

Supervised Learning for Link Prediction Using Similarity Indices

Sergey Korolev¹ and Leonid Zhukov

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

Abstract. The problem of link prediction gathered a lot of attention in the last few years, arising in different applications ranging from recommendation systems to social networks. In this paper, we will describe most popular similarity indices, compare their performance in their ability to show links with the highest probability of being removed from initial network and describe the approach that allows to use them to predict missing links using supervised machine learning. We will show the accuracy of prediction of this method on examples of real networks.

Keywords: network analysis, graph theory, link prediction

1 Introduction

Over the last few years the topic of link prediction has become very popular as a great number of systems can be described using networks. The range of different applications of network analysis made it crucial to understand the underlying structural properties of said networks and motivated the expansion of the theory.

The link prediction problem arises as part of the most recommendation systems. Social networks also use link prediction to show user's potential friends to grow the network and provide better user experience. These commercial applications of link prediction also favour algorithms with greater computational efficiency as such networks tend to be huge.

Other applications of link prediction include research as, for example, the ability to predict connections in food webs in biology can speed up the research by removing the necessity to explore all the connections experimentally.

The aim of this paper is to describe an approach to predict the missing links in the network using only structural information of said network. And also test for underlying properties of the index space of the links of the network.

2 Similarity Indices

In this section, we will describe most popular (as presented in [1]) 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 links E of existing graph G . Then we shuffle it and split it into n parts E_i , $i \in \{1, \dots, n\}$: $E = \cup_{i=1}^n E_i$, $E_i \cap E_j = \emptyset$, $i \neq j$. Then for each step we subtract E_i from E : $E'_i = E \setminus E_i$ and look at the complement of the remaining set of links $\overline{E'_i}$ to the set of links of the complete graph on these nodes, so that $E \cup \overline{E} = E'_i \cup E_i \cup \overline{E} = E'_i \cup \overline{E'_i}$, where \overline{E} is complement of graph G .

The main tool for measuring the accuracy of the ranking is AUC [2] score, which can be described as the probability that for a pair of links (e_a, e_b) $e_a \in E_i$, $e_b \in \overline{E}$, 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 a formula – $AUC = \frac{n' + 0.5n''}{n}$, where n is the number of pairs (e_a, e_b) $e_a \in E_i$, $e_b \in \overline{E}$, n' is the number of pairs such that $score(e_a) > score(e_b)$ and n'' is the number of pairs such that $score(e_a) = score(e_b)$.

For example, we have a graph on four nodes $\{a, b, c, d\}$ with links (a, c) , (a, d) , (b, d) and (c, d) . In order to calculate AUC score, we need to first remove some subset of existing links, for example $\{(c, d)\}$, then calculate the indices of the complement of the remaining set of links, so in this case its $\{(a, b)$, (b, c) , $(c, d)\}$ and then compare their indices. So if $score(a, b) = 1$, $score(b, c) = 2$ and $score(c, d) = 2$, then AUC is equal to $\frac{1+0.5 \cdot 1}{2} = 0.75$.

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 most obvious approach to the idea of capturing the similarity of two nodes in the network. As such, the 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

$$s_{xy}^{CN} = |\Gamma(x) \cap \Gamma(y)|, \quad (1)$$

where s_{xy} is the index. For symmetric graphs it is also obvious that $s_{xy} = (A^2)_{xy}$, where A is the adjacency matrix of said graph.

Salton Index (SaI) [3] This index is also called cosine similarity and is defined as

$$s_{xy}^{Salton} = \frac{|\Gamma(x) \cap \Gamma(y)|}{\sqrt{k_x \times k_x}}, \quad (2)$$

where k_x is the degree of node x .

