

Local diffusion algorithms for fast, personalized graph applications

David F. Gleich
CS Department

Kyle Kloster
Math Department

supported by
NSF CAREER CCF-1149756,
IIS-IS-1422918, DARPA SIMPLEX,

PURDUE
UNIVERSITY

Brief outline

Introduction & Application 20-30 min

Coordinate Relaxation for Strong Convergence 30 min

Break 15 min

Katz Diffusion 10 min

Weak Convergence for PageRank 20 min

Monte Carlo methods 15 min

Break 15 min

Implicit Regularization 25 min

Discussion 15 min

Brief outline

Introduction & Application 20-30 min

Coordinate Relaxation for Strong Convergence 30 min

Break 15 min

Katz Diffusion 10 min

Weak Convergence for PageRank 20 min

Monte Carlo methods 15 min

Break 15 min

Implicit Regularization 25 min

Discussion 15 min

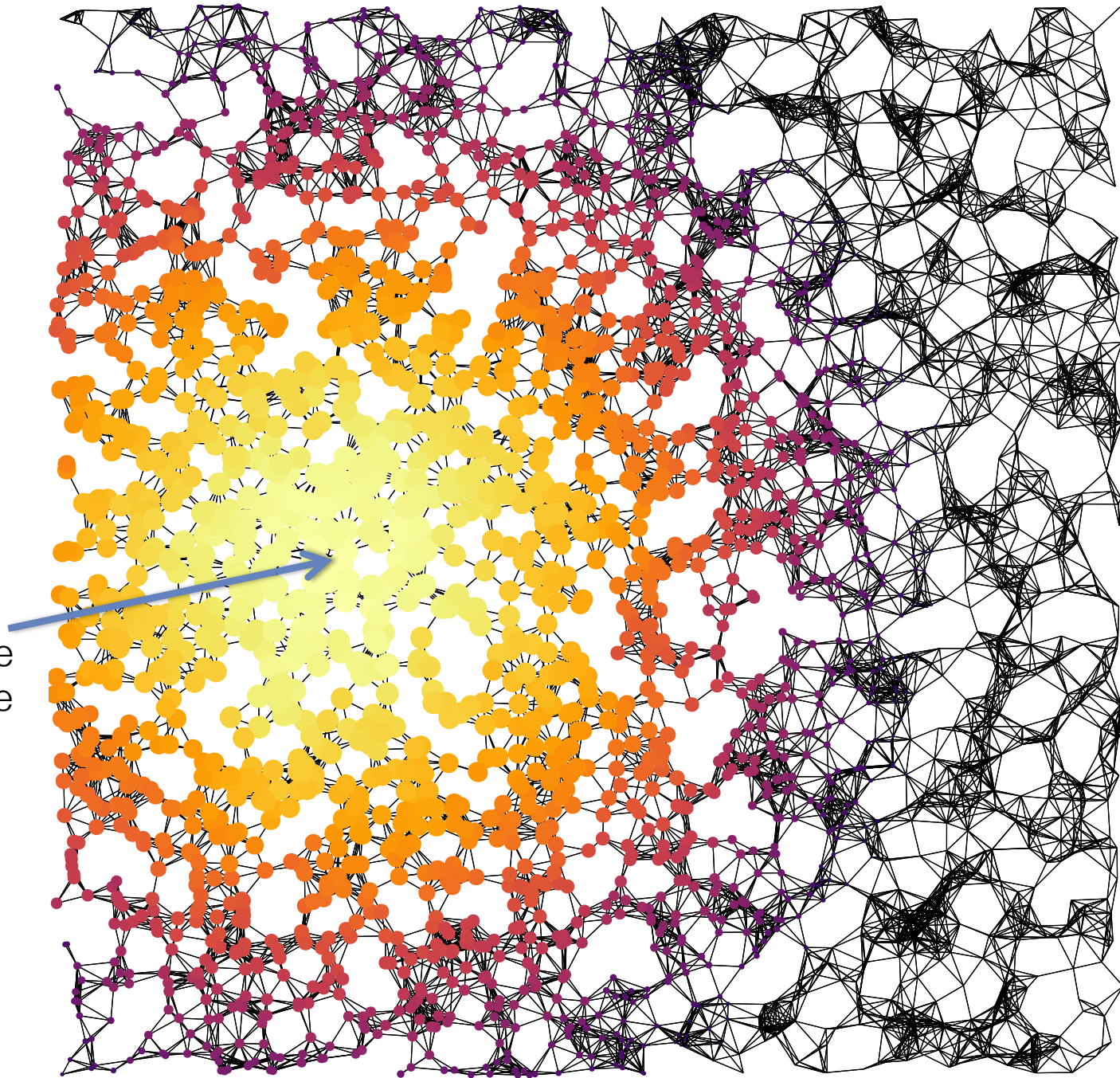


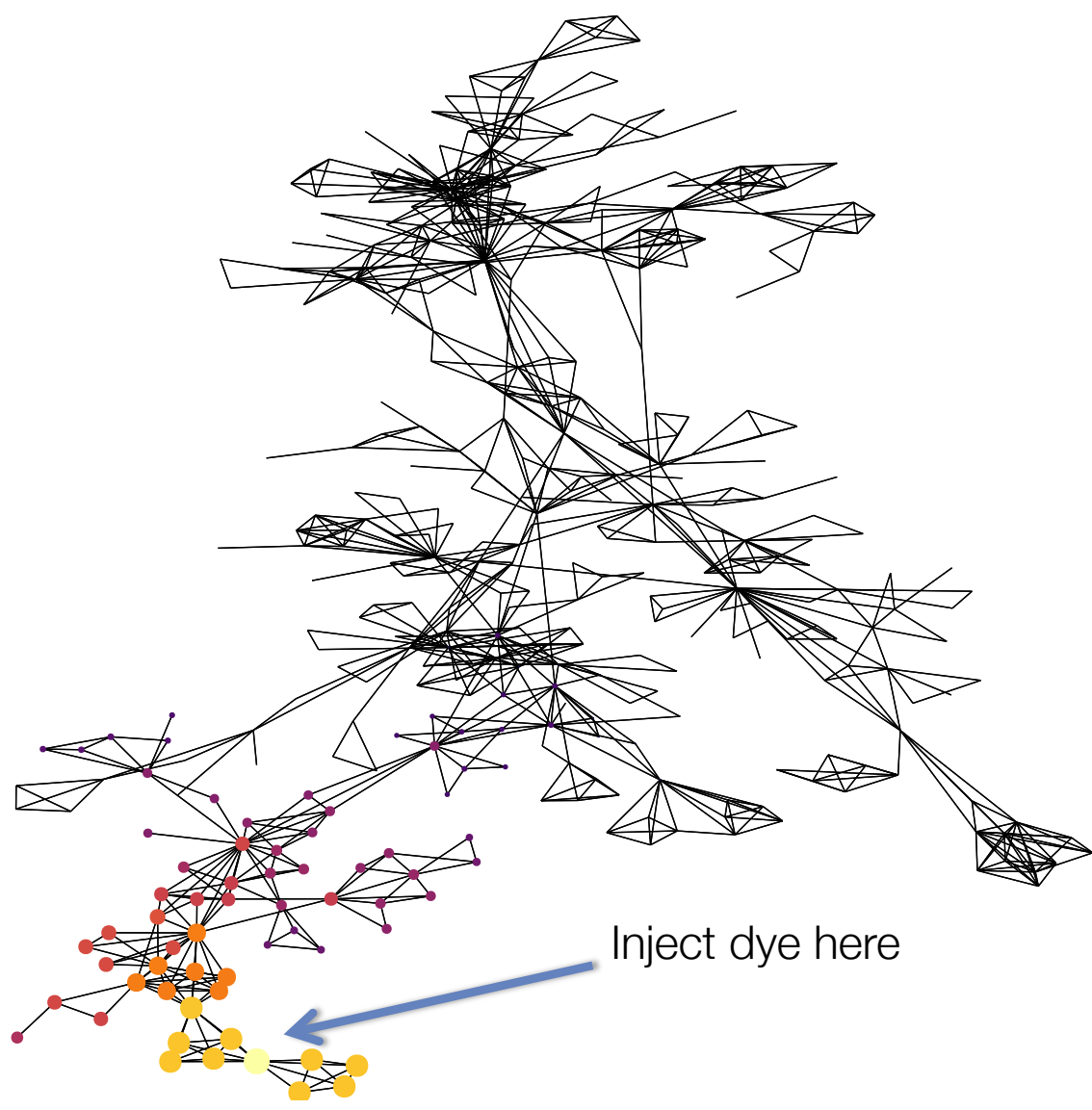


Wikipedia: Angiography

Let $G = (V, E) \dots$

Inject dye
here





Inject dye here

... work out sample diffusion on the board ...

Examples

Everything in the world can be explained by a matrix, and we see how deep the rabbit hole goes

The talk ends, you believe -- whatever you want to.



CHOOSE

Image from rockysprings, deviantart, CC share-alike



Gene
Golub



Charles
van Loan

MATRIX COMPUTATIONS

Matrices derived from networks

adjacency	\mathbf{A}		
diag degree matrix	\mathbf{D}	$d_{ii} = \text{deg}$	
random walk	\mathbf{P}	$= \mathbf{A}^T \mathbf{D}^{-1}$	Directed
Laplacian	\mathbf{L}	$= \mathbf{D} - \mathbf{A}$	
normalized Laplacian	\mathcal{L}	$= \mathbf{D}^{-1/2} (\mathbf{D} - \mathbf{A}) \mathbf{D}^{-1/2}$ $= \mathbf{I} - \mathbf{D}^{-1/2} \mathbf{A} \mathbf{D}^{-1/2}$	Undirected
normalized adjacency	\mathcal{A}	$= \mathbf{D}^{-1/2} \mathbf{A} \mathbf{D}^{-1/2}$	

A general diffusion definition

You'll see this again,
Kyle has more.

$$\mathbf{f} = \sum_{k=0}^{\infty} c_k \mathbf{M}^k \mathbf{s} \text{ Source}$$

PageRank

$$\mathbf{f} = \sum_{k=0}^{\infty} (1 - \alpha) \alpha^k \mathbf{P}^k \mathbf{s}$$

$$(\mathbf{I} - \alpha \mathbf{P}) \mathbf{f} = (1 - \alpha) \mathbf{s}$$

Katz scores

$$\mathbf{f} = \sum_{k=0}^{\infty} (1 - \alpha) \alpha^k \mathbf{A}^k \mathbf{s}$$

$$(\mathbf{I} - \alpha \mathbf{A}) \mathbf{f} = (1 - \alpha) \mathbf{s}$$

Heat kernel

$$\mathbf{f} = \sum_{k=0}^{\infty} e^{-t} \frac{t^k}{k!} \mathbf{A}^k \mathbf{s}$$

$$\mathbf{f} = \exp\{-t(\mathbf{I} - \mathbf{P})\} \mathbf{s}$$

... not just me ...

Adjacency Kernels

$$F_{\text{EXP}}(\mathbf{A}) = \exp(\alpha \mathbf{A})$$

$$F_{\text{EXP}}(\mathcal{A}) = \exp(\alpha \mathcal{A})$$

$$F_{\text{NEU}}(\mathbf{A}) = (\mathbf{I} - \alpha \mathbf{A})^{-1}$$

$$F_{\text{NEU}}(\mathcal{A}) = (\mathbf{I} - \alpha \mathcal{A})^{-1}$$

Laplacian Kernels

$$F_{\text{HEAT}}(\mathbf{L}) = \exp(-\alpha \mathbf{L})$$

$$F_{\text{HEAT}}(\mathcal{L}) = \exp(-\alpha \mathcal{L})$$

$$F_{\text{COMR}}(\mathbf{L}) = (\mathbf{I} + \alpha \mathbf{L})^{-1}$$

$$F_{\text{COMR}}(\mathcal{L}) = (\mathbf{I} + \alpha \mathcal{L})^{-1}$$

$$F_{\text{COM}}(\mathbf{L}) = \mathbf{L}^+$$

$$F_{\text{COM}}(\mathcal{L}) = \mathcal{L}^+$$

... demo ...

github.com/dgleich/diffusion-tutorial

Applications

Applications of localized and personalized diffusion

Seeds \rightarrow Scores \rightarrow “Learning”

Graph Kernels

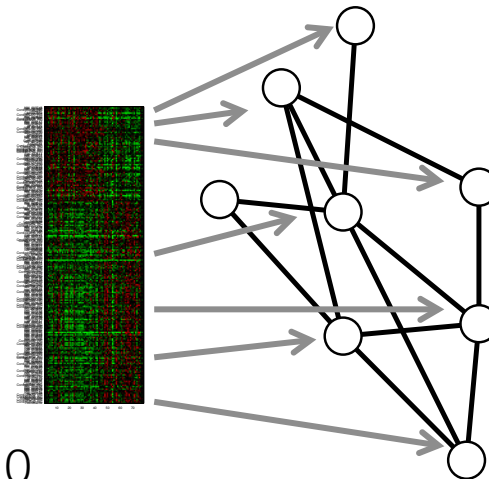
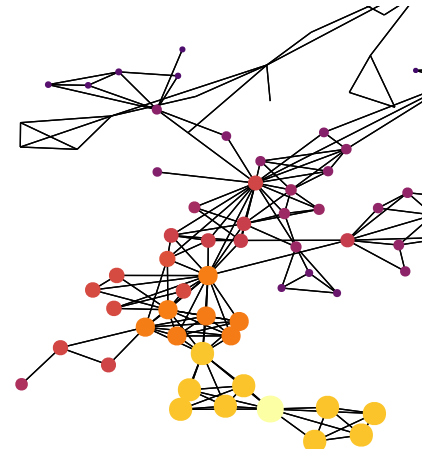
- Kondor & Lafferty, ICML 2002

Link prediction

- Liben-Nowell & Kleinberg, 2005, 2006
- Kunegis, Lommatzsch, 2009
- ... more ...

“Attribute prediction” “semi-dense vectors”

- GeneRank, Morrison et al. 2005
- ProteinRank, Freschi 2007
- Genes for cancer, Winter et al. 2012
- Cross-modal discovery, Pan et al. 2004
- Global information diffusion, Venner et al. 2010



PageRank beyond the Web

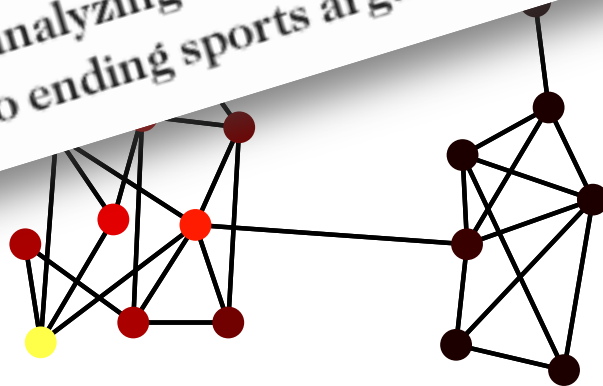
SIAM Review, Sept. 2015

$$(I - \alpha P)v$$

How Google's PageRank Quantifies Things (Like History's Best Tennis Player) Beyond The Web

by Jessica Leber
Fast Magazine

PageRank can be used for everything from analyzing the world's most important books to predicting traffic flow to ending sports arguments.

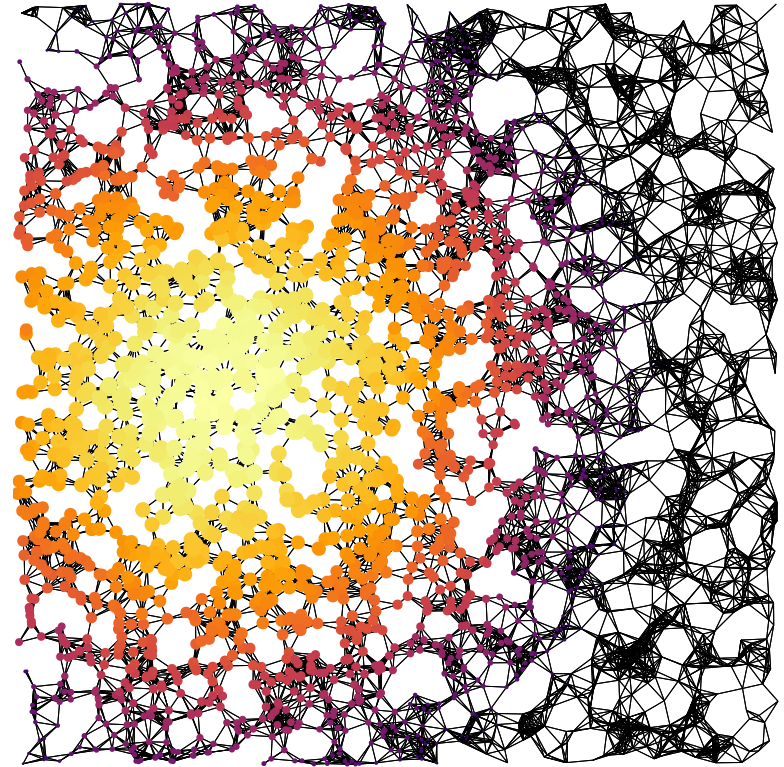
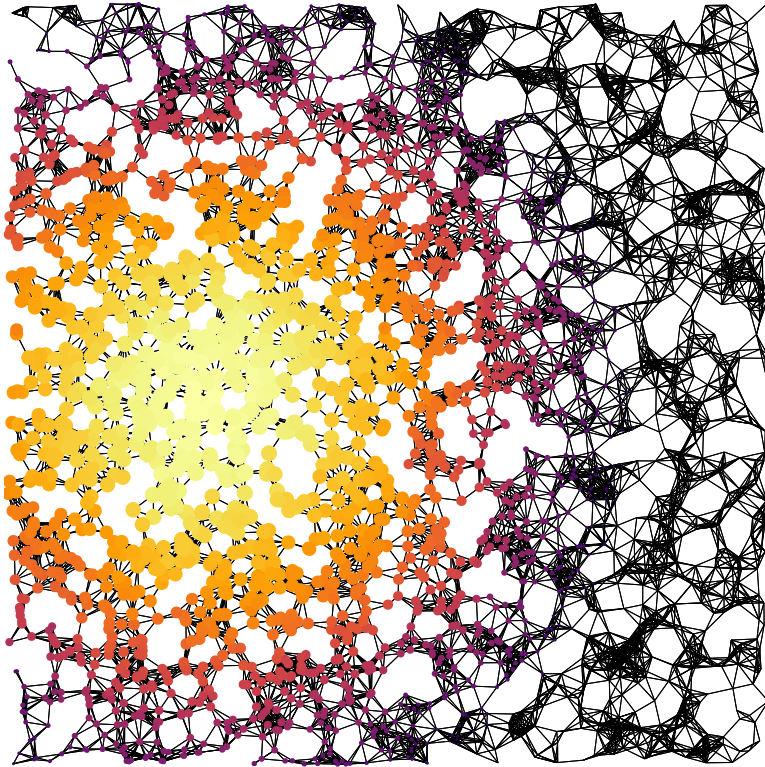


Vertex similarity

Diffusions in Open-Directory graph used for *semantic relatedness*

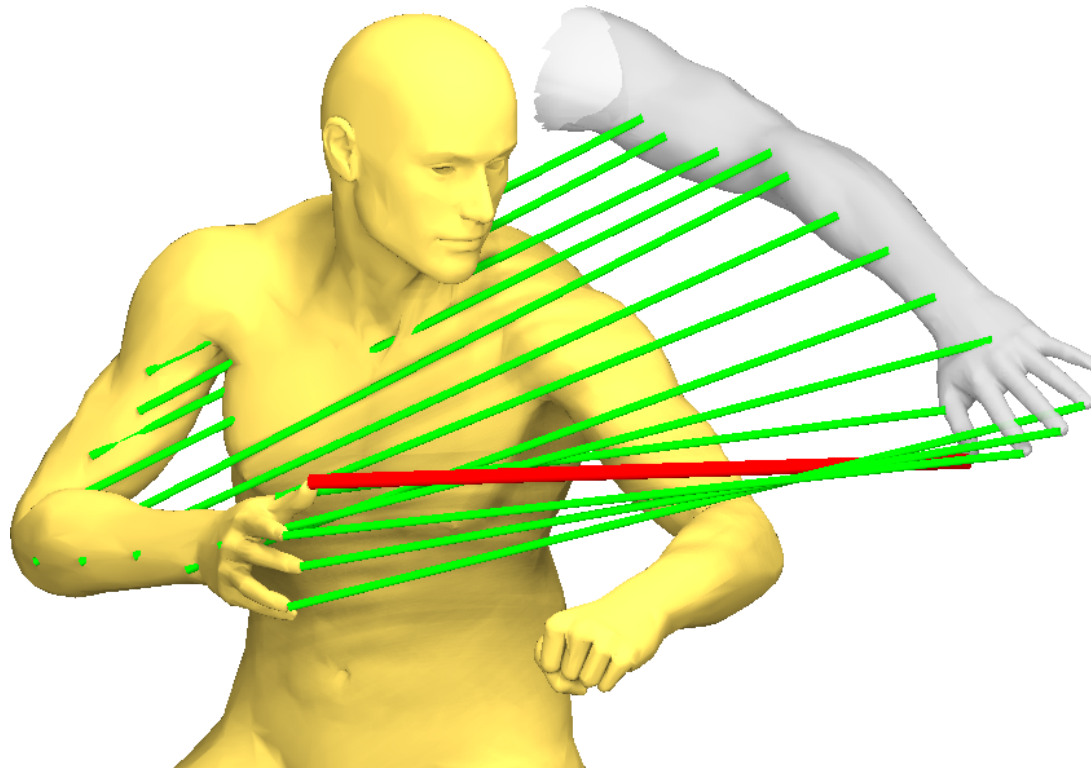
1. Reverse the direction of all edges
2. Compute the seeded PageRank matrix
e.g. “diffuse” from all individual seeds
3. Compute cosine distances between columns.

Vertex similarity



$$\|\mathbf{x} - \mathbf{y}\|_{\text{eye}} \approx 0$$

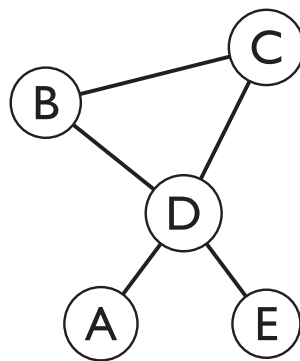
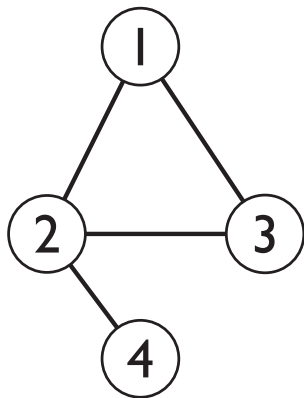
Vertex similarity



Network alignment & near isomorphisms

IsoRank

Diffuse from all potential matches through the Kronecker graph



$$P = \begin{bmatrix} 0 & 1/3 & 1/2 & 0 \\ 1/2 & 0 & 1/2 & 1 \\ 1/2 & 1/3 & 0 & 0 \\ 0 & 1/3 & 0 & 0 \end{bmatrix}$$

$$Q = \begin{bmatrix} 0 & 0 & 0 & 1/4 & 0 \\ 0 & 0 & 1/2 & 1/4 & 0 \\ 0 & 1/2 & 0 & 1/4 & 0 \\ 1 & 1/2 & 1/2 & 0 & 1 \\ 0 & 0 & 0 & 1/4 & 0 \end{bmatrix}$$

$$\begin{matrix} & A & B & C & D & E \\ \begin{matrix} 1 \\ 2 \\ 3 \\ 4 \end{matrix} & \begin{bmatrix} 0.03 & 0.05 & 0.05 & 0.09 & 0.03 \\ 0.04 & 0.07 & 0.07 & 0.15 & 0.04 \\ 0.03 & 0.05 & 0.05 & 0.09 & 0.03 \\ 0.02 & 0.03 & 0.03 & 0.05 & 0.02 \end{bmatrix} \end{matrix}$$

BIG PROBLEMS!

