

*MINING RESEARCH TOPIC-RELATED
INFLUENCE BETWEEN ACADEMIA
AND INDUSTRY*

Dan He

PKDD 2011, Athens, Greece

danhe@cs.ucla.edu



BACKGROUND

- Mining influence in networks, especially social networks has attracted tremendous attentions.



BACKGROUND

- Mining influence in networks, especially social networks has attracted tremendous attentions.
- Topic-level influence between academia and industry is useful in many cases:
 - Students whose career goal is to be a researcher in a company's research lab often seek advisor with high influence in industry.
 - Funding agencies may choose to grant certain awards to researchers who work closely with industry companies, on certain topics.
 - Companies may want to collaborate with academia researchers who have tight industry connection.



BACKGROUND

- Mining influence in networks, especially social networks has attracted tremendous attentions.
- Topic-level influence between academia and industry is useful in many cases:
 - Students whose career goal is to be a researcher in a company's research lab often seek advisor with high influence in industry.
 - Funding agencies may choose to grant certain awards to researchers who work closely with industry companies.
 - Companies may want to collaborate with academia researchers who have tight industry connection.
- To our knowledge, our work is the first one on mining influences between academia and industry.



PROPERTIES OF CO-AUTHORSHIP NETWORK

- Nodes: researchers.
- Edges: co-authorship between the pair of researchers.
- Edge weights: number of publications co-authored by the pair of researchers.
- Unlike social networks, there are no temporal “actions”: the pair of authors publish at the same time.



OUTLINE

- Topic-specific influence model in co-authorship network
- Three models for influence between academia and industry
 - Simple Additive model
 - Weighted Additive model
 - Clustering-based Additive model
- Evaluations of the three models
- Conclusions



OUTLINE

- Topic-specific influence model in co-authorship network
- Three models for influence between academia and industry
 - Simple Additive model
 - Weighted Additive model
 - Clustering-based Additive model
- Evaluations of the three models
- Conclusions



MODELS FOR PAIRWISE INFLUENCES IN CO-AUTHORSHIP NETWORK

○ Naïve Method:

- Based on the number of co-authored publications

- $$\text{Influence}(A \rightarrow B) = \frac{\# \text{co-publication}(A, B)}{\# \text{publication}(B)}$$

- Problem of the method: Considers A, B as independent to other authors, which is usually not true.



MODELS FOR PAIRWISE INFLUENCES IN CO-AUTHORSHIP NETWORK

- TAP Method (Tang et al., KDD 2009):
 - Based on topical affinity propagation (TAP)
 - A topical factor graph is built on top of the co-authorship network
 - A likelihood function for the graph is maximized by TAP.
 - Influences are computed via TAP.
 - Models influences among multiple researchers on specific topics.



OUTLINE

- Topic-specific influence model in co-authorship network
- Three models for influence between academia and industry
 - Simple Additive model
 - Weighted Additive model
 - Clustering-based Additive model
- Evaluations of the three models
- Conclusions



PREREQUISITES

- Assume we already obtain the topic-specific influences between researchers in the co-authorship network, via the TAP method.
- The influences are directed.
- The influences between two nodes are usually not symmetric.

$$\textit{Influence}(A \rightarrow B) \neq \textit{Influence}(B \rightarrow A)$$



SIMPLE ADDITIVE MODEL

- **Basic Idea:** Simply sum the influences on topic t from an academia researcher r to all the researchers in a company C

The diagram illustrates the Simple Additive Model equation: $I_t(r, C) = \sum_{i=1}^n I_t(r, C_i)$. Annotations include: 'company' pointing to C ; 'i-th researcher of the company' pointing to C_i ; 'Influence from research r to company C ' pointing to the left side of the equation; 'Influence from researcher r to the i-th researcher' pointing to the right side of the equation; 'Academia researcher' pointing to r ; and 'topic' pointing to the t in the equation.

company

i-th researcher of the company

Influence from research r to company C

$I_t(r, C) = \sum_{i=1}^n I_t(r, C_i)$

Influence from researcher r to the i-th researcher

Academia researcher

topic



SIMPLE ADDITIVE MODEL

- **Problem of the model:** not all the researchers are of equal importance to the company.
 - A manager should have higher weight than his/her group members do
 - A senior researcher should have higher weights than a junior researcher does



WEIGHTED ADDITIVE MODEL

- **Basic Idea:** Each industry researcher is weighted according to their internal influences in their company.

i-th researcher of the company

Internal influence of C_i

Weight of C_i → $W_t(C_i) = \frac{\sum_{j=1}^n I_t(C_i, C_j)}{\sum_{j=1, k=1}^{j=n, k=n} I_t(C_k, C_j)}$

$$I_t(r, C) = \sum_{i=1}^n I_t(r, C_i) \times W_t(C_i)$$

Total influence in the company



WEIGHTED ADDITIVE MODEL

- **Basic Idea:** Each industry research is weighted according to their internal influences in their company.
- **Problem of the model:** there might be researchers often publish together.
 - These researchers will have high influence to each other.
 - They may have high internal influences, leading to biased high weights.



