

Minireview

Link prediction techniques, applications, and performance: A survey



Ajay Kumar*, Shashank Sheshar Singh, Kuldeep Singh, Bhaskar Biswas

Department of Computer Science and Engineering, Indian Institute of Technology (BHU), Varanasi, 221-005, India

ARTICLE INFO

Article history:

Received 11 January 2019

Received in revised form 4 November 2019

Available online 8 February 2020

Keywords:

Link prediction

Similarity metrics

Probabilistic model

Embedding

Fuzzy logic

Deep learning

ABSTRACT

Link prediction finds missing links (in static networks) or predicts the likelihood of future links (in dynamic networks). The latter definition is useful in network evolution (Wang et al., 2011; Barabasi and Albert, 1999; Kleinberg, 2000; Leskovec et al., 2005; Zhang et al., 2015). Link prediction is a fast-growing research area in both physics and computer science domain. There exists a wide range of link prediction techniques like similarity-based indices, probabilistic methods, dimensionality reduction approaches, etc., which are extensively explored in different groups of this article. Learning-based methods are covered in addition to clustering-based and information-theoretic models in a separate group. The experimental results of similarity and some other representative approaches are tabulated and discussed. To make it general, this review also covers link prediction in different types of networks, for example, directed, temporal, bipartite, and heterogeneous networks. Finally, we discuss several applications with some recent developments and concludes our work with some future works.

© 2020 Elsevier B.V. All rights reserved.

Contents

1.	Introduction and background	2
2.	Existing methods	4
2.1.	Similarity-based methods	5
2.1.1.	Local similarity indices	5
2.1.2.	Global similarity indices	8
2.1.3.	Quasi-local indices	11
2.2.	Probabilistic and maximum likelihood models	12
2.2.1.	Local probabilistic model for link prediction	12
2.2.2.	Probabilistic relational model for link prediction (PRM)	14
2.2.3.	Hierarchical structure model (HSM) [72]	14
2.2.4.	Stochastic block model (SBM) [73]	15
2.2.5.	Exponential random graph model (ERGM) or P-star model	16
2.3.	Link prediction using dimensionality reduction	16
2.3.1.	Embedding-based link prediction	17
2.3.2.	Matrix factorization/decomposition-based link prediction	18
2.4.	Other approaches	19
2.4.1.	Learning-based frameworks for link prediction	19

* Corresponding author.

E-mail addresses: ajayk.rs.cse16@iitbhu.ac.in (A. Kumar), shashankrs.rs.cse16@iitbhu.ac.in (S.S. Singh), kuldeep.rs.cse13@iitbhu.ac.in (K. Singh), bhaskar.cse@iitbhu.ac.in (B. Biswas).

2.4.2.	Information theory-based link prediction.....	19
2.4.3.	Clustering-based link prediction.....	20
2.4.4.	Structural perturbation method (SPM) [159].....	21
3.	Experimental setup and results analysis.....	22
3.1.	Evaluation metrics.....	22
3.1.1.	Area under the receiver operating characteristics curve (AUROC).....	22
3.1.2.	Area under the precision–recall curve (AUPR).....	22
3.1.3.	Average precision.....	23
3.1.4.	Recall@k.....	23
3.2.	Datasets.....	23
3.3.	Accuracy results.....	23
3.3.1.	Recall@k.....	24
3.3.2.	Area under the precision–recall curve (AUPR).....	24
3.3.3.	Area under the receiver operating characteristics curve (AUROC).....	24
3.3.4.	Average precision.....	26
3.4.	Efficiency.....	27
4.	Variations of link prediction problem.....	27
4.1.	Link prediction in temporal networks.....	29
4.2.	Link prediction in bipartite networks.....	29
4.3.	Link prediction in heterogeneous networks.....	30
5.	Link prediction applications.....	30
5.1.	Network reconstruction.....	30
5.2.	Recommender system.....	31
5.3.	Network completion problem.....	31
5.4.	Spam mail detection.....	31
5.5.	Privacy control in social networks.....	32
5.6.	Identifying missing references in a publication.....	32
5.7.	Routing in networks.....	32
5.8.	Incorporating user's influence in link prediction.....	32
6.	Recent developments.....	33
6.1.	Link prediction using deep learning.....	33
6.2.	Fuzzy model-based link prediction.....	34
7.	Conclusion and future directions.....	37
	Declaration of competing interest.....	37
	References.....	38

1. Introduction and background

A social network (a more general term is a complex network) is a standard approach to model communication in a group or community of persons. Such networks can be represented as a graphical model in which a node maps to a person or social entity, and a link corresponds to an association or collaboration between corresponding persons or social entities. The relationships among individuals are continuously changing, so the addition and/or deletion of several links and vertices take place. It results in social networks to be highly dynamic and complex. Lots of issues arise when we study a social network, some of which are changing association patterns over time, factors that drive those associations, and the effects of those associations to other nodes. Here, we address a specific problem termed as link prediction.