But structured

Dataset	Size	Nonzeros
LCSH-2	59,849	227,464
WC-3	70,509	403,960
Product graph	4,219,893,141	91,886,357,440

Opinion dynamics

\mathbf{Y}_1 initial opinions

$$(\mathbf{I} - \alpha \mathbf{P}^T) \mathbf{Y}_\infty = \beta \mathbf{Y}_1$$

\mathbf{P} gives influences

Voting in social networks

Vicious democracy

Consider a social network

- Some small fraction of users express a vote 😊
- Other users delegate their vote with decay 😐
- The others don't vote at all 😞
- Determine a final vote by taking the expected diffusion where each non-voter picks a neighbor at random and then diffuse the vote

Communities & Clusters

Theory

- *Early* Fiedler, Anderson & Morley 1985, Mihail, Chung, Pothen et al., Simon et al., Lovász & Simonovits, FOCS 1990, Random. Struct. Alg. 1993
- Spielman & Teng, 2004, 2013
- Andersen, Chung, Lang, FOCS 2006
- Chung, PNAS 2007
- Ghosh et al. KDD 2014

Practice

- Andersen & Lang, WWW 2006
- Leskovec et al. Internet Math. 2009
- Gargi et al. 2011 (Google, YouTube communities)
- Epasto et al. 2014 (Google, Competing advertisers)
- ... so many ...

Andersen- Chung-Lang personalized PageRank community theorem

[Andersen et al. 2006]

Informally

Suppose the seeds are in a set of good conductance, then the personalized PageRank method *will find* a set with conductance that's nearly as good.

... also, it's really fast.

Examples

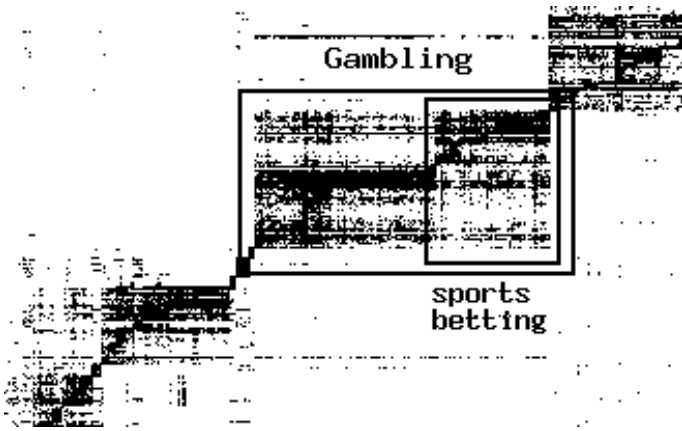


Figure 4: A low-resolution view of part of a very small version of the Yahoo sponsored search bipartite incidence matrix. The gambling co-cluster contains a sports betting subcluster.

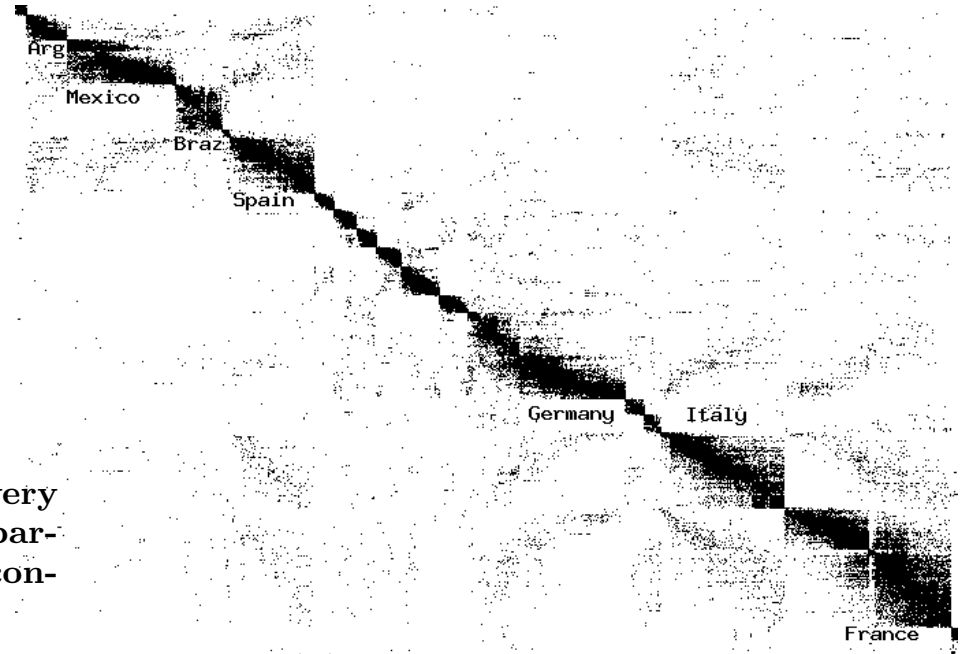
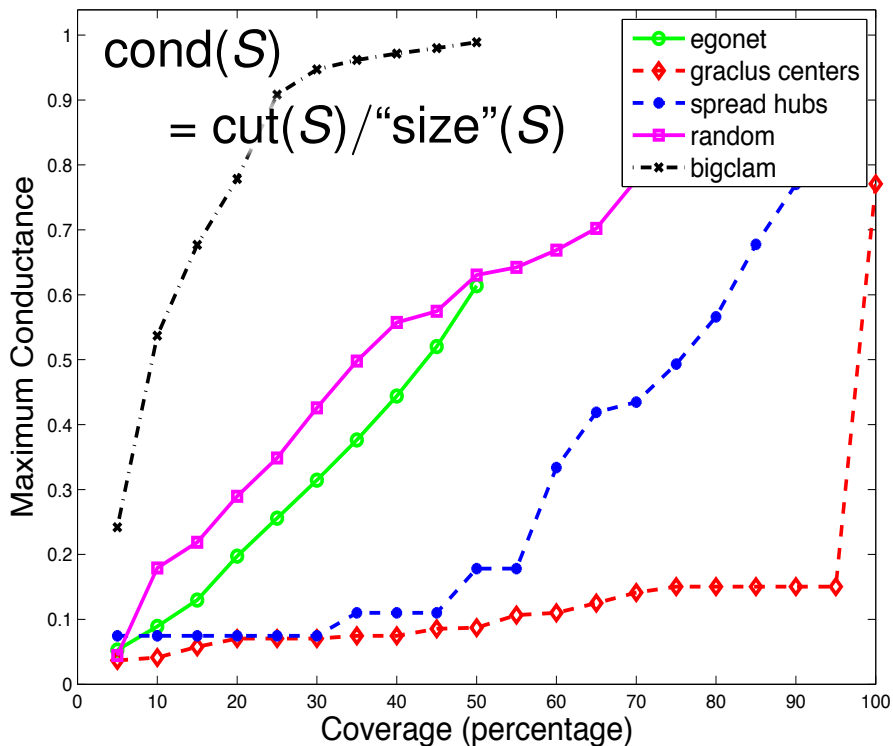
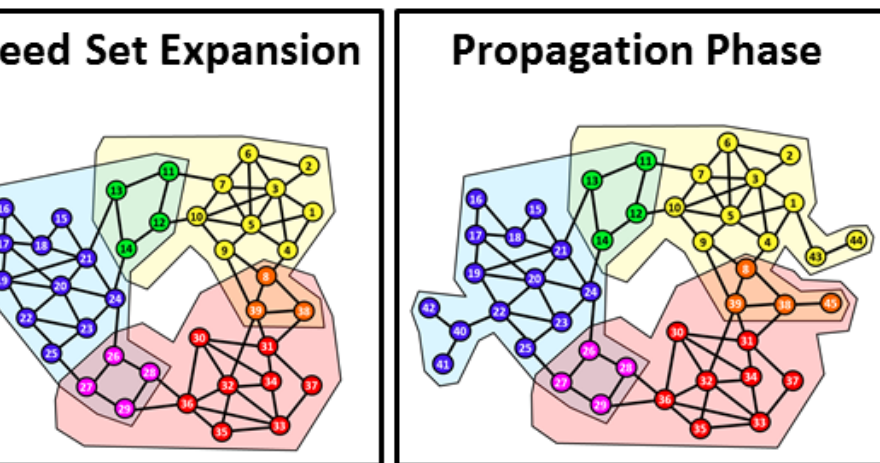
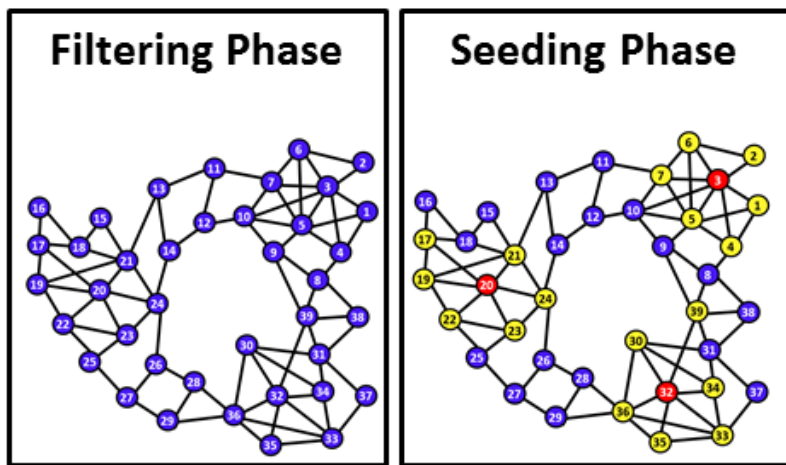


Figure 10: A low-resolution view of part of the movies vs actresses incidence matrix. The Spain co-cluster lies within a supercluster of Spanish- and Portuguese-language countries. Also there are many edges leading from Spain to other Romance-language countries. See section 3.5.

Overlapping communities via seed set expansion works nicely.

Flickr social network
2M vertices, 22M edges



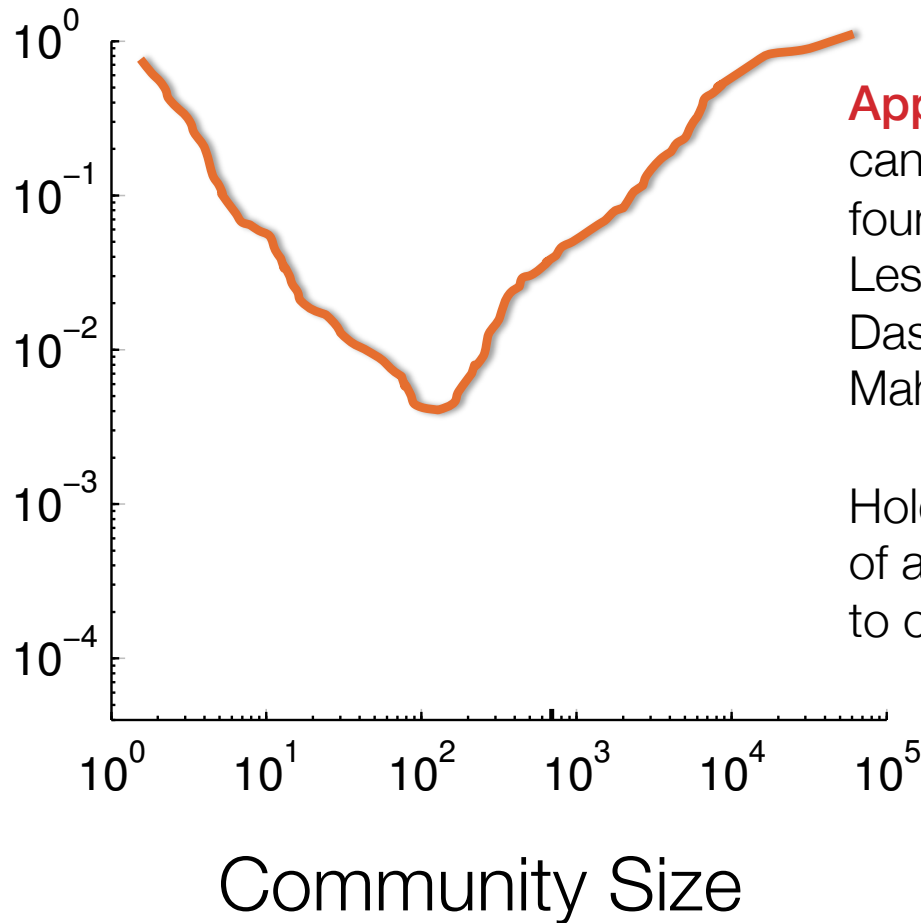
We can cover 95% of network with communities of cond. ~ 0.15 .

... demo ...

github.com/dgleich/diffusion-tutorial

Empirical Evaluation using Network Community Profiles

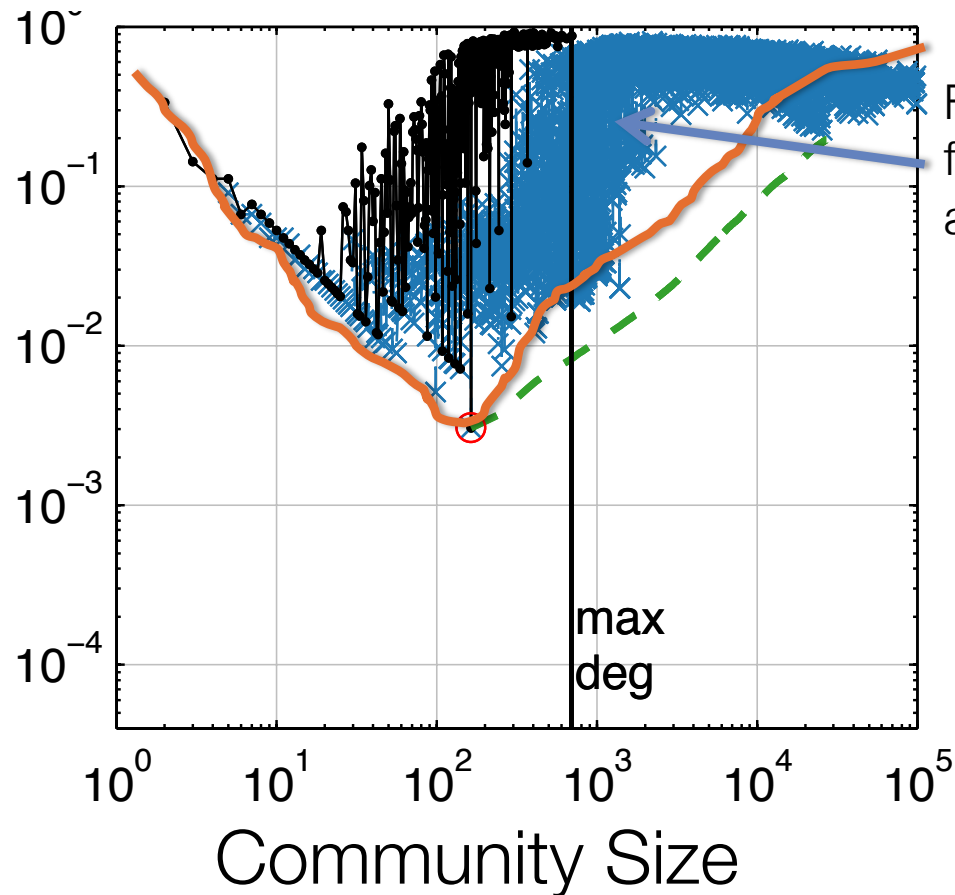
Minimum
conductance for
any community of
the given size



Network Community Profile

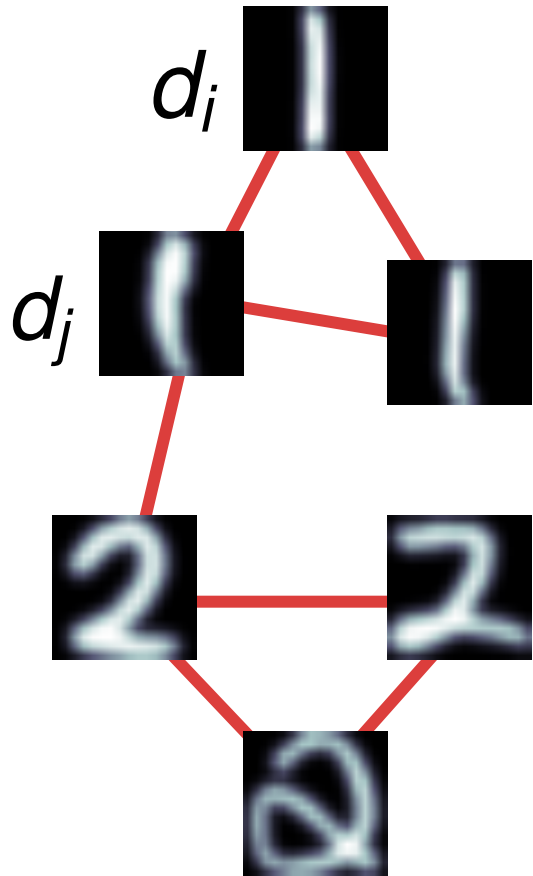
Minimum
conductance for
any community of
the given size

Facebook Sample - 1.1M verts, 4M edges

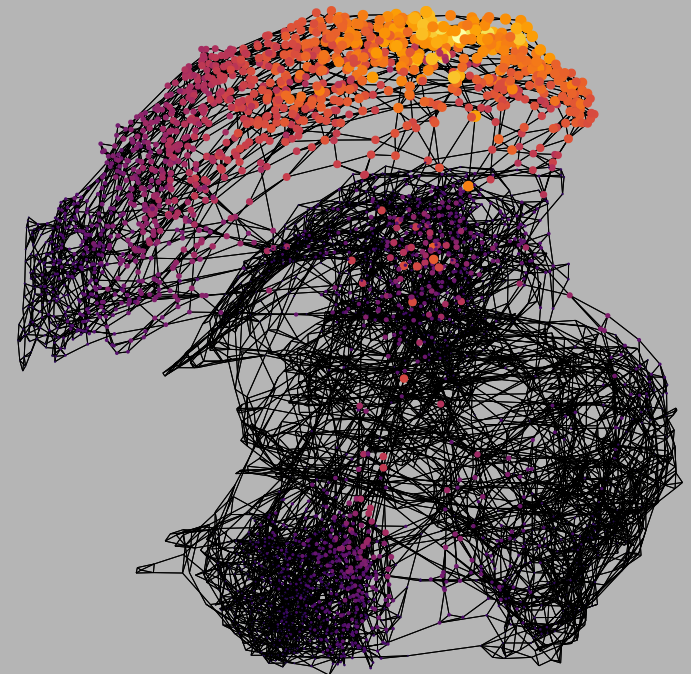
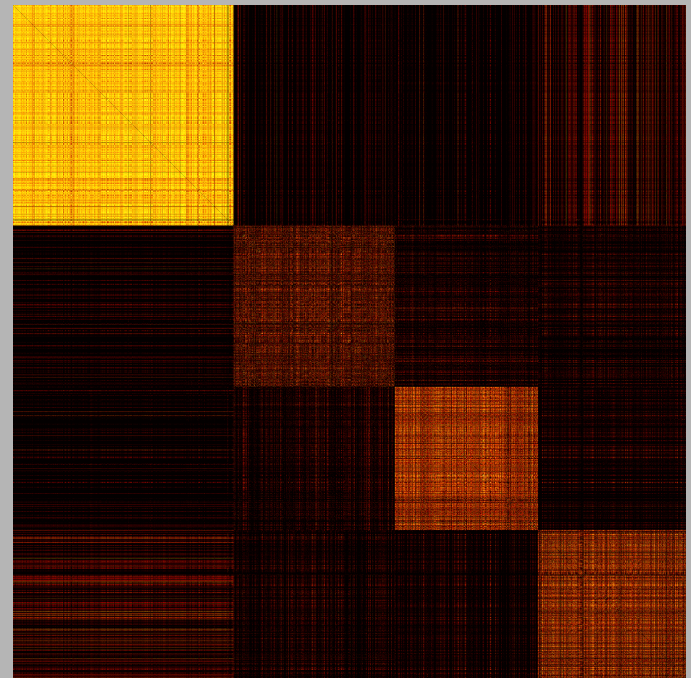


Facebook data
from Wilson et
al. 2009

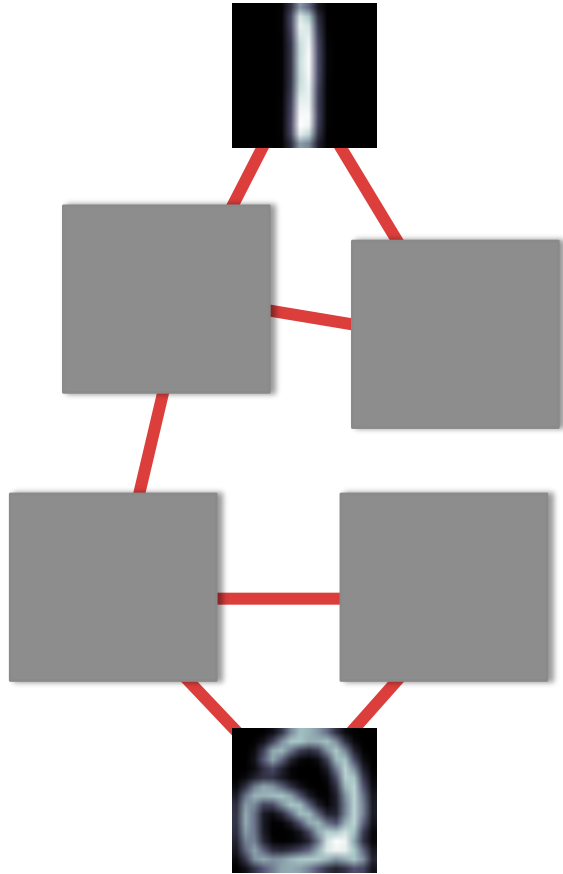
Semi-supervised Learning on Graphs



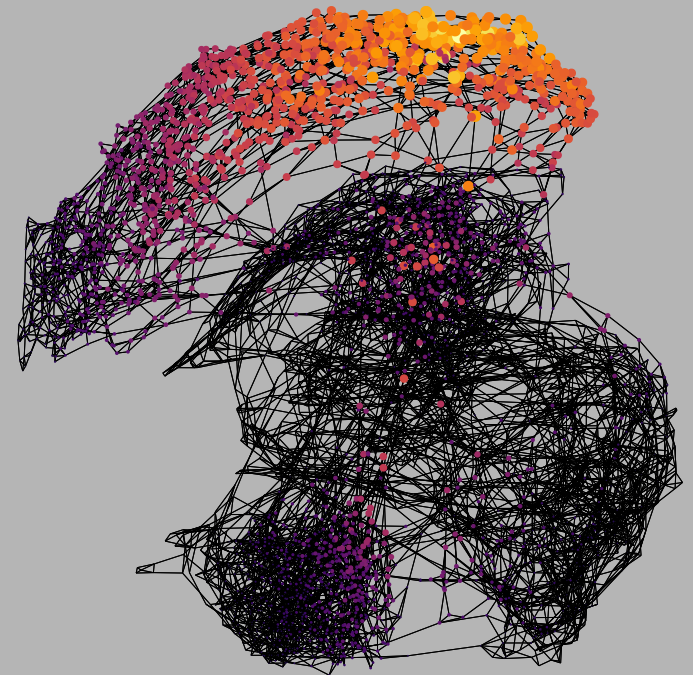
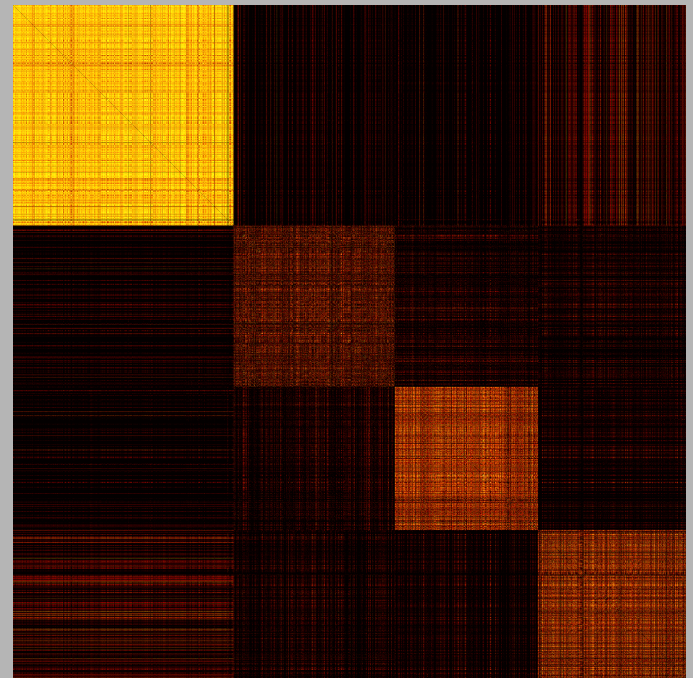
$$A_{i,j} = \exp\left(-\frac{\|d_i - d_j\|_2^2}{2\sigma^2}\right)$$



