

Syracuse University

**SURFACE**

---

Dissertations - ALL

SURFACE

---

January 2015

## Link Prediction in Dynamic Weighted and Directed Social Network using Supervised Learning

Ricky Laishram  
*Syracuse University*

Follow this and additional works at: <https://surface.syr.edu/etd>



Part of the [Engineering Commons](#)

---

### Recommended Citation

Laishram, Ricky, "Link Prediction in Dynamic Weighted and Directed Social Network using Supervised Learning" (2015). *Dissertations - ALL*. 355.

<https://surface.syr.edu/etd/355>

This Thesis is brought to you for free and open access by the SURFACE at SURFACE. It has been accepted for inclusion in Dissertations - ALL by an authorized administrator of SURFACE. For more information, please contact [surface@syr.edu](mailto:surface@syr.edu).

## **Abstract**

Link Prediction is an area of great interest in social network analysis. Previous works in the area of link prediction have only focused on networks where the links once created cannot be removed. In many real world social networks, the links should be assigned strengths; for example, the strength of a link should decrease over time, if there are no interactions between the two nodes for a long time and increase if the two nodes interact often. In this thesis we modify existing methods of link prediction to apply to weighted and directed networks. The features, developed in previous works for unweighted and undirected networks, are extended to apply to networks whose links have weight and direction, and algorithms are developed to calculate them efficiently. These network features are used to train an SVM classifier to predict which nodes will be connected by a link and which links will be broken in the future. The results obtained using Twitter @-mention network demonstrate that the method developed in this thesis is very effective.

# **Link Prediction in Dynamic Weighted and Directed Social Network using Supervised Learning**

by

**Ricky Laishram**

Submitted in partial fulfillment of the requirements for the degree of Master of  
Science in Computer Science

**Syracuse University**

**August 2015**

# **Copyright**

Copyright © Ricky Laishram, 2015

## **Acknowledgments**

I am very grateful to Dr. Chilukuri K. Mohan who introduced me to the area of social network analysis and encouraged me to work on the current research. Without his help and encouragements, I might not have even considered doing this research.

I am also extremely grateful to Dr. Kishan Mehrotra for guiding me in my research and providing valuable feedback while writing this thesis. Without his guidance and help this thesis would not have been possible.

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Previous Works</b>	<b>3</b>
2.1	Unsupervised Link Prediction . . . . .	3
2.2	Supervised Link Prediction . . . . .	5
2.3	Supervised Link Prediction in Sparse Network . . . . .	7
<b>3</b>	<b>Network Model with Weighted and Directed Links</b>	<b>8</b>
<b>4</b>	<b>Modification of Feature Scores for Weighted and Directed Network</b>	<b>10</b>
4.1	Common Neighbor Score . . . . .	10
4.2	Preferential Attachment Score . . . . .	11
4.3	Adamic-Adar Index . . . . .	12
4.4	Katz Score . . . . .	14
4.5	Shortest Path Score . . . . .	16
4.6	Rooted PageRank . . . . .	17
<b>5</b>	<b>Incremental Method for Calculating Feature Scores</b>	<b>19</b>
5.1	Common Neighbor Score . . . . .	19
5.1.1	Effect of decay of link weight . . . . .	20

5.1.2	Effect of removal of link . . . . .	20
5.1.3	Effect of interaction between nodes . . . . .	22
5.2	Preferential Attachment Score . . . . .	23
5.2.1	Effect of decay of link weight . . . . .	23
5.2.2	Effect of removal of link . . . . .	24
5.2.3	Effect of interaction between nodes . . . . .	25
5.3	Adamic Adar Index . . . . .	27
5.3.1	Effect of decay of link weight . . . . .	27
5.3.2	Effect of removal of links . . . . .	29
5.3.3	Effect of interaction between nodes . . . . .	30
5.4	Katz Score . . . . .	33
5.4.1	Effect of decay of link weight . . . . .	33
5.4.2	Effect of removal of link . . . . .	33
5.4.3	Effect of interaction between nodes . . . . .	36
5.5	Shortest Path Score . . . . .	37
5.5.1	Effect of decay of link weight . . . . .	38
5.5.2	Effect of interaction between nodes . . . . .	38
5.5.3	Effect of removal of link . . . . .	39
5.6	Rooted PageRank . . . . .	40

5.6.1	Effect of decay of link weight . . . . .	41
5.6.2	Effect of interactions between nodes and removal of links . . . . .	41
<b>6</b>	<b>Link Prediction as a Classification Problem</b>	<b>42</b>
6.1	Link Prediction Method . . . . .	42
6.2	Handling Class Imbalance . . . . .	43
<b>7</b>	<b>Data</b>	<b>45</b>
7.1	Data Source . . . . .	45
7.2	Data Preparation . . . . .	45
<b>8</b>	<b>Link Prediction Results</b>	<b>47</b>
<b>9</b>	<b>Analysis</b>	<b>54</b>
<b>10</b>	<b>Conclusion</b>	<b>54</b>
	<b>Appendices</b>	<b>55</b>
<b>A</b>	<b>Algorithms for Calculating Feature Scores</b>	<b>55</b>
A.1	Common Neighbor Score . . . . .	55
A.2	Adamic-Adar Index . . . . .	55
A.3	Preferential Attachment Score . . . . .	55

A.4	Katz Score . . . . .	59
A.5	Shortest Path Score . . . . .	61
A.6	Rooted PageRank . . . . .	61
<b>B</b>	<b>Run Time Analysis</b>	<b>63</b>
B.1	Common Neighbor Score . . . . .	63
B.2	Adamic-Adar Index . . . . .	63
B.3	Preferential Attachment Score . . . . .	64
B.4	Katz Score . . . . .	64
B.5	Shortest Path Score . . . . .	65
B.6	Rooted PageRank . . . . .	65
	<b>References</b>	<b>67</b>

# 1 Introduction

A social network is a structure containing a set of actors and relationships; actors are represented as nodes, and relationships are represented as edges. For example, the edges may represent communications between the nodes, spread of diseases from one node to another, etc. Two basic changes can happen in a social network - addition or removal of nodes, and addition or removal of edges.

In this thesis, I propose a method to predict addition or removal of edges. This particular problem is called the Link Prediction Problem in Social Network. David Liben-Nowell and Jon Kleinberg define the Link Prediction Problem as follows: "to what extent can the evolution of a social network be modeled using features intrinsic to the network itself?"[2].

In all the previous works addressing this problem, the network model is such that a link once created is not destroyed. However, in many real world networks, the links are not permanent. For example in a Twitter @-mention network, if two users have not communicated for a long time, the link between them should be removed. To address this problem, I consider the following:

1. A network model where the link weights change with time. In this model each new interaction between nodes strengthens the link between them, whereas lack of interaction decreases the strength of the link. An appropriate threshold is set such that if the weight of a link is below the threshold, the link is removed.
2. Link prediction in directed and weighted networks. I modify features, used in undirected and unweighted networks, to apply to weighted and directed networks.
3. Efficient algorithms to incrementally update the network feature scores are also developed.
4. The problem of class imbalance is handled by using the ensemble SVM method.

This thesis is organized as follows:

1. The problem of edge deletion is addressed in Chapter 3.

2. Feature modifications to address the weighted directed networks is addressed in Chapter 4.
3. Efficient algorithms to compute the feature scores is described in Chapter 5.
4. In Chapter 6, link prediction problem is formulated as a classification problem, and solved using SVM where class imbalance is also addressed.

In the next chapter, I examine the previous works done by other researchers and examine the shortcomings of the existing methods.

## 2 Previous Works

Many researchers have worked on the problem of link prediction in social networks in unsupervised as well as supervised context. In this chapter we look at some of the previous works done in the area of link prediction and their shortcomings.

### 2.1 Unsupervised Link Prediction

Liben-Nowell et al. [1] examined if links in a network can be predicted by considering only features intrinsic to the network. They observed a collaboration network at time  $[t_0, t_1, t_2, \dots T]$ . In this network the set of nodes  $V$  remains fixed but edges are added over time. The network over time are  $\langle V, E_0 \rangle, \langle V, E_1 \rangle, \langle V, E_2 \rangle \dots \langle V, E_T \rangle$  such that  $E_0 \subseteq E_1 \subseteq E_2 \dots \subseteq E_T$ . Out of these networks, they considered two networks  $G_1 = \langle V, E_t \rangle$  and  $G_2 = \langle V, E_T - E_t \rangle$  for time  $t$  such that  $0 < t < T$ , and let  $N = |E_T - E_t|$ . In other words,  $G_2$  contains only the links that were formed after time  $t$ . They used  $G_1$  as the training network and  $G_2$  is assumed to be unknown. Their aim was to use  $G_1$  to predict the edges in  $G_2$ .

They considered the following intrinsic features: common neighbors, preferential attachment count, Adamic-Adar Index, Katz Score, Hitting Time, Rooted PageRank and Commute Time for this evaluation.

Their method is described as below:

1. For all possible node pairs  $u$  and  $v$  such that  $u, v \in V \wedge (u, v) \notin E_1$ , the feature score  $s_i(u, v)$  is calculated for all the features,  $i = 1, 2, 3, \dots 7$ . Let  $S_i = \{s_i(u, v) : \forall (u, v) \in V \wedge (u, v) \notin E_1\}$ .
2. For  $i = 1, 2, \dots 7$ , the elements of  $S_i$  are sorted in descending order of magnitude and the top  $N$  pairs of nodes  $(u, v)$  with highest  $s_i(u, v)$  are selected. Let  $S_i^*$  represent this list of nodes. A link is predicted between  $(u, v)$  for all  $(u, v) \in S_i^*$ .

The performance  $P_i$  for a feature  $i$  is considered as the number of correctly predicted links. A link  $(u, v)$  in  $S_i^*$  is said to be correctly predicted if it is also in the set  $E_T - E_t$ .

$$P_i = |\{(u, v) : (u, v) \in S_i^* \wedge (u, v) \in (E_T - E_t)\}| \quad (1)$$

Liben-Nowell et al. compared the performances  $P_i$  against a random prediction. To generate randomly predicted set of  $N$  links in  $E_T - E_t$ , a set of  $N$  node pairs  $(u, v)$  is selected such that  $u, v \in V \wedge (u, v) \notin E_t$ . Let us denote this list by  $R$ , and  $Q$  be the number of correctly identified edges.

$$Q = |\{(u, v) : (u, v) \in R \wedge (u, v) \in (E_T - E_t)\}| \quad (2)$$

Liben-Nowell et al. applied their method to the collaboration network obtained from the ArXiv dataset. In this network,  $|V| = 1253$  and  $|E_T - E_t| = 1150$ .

Feature	$P_i/Q$
Adamic-Adar	54.8
Katz Score	54.8
Common Neighbors	41.1
Hitting Time	23.8
Rooted PageRank	42.3
Commute Time	15.5

Table 1:  $P_i/Q$  for different features

Liben-Nowell et al. found that the ratio  $P_i/Q$  is in the range  $[15, 55]$  (Table 1) for all the features. This demonstrates that the performance of the predictions based on network intrinsic features is significantly higher than a random predictor. From this observation, they concluded that a network contains intrinsic information to predict links.

The method used by Liben-Nowell et al. has two drawbacks that make it difficult to apply in real world networks.

1. In real world networks, we do not know the number of links that will be created in the future.

So calculation of  $N$  will not be possible.

2. Even though their method outperforms random prediction by a huge margin, the ratio of correctly predicted links to the actual number of new links ( $P_i/N$ ) is very low. In their work, Adamic-Adar Index and Katz Score were the best performing features, but the associated predictions are only 8% correct i.e.  $P_i/N \leq 0.08$ . The  $P_i/N$  ratios for all the features are shown in Table 2.

Feature	$P_i/N$
Adamic-Adar	0.080
Katz Score	0.080
Common Neighbors	0.060
Hitting Time	0.034
Jaccard's Coefficient	0.062
Rooted PageRank	0.060

Table 2:  $P_i/N$  for different features

## 2.2 Supervised Link Prediction

Al Hasan et al. [2] explored the use of supervised learning methods (Decision Tree, Naive Bayes, SVM, Multi-Layer Perceptron) to predict links in the co-authorship networks obtained from the BIOBASE and DBLP datasets.

As in Section 2.1, Al Hasan et al. considered two networks  $G_1$  and  $G_2$ . The network  $G_1$  is the training network and  $G_2$  is the network for which links are to be predicted. They transform the link prediction problem into a two-class classification problem by assigning class labels to all pairs of nodes as:

$$class(u, v|G_t) = \begin{cases} 1, & \text{if } (u, v) \in E_t & \text{(positive class)} \\ 0, & \text{if } (u, v) \notin E_t & \text{(negative class)} \end{cases} \quad (3)$$

Using  $G_1$ , they calculate the feature scores associated with each pair  $(u, v)$ . Let  $F_1$  denote the

feature scores for all n-features from  $G_1$ ,

$$F_1 = \{(s_1(u, v|G_1), s_2(u, v|G_1), \dots, s_n(u, v|G_1)) : u, v \in V\} \quad (4)$$

Here  $s_i(u, v|G_t)$  is a feature score between nodes  $u$  and  $v$  in network  $G_t$  for a feature  $i$ . They considered features that are both intrinsic to the network (sum of neighbors count, sum of secondary neighbors count, clustering index, shortest path) and extrinsic to the network (based on papers published - sum of keywords, sum of classification code, sum of papers, keyword match count).

Let  $C_1$  be the set of all the class labels for network  $G_1$ . Then  $(F_1, C_1)$  forms the training data. Likewise testing data is calculated from  $G_2$  for all pairs of nodes that do not have a link between them in  $G_1$ , i.e.

$$F_2 = \{(s_1(u, v|G_2), s_2(u, v|G_2), \dots, s_n(u, v|G_2)) : u, v \in V \wedge (u, v) \notin E_t\} \quad (5)$$

The classifiers are trained using the training data, and used to predict the class labels for the testing data  $F_2$ .

Al Hasan et al. achieved very good results with their method. All the supervised learning algorithms they used have accuracy of more than 85%. However, there are some shortcomings with their method.

1. In a real world network, we do not know the network  $G_2$ . So calculation of the testing data  $F_2$  (in equation 5) will not be possible.
2. Their method uses network extrinsic features, e.g., derived from similarity of the papers published. In other applications, it may not be possible to use extrinsic features.

## 2.3 Supervised Link Prediction in Sparse Network

The network that Al Hasan et al. [2] used in their work has positive and negative class of approximately equal size. However, many real world networks are sparse - the size of the positive class is much smaller compared to that that of the negative class. Lichtenwaltert et al. [3] improved upon the works of Al Hasan et al. as follows:

1. In their method, they do not use any extrinsic features. They used only network intrinsic features: common neighbors count, Katz Score, Jaccard's Coefficient, PropFlow and Preferential Attachment.
2. They overcame the problem of class imbalance by oversampling from the positive class and undersampling from the negative class to generate the training data. They used SVM to classify the test data into positive or negative class, and compared the results to those obtained using unsupervised method (Section 2.1).

Lichtenwaltert et al. [3] obtained very good results with their method ( $AUC > 70\%$ ) on sparse networks. This demonstrates two important ideas:

1. The network extrinsic features, like the ones used by Al Hasan et al. [2] are not always necessary in a supervised learning method of link prediction.
2. Supervised link prediction performs very well in sparse networks after adjusting for class imbalance.

### 3 Network Model with Weighted and Directed Links

In this chapter, we describe a network model where the links are assigned weights, representing its strength, and direction. Consider a network which has a link between nodes  $u$  and  $v$  at time  $t$ . If there is an interaction between  $u$  and  $v$  at time  $(t + 1)$ , the strength of the link between  $u$  and  $v$  should increase. If there is an interaction between two nodes, the weight of the link between them is increased. Otherwise, the weight of the link between them is decreased by a factor  $\delta$ . A weight threshold  $\theta$  is defined such that a link is considered broken if its weight is lower than  $\theta$ .

For example, in a Twitter @-mention network, an interaction is a tweet sent by one node to another. Each link in the network has a strength associated with it. Depending upon the strength of a link, from node  $u$  to  $v$ , the link could be destroyed assuming there are no more interactions between  $u$  and  $v$  [4, 5] for a long time.

Suppose that  $G_t = \langle V, E_t \rangle$  is the network at time  $t$ . If  $I_{t,t+1}$  is the set of node pairs that interact with each other in the time interval  $[t, t + 1]$ , the set of links at time  $t + 1$  is generated using the method described below.

1. A set of links is generated:

$$T_{t+1} = E_t \cup \{(u, v) : (u, v) \in I_{t,t+1} \wedge (u, v) \notin E_t\} \quad (6)$$

2. The weight of each link in  $T_{t+1}$  is calculated:

$$w(u, v|t + 1) = \begin{cases} \delta \cdot w(u, v|t) + 1, & \text{if } (u, v) \in I_{t,t+1} \\ \delta \cdot w(u, v|t), & \text{if } (u, v) \notin I_{t,t+1} \end{cases} \quad (7)$$

3. The links with weight less than  $\theta$  are removed from the set  $T_{t+1}$ . The resulting set  $E_{t+1}$  is

the set of links in the network at time  $(t + 1)$ .

$$E_{t+1} = T_{t+1} - \{(u, v) : (u, v) \in T_{t+1} \wedge w(u, v|t + 1) \leq \theta\} \quad (8)$$

The value of  $\theta$  and  $\delta$  depend on the network, and  $0 < \theta < 1$  and  $0 \leq \delta \leq 1$ .