Problem characterization. Consider a simple undirected network $G(V, E)$ (Refer to the Fig. 1), where V characterizes a vertex-set and E , the link-set. A simple graph is considered throughout the paper, i.e., parallel links and self-loops are not permitted. In this paper, we use (vertex \equiv node), (link \equiv edge) and (graph \equiv network) interchangeably. In the graph, a universal set U contains a total of $\frac{n(n-1)}{2}$ links (total node-pairs), where $n = |V|$ represents the number of total vertices of the graph. $(|U| - |E|)$ number of links are termed as the non-existing links, and some of these links may appear in the near future. Finding such missing links (i.e., AC, BD, and AD) is the aim of link prediction [1].

Formally, Liben-Nowell et al. [2] defined the link prediction problem as: suppose a graph $G_{t_0-t_1}(V, E)$ represents a snapshot of a network during time interval $[t_0, t_1]$ and $E_{t_0-t_1}$, a set of links present in that snapshot. The task of link prediction is to find set of links $E_{t'_0-t'_1}$ during the time interval $[t'_0, t'_1]$ where $[t_0, t_1] \leq [t'_0, t'_1]$. The link prediction idea is useful in several domains of application. Examples include automatic hyperlink creation [3], website hyperlink prediction [4] in the Internet and web science domain, and friend recommendation on Facebook. Building a recommendation system [5,6] in e-commerce is an essential task that uses link prediction as a basic building block. In Bioinformatics, protein–protein interactions (PPI) also have been implemented using link prediction [7]. In security concern areas, link prediction is used to distinguish hidden links among terrorists and their organizations.

Newman presented a paper on link prediction on collaboration networks in Physics and Biology [8]. In such networks, two authors are considered to be connected if they have at least one paper coauthored by them simultaneously. In the

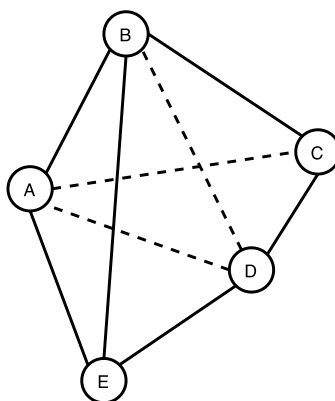


Fig. 1. The Link Prediction (LP) finds missing links (i.e., AD, AC, and BD) in this observed network.

empirical study, the author demonstrated that the likelihood of a pair of researchers teaming up increments with the numbers of different colleagues they have in mutual relation, and the likelihood of a specific researcher acquiring new partners increments with the number of his past teammates. The outcomes give experimental proof in favor of previously guessed mechanisms for clustering and power-law degree distributions in networks. Next, Liben-Nowell et al. [2] proposed a link prediction model explicitly for a social network. Each node in the network corresponds to a person or an entity, and a link between two nodes shows the interaction between them. The learning paradigm in this environment can be used to extract the similarities between two nodes by several similarity metrics. Ranks are assigned to each pair of nodes based on these similarities, then higher ranked node pairs are designated as predicted links. Further, Hasan et al. [9] expanded this work and demonstrated that there is a significant increase in prediction results when additional topological information about the network is available. They considered different similarity measures as features and performed a binary classification task using a supervised learning approach, which is similar to link prediction in this environment. In the relational context [10–12] and in the Internet domain [13], the link prediction problem has not used graph representation explicitly. The proposed frameworks can acknowledge any relational dataset where there is a relation among objects. In such frameworks, modeling paradigms like probabilistic relational models [14], graphical models [15], and stochastic relational models [7,16,17] have been used.