Semi-supervised Learning on Graphs



Experiment predict unlabeled images from the labeled ones



Semi-supervised Learning on Graphs

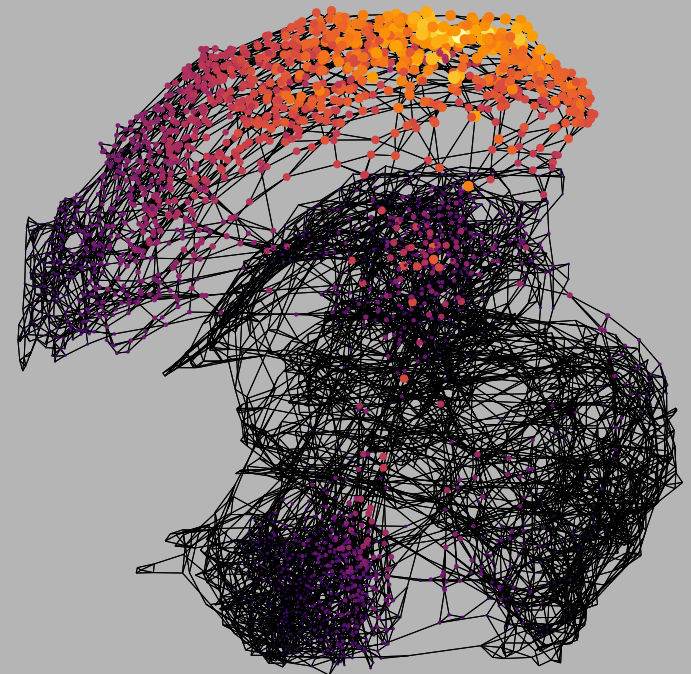
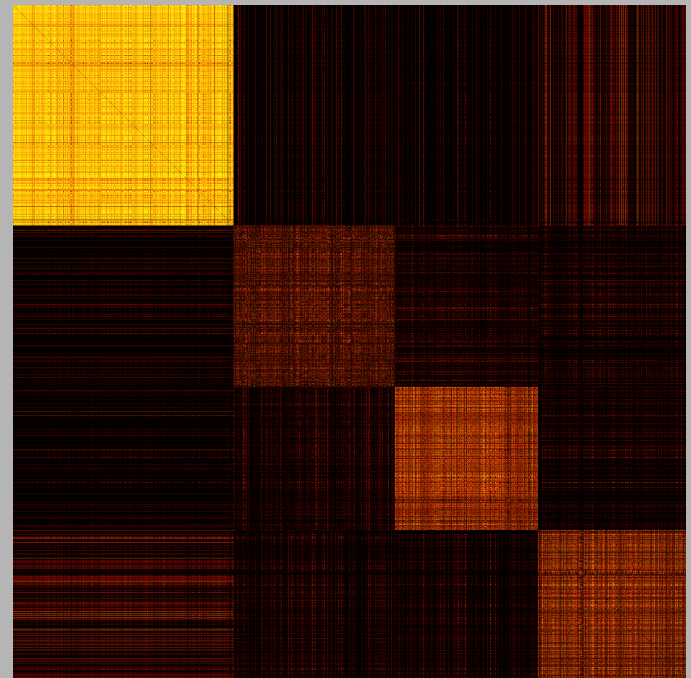
$$\mathbf{K}_1 = (\mathbf{I} - \beta \mathbf{A})^{-1}$$

$$\mathbf{K}_2 = (\mathbf{D} - \beta \mathbf{A})^{-1}$$

Predictions $\mathbf{Y} = \mathbf{K}_j \mathbf{L}$ ← Indicators on the revealed labels

$\mathbf{y} = \operatorname{argmax}_j \mathbf{Y}$

Experiment vary number of labeled images and track perf.



Even more diffusions

SIR and viral thresholds

- Wang et al. 2003
- Berger et al. 2005

Rumor spreading

- Chierichetti et al. 2010

Information cascades

- Farajtaba et al. 2015

Laplacian variations

- Bridle & Zhu, MLG 2013

Diffusions over semi-rings

- Kepner & Gilbert, 2011 (e.g. Peer-pressure clustering)

Generalized coefficients

- Boldi et al. 2005 – TotalRank
- Baeza-Yates et al. 2006 - Generalized
- Constantine & Gleich, 2007, 2010 – Random alpha PageRank

Brief outline

Introduction & Application 20-30 min

Coordinate Relaxation for Strong Convergence 30 min

Break 15 min

Katz Diffusion 10 min

Weak Convergence for PageRank 20 min

Monte Carlo methods 15 min

Break 15 min

Implicit Regularization 25 min

Discussion 15 min

Diffusion vectors

Given a graph G and a set of seed node(s)...

A **diffusion vector** assigns nodes values that quantify somehow the relationship of S to the rest of G , by propagating information from S to the rest of G .

$$\mathbf{f} = \sum_{k=0}^{\infty} c_k \mathbf{M}^k \mathbf{s}$$

Diffusion vector

Diffusion coefficients

Graph related matrix

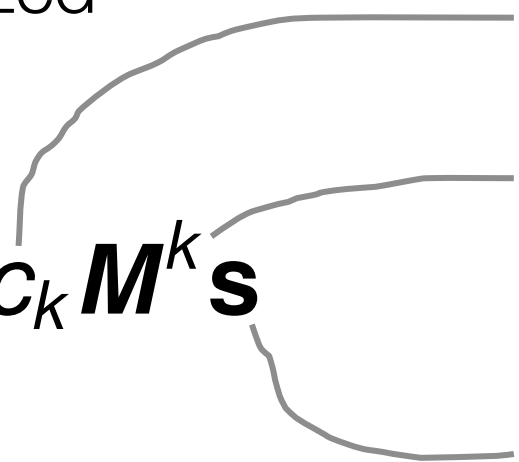
Seed vector

Another application, another diffusion

Different applications and sources of data call for different diffusions – change c_k , \mathbf{M}^k , \mathbf{s} , pre- and post-processing.

Algorithmic approaches

- deterministic
- randomized

$$\mathbf{f} = \sum_{k=0}^{\infty} c_k \mathbf{M}^k \mathbf{s}$$


Diffusion vector

- normalizations
- accuracy

Diffusion coefficients

- (many)

Graph related matrix

- (A,P,L, more)

Seed vector

- constructions
- normalizations

Another application, another diffusion

Different applications and sources of data call for different diffusions – change c_k , \mathbf{M}^k , \mathbf{s} , pre- and post-processing.

Algorithms – coordinate relaxation / push, monte carlo

Coefficients – probability, Katz, adaptive, many more

Graph matrix – adjacency, prob trans, Laplacians, more

Seed vector – seed set indicator vector, normalizations



Diffusion post-processing

- Degree scaling
- weak vs strong accuracy
- Sweep-procedure

Coordinate relaxation methods for graph diffusions

Local node rankings, similarity

Rank nodes with respect to S by taking the largest values from a diffusion \mathbf{f} . More accurate rankings require the solution to be precise enough that large diffusion values that are close to each other aren't mis-ranked:

Ranks		Diffusion values		Approx Ranks
 $\begin{bmatrix} 4 \\ 2 \\ 3 \\ 1 \end{bmatrix}$	$\mathbf{f} =$	$\begin{bmatrix} 0.139 \\ 0.251 \\ 0.25 \\ 0.36 \end{bmatrix}$, $\hat{\mathbf{f}} =$	$\begin{bmatrix} 4 \\ 3 \\ 2 \\ 1 \end{bmatrix}$ 

Local node rankings, similarity

Rank nodes with respect to S by taking the largest values from a diffusion \mathbf{f} . More accurate rankings require the solution to be precise enough ...

GOAL: compute \mathbf{f} with accuracy $\|\mathbf{f} - \hat{\mathbf{f}}\|_1 < \varepsilon$

We begin with the popular personalized PageRank diffusion vector, but will later apply this accuracy setting to a number of diffusions.

PageRank diffusion

The PageRank diffusion can be defined as the solution to

$$(\mathbf{I} - \alpha \mathbf{P})\mathbf{x} = (1 - \alpha)\mathbf{s}$$

for some alpha in (0, 1). The seed vector \mathbf{s} should be normalized to sum to 1. This linear system is equivalent to our definition for a diffusion:

$$\mathbf{f} = \sum_{k=0}^{\infty} \alpha^k \mathbf{P}^k \mathbf{s}$$

PageRank diffusion

The PageRank diffusion can be defined as the solution to

$$(\mathbf{I} - \alpha \mathbf{P})\mathbf{x} = (1 - \alpha)\mathbf{s}$$

for some alpha in (0, 1). The seed vector \mathbf{s} should be normalized to sum to 1. This linear system is equivalent to our definition for a diffusion:

$$\mathbf{f} = \sum_{k=0}^{\infty} \alpha^k \mathbf{P}^k \mathbf{s}$$

This holds when $\|\alpha \mathbf{P}\| < 1$
because of the Neumann series:

$$(\mathbf{I} - \alpha \mathbf{P})^{-1} = \sum_{k=0}^{\infty} \alpha^k \mathbf{P}^k \quad (\text{proof: geometric series})$$

Coordinate relaxation for PageRank

For a fast approximation, \mathbf{f} , to the following

$$(\mathbf{I} - \alpha \mathbf{P})\mathbf{x} = (1 - \alpha)\mathbf{s} = \tilde{\mathbf{s}}$$

we introduce a coordinate relaxation scheme:

Initial solution and residual: $\mathbf{x}^{(0)} = \mathbf{0}, \mathbf{r}^{(0)} = (1 - \alpha)\mathbf{s}$

Iterative updates: first pick entry of residual, j

(There are a number of ways of picking an entry

-- for now, assume just that the entry is non-zero)

Coordinate relaxation for PageRank

For a fast approximation, \mathbf{f} , to the following

$$(\mathbf{I} - \alpha \mathbf{P})\mathbf{x} = \tilde{\mathbf{s}}$$

we introduce a coordinate relaxation scheme:

Initial solution and residual: $\mathbf{x}^{(0)} = \mathbf{0}, \mathbf{r}^{(0)} = (1 - \alpha)\mathbf{s}$

Iterative updates: first pick entry of residual, j

- update solution: $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + r_j \cdot \mathbf{e}_j$

Coordinate relaxation for PageRank

For a fast approximation, \mathbf{f} , to the following

$$(\mathbf{I} - \alpha \mathbf{P})\mathbf{x} = \tilde{\mathbf{s}}$$

we introduce a coordinate relaxation scheme:

Initial solution and residual: $\mathbf{x}^{(0)} = \mathbf{0}, \mathbf{r}^{(0)} = (1 - \alpha)\mathbf{s}$

Iterative updates: first pick entry of residual, j

- update solution: $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + r_j \cdot \mathbf{e}_j$

- update residual: $\mathbf{r}^{(k+1)} = \tilde{\mathbf{s}} - (\mathbf{I} - \alpha \mathbf{P})\mathbf{x}^{(k+1)}$
 $= \mathbf{r}^{(k)} - r_j(\mathbf{I} - \alpha \mathbf{P})\mathbf{e}_j$
 $= \mathbf{r}^{(k)} - r_j\mathbf{e}_j + r_j\alpha \mathbf{P}\mathbf{e}_j$

Coordinate relaxation for PageRank

Approximating

$$(\mathbf{I} - \alpha \mathbf{P})\mathbf{x} = \tilde{\mathbf{s}}$$

Initial solution and residual: $\mathbf{x}^{(0)} = \mathbf{0}, \mathbf{r}^{(0)} = (1 - \alpha)\mathbf{s}$

Iterative updates: first pick entry of residual, j

- update solution: $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + r_j \cdot \mathbf{e}_j$

- update residual: $\mathbf{r}^{(k+1)} = \tilde{\mathbf{s}} - (\mathbf{I} - \alpha \mathbf{P})\mathbf{x}^{(k+1)}$
 $= \mathbf{r}^{(k)} - r_j \mathbf{e}_j + r_j \alpha \mathbf{P} \mathbf{e}_j$

Coordinate relaxation for PageRank

Approximating

$$(\mathbf{I} - \alpha \mathbf{P})\mathbf{x} = \tilde{\mathbf{s}}$$

Initial solution and residual: $\mathbf{x}^{(0)} = \mathbf{0}, \mathbf{r}^{(0)} = (\mathbf{1} - \alpha)\mathbf{s}$

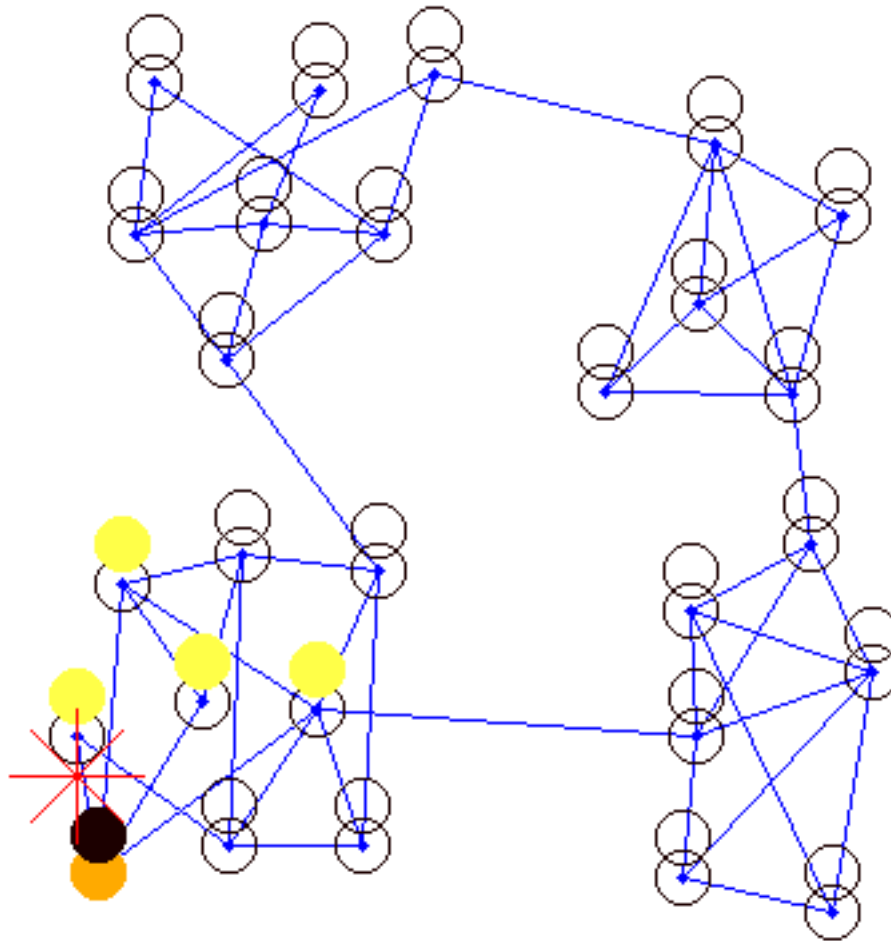
Iterative updates: first pick entry of residual, j

- update solution: $x_j^{(k+1)} = x_j^{(k)} + r_j$

- update residual:

$$r_i^{(k+1)} = \begin{cases} 0, & \text{if } i = j \\ r_i^{(k)} + r_j \alpha / d_j, & \text{if } i \sim j \\ r_i^{(k)}, & \text{else} \end{cases}$$

The push algorithm



Coordinate relaxation remarks, part 1

- This is the fundamental operation underlying many of the deterministic methods for diffusions
- Consists of single update to solution \mathbf{x} and a single column access of matrix \mathbf{P} to update \mathbf{r}
- Because of this, fast for sparse matrices (avoids $O(|E|)$ work required by full mat-vec!)
- Nothing special about PageRank, can be applied in other circumstances

Coordinate relaxation remarks, part 1

- *This is the fundamental operation underlying many of the deterministic methods for diffusions*
- Consists of single update to solution \mathbf{x} and a single column access of matrix \mathbf{P} to update \mathbf{r}
- Because of this, fast for sparse matrices (avoids $O(|E|)$ work required by full mat-vec!)
- Nothing special about PageRank, can be applied in other circumstances

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + r_j \mathbf{e}_j$$

$$\mathbf{r}^{(k+1)} = \mathbf{r}^{(k)} - r_j \mathbf{e}_j + r_j \alpha \mathbf{P} \mathbf{e}_j$$

Coordinate relaxation remarks, part 2

- Convergence depends on method of choosing entry, as well as the underlying matrix
 - Gauss-Southwell: choose r_j to be largest entry in \mathbf{r}
 - Gauss-Seidel: after j , choose $j+1$, then $j+2$, ...
 - Choose any entry \geq average magnitude of \mathbf{r}
 - Choose any entry above some threshold
 - Even random selection can work

Coordinate relaxation remarks, part 2

- Convergence depends on method of choosing entry, as well as the underlying matrix
 - Gauss-Southwell: choose r_j to be largest entry in \mathbf{r}
 - Gauss-Seidel: after j , choose $j+1$, then $j+2$, ...
 - Choose any entry \geq average magnitude of \mathbf{r}
 - Choose any entry above some threshold
 - Even random selection can work
- Many methods of entry selection converge for diagonally dominant and positive definite matrices (in particular, Gauss-Southwell does).
- Implementation requires intelligent choice of data structure for \mathbf{r} for fast entry selection/updates

PageRank Convergence: error & residual

Approximating a solution to

$$(\mathbf{I} - \alpha \mathbf{P})\mathbf{x} = \tilde{\mathbf{s}}$$

residual and error satisfy

$$\begin{aligned}\mathbf{r}^{(k)} &= \tilde{\mathbf{s}} - (\mathbf{I} - \alpha \mathbf{P})\mathbf{x}^{(k)} \\ &= (\mathbf{I} - \alpha \mathbf{P})\mathbf{x} - (\mathbf{I} - \alpha \mathbf{P})\mathbf{x}^{(k)}\end{aligned}$$

PageRank Convergence: error & residual

Approximating a solution to

$$(\mathbf{I} - \alpha \mathbf{P})\mathbf{x} = \tilde{\mathbf{s}}$$

residual and error satisfy

$$\mathbf{r}^{(k)} = \tilde{\mathbf{s}} - (\mathbf{I} - \alpha \mathbf{P})\mathbf{x}^{(k)}$$

$$= (\mathbf{I} - \alpha \mathbf{P})\mathbf{x} - (\mathbf{I} - \alpha \mathbf{P})\mathbf{x}^{(k)}$$

$$(\mathbf{I} - \alpha \mathbf{P})^{-1} \mathbf{r}^{(k)} = (\mathbf{x} - \mathbf{x}^{(k)})$$

$$\|\mathbf{x} - \mathbf{x}^{(k)}\| \leq \|(\mathbf{I} - \alpha \mathbf{P})^{-1}\| \|\mathbf{r}^{(k)}\|$$

for any sub-multiplicative matrix norm $\|\cdot\|$.