WEIGHTED ADDITIVE MODEL

- Basic Idea: Each industry research is weighted according to their internal influences in their company.
- Problem of the model: there might be researchers often publish together.
 - These researchers will have high influence to each other.
 - They may have high internal influence, leading to biased high weights.
- **Question:** How to normalize the weights of these researchers to address the bias?



CLUSTERING-BASED ADDITIVE MODEL

- **Basic Idea:** Put the researchers that collaborate a lot into a cluster.
 - Consider each cluster as a *super researcher*.
 - Average the influences in a cluster.
 - Conduct the weighted additive model on the cluster level instead of on the researcher level.
 - The weights of more correlated researchers, or researcher who collaborate often, can be adjusted appropriately.



CLUSTERING THE RESEARCHERS

- We do not know a reasonable number of clusters for clustering algorithms such as K-means, EM, spectral clustering etc.
- Modularity based clustering do not require a number of clusters, but it ignores the distances between the nodes.



CLUSTERING THE RESEARCHERS

- We do not know a reasonable number of clusters for clustering algorithms such as K-means, EM, spectral clustering etc.
- Modularity based clustering do not require a number of clusters, but it ignores the distances between the nodes.
- We propose a clustering algorithm which is based on distances between the nodes and does not require pre-knowledge on the number of clusters.



CLUSTERING THE RESEARCHERS

- Similarity of two researchers: Jaccard's coefficient on the number of publications.

$$\frac{\# \textit{co} - \textit{publication}}{\# \textit{total} \quad \textit{publication}}$$

- We use number of publications instead of influence since the similarity needs to be symmetric.



CLUSTERING THE RESEARCHERS

- The nodes in the same cluster are fully-connected, namely each cluster is a clique.
- Given a threshold t , an edge connects two nodes only when the similarity of the two nodes is no less than t
- We can generate a set of clusters, by searching for cliques in the graph. The nodes in the same clique are of similarity no less than t to each other.
- We want to select an appropriate t automatically.



CLUSTERING THE RESEARCHERS

- We propose an objective function for t

The set of clusters for threshold t

$$F(t) = \left\| \sum_{k=1}^{|cluster^t|} \sum_{i,j \in cluster_k^t} (1 - s_{i,j}) - |cluster^t| \right\|$$

Similarity between node i,j

The k -th cluster

- Find a t that minimizes $F(t)$
- Motivation of the objective function:
 - Minimize the sum of distances between the nodes in the clusters.
 - Also take into consideration the number of clusters



CLUSTERING THE RESEARCHERS

- Find t to minimize $F(t)$:
 - Compute pairwise similarity for all the nodes in the co-authorship network.
 - We prove at least one of the pairwise similarities is able to minimize $F(t)$.
 - We try all pairwise similarities and find the one that minimizes $F(t)$.
 - If we assume every researcher has only constant number of collaborators h , this step takes $O(nh)$ complexity, where n is the total number of nodes. Since h is a constant and usually speaking $h \ll n$, $O(nh) \approx O(n)$.



WEIGHTED ADDITIVE MODEL ON CLUSTERS

- Compute the weight of the researchers in the cluster according to the normal weighted additive model:

$$W_t(L_i) = \frac{\sum_{j=1}^n I_t(L_i, L_j)}{\sum_{j=1, k=1}^{j=n, k=n} I_t(L_k, L_j)}$$

- Compute the influence from a researcher r to a cluster L based on the weights of the researchers in the cluster:

$$I_t(r, L) = \sum_{i=1}^n I_t(r, L_i) \times W_t(L_i)$$



WEIGHTED ADDITIVE MODEL ON CLUSTERS

- Compute the influence between pair of clusters based on the weights of researchers in the clusters

$$I_t(L, K) = \sum_{i=1, j=1}^{i=|L|, j=|K|} I_t(L_i, K_j) \times W_t(L_i) \times W_t(K_j)$$

- Compute the weights of the clusters by considering each cluster as a “virtual researcher”

$$W_t(C_L) = \frac{\sum_{K \in \text{cluster}(C)} I_t(C_L, C_K)}{\sum_{L \in \text{cluster}(C), K \in \text{cluster}(C)} I_t(C_L, C_K)}$$



WEIGHTED ADDITIVE MODEL ON CLUSTERS

- Compute the influence of an academia researcher r to a company C according to the normal weighted additive model, by considering each cluster as a “virtual researcher”.

$$I_t(r, C) = \sum_{L \in \text{cluster}(C)}^n I_t(r, C_L) \times W_t(C_L)$$



CLUSTERING-BASED ADDITIVE MODEL

- As we will show later in the experiments, when there are bias towards the researchers who often publish together, the simple additive model and the weighted additive model are not able to address the bias.
- The clustering-based additive model is able to better address the bias.