The upsides of these methodologies incorporate genericity and simplicity, where the model can integrate attributes of the entities. On the downside, they are normally intricate and contain the excessive number of parameters, a large portion of which may be complex to the user. The research on social network evolution nearly takes after the task of link prediction where several seminal studies have been proposed like Barabasi and Albert work [18] on random network published in Science, Kleinberg work [19] in Nature Communication, [20] in KDD, [21] in EPL, and [22] in Scientific Reports. An evolution model considers some notable characteristics like the small world phenomenon [19] and the power-law degree distribution [18] during the evolution of the links of the underlying network. The essential contrast between the link prediction model and the evolution model is that the former spotlights on the overall characteristics of the network and the latter on the network's local characteristics to estimate missing links. The continuous growing size of social networks such as Myspace, Facebook, LinkedIn, Flickr, etc., has shown to be one of the key challenges in link prediction. Prior existing methodologies may not be implemented to such networks because of continuous evolving nature and their huge size, so some other methodologies are required to address these issues. As an example, Tylenka et al. [23] show that the timestamps of previous affiliations (that expressly use the genealogy of interactions) can be used to enhance the performance of link prediction. Song et al. [24] considered a social network consisting of around 2 million nodes and 90 million links and compute similarity measures among these nodes using matrix factorization. For such a huge network, a traditional algorithm would fail to calculate pair-wise similarities. Recently, Acar et al. [25] implemented tensor as the extension of matrix factorization, which is more richer and higher-order models.

We exhibit a review of previous methodologies shedding light on link prediction with the point of convergence mostly on social network graphs. We order these methodologies into several categories; one category of those calculates a similarity score between pairs of vertices in which higher scored pairs are assumed to have links between them. Another category of algorithms is based on probabilistic approaches in which Bayesian and relational models have been used. Dimensionality reduction approaches consisting of embedding and factorization-based methods have grouped into one, and some other approaches also have been studied.

Difference from the existing surveys. The present surveys on link prediction explore a large area of complex networks. It comprises several techniques ranging from classical structural and probabilistic ones to recent network embedding methods, fuzzy models, and deep learning models. It touches other areas like the information-theoretic model, clustering-based models, and factorization-based models. Link prediction in different types of networks like temporal, bipartite, and

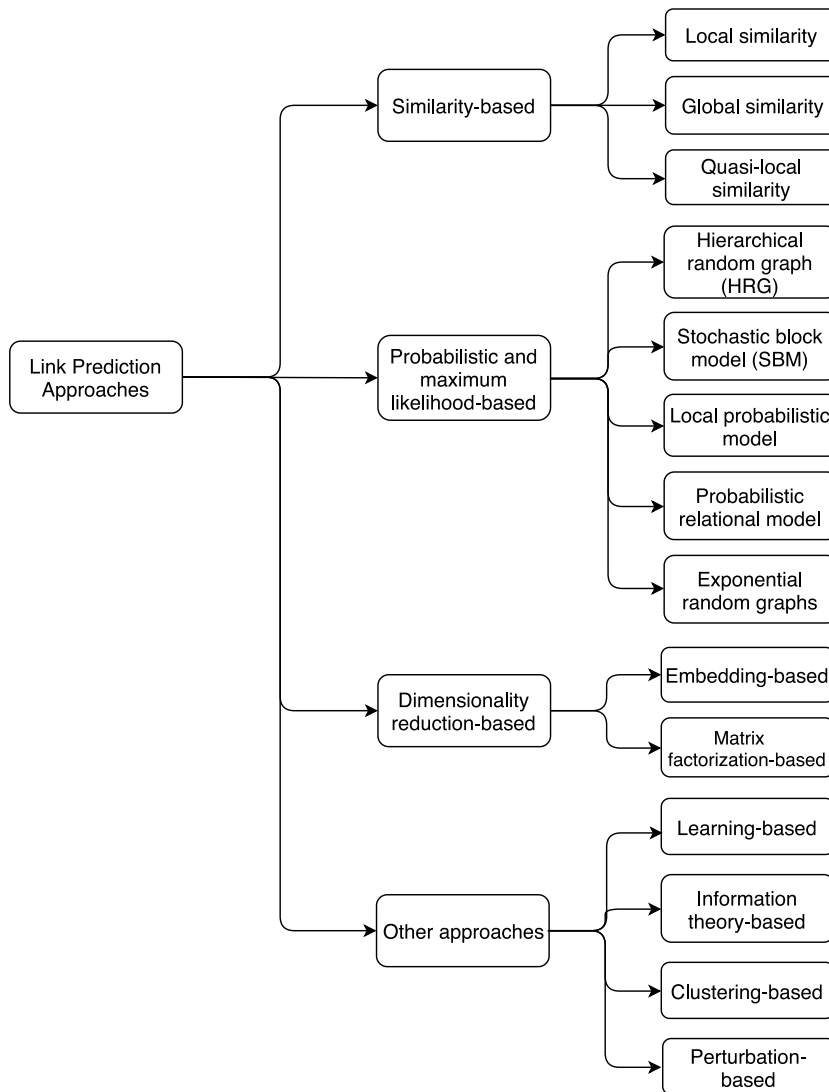


Fig. 2. Taxonomy of link prediction approaches.