PageRank Convergence: residual bound

Approximating a solution to $(\mathbf{I} - \alpha \mathbf{P})\mathbf{x} = \tilde{\mathbf{s}}$

Error satisfies $\|\mathbf{x} - \mathbf{x}^{(k)}\| \leq \|(\mathbf{I} - \alpha \mathbf{P})^{-1}\| \|\mathbf{r}^{(k)}\|$

Initial solution and residual: $\mathbf{x}^{(0)} = \mathbf{0}, \mathbf{r}^{(0)} = (1 - \alpha)\mathbf{s}$

Update residual: $\mathbf{r}^{(k+1)} = \mathbf{r}^{(k)} - r_j \mathbf{e}_j + r_j \alpha \mathbf{P} \mathbf{e}_j$

PageRank Convergence: residual bound

Approximating a solution to $(\mathbf{I} - \alpha \mathbf{P})\mathbf{x} = \tilde{\mathbf{s}}$

Error satisfies $\|\mathbf{x} - \mathbf{x}^{(k)}\| \leq \|(\mathbf{I} - \alpha \mathbf{P})^{-1}\| \|\mathbf{r}^{(k)}\|$

Initial solution and residual: $\mathbf{x}^{(0)} = \mathbf{0}, \mathbf{r}^{(0)} = (1 - \alpha)\mathbf{s}$

Update residual: $\mathbf{r}^{(k+1)} = \mathbf{r}^{(k)} - r_j \mathbf{e}_j + r_j \alpha \mathbf{P} \mathbf{e}_j$

$$\begin{aligned} \|\mathbf{r}^{(k+1)}\|_1 &\leq \|\mathbf{r}^{(k)} - r_j \mathbf{e}_j\|_1 + \|r_j \alpha \mathbf{P} \mathbf{e}_j\|_1 && \text{Triangle inequality} \\ &\leq \|\mathbf{r}^{(k)}\|_1 - r_j + |r_j \alpha| \|\mathbf{P} \mathbf{e}_j\|_1 && \text{Residual nonnegative} \\ &\leq \|\mathbf{r}^{(k)}\|_1 - r_j + |r_j \alpha| && \mathbf{P} \text{ is column-stochastic} \\ &\leq \|\mathbf{r}^{(k)}\|_1 - r_j(1 - \alpha) && \text{Residual nonnegative} \end{aligned}$$

PageRank Convergence: residual bound

Error satisfies $\|\mathbf{x} - \mathbf{x}^{(k)}\| \leq \|(\mathbf{I} - \alpha\mathbf{P})^{-1}\| \|\mathbf{r}^{(k)}\|$

Initial solution and residual: $\mathbf{x}^{(0)} = \mathbf{0}, \mathbf{r}^{(0)} = (1 - \alpha)\mathbf{s}$

Residual norm: $\|\mathbf{r}^{(k+1)}\|_1 \leq \|\mathbf{r}^{(k)}\|_1 - r_j(1 - \alpha)$

Assume we chose r_j to be at least as big as the average magnitude of the residual entries. Then

$$\begin{aligned} r_j &\geq \|\mathbf{r}^{(k)}\|_1 / \text{nnz}(\mathbf{r}^{(k)}) && \text{(definition of average)} \\ &\geq \|\mathbf{r}^{(k)}\|_1 / n && \text{(loose bound!)} \end{aligned}$$

PageRank Convergence: residual bound

Error satisfies $\|\mathbf{x} - \mathbf{x}^{(k)}\| \leq \|(\mathbf{I} - \alpha\mathbf{P})^{-1}\| \|\mathbf{r}^{(k)}\|$

Initial solution and residual: $\mathbf{x}^{(0)} = \mathbf{0}, \mathbf{r}^{(0)} = (1 - \alpha)\mathbf{s}$

Residual norm: $\|\mathbf{r}^{(k+1)}\|_1 \leq \|\mathbf{r}^{(k)}\|_1 - r_j(1 - \alpha)$

Assume we chose r_j to be at least as big as the average magnitude of the residual entries. Then

$$\begin{aligned} r_j &\geq \|\mathbf{r}^{(k)}\|_1 / \text{nnz}(\mathbf{r}^{(k)}) && \text{(definition of average)} \\ &\geq \|\mathbf{r}^{(k)}\|_1 / n && \text{(loose bound!)} \end{aligned}$$

$$\begin{aligned} \|\mathbf{r}^{(k+1)}\|_1 &\leq \|\mathbf{r}^{(k)}\|_1 - \|\mathbf{r}^{(k)}\|_1(1 - \alpha)/n \\ &\leq \|\mathbf{r}^{(k)}\|_1 \left(1 - \frac{(1 - \alpha)}{n}\right) \\ &\leq \|\mathbf{r}^{(0)}\|_1 \left(1 - \frac{(1 - \alpha)}{n}\right)^{k+1} = (1 - \alpha) \left(1 - \frac{(1 - \alpha)}{n}\right)^{k+1} \end{aligned}$$

PageRank Convergence: back to error

Error satisfies $\|\mathbf{x} - \mathbf{x}^{(k)}\| \leq \|(\mathbf{I} - \alpha\mathbf{P})^{-1}\| \|\mathbf{r}^{(k)}\|$

Substituting in for residual...

Residual norm: $\|\mathbf{r}^{(k)}\|_1 \leq (1 - \alpha) \left(1 - \frac{(1-\alpha)}{n}\right)^k$

$\|(\mathbf{I} - \alpha\mathbf{P})^{-1}\|_1 = \frac{1}{1-\alpha}$ (using Neumann series for inverse)

PageRank Convergence: back to error

Error satisfies $\|\mathbf{x} - \mathbf{x}^{(k)}\| \leq \|(\mathbf{I} - \alpha\mathbf{P})^{-1}\| \|\mathbf{r}^{(k)}\|$

Substituting in for residual...

Residual norm: $\|\mathbf{r}^{(k)}\|_1 \leq (1 - \alpha) \left(1 - \frac{(1-\alpha)}{n}\right)^k$

$\|(\mathbf{I} - \alpha\mathbf{P})^{-1}\|_1 = \frac{1}{1-\alpha}$ (using Neumann series for inverse)

Substitution gives $\|\mathbf{x} - \mathbf{x}^{(k)}\|_1 \leq \frac{1}{1-\alpha} \|\mathbf{r}^{(k)}\|$
 $\leq \left(1 - \frac{(1-\alpha)}{n}\right)^k$

Bounds number of iterations with $O(\log(1/\epsilon)n)$
(but the bound can be refined to give sublinear work,
depending on the underlying graph.)

Related work: strong coordinate relaxation

- Deterministic coordinate relaxation for diffusion
 - [Jeh & Widom '03] Scaling Personalized PageRank
 - [McSherry '05] Accelerated PageRank Computation
 - [Berkhin '07] Bookmark Coloring Algorithm for PPR
 - [Bonchi et al. '12] Fast Katz and Commute Times
 - [Kloster, Gleich WAW13] Coordinate relaxation for $\exp(\mathbf{P})e_j$
- Selecting entry to relax:
 - [Dhillon, Ravikumar, Tewari '11] Near neighbor-based greedy coordinate descent
 - [Nutini, et al. 2015] Gauss-Southwell better than random

Brief outline

Introduction & Application 20-30 min

Coordinate Relaxation for Strong Convergence 30 min

Break 15 min

Katz Diffusion 10 min

Weak Convergence for PageRank 20 min

Monte Carlo methods 15 min

Break 15 min

Implicit Regularization 25 min

Discussion 15 min

From PageRank
to the Katz diffusion

Coordinate relaxation for PageRank

Approximating a solution to

$$(\mathbf{I} - \alpha \mathbf{P})\mathbf{x} = \tilde{\mathbf{s}}$$

Initial solution and residual: $\mathbf{x}^{(0)} = \mathbf{0}$, $\mathbf{r}^{(0)} = (1 - \alpha)\mathbf{s}$

Iterative updates: first pick entry of residual, j

- update solution: $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + r_j \cdot \mathbf{e}_j$

- update residual: $\mathbf{r}^{(k+1)} = \tilde{\mathbf{s}} - (\mathbf{I} - \alpha \mathbf{P})\mathbf{x}^{(k+1)}$
 $= \mathbf{r}^{(k)} - r_j \mathbf{e}_j + r_j \alpha \mathbf{P} \mathbf{e}_j$

Coordinate relaxation for Katz diffusion

Approximating a solution to

$$(\mathbf{I} - \alpha \mathbf{A})\mathbf{x} = \mathbf{b}$$

Initial solution and residual: $\mathbf{x}^{(0)} = \mathbf{0}, \mathbf{r}^{(0)} = \mathbf{b}$

Iterative updates: first pick entry of residual, j

- update solution: $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + r_j \cdot \mathbf{e}_j$

- update residual: $\mathbf{r}^{(k+1)} = \mathbf{b} - (\mathbf{I} - \alpha \mathbf{A})\mathbf{x}^{(k+1)}$
 $= \mathbf{r}^{(k)} - r_j \mathbf{e}_j + r_j \alpha \mathbf{A} \mathbf{e}_j$

(Not much difference)

Katz diffusion

The Katz diffusion can be defined as the solution to

$$(\mathbf{I} - \alpha \mathbf{A})\mathbf{x} = \mathbf{s} - (\mathbf{I} - \alpha \mathbf{A})\mathbf{s}$$

The restrictions on alpha differ from PageRank:

PageRank: $0 < \alpha < 1$

Katz: $0 < \alpha < 1/d_{\max} \leq 1/\lambda_1(\mathbf{A})$

where λ_1 is the dominant eigenvalue, and d_{\max} is the largest degree in the graph.

Katz diffusion

The Katz diffusion can be defined as the solution to

$$(\mathbf{I} - \alpha \mathbf{A})\mathbf{x} = \mathbf{s} - (\mathbf{I} - \alpha \mathbf{A})\mathbf{s}$$

The restrictions on alpha differ from PageRank:

PageRank: $0 < \alpha < 1$

Katz: $0 < \alpha < 1/d_{\max} \leq 1/\lambda_1(\mathbf{A})$

where λ_1 is the dominant eigenvalue, and d_{\max} is the largest degree in the graph.

Note: $0 < \alpha < 1/d_{\max}$ guarantees $\|\alpha \mathbf{A}\|_1 < 1$ so

$$\mathbf{f} = \sum_{k=1}^{\infty} \alpha^k \mathbf{A}^k \mathbf{s} = \sum_{k=1}^{\infty} (\alpha d_{\max})^k \left(\frac{1}{d_{\max}} \mathbf{A}\right)^k \mathbf{s}$$

Coordinate relaxation for Katz diffusion

Approximating a solution to

$$(\mathbf{I} - \alpha \mathbf{A})\mathbf{x} = \mathbf{b}$$

Initial solution and residual: $\mathbf{x}^{(0)} = \mathbf{0}, \mathbf{r}^{(0)} = \mathbf{b}$

Iterative updates: first pick entry of residual, j

- update solution: $x_j^{(k+1)} = x_j^{(k)} + r_j$

- update residual:

$$r_i^{(k+1)} = \begin{cases} 0, & \text{if } i = j \\ r_i^{(k)} + r_j \alpha, & \text{if } i \sim j \\ r_i^{(k)}, & \text{else} \end{cases}$$

Convergence for PageRank (repeated)

Error satisfies $\|\mathbf{x} - \mathbf{x}^{(k)}\| \leq \|(\mathbf{I} - \alpha\mathbf{P})^{-1}\| \|\mathbf{r}^{(k)}\|$

Substituting in for residual...

Residual norm: $\|\mathbf{r}^{(k)}\|_1 \leq (1 - \alpha) \left(1 - \frac{(1-\alpha)}{n}\right)^k$

$\|(\mathbf{I} - \alpha\mathbf{P})^{-1}\|_1 = \frac{1}{1-\alpha}$ (using Neumann series for inverse)

Substitution gives $\|\mathbf{x} - \mathbf{x}^{(k)}\|_1 \leq \frac{1}{1-\alpha} \|\mathbf{r}^{(k)}\|$
 $\leq \left(1 - \frac{(1-\alpha)}{n}\right)^k$

Bounds number of iterations with $O(\log(1/\epsilon)n)$
(but the bound can be refined to give sublinear work,
depending on the underlying graph.)

Convergence for Katz

Error satisfies $\|\mathbf{x} - \mathbf{x}^{(k)}\| \leq \|(\mathbf{I} - \alpha\mathbf{A})^{-1}\| \|\mathbf{r}^{(k)}\|$

Differences:

Depending on which norm is desired, scale \mathbf{A} by $1/d_{\max}$ or $1/\lambda_1$. Then, rest of the convergence analysis for PageRank applies to the scaled adjacency matrix:

$$\|(\mathbf{I} - \alpha\mathbf{A})^{-1}\|_1 \leq \frac{1}{1 - \alpha d_{\max}} \quad \text{if} \quad \alpha < 1/d_{\max}$$

(If $\alpha < 1/\lambda_1$ the method still converges, but analysis is trickier.)

Brief outline

Introduction & Application 20-30 min

Coordinate Relaxation for Strong Convergence 30 min

Break 15 min

Katz Diffusion 10 min

Weak Convergence for PageRank 20 min

Monte Carlo methods 15 min

Break 15 min

Implicit Regularization 25 min

Discussion 15 min

Weak convergence coordinate
relaxation for communities and good
conductance sets.

Weak accuracy for community detection

Local community detection/ finding good conductance sets from a diffusion vector \mathbf{f} :

Seek the largest values in the diffusion vector; weak accuracy, because identifying the largest values is the goal, not the precise values themselves (or even the precise ranking of values).

Weak accuracy for community detection

Local community detection/ finding good conductance sets from a diffusion vector \mathbf{f} :

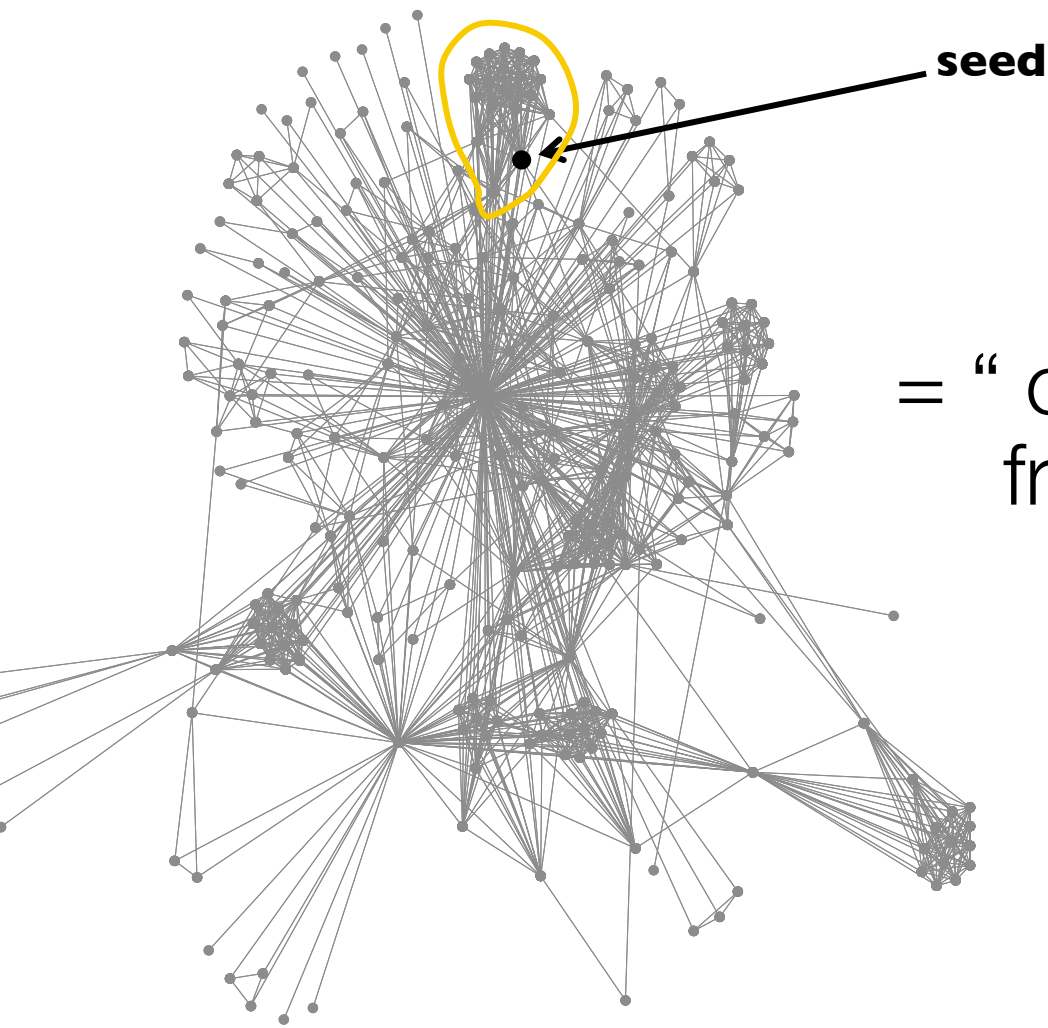
Seek the largest values in the diffusion vector; weak accuracy, because identifying the largest values is the goal, not the precise values themselves (or even the precise ranking of values).

GOAL: compute $\hat{\mathbf{f}}$ with accuracy $0 \leq f_j - \hat{f}_j \leq \varepsilon d_j$ (entry-wise). Equivalent to

$$\mathbf{f} \geq \hat{\mathbf{f}} \quad \text{and} \quad \|\mathbf{D}^{-1}(\mathbf{f} - \hat{\mathbf{f}})\|_{\infty} < \varepsilon$$

Low-conductance sets

$$\text{conductance}(T) = \frac{\# \text{ edges leaving } T}{\min(\text{vol}(T), \text{vol}(G-T))}$$



$$\text{vol}(S) = \sum_{v \in S} d(v)$$

= “ chance a random step from inside T exits T ”

Use a diffusion for good conductance sets

1. Approximate \mathbf{f} so $\|\mathbf{D}^{-1}(\mathbf{f} - \hat{\mathbf{f}})\|_{\infty} \leq \epsilon$
2. Scale by \mathbf{D} ,
3. Then “sweep” for best conductance set.

Sweep:

1. Sort diffusion vector so $f_1/d(1) \geq f_2/d(2) \geq \dots$
2. Consider the sweep sets $S(j) = \{1, 2, \dots, j\}$
3. Return the set $S(j)$ with the best conductance.

Diffusions used for conductance

Personalized PageRank (PPR)

$$\mathbf{f} = \sum_{k=0}^{\infty} \alpha^k \mathbf{P}^k \tilde{\mathbf{s}}$$

Heat Kernel (HK)

$$\mathbf{f} = \sum_{k=0}^{\infty} \frac{t^k}{k!} \mathbf{P}^k \tilde{\mathbf{s}}$$

Diffusions used for conductance

Personalized PageRank (PPR)

$$\mathbf{f} = \sum_{k=0}^{\infty} \alpha^k \mathbf{P}^k \tilde{\mathbf{s}}$$

Heat Kernel (HK)

$$\mathbf{f} = \sum_{k=0}^{\infty} \frac{t^k}{k!} \mathbf{P}^k \tilde{\mathbf{s}}$$

Time-dependent PageRank (TDPR)

$$\mathbf{f} = \sum_{k=0}^{\infty} \left[(1 - \alpha) \alpha^k \left(1 - e^{-\gamma} \sum_{r=0}^k \frac{\gamma^r}{r!} \right) + e^{-\gamma} \frac{\alpha^k \gamma^k}{k!} \right] \mathbf{P}^k \mathbf{s}$$

Comes from

$$\begin{aligned} \mathbf{x}'(t) &= (1 - \alpha) \mathbf{s} - (\mathbf{I} - \alpha \mathbf{P}) \mathbf{x}(t) \\ \mathbf{x}(0) &= \mathbf{s} \end{aligned}$$

Diffusions used for conductance

Personalized PageRank (PPR)

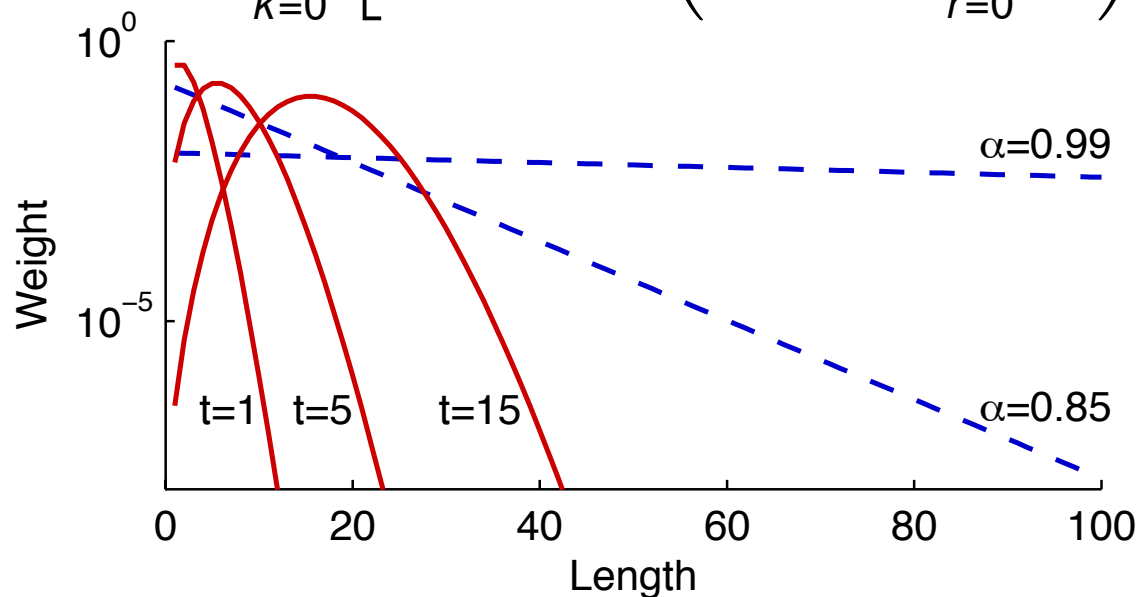
$$\mathbf{f} = \sum_{k=0}^{\infty} \alpha^k \mathbf{P}^k \tilde{\mathbf{s}}$$

Heat Kernel (HK)

$$\mathbf{f} = \sum_{k=0}^{\infty} \frac{t^k}{k!} \mathbf{P}^k \tilde{\mathbf{s}}$$

Time-dependent PageRank (TDPR)

$$\mathbf{f} = \sum_{k=0}^{\infty} \left[(1 - \alpha) \alpha^k \left(1 - e^{-\gamma} \sum_{r=0}^k \frac{\gamma^r}{r!} \right) + e^{-\gamma} \frac{\alpha^k \gamma^k}{k!} \right] \mathbf{P}^k \mathbf{s}$$



Various diffusions explore different aspects of graphs.

Diffusions used for conductance

Personalized PageRank (PPR)

$$\mathbf{f} = \sum_{k=0}^{\infty} \alpha^k \mathbf{P}^k \tilde{\mathbf{s}}$$

Heat Kernel (HK)

$$\mathbf{f} = \sum_{k=0}^{\infty} \frac{t^k}{k!} \mathbf{P}^k \tilde{\mathbf{s}}$$

Time-dependent PageRank (TDPR)

$$\mathbf{f} = \sum_{k=0}^{\infty} \left[(1 - \alpha) \alpha^k \left(1 - e^{-\gamma} \sum_{r=0}^k \frac{\gamma^r}{r!} \right) + e^{-\gamma} \frac{\alpha^k \gamma^k}{k!} \right] \mathbf{P}^k \mathbf{s}$$

Use other matrices, too:

$$\mathbf{f} = \sum_{k=0}^{\infty} c_k \mathbf{L}^k \mathbf{s}$$

(Various weightings and scalings

of the Laplacian have been explored, [Ghosh et al. '14])

Diffusions: conductance & algorithms

good
conductance

fast
algorithm

PR

Local Cheeger Inequality
[Andersen, Chung, Lang 06]

[Andersen Chung Lang 06]
“PPR-push” is $O(1/(\epsilon(1-\alpha)))$

HK

Local Cheeger Inequality
[Chung '07]

[Kloster, Gleich '14]
“HK-push” is $O(e^t C / \epsilon)$

TDPR

Open question

[Avron, Horesh '15]
Constant-time heuristically

Gen
Diff

[Ghosh et al. '14] on L;
open question for general f

In revision!

Weak convergence for PageRank

Approximating a solution to $(\mathbf{I} - \alpha \mathbf{P})\mathbf{x} = \tilde{\mathbf{s}}$

residual and error satisfy $\mathbf{r}^{(k)} = \tilde{\mathbf{s}} - (\mathbf{I} - \alpha \mathbf{P})\mathbf{x}^{(k)}$

$$(\mathbf{I} - \alpha \mathbf{P})^{-1} \mathbf{r}^{(k)} = (\mathbf{x} - \mathbf{x}^{(k)})$$

$$\|\mathbf{x} - \mathbf{x}^{(k)}\| \leq \|(\mathbf{I} - \alpha \mathbf{P})^{-1}\| \|\mathbf{r}^{(k)}\|$$

for any sub-multiplicative matrix norm $\|\cdot\|$. Scale by \mathbf{D} !

