

Link Prediction in Networks

Thanks to Jure Leskovec, Stanford and Panayiotis Tsaparas, Univ. of Ioannina for slides

Agenda

- Link Prediction in Networks
 - Estimating Scores for Missing Edges
 - Classification Approach (Omitted)
- Case studies:
 - Facebook: Supervised Random Walks for Link Prediction
 - Twitter: The who to follow service at Twitter

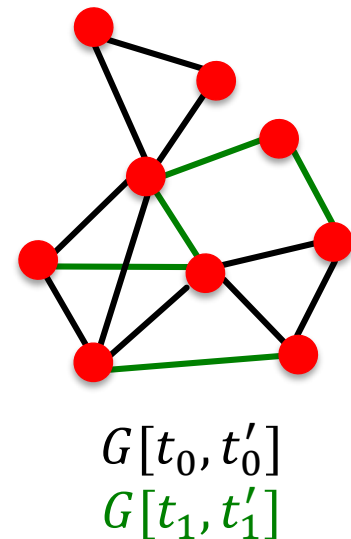
Link Prediction Motivation

- Recommending *new friends* in online social networks
- Predicting the participation of *actors* in events
- Suggesting *interactions* between the members of a company/organization that are external to the hierarchical structure of the organization itself
- Predicting *connections* between members of communities/organizations who have not been directly observed together
- Suggesting *collaborations* between researchers based on co-authorship
- Overcoming the data-sparsity problem *in recommender systems* using collaborative filtering

Link Prediction in Networks

■ The link prediction task:

- Given $G[t_0, t'_0]$ a graph on edges up to time t'_0 , **output a ranked list L** of links (not in $G[t_0, t'_0]$) that are predicted to appear in $G[t_1, t'_1]$



■ Evaluation:

- $n = |E_{new}|$: # new edges that appear during the test period $[t_1, t'_1]$
- Take top n elements of L and count correct edges

Link Prediction

- Predict links in a evolving collaboration network

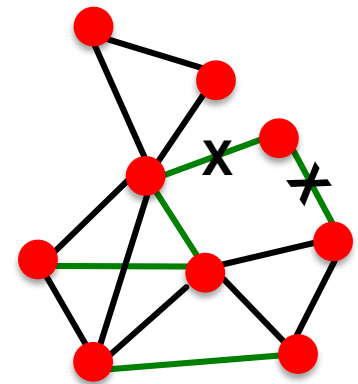
	training period			Core		
	authors	papers	collaborations ¹	authors	$ E_{old} $	$ E_{new} $
astro-ph	5343	5816	41852	1561	6178	5751
cond-mat	5469	6700	19881	1253	1899	1150
gr-qc	2122	3287	5724	486	519	400
hep-ph	5414	10254	47806	1790	6654	3294
hep-th	5241	9498	15842	1438	2311	1576

- **Core:** Because network data is very sparse
 - Consider only nodes with degree of at least 3
 - Because we don't know enough about these nodes to make good inferences

Link Prediction

■ Methodology:

- For each pair of nodes (x,y) compute score $c(x,y)$
 - For example, $c(x,y)$ could be the # of common neighbors of x and y
- Sort pairs (x,y) by the decreasing score $c(x,y)$
 - **Note:** Only consider/predict edges where both endpoints are in the core ($\text{deg.} \geq 3$)
- **Predict top n pairs as new links**
- **See which of these links actually appear in $G[t_1, t'_1]$**



Link Prediction

- **Different scoring functions** $c(x, y) =$
 - **Graph distance:** (negated) Shortest path length
 - **Common neighbors:** $|\Gamma(x) \cap \Gamma(y)|$
 - **Jaccard's coefficient:** $|\Gamma(x) \cap \Gamma(y)| / |\Gamma(x) \cup \Gamma(y)|$
 - **Adamic/Adar:** $\sum_{z \in \Gamma(x) \cap \Gamma(y)} 1 / \log |\Gamma(z)|$
 - **Preferential attachment:** $|\Gamma(x)| \cdot |\Gamma(y)|$ $\Gamma(x)$... neighbors of node x
 - **PageRank:** $r_x(y) + r_y(x)$
 - $r_x(y)$... stationary distribution score of y under the random walk:
 - with prob. 0.15, jump to x
 - with prob. 0.85, go to random neighbor of current node
- **Then, for a particular choice of $c(\cdot)$**
 - For every pair of nodes (x, y) compute $c(x, y)$
 - Sort pairs (x, y) by the decreasing score $c(x, y)$
 - **Predict top n pairs as new links**

Link Prediction Methods

Link Prediction Methods

- How to assign the score $c(x, y)$ for each pair (x, y) ?
 - Some form of **similarity** between x and y
 - Some form of **node proximity** between x and y
- Methods
 - Neighborhood-based (shared neighbors)
 - Network proximity based (paths between x and y)
 - Other

Methods for Link Prediction

Neighborhood-based

Neighborhood-based Methods

Let $\Gamma(x)$ be the set of neighbors of x in G_{old}

- Methods
 - Common Neighbors Overlap
 - Jaccard
 - Adamic/Adar
 - Preferential Attachment

Neighborhood-based Methods

Intuition: The **larger the overlap** of the neighbors of two nodes, the **more likely** the nodes to be linked in the future

- **Common neighbors**

- **A:** adjacency matrix, $\mathbf{A}_{x,y}^2$: #paths of length 2

$$\text{score}(x, y) = |\Gamma(x) \cap \Gamma(y)|$$

- **Jaccard coefficient**

- The probability that both **x** and **y** have a feature for a randomly selected feature that either **x** or **y** has

$$\text{score}(x, y) = \frac{|\Gamma(x) \cap \Gamma(y)|}{|\Gamma(x) \cup \Gamma(y)|}$$

Neighborhood-based Methods

■ Adamic/Adar

- Assigns large weights to common neighbors z of x and y which themselves have few neighbors (weight rare features more heavily)

$$\text{score}(x, y) = \sum_{z \in \Gamma(x) \cap \Gamma(y)} \frac{1}{\log |\Gamma(z)|}$$

■ Preferential attachment

- Based on the premise that the probability that a new edge has node x as its endpoint is proportional to $|\Gamma(x)|$, i.e., nodes like to form ties with ‘popular’ nodes

$$\text{score}(x, y) = |\Gamma(x)| |\Gamma(y)|$$

Methods for Link Prediction

Network proximity based

Network Proximity Methods

Intuition: The “**closer**” two nodes are in the network, the **more likely** are to be linked in the future

- **Methods**
 - based on **shortest path length** between **x** and **y**
 - based on **all paths** between **x** and **y**
 - Katz _{β} measure (unweighted, weighted)
 - Random walk-based
 - hitting time
 - commute time
 - Rooted PageRank
 - SimRank

Shortest Path Based

For $x, y \in V \times V - E_{\text{old}}$,

score(x, y) = (negated) **length of shortest path** between **x** and **y**

If there are more than **n** pairs of nodes tied for the shortest path length, order them at random

Ensemble of All Paths

- **Katz _{β} measure**

$$\text{score}(x, y) := \sum_{\ell=1}^{\infty} \beta^{\ell} \cdot |\text{paths}_{x,y}^{(\ell)}|$$

- Sum over all paths of length ℓ
- $0 < \beta < 1$: a parameter of the predictor, exponentially damped to count short paths more heavily

Ensemble of All Paths

- **Katz_β measure**

$$\text{score}(x, y) := \sum_{\ell=1}^{\infty} \beta^{\ell} \cdot |\text{paths}_{x,y}^{(\ell)}|$$

$$\sum_{l=1}^{\infty} \beta^l \cdot |\text{paths}_{xy}^{(l)}| = \beta A_{xy} + \beta^2 (A^2)_{xy} + \beta^3 (A^3)_{xy} + \dots$$

- **Unweighted** version: $\text{path}_{x,y}(1) = 1$, if \mathbf{x} and \mathbf{y} have collaborated, **0** otherwise
- **Weighted** version: $\text{path}_{x,y}(1) = \text{\#times}$ \mathbf{x} and \mathbf{y} have collaborated

Random Walk Based

Consider a random walk on \mathbf{G}_{old} that starts at \mathbf{x} and iteratively moves to a neighbor of \mathbf{x} chosen uniformly at random from $\Gamma(\mathbf{x})$

- **Hitting $H_{\mathbf{x},\mathbf{y}}$** (from \mathbf{x} to \mathbf{y}): the expected number of steps it takes for the random walk starting at \mathbf{x} to reach \mathbf{y}

$$\text{score}(\mathbf{x}, \mathbf{y}) = -H_{\mathbf{x},\mathbf{y}}$$

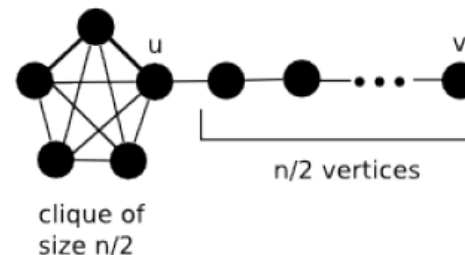
- **Commute Time $C_{\mathbf{x},\mathbf{y}}$** (from \mathbf{x} to \mathbf{y}): the expected number of steps to travel from \mathbf{x} to \mathbf{y} and from \mathbf{y} to \mathbf{x}

$$\text{score}(\mathbf{x}, \mathbf{y}) = -(H_{\mathbf{x},\mathbf{y}} + H_{\mathbf{y},\mathbf{x}})$$

Not symmetric, can be shown

$$h_{vu} = \Theta(n^2)$$

$$h_{uv} = \Theta(n^3)$$



Random Walk Based

- The hitting time and commute time measures are sensitive to parts of the graph far away from x and y → periodically **reset the walk**
- Random walk on G_{old} that starts at x and has a probability α of returning to x at each step
- **Rooted PageRank**
 - Starts from x
 - with probability $(1 - \alpha)$ moves to a random neighbor
 - with probability α returns to x

score(x, y) = stationary probability of y in a
rooted PageRank

SimRank

Intuition: Two objects are *similar*, if they are *related to similar objects*

- Two objects x and y are *similar*, if they are related to objects a and b respectively and a and b are themselves similar

$$\text{similarity}(x, y) := \gamma \cdot \frac{\sum_{a \in \Gamma(x)} \sum_{b \in \Gamma(y)} \text{similarity}(a, b)}{|\Gamma(x)| \cdot |\Gamma(y)|}$$

- Expresses the average similarity between neighbors of x and neighbors of y
score(x, y) = similarity(x, y)

Methods for Link Prediction

Other Methods

Other Methods

- Low-rank Approximations
- Unseen bigrams
- High-level Clustering

Low Rank Approximations

Intuition: represent the adjacency matrix \mathbf{M} with a lower rank matrix \mathbf{M}_k

- Method

- Apply SVD (singular value decomposition)
- Obtain the **rank-k** matrix that best approximates \mathbf{M}