heterogeneous networks are also explored. The accuracy of structural-based and some other representative models are experimentally compared with four well-known evaluation metrics on seven real-world datasets. Existing surveys [1,2,26,27] present a good effort in link prediction with their limitations, such as Nowell and Kleinberg [2] in 2007, experimentally explored structural and high level approaches lacking diversity. One more experimental survey by Martinez et al. [26] extensively visits several areas but lacking recent models link embedding and deep learning, etc. Hasan et al. [27] present a theoretical survey mainly focused on machine learning approaches.

Organization. The existing methods in the literature have been reported in Section 2. Section 3 discusses an experimental study consisting of evaluation strategies and basic topological information of several real network datasets. Moreover, the experimental results of the accuracy and efficiency of similarity-based methods also have been shown in this section. Variations of link prediction problem are reported in Section 4 and different applications are described in Section 5. Section 6 shows some recent developments. Finally, Section 7 concludes this work with some future directions.

2. Existing methods

Recently, numerous methodologies of link prediction have been implemented. These methods can be grouped into several categories, like similarity-based, probabilistic models, learning-based models, etc as shown in the Fig. 2.

2.1. Similarity-based methods

Similarity-based metrics are the simplest one in link prediction, in which for each pair x and y , a similarity score $S(x, y)$ is calculated. The score $S(x, y)$ is based on the structural or node's properties of the considered pair. The non-observed links (i.e., $U - E^T$) are assigned scores according to their similarities. The pair of nodes having a higher score represents the predicted link between them. The similarity measures between every pair can be calculated using several properties of the network, one of which is structural property. Scores based on this property can be grouped in several categories like local and global, node-dependent and path-dependent, parameter-dependent and parameter-free, and so on.

2.1.1. Local similarity indices

Local indices are generally calculated using information about common neighbors and node degree. These indices consider immediate neighbors of a node. Examples of such indices contains common neighbor [8], preferential attachment [28], Adamic/Adar [29], resource allocation [30], etc.

- (i) Common Neighbors (CN) [8] In a given network or graph, the size of common neighbors for a given pair of nodes x and y is calculated as the size of the intersection of the two nodes neighborhoods.

$$S(x, y) = |\Gamma(x) \cap \Gamma(y)|, \quad (1)$$

where $\Gamma(x)$ and $\Gamma(y)$ are neighbors of the node x and y respectively. The likelihood of the existence of a link between x and y increases with the number of common neighbors between them. In a collaboration network, Newman calculated this quantity and demonstrated that the probability of collaboration between two nodes depends upon the common neighbors of the selected nodes. Kossinets and Watts [31,32] investigated a large social network and recommended that two students are more likely to be friends who are having numerous common friends. It has been observed that the common neighbor approach performs well on most real-world networks and beats other complex methods.

- (ii) Jaccard Coefficient [33] This metric is similar to the common neighbor. Additionally, it normalizes the above score, as given below.

$$S(x, y) = \frac{|\Gamma(x) \cap \Gamma(y)|}{|\Gamma(x) \cup \Gamma(y)|}. \quad (2)$$

i.e., the Jaccard coefficient is defined as the probability of selection of common neighbors of pairwise vertices from all the neighbors of either vertex. The pairwise Jaccard score increases with the number of common neighbors between the two vertices considered. Liben-Nowell et al. [2] demonstrated that this similarity metric performs worse as compared to Common Neighbors.

- (iii) Adamic/Adar Index [29] Adamic and Adar presented a metric to calculate a similarity score between two web pages based on shared features, which are further used in link prediction after some modification by Liben-Nowell et al. [2].

$$S(x, y) = \sum_{z \in \Gamma(x) \cap \Gamma(y)} \frac{1}{\log k_z}, \quad (3)$$

where k_z is the degree of the node z . It is clear from the equation that more weights are assigned to the common neighbors having smaller degrees. This is also intuitive in the real-world scenario, for example, a person with more number of friends spend less time/resource with an individual friend as compared to the less number of friends.