Weak convergence for PageRank

Approximating a solution to $(\mathbf{I} - \alpha\mathbf{P})\mathbf{x} = \tilde{\mathbf{s}}$

residual and error satisfy $\mathbf{r}^{(k)} = \tilde{\mathbf{s}} - (\mathbf{I} - \alpha\mathbf{P})\mathbf{x}^{(k)}$

$$(\mathbf{I} - \alpha\mathbf{P})^{-1}\mathbf{r}^{(k)} = (\mathbf{x} - \mathbf{x}^{(k)})$$

$$\|\mathbf{x} - \mathbf{x}^{(k)}\| \leq \|(\mathbf{I} - \alpha\mathbf{P})^{-1}\| \|\mathbf{r}^{(k)}\|$$

for any sub-multiplicative matrix norm $\|\cdot\|$. Scale by \mathbf{D} !

$$\mathbf{D}^{-1}\mathbf{r}^{(k)} = \mathbf{D}^{-1}\tilde{\mathbf{s}} - (\mathbf{I} - \alpha\mathbf{P}^T)\mathbf{D}^{-1}\mathbf{x}^{(k)}$$

$$(\mathbf{I} - \alpha\mathbf{P}^T)^{-1}\mathbf{D}^{-1}\mathbf{r}^{(k)} = \mathbf{D}^{-1}(\mathbf{x} - \mathbf{x}^{(k)})$$

$$\|\mathbf{D}^{-1}(\mathbf{x} - \mathbf{x}^{(k)})\| \leq \|(\mathbf{I} - \alpha\mathbf{P}^T)^{-1}\| \|\mathbf{D}^{-1}\mathbf{r}^{(k)}\|$$

This requires \mathbf{A} is symmetric:

$$\mathbf{D}^{-1}\mathbf{P} = \mathbf{D}^{-1}(\mathbf{A}\mathbf{D}^{-1}) = (\mathbf{D}^{-1}\mathbf{A})\mathbf{D}^{-1} = \mathbf{P}^T\mathbf{D}^{-1}$$

Weak coordinate relaxation

Approximating a solution to $(\mathbf{I} - \alpha \mathbf{P})\mathbf{x} = \tilde{\mathbf{s}}$

Residual and error satisfy

$$\|\mathbf{D}^{-1}(\mathbf{x} - \mathbf{x}^{(k)})\|_{\infty} \leq \|(\mathbf{I} - \alpha \mathbf{P}^T)^{-1}\|_{\infty} \|\mathbf{D}^{-1} \mathbf{r}^{(k)}\|_{\infty}$$

$$\|\mathbf{D}^{-1}(\mathbf{x} - \mathbf{x}^{(k)})\|_{\infty} \leq \frac{1}{1-\alpha} \max\{r_j/d_j\}$$

Weak coordinate relaxation

Approximating a solution to $(\mathbf{I} - \alpha \mathbf{P})\mathbf{x} = \tilde{\mathbf{s}}$

Residual and error satisfy

$$\|\mathbf{D}^{-1}(\mathbf{x} - \mathbf{x}^{(k)})\|_{\infty} \leq \|(\mathbf{I} - \alpha \mathbf{P}^T)^{-1}\|_{\infty} \|\mathbf{D}^{-1} \mathbf{r}^{(k)}\|_{\infty}$$

$$\|\mathbf{D}^{-1}(\mathbf{x} - \mathbf{x}^{(k)})\|_{\infty} \leq \frac{1}{1-\alpha} \max\{r_j/d_j\}$$

Contrast with 1-norm version: here simply track residual entries (degree normalized) that exceed a threshold.

This suggests a new method of choosing the coordinate j .

The rest of the update operation stays the same.

Weak coordinate relaxation, operation

Approximating a solution to $(\mathbf{I} - \alpha\mathbf{P})\mathbf{x} = \tilde{\mathbf{s}}$

Initial solution and residual: $\mathbf{x}^{(0)} = \mathbf{0}, \mathbf{r}^{(0)} = (1 - \alpha)\mathbf{s}$

Iterative update: a queue stores big entries: $r_j \geq \varepsilon d_j$
- pick top entry off $Q(\mathbf{r}), j$.

Weak coordinate relaxation

Approximating a solution to $(\mathbf{I} - \alpha\mathbf{P})\mathbf{x} = \tilde{\mathbf{s}}$

Initial solution and residual: $\mathbf{x}^{(0)} = \mathbf{0}, \mathbf{r}^{(0)} = (1 - \alpha)\mathbf{s}$

Iterative update: a queue stores big entries: $r_j \geq \varepsilon d_j$

- pick top entry off $Q(\mathbf{r})$, j .

- update solution: $x_j^{(k+1)} = x_j^{(k)} + r_j$

- update residual:
$$r_i^{(k+1)} = \begin{cases} 0, & \text{if } i = j \\ r_i^{(k)} + r_j \alpha / d_j, & \text{if } i \sim j \\ r_i^{(k)}, & \text{else} \end{cases}$$

Weak coordinate relaxation

Approximating a solution to $(\mathbf{I} - \alpha\mathbf{P})\mathbf{x} = \tilde{\mathbf{s}}$

Initial solution and residual: $\mathbf{x}^{(0)} = \mathbf{0}, \mathbf{r}^{(0)} = (1 - \alpha)\mathbf{s}$

Iterative update: a queue stores big entries: $r_j \geq \varepsilon d_j$

- pick top entry off $Q(\mathbf{r})$, j .

- update solution: $x_j^{(k+1)} = x_j^{(k)} + r_j$

- update residual:
$$r_i^{(k+1)} = \begin{cases} 0, & \text{if } i = j \\ r_i^{(k)} + r_j \alpha / d_j, & \text{if } i \sim j \\ r_i^{(k)}, & \text{else} \end{cases}$$

- for $i \sim j$: $r_i \geq \varepsilon d_i$, add r_i to $Q(\mathbf{r})$ if not present

Weak coordinate relaxation, work bound

1. Every update is on an entry in $Q(\mathbf{r})$ satisfying $\varepsilon d_i \leq r_i$

2. Sum of updates is $\sum_{t=1} r_{i(t)} = \sum_{k=1}^n \hat{f}_k \leq 1$

3. Total work is $\sum_{t=1} d_{i(t)}$

Weak coordinate relaxation, work bound

1. Every update is on an entry in $Q(\mathbf{r})$ satisfying $\varepsilon d_i \leq r_i$

2. Sum of updates is $\sum_{t=1} r_{i(t)} = \sum_{k=1}^n \hat{f}_k \leq 1$

3. Total work is $\sum_{t=1} d_{i(t)}$

All together: $\varepsilon \sum_{t=1} d_{i(t)} = \sum_{t=1} \varepsilon d_{i(t)} \leq \sum_{t=1} r_{i(t)} \leq 1$

Work is bounded by $1 / ((1 - \alpha)\varepsilon)$
a constant independent of the graph size!

Weak coordinate relaxation, work bound

1. Every update is on an entry in $Q(\mathbf{r})$ satisfying $\varepsilon d_i \leq r_i$

2. Sum of updates is $\sum_{t=1} r_{i(t)} = \sum_{k=1}^n \hat{f}_k \leq 1$

3. Total work is $\sum_{t=1} d_{i(t)}$

All together: $\varepsilon \sum_{t=1} d_{i(t)} = \sum_{t=1} \varepsilon d_{i(t)} \leq \sum_{t=1} r_{i(t)} \leq 1$

Work is bounded by $1 / ((1 - \alpha)\varepsilon)$ (comes from $\mathbf{r}^{(0)} = (1 - \alpha)\mathbf{s}$)
a constant independent of the graph size!

Weak coordinate relaxation remarks

- [Ghosh et al '14] proved Cheeger inequalities for related diffusions that use weighted Laplacians
- Is there a related Cheeger inequality, and a constant-time algorithm, for the degree-normalized Katz diffusion?
- Is there a “best” set of diffusion coefficients for identifying particular structures?
- Can we improve on the sweep procedure? Or bound its performance?
 - [Kenter et al. '15] introduced a randomized subroutine that improves on sweep in certain conditions

Related work: weak coordinate relaxation

- Deterministic coordinate descent
 - [Andersen, Chung, Lang '06] Local Graph Partitioning
 - [Andersen, Lang '06] Communities from Seed Sets
 - [Kloster, Gleich '14] Heat Kernel clustering
 - [Kloumann, Kleinberg '14] Community membership from seed
 - [Ghosh et al. '14] Interplay between dynamics and networks
 - [Avron, Horesh '15] Time-Dependent PageRank clustering

Brief outline

Introduction & Application 20-30 min

Coordinate Relaxation for Strong Convergence 30 min

Break 15 min

Katz Diffusion 10 min

Weak Convergence for PageRank 20 min

Monte Carlo methods 15 min

Break 15 min

Implicit Regularization 25 min

Discussion 15 min

Monte Carlo methods for diffusion vectors

Monte Carlo motivation

Benefits of MC over deterministic?

- strong convergence method drastically slows down as it operates on nodes of large degree
- weak convergence method gets no accuracy when it encounters nodes of large degree
- MC avoids out-link accesses best

Monte Carlo method for Matrix inversion

[Forsyth & Liebler, 1950]

Matrix Inversion by a Monte Carlo method:

Want $(\mathbf{B}^{-1})_{ij}$, so design a game such that the expected value is exactly $(\mathbf{B}^{-1})_{ij}$. It has inspired other work:

Monte Carlo method for Matrix inversion

[Forsyth & Liebler, 1950]

Matrix Inversion by a Monte Carlo method:

Want $(\mathbf{B}^{-1})_{ij}$, so design a game such that the expected value is exactly $(\mathbf{B}^{-1})_{ij}$. It has inspired other work:

- [K. Avrachenkov '05] MC methods in PageRank
- [Fogaras et al. '05] Fully scaling personalized PageRank
- [Das Sarma et al. '08] Estimating PageRank on graph streams
- [Bahmani '10] Fast incremental Personalized PageRank
- [Bahmani '10] PageRank & MapReduce
- [Borgs '12] Sublinear PageRank
- [Chung, Simpson WAW13] Solving systems w/ heat kernel

Designing a game for PageRank

We want a specific entry of
say, \mathbf{f}_i .

$$\mathbf{f} = (1 - \alpha) \sum_{k=0}^{\infty} \alpha^k \mathbf{P}^k \mathbf{e}_j$$

GOAL: design a random process for producing $\hat{\mathbf{f}}$
so that the expected value of each entry is the true value.

Designing a game for PageRank

We want a specific entry of say, \mathbf{f}_i .

$$\mathbf{f} = (1 - \alpha) \sum_{k=0}^{\infty} \alpha^k \mathbf{P}^k \mathbf{e}_j$$

GOAL: design a random process for producing $\hat{\mathbf{f}}$ so that the expected value of each entry is the true value.

Observe that

$$\begin{aligned} \mathbf{f}_i &= (1 - \alpha) \sum_{k=0}^{\infty} \alpha^k (\mathbf{P}^k)_{ij} \\ &= (1 - \alpha) \sum_{k=0}^{\infty} \alpha^k \left(\sum_{p_k(i,j)} \sum_{i_1 \in N(j)} \sum_{i_2 \in N(i_1)} \cdots \sum_{i_k \in N(i_{k-1})} P_{i_k, i_{k-1}} \cdots P_{i_2, i_1} P_{i_1, j} \right) \end{aligned}$$

where $p_k(i,j)$ is the set of all k -walks from j to i (and $i_k = i$).

Designing a game for PageRank

We'll convert this so it looks like an expected value:

$$\mathbf{f}_j = (1 - \alpha) \sum_{k=0}^{\infty} \alpha^k \left(\sum_{p_k(i,j)} \sum_{i_1 \in N(j)} \sum_{i_2 \in N(i_1)} \cdots \sum_{i_k \in N(i_{k-1})} P_{i_k, i_{k-1}} \cdots P_{i_2, i_1} P_{i_1, j} \right)$$

Note that $P(i_k, i_{k-1}) \cdots P(i_1, j)$, is the probability of taking a specific walk, $w_k(i, j)$.

Designing a game for PageRank

We'll convert this so it looks like an expected value:

$$\mathbf{f}_j = (1 - \alpha) \sum_{k=0}^{\infty} \alpha^k \left(\sum_{p_k(i,j)} \sum_{i_1 \in N(j)} \sum_{i_2 \in N(i_1)} \cdots \sum_{i_k \in N(i_{k-1})} P_{i_k, i_{k-1}} \cdots P_{i_2, i_1} P_{i_1, j} \right)$$

Note that $P(i_k, i_{k-1}) \dots P(i_1, j)$, is the probability of taking a specific walk, $w_k(i, j)$. We can rewrite...

$$\begin{aligned} \mathbf{f}_j &= (1 - \alpha) \sum_{k=0}^{\infty} \alpha^k \left(\sum_{w_k(i,j) \in p_k(i,j)} \mathbb{P}(w_k(i, j)) \right) \\ &= \sum_{k=0}^{\infty} \left(\sum_{w_k(i,j) \in p_k(i,j)} (1 - \alpha) \alpha^k \mathbb{P}(w_k(i, j)) \right) \end{aligned}$$

Designing a game for PageRank

GOAL: convert so it looks like an expected value.

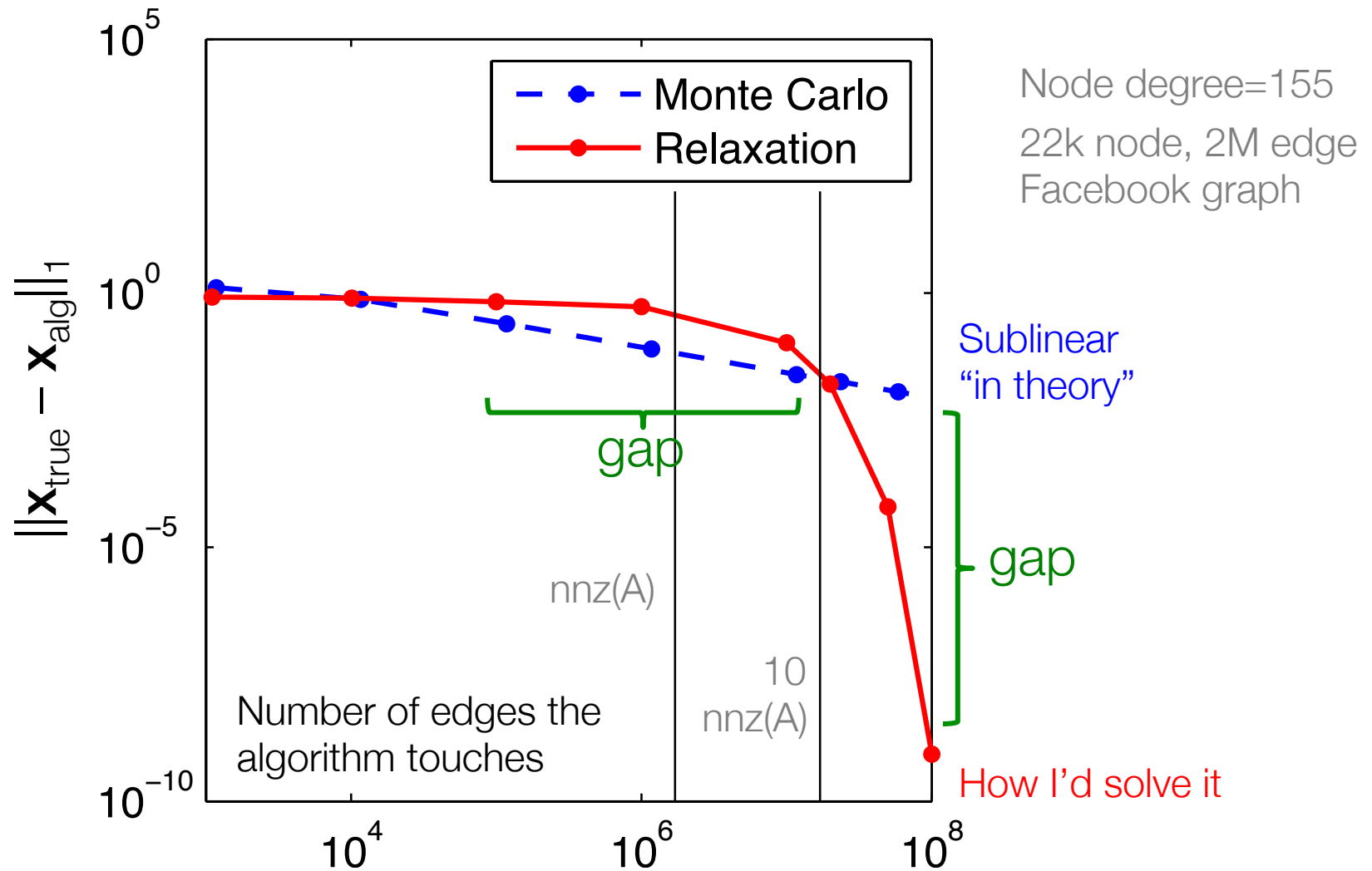
$$\mathbf{f}_i = \sum_{k=0}^{\infty} \left(\sum_{w_k(i,j) \in P_k(i,j)} (1 - \alpha) \alpha^k \mathbb{P}(w_k(i,j)) \cdot 1 \right)$$

This is the expected value of:

1. Choose a length k with probability $(1 - \alpha) \alpha^k$
2. Make a random k -walk from j . It lands at node i_k
3. Update solution where walk ends, i_k , by adding 1.

Hence for a single iteration of this, $\mathbb{E}(\hat{\mathbf{f}}_i) = \mathbf{f}_i$

Monte Carlo vs Deterministic



Monte Carlo Remarks

Accuracy, convergence are problematic (previous slide)

OPEN QUESTION: can we improve on number of samples / random walks required, or the accuracy attained?

Related work: Monte Carlo

- Monte Carlo methods:
 - [Forsyth Liebler '50] fore-runner
 - [K. Avrachenkov '05] MC methods in PageRank
 - [Fogaras et al. '05] Fully scaling personalized PageRank
 - [Das Sarma et al. '08] Estimating PageRank on graph streams
 - [Bahmani '10] Fast incremental Personalized PageRank
 - [Bahmani '10] PageRank & MapReduce
 - [Borgs '12] Sublinear PageRank
 - [Chung, Simpson WAW13] Solving systems w/ heat kernel
- Hybrid monte carlo / coordinate relaxation
 - [Lofgren et al. 2014], node-to-node PPR estimate

Brief outline

Introduction & Application 20-30 min

Coordinate Relaxation for Strong Convergence 30 min

Break 15 min

Katz Diffusion 10 min

Weak Convergence for PageRank 20 min

Monte Carlo methods 15 min

Break 15 min

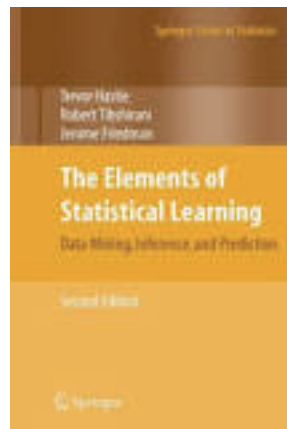
Implicit Regularization 25 min

Discussion 15 min

Statistical regularization

Least squares	minimize \mathbf{x}	$\ \mathbf{Ax} - \mathbf{b}\ _2^2$
Lasso	minimize \mathbf{x}	$\ \mathbf{Ax} - \mathbf{b}\ _2^2 + \lambda\ \mathbf{x}\ _1$

Choose the regularizer to counter a given noise type



See Hastie, Tibshirani, Friedman 2009
The Elements of Statistical Learning: Data Mining, Inference, and Prediction

Implicit regularization

1. Run an algorithm procedure
2. Show that your algorithm implicitly is tolerant to a type of noise.

Implementing Regularization Implicitly Via Approximate Eigenvector Computation

Michael W. Mahoney

MMAHONEY@CS.STANFORD.EDU

Department of Mathematics, Stanford University, Stanford, CA 94305

Lorenzo Orecchia

ORECCHIA@EECS.BERKELEY.EDU

Computer Science Division, UC Berkeley, Berkeley, CA 94720

Consider a graph diffusion (of our type) these diffusions implicitly regularize the resulting solution vectors to be tolerant to noise.

An example theorem

SPECTRAL CLUSTERING

$$\text{minimize } \sum_{ij} \mathcal{L}_{ij} X_{ij}$$

$$\text{subject to } \text{trace}(\mathbf{X}) = 1 \\ \mathbf{X} \succeq 0$$

REGULARIZED

SPECTRAL CLUSTERING

$$\text{minimize } \sum_{ij} \mathcal{L}_{ij} X_{ij} + \lambda F(\mathbf{X})$$

$$\text{subject to } \text{trace}(\mathbf{X}) = 1 \\ \mathbf{X} \succeq 0$$

Let $F(\mathbf{X}) = \text{trace}(\mathbf{X} \log \mathbf{X}) - \text{trace}(\mathbf{X})$
then the solution of regularized spectral is
 $\mathbf{X} = C \exp(-(1/\lambda)\mathcal{L})$

You do nothing special!

Using the heat kernel implicitly regularizes solutions against noise characterized by the generalized entropy function

More examples

SPECTRAL CLUSTERING

$$\text{minimize } \sum_{ij} \mathcal{L}_{ij} X_{ij}$$

$$\text{subject to } \text{trace}(\mathbf{X}) = 1 \\ \mathbf{X} \succeq 0$$

REGULARIZED

SPECTRAL CLUSTERING

$$\text{minimize } \sum_{ij} \mathcal{L}_{ij} X_{ij} + \lambda F(\mathbf{X})$$

$$\text{subject to } \text{trace}(\mathbf{X}) = 1 \\ \mathbf{X} \succeq 0$$

PageRank

Let $F(\mathbf{X}) = \log \det \mathbf{X}$

then the solution of regularized spectral is

$$\mathbf{X} = \mathbf{C}(\mathbf{I} - \alpha \mathcal{L})^{-1}$$

Truncation

Let $F_p(\mathbf{X}) = 1/p \text{ trace}(\mathbf{X}^p)$

then the solution of regularized spectral is

$$\mathbf{X} = \mathbf{C}(\mathcal{A})^{q-1} \text{ where } 1/p + 1/q = 1$$

These results should be true up to degree normalization on the solution.

Our question

Why does the “push method” have such incredible empirical utility?

Answer

Gleich & Mahoney, ICML 2014.

Anti-differentiating approximation algorithms, a case study with min-cuts, spectral, and flow.

Algorithmic Anti-differentiation

Understanding how and why heuristic procedures

- Early stopping
- Truncating small entries
- etc

are actually *algorithms* for implicit objectives.

The ideal world



Given Problem P

Derive solution
characterization C

Show algorithm A
finds a solution where C
holds

Profit?

Given “min-cut”

Derive “max-flow is
equivalent to min-cut”

Show push-relabel
solves max-flow

Profit!

(The ideal world)'

Given Problem P

Derive solution approx. characterization C'

Show algorithm A' quickly finds a solution where C' holds

Profit?

Given “sparest-cut”

Derive Rayleigh-quotient approximation

Show power-method finds a good Rayleigh-quotient

Profit?

The real world?

Given Task P

Hack around until you find something useful

Write paper presenting “novel heuristic” H for P and ...

Profit!

Given “find-communities”

Hack around
??? (hidden) ???

Write paper presenting “three matvecs finds real-world communities”

Profit!

Algorithmic Anti-differentiation

Given heuristic H , is there a problem P' such that H is an algorithm for P' ?

Understand why H works **Given** “find-communities”

Show heuristic H solves P' **Hack around**

Guess and check

until you find something H solves

Derive characterization of heuristic H

Write paper presenting “three matvecs finds real-world communities”

Profit!

e.g. Mahoney & Orecchia

Algorithmic Anti-differentiation

Given heuristic H , is there a problem P' such that H is an algorithm for P' ?

If your algorithm is related to optimization, this is:

Given a procedure X , what objective does it optimize?

In an unconstrained case, this is just “anti-differentiation!”