Jaccard Index (JI) [4] This index was proposed by Jaccard and measures the number of common neighbours weighted by the size of the union of the sets of neighbours:

$$s_{xy}^{Jaccard} = \frac{|\Gamma(x) \cap \Gamma(y)|}{|\Gamma(x) \cup \Gamma(y)|}. \quad (3)$$

Sørensen Index (SoI) [5]

$$s_{xy}^{Sørensen} = \frac{2|\Gamma(x) \cap \Gamma(y)|}{k_x + k_y}. \quad (4)$$

Hub Promoted Index (HPI) [6] This index gives higher scores to links adjacent to hubs since it is weighted by the minimum of node degrees:

$$s_{xy}^{HPI} = \frac{|\Gamma(x) \cap \Gamma(y)|}{\min\{k_x, k_y\}}. \quad (5)$$

Hub Depressed Index (HDI) This index penalises scores for the links adjacent to hubs:

$$s_{xy}^{HDI} = \frac{|\Gamma(x) \cap \Gamma(y)|}{\max\{k_x, k_y\}}. \quad (6)$$

LeichtHolmeNewman Index (LHN1) [7] This index promotes links adjacent to the nodes with lower degrees:

$$s_{xy}^{LHN1} = \frac{|\Gamma(x) \cap \Gamma(y)|}{k_x \times k_y}. \quad (7)$$

Preferential Attachment Index (PAI) This index is derived from the preferential attachment model of network growth [8]. It promotes links adjacent to the nodes of higher degrees:

$$s_{xy}^{PA} = k_x \times k_y. \quad (8)$$

AdamicAdar Index (AAI) [9] This index is a modified version of counting the amount of common neighbours as it values neighbours with lower degree higher:

$$s_{xy}^{AA} = \sum_{z \in \Gamma(x) \cap \Gamma(y)} \frac{1}{\log k_z}. \quad (9)$$

Resource Allocation Index (RAI) [10] This index is similar to the (9) and is derived from the model of resource allocation in the networks, assuming that node's ability to transfer something is inversely proportional to its degree:

$$s_{xy}^{RA} = \sum_{z \in \Gamma(x) \cap \Gamma(y)} \frac{1}{k_z}. \quad (10)$$

2.2 Global Similarity Indices

Katz Index (KI) [11] This index measures the number of paths of length l weighted by β^l , where β is the dampening factor. Since the number of paths in a symmetric graph can be expressed as the power of the adjacency matrix, it can be expressed in the matrix form as:

$$S^{Katz} = (I - \beta A)^{-1} - I, \quad (11)$$

where β should be less than the reciprocal of the largest eigenvalue of A in order for the equation to converge.

LeichtHolmeNewman Index (LHN2) [7] This index is the modification of Katz index that takes into account that two nodes tend to be similar if their neighbours are similar. It can be expressed in the matrix form [1] as:

$$S^{LHN2} = 2m\lambda_1 D^{-1} \left(I - \frac{\phi A}{\lambda_1} \right)^{-1} D^{-1}, \quad (12)$$

where m is the number of links in the network, λ_1 is the largest eigenvalue of adjacency matrix A , D is a diagonal degree matrix of form $diag\{k_1, \dots, k_n\}$, $\phi: 0 < \phi < 1$ is a free parameter that assigns higher weights to shorter paths if it's closer to 0 and to longer paths if it's closer to 1.

Average Commute Time (ACT) This index assigns higher scores to the pair of nodes for which the average length of the random walk from one node to another is shorter. It can be expressed as:

$$s_{xy}^{ACT} = \frac{1}{l_{xx}^+ + l_{yy}^+ - 2l_{xy}^+}, \quad (13)$$

where l_{ij}^+ is the (i, j) entry of pseudoinverse of the Laplacian matrix $L = D - A$.