Singular Value Decomposition

$$A = U \Sigma V^T = \begin{bmatrix} \vec{u}_1 & \vec{u}_2 & \cdots & \vec{u}_r \end{bmatrix} \begin{bmatrix} \sigma_1 & & & \\ & \sigma_2 & & \\ & & \ddots & \\ & & & \sigma_r \end{bmatrix} \begin{bmatrix} \vec{v}_1 \\ \vec{v}_2 \\ \vdots \\ \vec{v}_r \end{bmatrix}$$

$[n \times r] \quad [r \times r] \quad [r \times n]$

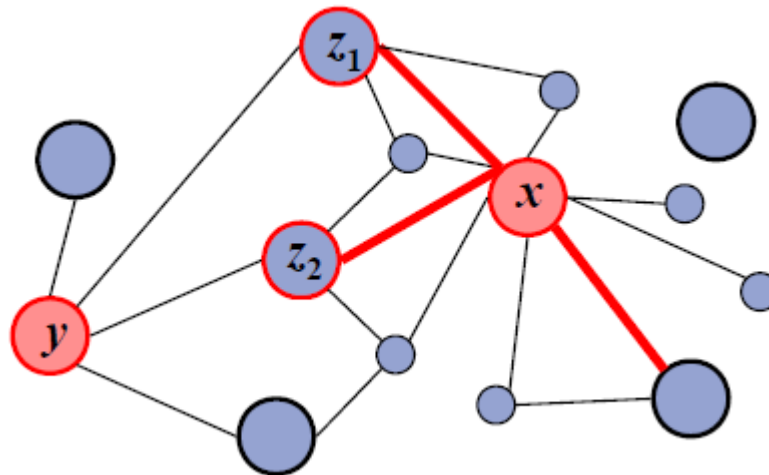
- r : rank of matrix A
- $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r$: singular values (square roots of eig-vals AA^T, A^TA)
- $\vec{u}_1, \vec{u}_2, \dots, \vec{u}_r$: left singular vectors (eig-vectors of AA^T)
- $\vec{v}_1, \vec{v}_2, \dots, \vec{v}_r$: right singular vectors (eig-vectors of A^TA)
- $A = \sigma_1 \vec{u}_1 \vec{v}_1^T + \sigma_2 \vec{u}_2 \vec{v}_2^T + \dots + \sigma_r \vec{u}_r \vec{v}_r^T$
- $A_k = \sigma_1 \vec{u}_1 \vec{v}_1^T + \sigma_2 \vec{u}_2 \vec{v}_2^T + \dots + \sigma_k \vec{u}_k \vec{v}_k^T, k \in \{1, 2, \dots, r\}$

Unseen bigrams

- **Unseen bigrams**: Predict pairs of words that co-occur in a test corpus, but not in the corresponding training corpus
- Not just $\text{score}(\mathbf{x}, \mathbf{y})$ but $\text{score}(\mathbf{z}, \mathbf{y})$ for nodes \mathbf{z} that are similar to \mathbf{x} --- $\mathbf{S}_x^{(\delta)}$: the δ nodes *most related to x*

$$\text{score}_{unweighted}^*(x, y) := \left| \{z : z \in \Gamma(y) \cap S_x^{(\delta)}\} \right|$$

$$\text{score}_{weighted}^*(x, y) := \sum_{z \in \Gamma(y) \cap S_x^{(\delta)}} \text{score}(x, z)$$



Clustering

- Compute **score(x, y)** for all edges in E_{old}
- Delete the **(1-p)** fraction of the edges whose score is the lowest, for some parameter p
- Re-compute **score(x, y)** for all pairs in the subgraph

Evaluation & Results

How to Evaluate the Prediction

- Each link predictor p outputs a ranked list L_p of pairs in $V \times V - E_{\text{old}}$ in decreasing order of confidence
 - focus on **Core** network, ($d > 3$)
$$E^*_{\text{new}} = E_{\text{new}} \cap (\text{Core} \times \text{Core}) = |E^*_{\text{new}}|$$
- Evaluation method: **Size of intersection** of
 - the first n edge predictions from L_p that are in $\text{Core} \times \text{Core}$, and
 - the actual set E^*_{new}

How many of the (relevant) top- n predictions are correct (precision?)

Evaluation: Baseline Predictor

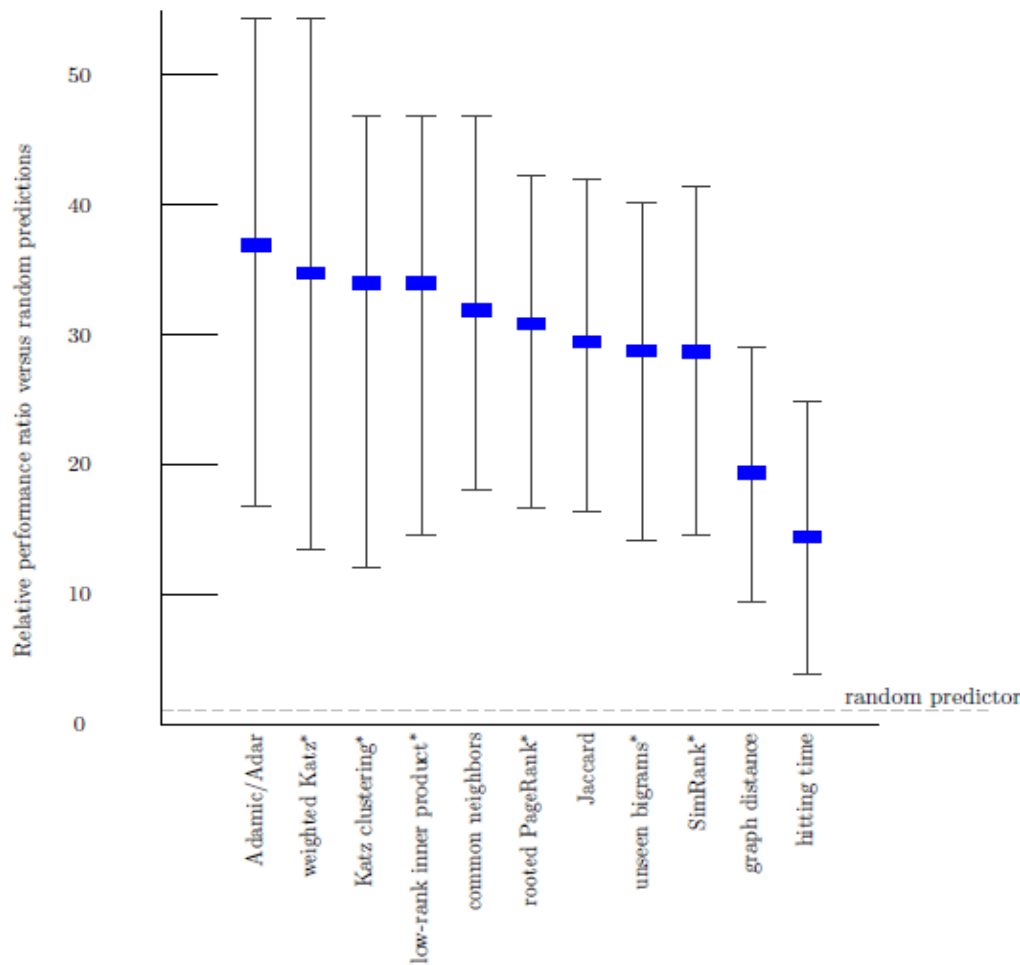
- **Random Predictor:** Randomly select pairs of authors who did not collaborate in the training interval
 - Probability that a random prediction is correct:

$$\frac{|E_{new}|}{\binom{|Core|}{2} - |E_{old}|}$$

In the datasets, from 0.15% (cond-mat) to 0.48% (astro-ph)

Average Relevance Performance

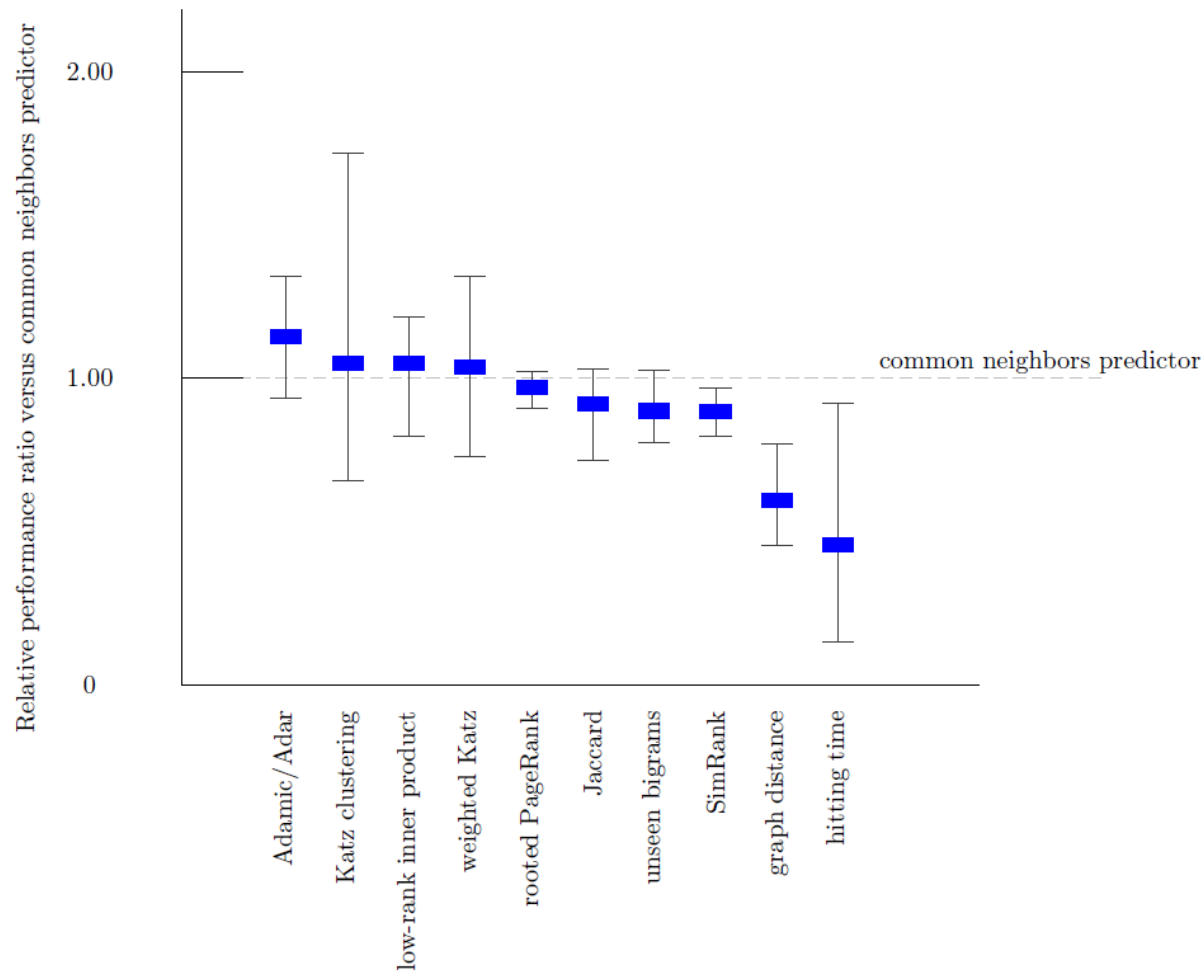
■ Improvement over random predictor



- average ratio over the five datasets of the given predictor's performance *versus a baseline* predictor's performance.
- the error bars indicate the minimum and maximum of this ratio over the five datasets.
- the parameters for the starred predictors are: (1) for weighted Katz, $\beta = 0.005$; (2) for Katz clustering, $\beta_1 = 0.001$; $\rho = 0.15$; $\beta_2 = 0.1$; (3) for low-rank inner product, rank = 256; (4) for rooted Pagerank, $\alpha = 0.15$; (5) for unseen bigrams, unweighted, common neighbors with $\delta = 8$; and (6) for SimRank, $C(\gamma) = 0.8$.