Our question

Why does the “push method” have such incredible empirical utility?

Answer

Gleich & Mahoney, ICML 2014.

Anti-differentiating approximation algorithms, a case study with min-cuts, spectral, and flow.

The O (correct) answer

1. PageRank related to Laplacian
2. Laplacian related to cuts
3. Andersen, Chung, Lang provides the “right” bounds and “localization”

Now the θ (correct) answer?

A deeper insight into the relationship

Intellectually indebted to ...

Chin, Mađry, Miller & Peng [2013]

Orecchia & Zhu [2014]

The s-t min-cut problem

Unweighted incidence matrix

Diagonal capacity matrix

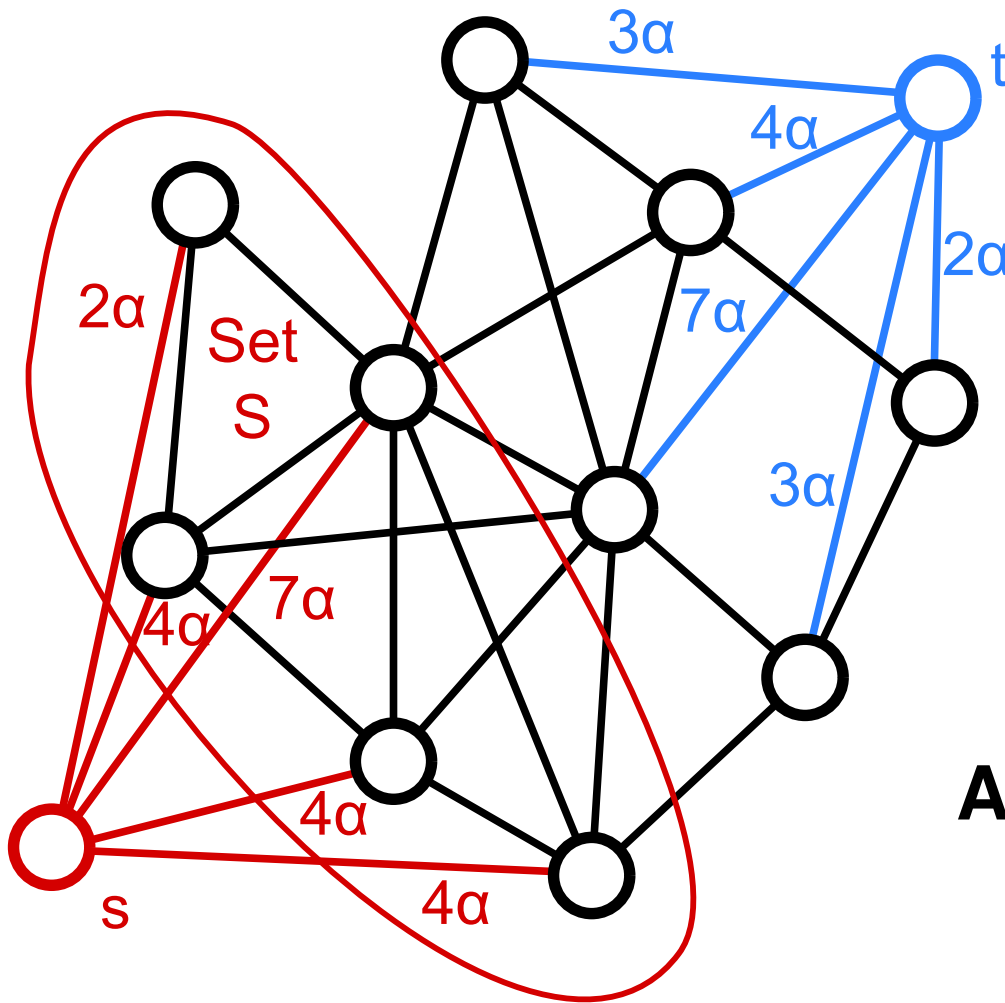
minimize

$$\|\mathbf{B}\mathbf{x}\|_{C,1} = \sum_{ij \in E} C_{i,j} |x_i - x_j|$$

subject to

$$x_s = 1, x_t = 0, \mathbf{x} \geq 0.$$

The localized cut graph

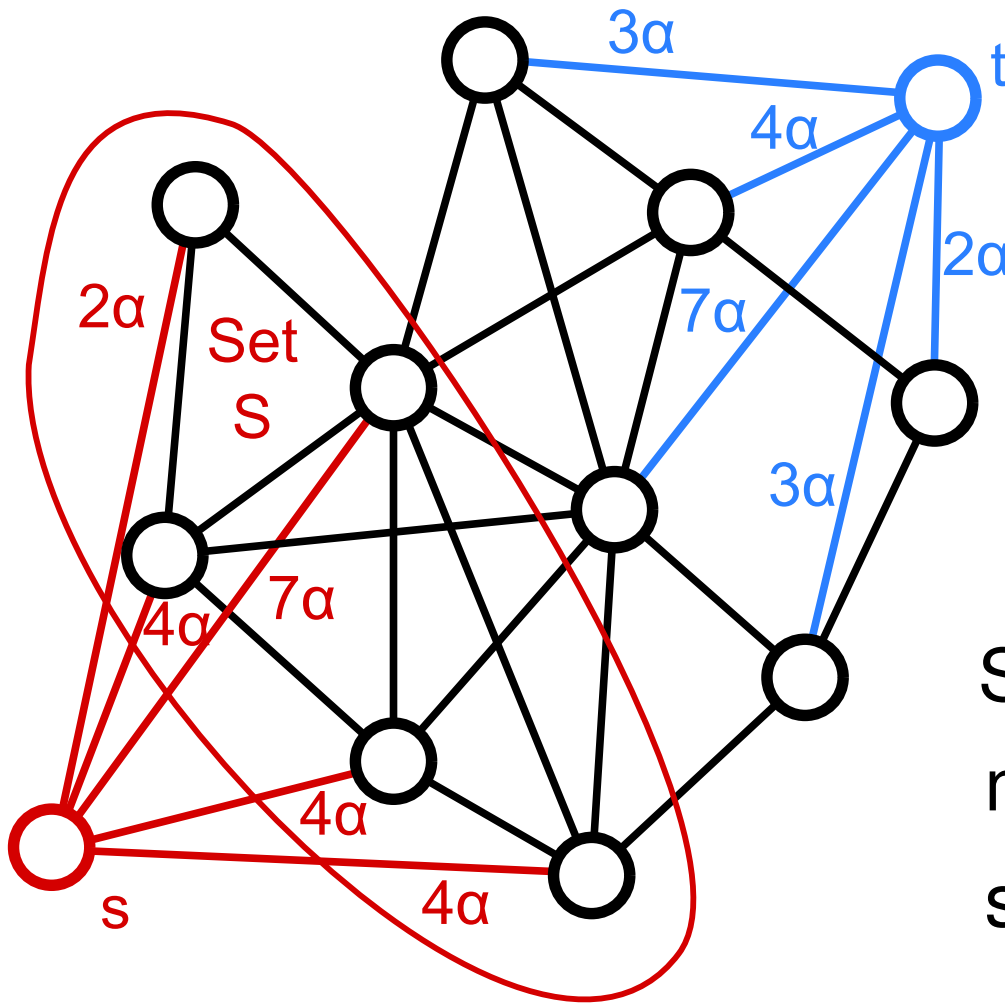


Connect **s** to vertices in **S** with weight $\alpha \cdot \text{degree}$
 Connect **t** to vertices in \bar{S} with weight $\alpha \cdot \text{degree}$

Related to a construction used in “FlowImprove”
 Andersen & Lang (2007); and
 Orecchia & Zhu (2014)

$$\mathbf{A}_S = \begin{bmatrix} 0 & \alpha \mathbf{d}_S^T & 0 \\ \alpha \mathbf{d}_S & \mathbf{A} & \alpha \mathbf{d}_{\bar{S}} \\ 0 & \alpha \mathbf{d}_{\bar{S}}^T & 0 \end{bmatrix}$$

The localized cut graph



Connect **s** to vertices in **S** with weight $\alpha \cdot \text{degree}$
 Connect **t** to vertices in \bar{S} with weight $\alpha \cdot \text{degree}$

$$\mathbf{B}_S = \begin{bmatrix} \mathbf{e} & -\mathbf{I}_S & 0 \\ 0 & \mathbf{B} & 0 \\ 0 & -\mathbf{I}_{\bar{S}} & \mathbf{e} \end{bmatrix}$$

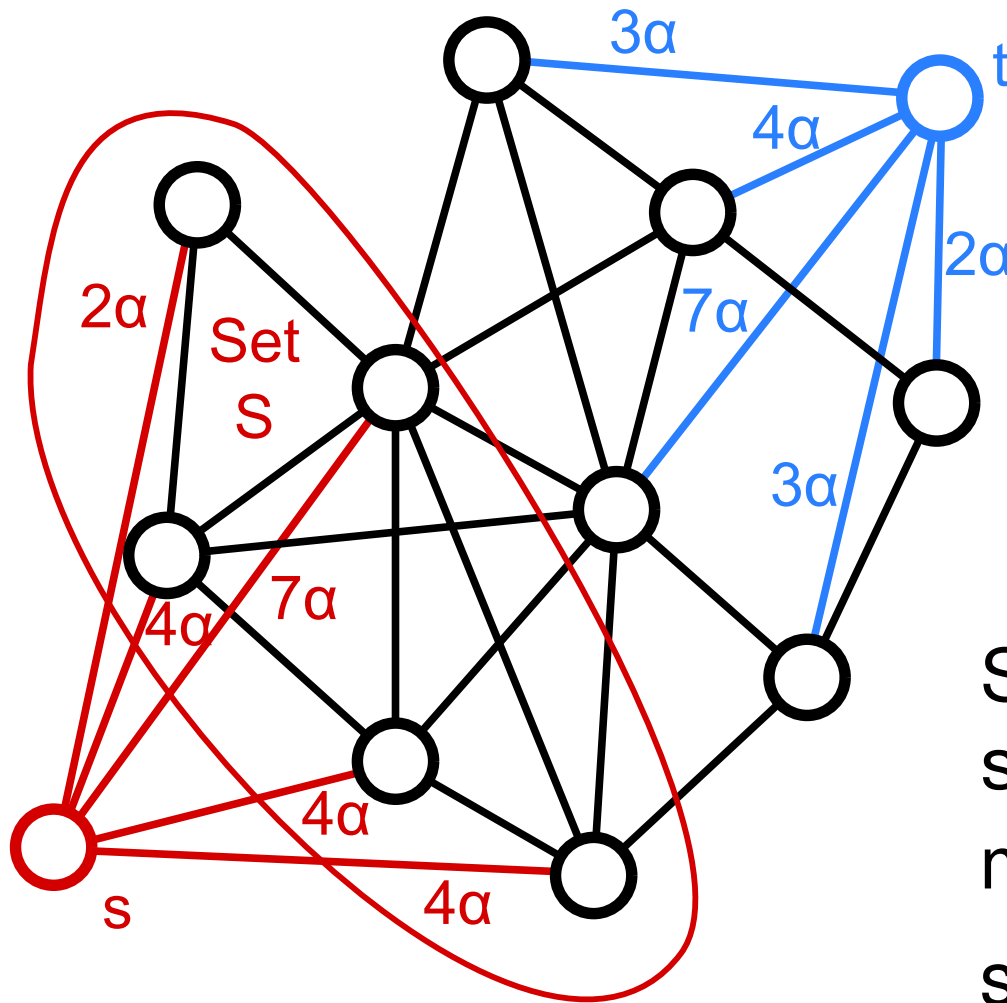
Solve the s-t min-cut

minimize $\|\mathbf{B}_S \mathbf{x}\|_{C(\alpha), 1}$

subject to $x_s = 1, x_t = 0$

$\mathbf{x} \geq 0.$

The localized cut graph



Connect **s** to vertices in **S** with weight $\alpha \cdot \text{degree}$
 Connect **t** to vertices in \bar{S} with weight $\alpha \cdot \text{degree}$

$$\mathbf{B}_S = \begin{bmatrix} \mathbf{e} & -\mathbf{I}_S & 0 \\ 0 & \mathbf{B} & 0 \\ 0 & -\mathbf{I}_{\bar{S}} & \mathbf{e} \end{bmatrix}$$

Solve the “electrical flow”
 s-t min-cut

minimize $\|\mathbf{B}_S \mathbf{x}\|_{C(\alpha), 2}$

subject to $x_s = 1, x_t = 0$

s-t min-cut \rightarrow PageRank

The PageRank vector \mathbf{z} that solves

$$(\alpha \mathbf{D} + \mathbf{L})\mathbf{z} = \alpha \mathbf{v}$$

with $\mathbf{v} = \mathbf{d}_S / \text{vol}(S)$ is a renormalized solution of the electrical cut computation:

$$\text{minimize} \quad \|\mathbf{B}_S \mathbf{x}\|_{C(\alpha), 2}$$

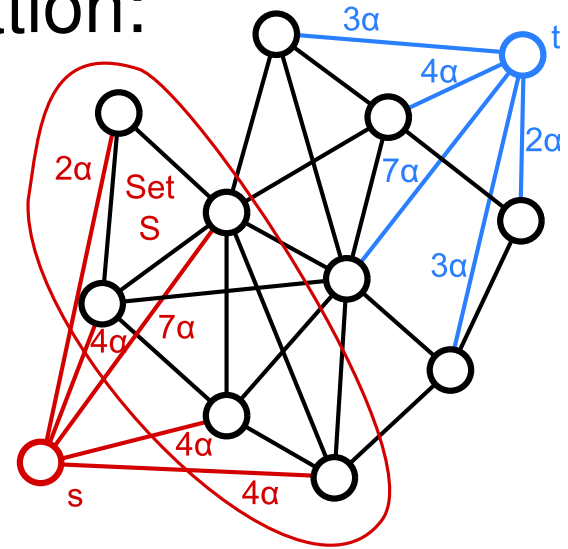
$$\text{subject to} \quad x_s = 1, x_t = 0.$$

Specifically, if \mathbf{x} is the solution, then

$$\mathbf{x} = \begin{bmatrix} 1 \\ \text{vol}(S)\mathbf{z} \\ 0 \end{bmatrix}$$

Proof

Square and expand the objective into a Laplacian, then apply constraints.



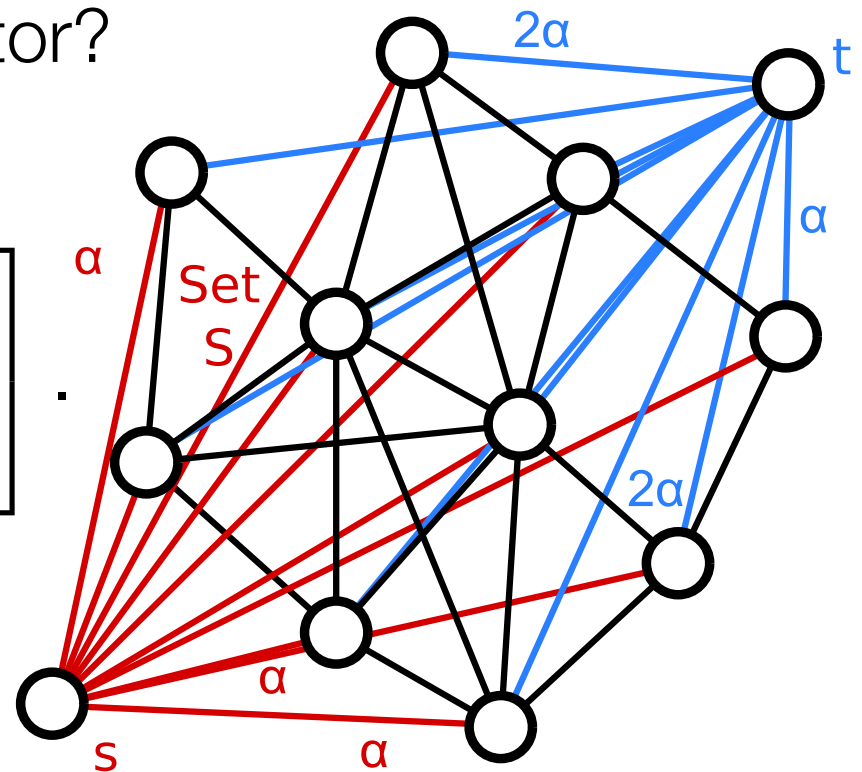
PageRank \rightarrow s-t min-cut

That equivalence works if \mathbf{v} is degree-weighted.

What if \mathbf{v} is the uniform vector?

$\mathbf{A}(\mathbf{s}) =$

$$\begin{bmatrix} 0 & \alpha \mathbf{s}^T & 0 \\ \alpha \mathbf{s} & \mathbf{A} & 0 \\ 0 & \alpha(\mathbf{d} - \mathbf{s})^T & 0 \end{bmatrix} \cdot \alpha(\mathbf{d} - \mathbf{s})$$



And beyond ...

$$\begin{bmatrix} 0 & \mathbf{e}_S^T & 0 \\ \mathbf{e}_S & \theta \mathbf{A} & \mathbf{e}_{\bar{S}} \\ 0 & \mathbf{e}_{\bar{S}} & 0 \end{bmatrix} \cdot (\mathbf{I} + \theta \mathbf{L}) \mathbf{x} = \mathbf{e}_S$$

Easy to cook up interesting diffusion-like problems and adapt them to this framework. In particular, Zhou et al. (2004) gave a semi-supervised learning diffusion we study soon.

The Push Algorithm for PageRank

Proposed (in closest form) in Andersen, Chung, Lang
(also by McSherry, Jeh & Widom) for *personalized PageRank*

Strongly related to Gauss-Seidel (as Kyle mentioned!)

Derived to show improved runtime for balanced solvers

1. $\mathbf{x}^{(1)} = \mathbf{0}, \mathbf{r}^{(1)} = (1 - \beta)\mathbf{e}_i, k = 1$
2. *while any $r_j > \tau d_j$ (d_j is the degree of node j)*
3. $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + (r_j - \tau d_j \rho)\mathbf{e}_j$
4.
$$\mathbf{r}_i^{(k+1)} = \begin{cases} \tau d_j \rho & i = j \\ r_i^{(k)} + \beta(r_j - \tau d_j \rho)/d_j & i \sim j \\ r_i^{(k)} & \text{otherwise} \end{cases}$$
5. $k \leftarrow k + 1$

The
Push
Method
 τ, ρ

Back to the push method

Let \mathbf{x} be the output from the push method
with $0 < \beta < 1$, $\mathbf{v} = \mathbf{d}_S / \text{vol}(S)$,
 $\rho = 1$, and $\tau > 0$.

Set $\alpha = \frac{1-\beta}{\beta}$, $\kappa = \tau \text{vol}(S) / \beta$, and let \mathbf{z}_G solve:

minimize $\frac{1}{2} \|\mathbf{B}_S \mathbf{z}\|_{C(\alpha), 2}^2 + \kappa \|\mathbf{Dz}\|_1$

subject to $z_S = 1, z_t = 0, \mathbf{z} \geq 0$

Need for
normalization

Regularization
for sparsity

where $\mathbf{z} = \begin{bmatrix} 1 \\ \mathbf{z}_G \\ 0 \end{bmatrix}$.

Then $\mathbf{x} = \mathbf{Dz}_G / \text{vol}(S)$.

Proof Write out KKT conditions
Show that the push method
solves them. Slackness was “tricky”

Some reflections on algorithmic anti-differentiating

Differentiation

Given $f(x)$, computing $g(x) = f'(x)$ analytically usually isn't too hard.

Anti-differentiation

Given $f(x)$, computing $F(x)$ where $f(x) = F'(x)$ can be very hard and or impossible

Algorithms solve for the KKT conditions (e.g. $f(x) = 0$)

Algorithmic anti-differentiation finds the objective that the algorithm solves (e.g. minimize $F(x)$)

For simple optimization, this analogy is precise.

*In general, this is very hard or impossible.
(3-4 years for PageRank)*

Your question?

So what? Why does this matter?

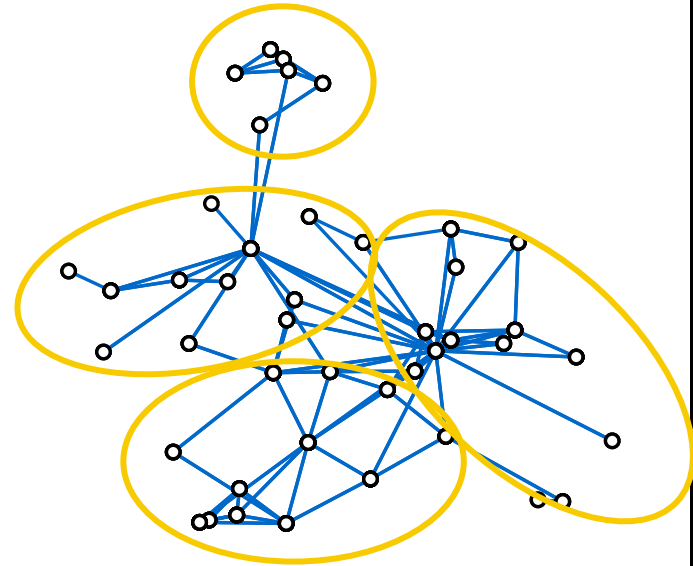
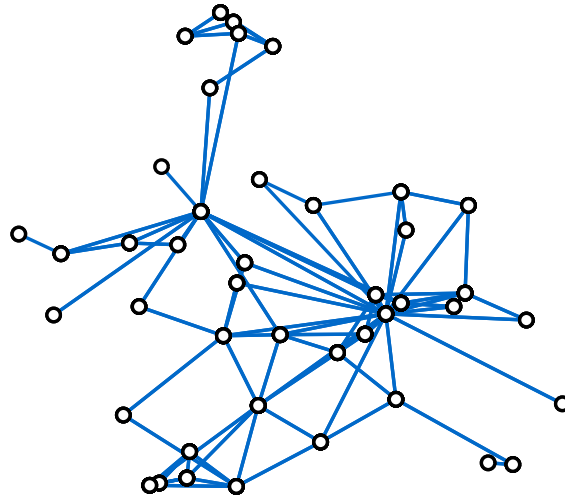
Answer

Gleich & Mahoney, KDD 2015

Using local spectral methods to robustify graph-based learning.

The graph-based data analysis pipeline

```
1 0 0 0 1 0 0 1
0 1 0 1 0 0 1 1
0 1 0 1 0 0 0 1
1 0 0 0 0 0 1 1
1 1 0 1 1 1 0 1
1 0 1 1 0 0 0 1
1 0 1 1 1 0 1 0
1 1 1 1 0 1 0 0
1 1 1 0 0 1 1 1
1 1 0 1 1 1 1 1
```



Raw data

- Relationships
- Images
- Text records
- Etc.

Convert to a graph

- Nearest neighs
- Kernels
- 2-mode to 1-mode
- Etc.

Algorithm/Learning

- Important nodes
- Infer features
- Clustering
- Etc.

“Noise” in the initial data modeling decisions

Explicit graphs

are those that are given to a data analyst.

“A social network”

- Known spam accounts included?
- Users not logged in for a year?
- Etc.

A type of noise

Constructed graphs

are built based on some other primary data.

“nearest neighbor graphs”

- K-NN or ϵ -NN
- Thresholding correlations to zero

Often made for computational convenience! (Graph too big.)

A different type a noise!

