# Link Prediction in Networks

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



- 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 coauthorship
- Overcoming the data-sparsity problem *in recommender* systems using collaborative filtering

# Link Prediction in Networks

#### The link prediction task:

Given G[t<sub>0</sub>, t'<sub>0</sub>] a graph on edges up to time t'<sub>0</sub>, output a ranked list L of links (not in G[t<sub>0</sub>, t'<sub>0</sub>]) that are predicted to appear in G[t<sub>1</sub>, t'<sub>1</sub>]



 $\begin{array}{l} G[t_0,t_0']\\ G[t_1,t_1'] \end{array}$ 

#### Evaluation:

- n = /E<sub>new</sub>/: # new edges that appear during the test period [t<sub>1</sub>, t'<sub>1</sub>]
- 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 (*deg.* ≥ 3)
- Predict top n pairs as new links
- See which of these links actually appear in G[t<sub>1</sub>, t'<sub>1</sub>]



of node *x* 

# **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:  $|\Gamma(x)| \cdot |\Gamma(y)|$   $\Gamma(x) \dots$  neighbors
  - 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(·)
  - 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**

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

## **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 **r(x)** be the set of neighbors of **x** in **G**<sub>old</sub>

#### 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, A<sub>x,y</sub><sup>2</sup>: #paths of length 2

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

#### Jaccard coefficient

• The probability that both **x** and **y** have a feature for a randomly selected feature that either **x** or **y** has  $\operatorname{score}(x, y) = \frac{|\Gamma(x) \cap \Gamma(y)|}{|\Gamma(x) \cup \Gamma(y)|}$ 

# **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)

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

#### Preferential attachment

 Based on the premise that the probability that a new edge has node x as its endpoint is proportional to |Γ(x)|, i.e., nodes like to form ties with 'popular' nodes

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

# **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<sub>β</sub> measure (unweighted, weighted)
    - Random walk-based
      - hitting time
      - commute time
      - Rooted PageRank
      - SimRank

#### **Shortest Path Based**

For  $\mathbf{x}, \mathbf{y} \in \mathbf{V} \times \mathbf{V} - \mathbf{E}_{old}$ ,

# 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**

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

- Sum over all paths of length *l*
- O<B<1: a parameter of the predictor, exponentially damped to count short paths more heavily

#### **Ensemble of All Paths**

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

$$\sum_{l=1}^{\infty} \beta^l \cdot |\mathsf{paths}_{xy}^{\langle l \rangle}| = \beta A_{xy} + \beta^2 (A^2)_{xy} + \beta^3 (A^3)_{xy} + \cdots$$

- Unweighted version: path<sub>x,y</sub>(1) = 1, if x and y have collaborated, 0 otherwise
- Weighted version: path<sub>x,y</sub>(1) = #times x and y have collaborated

## **Random Walk Based**

Consider a random walk on **G**<sub>old</sub> that starts at **x** and iteratively moves to a neighbor of **x** chosen uniformly at random from **Γ(x)** 

- Hitting H<sub>x,y</sub> (from x to y): the expected number of steps it takes for the random walk starting at x to reach y
   score(x, y) = H<sub>x,y</sub>
- Commute Time C<sub>x,y</sub> (from x to y): the expected number of steps to travel from x to y and from y to x

**score(x, y)** = 
$$-(H_{x,y} + H_{y,x})$$

Not symmetric, can be shown  $h_{vu} = \Theta(n^2)$   $h_{uv} = \Theta(n^3)$   $h_{uv} = \Theta(n^3)$   $h_{uv} = \Theta(n^3)$ 

Jure Leskovec, Stanford CS224W: Social and Information Network Analysis, http://cs224w.stanford.edu

## **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<sub>old</sub> that starts at x and has a probability α of returning to x at each step
- Rooted PageRank
  - Starts from x
  - with probability (1 a) moves to a random neighbor
  - with probability *a* returns to x

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

**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

$$\label{eq:similarity} \begin{split} \mathsf{similarity}(x,y) := \gamma \cdot \frac{\sum_{a \in \Gamma(x)} \sum_{b \in \Gamma(y)} \mathsf{similarity}(a,b)}{|\Gamma(x)| \cdot |\Gamma(y)|} \end{split}$$