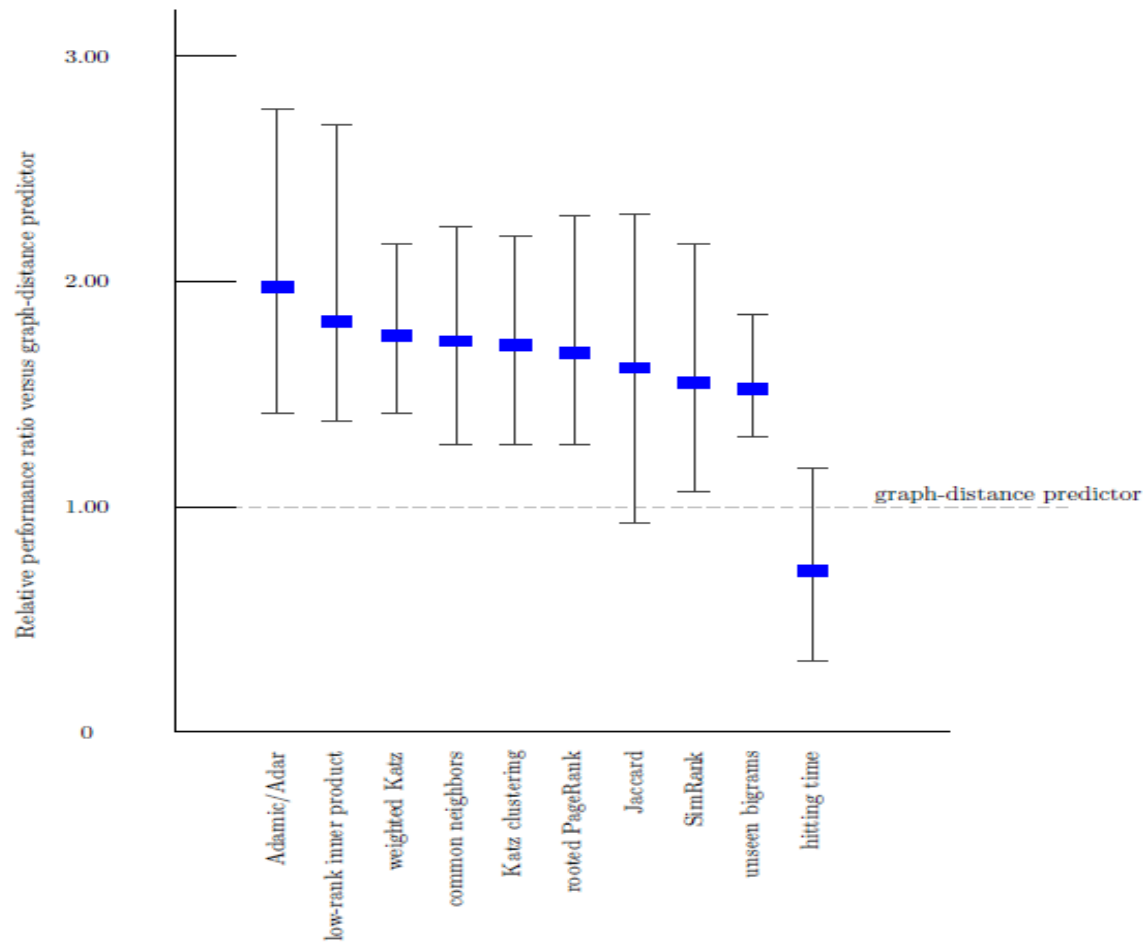
Results: Common Neighbors

- Improvement over #common neighbors

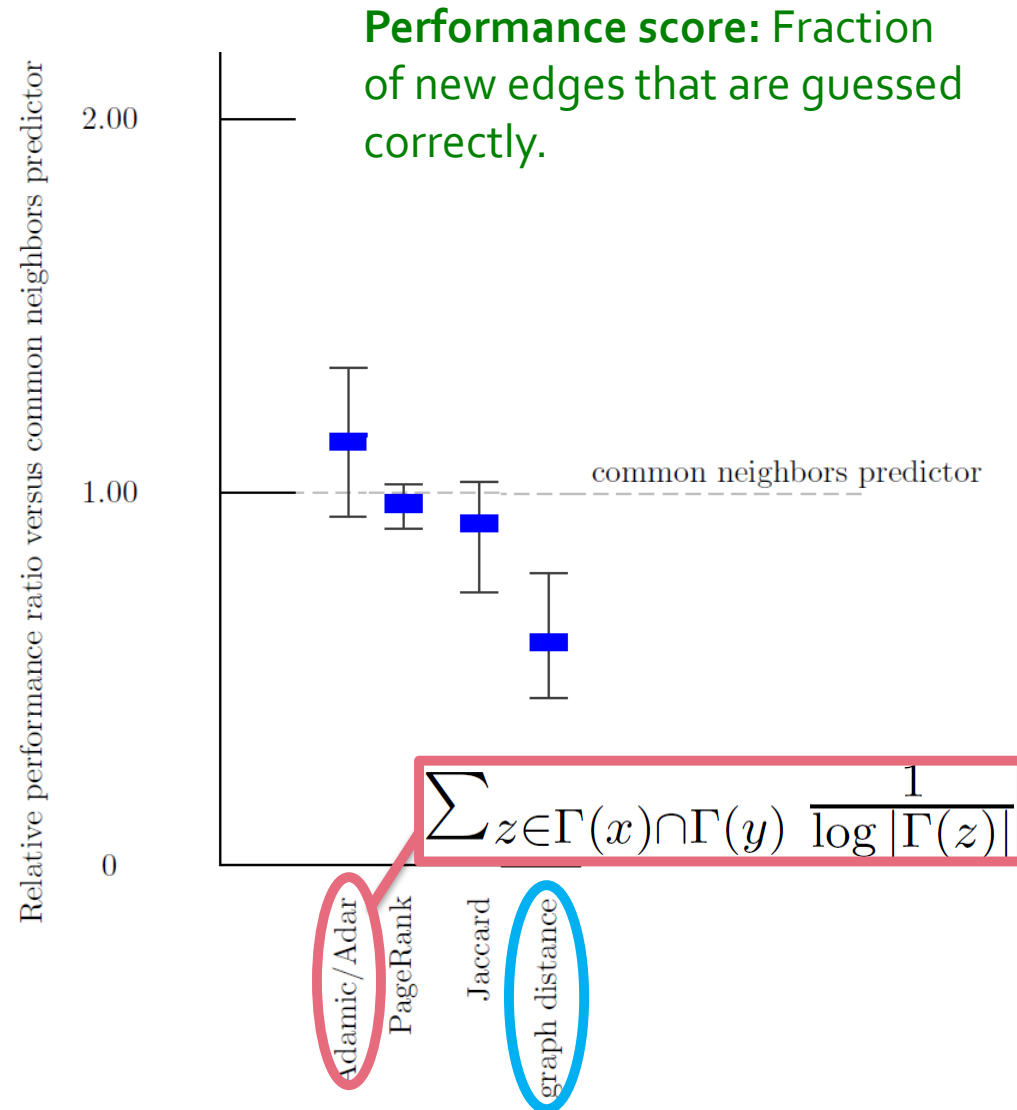
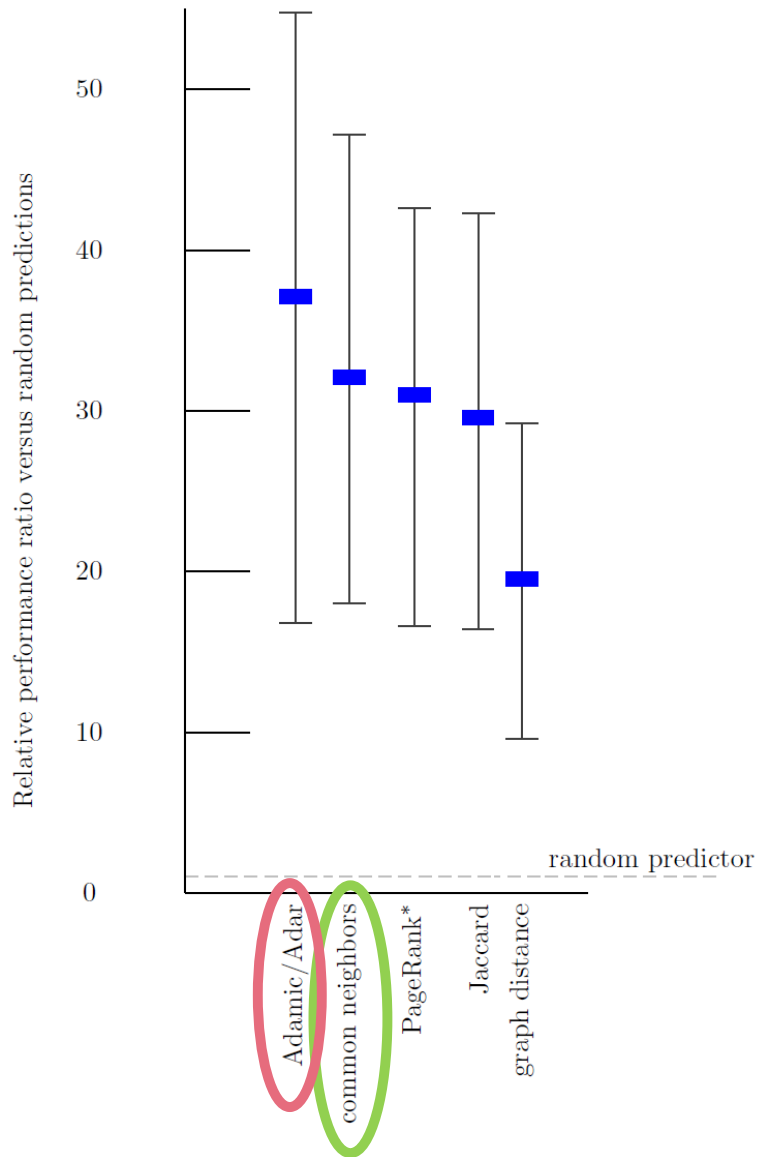


