

# Supervised Learning for Link Prediction Using Similarity Indices

Sergey Korolev<sup>1</sup> and Leonid Zhukov

Higher School of Economics, Moscow, Russia,  
sokorolev@edu.hse.ru

**Abstract.** The abstract should summarize the contents of the paper using at least 70 and at most 150 words. It will be set in 9-point font size and be inset 1.0 cm from the right and left margins. There will be two blank lines before and after the Abstract. . . .

**Keywords:** network analysis, graph theory, link prediction

## 1 Introduction

Over the last few years the topic of link prediction has become very popular due to the rise of recommendation systems and social networks. The

The aim of this paper is to try to devise an approach to predict the missing links in the network using only structural information of said network.

## 2 Similarity Indices

In this section, we will describe most popular local and global similarity indices and measure their accuracy in ranking links according to the probability that said links were removed from the network.

We'll start with describing the method for measuring said accuracy. We start with the set of edges  $E$  of existing graph  $G$ . Then we shuffle it and split it into  $n$  parts  $E'_i$ ,  $i \in \{1, \dots, n\}$ . Then for each step we subtract  $E'_i$  from  $E$  and look at the complement of the remaining set of edges  $\bar{E}_i$  to the set of edges of complete graph on these nodes, so that  $E^* = E \cup \bar{E} = E_i \cup E'_i \cup \bar{E} = E_i \cup \bar{E}_i$ , where  $E^*$  is the set of edges of complete graph and  $\bar{E}$  is complement of graph  $G$ .

The main tool for measuring the accuracy of ranking is AUC score, which can be described as the probability that for pair of edges  $(e_a, e_b)$   $e_a \in E'_i$ ,  $e_b \in \bar{E}$ , the ranking is such that  $score(e_a) > score(e_b)$  plus 0.5 times the probability that  $score(e_a) = score(e_b)$ . Or as formula –  $AUC = \frac{n' + 0.5n''}{n}$ , where  $n$  is number of pairs  $(e_a, e_b)$   $e_a \in E'_i$ ,  $e_b \in \bar{E}$ ,  $n'$  is number of pairs such that  $score(e_a) > score(e_b)$  and  $n''$  is number of pairs such that  $score(e_a) = score(e_b)$ .

We'll first describe all indices used for the prediction of links and then compare their AUC scores on different networks.

## 2.1 Local Similarity Indices

**Common neighbours (CN)** This is the easiest and most obvious approach to the idea of capturing the similarity of two nodes in the network. As such, a number of indices described below will use this measure with different weights. So if we denote  $\Gamma(x)$  the number of neighbours of node  $x$  in the graph, the formula looks like

$$s_{xy}^{CN} = |\Gamma(x) \cap \Gamma(y)|,$$

where  $s_{xy}$  is the index.

**Salton Index (SaI)** This index is also called cosine similarity and is defined as

$$s_{xy}^{Salton} = \frac{|\Gamma(x) \cap \Gamma(y)|}{\sqrt{k_x \times k_y}},$$

where  $k_x$  is the degree of node  $x$ .

**Jaccard Index (JI)**

$$s_{xy}^{Jaccard} = \frac{|\Gamma(x) \cap \Gamma(y)|}{|\Gamma(x) \cup \Gamma(y)|}.$$