In the next chapter, we will describe how the existing feature scores are modified to apply to the network described in this chapter.

## 4 Modification of Feature Scores for Weighted and Directed Network

As mentioned in Chapter 2, previous works on link prediction have been performed on unweighted and undirected networks. So, most of the feature scores are independent of the link weight and direction. In this chapter we will extend the existing feature scores to apply to weighted and directed networks.

Notations:

1. In an undirected network, the set of neighbors of a node  $u$  is represented by  $\Gamma(u)$ , and  $w(u, v)$  is the weight of the link from  $u$  to  $v$ .
2. In a directed network, a node  $u$  has two types of neighbors - neighbors with links that are directed away from  $u$ , represented by  $\Gamma_o(u)$ , and neighbors with links that are directed to  $u$ , represented by  $\Gamma_i(u)$ .
3. The set  $\rho(u, v|l)$  denotes the set of paths of length  $l$  from node  $u$  to  $v$ ,

In the following sections, we modify the features discussed in Chapter 2 for weighted and directed networks:

### 4.1 Common Neighbor Score

Newman's work [8] on collaboration networks shows that there is a positive correlation between the number of common neighbors of two nodes and the probability that there will be a link between them in the future.

In an undirected network, the Common Neighbor Score of two nodes  $u$  and  $v$  is given by,

$$CN(u, v) = |\Gamma(u) \cap \Gamma(v)| \quad (9)$$

In a directed network, there are two Common Neighbor Scores - one based on the in-neighbors and the other based on the out-neighbors.

$$CN_o = |\Gamma_o(u) \cap \Gamma_o(v)| \quad (10)$$

$$CN_i = |\Gamma_i(u) \cap \Gamma_i(v)| \quad (11)$$

To take into account the weight of the links, the sum of the average weights of the links from the nodes  $u$  and  $v$  to the common neighbors is taken instead of simply the number of neighbors. So we propose to extend the Common Neighbor Score as follows:

$$CN_{o,weighted} = \sum_{z \in (\Gamma_o(u) \cap \Gamma_o(v))} \frac{w(u, z) + w(v, z)}{2} \quad (12)$$

$$CN_{i,weighted} = \sum_{z \in (\Gamma_i(u) \cap \Gamma_i(v))} \frac{w(z, u) + w(z, v)}{2} \quad (13)$$

## 4.2 Preferential Attachment Score

Work by Barabasi et al. [10] on collaboration networks suggest that there is a positive correlation between the probability of collaboration between two nodes and the product of the number of neighbors of the two nodes.

In an undirected and unweighted network, the Preferential Attachment Score between two nodes  $u$  and  $v$  is given by,

$$PA(u, v) = |\Gamma(u)| \cdot |\Gamma(v)| \quad (14)$$

In a directed network, since there are two neighbor sets, the two Preferential Attachment Scores will be given by,

$$PA_i(u, v) = |\Gamma_i(u)| \cdot |\Gamma_i(v)| \quad (15)$$

$$PA_o(u, v) = |\Gamma_o(u)| \cdot |\Gamma_o(v)| \quad (16)$$

In a directed weighted network, to incorporate the weights of the links into the Preferential Attachment Score, it can be redefined as the product of the average link weights of the two neighbors of the two nodes, written as,

$$PA_{i,weighted}(u, v) = \left( \frac{\sum_{z \in \Gamma_i(u)} w(z, u)}{|\Gamma_i(u)|} \right) \cdot \left( \frac{\sum_{z \in \Gamma_i(v)} w(z, v)}{|\Gamma_i(v)|} \right) \quad (17)$$

$$PA_{o,weighted}(u, v) = \left( \frac{\sum_{z \in \Gamma_o(u)} w(u, z)}{|\Gamma_o(u)|} \right) \cdot \left( \frac{\sum_{z \in \Gamma_o(v)} w(v, z)}{|\Gamma_o(v)|} \right) \quad (18)$$

### 4.3 Adamic-Adar Index

Adamic et al. [9] studied web pages and developed a measure to calculate the similarity between web pages. They found that rare features are more important than common ones while calculating the similarity between web pages. Their work can be easily applied to social networks - nodes take the place of web pages and links take the place of web links.

In an undirected network, the Adamic-Adar Index of two nodes  $u$  and  $v$  is given by,

$$AA(u, v) = \sum_{z \in (\Gamma(u) \cap \Gamma(v))} \frac{1}{\log(|\Gamma(z)|)} \quad (19)$$

If the network is directed, we will have two Adamic-Adar Index values based on the two neighbor sets.

$$AA_o(u, v) = \sum_{z \in (\Gamma_o(u) \cap \Gamma_o(v))} \frac{1}{\log |\Gamma_o(z)|} \quad (20)$$

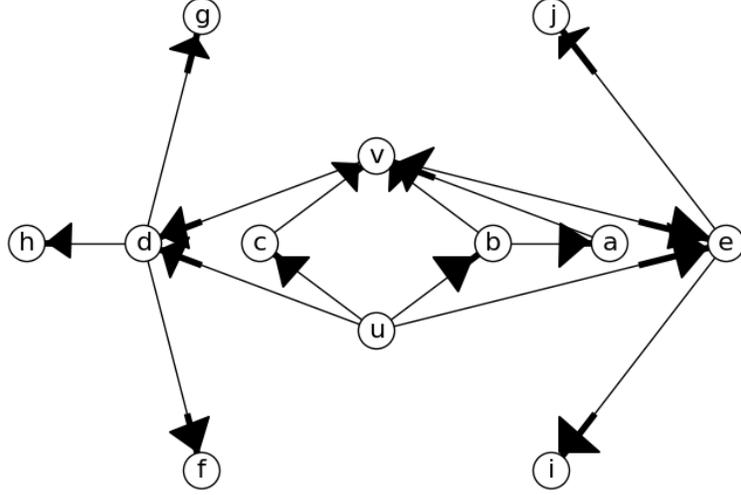


Figure 1: An example graph

$$AA_i(u, v) = \sum_{z \in (\Gamma_i(u) \cap \Gamma_i(v))} \frac{1}{\log |\Gamma_i(z)|} \quad (21)$$

To illustrate, consider the network in Figure 1. If we want to calculate the out-Adamic-Adar Index between nodes  $u$  and  $v$ , the contributions of the common nodes  $d$  and  $e$  will be given by,

$$\gamma(d) = \frac{1}{\log(w(d, f) + w(d, g) + w(d, h))} \quad (22)$$

$$\gamma(e) = \frac{1}{\log(w(e, i) + w(e, j))} \quad (23)$$

The contributions of  $\gamma(d)$  and  $\gamma(e)$  towards the final out-Adamic-Adar Index need to take into account the average weights of links to node  $d$  and  $e$  respectively.

$$AA_{o,weighted}(u, v) = \frac{1}{2} \left( \frac{w(u, d) + w(v, d)}{\gamma(d)} + \frac{w(u, e) + w(v, e)}{\gamma(e)} \right) \quad (24)$$

Equation 24 can be generalized as,

$$AA_{o,weighted}(u, v) = \frac{1}{2} \left( \sum_{z \in (\Gamma_o(u) \cap \Gamma_o(v))} \frac{w(u, z) + w(v, z)}{\log(\sum_{x \in \Gamma_o(z)} w(z, x))} \right) \quad (25)$$

Similarly, the in-Adamic-Adar Index can be derived to be,

$$AA_{i,weighted}(u, v) = \frac{1}{2} \left( \sum_{z \in (\Gamma_i(u) \cap \Gamma_i(v))} \frac{w(z, u) + w(z, v)}{\log(\sum_{x \in \Gamma_i(z)} w(x, z))} \right) \quad (26)$$

#### 4.4 Katz Score

Katz [11] defined a similarity measure based on the paths between two nodes. In the measure that Katz defined, the contribution of shorter paths to the final score is higher than that of the longer paths.

Mathematically, the Katz Score of  $(u, v)$  is given by,

$$KS(u, v) = \sum_{l=1}^{\infty} \beta^l \cdot |\rho(u, v|l)| \quad (27)$$

where  $\beta$  is the damping factor ( $0 \leq \beta \leq 1$ ) and  $|\rho(u, v|l)|$  is the number of paths of length  $l$  from node  $u$  to  $v$ . Liben-Nowell et al. [1] have shown that performance of the Katz Score is best for  $\beta \leq 0.005$ .

As seen in equation 27, the Katz Score depends on only the path between the nodes. So the same equation can be applied to directed networks. However, it still needs to be adjusted for link weights. To account for the link weights, the number of paths in equation 27 is replaced by the

sum of average link weights.

Consider the network in Figure 1. The sets of paths from node  $u$  to  $v$  are,

$$\rho(u, v|1) = \phi \quad (28)$$

$$\rho(u, v|2) = \{\{(u, b), (b, v)\}, \{(u, c), (c, v)\}\} \quad (29)$$

$$\rho(u, v|3) = \{\{(u, b), (b, a), (a, v)\}\} \quad (30)$$

The contribution of each path set towards the final Katz Score is as follows:

$$\gamma(1) = \beta^1 \cdot 0 \quad (31)$$

$$\gamma(2) = \beta^2 \cdot \left( \frac{w(u, b) + w(b, v)}{2} + \frac{w(u, c) + w(c, v)}{2} \right) \quad (32)$$

$$\gamma(3) = \beta^3 \cdot \left( \frac{w(u, b) + w(b, a) + w(a, v)}{3} \right) \quad (33)$$

The Katz Score from node  $u$  to  $v$  is then given by:

$$KS_{weighted}(u, v) = \gamma(1) + \gamma(2) + \gamma(3) \quad (34)$$

Equation 34 can be generalized as,

$$KS_{weighted}(u, v) = \sum_{l=1}^{\infty} \beta^l \cdot \sum_{p \in \rho(u, v|l)} \frac{\sum_{(x, y) \in p} w(x, y)}{l} \quad (35)$$

## 4.5 Shortest Path Score

In an undirected and unweighted network, the Shortest Path score of  $(u, v)$  is simply the reciprocal of the length of the shortest path from  $u$  to  $v$ .

To extend the Shortest Path Score to weighted and directed network, let us take the network in Fig. 1 as an example. Let  $l_{min}$  represent the length of the shortest path. In the example network, the length of the shortest path from  $u$  to  $v$  is 2.

The set of paths of length  $l_{min}$  is,

$$\rho(u, v|2) = \{ \{(u, b), (b, v)\}, \{(u, c), (c, v)\} \} \quad (36)$$

The set of average link weights in the path set  $\rho(u, v|2)$  is given by,

$$\bar{\rho}(u, v|2) = \left\{ \frac{w(u, b) + w(b, v)}{2}, \frac{w(u, c) + w(c, v)}{2} \right\} \quad (37)$$

The Shortest Path Score of  $(u, v)$  is then given by the average of the elements of the set  $\bar{\rho}(u, v|2)$ ,

$$SP_{weighted}(u, v) = \frac{1}{2} \left( \frac{w(u, b) + w(b, v)}{2} + \frac{w(u, c) + w(c, v)}{2} \right) \quad (38)$$

Equation 38 can be generalized as follows:

$$SP_{weighted}(u, v) = \frac{1}{|\rho(u, v|l_{min})|} \sum_{p \in \rho(u, v|l_{min})} \frac{\sum_{(x, y) \in p} w(x, y)}{l_{min}} \quad (39)$$

## 4.6 Rooted PageRank

Rooted PageRank is a modification of the PageRank algorithm developed by Brin et al. [12]. The Rooted PageRank between nodes  $u$  and  $v$  is defined as the probability that we will land on the node  $v$  when we perform a random walk starting from node  $u$  under the following conditions:

1. There is a  $(1 - \alpha)$  probability of moving to a randomly chosen neighbor node from the current node.
2. There is an  $\alpha$  probability of the random walk restarting from the root  $u$  again.

To account for the link weights the probability of selecting the next neighbor node is proportional to the associated link weight. This is shown in more detail in the algorithm for calculating the Rooted PageRank (Algorithm 1).

---

**Algorithm 1** Weighted PageRank Algorithm

---

```
1: procedure WEIGHTED-PAGERANK( $u, v$ )
2:    $N = 1000$  ▷ The number of walks to perform
3:    $i = 0$ 
4:    $n = 0$ 
5:    $currentNode = u$ 
6:   while  $i < N$  do
7:     if  $random(0,1) > \alpha$  then ▷ Reset the random walk
8:        $currentNode = u$ 
9:     else ▷ Select a neighbor node
10:       $probDist = []$ 
11:      for all  $n \in \Gamma_o(currentNode)$  do ▷ Probability of next node
12:         $probDist[n] = w(currentNode, n) / \sum_{z \in \Gamma_o(currentNode)} w(currentNode, z)$ 
13:      end for
14:       $currentNode = selectNode(\Gamma_o(currentNode), probDist)$ 
15:      if  $currentNode == v$  then ▷ Check if node v is selected
16:         $n = n+1$ 
17:      end if
18:       $i = i+1$ 
19:    end if
20:  end while
21:  return  $n/i$ 
22: end procedure
```

---

## 5 Incremental Method for Calculating Feature Scores

When the network feature scores are modified to take into account the link weights and direction (Chapter 4), the time complexity of calculating them is higher compared to the unweighted and undirected feature scores. In this chapter, we describe methods to calculate the weighted features by iteratively updating them as new interactions between nodes and removal of links are observed.

Notations:

1. The network is observed at time  $[1, 2, 3, \dots, T]$ .
2. The decay factor of the network is  $\delta$  per time interval.
3.  $I_{t,t+1}$  is the set of node pairs that interact in the time interval  $[t, t + 1]$ .
4.  $S$  is the set of all possible node pairs.

### 5.1 Common Neighbor Score

Let  $CN_o(u, v|t)$  and  $CN_i(u, v|t)$  be the Common Neighbor Score of node pair  $(u, v)$  at time  $t$  based on out-neighbor (equation 10) and in-neighbor (equation 11).

To update Common Neighbor Scores, we need to consider three cases:

1. Decay of link weight
2. Removal of link
3. Interaction between nodes

### 5.1.1 Effect of decay of link weight

If there were no interactions between the common neighbors of  $x$  and  $y$  at time interval  $[t, t + 1]$ , the weight of the link between two nodes  $a$  and  $b$  is,

$$w(a, b|t + 1) = \delta \cdot w(a, b|t) \quad (40)$$

Consider the in-Common Neighbor Score between nodes  $x$  and  $y$  at time  $t$ .

$$CN_i(x, y|t) = \sum_{z \in (\Gamma_i(x) \cap \Gamma_i(y))} \frac{w(z, x|t) + w(z, y|t)}{2} \quad (41)$$

The in-Common Neighbor Score of  $(x, y)$  at time  $t + 1$  is

$$CN_i(x, y|t + 1) = \sum_{z \in (\Gamma_i(x) \cap \Gamma_i(y))} \frac{w(z, x|t + 1) + w(z, y|t + 1)}{2} \quad (42)$$

$$= \delta \cdot \sum_{z \in (\Gamma_i(x) \cap \Gamma_i(y))} \frac{w(z, x|t) + w(z, y|t)}{2} \quad (43)$$

$$\boxed{CN_i(x, y|t + 1) = \delta \cdot CN_i(x, y|t)} \quad (44)$$

Similarly, the out-Common Neighbor Score is,

$$\boxed{CN_o(x, y|t + 1) = \delta \cdot CN_o(x, y|t)} \quad (45)$$

### 5.1.2 Effect of removal of link

Suppose that there was a link between nodes  $u$  and  $v$  at time  $t$  but it has been removed at time  $(t + 1)$ . The removal of the link  $(u, v)$  will affect  $CN_o(u, a)$  and  $CN_i(v, a)$  for a node  $a$  that was connected to  $u$  or  $v$  as shown in Figures 2a and 2b. Consider the node  $a$  connected as shown in

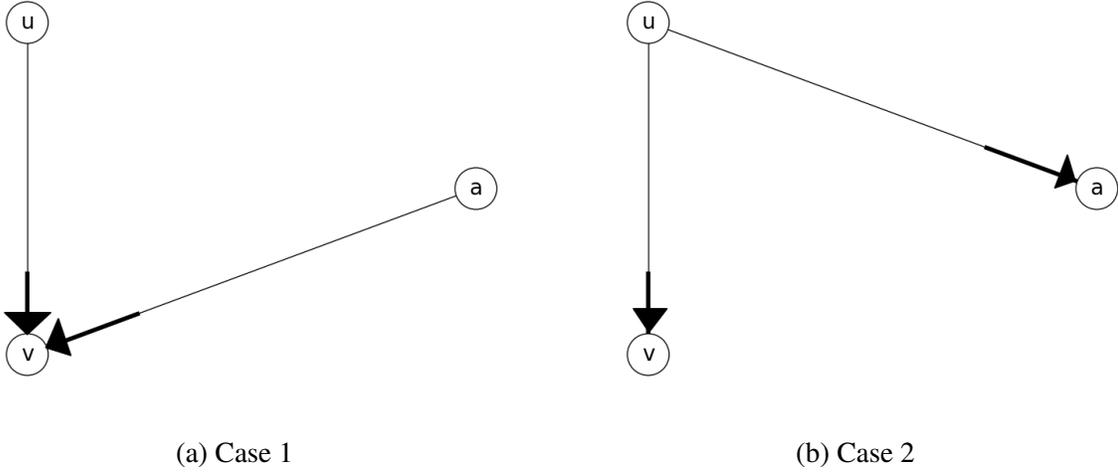


Figure 2: Node pairs that interaction between  $(u, v)$  will affect

Figure 2a and recall that,

$$CN_o(u, a|t) = \frac{w(u, v|t) + w(a, v|t)}{2} + \kappa, \quad (46)$$

where  $\kappa$  is the component of  $CN_o(u, a|t)$  that is independent of node  $v$ .