Results: Common Neighbors

- Improvement over graph distance predictor



Results: Improvement



Factor Improvement Over Random

predictor		astro-ph	cond-mat	gr-qc	hep-ph	hep-th
probability that a random prediction is correct		0.475%	0.147%	0.341%	0.207%	0.153%
graph distance (all distance-two pairs)		9.4	25.1	21.3	12.0	29.0
common neighbors		18.0	40.8	27.1	26.9	46.9
preferential attachment		4.7	6.0	7.5	15.2	7.4
Adamic/Adar		16.8	54.4	30.1	33.2	50.2
Jaccard		16.4	42.0	19.8	27.6	41.5
SimRank $\gamma = 0.8$		14.5	39.0	22.7	26.0	41.5
hitting time		6.4	23.7	24.9	3.8	13.3
hitting time—normed by stationary distribution		5.3	23.7	11.0	11.3	21.2
commute time		5.2	15.4	33.0	17.0	23.2
commute time—normed by stationary distribution		5.3	16.0	11.0	11.3	16.2
rooted PageRank	$\alpha = 0.01$	10.8	27.8	33.0	18.7	29.1
	$\alpha = 0.05$	13.8	39.6	35.2	24.5	41.1
	$\alpha = 0.15$	16.6	40.8	27.1	27.5	42.3
	$\alpha = 0.30$	17.1	42.0	24.9	29.8	46.5
	$\alpha = 0.50$	16.8	40.8	24.2	30.6	46.5
Katz (weighted)	$\beta = 0.05$	3.0	21.3	19.8	2.4	12.9
	$\beta = 0.005$	13.4	54.4	30.1	24.0	51.9
	$\beta = 0.0005$	14.5	53.8	30.1	32.5	51.5
Katz (unweighted)	$\beta = 0.05$	10.9	41.4	37.4	18.7	47.7
	$\beta = 0.005$	16.8	41.4	37.4	24.1	49.4
	$\beta = 0.0005$	16.7	41.4	37.4	24.8	49.4

Factor Improvement Over Random

predictor		astro-ph	cond-mat	gr-qc	hep-ph	hep-th
probability that a random prediction is correct		0.475%	0.147%	0.341%	0.207%	0.153%
graph distance (all distance-two pairs)		9.4	25.1	21.3	12.0	29.0
common neighbors		18.0	40.8	27.1	26.9	46.9
Low-rank approximation: Inner product	rank = 1024	15.2	53.8	29.3	34.8	49.8
	rank = 256	14.6	46.7	29.3	32.3	46.9
	rank = 64	13.0	44.4	27.1	30.7	47.3
	rank = 16	10.0	21.3	31.5	27.8	35.3
	rank = 4	8.8	15.4	42.5	19.5	22.8
	rank = 1	6.9	5.9	44.7	17.6	14.5
Low-rank approximation: Matrix entry	rank = 1024	8.2	16.6	6.6	18.5	21.6
	rank = 256	15.4	36.1	8.1	26.2	37.4
	rank = 64	13.7	46.1	16.9	28.1	40.7
	rank = 16	9.1	21.3	26.4	23.1	34.0
	rank = 4	8.8	15.4	39.6	20.0	22.4
	rank = 1	6.9	5.9	44.7	17.6	14.5
Low-rank approximation: Katz ($\beta = 0.005$)	rank = 1024	11.4	27.2	30.1	27.0	32.0
	rank = 256	15.4	42.0	11.0	34.2	38.6
	rank = 64	13.1	45.0	19.1	32.2	41.1
	rank = 16	9.2	21.3	27.1	24.8	34.9
	rank = 4	7.0	15.4	41.1	19.7	22.8
	rank = 1	0.4	5.9	44.7	17.6	14.5
unseen bigrams (weighted)	common neighbors, $\delta = 8$	13.5	36.7	30.1	15.6	46.9
	common neighbors, $\delta = 16$	13.4	39.6	38.9	18.5	48.6
	Katz ($\beta = 0.005$), $\delta = 8$	16.8	37.9	24.9	24.1	51.1
	Katz ($\beta = 0.005$), $\delta = 16$	16.5	39.6	35.2	24.7	50.6
unseen bigrams (unweighted)	common neighbors, $\delta = 8$	14.1	40.2	27.9	22.2	39.4
	common neighbors, $\delta = 16$	15.3	39.0	42.5	22.0	42.3
	Katz ($\beta = 0.005$), $\delta = 8$	13.1	36.7	32.3	21.6	37.8
	Katz ($\beta = 0.005$), $\delta = 16$	10.3	29.6	41.8	12.2	37.8
clustering: Katz ($\beta_1 = 0.001, \beta_2 = 0.1$)	$\rho = 0.10$	7.4	37.3	46.9	32.9	37.8
	$\rho = 0.15$	12.0	46.1	46.9	21.0	44.0
	$\rho = 0.20$	4.6	34.3	19.8	21.2	35.7
	$\rho = 0.25$	3.3	27.2	20.5	19.4	17.4

Evaluation: Prediction Overlap

[illegible]

How similar are the predictions made by the different methods?

common predictions

[illegible]

correct predictions

Extensions

- Improve **performance**. Even the best (Katz clustering on gr-qc) correct on only about 16% of its prediction
- Improve **efficiency** on very large networks (approximation of distances)
- Treat more **recent** links (e.g., collaborations) as more important
- Additional information (paper titles, author institutions, etc) latently present in the graph

Facebook: Supervised Random Walks for Link Prediction

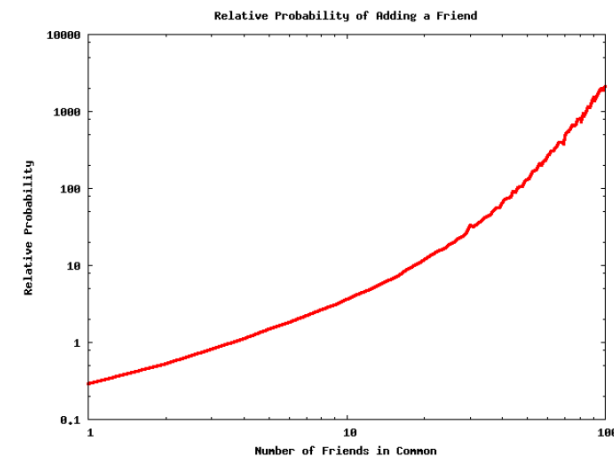
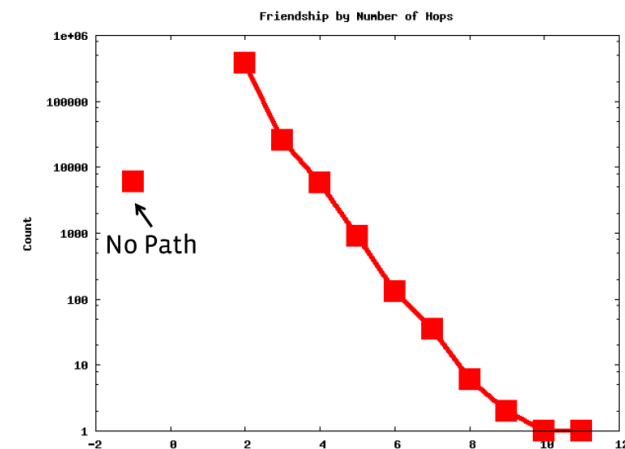
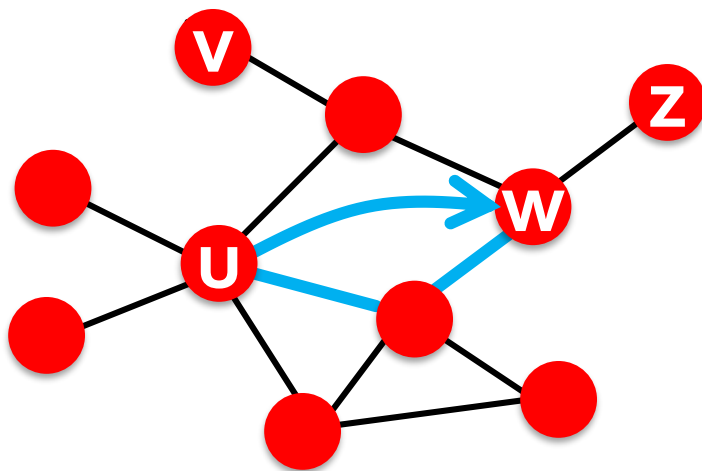
Supervised Link Prediction

■ Can we learn to predict new friends?

■ Facebook's People You May Know

■ Let's look at the FB data:

- 92% of new friendships on FB are friend-of-a-friend
- More mutual friends helps



Supervised Link Prediction

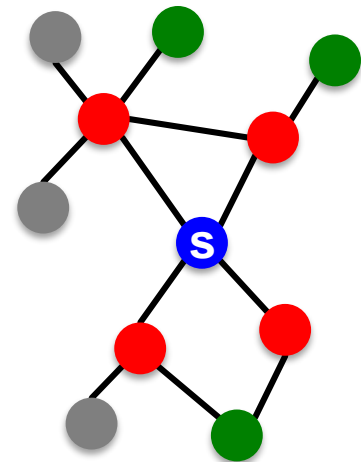
- **Goal:** Recommend a list of possible friends
- **Supervised machine learning setting:**

- **Labeled training examples:**

- For every user s have a list of others she will create links to $\{d_1 \dots d_k\}$ **in the future**
 - Use FB network from May 2012 and $\{d_1 \dots d_k\}$ are the new friendships you created since then
 - These are the “positive” training examples
- Use all other users as “negative” example

- **Task:**

- For a given node s , **score** nodes $\{d_1 \dots d_k\}$ **higher** than any other node in the network