- (iv) Preferential Attachment [28] (PA) The idea of preferential attachment is applied to generate a growing scale-free network. The term growing represents the incremental nature of nodes over time in the network. The likelihood incrementing new connection associated with a node x is proportional to k_x , the degree of the node. Preferential attachment score between two nodes x and y can be computed as

$$S(x, y) = k_x \cdot k_y. \quad (4)$$

This index shows the worst performance on most networks, as reported in the result section. The simplicity (as it requires the least information for the score calculation) and the computational time of this metric are the main advantages. Also, it can be used in a non-local context as it requires only degree as information and not the common neighbors. In assortative networks, the performance of the PA improves, while very bad for disassortative networks. In other words, PA shows better results if larger degree nodes are densely connected, and lower degree nodes are rarely connected.

In a supervised learning framework, Hasan et al. [9] showed that aggregate functions (e.g., sum, multiplication, etc.) over feature values of vertices could be applied to compute link feature value. In the above equation, summation can also be used instead of multiplication as an aggregate function, and in fact, it has been proved to be quite useful. [9] showed the preferential attachment with aggregate function "sum" performs well for the link prediction in coauthorship network.

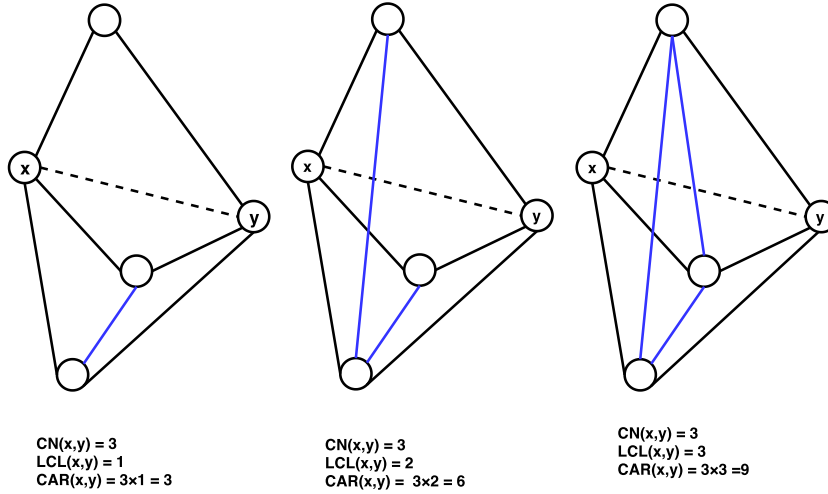


Fig. 3. CAR Index = (Number of CNs) \times (Number of LCLs).

- (v) Resource Allocation Index (RA) [30] The original dynamics of this similarity index is originated from Ou et al. [34] work published in “Physical Review E” on resource allocation dynamics on complex networks. Consider two non-adjacent vertices x and y . Suppose node x sends some resources to y through the common nodes of both x and y then the similarity between the two vertices is computed in terms of resources sent from x to y . This is expressed mathematically as

$$S(x, y) = \sum_{z \in \Gamma(x) \cap \Gamma(y)} \frac{1}{k_z}. \quad (5)$$

This similarity measure and the Adamic/Adar are very similar to each other, as shown by the Eqs. (3) and (5), respectively. The difference is that the RA index heavily penalizes to higher degree nodes compared to the AA index. Prediction results of these indices become almost the same for smaller average degree networks. This index shows good performance on heterogeneous networks with a high clustering coefficient, especially on transportation networks (e.g., usair97 as reported in the result section).

- (vi) Cosine similarity or Salton Index (SI) In a vector space, document similarities can be computed using the Salton index, also known as Cosine similarity [35]. This similarity index between two records (documents) is measured by calculating the Cosine of the angle between them. The metric is all about the orientation and not magnitude. The Cosine similarity can be computed as

$$S(x, y) = \frac{|\Gamma(x) \cap \Gamma(y)|}{\sqrt{(k_x \cdot k_y)}}. \quad (6)$$

- (vii) Sorensen Index [36] This index of similarity was applied mainly to the ecological data samples and given by Thorvald Sorensen in 1948. It is very similar to the Jaccard index, as we can observe in Eq. (7). McCune et al. show that it is more robust than Jaccard against the outliers [37].