Cosine based on L^+ (CBL) This index is defined as the cosine of two node vectors:

$$s_{xy}^{cos+} = \cos(x, y)^+ = \frac{v_x^T v_y}{|v_x| \cdot |v_y|} = \frac{l_{xy}^+}{\sqrt{l_{xx}^+ \cdot l_{yy}^+}}. \quad (14)$$

Random Walk with Restart (RWR) This index is derived from the PageRank algorithm and can be expressed as:

$$s_{xy}^{RWR} = q_{xy} + q_{yx}, \quad (15)$$

where q_{xy} is the y th element of vector \vec{q}_x . $\vec{q}_x = (1 - c)(I - cP^T)^{-1}\vec{e}_x$, where c is probability in the PageRank algorithm, P is the transition matrix with $P_{xy} = \frac{1}{k_x}$ if x and y are connected and 0 otherwise and \vec{e}_x is the vector of form (δ_{ix}) , where δ_{ix} is the Kronecker function.

Matrix Forest Index (MFI) [12] This index can be described as the ratio of the number of spanning rooted forests such that y belongs to the tree rooted in x to the number of all spanning rooted forests of the network:

$$S = (I + L)^{-1}. \quad (16)$$

We can now implement and compare AUC scores of all indices described above on different networks. (Results were achieved using Python's NetworkX [13], Pandas [14] and NumPy [15] in IPython Notebook [16]).

Table 1. AUC scores of local indices. AN - Word adjacencies network[18], CN - Neural network[19], Dp - Dolphin social network network[20], FB - American College football network[21], LM - Les Miserables network[22], PB - Books about US politics network[23], KC - Zachary’s karate club network[24], UA - US air transportation system network[25]. Best results for given network are in bold.

| | 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 | 0.744 | 0.662 | 0.659 |
| CN | 0.844 | 0.797 | 0.790 | 0.790 | 0.805 | 0.779 | 0.725 | 0.750 | 0.861 | 0.866 |
| Dp | 0.779 | 0.774 | 0.779 | 0.779 | 0.763 | 0.781 | 0.762 | 0.619 | 0.781 | 0.781 |
| FB | 0.846 | 0.858 | 0.858 | 0.858 | 0.856 | 0.857 | 0.859 | 0.270 | 0.846 | 0.846 |
| LM | 0.911 | 0.882 | 0.880 | 0.880 | 0.847 | 0.878 | 0.820 | 0.776 | 0.918 | 0.919 |
| PB | 0.887 | 0.884 | 0.875 | 0.875 | 0.894 | 0.863 | 0.848 | 0.653 | 0.897 | 0.900 |
| KC | 0.700 | 0.636 | 0.607 | 0.607 | 0.712 | 0.593 | 0.600 | 0.712 | 0.726 | 0.733 |
| UA | 0.935 | 0.909 | 0.898 | 0.898 | 0.870 | 0.892 | 0.768 | 0.887 | 0.946 | 0.952 |

Table 2. AUC scores of global indices.

| | KI | LHN2 | ACT | CBL | RWR | MFI |
|----|--------------|-------|--------------|--------------|-------|--------------|
| AN | 0.714 | 0.516 | 0.742 | 0.583 | 0.737 | 0.667 |
| CN | 0.852 | 0.547 | 0.738 | 0.848 | 0.725 | 0.865 |
| Dp | 0.799 | 0.662 | 0.760 | 0.785 | 0.647 | 0.804 |
| FB | 0.857 | 0.877 | 0.588 | 0.885 | 0.272 | 0.878 |
| LM | 0.884 | 0.568 | 0.859 | 0.812 | 0.794 | 0.867 |
| PB | 0.891 | 0.862 | 0.729 | 0.891 | 0.618 | 0.899 |
| KC | 0.755 | 0.612 | 0.666 | 0.736 | 0.620 | 0.749 |
| UA | 0.920 | 0.500 | 0.892 | 0.905 | 0.832 | 0.913 |