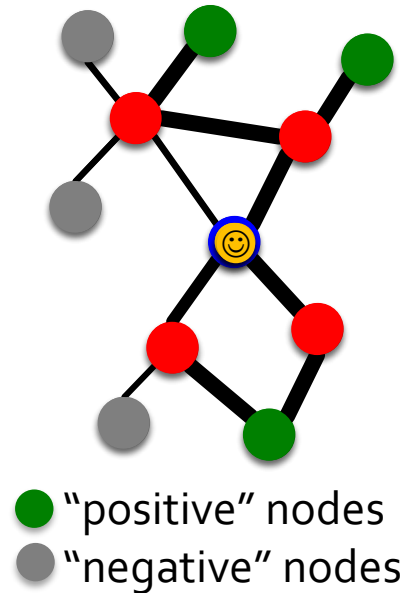


● “positive” nodes
● “negative” nodes

Green nodes
are the nodes
to which **s**
creates links in
the future

Supervised Link Prediction

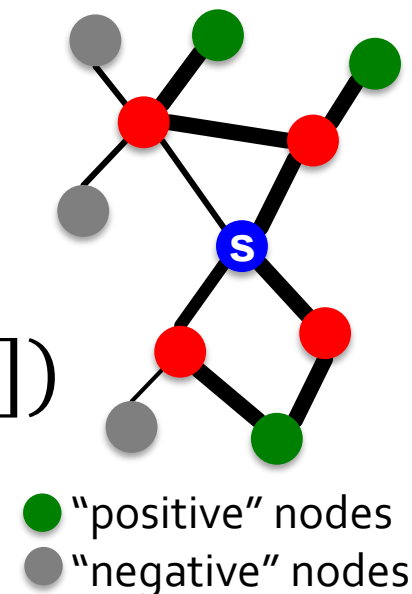
- How to combine node/edge features and the network structure?
 - Estimate **strength** of each friendship (u, v) using:
 - Profile of user u , profile of user v
 - Interaction history of users u and v
 - This creates a **weighted graph**
 - Do **Personalized PageRank from s** and measure the “**proximity**” (the visiting prob.) of any other node w from s
 - Sort nodes w by decreasing “**proximity**”



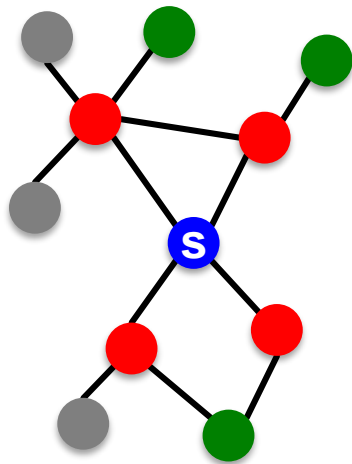
Supervised Random Walks

- Let s be the starting node
- Let $f_{\beta}(u, v)$ be a function that assigns **strength a_{uv} to edge (u, v)**

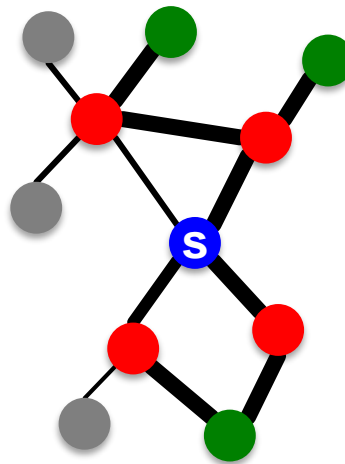
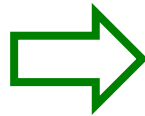
$$a_{uv} = f_{\beta}(u, v) = \exp(-\sum_i \beta_i \cdot x_{uv}[i])$$
 - x_{uv} is a feature vector of (u, v)
 - Features of node u
 - Features of node v
 - Features of edge (u, v)
 - **Note: β is the weight vector we will later estimate!**
- **Do Random Walk with Restarts** from s where transitions are according to edge strengths a_{uv}



SRW: Prediction



Network



Set edge
strengths
 $a_{uv} = f_{\beta}(u, v)$



Random Walk with Restarts on the weighted graph. Each node w has a PageRank proximity p_w



Sort nodes w by the decreasing PageRank score p_w



Recommend top k nodes with the highest proximity p_w to node s

- **How to estimate edge strengths?**
 - How to set parameters β of $f_{\beta}(u, v)$?
- **Idea:** Set β such that it (correctly) predicts the known future links

Personalized PageRank

- a_{uv} Strength of edge (u, v)
- Random walk transition matrix:

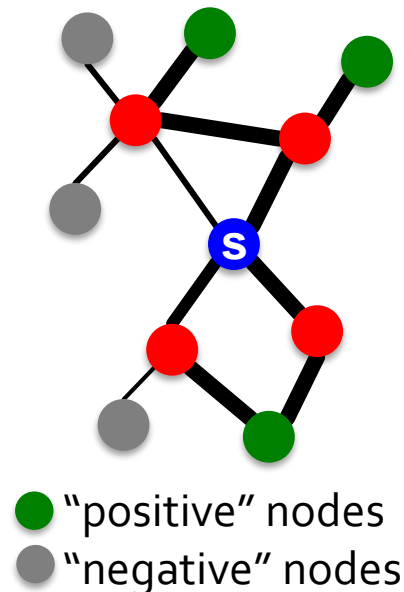
$$Q'_{uv} = \begin{cases} \frac{a_{uv}}{\sum_w a_{uw}} & \text{if } (u, v) \in E, \\ 0 & \text{otherwise} \end{cases}$$

- PageRank transition matrix:

$$Q_{ij} = (1 - \alpha)Q'_{ij} + \alpha \mathbf{1}(j = s)$$

- Where with prob. α we jump back to node s

- Compute PageRank vector: $p = p^T Q$
- Rank nodes w by decreasing p_w



The Optimization Problem

- **Positive** examples
 $D = \{d_1, \dots, d_k\}$
- **Negative** examples
 $L = \{\textit{other nodes}\}$
- **What do we want?**

$$\min_{\beta} F(\beta) = ||\beta||^2$$

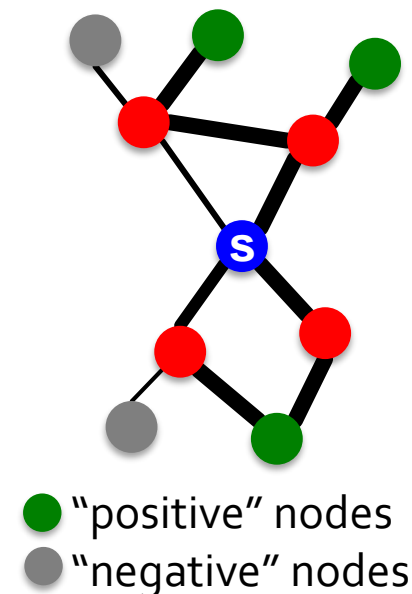
such that

$$\forall d \in D, l \in L : p_l < p_d$$

- **Note:**

- Exact solution to this problem may not exist
- So we make the constraints “soft” (i.e., optional)

We prefer small weights β to prevent overfitting



Every positive example has to have higher PageRank score than every negative example

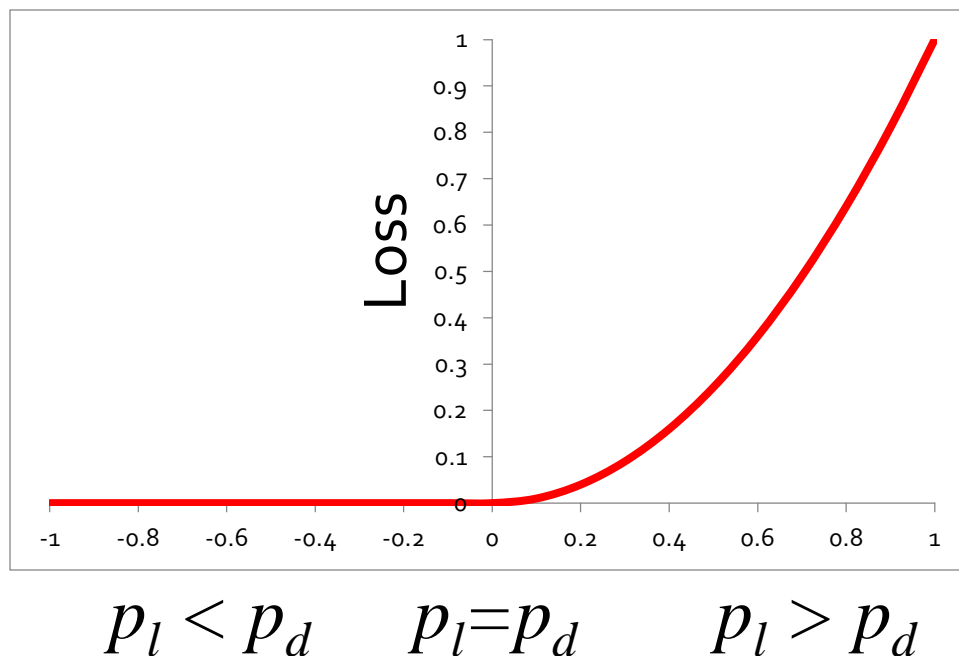
Making Constraints “Soft”

- **Want to minimize:**

$$\min_{\beta} F(\beta) = \sum_{d \in D, l \in L} h(p_l - p_d) + \lambda ||\beta||^2$$

- **Loss:** $h(x) = 0$ if $x < 0$, or x^2 else

Penalty for violating the constraint that $p_d > p_l$



Solving the Problem: Intuition

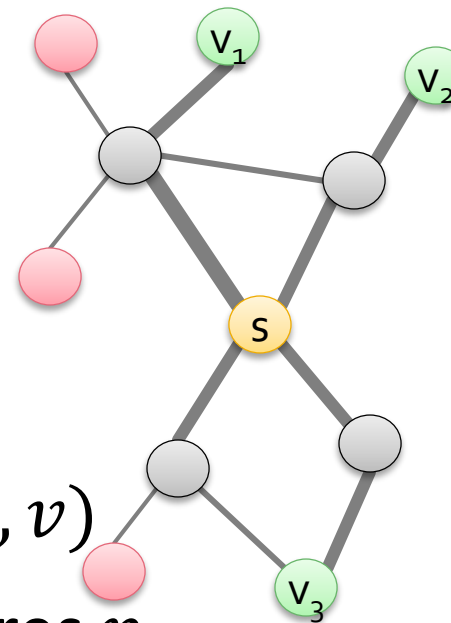
- **How to minimize F ?**

$$\min_{\beta} F(\beta) = \sum_{d \in D, l \in L} h(p_l - p_d) + \lambda ||\beta||^2$$

- **Both p_l and p_d depend on β**

- Given β assign edge weights $a_{uv} = f_{\beta}(u, v)$
- Using $Q = [a_{uv}]$ compute PageRank scores p_{β}
- Rank nodes by the decreasing score

- **Goal: Want to find β such that $p_l < p_d$**



Solving the Problem: Intuition

■ How to minimize $F(\beta)$?