Labeled graphs

occur in information diffusion/propagation

“function prediction”

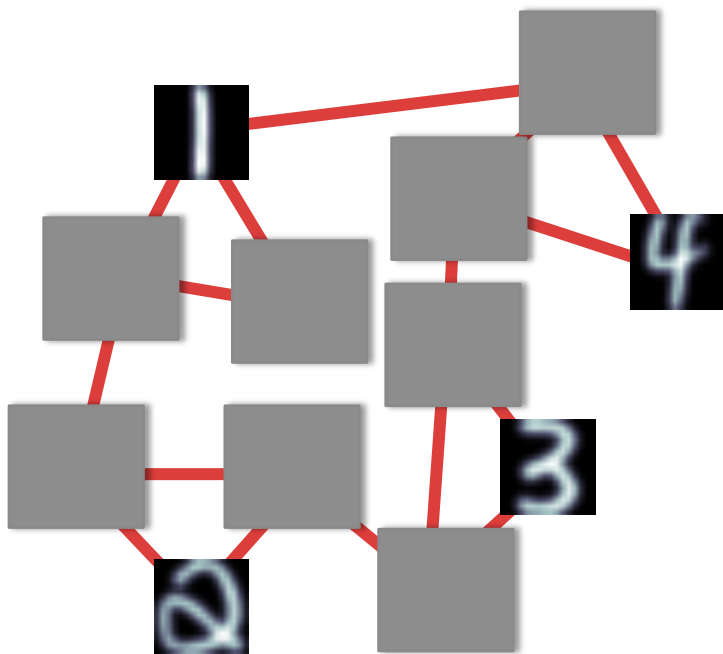
- Labeled nodes
- Labeled edges
- Some are wrong

A direct type of noise!

Do these decisions matter?
Our experience Yes! Dramatically so!

Semi-supervised graph-based learning

Given a graph, and a few labeled nodes, predict the labels on the rest of the graph.

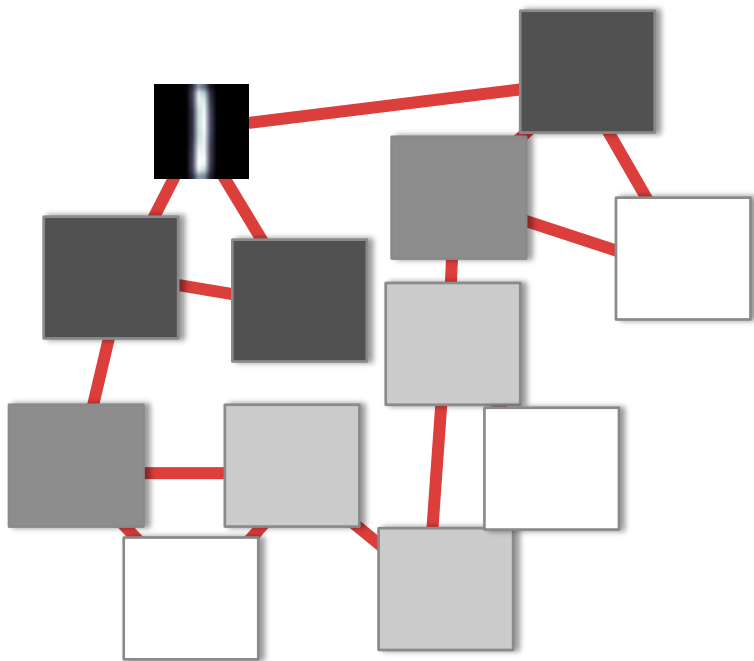


Algorithm

1. Run a diffusion for each label (possibly with neg. info from other classes)
2. Assign new labels based on the value of each diffusion

Semi-supervised graph-based learning

Given a graph, and a few labeled nodes, predict the labels on the rest of the graph.

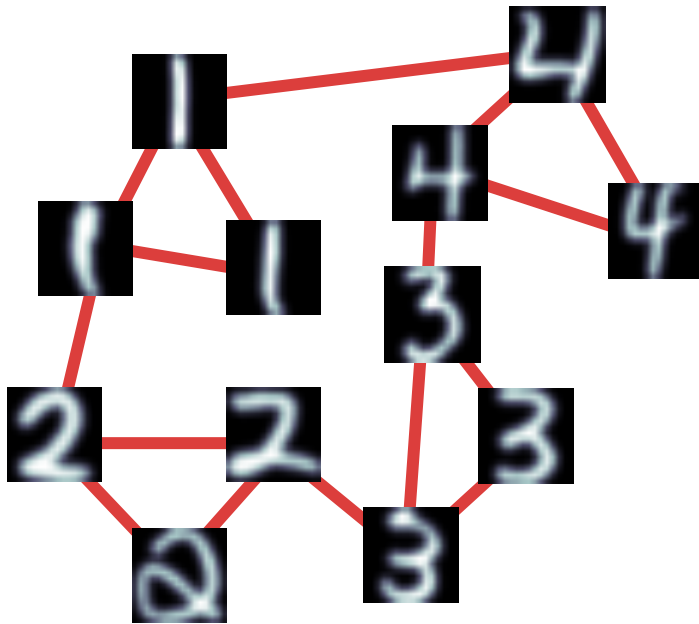


Algorithm

1. Run a diffusion for each label (possibly with neg. info from other classes)
2. Assign new labels based on the value of each diffusion

Semi-supervised graph-based learning

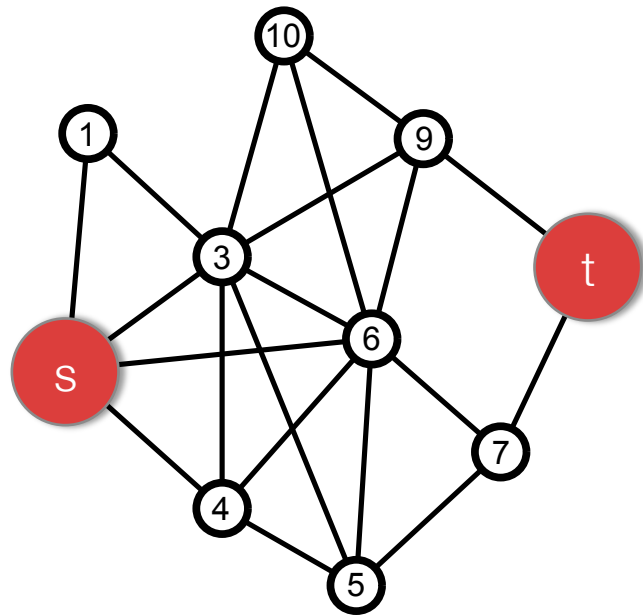
Given a graph, and a few labeled nodes, predict the labels on the rest of the graph.



Algorithm

1. Run a diffusion for each label (possibly with neg. info from other classes)
2. Assign new labels based on the value of each diffusion

The diffusions proposed for semi-supervised learning are s,t-cut minorants



In the unweighted case,
solve via max-flow.

In the weighted case,
solve via network simplex
or industrial LP.

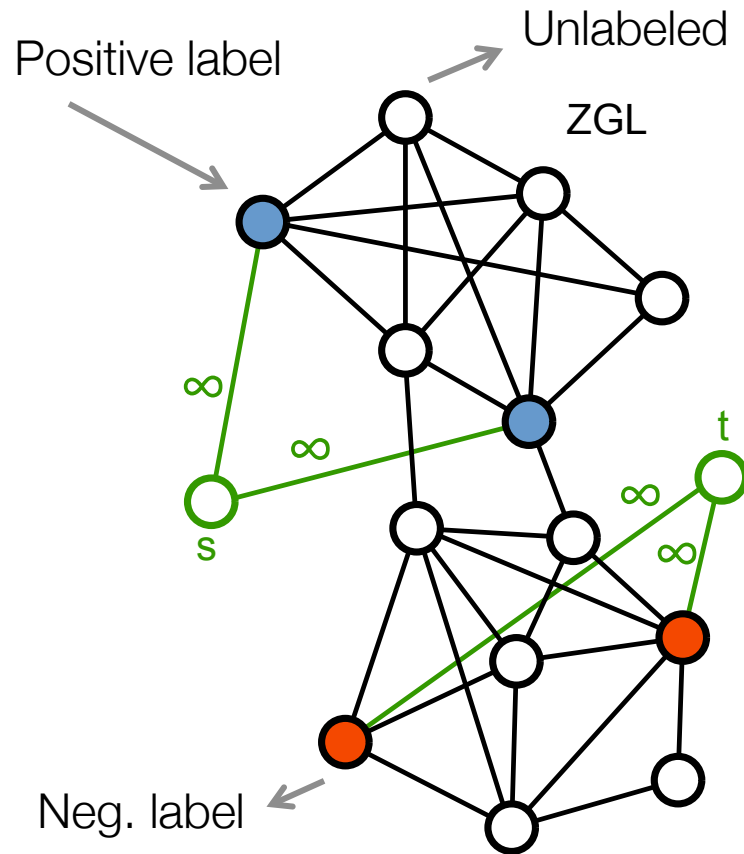
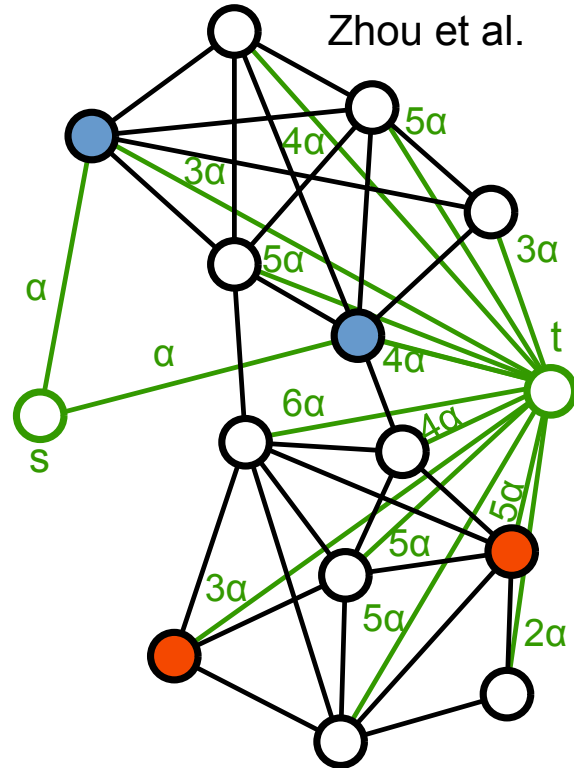
minimize $\sum_{ij \in E} C_{i,j} |x_i - x_j|$
subject to $x_s = 1, x_t = 0.$

MINCUT LP

minimize $\sqrt{\sum_{ij \in E} C_{i,j} |x_i - x_j|^2}$
subject to $x_s = 1, x_t = 0.$

Spectral minorant – lin. sys.

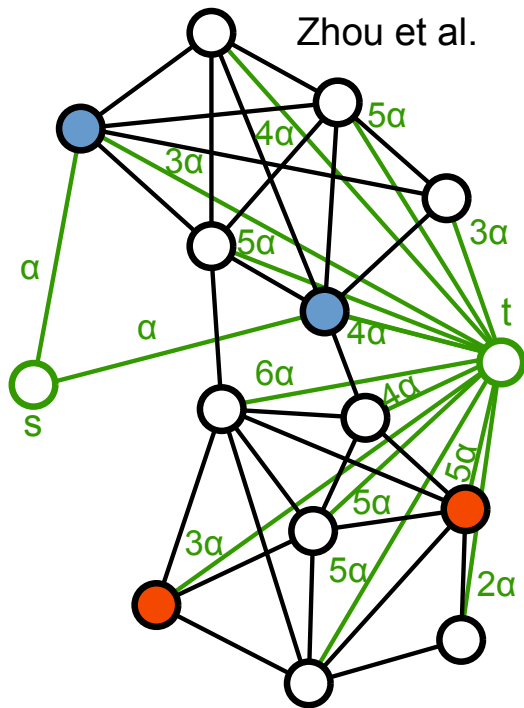
Representative cut problems



These help our intuition about the solutions
All spectral minorants are linear systems.

Implicit regularization views on the Zhou et al. diffusion

The Mahoney-Orecchia-Vishnoi (MOV) vector is a localized variation on the Fiedler vector to find a small conductance set nearby a seed set.



$$\begin{aligned} &\text{minimize} && \sqrt{\sum_{ij \in E} C_{i,j} |x_i - x_j|^2} \\ &\text{subject to} && x_s = 1, x_t = 0. \end{aligned}$$

RESULT

The spectral minorant of Zhou is equivalent to the weakly-local MOV solution.

PROOF

The two linear systems are the same (after working out a few equivalences).

IMPORTANCE

We'd expect Zhou to be "more robust"

A scalable, localized algorithm for Zhou et al's diffusion.

$$\begin{aligned} &\text{minimize} && \sqrt{\sum_{ij \in E} C_{i,j} |x_i - x_j|^2} \\ &\text{subject to} && x_s = 1, x_t = 0. \end{aligned}$$

RESULT

We can use a variation on coordinate descent methods related to the Andersen-Chung-Lang PUSH procedure to solve Zhou's diffusion in a scalable manner.

PROOF. See Gleich-Mahoney ICML '14

IMPORTANCE (1)

We should be able to make Zhou et al. scale.

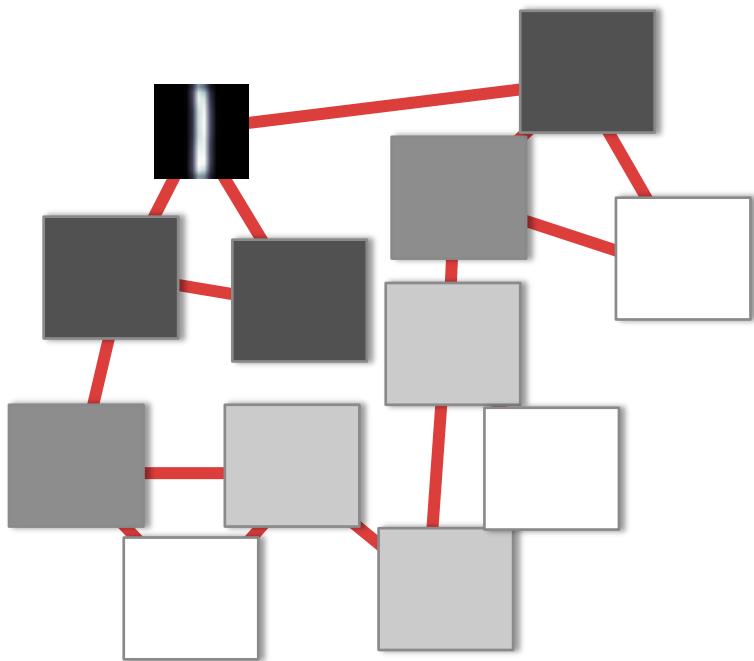
IMPORTANCE (2)

Using this algorithm adds another implicit regularization term that should further improve robustness!

$$\begin{aligned} &\text{minimize} && \sum_{ij \in E} C_{i,j} |x_i - x_j|^2 + \tau \sum_{i \in V} d_i x_i \\ &\text{subject to} && x_s = 1, x_t = 0, x_i \geq 0. \end{aligned}$$

Semi-supervised graph-based learning

Given a graph, and a few labeled nodes, predict the labels on the rest of the graph.



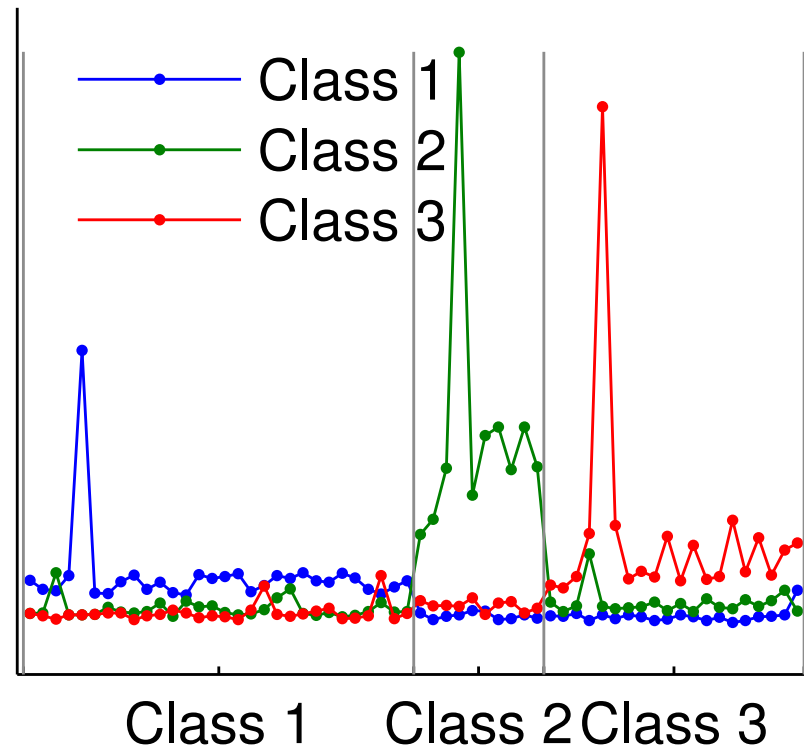
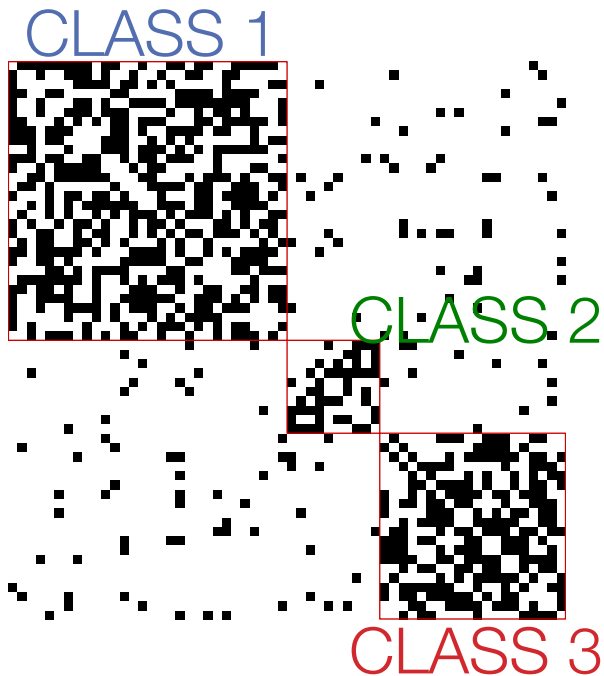
Algorithm

1. Run a diffusion process for each label (positive and negative) with neg. info from other classes
2. Assign new labels based on the value of each diffusion

New interpretation of the diffusion process: scalable and more robust!

Traditional rounding methods for SSL are value-based

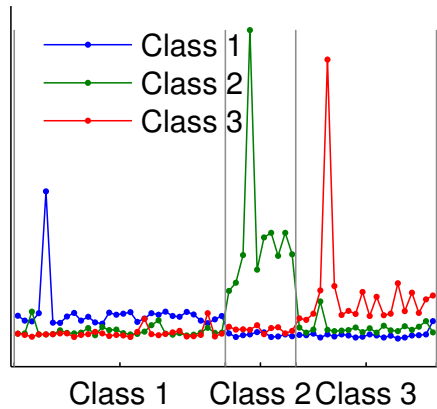
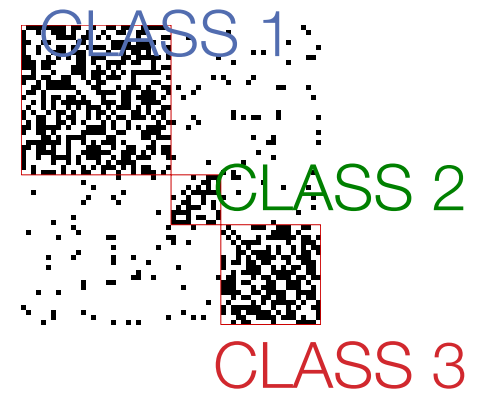
Zhou's diffusion



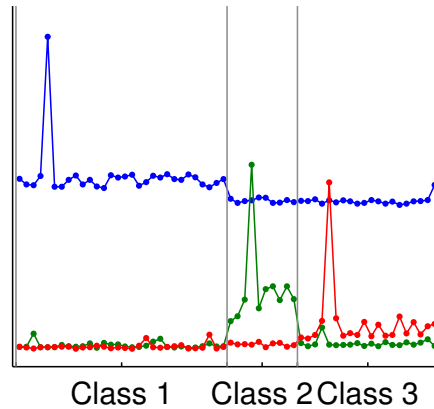
VALUE-BASED

Use the largest value of the diffusion to pick the label.

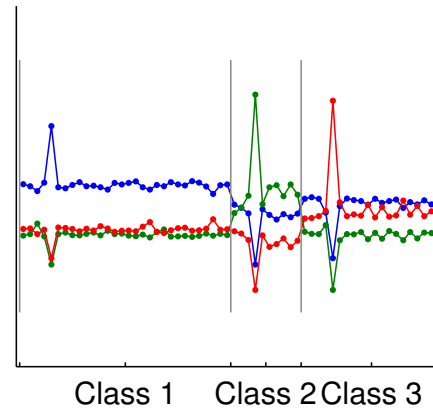
But value based rounding doesn't work for all diffusions



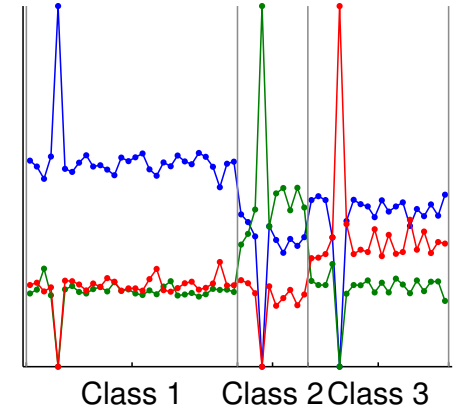
(b) Zhou et al., $l = 3$



(c) Andersen-Lang, $l = 3$



(d) Joachims, $l = 3$



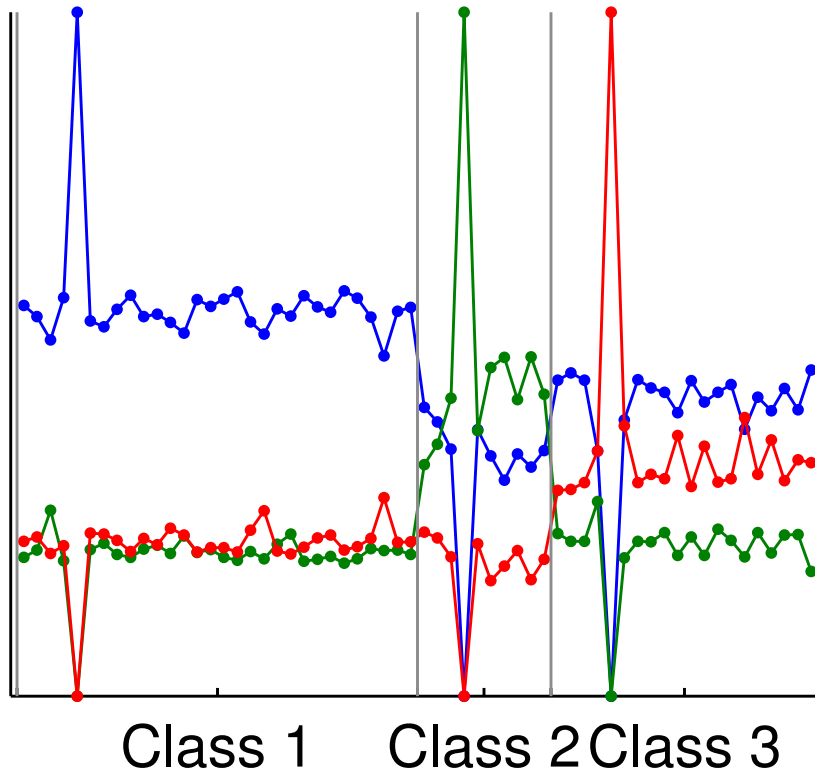
(e) ZGL, $l = 3$

VALUE-BASED rounding fails for most of these diffusions

BUT

There is still a signal there!