If the link  $(u, v)$  is removed at time  $(t + 1)$ ,

$$CN_o(u, a|t + 1) = \delta \cdot \kappa, \quad (47)$$

which can be written as,

$$CN_o(u, a|t + 1) = \delta \cdot \left( \frac{w(u, v|t) + w(a, v|t)}{2} + \kappa \right) - \delta \cdot \frac{w(u, v|t) + w(a, v|t)}{2} \quad (48)$$

$$\boxed{CN_o(u, a|t + 1) = \delta \cdot CN_o(u, a|t) - \delta \cdot \frac{w(u, v|t) + w(a, v|t)}{2}} \quad (49)$$

Removal of the link  $(u, v)$  will not affect the in-Common Neighbor Score of  $(u, a)$ .

Similarly, for the case where node  $a$  is connected as shown in Figures 2a, it can be shown that

the out-Common Neighbor Score of  $(a, v)$  will not be affected by the removal of the link  $(u, v)$  and,

$$\boxed{CN_i(a, v|t+1) = \delta \cdot CN_i(a, v|t) - \delta \cdot \frac{w(u, v|t) + w(u, a|t)}{2}} \quad (50)$$

### 5.1.3 Effect of interaction between nodes

If a pair of nodes  $(u, v)$  interact with each other in  $[t, t+1]$ , it will affect  $CN_o(u, a)$  and  $CN_i(v, a)$  for node  $a$  that is connected to  $u$  or  $v$ , as shown in Figures 2a and 2b, consider the node  $a$ , as shown in Figure 2a, at time  $(t+1)$ .

If there is a link from  $u$  to  $v$  at time  $t$ , then,

$$CN_o(u, a|t) = \frac{w(u, v|t) + w(a, v|t)}{2} + \kappa, \quad (51)$$

where  $\kappa$  is the component of  $CN_o(u, a|t)$  that is independent of node  $v$ .

At time  $t+1$ , the weights  $w(u, v)$  and  $w(a, v)$  are updated, and will be,

$$w(u, v|t+1) = \delta \cdot w(u, v) + 1 \quad (52)$$

$$w(a, v|t+1) = \delta \cdot w(a, v) \quad (53)$$

Consequently, at time  $(t+1)$ ,

$$CN_o(u, a|t+1) = \frac{w(u, v|t+1) + w(a, v|t+1)}{2} + \delta \cdot \kappa \quad (54)$$

$$= \frac{\delta \cdot w(u, v|t) + 1 + \delta \cdot w(a, v|t)}{2} + \delta \cdot \kappa \quad (55)$$

$$= \delta \cdot \left( \frac{w(u, v) + w(a, v)}{2} + \kappa \right) + \frac{1}{2} \quad (56)$$

$$= \delta \cdot CN_o(u, a|t) + 0.5 \quad (57)$$

On the other hand, if there was no link from  $u$  to  $v$  at time  $t$ ,

$$CN_o(u, a|t) = \kappa \quad (58)$$

$$CN_o(u, a|t+1) = \frac{w(u, v|t+1) + w(a, v|t+1)}{2} + \delta \cdot \kappa \quad (59)$$

$$\boxed{CN_o(u, a|t+1) = \frac{w(u, v|t+1) + w(a, v|t+1)}{2} + \delta \cdot CN_o(u, a|t)} \quad (60)$$

Likewise, if there is a node  $a$  as shown in Figure 2b at time  $(t+1)$ , if there is a link from  $u$  to  $v$  at time  $t$ , it can be similarly shown that,

$$CN_i(v, a|t+1) = \delta \cdot CN_i(v, a|t) + 0.5 \quad (61)$$

On the other hand, if there is no link from  $u$  to  $v$  at time  $t$ , then

$$\boxed{CN_i(a, v|t+1) = \frac{w(u, v|t+1) + w(u, a|t+1)}{2} + \delta \cdot CN_i(a, v|t)} \quad (62)$$

## 5.2 Preferential Attachment Score

Let  $PF_o(u, v|t)$  and  $PF_i(u, v|t)$  be the Preferential Attachment Score of node pair  $(u, v)$  at time  $t$ . To update  $PF_o(u, v|t)$  and  $PF_i(u, v|t)$ , we need to take into account the decay in the link weight, changes due to links that are removed and the effect of interactions between two nodes.

### 5.2.1 Effect of decay of link weight

Consider the out-Preferential Attachment Score between nodes  $u$  and  $v$  at time  $t$ .

$$PF_o(u, v|t) = \left( \frac{\sum_{z \in \Gamma_o(u|t)} w(u, z|t)}{|\Gamma_o(u|t)|} \right) \cdot \left( \frac{\sum_{z \in \Gamma_o(v|t)} w(v, z|t)}{|\Gamma_o(v|t)|} \right) \quad (63)$$

If there was no interactions between  $u$  and  $v$  and any of their neighbors in time interval  $[t, t + 1]$ ,

$$\Gamma_o(u|t + 1) = \Gamma_o(x|t) \quad (64)$$

$$\Gamma_o(v|t + 1) = \Gamma_o(x|t) \quad (65)$$

$$w(u, x|t + 1) = \delta \cdot w(u, x|t) \quad (66)$$

$$w(v, x|t + 1) = \delta \cdot w(v, x|t) \quad (67)$$

The out-Preferential Attachment Score at time  $(t + 1)$  is,

$$PF_o(u, v|t + 1) = \left( \frac{\sum_{z \in \Gamma_o(u|t+1)} w(u, z|t + 1)}{|\Gamma_o(u|t + 1)|} \right) \cdot \left( \frac{\sum_{z \in \Gamma_o(v|t+1)} w(v, z|t + 1)}{|\Gamma_o(v|t + 1)|} \right) \quad (68)$$

$$= \left( \frac{\sum_{z \in \Gamma_o(u|t)} \delta \cdot w(u, z|t)}{|\Gamma_o(u|t)|} \right) \cdot \left( \frac{\sum_{z \in \Gamma_o(v|t)} \delta \cdot w(v, z|t)}{|\Gamma_o(v|t)|} \right) \quad (69)$$

$$= \delta^2 \cdot \left( \frac{\sum_{z \in \Gamma_o(u|t+1)} w(u, z|t + 1)}{|\Gamma_o(u|t + 1)|} \right) \cdot \left( \frac{\sum_{z \in \Gamma_o(v|t+1)} w(v, z|t + 1)}{|\Gamma_o(v|t + 1)|} \right) \quad (70)$$

$$\boxed{PF_o(u, v|t + 1) = \delta^2 \cdot PF_o(u, v|t)} \quad (71)$$

Similarly, the in-Preferential Attachment Score at time  $(t + 1)$  can be show to be,

$$\boxed{PF_i(u, v|t + 1) = \delta^2 \cdot PF_i(u, v|t)} \quad (72)$$

### 5.2.2 Effect of removal of link

Suppose that there was a link between nodes  $u$  and  $v$  at time  $t$  but it is broken at  $(t + 1)$ .

Consider the Preferential Attachment Score between nodes  $u$  and  $x$  at time  $t$ .

$$PF_o(u, x|t) = \left( \frac{\sum_{z \in \Gamma_o(u|t)} w(u, z|t)}{|\Gamma_o(u|t)|} \right) \cdot \left( \frac{\sum_{z \in \Gamma_o(x|t)} w(x, z|t)}{|\Gamma_o(x|t)|} \right) \quad (73)$$

The Preferential Attachment Score of  $(u, x)$  at time  $(t + 1)$ ,

$$PF_o(u, x|t + 1) = \left( \frac{\sum_{z \in \Gamma_o(u|t+1)} w(u, z|t + 1)}{|\Gamma_o(u|t + 1)|} \right) \cdot \left( \frac{\sum_{z \in \Gamma_o(x|t+1)} w(x, z|t + 1)}{|\Gamma_o(x|t + 1)|} \right) \quad (74)$$

Since the link  $(u, v)$  is broken at time  $t + 1$ ,

$$PF_o(u, x|t + 1) = \left( \frac{\sum_{z \in \Gamma_o(x|t)} \delta \cdot w(x, z|t)}{|\Gamma_o(x|t)|} \right) \quad (75)$$

$$\begin{aligned} & \cdot \left( \frac{(\sum_{z \in \Gamma_o(u|t)} \delta \cdot w(u, z|t)) - \delta \cdot w(u, v|t)}{|\Gamma_o(u|t) - 1|} \right) \\ & = \delta^2 \cdot \frac{|\Gamma_o(u|t)|}{|\Gamma_o(u|t) - 1|} \cdot \left( \frac{\sum_{z \in \Gamma_o(u|t)} w(u, z|t)}{|\Gamma_o(u|t)|} \right) \cdot \left( \frac{\sum_{z \in \Gamma_o(x|t)} w(x, z|t)}{|\Gamma_o(x|t)|} \right) \\ & \quad - \delta \cdot \left( \frac{w(u, v|t)}{|\Gamma_o(u|t) - 1|} \right) \cdot \left( \frac{\sum_{z \in \Gamma_o(x|t)} w(x, z|t)}{|\Gamma_o(x|t)|} \right) \end{aligned} \quad (76)$$

$$\boxed{PF_o(u, x|t + 1) = \delta^2 \cdot \frac{|\Gamma_o(u|t)|}{|\Gamma_o(u|t) - 1|} \cdot PF_o(u, x|t) - \delta \cdot \left( \frac{w(u, v|t)}{|\Gamma_o(u|t) - 1|} \right) \cdot \left( \frac{\sum_{z \in \Gamma_o(x|t)} w(x, z|t)}{|\Gamma_o(x|t)|} \right)} \quad (77)$$

Similarly,

$$\boxed{PF_i(x, v|t + 1) = \delta^2 \cdot \frac{|\Gamma_i(v|t)|}{|\Gamma_i(v|t) - 1|} \cdot PF_i(x, v|t) - \delta \cdot \left( \frac{w(u, v|t)}{|\Gamma_i(v|t) - 1|} \right) \cdot \left( \frac{\sum_{z \in \Gamma_i(x|t)} w(z, x|t)}{|\Gamma_i(x|t)|} \right)} \quad (78)$$

### 5.2.3 Effect of interaction between nodes

Consider the out-Preferential Attachment Score between nodes  $u$  and  $x$  at time  $t$ .

$$PF_o(u, x|t) = \left( \frac{\sum_{z \in \Gamma_o(u|t)} w(u, z|t)}{|\Gamma_o(u|t)|} \right) \cdot \left( \frac{\sum_{z \in \Gamma_o(x|t)} w(x, z|t)}{|\Gamma_o(x|t)|} \right) \quad (79)$$

Let,

$$P = \left( \frac{\sum_{z \in \Gamma_o(x|t)} w(x, z|t)}{|\Gamma_o(x|t)|} \right) \quad (80)$$

Assume there is an interaction between nodes  $u$  and  $v$  at time interval  $[t, t + 1]$  and the link  $(u, v)$  existed at time  $t$ ,

$$PF_o(u, x|t + 1) = \left( \frac{(\sum_{z \in \Gamma_o(u|t)} \delta \cdot w(u, z|t)) + 1}{|\Gamma_o(u|t)|} \right) \cdot \left( \frac{\sum_{z \in \Gamma_o(v|t)} \delta \cdot w(x, z|t)}{|\Gamma_o(x|t)|} \right) \quad (81)$$

$$= \delta^2 \cdot \left( \frac{\sum_{z \in \Gamma_o(u|t)} w(u, z|t)}{|\Gamma_o(u|t)|} \right) \cdot P + \frac{\delta \cdot P}{|\Gamma_o(u|t)|} \quad (82)$$

$$\boxed{PF_o(u, x|t + 1) = \delta^2 \cdot PF_o(u, x|t) + \delta \cdot \left( \frac{\sum_{z \in \Gamma_o(x|t)} w(x, z|t)}{|\Gamma_o(u|t)| \cdot |\Gamma_o(x|t)|} \right)} \quad (83)$$

If the link  $(u, v)$  did not exist at time  $t$ ,

$$PF_o(u, x|t + 1) = \left( \frac{(\sum_{z \in \Gamma_o(u|t)} \delta \cdot w(u, z|t)) + 1}{|\Gamma_o(u|t)| + 1} \right) \cdot \left( \frac{\sum_{z \in \Gamma_o(v|t)} \delta \cdot w(x, z|t)}{|\Gamma_o(x|t)|} \right) \quad (84)$$

$$= \left( \frac{(\sum_{z \in \Gamma_o(u|t)} \delta \cdot w(u, z|t)) + 1}{|\Gamma_o(u|t)|} \right) \cdot \delta \cdot P \cdot \left( \frac{|\Gamma_o(u|t)|}{|\Gamma_o(u|t)| + 1} \right) \quad (85)$$

$$= \left( \frac{\delta^2 \cdot P \cdot (\sum_{z \in \Gamma_o(u|t)} w(u, z|t)) + \delta \cdot P}{|\Gamma_o(u|t)|} \right) \cdot \left( \frac{|\Gamma_o(u|t)|}{|\Gamma_o(u|t)| + 1} \right) \quad (86)$$

$$\boxed{PF_o(u, x|t + 1) = \delta \cdot \left( \frac{\sum_{z \in \Gamma_o(x|t)} w(x, z|t)}{|\Gamma_o(x|t)| \cdot (|\Gamma_o(u|t)| + 1)} \right) + \delta^2 \cdot \left( \frac{|\Gamma_o(u|t)|}{|\Gamma_o(u|t)| + 1} \right) \cdot PF_o(u, x|t)} \quad (87)$$

Similarly, for the node pair  $(v, x)$  if  $(u, v)$  existed at time  $t$ , the in-Preferential Attachment is,

$$\boxed{PF_i(v, x|t + 1) = \delta^2 \cdot PF_i(v, x|t) + \delta \cdot \left( \frac{\sum_{z \in \Gamma_i(x|t)} w(z, x|t)}{|\Gamma_i(v|t)| \cdot |\Gamma_i(x|t)|} \right)} \quad (88)$$

If  $(u, v)$  did not exist at time  $t$ ,

$$PF_i(v, x|t+1) = \delta \cdot \left( \frac{\sum_{z \in \Gamma_i(x|t)} w(z, x|t)}{|\Gamma_i(x|t)| \cdot (|\Gamma_i(v|t)| + 1)} \right) + \delta^2 \cdot \left( \frac{|\Gamma_i(v|t)|}{|\Gamma_i(v|t)| + 1} \right) \cdot PF_i(v, x|t) \quad (89)$$

### 5.3 Adamic Adar Index

Let  $AA_o(u, v|t)$  and  $AA_i(u, v|t)$  be the Adamic-Adar Index of node pair  $(u, v)$  at time  $t$ . To update  $AA_o(u, v|t)$  and  $AA_i(u, v|t)$ , we need to take into account the decay in the link weight, the changes due to links that are removed and the effect of interactions between two nodes.

#### 5.3.1 Effect of decay of link weight

Assume that for a pair of nodes  $u$  and  $v$ , there are no interactions between nodes that affect the Adamic-Adar Index of  $(u, v)$  at time  $(t+1)$ .