3 Method

We begin by defining the test dataset to check the performance of the algorithm. We shuffle the set of existing links and split it into n parts E'_i , $i \in \{1, \dots, n\}$. It is said that $n = 10$ achieves the best complexity-precision trade-off. 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 a training dataset to train our classifier on it and then check its performance on the test dataset. We take the remaining set of links E'_i , shuffle it and split it into n parts E'_{ij} : $E_i = \cup_{j=1}^n E'_{ij}$. To achieve balanced training dataset with the same amount of links with class 0 and 1 we split the complement \bar{E}'_i into subsets of the same size as E'_{ij} . Now for each subset E'_{ij} we take this subset out of E'_i , then calculate indices on the complement of the remaining set \bar{E}'_i using graph with edges $E'_i \setminus E'_{ij}$ and mark classes as $class_e$, where if $e \in E'_{ij}$ 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 links from only one subset of \bar{E}'_i of the same size as E'_{ij} 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 with linear kernel, decision trees, k nearest neighbours and random forests it was observed that the best performance is achieved with classification using random forests (as suggested in [17] and also tested in [28]).

4 Results

We test the performance our method on proposed test datasets and measure its F1 score, precision, recall, AUC score and accuracy on the class 1 of links removed from the network. We then compile all these results and take the mean and standard deviation and compare these scores for different networks. (Results were achieved using Python's Scikit-learn [26]).

The results are interesting as they show very low precision and F1 scores of class of removed links, but high recall, AUC and accuracy. It means that the classifier correctly marks most of removed links, but also classifies a lot of links that were not present in the network in the first place as the ones that were removed.

These results also show that there is difference in index space between the links removed and links from complement that differentiates removed links within the complement and that difference is enough to separate most of the initially removed links from others. This difference can be easily explored on the train set using different cross validation techniques.

Table 3. Scores of the algorithm performance on the class 1 of marked links (mean \pm std).

| | F1 | Precision | Recall | ROCAUC | Accuracy |
|----|-------------------|-------------------|-------------------|-------------------|-------------------|
| AN | 0.028 ± 0.003 | 0.014 ± 0.002 | 0.563 ± 0.077 | 0.637 ± 0.034 | 0.710 ± 0.026 |
| CN | 0.042 ± 0.002 | 0.021 ± 0.001 | 0.806 ± 0.024 | 0.808 ± 0.012 | 0.810 ± 0.007 |
| Dp | 0.046 ± 0.008 | 0.024 ± 0.004 | 0.605 ± 0.081 | 0.688 ± 0.044 | 0.770 ± 0.029 |
| FB | 0.203 ± 0.029 | 0.119 ± 0.020 | 0.710 ± 0.036 | 0.827 ± 0.018 | 0.942 ± 0.010 |
| LM | 0.163 ± 0.031 | 0.091 ± 0.020 | 0.807 ± 0.081 | 0.863 ± 0.039 | 0.919 ± 0.017 |
| PB | 0.075 ± 0.008 | 0.039 ± 0.004 | 0.757 ± 0.081 | 0.797 ± 0.039 | 0.837 ± 0.014 |
| KC | 0.078 ± 0.022 | 0.042 ± 0.012 | 0.641 ± 0.189 | 0.698 ± 0.092 | 0.752 ± 0.057 |
| UA | 0.085 ± 0.005 | 0.044 ± 0.003 | 0.896 ± 0.026 | 0.909 ± 0.014 | 0.922 ± 0.004 |

5 Conclusion

The 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.

Such kind of "optimistic" classifier can be useful in social network applications as most of the time it's easier to display lots of potential friends to user and let them decide whether or not the presented links are correct. But it can be useless in applications that require high precision and can settle down for lower recall.

This approach also allows to build fast evolving network models and evaluate change in parameters of said network with growth.

