for PPR
on a
graph

$\varepsilon, \rho$

1. $\mathbf{x}^{(1)} = 0, \mathbf{r}^{(1)} = (1 - \beta)\mathbf{e}_i, k = 1$

2. *while any* $r_j > \varepsilon d_j$     ($d_j$ *is the degree of node j*)

3. $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + (r_j - \varepsilon d_j \rho)\mathbf{e}_j$

4. $r_i^{(k+1)} = \begin{cases} \varepsilon d_j \rho & i = j \\ r_j^{(k)} + \beta(r_j - \varepsilon d_j \rho)/d_j & i \sim j \\ r_i^{(k)} & \text{otherwise} \end{cases}$

5. $k \leftarrow k + 1$

David Gleich · Purdue

# Push is fast!

PageRank
$(I - \alpha P)x$
$= (1 - \alpha)e_s$

Katz
$(I - \alpha A)x$
$= (1 - \alpha)e_s$

For the PageRank diffusion, Push
gives constant work (entry-wise)
Andersen, Chung, Lang FOCS 2006

1. For the Katz diffusion
   Push works empirically fast
   Bonchi, Gleich, et al., 2012, internet Math

2. For the exponential $x = \exp(P)e_s$
   Push gives uniform localization
   on power-law graphs and fast
   runtimes
   Gleich and Kloster 2014, Kloster Gleich

3. For the heat-kernel diffusion
   Push gives constant work
   (entry-wise) $x = \exp(tP)e_s$
   Kloster and Gleich 2014, KDD

4. For the PageRank diffusion
   Push yields sparsity
   regularization
   Gleich and Mahoney ICML 2014

5. For a general class of diffusions
   There is a Cheeger inequality
   like before
   Ghosh, Teng, et al. KDD 2014

6. For the PageRank diffusion
   Push gives the solution path in
   constant work (entry-wise)
   Kloster and Gleich arXiv 2015 WAW

# Push is useful!



1. Push implicitly regularizes semi-supervised learning
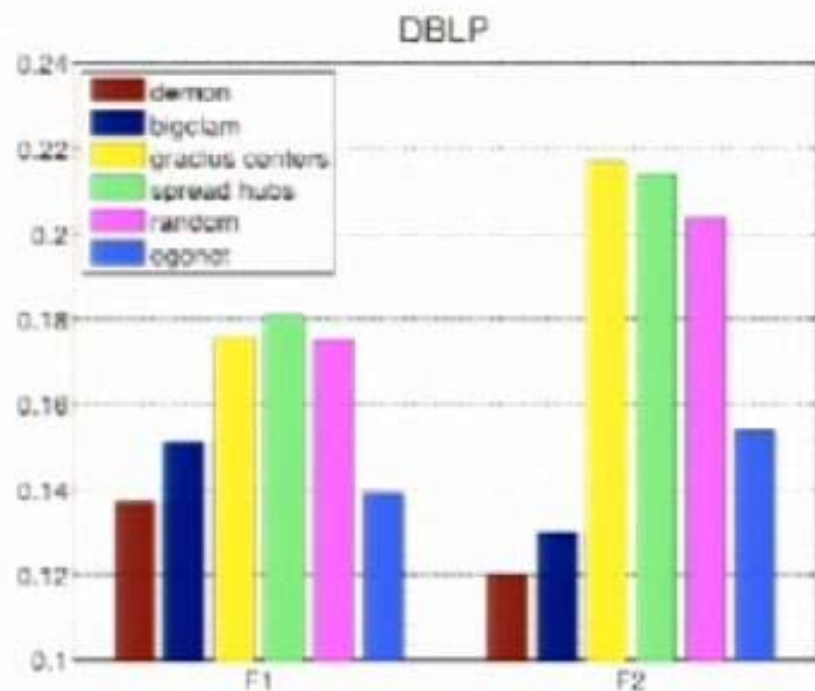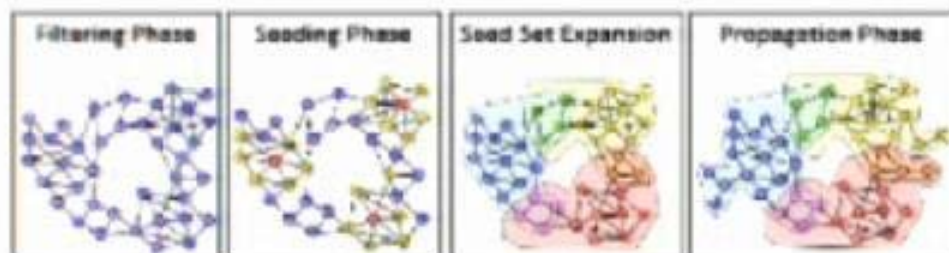   *Gleich and Mahoney; submitted*