Consider the out-Adamic-Adar Index between nodes  $u$  and  $v$  at time  $t$ .

$$AA_o(u, v|t) = \frac{1}{2} \left( \sum_{z \in (\Gamma_o(u|t) \cap \Gamma_o(v|t))} \frac{w(u, z|t) + w(v, z|t)}{\log(\sum_{x \in \Gamma_o(z|t)} w(z, x|t))} \right) \quad (90)$$

Let,

$$A(u, v, z|t) = w(u, z|t) + w(v, z|t) \quad (91)$$

$$B(z|t) = \log\left(\sum_{x \in \Gamma_o(z|t)} w(z, x|t)\right) \quad (92)$$

At time  $(t + 1)$ ,

$$AA_o(u, v|t + 1) = \frac{1}{2} \left( \sum_{z \in (\Gamma_o(u|t+1) \cap \Gamma_o(v|t+1))} \frac{(w(u, z|t + 1) + w(v, z|t + 1))}{\log(\sum_{x \in \Gamma_o(z|t+1)} w(z, x|t + 1))} \right) \quad (93)$$

$$= \frac{1}{2} \left( \sum_{z \in (\Gamma_o(u|t) \cap \Gamma_o(v|t))} \frac{\delta \cdot (w(u, z|t) + w(v, z|t))}{\log(\delta) + \log(\sum_{x \in \Gamma_o(z|t)} w(z, x|t))} \right) \quad (94)$$

$$= \frac{1}{2} \left( \sum_{z \in (\Gamma_o(u|t) \cap \Gamma_o(v|t))} \frac{\delta \cdot A(u, v, z|t)}{\log(\delta) + B(z|t)} \right) \quad (95)$$

Let us assume that there is an  $F$  such that,

$$\frac{\delta \cdot A(u, v, z|t)}{\log(\delta) + B(z|t)} = F \cdot \frac{A(u, v, z|t)}{B(z|t)} \quad (96)$$

Then,

$$F = \frac{\delta}{\log(\delta)/B(z|t) + 1} \quad (97)$$

For  $\delta \rightarrow 1$ ,

$$\frac{\log(\delta)}{B(z|t)} \approx 0 \quad (98)$$

$$F \approx \delta \quad (99)$$

So,

$$\boxed{AA_o(u, v|t + 1) \approx \delta \cdot AA_o(u, v|t)} \quad (100)$$

Similarly,

$$\boxed{AA_i(u, v|t + 1) \approx \delta \cdot AA_i(u, v|t)} \quad (101)$$

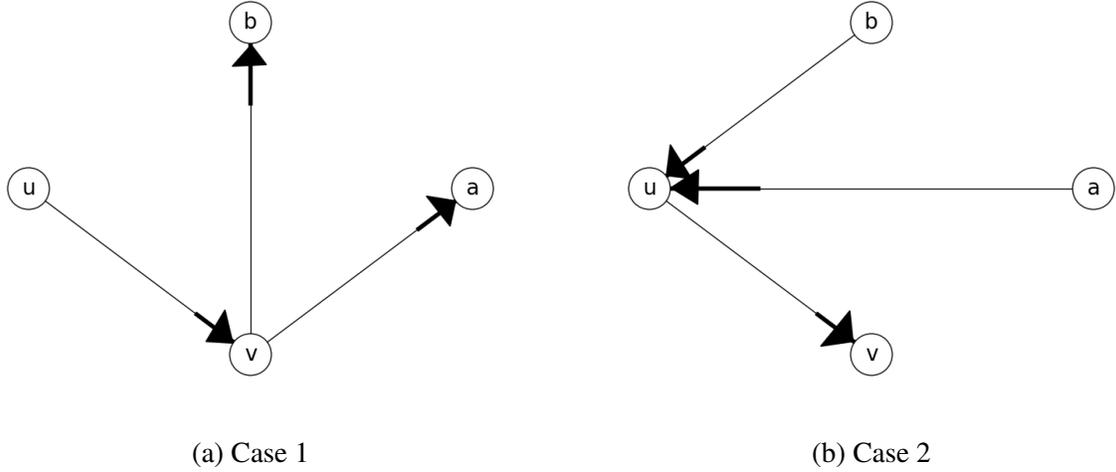


Figure 3: Node pairs that interaction between  $(u, v)$  will affect

### 5.3.2 Effect of removal of links

Suppose that there was a link between  $u$  and  $v$  at time  $t$ , but it is removed at time  $(t + 1)$ . The removal of  $(u, v)$  will affect the Adamic-Adar Index of  $(u, a)$  and  $(v, a)$  for a node  $a$  which was connected to  $u$  or  $v$  as shown in Figure 2a and 2b at time  $t$ . If there are two nodes  $a$  and  $b$  that was connected to  $u$  or  $v$  at time  $t$  as shown in Figure 3a and 3b, it will affect the Adamic-Adar Index of  $(a, b)$ .

1. Let us assume that the nodes  $u, v$  and  $a$  were connected to each other as shown in Figure 2a at time  $t$ .

At time  $t$ ,

$$AA_o(u, a|t) = \frac{w(u, v|t) + w(a, v|t)}{2 \cdot \log(\sum_{z \in \Gamma_o(v|t)} w(v, z|t))} + \kappa \quad (102)$$

where  $\kappa$  is the component of the  $AA_o(u, a|t)$  that is independent of node  $v$ .

At time  $(t + 1)$ ,

$$AA_o(u, a|t + 1) = \delta \cdot \kappa \quad (103)$$

$$AA_o(u, a|t+1) = \delta \cdot \left( \frac{w(u, v|t) + w(a, v|t)}{2 \cdot \log \left( \sum_{z \in \Gamma_o(v|t)} w(v, z|t) \right)} + \kappa \right) - \frac{\delta \cdot (w(u, v|t) + w(a, v|t))}{2 \cdot \log \left( \sum_{z \in \Gamma_o(v|t)} w(v, z|t) \right)} \quad (104)$$

$$AA_o(u, a|t+1) = \delta \cdot AA_o(u, a|t) - \frac{\delta \cdot (w(u, v|t) + w(a, v|t))}{2 \cdot \log \left( \sum_{z \in \Gamma_o(v|t)} w(v, z|t) \right)} \quad (105)$$

2. If the nodes  $u$ ,  $v$  and  $a$  were connected to each other as shown in Figure 2b at time  $t$ , it can be similarly shown that,

$$AA_i(v, a|t+1) = \delta \cdot AA_i(v, a|t) - \frac{\delta \cdot (w(u, v|t) + w(u, a|t))}{2 \cdot \log \left( \sum_{z \in \Gamma_i(u|t)} w(z, u|t) \right)} \quad (106)$$

3. If the nodes  $u$ ,  $v$ ,  $a$  and  $b$  were connected to each other as shown in Figure 3b at time  $t$ , it can be shown that,

$$AA_o(a, b|t+1) = \delta \cdot AA_o(a, b|t) - \frac{\delta \cdot (w(a, u|t) + w(b, u|t))}{2 \cdot \log \left( \sum_{z \in \Gamma_o(u|t)} w(u, z|t) \right)} + \frac{\delta \cdot (w(a, u|t) + w(b, u|t))}{2 \cdot \log \left( \sum_{z \in \Gamma_o(u|t+1)} w(u, z|t+1) \right)} \quad (107)$$

4. If the nodes  $u$ ,  $v$ ,  $a$  and  $b$  were connected to each other as shown in Figure 3a at time  $t$ ,

$$AA_i(a, b|t+1) = \delta \cdot AA_i(a, b|t) - \frac{\delta \cdot (w(v, a|t) + w(v, b|t))}{2 \cdot \log \left( \sum_{z \in \Gamma_i(v|t)} w(z, v|t) \right)} + \frac{\delta \cdot (w(v, a|t) + w(v, b|t))}{2 \cdot \log \left( \sum_{z \in \Gamma_i(v|t+1)} w(z, v|t+1) \right)} \quad (108)$$

### 5.3.3 Effect of interaction between nodes

If a pair of nodes  $(u, v)$  interact in  $[t, t+1]$ , it will affect the Adamic-Adar Index of  $(u, a)$  and  $(v, a)$  where  $a$  is a node connected to  $u$  or  $v$  as shown in Figure 2a and 2b. If there are two nodes  $a$  and  $b$

that are connected to  $u$  or  $v$  as shown in Figure 3a and 3b, it will affect the Adamic-Adar Index of  $(a, b)$

1. Let us consider the nodes  $u, v$  and  $a$  connected to each other as shown in Figure 2a at time  $(t + 1)$ . If the link  $(u, v)$  exist at time  $t$ ,

$$AA_o(u, a|t) = \frac{w(u, v|t) + w(a, v|t)}{2 \cdot \log \left( \sum_{z \in \Gamma_o(v|t)} w(v, z|t) \right)} + \kappa \quad (109)$$

Then, at time  $(t + 1)$ ,

$$AA_o(u, a|t + 1) = \frac{\delta \cdot w(u, v|t) + \delta \cdot w(a, v|t) + 1}{2 \cdot \log \left( \sum_{z \in \Gamma_o(v|t)} \delta \cdot w(v, z|t) \right)} + \delta \cdot \kappa \quad (110)$$

$$\boxed{AA_o(u, a|t + 1) = \delta \cdot AA_o(u, a|t) + \frac{1}{2 \cdot \log \left( \sum_{z \in \Gamma_o(v|t)} \delta \cdot w(v, z|t) \right)}} \quad (111)$$

If the link  $(u, v)$  did not exist at time  $t$ , at time  $t$ ,

$$AA_o(u, a|t) = \kappa \quad (112)$$

Then, at time  $(t + 1)$ ,

$$AA_o(u, a|t + 1) = \frac{1 + \delta \cdot w(a, v|t)}{2 \cdot \log \left( \sum_{z \in \Gamma_o(v|t)} \delta \cdot w(v, z|t) \right)} + \delta \cdot \kappa \quad (113)$$

$$\boxed{AA_o(u, a|t + 1) = \frac{1 + \delta \cdot w(a, v|t)}{2 \cdot \log \left( \sum_{z \in \Gamma_o(v|t)} \delta \cdot w(v, z|t) \right)} + \delta \cdot AA_o(u, v|t)} \quad (114)$$

2. Suppose the nodes  $u, v$  and  $a$  are connected to each other as shown in Figure 2b at time  $(t + 1)$ . If  $(u, v)$  existed in time  $t$ , it can be similarly shown that,

$$\boxed{AA_i(a, v|t + 1) = \delta \cdot AA_i(a, v|t) + \frac{1}{2 \cdot \log \left( \sum_{z \in \Gamma_i(u|t)} \delta \cdot w(z, u|t) \right)}} \quad (115)$$

If  $(u, v)$  did not exist in time  $t$ ,

$$\boxed{AA_i(a, v|t+1) = \delta \cdot AA_i(a, v|t) + \frac{1 + \delta \cdot w(a, u|t)}{2 \cdot \log(\sum_{z \in \Gamma_i(u|t)} \delta \cdot w(z, u|t))}} \quad (116)$$

3. Let us consider the nodes  $u, v, a$  and  $b$  connected as shown in Figure 3b. At time  $t$ ,

$$AA_o(a, b|t) = \left( \frac{w(a, u|t) + w(b, u|t)}{\log(\sum_{z \in \Gamma_o(u|t)} w(u, z|t))} \right) + \kappa, \quad (117)$$

where  $\kappa$  is the component of  $AA_o(a, b|t)$  that is independent of the node  $u$ .

At time  $(t+1)$ ,

$$AA_o(a, b|t+1) = \left( \frac{w(a, u|t+1) + w(b, u|t+1)}{\log(\sum_{z \in \Gamma_o(u|t+1)} w(u, z|t+1))} \right) + \delta \cdot \kappa \quad (118)$$

$$= \left( \frac{\delta \cdot w(a, u|t) + \delta \cdot w(b, u|t)}{\log(\sum_{z \in \Gamma_o(u|t+1)} w(u, z|t+1))} \right) + \delta \cdot \kappa \quad (119)$$

Then,

$$\boxed{AA_o(a, b|t+1) = \delta \cdot AA_o(a, b|t) + \delta \cdot \left( \frac{w(a, u|t) + w(b, u|t)}{\log(\sum_{z \in \Gamma_o(u|t+1)} w(u, z|t+1))} - \frac{w(a, u|t) + w(b, u|t)}{\log(\sum_{z \in \Gamma_o(u|t)} w(u, z|t))} \right)} \quad (120)$$

4. If the nodes  $u, v, a$  and  $b$  are connected as show in Figure 3a at time  $(t+1)$ .

$$\boxed{AA_i(a, b|t+1) = \delta \cdot AA_i(a, b|t) + \delta \cdot \left( \frac{w(v, a|t) + w(v, b|t)}{\log(\sum_{z \in \Gamma_i(v|t+1)} w(z, v|t+1))} - \frac{w(a, u|t) + w(b, u|t)}{\log(\sum_{z \in \Gamma_i(u|t)} w(z, v|t))} \right)} \quad (121)$$

## 5.4 Katz Score

Let us assume that  $KS(u, v|t)$  is the Katz Score of  $(u, v)$  at time  $t$  and  $\beta^l$  is the damping factor for the Katz Score. Work by Libel-Nowell et al. [1] have shown that  $\beta < 0.005$ . So for paths of length four or above, the damping factor becomes  $6.25 \times 10^{-10}$  or lower. So calculations of Katz Score in real world networks are approximated by considering only paths with length three or less.

### 5.4.1 Effect of decay of link weight

Assume that  $KS(u, v|t)$  is the Katz Score at time  $t$  and there is no interactions among any nodes that affect  $KS(u, v|t)$ .

$$KS(u, v|t) = \sum_{l=1}^3 \beta^l \cdot \sum_{p \in \rho(u, v|l|t)} \left( \frac{\sum_{(x, y) \in p} w(x, y|t)}{l} \right) \quad (122)$$

Then, at time  $(t + 1)$ ,

$$KS(u, v|t + 1) = \sum_{l=1}^3 \beta^l \cdot \sum_{p \in \rho(u, v|l|t)} \left( \frac{\sum_{(x, y) \in p} w(x, y|t + 1)}{l} \right) \quad (123)$$

$$= \delta \cdot \sum_{l=1}^3 \beta^l \cdot \sum_{p \in \rho(u, v|l|t)} \left( \frac{\sum_{(x, y) \in p} w(x, y|t)}{l} \right) \quad (124)$$

$$\boxed{KS(u, v|t + 1) = \delta \cdot KS(u, v|t)} \quad (125)$$

### 5.4.2 Effect of removal of link

Let us assume that there was a link between nodes  $u$  and  $v$  at time  $t$ , but it is removed at time  $(t + 1)$ . This will affect the Katz score between  $(u, v)$ ,  $(u, a)$ ,  $(a, v)$  and  $(a, b)$  if they were connected as shown in Figure 4 at time  $t$ .

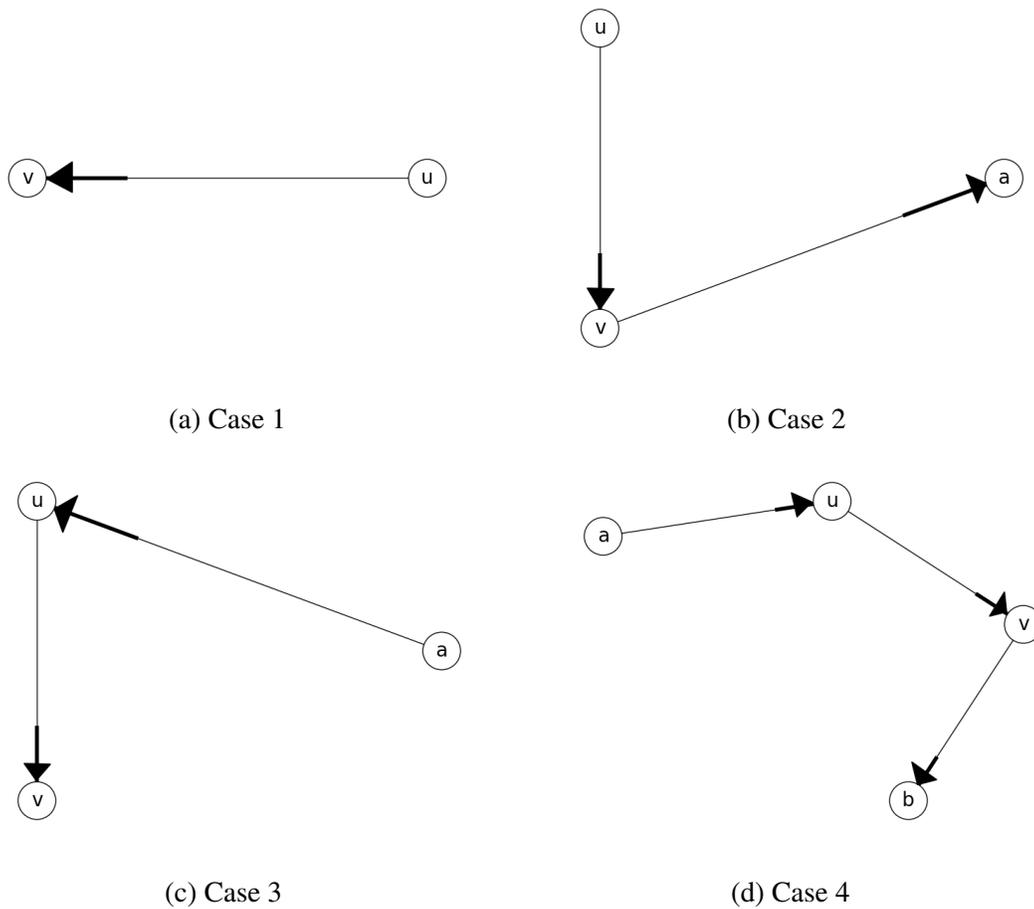


Figure 4: Node pairs for which interaction between  $(u, v)$  will affect the Katz Score

1. Case 1. (Figure: 4a)

If the Katz Score of  $(u, v)$  at time  $t$  is,

$$KS(u, v|t) = \beta \cdot w(u, v) + \kappa, \quad (126)$$

where  $\kappa$  is the component of  $KS(u, v|t)$  that is independent of the path  $(u, v)$ .