**Sørensen Index (SoI)**

$$s_{xy}^{Sørensen} = \frac{2|\Gamma(x) \cap \Gamma(y)|}{k_x + k_y}.$$

**Hub Promoted Index (HPI)**

$$s_{xy}^{HPI} = \frac{|\Gamma(x) \cap \Gamma(y)|}{\min\{k_x, k_y\}}.$$

**Hub Depressed Index (HDI)**

$$s_{xy}^{HDI} = \frac{|\Gamma(x) \cap \Gamma(y)|}{\max\{k_x, k_y\}}.$$

**LeichtHolmeNewman Index (LHN1)**

$$s_{xy}^{LHN1} = \frac{|\Gamma(x) \cap \Gamma(y)|}{k_x \times k_y}.$$

**Preferential Attachment Index (PAI)**

$$s_{xy}^{PA} = k_x \times k_y.$$

### AdamicAdar Index (AAI)

$$s_{xy}^{AA} = \sum_{z \in \Gamma(x) \cap \Gamma(y)} \frac{1}{\log k_z}.$$

### Resource Allocation Index (RAI)

$$s_{xy}^{RA} = \sum_{z \in \Gamma(x) \cap \Gamma(y)} \frac{1}{k_z}.$$

## 2.2 Global Similarity Indices

### Katz Index (KI)

$$S^{Katz} = (I - \beta A)^{-1} - I.$$

### LeichtHolmeNewman Index (LHN2)

$$S^{LHN2} = 2m\lambda_1 D^{-1} (I - \frac{\phi A}{\lambda_1})^{-1} D^{-1}.$$

### Average Commute Time (ACT)

$$S_{xy}^{ACT} = \frac{1}{l_{xx}^+ + l_{yy}^+ - 2l_{xy}^+}.$$

### Cosine based on $L^+$ (CBL)

$$s_{xy}^{cos^+} = \cos(x, y)^+ = \frac{v_x^T v_y}{|v_x| \cdot |v_y|} = \frac{l_{xy}^+}{\sqrt{l_{xx}^+ \cdot l_{yy}^+}}.$$

### Random Walk with Restart (RWR)

$$s_{xy}^{RWR} = q_{xy} + q_{yx}.$$

### Matrix Forest Index (MFI)

$$S = (I + L)^{-1}.$$

**Table 1.** AUCs of local indices. AN - Word adjacencies network[1], CN - Neural network[2], Dp - Dolphin social network network[3], FB - American College football network[4], LM - Les Miserables network[5], PB - Books about US politics network[6], KC - Zachary’s karate club network[7], UA - US air transportation system network[8]. Best results for given network in bold.

Net	CN	SaI	JI	SoI	HPI	HDI	LHN1	PAI	AAI	RAI
AN	0.662	0.606	0.603	0.603	0.618	0.603	0.566	<b>0.744</b>	0.662	0.659
CN	0.844	0.797	0.790	0.790	0.805	0.779	0.725	0.750	0.861	<b>0.866</b>
Dp	0.779	0.774	0.779	0.779	0.763	0.780	0.762	0.619	<b>0.781</b>	<b>0.781</b>
FB	0.846	0.856	0.858	0.858	0.856	0.857	<b>0.859</b>	0.270	0.846	0.846
LM	0.910	0.882	0.880	0.880	0.847	0.878	0.820	0.776	0.918	<b>0.919</b>
PB	0.887	0.884	0.875	0.875	0.894	0.863	0.848	0.653	<b>0.897</b>	0.890
KC	0.700	0.636	0.607	0.607	0.712	0.593	0.600	0.712	0.726	<b>0.733</b>
UA	0.934	0.908	0.897	0.897	0.869	0.891	0.767	0.885	0.945	<b>0.951</b>

**Table 2.** AUCs of global indices.

Net	KI	LHN2	ACT	CBL	RWR	MFI
AN	0.714	0.549	<b>0.742</b>	0.586	0.738	0.667
CN	0.852	0.717	0.738	0.848	0.725	<b>0.865</b>
Dp	0.799	<b>0.825</b>	0.760	0.791	0.649	0.804
FB	0.857	0.879	0.588	<b>0.885</b>	0.273	0.878
LM	<b>0.884</b>	0.813	0.863	0.825	0.794	0.867
PB	0.891	0.858	0.729	0.891	0.618	<b>0.899</b>
KC	<b>0.755</b>	0.610	0.666	0.739	0.618	0.749
UA	<b>0.920</b>	NaN	0.892	0.913	0.862	0.913

### 3 Method

We begin by defining the test dataset to check the performance of the algorithm. We shuffle the set of existing edges and split it into  $n$  parts  $E'_i$ ,  $i \in \{1, \dots, n\}$ . It is said that  $n = 10$  achieves the best complexity-precision tradeoff. Now we can use the set  $\bar{E}_i = E'_i \cup \bar{E}$  as test dataset, where if  $e \in E'_i$  then  $class_e = 1$  and if  $e \in \bar{E}$  then  $class_e = 0$ .