2. Push gives state of the art results for overlapping community detection
   *Whang, Gleich, Dhillon, CIKM 2013*
   *Whang, Gleich, Dhillon, in prep.*

3. Push for overlapping clusters decrease communication in parallel solutions
   *Andersen, Gleich, Mirrokni, WSDM 2012*

# Heat kernel localization

## General recipe

1. Take problem X, convert into a linear system

2. Apply "push" to that linear system

3. Analyze and bound total work

## Heat kernel recipe

1. Convert $X = \exp(tP)e_i$ into

$$
\begin{bmatrix}
I & & & & \\
-tP/1 & I & & & \\
& -tP/2 & \ddots & & \\
& & \ddots & I & \\
& & & -tP/N & I
\end{bmatrix}
\begin{bmatrix}
v_0 \\ v_1 \\ \vdots \\ \vdots \\ v_N
\end{bmatrix}
=
\begin{bmatrix}
e_i \\ 0 \\ \vdots \\ \vdots \\ 0
\end{bmatrix}
$$

2. Apply "push"

3. Analyze work bound

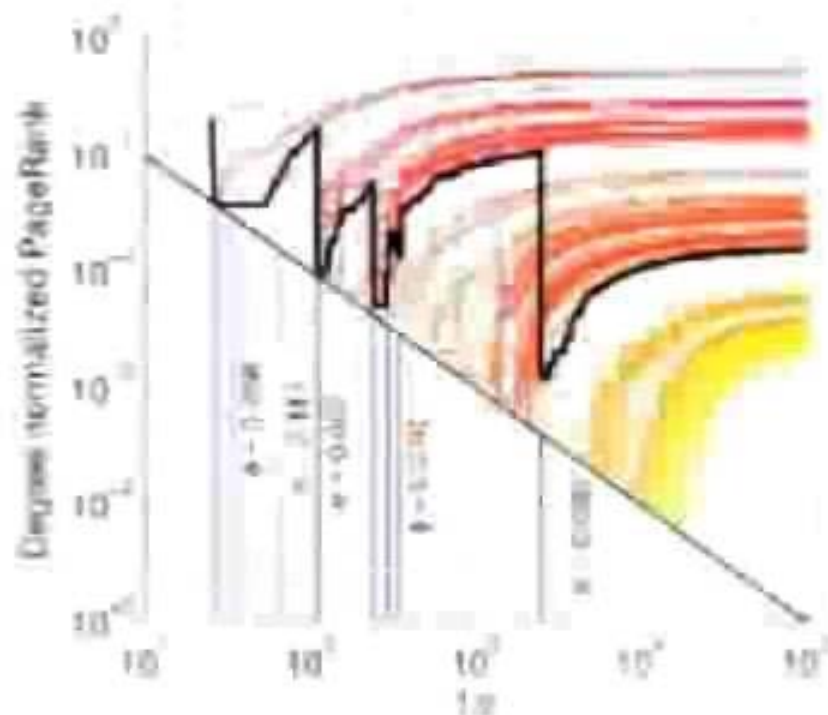| data | $F_1$ | | precision | | set size | | comm |
|---|---|---|---|---|---|---|---|
| | HK | PR | HK | PR | HK | PR | size |
| amazon | 0.325 | 0.140 | 0.244 | 0.107 | 193 | 15293 | 495 |
| dblp | 0.257 | 0.115 | 0.208 | 0.081 | 44 | 16026 | 1429 |
| youtube | 0.177 | 0.136 | 0.135 | 0.098 | 1010 | 6079 | 1615 |
| lj | 0.131 | 0.107 | 0.102 | 0.086 | 283 | 738 | 662 |
| orkut | 0.055 | 0.044 | 0.036 | 0.031 | 537 | 1989 | 4526 |
| friendster | 0.078 | 0.090 | 0.066 | 0.075 | 229 | 333 | 724 |



Seed

hkgrow community

hkgrow mistakes

PR achieves high recall by "guessing" a huge set

HK identifies a tighter cluster, so attains better precision

David Gleich · Purdue

26

# PageRank solution paths

$x = 0.0219068$

Compute one diffusion, and all sweep-cuts, for all values of epsilon

# An algorithm to find overlapping communities using local diffusions

1. Extract part of the graph that might have overlapping communities.

2. Compute a partitioning of the network into many pieces (think sqrt(n)) using Graclus.

3. Find the center of these partitions.

4. Use "push" seeded with the egonets of these partitions

5. Add back any missing pieces.

# Recap

- Local methods give rapid insight into massive graphs

- Seeded diffusions are usually localized

QUESTIONS?