 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 **M** with a lower rank matrix  $\mathbf{M}_{\mathbf{k}}$ 

- Method
  - Apply SVD (singular value decomposition)
  - Obtain the rank-k matrix that best approximates M

# **Singular Value Decomposition**

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

## **Unseen bigrams**

- Unseen bigrams: pairs of words that co-occur in a test corpus, but not in the corresponding training corpus
- Not just score(x, y) but score(z, y) for nodes z that are similar to x --- S<sub>x</sub><sup>(δ)</sup>: the δ nodes most related to x

$$\operatorname{score}_{unweighted}^{*}(x,y) := \left| \{z : z \in \Gamma(y) \cap S_{x}^{\langle \delta \rangle} \} \right|$$
$$\operatorname{score}_{weighted}^{*}(x,y) := \sum_{z \in \Gamma(y) \cap S_{x}^{\langle \delta \rangle}} \operatorname{score}(x,z)$$

# Clustering

- Compute score(x, y) for all edges in E<sub>old</sub>
- 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<sub>p</sub>* of pairs in V×V–E<sub>old</sub> in decreasing order of confidence
  - focus on Core network, (d > 3)

 $E*_{new} = E_{new} \cap (Core \times Core) = |E*_{new}|$ 

- Evaluation method: Size of intersection of
  - the first n edge predictions from L<sub>p</sub> that are in Core
     × Core, and
  - the actual set E\*<sub>new</sub>

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:

$$|E_{new}|$$

$$\binom{|\mathsf{Core}|}{2} - |E_{old}|$$

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



3/2/2017

## **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
lpha=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
lpha = 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
eta=0.005	13.4	54.4	30.1	24.0	51.9
eta=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
eta=0.005	16.8	41.4	37.4	24.1	49.4
eta=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	0.475%	0.147%	0.341%	0.207%	0.153%	
graph distance (all distan	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:	rank = 1024	15.2	53.8	29.3	34.8	49.8
Inner product	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:	rank = 1024	8.2	16.6	6.6	18.5	21.6
Matrix entry	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:	rank = 1024	11.4	27.2	30.1	27.0	32.0
Katz ( $\beta = 0.005$ )	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	common neighbors, $\delta = 8$	13.5	36.7	30.1	15.6	46.9
(weighted)	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	common neighbors, $\delta = 8$	14.1	40.2	27.9	22.2	39.4
(unweighted)	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:	ho=0.10	7.4	37.3	46.9	32.9	37.8
Katz ( $\beta_1 = 0.001, \beta_2 = 0.2$	1) $\rho = 0.15$	12.0	46.1	46.9	21.0	44.0
	ho = 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**

	Adamic/Adar	Katz clustering	common neighbors	hitting time	Jaccard's coefficient	weighted Katz	low-rank inner product	rooted Pagerank	SimRank	unseen bigrams
Adamic/Adar	1150	638	520	193	442	1011	905	528	372	486
Katz clustering		1150	411	182	285	630	623	347	245	389
 common neighbors			1150	135	506	494	467	305	332	489
hitting time				1150	87	191	192	247	130	156
Jaccard's coefficient					1150	414	382	504	845	458
weighted Katz						1150	1013	488	344	474
low-rank inner product							1150	453	320	448
rooted Pagerank								1150	678	461
SimRank									1150	423
unseen bigrams										1150

How similar are the predictions made by the different methods? # common predictions

		Adamic/Adar	Katz clustering	common neighbors	hitting time	Jaccard's coefficient	weighted Katz	low-rank inner product	rooted Pagerank	SimRank	unseen bigrams	
	Adamic/Adar	92	65	53	22	43	87	72	44	36	49	ĺ
	Katz clustering		78	41	20	29	66	60	31	22	37	
	common neighbors			69	13	43	52	43	27	26	40	
	hitting time				40	8	22	19	17	9	15	
-	Jaccard's coefficient					71	41	32	39	51	43	
	weighted Katz						92	75	44	32	51	
	low-rank inner product							79	39	26	46	
	rooted Pagerank								69	48	39	
	SimRank									66	34	
	unseen bigrams										68	

# # correct predictions

Jure Leskovec, Stanford CS224W: Social and Information Network Analysis, http://cs224w.stanford.edu