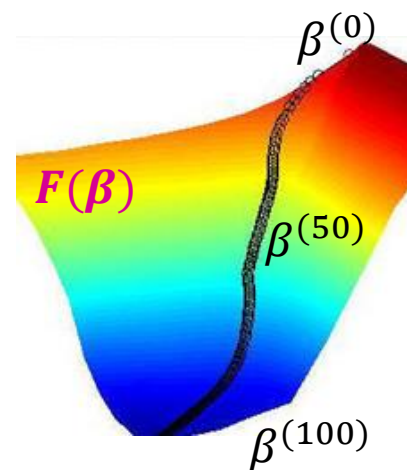
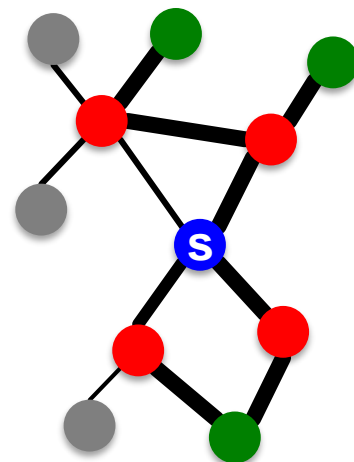
$$\min_{\beta} F(\beta) = \sum_{d \in D, l \in L} h(p_l - p_d) + \lambda \|\beta\|^2$$

■ Idea:

- Start with some random $\beta^{(0)}$
- Evaluate the derivative of $F(\beta)$ and do a small step in the opposite direction

$$\beta^{(t+1)} = \beta^{(t)} - \eta \frac{\partial F(\beta^{(t)})}{\partial \beta}$$

- Repeat until convergence



Gradient Descent

- What's the derivative $\frac{\partial F(\beta^{(t)})}{\partial \beta}$?

$$F(\beta) = \sum_{d \in D, l \in L} h(p_l - p_d) + \lambda \|\beta\|^2$$

$$\frac{\partial F(\beta)}{\partial \beta} = \sum_{l,d} \frac{\partial h(p_l - p_d)}{\partial \beta} + 2\lambda\beta$$

$$= \sum_{l,d} \frac{\partial h(p_l - p_d)}{\partial (p_l - p_d)} \left(\frac{\partial p_l}{\partial \beta} - \frac{\partial p_d}{\partial \beta} \right) + 2\lambda\beta$$

$$h(x) = \max\{x, 0\}^2$$

Easy!

- We know:

$$p = p^T Q \text{ that is } p_u = \sum_j p_j Q_{ju}$$

- So:

$$\frac{\partial p_u}{\partial \beta} = \sum_j Q_{ju} \frac{\partial p_j}{\partial \beta} + p_j \frac{\partial Q_{ju}}{\partial \beta}$$

Gradient Descent

■ **We just got:**
$$\frac{\partial p_u}{\partial \beta} = \sum_j Q_{ju} \frac{\partial p_j}{\partial \beta} + p_j \frac{\partial Q_{ju}}{\partial \beta}$$

■ **Few details:**

■ Computing $\partial Q_{ju}/\partial \beta$ is easy. **Remember:**
$$Q'_{uv} = \begin{cases} \frac{a_{uv}}{\sum_w a_{uw}} & \text{if } (u, v) \in E, \\ 0 & \text{otherwise} \end{cases}$$

■ We want $\frac{\partial p_j}{\partial \beta}$ but it appears on both sides of the equation. Notice the whole thing looks like a PageRank equation: $x = Q \cdot x + z$

$$a_{uv} = f_{\beta}(u, v) = \exp\left(-\sum_i \beta_i \cdot x_{uv}[i]\right)$$

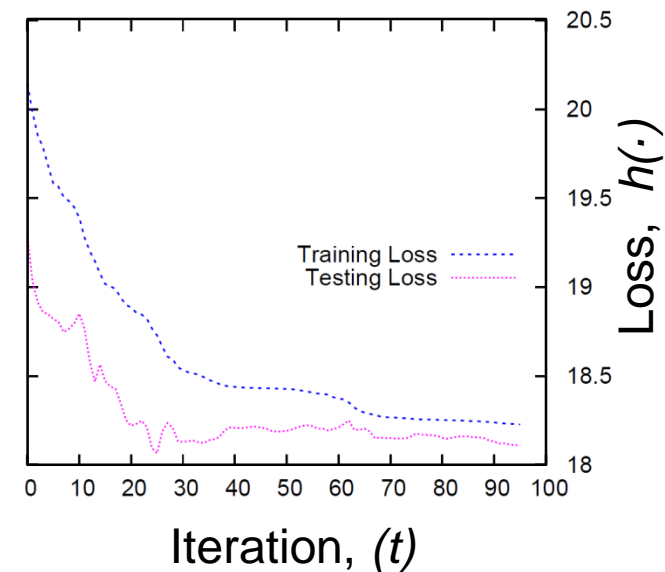
■ **As with PageRank we can use the power-iteration to solve it:**

■ Start with a random $\frac{\partial p}{\partial \beta}^{(0)}$

■ Then iterate:
$$\frac{\partial p}{\partial \beta}^{(t+1)} = Q \cdot \frac{\partial p}{\partial \beta}^{(t)} + \frac{\partial Q_{ju}}{\partial \beta} \cdot p$$

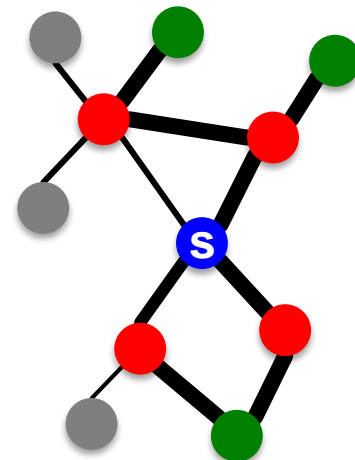
Optimizing $F(\beta)$

- To optimize $F(\beta)$, use gradient descent:
 - Pick a random starting point $\beta^{(0)}$
 - Using current $\beta^{(t)}$ compute edge strenghts and the transition matrix Q
 - Compute PageRank scores p
 - Compute the gradient with respect to weight vector $\beta^{(t)}$
 - Update $\beta^{(t+1)}$



Data: Facebook

- **Facebook Iceland network**
 - 174,000 nodes (55% of population)
 - Avg. degree 168
 - Avg. person added 26 friends/month
- **For every node s :**
 - **Positive examples:**
 - $D = \{ \text{new friendships } s \text{ created in Nov '09} \}$
 - **Negative examples:**
 - $L = \{ \text{other nodes } s \text{ did not create new links to} \}$
 - **Limit to friends of friends:**
 - On avg. there are 20,000 FoFs (maximum is 2 million)!



Experimental Setting

- **Node and Edge features for learning:**
 - **Node:** Age, Gender, Degree
 - **Edge:** Age of an edge, Communication, Profile visits, Co-tagged photos
- **Evaluation:**
 - **Precision at top 20**
 - We produce a list of 20 candidates
 - By taking top 20 nodes x with highest PageRank score p_x
 - Measure to what fraction of these nodes s actually links to

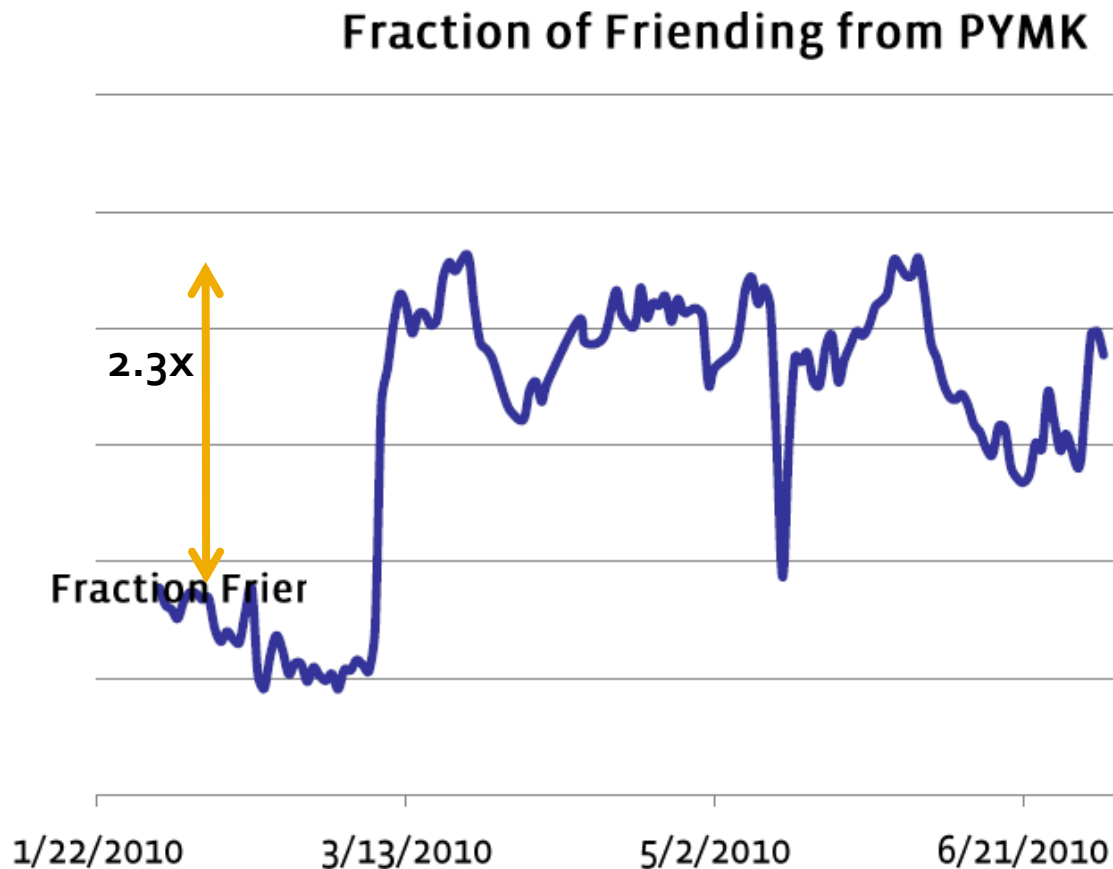
Results: Facebook Iceland

- Facebook: Predict future friends
 - Adamic-Adar already works great
 - Supervised Random Walks (SRW) gives slight improvement

Learning Method	Prec@Top20
Random Walk with Restart	6.80
Adamic-Adar	7.35
Common Friends	7.35
Degree	3.25
SRW: one edge type	6.87
SRW: multiple edge types	7.57

Results: Facebook

- 2.3x improvement over previous FB-PYMK (People You May Know)