$$S(x, y) = \frac{2|\Gamma(x) \cap \Gamma(y)|}{k_x + k_y}. \quad (7)$$

- (viii) CAR-based Common Neighbor Index (CAR) [38] CAR-based indices are presented based on the assumption that the link existence between two nodes is more likely if their common neighbors are members of a local community (local-community-paradigm (LCP) theory) [38]. In other words, the likelihood existence increases with the number of links among the common neighbors (local community links (LCLs)) of the seed node pair as described in Fig. 3.

$$\begin{aligned}
 S(x, y) &= CN(x, y) \times LCL(x, y) \\
 &= CN(x, y) \times \sum_{z \in \Gamma(x) \cap \Gamma(y)} \frac{|\gamma(z)|}{2},
 \end{aligned} \quad (8)$$

where $CN(x, y) = |\Gamma(x) \cap \Gamma(y)|$ is number of common neighbors $LCL(x, y)$ refers to the number of local community links which are defined as the links among the common neighbors of seed nodes x and y [38]. $\gamma(z)$ is the subset of neighbors of node z that are also common neighbors of x and y .

- (ix) CAR-based Adamic/Adar Index (CAA) [38] If LCL is considered as an accuracy enhancer, then the CAA index is obtained by incorporating the LCL theory to the well-known AA index and mathematically expressed by the equation given below.

$$S(x, y) = \sum_{z \in \Gamma(x) \cap \Gamma(y)} \frac{|\gamma(z)|}{\log_2(k_z)}, \tag{9}$$

- (x) CAR-based Resource Allocation Index (CRA) [38] The authors show the general application of the LCL theory to other indices and generate the CRA index by incorporating this concept into the existing RA index of the literature. Mathematically, the CRA can be expressed as

$$S(x, y) = \sum_{z \in \Gamma(x) \cap \Gamma(y)} \frac{|\gamma(z)|}{k_z}, \tag{10}$$

- (xi) CAR-based Preferential Attachment Index (CPA) [38] This is the preferential attachment index based on the CAR index. CPA is obtained by incorporating the LCL theory to the original PA method and expressed mathematically by

$$S(x, y) = e_x \cdot e_y + e_x \cdot \text{CAR}(x, y) + e_y \cdot \text{CAR}(x, y) + \text{CAR}(x, y)^2, \tag{11}$$

where e_x is the number of neighbors of x not shared by y and $\text{CAR}(x, y)$ is the similarity score of the node pair x and y using CAR index.

CAR-based methods listed above show the best performance on LCP networks. The LCP networks are related to dynamic and heterogeneous systems and facilitate network evolution of social and biological networks.

- (xii) Hub Promoted Index (HPI) [39] Ravasz et al. published a paper on a cellular organization in metabolic networks. They show that the metabolic networks are composed of several small and highly connected topological modules and are combined into larger and less cohesive hierarchical structures. The number of such modules and their degree of clustering follow the power law. This similarity index promotes the formation of links between the sparsely connected nodes and hubs. It also tries to prevent links formation between the hub nodes. This similarity metric can be expressed mathematically as

$$S(x, y) = \frac{|\Gamma(x) \cap \Gamma(y)|}{\min(k_x, k_y)}. \tag{12}$$

- (xiii) Hub Depressed Index (HDI) [39] This index is the same as the previous one but with the opposite goal as it avoids the formation of links between hubs and low degree nodes in the networks. The Hub depressed index promotes the links evolution between the hubs as well as the low degree nodes. The mathematical expression for this index is given below.

$$S(x, y) = \frac{|\Gamma(x) \cap \Gamma(y)|}{\max(k_x, k_y)}. \tag{13}$$

- (xiv) Local Naive Bayes-based Common Neighbors (LNBCN) [40] The above similarity indices are somehow based on common neighbors of the node pair where each of the which are equally weighted. This method is based on the Naive Bayes theory and arguments that different common neighbors play different role in the network and hence contributes differently to the score function computed for non-observed node pairs.

$$S(x, y) = \sum_{z \in \Gamma(x) \cap \Gamma(y)} [\log(\frac{C(z)}{1 - C(z)}) + \log(\frac{1 - \rho}{\rho})], \tag{14}$$

where $C(z)$ is node clustering coefficient and ρ is the network density expressed as

$$\rho = \frac{|E|}{n(n - 1)/2}.$$

- (xv) Leicht-Holme-Newman Local Index (LHNL) [41] Leicht et al. [41] presented a paper on vertex similarity in networks. Their work is based on the concept of self-similarity, i.e., two vertices are similar to each other if their corresponding neighbors are self-similar to themselves. This score is defined by the ratio of the path of length two that exits between two vertices and the expected path of the same length between them.