At time  $t$ , if the link  $(u, v)$  has been removed,

$$KS(u, v|t+1) = \delta \cdot \kappa \quad (127)$$

$$= \delta \cdot (\beta \cdot w(u, v|t) + \kappa) - \delta \cdot \beta \cdot w(u, v|t) \quad (128)$$

$$\boxed{KS(u, v|t+1) = \delta \cdot KS(u, v|t) - \delta \cdot \beta \cdot w(u, v|t)} \quad (129)$$

2. Similarly in Case 2. (Figure: 4b),

$$\boxed{KS(u, a|t+1) = \delta \cdot KS(u, a|t) - \delta \cdot \beta^2 \cdot \frac{w(u, v|t) + w(v, a|t)}{2}} \quad (130)$$

3. Similarly in Case 3. (Figure: 4c),

$$\boxed{KS(a, v|t+1) = \delta \cdot KS(a, v|t) - \delta \cdot \beta^2 \cdot \frac{w(u, v|t) + w(a, u|t)}{2}} \quad (131)$$

4. Similarly in Case 4. (Figure: 4d),

$$\boxed{KS(a, b|t+1) = \delta \cdot KS(a, b|t) - \delta \cdot \beta^3 \cdot \frac{w(a, u|t) + w(u, v|t) + w(v, b|t)}{3}} \quad (132)$$

### 5.4.3 Effect of interaction between nodes

Let us assume that interaction between  $(u, v)$  affects the Katz score between two nodes. The interaction will affect the Katz score between  $(u, v)$ ,  $(u, a)$ ,  $(a, v)$  and  $(a, b)$  if they are connected as shown in Figure 4.

1. If the nodes  $u$  and  $v$  are connected as shown in Figure 4a,

$$w(u, v|t+1) = \delta \cdot w(u, v|t) + 1 \quad (133)$$

$$KS(u, v|t+1) = \sum_{l=1}^3 \beta^l \cdot \sum_{p \in \rho(u, v|l|t+1)} \left( \frac{\sum_{(x, y) \in p} w(x, y|t+1)}{l} \right) \quad (134)$$

$$= \left( \sum_{l=1}^3 \beta^l \cdot \sum_{p \in \rho(u, v|l)} \left( \frac{\sum_{(x, y) \in p} \delta \cdot w(x, y|t)}{l} \right) \right) + \beta \quad (135)$$

$$\boxed{KS(u, v|t+1) = \delta \cdot KS(u, v|t) + \beta} \quad (136)$$

2. If the nodes are connected as shown in Figure 4b, and there is a link from  $u$  to  $v$  at time  $t$ ,

$$\rho(u, v|l|t+1) = \rho(u, v|l|t) \quad (137)$$

$$w(u, v|t+1) = \delta \cdot w(u, v|t) + 1 \quad (138)$$

$$KS(u, a|t+1) = \delta \cdot \left( \sum_{l=1}^3 \beta^l \cdot \sum_{p \in \rho(u, a|l|t)} \left( \frac{\sum_{(x, y) \in p} w(x, y|t+1)}{l} \right) \right) + \frac{\beta^2}{2} \quad (139)$$

$$\boxed{KS(u, a|t+1) = \delta \cdot KS(u, a|t) + \frac{\beta^2}{2}} \quad (140)$$

If there was no link from  $u$  to  $v$  at time  $t$ ,

$$\rho(u, a|l|t+1) = \rho(u, a|l|t) + \{(u, v), (v, a)\} \quad (141)$$

$$\boxed{KS(u, a|t+1) = \delta \cdot KS(u, a|t) + \beta^2 \cdot \frac{1 + \delta \cdot w(v, a)}{2}} \quad (142)$$

3. Similarly, if the nodes are connected as shown in Figure 4c and there is a link from  $u$  to  $v$  at time  $t$ ,

$$\boxed{KS(a, v|t+1) = \delta \cdot KS(a, v|t) + \frac{\beta^2}{2}} \quad (143)$$

If there was no link from  $u$  to  $v$  at time  $t$ ,

$$\boxed{KS(a, v|t+1) = \delta \cdot KS(a, v|t) + \beta^2 \cdot \frac{1 + \delta \cdot w(a, u|t)}{2}} \quad (144)$$

4. Similarly, if the nodes are connected as shown in Figure 4d, and there is a link from  $u$  to  $v$  at time  $t$ ,

$$\boxed{KS(a, b|t+1) = \delta \cdot KS(a, b|t) + \frac{\beta^3}{3}} \quad (145)$$

If there was no link from  $u$  to  $v$  at time  $t$ ,

$$\boxed{KS(a, b|t+1) = \beta^2 \cdot \frac{1 + \delta \cdot w(a, u|t) + \delta \cdot w(v, b|t)}{3} + \delta \cdot KS(a, b|t)} \quad (146)$$

## 5.5 Shortest Path Score

Let  $SP(u, v|t)$  be the shortest path score between two nodes at time  $t$ ,  $l_{min}(u, v|t)$  be the length of the shortest path from node  $u$  to  $v$  at time  $t$  and  $\rho_{min}(u, v|t)$  be the set of paths of length  $l_{min}(u, v|t)$  from  $u$  to  $v$  at  $t$ .

### 5.5.1 Effect of decay of link weight

The Shortest Path Score of  $(u, v)$  at time  $t$  is,

$$SP(u, v|t) = \frac{1}{|\rho_{min}(u, v|t)|} \sum_{p \in \rho_{min}(u, v|t)} \frac{\sum_{(x,y) \in p(x,y|t)} w(x, y|t)}{l_{min}(u, v|t)} \quad (147)$$

If the interactions between the nodes in time interval  $[t, t + 1]$  do not affect  $l_{min}(u, v|t)$  or  $\rho_{min}(u, v|t)$ , the Shortest Path Score at time  $(t + 1)$  is,

$$SP(u, v|t + 1) = \frac{1}{|\rho_{min}(u, v|t + 1)|} \cdot \sum_{p \in \rho_{min}(u, v|t + 1)} \frac{\sum_{(x,y) \in p(x,y|t + 1)} w(x, y|t + 1)}{l_{min}(u, v|t + 1)} \quad (148)$$

$$= \frac{1}{|\rho_{min}(u, v|t)|} \sum_{p \in \rho_{min}(u, v|t)} \frac{\sum_{(x,y) \in p(x,y|t + 1)} \delta \cdot w(x, y|t)}{l_{min}(u, v|t)} \quad (149)$$

$$\boxed{SP(u, v|t + 1) = \delta \cdot SP(u, v|t)} \quad (150)$$

### 5.5.2 Effect of interaction between nodes

If there is an interaction between nodes  $(u, v)$  at time interval  $[t, t + 1]$ , the interaction will change the Shortest Path Score of  $(x, y)$  under two cases.

1. The link  $(u, v)$  did not exist in  $E_t$ , and the new interaction creates a new shortest path.
2. The link  $(u, v)$  exist in  $E_t$ , nodes  $u$  and  $v$  lies in the shortest path from  $x$  to  $y$ .

Suppose the link  $(u, v)$  did not exist in  $E_t$ , if the interaction between  $(u, v)$  forms a new shortest path from  $x$  to  $y$ :

1. The new shortest path should contain  $(u, v)$ .

2. The path from  $x$  to  $u$  in the new shortest path should be the shortest path from  $x$  to  $u$ .
3. The path from  $v$  to  $y$  in the new shortest path should be the shortest path from  $v$  to  $y$ .

So, the interaction between  $u$  and  $v$  will create a new shortest path from  $x$  to  $y$  if,

$$l_{min}(x, y|t) \geq l_{min}(x, u|t) + l_{min}(v, y|t) + 1 \quad (151)$$

If a new shortest path is created from  $x$  to  $y$ ,

$$l_{min}(x, y|t+1) = l_{min}(x, u|t) + l_{min}(v, y|t) + 1 \quad (152)$$

$$\rho_{min}(x, y|t+1) = \{\rho_{min}(x, v|t) + (u, v) + \rho_{min}(v, y|t)\} \quad (153)$$

However, if the link  $(u, v)$  exist in  $E_t$  and if  $(u, v)$  is in the shortest path from  $x$  to  $y$ ,

$$w(u, v|t+1) = \delta \cdot w(u, v|t) + 1 \quad (154)$$

$$\boxed{SP(x, y|t+1) = \delta \cdot SP(x, y|t) + \frac{1}{|\rho_{min}(x, y|t)| \cdot l_{min}(x, y|t)}} \quad (155)$$

If  $(u, v)$  is in not the shortest path from  $x$  to  $y$ ,

$$\boxed{SP(x, y|t+1) = \delta \cdot SP(x, y|t)} \quad (156)$$

### 5.5.3 Effect of removal of link

Let us assume that there was a link between nodes  $u$  and  $v$  at time  $t$ , but not at time  $t$ . Let us assume that  $(u, v)$  was in a shortest path  $p$  from node  $x$  to  $y$ .

Then removal of  $(u, v)$  will break the path  $p$ , and its contribution towards the Shortest Path

Score from  $x$  to  $y$  should be removed.

At time  $t$ ,

$$SP(x,y|t) = \frac{(\sum_{(a,b) \in p} w(a,b|t)) + \kappa}{|\rho(x,y|t)| \cdot l_{min}(x,y|t)} \quad (157)$$

where  $\kappa$  is the contribution from the other paths that do not contain  $(u, v)$ .

At time  $(t + 1)$ , if  $(u, v)$  is broken, the path  $p$  does not exist anymore. So,

$$|\rho(x,y|t+1)| = |\rho(x,y|t)| - 1 \quad (158)$$

$$SP(x,y|t) = \frac{\kappa}{(|\rho(x,y|t)| - 1) \cdot l_{min}(x,y|t)} \quad (159)$$

$$= \delta \cdot \frac{|\rho(x,y|t)|}{|\rho(x,y|t)| - 1} \cdot \left( \frac{(\sum_{(a,b) \in p} w(a,b|t)) + \kappa}{|\rho(x,y|t)| \cdot l_{min}(x,y|t)} \right) \quad (160)$$

$$- \delta \cdot \frac{\sum_{(a,b) \in p} w(a,b|t)}{(|\rho(x,y|t)| - 1) \cdot l_{min}(x,y|t)}$$

$$\boxed{SP(x,y|t) = \delta \cdot \frac{|\rho(x,y|t)|}{|\rho(x,y|t)| - 1} \cdot SP(x,y|t) - \delta \cdot \frac{\sum_{(a,b) \in p} w(a,b|t)}{(|\rho(x,y|t)| - 1) \cdot l_{min}(x,y|t)}} \quad (161)$$

## 5.6 Rooted PageRank

Consider the Rooted PageRank algorithm given in 1. The teleportation probability  $\alpha$  determines the probability that the random walk with reach a certain length.

We can approximate the Rooted PageRank by considering only paths that has a probability of at least  $p$  of being reached by a random walk. Suppose that  $l_c$  is the maximum path length that has a probability  $p$  or more of being reached.

$$p \geq (1 - \alpha)^{l_c} \quad (162)$$

$$l_c \leq \left\lceil \frac{\log p}{\log(1 - \alpha)} \right\rceil \quad (163)$$

This means that if there is an interaction between  $(u, v)$ , it will effect the Rooted PageRank of nodes  $x$  and  $y$  if they are at a maximum distance of  $l_c$  away from  $u$  or  $v$ .

Work by Liben-Nowell et al [1] have shown that the optimum value of  $\alpha$  is 0.5. So, to approximate the Rooted PageRank by taking only paths of probability  $p = 0.1$  or greater,

$$l_c \leq \left\lceil \frac{\log(0.1)}{\log(1 - 0.5)} \right\rceil \quad (164)$$

$$l_c \leq \lfloor 3.21 \rfloor \quad (165)$$

$$l_c \leq 3 \quad (166)$$

### 5.6.1 Effect of decay of link weight

Let  $I$  be the set of node pairs that interact in time interval  $[t, t + 1]$ ,  $R$  be the set of node pairs that had a link at time  $t$  but is broken in time  $(t + 1)$ ,  $C$  be the set of all nodes in the set of node pairs  $(I \cup R)$  and  $V$  be the set of all node pairs.

In Algorithm 1, it can be seen that if the link weights are changed by the same factor, the value of the Rooted PageRank will not change. So, for any node pair  $(u, v)$  such that  $u \notin C \wedge v \notin C$ , the decay in link weight will not change the Rooted PageRank.

### 5.6.2 Effect of interactions between nodes and removal of links

Let  $D$  be the set of all node pairs that are within  $l_c$  distance from the nodes in  $C$ .

So, to update the Rooted PageRank, we need to only recalculate it for nodes in  $D$ . The algorithm for updating the Rooted PageRank is given in Algorithm 8.

## 6 Link Prediction as a Classification Problem

In this chapter, we discuss the problem of link prediction as a classification problem. We will also address the problem that arises naturally due to the imbalance between the negative and positive classes.

### 6.1 Link Prediction Method

As discussed earlier, the network observed over time are  $\langle V, E_0 \rangle, \langle V, E_1 \rangle, \dots, \langle V, E_T \rangle$ . However, we consider only three snapshots; that is we consider the networks  $G_1^* = \langle V, E_{t_1} \rangle$ ,  $G_2^* = \langle V, E_{t_2} \rangle$  and  $G_3^* = \langle V, E_{t_3} \rangle$  at time  $t_1, t_2$  and  $t_3$  such that  $t_1 < t_2 < t_3$ .<sup>1</sup>

Recall that a deficiency of the method described in Chapter 2.2 is that the network  $G_2$  cannot be observed in observed in real world cases, and consequently the feature set cannot be used to predict edges. For this reason, in our method, the features for the training data are generated from  $G_1^*$  and the class labels for the training data are generated from  $G_2^*$ . In other words, the training set for a classifier consists of  $(F_1, C_2)$ . The features for the testing data are generated from  $G_2^*$  to predict the links in  $G_3^*$ .

The training data and the testing data are generated as described below:

1. For all possible node pairs  $(u, v)$  the features for training are calculated from  $G_1^*$ .

$$F_{train} = \{(s_1(u, v|G_1^*), s_2(u, v|G_1^*), \dots, s_n(u, v|G_1^*)) : u, v \in V\} \quad (167)$$

where,  $s_f(u, v|G_t^*)$  is a score associated with node pair  $(u, v)$  in network  $G_t^*$  for a feature  $f$ .

2. The class label for the training data is calculated from  $G_2^*$ . A pair of nodes  $(u, v)$  is considered to be in the positive class if there is a link from  $u$  to  $v$  in  $G_2^*$ ; otherwise it is considered to be

---

<sup>1</sup>It should be noted here that  $t_1$  should not be very small. If  $t_1$  is very small, the missing past effect [13] will come into play and the learning data will be inaccurate.

in the negative class.

$$C_{train} = \{class(u, v | G_2^*) : u, v \in V\} \quad (168)$$

3. The feature set  $F_{test}$  is calculated for all possible node pairs  $(u, v)$  using  $G_2^*$ . This is the testing data.

$$F_{test} = \{(s_1(u, v | G_2^*), s_2(u, v | G_2^*), \dots, s_n(u, v | G_2^*)) : u, v \in V\} \quad (169)$$

There are two differences between our method and the earlier methods:

1. While the method by Al Hasan et al. [2] uses two networks,  $G_1, G_2$ , our method uses three ( $G_1^*, G_2^*$  and  $G_3^*$ ). This is to overcome the problem with the training network used by Al Hasan et al., where  $G_2$  will be unknown in a real world network, as described in Section 2.2. In our method  $G_3^*$  is the unknown and  $G_2^*$  is used to calculate the testing data.
2. In the earlier methods of link prediction (Chapter 2), the size of the testing data is always smaller than the training data because existing links are not included in the testing data. If the links are permanent, the node pairs that already have a link between them in the training network are not included in the testing data. However, in our method, since links are not permanent, the training data and the testing data have the same size.

## 6.2 Handling Class Imbalance

As mentioned in Section 2.3, in many real world networks, the size of the positive class is much smaller than that of the negative class. So, when we look at link prediction as a classification problem, the training data will be highly imbalanced.

Previous works done in classification of extremely imbalanced data [14, 15, 16, 17] have shown that good results can be obtained using an ensemble of SVM classifiers trained with the minority class and different samples of the majority class. We implement the proposed ensemble approach as described below:

1. Suppose the training data is  $D_{train}$ . Consider its two components  $D_+$  and  $D_-$  associated with the positive and negative class respectively, where  $|D_-| \gg |D_+|$ .
2. To use an ensemble of  $m$  SVM classifiers,  $n$  samples  $[S_1, S_2, \dots, S_n]$  each of size  $|D_+|$ , are taken from  $D_-$  without replacement. In our simulations,  $m = 10$ .
3. For  $i = 1, 2, \dots, m$ , we train the classifier  $C_i$  using the sample  $S_i$  and  $D_+$ , and class membership represented by 1 for  $D_+$  and 0 for  $S_i$ .
4. Let  $P_i$  consists of the predictions by classifier  $C_i$ ,  $i = 1, 2, \dots, m$ .
5. The predictions  $P_1, P_2, \dots, P_n$  are combined using majority voting to generate the final prediction  $P$ .

$$P(u, v) = \begin{cases} 1 & \text{if } (\sum_{i=1}^m P_i(u, v)) > n/2 \\ 0 & \text{otherwise} \end{cases} \quad (170)$$

## 7 Data

In this chapter, we describe the data used for simulation and testing, and the pre-processing done to the data.