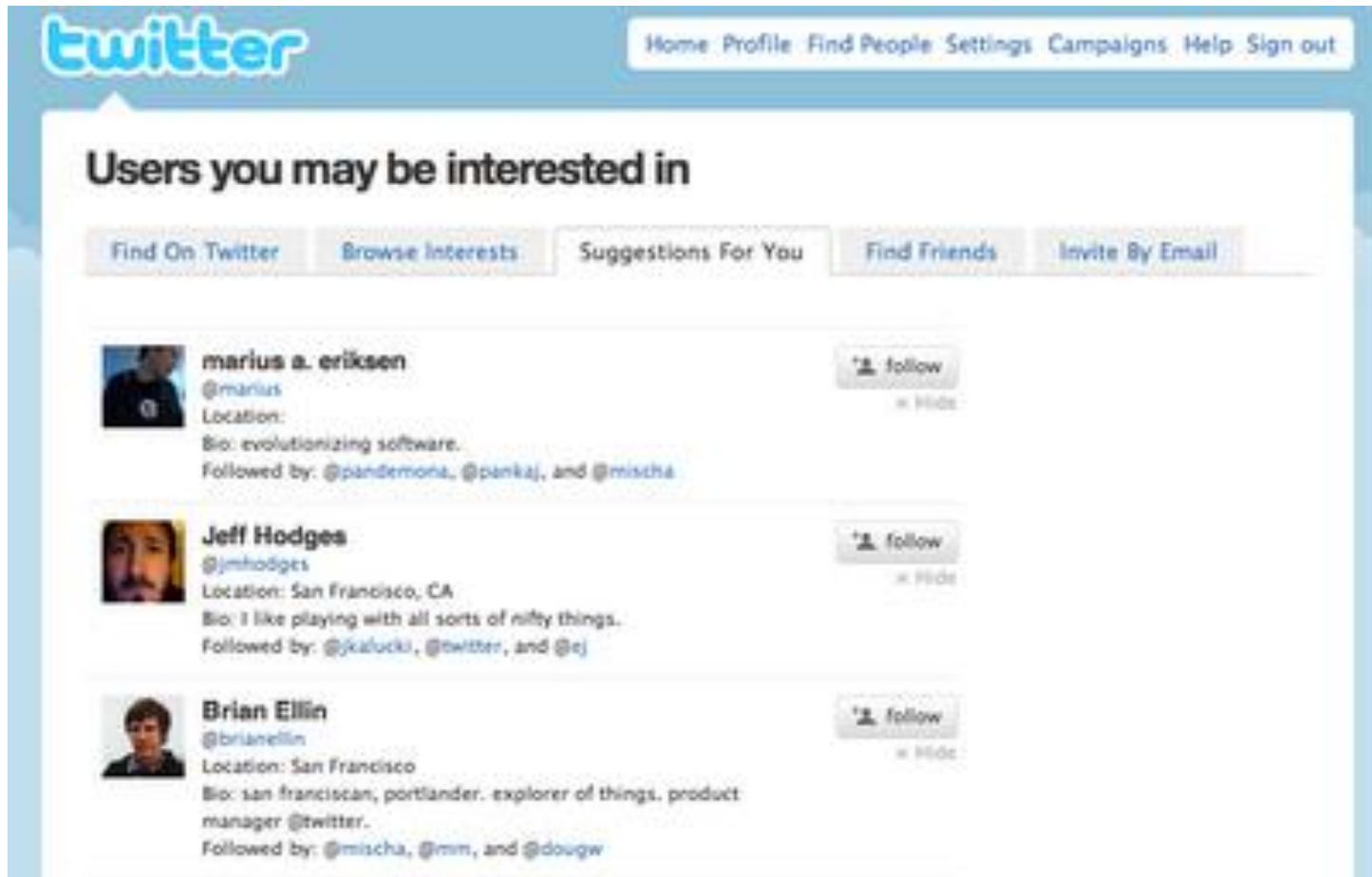
Results: Co-Authorship

- **Arxiv Hep-Ph collaboration network:**
 - Poor performance of unsupervised methods
 - SRW gives a boost of 25%!

Learning Method	Prec@Top20
Random Walk with Restart	3.41
Adamic-Adar	3.13
Common Friends	3.11
Degree	3.05
SRW: one edge type	4.24
SRW: multiple edge types	4.25

Wtf: The Who to Follow Service at Twitter

Introduction



Semantic differences between “interested in” and “similar to”

WtF ("Who to Follow")

- **WtF ("Who to Follow")**: the Twitter user recommendation service
 - help existing and new users to **discover connections** to sustain and grow
 - used for **search relevance, content discovery, promoted products**, etc.
 - **Twitter Data**:
 - 200 million users
 - 400 million tweets every day (as of early 2013)
 - <http://www.internetlivestats.com/twitter-statistics/>

The Twitter Graph

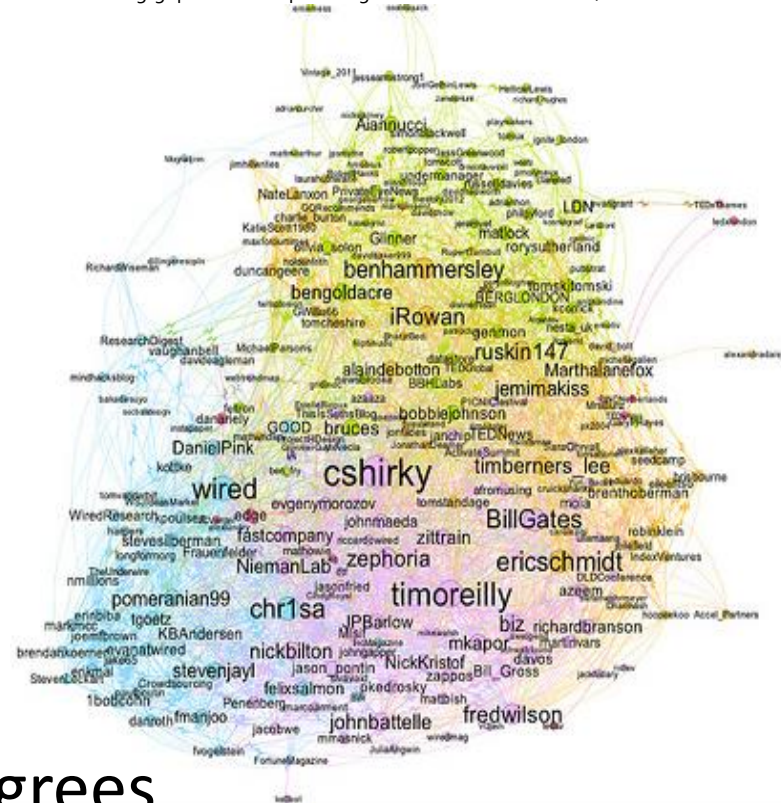
<http://blog.ouseful.info/2011/07/07/visualising-twitter-friend-connections-using-gephi-an-example-using-wired-uk-friends-network/>

■ Graph

- **Node:** user
- **(Directed) Edge:** follows

■ Graph Statistics (Aug'12)

- Over 20 billion edges
- Power law of in- and out-degrees
- Over 1000 with more than 1 million followers
- 25 users with more than 10 million followers



Algorithms: Circle of Trust

- **Circle of Trust**

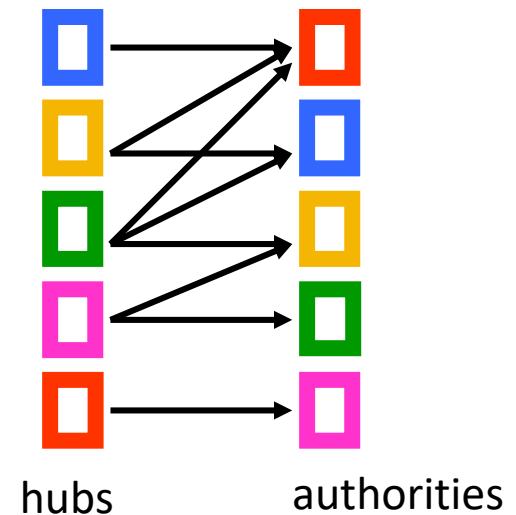
- Based on an **egocentric random walk** (similar to **personalized (rooted) PageRank**)
- Computed in an **online fashion** (from scratch each time) given **a set of parameters**
 - # of random walk steps
 - reset probability
 - pruning settings to discard low probability vertices
 - parameters to control sampling of outgoing edges at vertices with large out-degrees

Algorithms

- Directed edge
 - **Asymmetric** nature of the **follow relationship**
 - Friendships in other social networks such as Facebook or LinkedIn are symmetric/**reciprocal**
 - Similar to the **user-item recommendations** problem where the “item” is also a user

Algorithms: SALSA

- **SALSA** (Stochastic Approach for Link-Structure Analysis)
 - a variation of HITS
- **HITS**
 - Intuition:
 - Good **hubs** point to good **authorities**
 - Good **auth.** are pointed by good **hubs**
 - Recurs. comput. of **hub score**
 - Recurs. comput. of **auth. score**



$$h_i = \sum_{j:i \rightarrow j} a_j$$

$$a_i = \sum_{j:j \rightarrow i} h_j$$

Algorithms: SALSA

- Random walks to rank **hubs** and **authorities**
 - Two different random walks (Markov chains): **a chain of hubs** and **a chain of authorities**
 - Each walk traverses nodes only in one side by traversing two links in each step **$h \rightarrow a \rightarrow h$** , **$a \rightarrow h \rightarrow a$**

Transition matrices of each chain: **H** and **A**

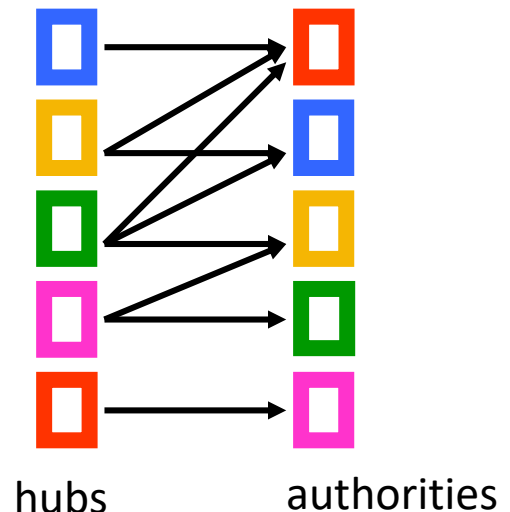
W: the adjacency of the directed graph

W_r: divide each entry by the sum of its row

W_c: divide each entry by the sum of its column

$$\mathbf{H} = \mathbf{W}_r \mathbf{W}_c^T$$

$$\mathbf{A} = \mathbf{W}_c^T \mathbf{W}_r$$

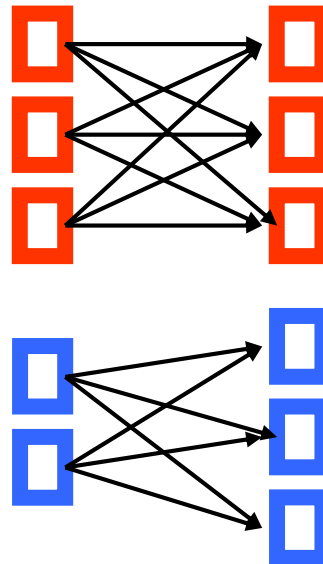


Algorithms: SALSA

- Reduces to the problem of HITS with **tightly knit communities**
 - TKC effect
- Better for single-topic communities
- More efficient implementation