$$S(x, y) = \frac{|\Gamma(x) \cap \Gamma(y)|}{k_x \cdot k_y}. \tag{15}$$

- (xvi) Node Clustering Coefficient (CCLP) [42] This index is also based on the clustering coefficient property of the network in which the clustering coefficients of all the common neighbors of a seed node pair are computed and summed to find the final similarity score of the pair. Mathematically, this index can be expressed as follows.

$$S(x, y) = \sum_{z \in \Gamma(x) \cap \Gamma(y)} C(z), \tag{16}$$

where

$$C(z) = \frac{t(z)}{k_z(k_z - 1)}$$

is clustering coefficient of the node z and $t(z)$ is the total triangles passing through the node z .

(xvii) Node and Link Clustering coefficient (NLC) [43] This similarity index is based on the basic topological feature of a network called "Clustering Coefficient" [44,45]. The clustering coefficients of both nodes and links are incorporated to compute the similarity score.

$$S(x, y) = \sum_{z \in \Gamma(x) \cap \Gamma(y)} \frac{|\Gamma(x) \cap \Gamma(z)|}{k_z - 1} \times C(z) + \frac{|\Gamma(y) \cap \Gamma(z)|}{k_z - 1} \times C(z), \quad (17)$$

2.1.2. Global similarity indices

Global indices are computed using entire topological information of a network. The computational complexities of such methods are higher and seem to be infeasible for large networks.

(i) Katz Index [46] This index can be considered as a variant of the shortest path metric. It directly aggregates over all the paths between x and y and damps exponentially for longer paths to penalize them. It can be expressed mathematically as

$$S(x, y) = \sum_{l=1}^{\infty} \beta^l |\text{paths}_{x,y}^{(l)}| = \sum_{l=1}^{\infty} \beta^l (A^l)_{x,y}, \quad (18)$$

where, $\text{paths}_{x,y}^{(l)}$ is considered as the set of total l length paths between x and y , β is a damping factor that controls the path weights and A is the adjacency matrix. For the convergence of above equation,

$$\beta < \frac{1}{\lambda_1},$$

where λ_1 is the maximum eigen value of the matrix A . If 1 is added to each element of the diagonal of the resulting similarity matrix S , this expression can be written in matrix terms as

$$S = \beta AS + I, \quad (19)$$

where I is the identity matrix of the proper dimension. The similarity between all pairs of nodes can be directly computed using the closed-form by rearranging for S in the previous expression and subtracting the previously added 1 to the elements in the diagonal. Katz score for each pair of nodes in the network is calculated by finding the similarity matrix as

$$S = (I - \beta A)^{-1} - I. \quad (20)$$

The computational complexity of the given metric is high, and it can be roughly estimated to be cubic complexity which is not feasible for a large network.

(ii) Random Walk with Restart (RWR) [47] Let α be a probability that a random walker iteratively moves to an arbitrary neighbor and returns to the same starting vertex with probability $(1 - \alpha)$. Consider q_{xy} to be the probability that a random walker who starts walking from vertex x and located at the vertex y in steady-state. Now, this probability of walker to reach the vertex y is expressed mathematically as

$$\vec{q}_x = \alpha P^T \vec{q}_x + (1 - \alpha) \vec{e}_x, \quad (21)$$

where \vec{e}_x is the seed vector of length $|V|$ (i.e., the total number of vertices in the graph). This vector consists of zeros for all components except the elements x itself. The transition matrix P can be expressed as

$$P_{xy} = \begin{cases} \frac{1}{k_x} & \text{if } x \text{ and } y \text{ are connected,} \\ 0 & \text{otherwise.} \end{cases}$$

Simplifying the above equation we get,

$$\vec{q}_x = (1 - \alpha)(I - \alpha P^T)^{-1} \vec{e}_x. \quad (22)$$

Since this similarity is not symmetric, the final score between the node pair (x, y) can be computed as

$$S(x, y) = q_{xy} + q_{yx}. \quad (23)$$

It is clear from Eq. (22) that matrix inversion is required to solve, which is quite expensive and prohibitive for large networks. A faster version of this index is implemented in [47].