Rank-based rounding is far more robust.



NEW IDEA

Look at the RANK of the item in each diffusion instead of its VALUE.

JUSTIFICATION

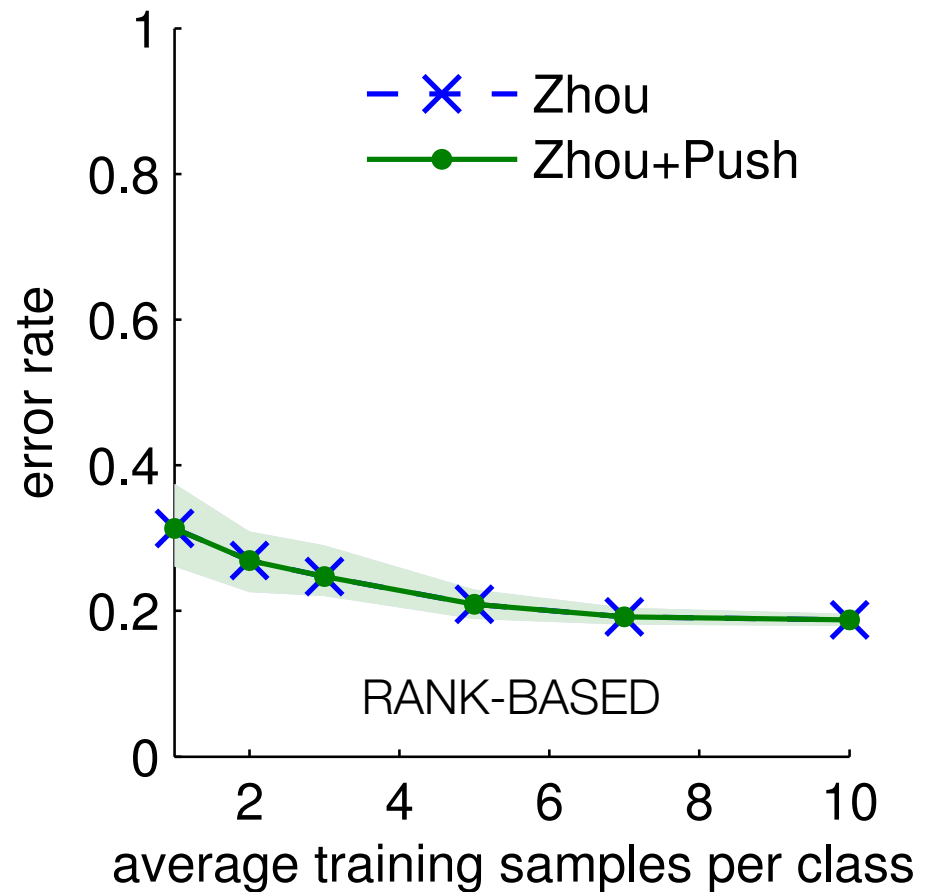
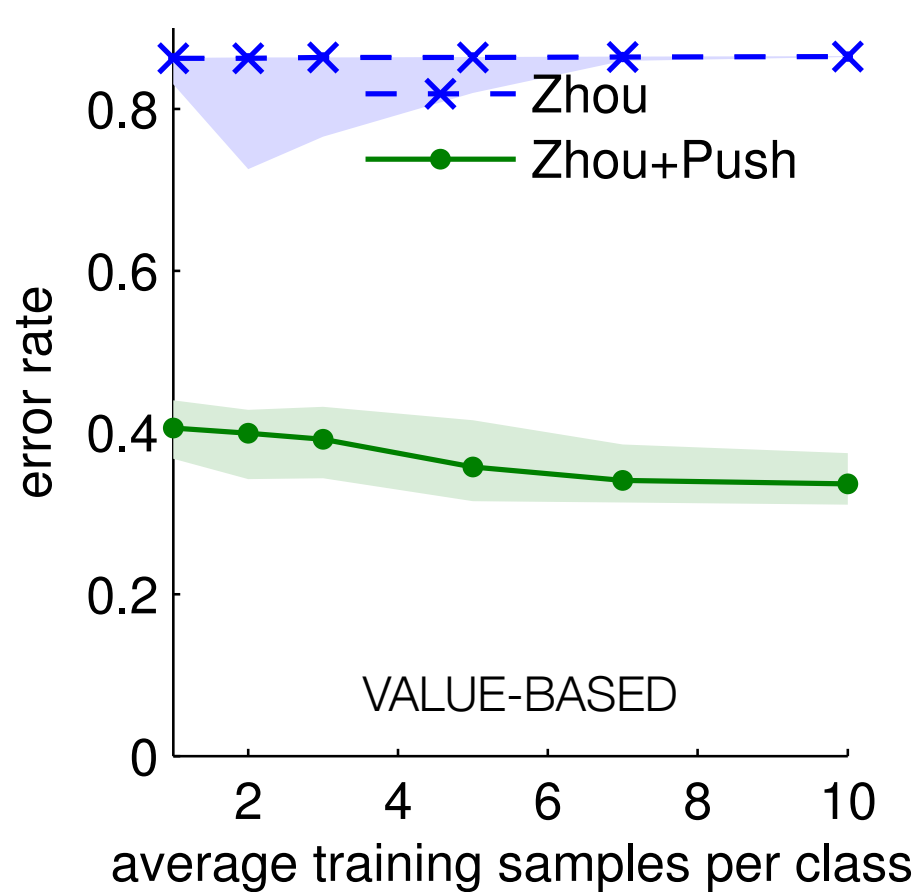
Based on the idea of sweep-cut rounding in spectral methods (use the order induced by the eigenvector, not its values)

IMPACT

Much more robust rounding to labels

Rank-based rounding has a big impact on a real-study.

We used the digit prediction task out of Zhou's paper and added just a bit of noise as label errors and switched parameters.



Solutions Paths

One benefit of the weak convergence algorithms...

A direct decomposition is a black box:
Feed in input, get output.

In contrast, the iterative nature of “push” means running the algorithm is essentially “watching” the diffusion process occur.

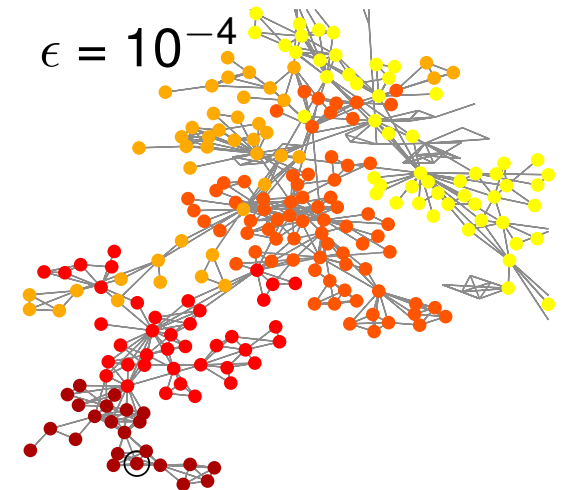
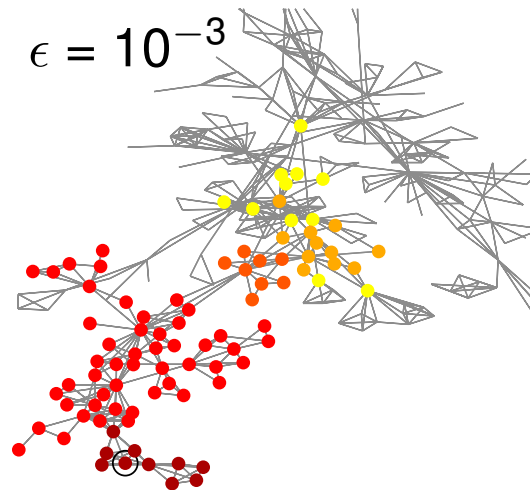
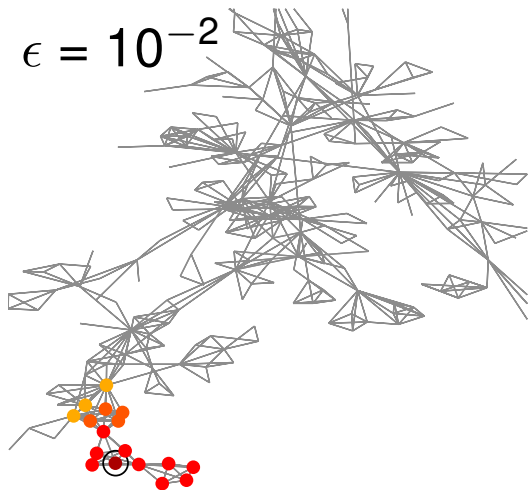
We get more information this way!

Solutions Paths

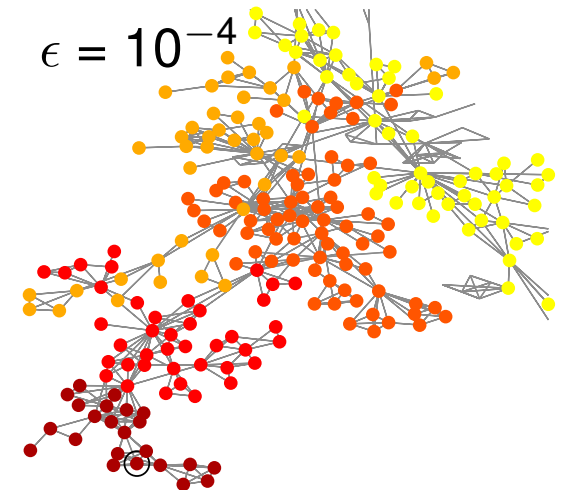
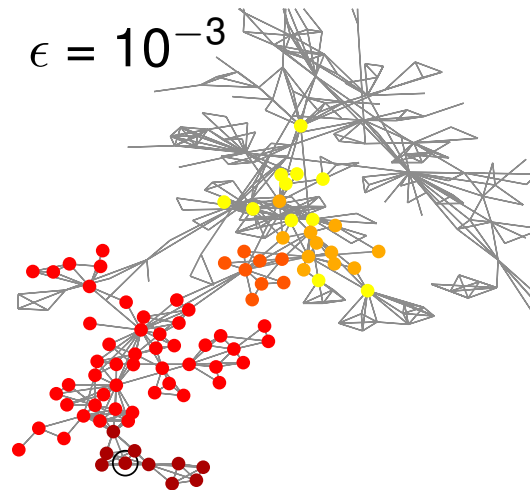
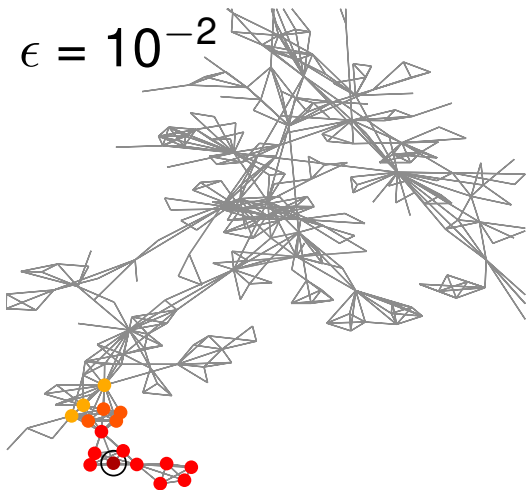
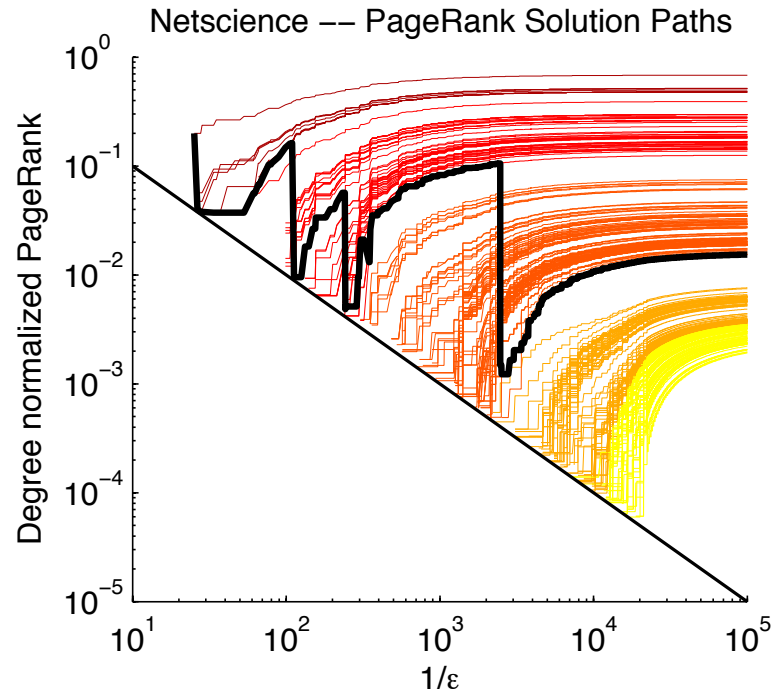
One benefit of the weak convergence algorithms...

A direct decomposition is a black box:
Feed in input, get output.

In contrast, the iterative nature of “push” means
running the algorithm is essentially “watching” the
diffusion process occur.



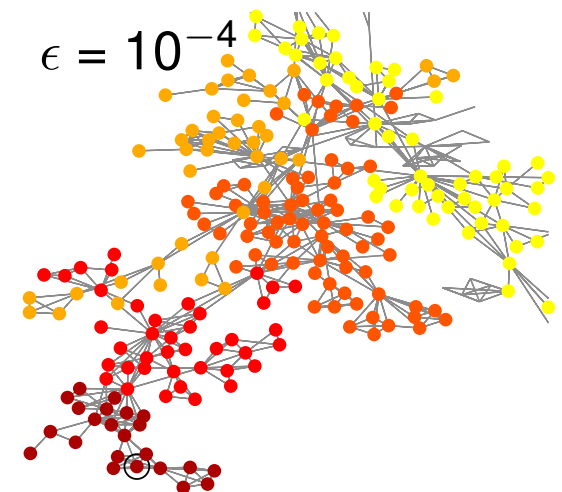
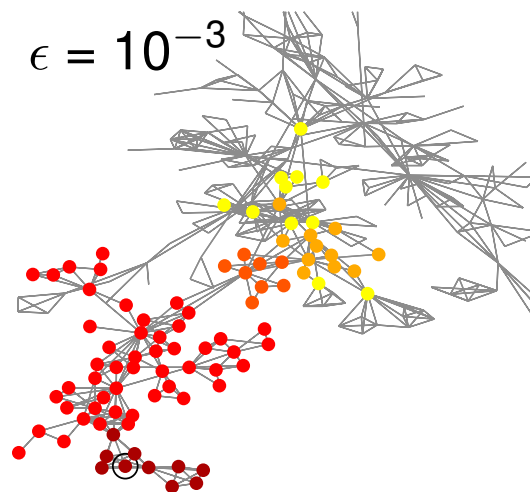
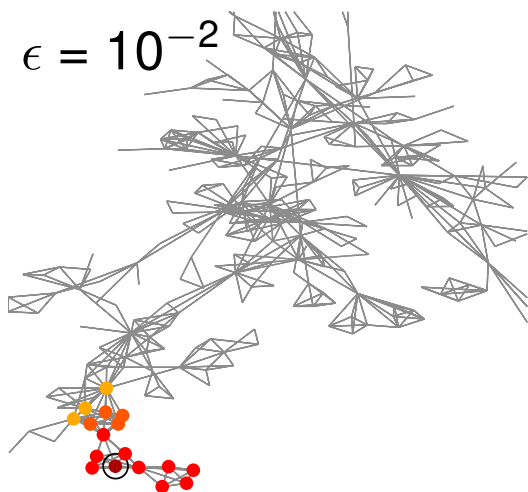
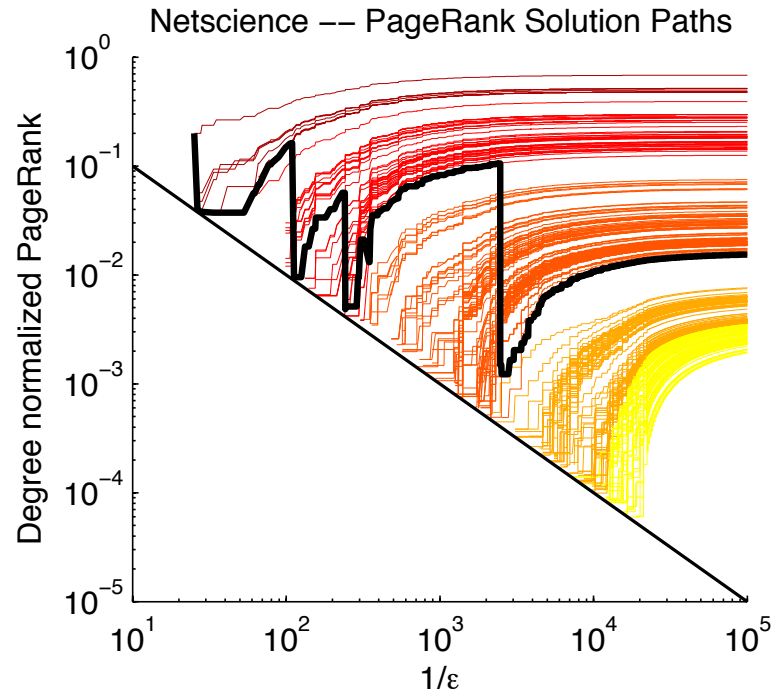
Solutions Paths



Solutions Paths

Each curve is a node. Its value increases as ϵ goes to 0.

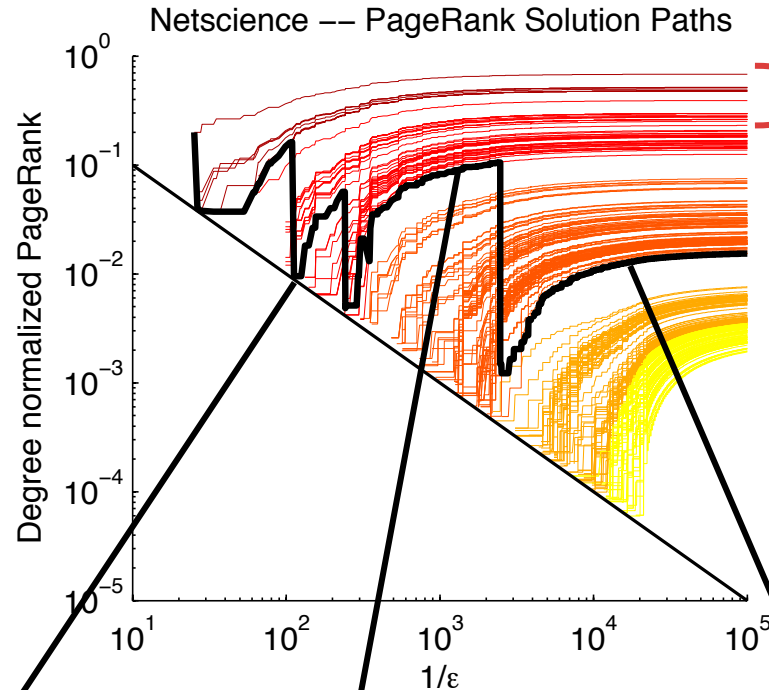
Thick black line shows set of best conductance.



Solutions Paths

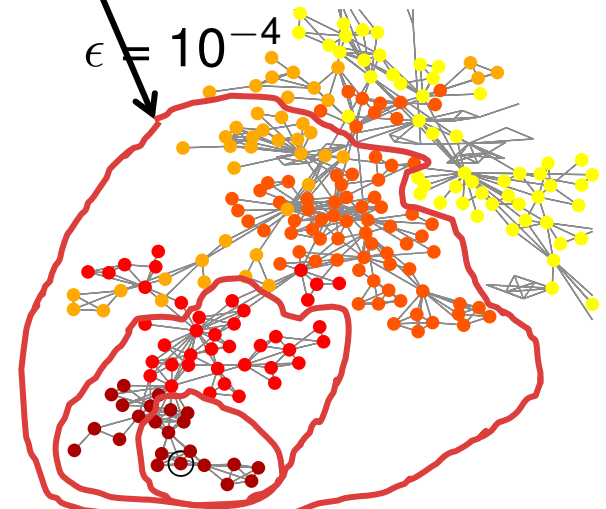
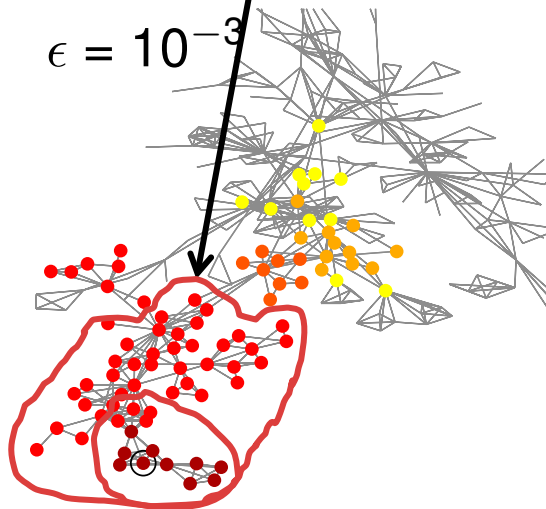
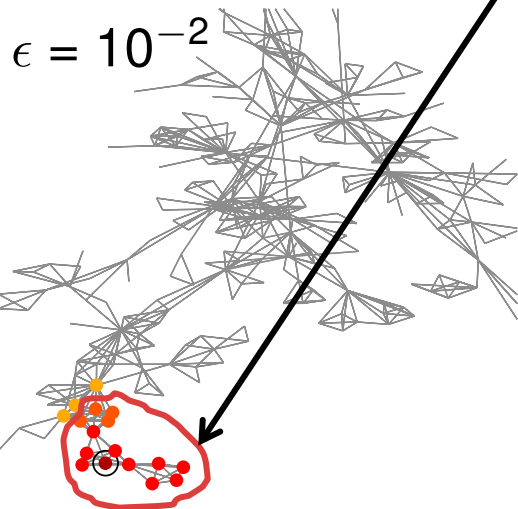
Each curve is a node. Its value increases as ϵ goes to 0.

Thick black line shows set of best conductance.



Bundles of curves are good clusters

Paths identify *nested* clusters



Solutions Paths

Locate nested, good-conductance sets that a single diffusion + sweep could miss.

Can be done efficiently because the constant-time approach to computing diffusions enables efficient storage and analysis of the push process

Total Paths work (for PageRank): $O\left(\frac{1}{\epsilon(1-\alpha)}\right)^2$
Still efficient!

Recent directions

1. Diffusions in time-dependent networks
 - Grindrod et al. 2011 (Katz); Gleich & Rossi, 2014 (PageRank); Grindrod & Higham, 2014 (Katz again...)

Open issues

1. Parameter selection
2. Tuned diffusions