- 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<sub>1</sub> ... d<sub>k</sub>} in the future
      - Use FB network from May 2012 and {d<sub>1</sub> ... d<sub>k</sub>} 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<sub>1</sub> ... d<sub>k</sub>}
       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"



"positive" nodes

"negative" nodes

# **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<sub>uv</sub>

## **SRW: Prediction**



#### [WSDM '11]

# Personalized PageRank

*a<sub>uv</sub>* .... 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.  $\alpha$  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<sub>1</sub>,...,d<sub>k</sub>}

  Negative examples

  L = {other nodes}

  What do we want?

  min F(β) = ||β||<sup>2</sup>
  - such that
  - $\forall d \in D, l \in L : p_l < p_d$

#### Note:

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

We prefer small

overfitting

weights  $\beta$  to prevent

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



## **Making Constraints "Soft"**

#### Want to minimize:



[WSDM '11]

# **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 $\beta$

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

#### • Goal: Want to find $\beta$ such that $p_l < p_d$

[WSDM '11]

S

V-

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

Jure Leskovec, Stanford CS224W: Social and Information Network Analysis, http://cs224w.stanford.edu

Repeat until convergence





[WSDM '11]

### **Gradient Descent**

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

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

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

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

#### We know:

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

So:  

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

Details!

#### 3/2/2017

# **Gradient Descent**

# Few details:

- We just got:  $\frac{\partial p_u}{\partial \beta} = \sum_j Q_{ju} \frac{\partial p_j}{\partial \beta} + p_j \frac{\partial Q_{ju}}{\partial \beta}$ 
  - 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:  $x = Q \cdot x + z$

$$a_{uv} = f_{\beta}(u, v)$$
$$= \exp\left(-\sum_{i} \beta_{i} \cdot \mathbf{x}_{uv}[i]\right)$$

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



# Optimizing $F(\beta)$

#### • To optimize $F(\beta)$ , use gradient descent:

- Pick a random starting point  $\beta^{(0)}$
- Using current β<sup>(t)</sup> compute edge strenghts and the transition matrix Q
- Compute PageRank scores p
- Compute the gradient with respect to weight vector  $\beta^{(t)}$
- Update  $\beta^{(t+1)}$



**Details**!

### 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 = { new friendships s created in Nov '09 }
- Negative examples:
  - L = { other nodes s did not create new links to }
- Limit to friends of friends:
  - On avg. there are 20,000 FoFs (maximum is 2 million)!



[WSDM '11]

# **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<sub>x</sub>
  - 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)

Fraction of Friending from PYMK



# **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

sers you r	nay be intere	ested in		
nd On Twitter	Browse Interests	Suggestions For You	Find Friends	Invite By Email
marius a Ornarius Location: Bio: evoluti Followed by	. eriksen onizing software. r: @pandemona, @pankaj.	and @mischa	12. follow	
Jeff Hod Opiniodoti Location: Si Bio: 1 like pi Followed by	ges an Francisco, CA laying with all sorts of nife r: @jkalucki, @twitter, and	r things.	*1. follow a Hide	
Brian Ell	in		'1. follow	

#### 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**

### Graph

- Node: user
- (Directed) Edge: follows

#### Graph Statistics (Aug'12)

Over 20 billion edges

Alainucci Alainucci Natelanoo Protective 

http://blog.ouseful.info/2011/07/07/visualising-twitter-friend-connections-

using-gephi-an-example-using-wireduk-friends-network/

- 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

hubs autho

authorities



# **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→a→h, a→h→a

Transition matrices of each chain: H and A W: the adjacency of the directed graph W<sub>r</sub>: divide each entry by the sum of its row W<sub>c</sub>: divide each entry by the sum of its column



# **Algorithms: SALSA**

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

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



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



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



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



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

- 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.)

#### Two phase algorithm

- 1<sup>st</sup> Candidate generation: produce a list of promising recommendations for each user, using any algorithm
- 2<sup>nd</sup> Rescoring: apply a machine-learned model to the candidates, binary classification problem (logistic regression)

#### Evaluation

- 1<sup>st</sup> Phase: recall + diversity
- 2<sup>nd</sup> Phase: precision + maintain diversity