### 7.1 Data Source

In this thesis, the data used for simulation is the Twitter @-mention network. This network is obtained from the Twitter Streaming API [6] by observing it for 28 days. It contains 91169 nodes and 13087872 interactions between node pairs over the 28 days.

In this network, nodes represent users and there is a link from node  $u$  to  $v$  if  $u$  sends a tweet mentioning  $v$ . The links in this network are not permanent and they are directed. In the Twitter network, the size of the positive class is much smaller than that of the negative class.

The dataset consists of networks observed every hour  $\langle V, E_1 \rangle, \langle V, E_2 \rangle, \dots, \langle V, E_{672} \rangle$  which are generated using the network model described in Chapter 3. In this network the values of  $\theta$  and  $\delta$  are 0.2 and 0.9 per hour respectively.

Spammers and bots are removed from this network by using the Local Outlier Factor density-based outlier detection [20] using the in-degree and out-degree as features.

### 7.2 Data Preparation

As mentioned in Section 7.1, the Twitter data has approximately  $10^5$  nodes. So, there are approximately  $10^{10}$  possible links.

Because of the size of the Twitter network, it is computationally very expensive to perform link prediction over the entire network. So five sets of sample nodes  $V_1, V_2, \dots, V_5$  are taken from the set

$V$ , each containing 500 nodes. The five samples are taken from areas of the network with different densities.<sup>2</sup>

Sample	$ E_{t_1} $	$ E_{t_2} $	$ E_{t_3} $	Density at $t_3$
$V_1$	729	810	824	0.0033
$V_2$	780	893	973	0.0039
$V_3$	1251	1314	1298	0.0056
$V_4$	1407	1463	1497	0.0060
$V_5$	1564	1768	1699	0.0068

Table 3: Sizes of the sample networks

To apply the link prediction method described in Section 6.1,  $t_1 = 336$ ,  $t_2 = 504$  and  $t_3 = 672$ . The observations were made every hour, so the value of  $t_i$  is in hours. Thus we are using the networks observed on day 14 and day 21 to predict the links in the network on day 28.

---

<sup>2</sup>The density of a set of nodes is defined as the ratio of the number of connections between the nodes to the number of potential connections. The density of a set of nodes  $V_i$  at time  $t$  will be given by,

$$Density(V_i|t) = \frac{|\{(u,v) : (u,v \in V_i) \wedge ((u,v) \in E_t)\}|}{|V_i| \cdot (|V_i| - 1)} \quad (171)$$

## 8 Link Prediction Results

In this chapter, we present the results of from the simulation done using the data described in Chapter 7. As mentioned in Section 7.1, the positive class is much smaller than the negative class. So, the ensemble method described in Section 6.2 is used.

The link prediction (Section 6.1) is performed on the five sets of nodes mentioned in Section 7.2. Let  $P_1, P_2, P_3, P_4$  and  $P_5$  represent the predictions for node sets  $V_1, V_2, V_3, V_4$  and  $V_5$ .

Work by Fawcett [7] suggests that for datasets with imbalanced classes, the area under curve (AUC) of the receiver operating curve (ROC) is a more accurate measure of performance than accuracy and recall. So the predictions made using our method ( $P_i$ ) is compared to the baseline prediction ( $Q_i$ ) using the AUC of the ROC.

To provide a baseline to compare for comparison with our method, the Twitter network is assumed to be unweighted and undirected, and the link prediction method described by Leichtenwalter et al. [3] is applied. The baseline link prediction is performed for node sets  $V_1, V_2, V_3, V_4$  and  $V_5$ . Let  $Q_1, Q_2, Q_3, Q_4$  and  $Q_5$  represent these predictions.

The comparison between the AUC of  $P_i$  and  $Q_i$  is given in Table 4. The comparison of the ROC for  $V_1, V_2, V_3, V_4$  and  $V_5$  are shown in Figures 5, 6, 7, 8 and 9.

Sample	AUC of $Q_i$	AUC of $P_i$
$V_1$	0.52	0.77
$V_2$	0.73	0.79
$V_3$	0.74	0.82
$V_4$	0.78	0.83
$V_5$	0.83	0.87

Table 4: AUC comparison between prediction using our method  $P_i$  and the baseline method  $Q_i$

Sample	Number of correct predictions in $Q_i$	Number of correct predictions in $P_i$
$V_1$	190	577
$V_2$	195	632
$V_3$	844	1064
$V_4$	1168	1272
$V_5$	1393	1529

Table 5: Number of correct predictions comparison between prediction using our method  $P_i$  and the baseline method  $Q_i$