Now we need to construct training dataset to train our classifier on it and then check its performance on the test dataset. We take the remaining set of edges  $E_i$ , shuffle it and split it into  $n$  parts  $(E'_i)'_j$ . To achieve balanced training dataset with the same amount of edges with class 0 and 1 we split the complement  $\bar{E}_i$  into subsets of the same size as  $(E'_i)'_j$ . Now for each subset  $(E'_i)'_j$  we take this subset out of  $E_i$ , then calculate indices on the complement of the remaining set and mark classes as  $class_e$ , where if  $e \in (E'_i)'_j$  then  $class_e = 1$  and if  $e \in \bar{E}_i$  then  $class_e = 0$ . And to get balanced dataset, on each step we take edges from only one subset of  $\bar{E}_i$  as examples of class 0.

After  $n$  steps of the above process we get the training dataset of size  $2|E_i|$ . Now we can train our classifier on this dataset and then test its performance on the test dataset. After testing SVM, decision trees, k nearest neighbours, random forests, LDA and QDA it was observed that the best performance is achieved with classification using random forests ().

### 4 Results

**Table 3.** Scores of the algorithm performance on the class 1 of marked links.

Net	F1	Precision	Recall	ROCAUC	Accuracy
AN	$0.027 \pm 0.003$	$0.014 \pm 0.002$	$0.579 \pm 0.072$	$0.638 \pm 0.032$	$0.696 \pm 0.023$
CN	$0.041 \pm 0.001$	$0.021 \pm 0.000$	$0.811 \pm 0.018$	$0.809 \pm 0.008$	$0.808 \pm 0.004$
Dp	$0.049 \pm 0.008$	$0.025 \pm 0.005$	$0.635 \pm 0.082$	$0.704 \pm 0.040$	$0.771 \pm 0.035$
FB	$0.199 \pm 0.025$	$0.116 \pm 0.017$	$0.721 \pm 0.044$	$0.831 \pm 0.019$	$0.939 \pm 0.011$
LM	$0.163 \pm 0.021$	$0.090 \pm 0.013$	$0.831 \pm 0.068$	$0.875 \pm 0.032$	$0.918 \pm 0.013$
PB	$0.080 \pm 0.009$	$0.042 \pm 0.005$	$0.780 \pm 0.080$	$0.812 \pm 0.040$	$0.842 \pm 0.012$
KC	$0.092 \pm 0.026$	$0.049 \pm 0.014$	$0.641 \pm 0.145$	$0.717 \pm 0.082$	$0.790 \pm 0.043$
UA	$0.083 \pm 0.008$	$0.044 \pm 0.004$	$0.901 \pm 0.024$	$0.910 \pm 0.011$	$0.919 \pm 0.009$

## 5 Conclusion

Proposed approach of supervised learning for link prediction showed interesting results on the selected networks. It appears that the algorithm classifies most of the removed links correctly, but also classifies a lot of links that were not present in the network in the first place as the ones that were removed.

## References

1. M. E. J. Newman, Phys. Rev. E 74, 036104 (2006). Available at: <http://www-personal.umich.edu/~mejn/netdata/>.
2. D. J. Watts and S. H. Strogatz, Nature 393, 440-442 (1998). Available at: <http://www-personal.umich.edu/~mejn/netdata/>.
3. D. Lusseau, K. Schneider, O. J. Boisseau, P. Haase, E. Slooten, and S. M. Dawson, Behavioral Ecology and Sociobiology 54, 396-405 (2003). Available at: <http://www-personal.umich.edu/~mejn/netdata/>.
4. M. Girvan and M. E. J. Newman, Proc. Natl. Acad. Sci. USA 99, 7821-7826 (2002). Available at: <http://www-personal.umich.edu/~mejn/netdata/>.
5. D. E. Knuth, The Stanford GraphBase: A Platform for Combinatorial Computing, Addison-Wesley, Reading, MA (1993). Available at: <http://www-personal.umich.edu/~mejn/netdata/>.
6. V. Krebs, Datasets. Available at: <http://www-personal.umich.edu/~mejn/netdata/>.
7. W. W. Zachary, An information flow model for conflict and fission in small groups, Journal of Anthropological Research 33, 452-473 (1977). Available at: <http://www-personal.umich.edu/~mejn/netdata/>.
8. V. Batageli, A. Mrvar, Pajek datasets. Available at: <http://vlado.fmf.uni-lj.si/pub/networks/data/default.htm>.