OUTLINE

- Topic-specific influence model in co-authorship network
- Three models for influence between academia and industry
 - Simple Additive model
 - Weighted Additive model
 - Clustering-based Additive model
- Evaluations of the three models
- Conclusions



EXPERIMENTS ON REAL DATA

- Experiments Set up:
 - The same data set as used in the work of Tang et al, KDD 2009, with 8 different topics.
 - Researchers are affiliated with either companies or universities
 - One researcher only has one affiliation.

Topics	#Researchers	#Companies	#Universities
<i>Data Mining</i>	679	33	99
<i>Machine Learning</i>	976	48	97
<i>Database System</i>	1127	66	116
<i>Information Retrieval</i>	657	49	87
<i>Web Services</i>	400	27	48
<i>Semantic Web</i>	671	38	35
<i>Bayesian Network</i>	554	24	45
<i>Web Mining</i>	348	25	47



EXPERIMENTS ON REAL DATA

○ Experiments Set up:

- The researchers with missing affiliation are excluded.
- Influences on the 8 topics among the researchers are computed via the TAP method.

Topics	#Researchers	#Companies	#Universities
<i>Data Mining</i>	679	33	99
<i>Machine Learning</i>	976	48	97
<i>Database System</i>	1127	66	116
<i>Information Retrieval</i>	657	49	87
<i>Web Services</i>	400	27	48
<i>Semantic Web</i>	671	38	35
<i>Bayesian Network</i>	554	24	45
<i>Web Mining</i>	348	25	47



SAMPLE RESULTS

- The top-5 most influential researchers to Microsoft and IBM from the three models on the topic “Data Mining”

Data Mining (Microsoft)			Data Mining (IBM)		
<i>simple additive</i>	<i>weighted additive</i>	<i>clustering-based additive</i>	<i>simple additive</i>	<i>weighted additive</i>	<i>clustering-based additive</i>
Jiawei Han	Huan Liu	Clark Glymour	Jiawei Han	Jiawei Han	Jiawei Han
Huan Liu	Clark Glymour	Huan Liu	Philip S. Yu	Philip S. Yu	Philip S. Yu
Xifeng Yan	Michail Vlachos	Michail Vlachos	Michail Vlachos	Michail Vlachos	Michail Vlachos
Clark Glymour	Xuanhui Wang	Padhraic Smyth	Tao Tao	Tao Tao	Tao Tao
Philip S. Yu	Padhraic Smyth	Bing Liu	Ricardo Vilalta	Ricardo Vilalta	Ricardo Vilalta



SAMPLE RESULTS

- The top-5 most influential researchers to Microsoft and IBM from the three models on the topic “Data Mining”
- The ranks may not be accurate due to missing affiliation information.

Data Mining (Microsoft)			Data Mining (IBM)		
<i>simple additive</i>	<i>weighted additive</i>	<i>clustering-based additive</i>	<i>simple additive</i>	<i>weighted additive</i>	<i>clustering-based additive</i>
Jiawei Han	Huan Liu	Clark Glymour	Jiawei Han	Jiawei Han	Jiawei Han
Huan Liu	Clark Glymour	Huan Liu	Philip S. Yu	Philip S. Yu	Philip S. Yu
Xifeng Yan	Michail Vlachos	Michail Vlachos	Michail Vlachos	Michail Vlachos	Michail Vlachos
Clark Glymour	Xuanhui Wang	Padhraic Smyth	Tao Tao	Tao Tao	Tao Tao
Philip S. Yu	Padhraic Smyth	Bing Liu	Ricardo Vilalta	Ricardo Vilalta	Ricardo Vilalta



SAMPLE RESULTS

- The top-5 most influential researchers to Microsoft and IBM from the three models on the topic “Machine Learning”

Machine Learning (Microsoft)		
<i>simple additive</i>	<i>weighted additive</i>	<i>clustering-based additive</i>
Brendan J. Frey	Aaron Hertzmann	Aaron Hertzmann
Aaron Hertzmann	Michael I. Jordan	Michael I. Jordan
Andrew McCallum	Andrew McCallum	Andrew McCallum
Michael I. Jordan	Brendan J. Frey	William T. Freeman
William T. Freeman	William T. Freeman	Yoav Freund