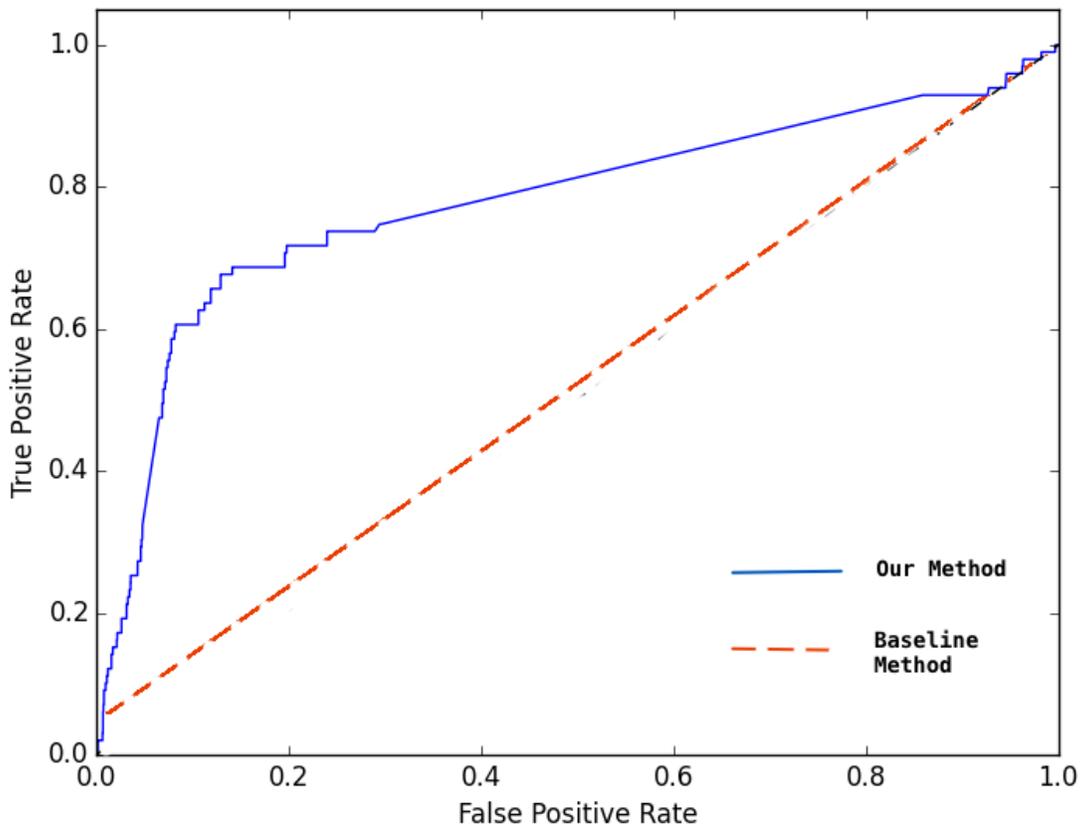


Figure 5: ROC for  $V_1$

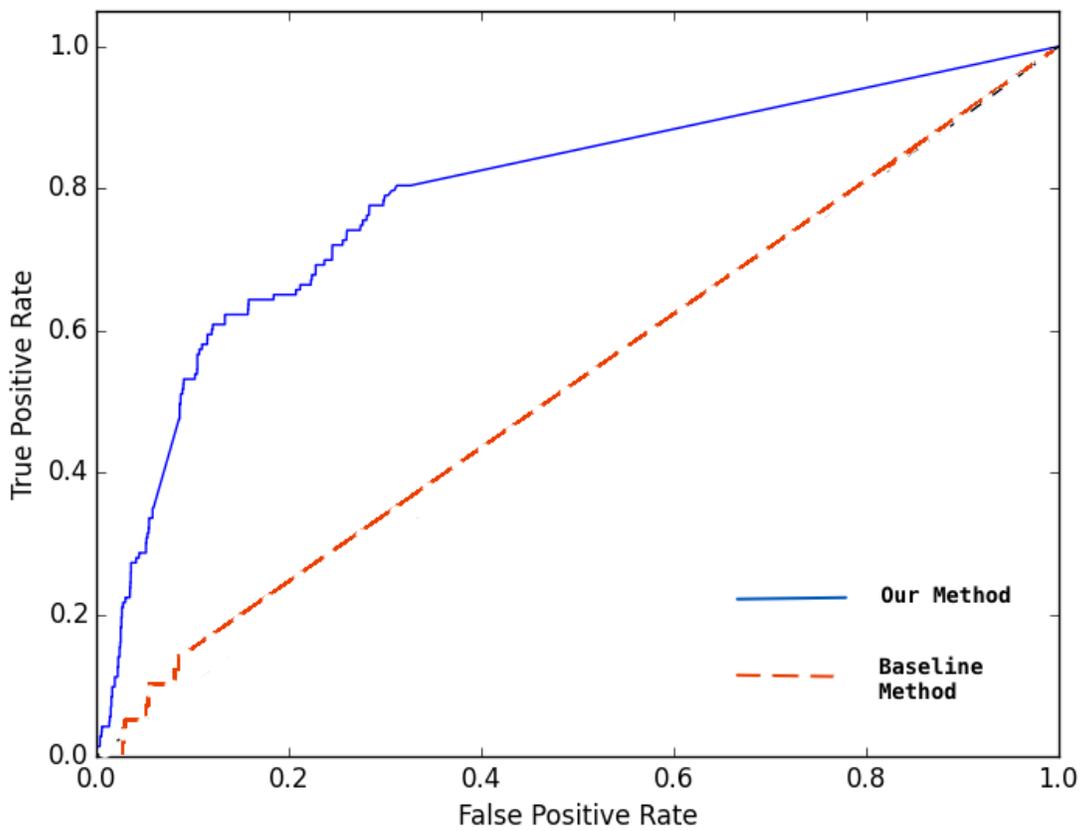


Figure 6: ROC for  $V_2$

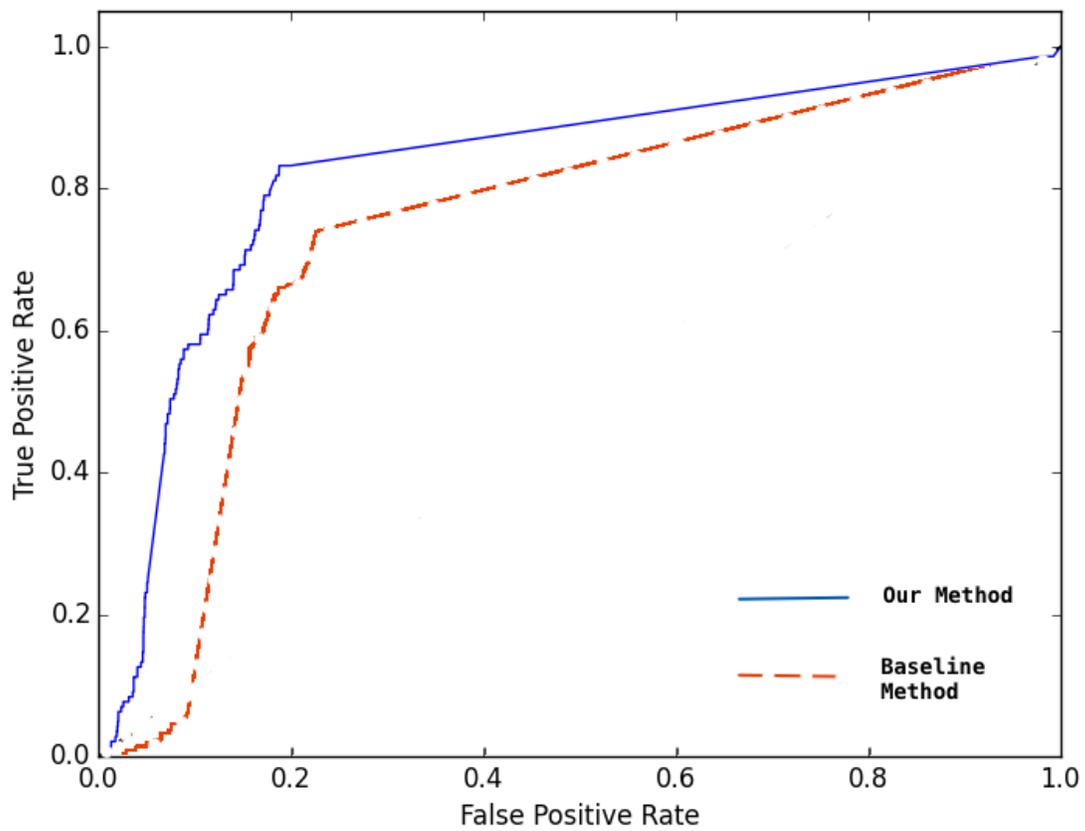


Figure 7: ROC for  $V_3$

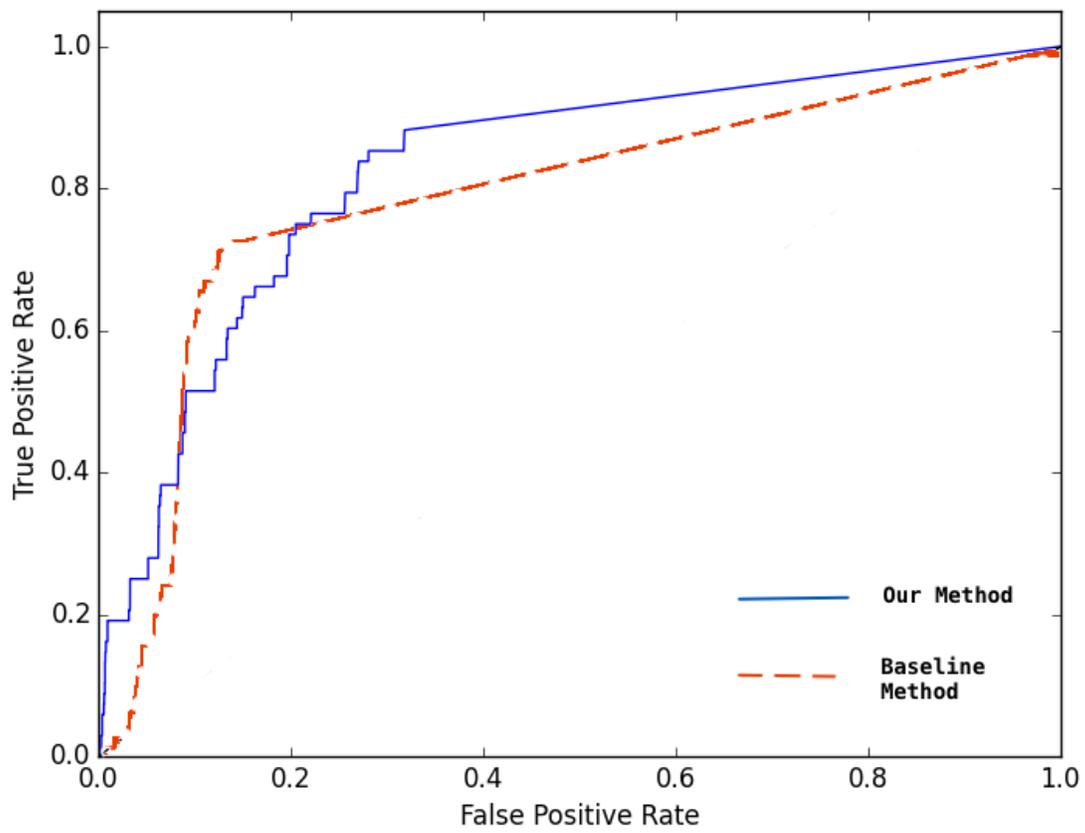


Figure 8: ROC for  $V_4$

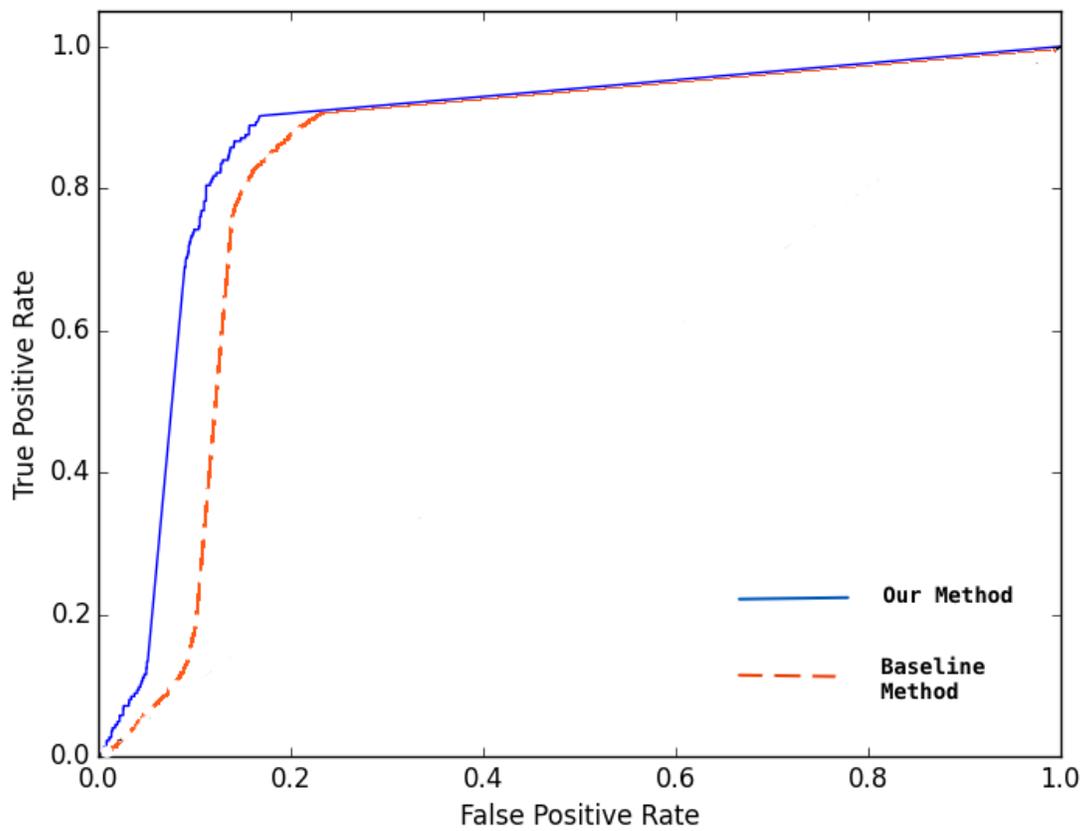


Figure 9: ROC for  $V_5$

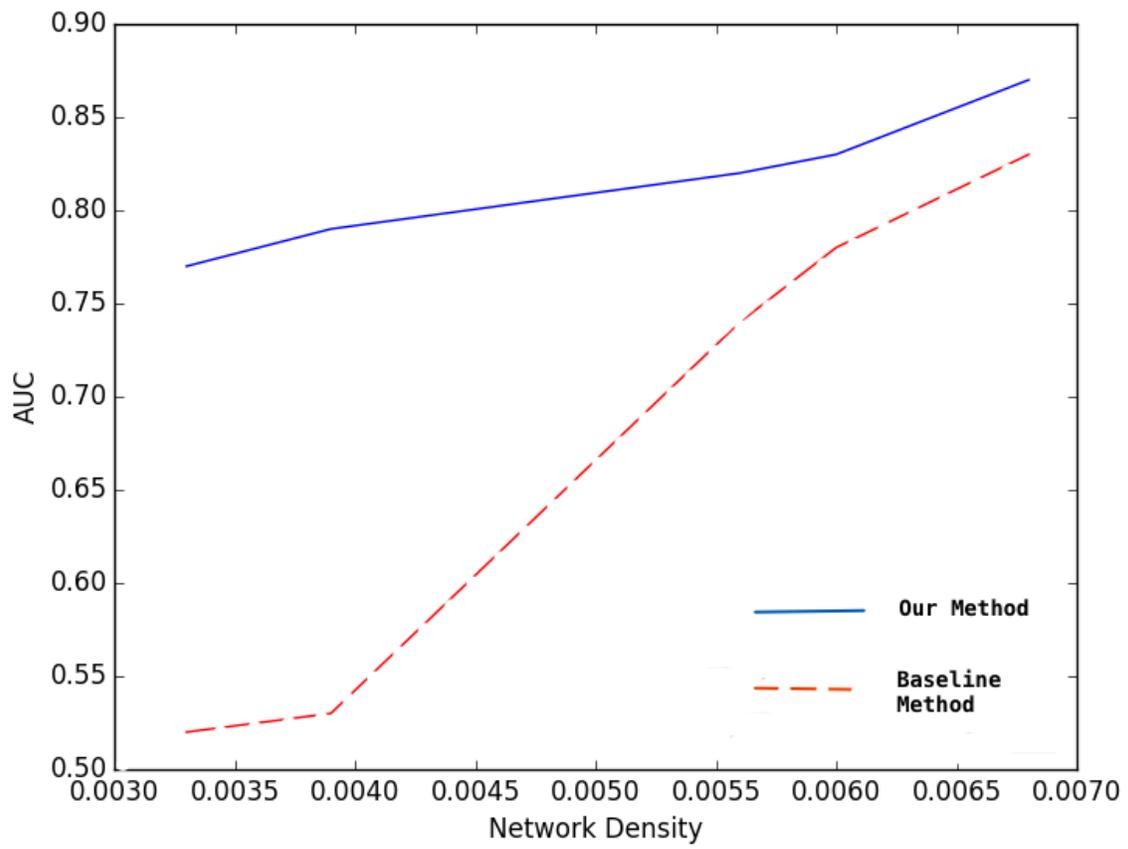


Figure 10: Variation of AUC with Density

## 9 Analysis

The results in Chapter 8 show that our method performs better than the baseline method in all cases. Our method takes the link weight and directions into consideration, but the baseline method assumes that the network is undirected and unweighted.

The graph in Figure 10 shows the comparison of AUC for our method and the baseline method against the density of the sample nodes. It can be seen that the performance of our method is much better than the baseline method if the density is lower, i.e., the network is sparse.

## 10 Conclusion

When we consider networks in which the links are not permanent, the method of link prediction described in this thesis can achieve performance better over other supervised learning methods. While the improvement over the baseline method is around 3% AUC for a dense network, the difference is very substantial (around 25% AUC) when the network is very sparse.

As mentioned in Chapter 2, most of the real world networks are sparse, and existing methods do not take into consideration links that can be removed.

The features (4) for classification used in this thesis depend on only the network topology. This means that this method can be applied easily to other type of networks.

# Appendices

## A Algorithms for Calculating Feature Scores

In Chapter 5, methods to incrementally update the feature scores are described. In this section, we describe the algorithms to calculate the feature scores by using the methods described in Chapter 5.

### A.1 Common Neighbor Score

Let  $I$  be the set of node pairs that interact with each other in time interval  $[t, t + 1]$  and let  $R$  be the set of node pairs whose links are broken in time interval  $[t, t + 1]$ .

Equations 44, 45, 49, 50, 57, 60, 61 and 62 are used to develop an algorithm to iteratively update the common neighbor score for all possible node pairs. This is given in Algorithm 2.

### A.2 Adamic-Adar Index

Let us assume that  $I$  is the set of node pairs that interact in time interval  $[t, t + 1]$ ,  $R$  be the set of node pairs which had a link at time  $t$  but is removed at time  $(t + 1)$  and let  $S$  be the set of all possible node pairs. The algorithm for updating the Adamic Adar Index is given in Algorithm 3 and 4.

### A.3 Preferential Attachment Score

Let  $I$  be the set of node pairs that interact in time interval  $[t, t + 1]$ ,  $R$  be the set of node pairs which had a link at time  $t$  but do not at time  $(t + 1)$ ,  $N$  be the set of all nodes,  $S$  be the set of all possible

---

**Algorithm 2** Common Neighbor Score Update Algorithm

---

```
1: procedure COMMONNEIGHBORSOREUPDATE(I)
2:   for all  $(u, v) \in S$  do                                     ▷ Update all scores for decay
3:      $CN((u, v)|t + 1) = \delta \cdot CN((u, v)|t)$ 
4:   end for
5:   for all  $(u, v) \in I$  do                                     ▷ Update for  $I$ 
6:     for all  $x \in (\Gamma_i(v|t + 1) - \{u\})$  do               ▷ Update  $CN_o$ 
7:       if  $(u, v) \notin E_t$  then
8:          $CN_o(u, x|t + 1) = CN_o(u, x|t + 1) + 0.5(1 + w(x, v|t + 1))$ 
9:       else
10:         $CN_o(u, x|t + 1) = CN_o(u, x|t + 1) + 0.5$ 
11:      end if
12:    end for
13:    for all  $x \in (\Gamma_o(u|t + 1) - \{v\})$  do                 ▷ Update  $CN_i$ 
14:      if  $(u, v) \notin E_t$  then
15:         $CN_i(x, v|t + 1) = CN_i(x, v|t + 1) + 0.5(1 + w(u, x|t + 1))$ 
16:      else
17:         $CN_i(x, v|t + 1) = CN_i(x, v|t + 1) + 0.5$ 
18:      end if
19:    end for
20:  end for
21:  for all  $(u, v) \in R$  do                                     ▷ Update for  $R$ 
22:    for all  $x \in (\Gamma_i(v|t) - \{u\})$  do                   ▷ Update  $CN_o$ 
23:       $CN_o(u, x|t + 1) = CN_o(u, x|t + 1) - 0.5 \cdot \delta((w(u, v|t) + w(x, v|t)))$ 
24:    end for
25:    for all  $x \in (\Gamma_o(u|t) - \{v\})$  do                   ▷ Update  $CN_i$ 
26:       $CN_i(x, v|t + 1) = CN_i(x, v|t + 1) - 0.5 \cdot \delta(w(u, v|t) + w(u, x|t))$ 
27:    end for
28:  end for
29: end procedure
```

---

---

**Algorithm 3** out-Adamic Adar Index Update Algorithm
 

---

```

1: procedure OUTADAMICADARUPDATE(I)
2:   for all  $(u, v) \in I$  do ▷ Update for interactions
3:     if  $(u, v) \in E_t$  then
4:       for all  $x \in \Gamma_i(v|t+1)$  do
5:          $AA_o(u, x|t+1) = \delta \cdot AA_o(u, x|t) + \frac{1}{2 \cdot \log(\sum_{z \in \Gamma_o(v|t+1)} w(v, z|t+1))}$ 
6:       end for
7:     else
8:       for all  $x \in \Gamma_i(v|t+1)$  do
9:          $AA_o(u, x|t+1) = \delta \cdot AA_o(u, x|t) + \frac{1 + \delta \cdot w(a, v|t)}{2 \cdot \log(\sum_{z \in \Gamma_o(v|t+1)} w(v, z|t+1))}$ 
10:      end for
11:    end if
12:    for all  $x, y \in \Gamma_i(u|t+1)$  do
13:       $AA_o(x, y|t+1) = \delta \cdot \left( \frac{w(x, u|t) + w(y, u|t)}{\log(\sum_{z \in \Gamma_o(u|t+1)} w(u, z|t+1))} - \frac{w(x, u|t) + w(y, u|t)}{2 \cdot \log(\sum_{z \in \Gamma_o(u|t)} w(u, z|t))} \right) + \delta \cdot AA_o(a, b|t)$ 
14:    end for
15:  end for
16:  for all  $(u, v) \in R$  do ▷ Update for removed links
17:    for all  $x \in \Gamma_i(v|t+1)$  do
18:       $AA_o(u, x|t+1) = \delta \cdot AA_o(u, x|t) - \frac{\delta \cdot (w(u, v) + w(x, v|t))}{2 \cdot \log(\sum_{z \in \Gamma_o(v|t)} w(v, z|t))}$ 
19:    end for
20:    for all  $x, y \in \Gamma_i(u|t+1)$  do
21:       $AA_o(a, y|t+1) = \delta \cdot AA_o(x, y|t) - \frac{\delta \cdot (w(x, u|t) + w(y, u|t))}{2 \cdot \log(\sum_{z \in \Gamma_o(u|t)} w(u, z|t))} + \frac{\delta \cdot (w(x, u|t) + w(y, u|t))}{2 \cdot \log(\sum_{z \in \Gamma_o(u|t+1)} w(u, z|t+1))}$ 
22:    end for
23:  end for
24:  for all  $(u, v) \in S$  do ▷ Update all the remaining node pairs
25:    if  $AA_o(u, v|t+1)$  do not exist then
26:       $AA_o(u, v|t+1) = \delta \cdot AA_o(u, v|t)$ 
27:    end if
28:  end for
29: end procedure

```

---

---

**Algorithm 4** in-Adamic Adar Index Update Algorithm
 

---

```

1: procedure INADAMICADARUPDATE(I)
2:   for all  $(u, v) \in I$  do ▷ Update for interactions
3:     if  $(u, v) \in E_t$  then
4:       for all  $x \in \Gamma_o(u|t+1)$  do
5:          $AA_i(v, x|t+1) = \delta \cdot AA_i(v, x|t) + \frac{1}{2 \cdot \log(\sum_{z \in \Gamma_i^-(u|t+1)} w(z, u|t+1))}$ 
6:       end for
7:     else
8:       for all  $x \in \Gamma_o(u|t+1)$  do
9:          $AA_i(v, x|t+1) = \delta \cdot AA_i(v, x|t) + \frac{1 + \delta \cdot w(u, a|t)}{2 \cdot \log(\sum_{z \in \Gamma_i^-(u|t+1)} w(z, u|t+1))}$ 
10:      end for
11:    end if
12:    for all  $a, b \in \Gamma_o(v|t+1)$  do
13:       $AA_i(a, b|t+1) = \delta \cdot AA_i(a, b|t) + \frac{\delta \cdot (w(v, a|t) + w(v, b|t))}{\log(\sum_{z \in \Gamma_i^-(v|t+1)} w(z, v|t+1))} - \frac{\delta \cdot (w(v, a|t) + w(v, b|t))}{2 \cdot \log(\sum_{z \in \Gamma_i^-(v|t)} w(z, v|t))}$ 
14:    end for
15:  end for
16:  for all  $(u, v) \in R$  do ▷ Update for removed links
17:    for all  $x \in \Gamma_o(u|t+1)$  do
18:       $AA_i(v, x|t+1) = \delta \cdot AA_i(v, x|t) - \frac{\delta \cdot (w(u, v) + w(v, x|t))}{2 \cdot \log(\sum_{z \in \Gamma_i^-(u|t)} w(z, u|t))}$ 
19:    end for
20:    for all  $x, y \in \Gamma_o(v|t+1)$  do
21:       $AA_i(a, b|t+1) = \delta \cdot AA_i(a, b|t) - \frac{\delta \cdot (w(v, a|t) + w(v, b|t))}{2 \cdot \log(\sum_{z \in \Gamma_i^-(v|t)} w(z, v|t))} + \frac{\delta \cdot (w(v, a|t) + w(v, b|t))}{2 \cdot \log(\sum_{z \in \Gamma_i^-(v|t+1)} w(z, v|t+1))}$ 
22:    end for
23:  end for
24:  for all  $(u, v) \in S$  do ▷ Update all the remaining node pairs
25:    if  $AA_i(u, v|t+1)$  do not exist then
26:       $AA_i(u, v|t+1) = \delta \cdot AA_i(u, v|t)$ 
27:    end if
28:  end for
29: end procedure

```

---

node pairs and  $SUM_o(x|t)$  and  $SUM_i(x|t)$  be the sum of the weights of the out-links and in-links of node  $x$ . The algorithm for updating the Preferential Attachment Score is given in Algorithm 5.

---

**Algorithm 5** Preferential Attachment Score Update Algorithm

---

```

1: procedure PREFERENTIALATTACHMENTUPDATE(I)
2:   for all  $x \in N$  do
3:     for all  $(u, v) \in I$  do ▷ Update for node interactions
4:       if  $(u, v) \in E_t$  then
5:          $PF_o(u, x|t+1) = \delta^2 \cdot PF_o(u, x|t) + \frac{\delta \cdot SUM_o(x|t)}{|\Gamma_o(u|t)| \cdot |\Gamma_o(x|t)|}$ 
6:          $PF_i(v, x|t+1) = \delta^2 \cdot PF_i(v, x|t) + \frac{\delta \cdot SUM_i(x|t)}{|\Gamma_i(v|t)| \cdot |\Gamma_i(x|t)|}$ 
7:       else
8:          $PF_o(u, x|t+1) = \frac{|\Gamma_o(u|t)|}{|\Gamma_u(u|t)|+1} \cdot \delta^2 \cdot PF_o(u, x|t) + \frac{\delta \cdot SUM_o(x|t)}{(|\Gamma_o(u|t)|+1) \cdot |\Gamma_o(x|t)|}$ 
9:          $PF_i(v, x|t+1) = \frac{|\Gamma_i(v|t)|}{|\Gamma_i(v|t)|+1} \cdot \delta^2 \cdot PF_i(v, x|t) + \frac{\delta \cdot SUM_i(x|t)}{(|\Gamma_i(v|t)|+1) \cdot |\Gamma_i(x|t)|}$ 
10:      end if
11:    end for
12:    for all  $(u, v) \in R$  do ▷ Update for links removed
13:       $PF_o(u, x|t+1) = \frac{\delta^2 \cdot |\Gamma_o(u|t)| \cdot PF_o(u, x|t)}{|\Gamma_o(u|t)-1} - \frac{\delta \cdot w(u, v|t) \cdot SUM_o(x|t)}{(|\Gamma_o(u|t)-1) \cdot |\Gamma_o(x|t)|}$ 
14:       $PF_i(v, x|t+1) = \frac{\delta^2 \cdot |\Gamma_i(v|t)| \cdot PF_i(v, x|t)}{|\Gamma_i(v|t)-1} - \frac{\delta \cdot w(u, v|t) \cdot SUM_i(x|t)}{(|\Gamma_i(v|t)-1) \cdot |\Gamma_i(x|t)|}$ 
15:    end for
16:  end for
17:  for all  $(u, v) \in S$  do ▷ Update all the remaining node pairs
18:    if  $PF_i(u, v|t+1)$  do not exist then
19:       $PF_i(u, v|t+1) = \delta^2 \cdot PF_i(u, v|t)$ 
20:    end if
21:    if  $PF_o(u, v|t+1)$  do not exist then
22:       $PF_o(u, v|t+1) = \delta^2 \cdot PF_o(u, v|t)$ 
23:    end if
24:  end for
25: end procedure

```

---

## A.4 Katz Score

Let  $I$  be the set of node pairs that interact in time interval  $[t, t+1]$ ,  $R$  be the set of node pairs that had a link between them at time  $t$  but not at  $(t+1)$  and  $S$  be the set of all possible node pairs. The algorithm for updating the Katz Score is given in Algorithm 6.