References

1. Lü, L., Zhou, T.: Link prediction in complex networks: A survey. *Physica A* 390, 1150-1170 (2011)
2. Hanely, J. A., McNeil, B. J.: The meaning and use of the area under a receiver operating characteristic (ROC) curve. *Radiology* 143, 29 (1982)
3. Salton, G., McGill, M. J.: *Introduction to Modern Information Retrieval*. McGraw-Hill, Auckland (1983)
4. Jaccard, P.: Étude comparative de la distribution florale dans une portion des Alpes et des Jura. *Bull. Soc. Vaud. Sci. Nat.* 37, 547 (1901)
5. Sørensen, T.: A method of establishing groups of equal amplitude in plant sociology based on similarity of species content and its application to analyses of the vegetation on Danish commons. *Biol. Skr.* 5, 1 (1948)

6. Ravasz, E., Somera, A. L., Mongru, D. A., Oltvai, Z. N., Barabási, A.-L.: Hierarchical organization of modularity in metabolic networks. *Science* 297, 1551 (2002)
7. Leicht, E. A., Holme, P., Newman, M. E. J.: Vertex similarity in networks. *Phys. Rev. E* 73, 026120 (2006)
8. Barabási, A.-L., Albert, R.: Emergence of scaling in random networks. *Science* 286, 509 (1999)
9. Adamic, L. A., Adar, E.: Friends and neighbors on the web. *Soc. Networks* 25, 211 (2003)
10. Zhou, T., Lü, L., Zhang, Y.-C.: Predicting missing links via local information. *Eur. Phys. J. B* 71, 623 (2009)
11. Katz, L.: A new status index derived from sociometric analysis. *Psychometrika* 18, 39 (1953)
12. Chebotarev, P., Shamis, E. V.: The matrix-forest theorem and measuring relations in small social groups. *Autom. Remote Control* 58, 1505 (1997)
13. Hagberg, A. A., Schult, D. A., Swart, P. J.: Exploring network structure, dynamics, and function using NetworkX. *Proceedings of the 7th Python in Science Conference* 11–15 (2008)
14. McKinney, W.: Data Structures for Statistical Computing in Python. *Proceedings of the 9th Python in Science Conference*, 51–56 (2010)
15. van der Walt, S., Colbert, S. C., Varoquaux, G.: The NumPy Array: A Structure for Efficient Numerical Computation. *Computing in Science & Engineering*, 13, 22–30 (2011)
16. Pérez, F., Granger, B. E.: IPython: A System for Interactive Scientific Computing. *Computing in Science & Engineering*, 9, 21–29 (2007)
17. Cukierski, W., Hamner, B., Yang, B.: Graph-based Features for Supervised Link Prediction. *Proceedings of International Joint Conference on Neural Networks*, San Jose, California, USA (2011)
18. Newman, M. E. J. *Phys. Rev. E* 74, 036104 (2006) Available at: <http://www-personal.umich.edu/~mejn/netdata/>.
19. Watts, D. J., Strogatz, S. H. *Nature* 393, 440–442 (1998) Available at: <http://www-personal.umich.edu/~mejn/netdata/>.
20. Lusseau, D., Schneider, K., Boisseau, O. J., Haase, P., Slooten, E., Dawson, S. M. *Behavioral Ecology and Sociobiology* 54, 396–405 (2003) Available at: <http://www-personal.umich.edu/~mejn/netdata/>.
21. Girvan, M., Newman, M. E. J. *Proc. Natl. Acad. Sci. USA* 99, 7821–7826 (2002) Available at: <http://www-personal.umich.edu/~mejn/netdata/>.
22. Knuth, D. E. *The Stanford GraphBase: A Platform for Combinatorial Computing*, Addison-Wesley, Reading, MA (1993) Available at: <http://www-personal.umich.edu/~mejn/netdata/>.
23. Krebs, V. *Datasets*. Available at: <http://www-personal.umich.edu/~mejn/netdata/>.
24. Zachary, W. W.: 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/>.
25. Batageli, V., Mrvar, A. *Pajek datasets*. Available at: <http://vlado.fmf.uni-lj.si/pub/networks/data/default.htm>.
26. Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., Duchesnay, É.: Scikit-learn: Machine Learning in Python. *Journal of Machine Learning Research*, 12, 2825–2830 (2011)