Machine Learning(IBM)		
<i>simple additive</i>	<i>weighted additive</i>	<i>clustering-based additive</i>
Inderjit S. Dhillon	Inderjit S. Dhillon	Inderjit S. Dhillon
Adam R. Klivans	Manfred K. Warmuth	Manfred K. Warmuth
Nader H. Bshouty	Roni Khardon	Roni Khardon
Joydeep Ghosh	Geoffrey J. Gordon	Geoffrey J. Gordon
Manfred K. Warmuth	Gerald Tesauro	Gerald Tesauro



EVALUATION ON THE CLUSTERING ALGORITHM

- Our clustering algorithm found clusters of researchers who collaborated quite often.
- We show 4 pairs of researchers who co-authored many papers. According to the total number of their publications, we guess one of them is relatively senior and the other is relatively junior.
- We validate the titles of the senior researchers and it is true that all of them are manager/principal

Senior	Junior	Co-authored
47	23	11
64	16	9
62	18	11
94	28	22



EXPERIMENTS ON SIMULATED DATA

- Due to lack of ground-truth and missing affiliation, we can not effectively validate our models on real data.



EXPERIMENTS ON SIMULATED DATA

- We validate our models on simulated data:
 - Generate a company with 200 researchers
 - Randomly select 20 of them as manager
 - Assign the remaining researchers to one of the 20 groups randomly
 - Generate 400 researchers in academia
 - Ground-truth of two types of most influential academia researchers to the company:
 - *Influence many researcher*: Academia researcher who influences many researchers of the company
 - *Influence important researchers*: Academia researcher who influences important researchers, such as manager, of the company



EXPERIMENTS ON SIMULATED DATA

- We select 10 *influence many* researchers and 10 *influence important* researchers
- We also select 20 *influence same group* researchers, who influence researchers in the same group of the company.
- Rank the academia researchers according to their influences to the company by the three models.
- We expect the top-20 most influential researchers are the 10 *influence many* researchers + the 10 *influence important* researchers



EXPERIMENTS ON SIMULATED DATA

- We conduct two sets of experiments:
 - The researchers in the same group have LOW influence to each other (no bias)
 - The researchers in the same group have HIGH influence to each other (with bias)



EXPERIMENTS ON SIMULATED DATA

- We conduct two sets of experiments:
 - The researchers in the same group have LOW influence to each other (no bias)
 - The researchers in the same group have HIGH influence to each other (with bias)
- What we observe:
 - Simple Additive model tends to assign high ranks to the *influence many* researchers



EXPERIMENTS ON SIMULATED DATA

- We conduct two sets of experiments:
 - The researchers in the same group have LOW influence to each other (no bias)
 - The researchers in the same group have HIGH influence to each other (with bias)
- What we observe:
 - Simple Additive model tends to assign high ranks to the *influence many* researchers
 - Weighted Additive model tends to assign high ranks to the *influence important* researchers



EXPERIMENTS ON SIMULATED DATA

- What we observe:
 - Clustering-based additive model
 - Is similar to the weighted additive model, if researchers in the same group have low influence to each other (no bias)



EXPERIMENTS ON SIMULATED DATA

- What we observe:
 - Clustering-based additive model
 - Is similar to the weighted additive model, if researchers in the same group have low influence to each other (no bias)
 - Tends to assign higher ranks to the *influence important* researchers and lower ranks to the *influence same group* researchers, if researchers in the same group have high influence to each other (with bias)



EXPERIMENTS ON SIMULATED DATA

- What we observe:
 - Clustering-based additive model
 - Is similar to the weighted additive model, if researchers in the same group have low influence to each other (no bias)
 - Tends to assign higher ranks to the *influence important* researchers and lower ranks to the *influence same group* researchers, if researchers in the same group have high influence to each other (with bias)
 - When the researchers in the same group have high influence to each other, clustering-based additive model is able to better distinguish the *influence same group* researchers from the *influence important* and the *influence many* researchers.



OUTLINE

- Topic-specific influence model in co-authorship network
- Three models for influence between academia and industry
 - Simple Additive model
 - Weighted Additive model
 - Clustering-based Additive model
- Evaluations of the three models
- **Conclusions**



TAKE HOME MESSAGE

- We addressed the problem of mining research topic-specific influence between academia and industry.
- We proposed three models to learn the influence:
 - Simple additive model
 - Weighted additive model
 - Clustering-based additive model
- The influence from industry to academia can be obtained using the same models but influence with reverse directions.



THANKS

- Questions?