---

**Algorithm 6** Katz Score Update Algorithm

---

```
1: procedure KATZUPDATE(I)
2:   for all  $(u, v) \in I$  do ▷ Update for interaction
3:     if  $(u, v) \in E_t$  then
4:        $KS(u, v|t+1) = \delta \cdot KS(u, v|t) + \beta$ 
5:       for all  $x \in \Gamma_o(v)$  do
6:          $KS(u, x|t+1) = \delta \cdot KS(u, x|t) + 0.5 \cdot \beta^2$ 
7:       end for
8:       for all  $x \in \Gamma_i(u)$  do
9:          $KS(x, v|t+1) = \delta \cdot KS(x, v|t) + 0.5 \cdot \beta^2$ 
10:      end for
11:      for all  $y \in \Gamma_o(v)$  do
12:        for all  $x \in \Gamma_i(u)$  do
13:           $KS(x, y|t+1) = \delta \cdot KS(x, y|t) + \frac{\beta^2}{3}$ 
14:        end for
15:      end for
16:    else
17:       $KS(u, v|t+1) = \delta \cdot KS(u, v|t) + \beta$ 
18:      for all  $x \in \Gamma_o(v)$  do
19:         $KS(u, x|t+1) = \delta \cdot KS(u, x|t) + \beta^2 \cdot \frac{1+\delta \cdot w(v, x|t)}{2}$ 
20:      end for
21:      for all  $x \in \Gamma_i(u)$  do
22:         $KS(x, v|t+1) = \delta \cdot KS(x, v|t) + \beta^2 \cdot \frac{1+\delta \cdot w(x, u|t)}{2}$ 
23:      end for
24:      for all  $y \in \Gamma_o(v)$  do
25:        for all  $x \in \Gamma_i(u)$  do
26:           $KS(x, y|t+1) = \delta \cdot KS(x, y|t) + \beta^3 \cdot \frac{1+\delta \cdot w(x, u)+\delta \cdot w(v, y)}{3}$ 
27:        end for
28:      end for
29:    end if
30:  end for
```

---

---

```

31:   for all  $(u, v) \in R$  do ▷ Update for removal of links
32:      $KS(u, v|t+1) = KS(u, v|t) - \delta \cdot \beta w(u, v|t)$ 
33:     for all  $x \in \Gamma_o(v)$  do
34:        $KS(u, x|t+1) = KS(u, x|t) - \delta \cdot \beta^2 \cdot \frac{w(u, v|t) + w(v, x|t)}{2}$ 
35:     end for
36:     for all  $x \in \Gamma_i(u)$  do
37:        $KS(x, v|t+1) = KS(x, v|t) - \delta \cdot \beta^2 \cdot \frac{w(u, v|t) + w(x, u|t)}{2}$ 
38:     end for
39:     for all  $y \in \Gamma_o(v)$  do
40:       for all  $x \in \Gamma_i(u)$  do
41:          $KS(x, y|t+1) = KS(x, y|t) - \delta \cdot \beta^3 \cdot \frac{w(u, v|t) + w(x, u|t) + w(v, y|t)}{3}$ 
42:       end for
43:     end for
44:   end for
45:   for all  $(u, v) \in S$  do ▷ Update remaining node pairs
46:     if  $KS(u, v|t+1)$  do not exist then
47:        $KS(u, v|t+1) = \delta \cdot KS(u, v|t)$ 
48:     end if
49:   end for
50: end procedure

```

---

## A.5 Shortest Path Score

Let us assume that  $P_s(a)$  is the set of all nodes  $b$  such that there is a valid shortest path from  $a$  to  $b$ . Also suppose that  $P_e(a)$  is the set of all nodes  $b$  such that there is a valid shortest path from  $b$  to  $a$ . Let  $Q(u, v)$  be the set of all node pairs whose shortest paths passes through  $(u, v)$ .

Let  $I$  be the set of node pairs that interact with each other in time interval  $[t, t+1]$ ,  $R$  be the set of all node pairs which had a link at time  $t$  but not at time  $(t+1)$  and  $S$  be the set of all possible node pairs. The algorithm to update the Shortest Path Score is given in Algorithm 7.

## A.6 Rooted PageRank

To update the Rooted PageRank we only need to recalculate the Rooted PageRank for the nodes that are within a distance of  $l_c$  from the nodes which interacted or which had links that are broken.

---

**Algorithm 7** Shortest Path Score Update Algorithm

---

```
1: procedure SHORTESTPATHUPDATE(I,R)
2:   for all  $(u, v) \in I$  do ▷ Update for node interactions
3:     for all  $x \in P_e(u)$  do
4:       for all  $y \in P_s(v)$  do
5:         if  $(u, v) \in E_t$  then
6:            $SP(x, y|t+1) = \delta \cdot SP(x, y|t) + \frac{1}{|\rho_{min}(x, y|t)| \cdot l_{min}(u, v|t)}$ 
7:         else
8:           if  $l_{min}(x, y|t) \geq (l_{min}(x, u|t) + l_{min}(v, y|t) + 1)$  then
9:              $l_{min}(x, y|t+1) \geq l_{min}(x, u|t) + l_{min}(v, y|t) + 1$ 
10:             $\rho_{min}(x, y|t+1) = \{\rho_{min}(x, v|t) + (u, v) + \rho_{min}(v, y|t)\}$ 
11:            Recalculate  $SP(x, y|t+1)$  using  $\rho_{min}(x, y|t+1)$ 
12:            Update  $Q(x, y)$ 
13:          end if
14:        end if
15:      end for
16:    end for
17:  end for
18:  for all  $(u, v) \in I$  do ▷ Update remaining node pairs
19:    if  $SP(u, v|t+1)$  do not exist then
20:       $SP(u, v|t+1) = \delta \cdot SP(u, v|t)$ 
21:    end if
22:  end for
23:  for all  $(u, v) \in R$  do ▷ Update for link removal
24:    for all  $(x, y) \in Q(u, v)$  do
25:       $SP(x, y|t+1) = \frac{|\rho(x, y|t)|}{|\rho(x, y|t)|-1} \cdot SP(x, y|t+1)$ 
26:      for all  $p \in \rho_{min}(x, y|t)$  do
27:        if  $(u, v) \in p$  then
28:           $SP(x, y|t+1) = SP(x, y|t+1) - \delta \cdot \frac{\sum_{(a,b) \in p} w(a, b|t)}{(|\rho(x, y|t)|-1) \cdot l_{min}(x, y|t)}$ 
29:           $\rho(x, y|t+1) = \rho(x, y|t) - \{p\}$ 
30:        end if
31:      end for
32:      Update  $Q(x, y) = Q(x, y) - \{(x, y)\}$ 
33:    end for
34:  end for
35: end procedure
```

---

The algorithm to recalculate the Rooted PageRank is given in Algorithm 8

---

**Algorithm 8** Rooted PageRank Update Algorithm

---

```

1: procedure PAGERANKUPDATE(I)
2:   for all  $u \in N$  do
3:     for all  $v \in D$  do
4:        $PR(u, v|t + 1) = WeighedPageRank(u, v)$ 
5:        $PR(v, u|t + 1) = WeighedPageRank(v, v)$ 
6:     end for
7:   end for
8: end procedure

```

---

## B Run Time Analysis

### B.1 Common Neighbor Score

Let us assume  $n$  is the number nodes,  $i$  is the average number of node pairs that interact in time interval  $[t, t + 1]$ ,  $r$  is the number of links that were broken on  $[t, t + 1]$  and  $c$  is the average number of common neighbors between two nodes.

Then the time it takes to calculate the common neighbor scores for all possible node pairs using equation 12 and 13 is  $(c \cdot n^2)$ , and the time it takes to update the Common Neighbor Scores for all possible node pairs using Algorithm 2 is  $n^2 + 2 \cdot c \cdot (i + r)$ . For most real world networks,  $n^2 \gg (i + r)$  and  $n \gg c$ . So,  $n^2 \gg 2 \cdot c \cdot (i + r)$ . The running time for Algorithm 2 can then be approximated as  $O(n^2)$ .

### B.2 Adamic-Adar Index

Let us assume  $n$  is the number nodes,  $i$  is the number of node pairs that interact in time interval  $[t, t + 1]$ ,  $r$  is the number of links that were broken in time interval  $[t, t + 1]$  and  $g$  is the average

number of neighbors for a node.

Then the time it takes to calculate the Adamic Adar Index for all possible node pairs using equation 25 and 26 is  $(gcn^2)$ .

If we use Algorithm 3 and 4, the time it takes to update the Adamic Adar Index for all possible node pairs is  $(n^2 + (i + r) \cdot (g + g^2))$ . For a sparse network,  $g$  is small, and we can choose the time intervals such that the amount of changes in the network  $(i + r)$  is not very high. Then,  $n^2 \gg (i + r)(g + g^2)$ . So, the time it will take to run Algorithm 3 and 4 can be approximated as  $O(n^2)$ .

### B.3 Preferential Attachment Score

Let us assume  $n$  is the number nodes,  $g$  is the average number of neighbors for a node,  $r$  is the number of links that were broken in in time interval  $[t, t + 1]$  and  $i$  is the number of node pairs that interact in time interval  $[t, t + 1]$ .

Then the time it takes to calculate the Preferential Attachment Score for all possible node pairs using equation 17 and 18 is  $(2 \cdot g \cdot n^2)$ .

If we use Algorithm 5, the time it takes to update the Preferential Attachment Score for all possible node pairs is  $n^2 + 2 \cdot n \cdot g \cdot (i + r)$ . The interval can be chosen such  $n > (i + r)$ . For a sparse network,  $n \gg g$ . So, the time it will take to run Algorithm 5 can be approximated as  $O(n^2)$ .

### B.4 Katz Score

Let us assume  $n$  is the number nodes,  $g$  is the average number of neighbors for a node,  $r$  is the number of node pairs that had a link between them in time  $t$  but not in  $(t + 1)$  and  $i$  is the average number of node pairs that interact in time interval  $[t, t + 1]$ .

Then the time it takes to calculate the Katz Score for all possible node pairs using equation 35 and paths of length three or less is  $(g^3 \cdot n^2)$ .

If we use Algorithm 6, the time it takes to update the Katz Score for all possible node pairs is  $n^2 + (i + r) \cdot (2 \cdot g + g^2)$ .

The time interval can be chosen so that  $n \gg (i + r)$ . And for a sparse network,  $n \gg g$ . So, for a real world network, the time  $n^2 \gg (i + r) \cdot (2 \cdot g + g^2)$ . So, the it will take to run Algorithm 6 can be approximated as  $O(n^2)$ .

## B.5 Shortest Path Score

Let us assume  $n$  is the number nodes,  $p$  is the average number of nodes that can be reached from each node,  $c$  is the average number of paths of shortest length between two nodes,  $r$  be the number of node pairs that had a link at time  $t$  but not at  $(t + 1)$  and  $i$  is the average number of node pairs that interact in time interval  $[t, t + 1]$ .

Then the time it takes to calculate the Shortest Path Score for all possible node pairs using equation 39 and paths of length three or less is  $c \cdot n^3$ .

If we use Algorithm 6, the time it takes to update the Shortest Path Score for all possible node pairs is  $(n^2 + (i + r) \cdot p^2)$ . The time interval can be chosen such that  $(i + r)$  is small. And in a sparse network,  $(n \gg p)$ . So, in a real world network,  $(n^2 \gg (i + r) \cdot p^2)$ . So, the running time of the algorithm can be approximated as  $O(n^2)$ .

## B.6 Rooted PageRank

Since the algorithm for calculating the Rooted PageRank between two nodes given in Algorithm 1 depends on random walk, it is difficult to calculate the running time for this algorithm. So, let us

assume that this running time is  $f$ .

If we are to calculate the Rooted PageRank for all node pairs at each time interval, the time it will take is  $n^2 \cdot f$ .

If we calculate the Rooted PageRank for only the nodes that are in the set  $D$ . To calculate the size of the set  $D$ , let us assume that each node has  $m$  neighbors on average. Then the size of the set  $D$  is,

$$|D| < 2(i+r) \cdot (m^{l_c} + m^{l_c-1} + \dots + 1) \quad (172)$$

$$|D| < 2(i+r) \cdot \frac{(m-1)^{l_c+1} - 1}{m-2} \quad (173)$$

Then the approximate time it will take to update the Rooted PageRank using Algorithm 8 is,  $|D| \cdot n \cdot f$ .

Since  $|D| \leq n$ , Algorithm 8 will run in time less than  $n^2 \cdot f$ .

In most sparse networks, the average number of neighbors is low, and if we select  $l_c \leq 3$ ,  $|D| \ll n$ . So in real world networks, running time of the Algorithm 8 will be much lower than  $n^2 \cdot f$ .

## References

- [1] David Liben-Nowell, Jon Kleinberg "The Link Prediction Problem in Social Networks" *Journal of the American Society for Information Science and Technology*, 2007.
- [2] Mohammad Al Hasan, Vineet Chaoji, Saeed Salem, Mohammed Zaki "Link Prediction using Supervised Learning" *SIAM International Conference on Data Mining*, 2011.
- [3] Ryan N Lichtenwalter, Jake T Lussier, Nitesh V Chawla "New Perspectives and Methods in Link Prediction" *Proceeding of the 16th ACM*, 2010
- [4] Julia Preusse, Jörn Kleinberg, Matthias Thimm, Sergej Sizov "DecLiNe - Models for Decay of Links in Networks", 2014
- [5] Sitaram Asur, Bernardo A. Huberman, Gabor Szabo, Chunyan Wang "Trends in Social Media : Persistence and Decay" *Fifth International AAAI Conference on Weblogs and Social Media*, 2011
- [6] Twitter Streaming API <https://dev.twitter.com/streaming/overview>
- [7] Tom Fawcett "An Introduction to ROC Analysis" *Pattern Recognition Letters*, 2006
- [8] M E J Newman "Clustering and Preferential Attachment in growing networks" *Physical Review Letters* E, 64(025102), 2001
- [9] Lada A Adamic, Eytan Adar "Friends and neighbors on the web" *Social Networks*, 2003
- [10] A L Barabasi, H Jeong, Z Neda, E Ravasz, A Schubert, T Vicsek "Evolution of the social network of scientific collaboration" *Physica*, 2002
- [11] Leo Katz "A new status index derived from sociometric analysis" *Psychometrika*, 1953
- [12] Sergey Brin, Larry Page "The anatomy of a large-scale hypertextual Web search engine" *Computer Networks and ISDN Systems*, 1998

- [13] Jure Leskovec, Jon Kleinberg, Christos Faloutsos "Graph Evolution: Densification and Shrinking Diameters" *ACM Transactions on Knowledge Discovery from Data*, 2007
- [14] Z. Lin, Z. Hao, X. Yang, X. Liu, "Several svm ensemble methods integrated with under-sampling for imbalanced data learning" *Proceedings of the 5th International Conference on Advanced Data Mining and Applications*, 2009
- [15] P. Kang, S. Cho "Eus svms: ensemble of under-sampled svms for data imbalance problems" *Proceedings of the 13th international conference on Neural Information Processing*, 2006
- [16] Y. Liu, A. An, X. Huang "Boosting prediction accuracy on imbalanced datasets with svm ensembles" *Proceedings of the 10th Pacific-Asia conference on Advances in Knowledge Discovery and Data Mining*, 2006
- [17] B. Wang, N. Japkowicz "Boosting support vector machines for imbalanced data sets" *Knowledge and Information Systems*, 2010
- [18] John C. Platt "Probabilistic Outputs for Support Vector Machines and Comparisons to Regularized Likelihood Methods" *Advances in Large Margin Classifiers*, 1999
- [19] David J Rogers, Taffee T Tanimoto "A Computer Program for Classifying Plants" *Science*, 1960
- [20] MM Breunig, HP Kriegel, RT Ng, J Sander, "LOF: identifying density-based local outliers" *SIGMOD*, 2000
- [21] KL Li, HK Huang, SF Tian, W Xu "Improving one-class SVM for anomaly detection" *Proceedings of the Second International Conference on Machine Learning and Cybernetics*, 2003

## **Vita**

Ricky Laishram received a Bachelor of Engineering in Electronics and Communication Engineering from Birla Institute of Technology in December 2010. In 2013, he joined the Master of Science in Computer and Information Science program in Syracuse University.