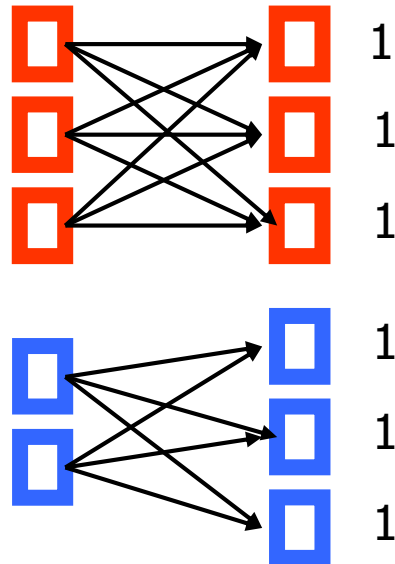
HITS and the TKC effect

- The HITS algorithm favors the most **dense community** of hubs and authorities
 - Tightly Knit Community (TKC) effect



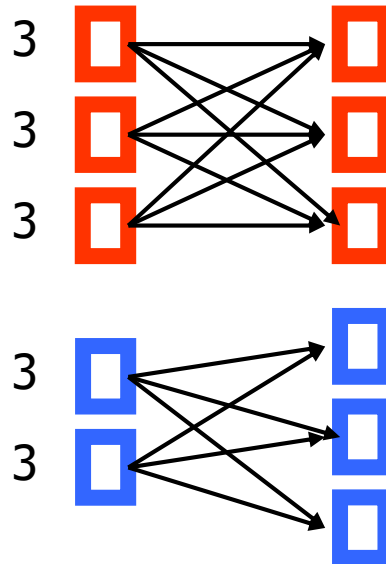
HITS and the TKC effect

- The HITS algorithm favors the most **dense community** of hubs and authorities
 - Tightly Knit Community (TKC) effect



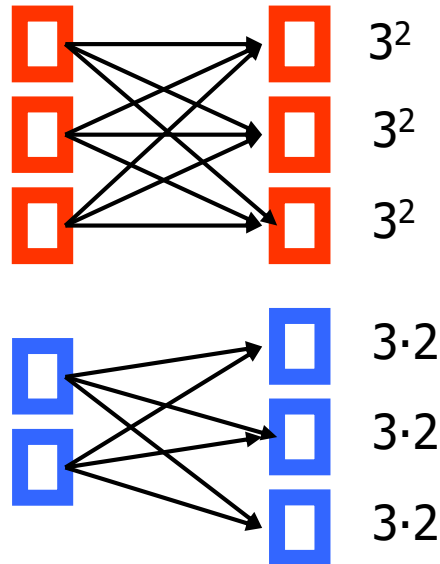
HITS and the TKC effect

- The HITS algorithm favors the most **dense community** of hubs and authorities
 - Tightly Knit Community (TKC) effect



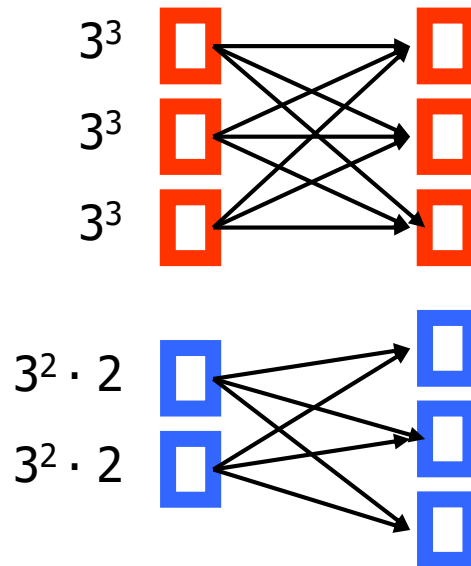
HITS and the TKC effect

- The HITS algorithm favors the most **dense community** of hubs and authorities
 - Tightly Knit Community (TKC) effect



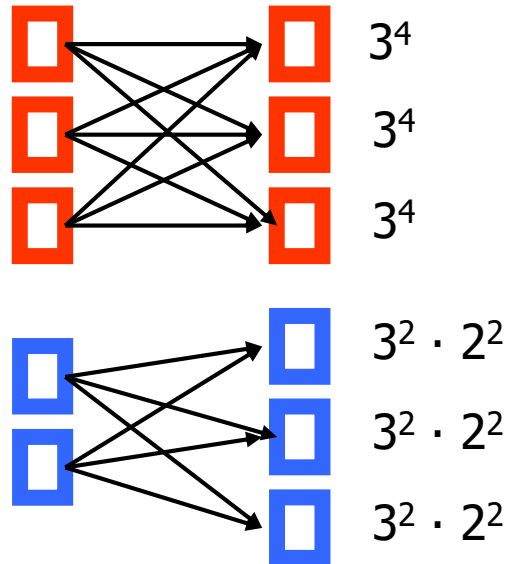
HITS and the TKC effect

- The HITS algorithm favors the most **dense community** of hubs and authorities
 - Tightly Knit Community (TKC) effect



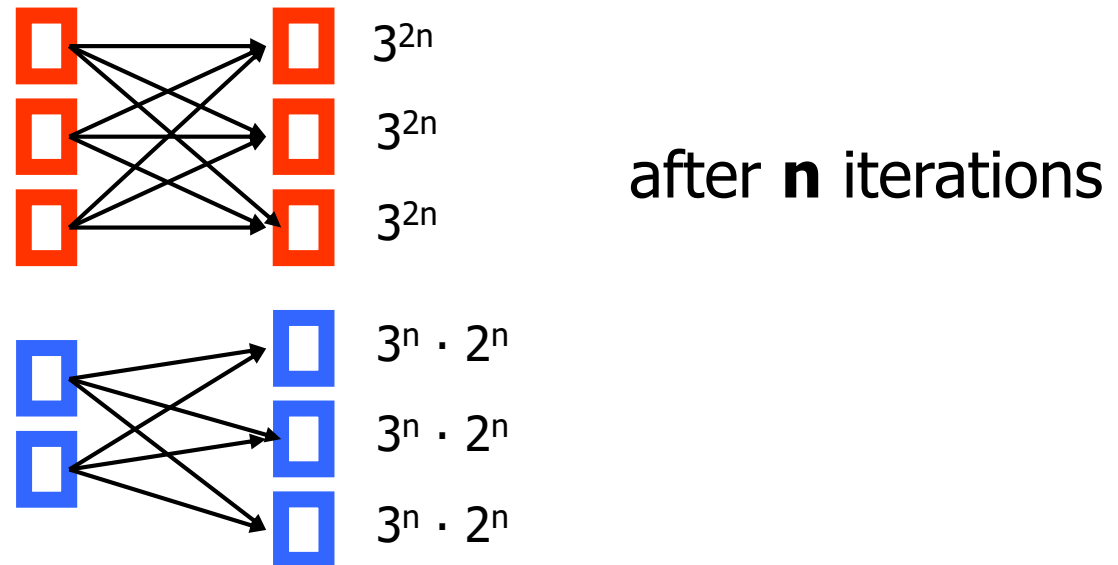
HITS and the TKC effect

- The HITS algorithm favors the most **dense community** of hubs and authorities
 - Tightly Knit Community (TKC) effect



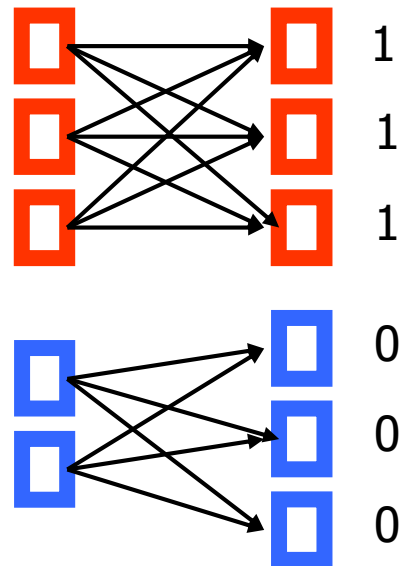
HITS and the TKC effect

- The HITS algorithm favors the most **dense community** of hubs and authorities
 - Tightly Knit Community (TKC) effect



HITS and the TKC effect

- The HITS algorithm favors the most **dense community** of hubs and authorities
 - Tightly Knit Community (TKC) effect



after normalization
with the max
element as $n \rightarrow \infty$

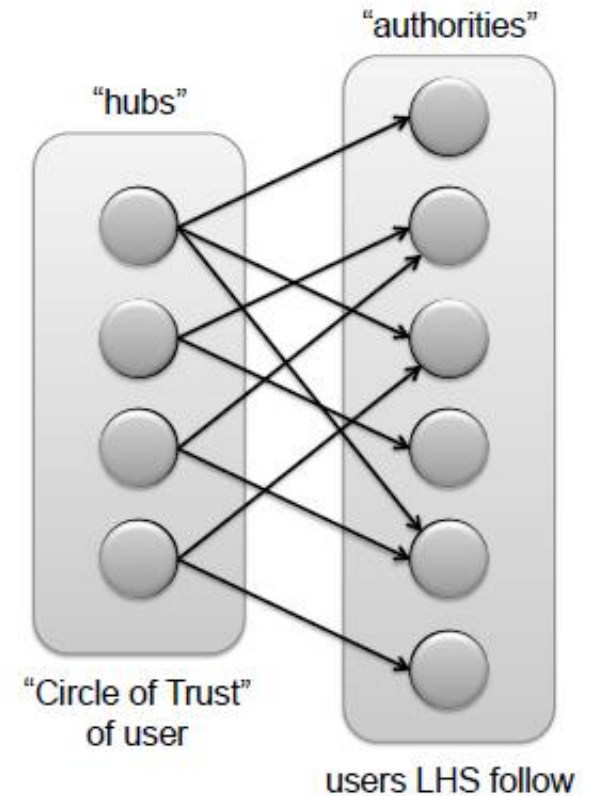
Algorithms: SALSA

■ Hubs:

- 500 top-ranked nodes from a user's **circle of trust**
- user similarity (based on homophily, also useful)

■ Authorities:

- users that the hubs follow
- “interested in” user recommendations



Algorithms: SALSA

- **SALSA**'s recursive nature
 - Two users are similar if they **follow** the same (or similar) users (**LHS**)
 - A user **u** is likely to follow those who are **followed by** users that are similar to **u** (**RHS**)
- The random walk ensures **fair distribution** of scores in both directions
- Similar users are selected from the **circle of trust** of a user (via **personalized PageRank**)

Evaluation

- Approaches
 - Offline experiments on retrospective data
 - Online **A/B testing** on live traffic
- Various parameters may interfere:
 - How the results are rendered
 - Platform (mobile, etc.)
 - New vs old users

Extensions

- Add **metadata to vertices** (e.g., user profile information) and **edges** (e.g., edge weights, timestamp, etc.)
- Consider **interaction graphs** (e.g., graphs defined in terms of retweets, favorites, replies, etc.)

Extensions

- **Two phase algorithm**
 - **1st - Candidate generation:** produce a list of promising recommendations for each user, using any algorithm
 - **2nd - Rescoring:** apply a machine-learned model to the candidates, binary classification problem (logistic regression)
- **Evaluation**
 - **1st Phase:** recall + diversity
 - **2nd Phase:** precision + maintain diversity