- (iii) Shortest Path [2] Lots of algorithms [48–50] are available to compute the shortest path between a vertex pair in a graph that applies to a different scenario. Liben-Nowell et al. [2] provided the shortest path with its negation as a metric to link prediction. The inverse relation between the similarity and length of the shortest path is captured by the following mathematical equation given below [2].

$$S(x, y) = -|d(x, y)|, \tag{24}$$

where Dijkstra algorithm [48] is applied to efficiently compute the shortest path $d(x, y)$ between the node pair (x, y) . The prediction accuracy of this index is low compared to most local indices that make room for the consideration of indirect path in link prediction.

Several paths of different lengths can exist between a vertex pair, the similarity between such pair is computed by several other methods like Katz index, local path index, etc.

- (iv) Leicht–Holme–Newman Global Index (LHNG) [41] This global index, proposed by Leicht et al. [41], is based on the principle that two nodes are similar if either of them has an immediate neighbor, which is similar to the other node. This is a recursive definition of similarity where a termination condition is needed. The termination condition is introduced in terms of self-similarity, i.e., a node is similar to itself. Thus, the similarity score equation consists of two terms: first, the neighbor similarity, and the second, self-similarity, as given below.

$$S(x, y) = \phi \sum_z A_{x,z} S_{z,y} + \psi \delta_{x,y}. \tag{25}$$

Here, the first term is neighborhood similarity and the second term is self-similarity. ϕ and ψ are free parameters that make a balance between these two terms. In matrix form [1,26]

$$\begin{aligned} S &= \phi AS + \psi I = \psi(I - \phi A)^{-1} \\ &= \psi(I + \phi A + \phi^2 A^2 + \dots) \end{aligned} \tag{26}$$

When the free parameter $\psi = 1$, this index resembles to the Katz index [46]. Moreover, we note that $A^1(x, y)$, $A^2(x, y)$, etc, represent number of paths of length 1, 2, and so on respectively. After some calculation, the final similarity score can be expressed as given below [41].

$$S = 2m\lambda_1 D^{-1} (I - \frac{\alpha}{\lambda_1} A)^{-1} D^{-1}, \tag{27}$$

where D is the diagonal matrix, and β is dumping factor that penalizes the longer path contribution. Dropping the constant term $2m\lambda_1$ and rearranging Eq. (27), it becomes

$$DSD = \frac{\beta}{\lambda_1} A(DSD) + I. \tag{28}$$

The Eq. (28) solved by iterating this equation repeatedly with the initial value of $(DSD) = 0$ and converges normally in 100 iterations as claimed by the authors [41].

- (v) Cosine based on L^+ (Cos^+) [51] Laplacian matrix is extensively used as an alternative representation of graphs in spectral graph theory [52]. This matrix can be defined as $L = D - A$, where, D is the diagonal matrix consisting of the degrees of each node of the matrix and A is the adjacency matrix of the graph. The pseudoinverse of the matrix defined by Moore–Penrose is represented as L^+ and each entry of this matrix is used to represent the similarity score [51] between the two corresponding nodes. The most common way to compute this pseudoinverse is by computing the singular value decomposition (SVD) of the Laplacian matrix ($L = \mathcal{U} \Sigma \mathcal{V}^T$), where \mathcal{U} and \mathcal{V} are left and right singular vectors of SVD as follows

$$L^+ = \mathcal{V} \Sigma^+ \mathcal{U}^T, \tag{29}$$

Σ^+ is obtained by taking the inverse of each nonzero element of the Σ . Further, the similarity between two nodes x and y can be computed using any inner product measure such as Cosine similarity given as

$$S(x, y) = \frac{L^+_{x,y}}{\sqrt{L^+_{x,x} L^+_{y,y}}}. \tag{30}$$

- (vi) Average Commute Time (ACT) [53] This index is based on the random walk concept. A random walk is a Markov chain [54,55] which describes the movements of a walker. ACT is first coined by Gobel and Jagers [56] and applied in link prediction by [53]. It defined as the average number of movements/steps required by a random walker to reach the destination node y , and come back to the starting node x . If $m(x, y)$ be the number of steps required by the walker to reach y from x , then the following expression captures this concept.

$$n(x, y) = m(x, y) + m(y, x). \tag{31}$$

This equation can be simplified using the pseudoinverse of the Laplacian matrix L^+ [51,57] as

$$n(x, y) = |E|(L^+_{xx} + L^+_{yy} - 2L^+_{xy}), \tag{32}$$

where l_{xy}^+ denotes the (x, y) entry of the matrix L^+ . Pseudoinverse of the Laplacian, L^+ can be computed as [51]

$$L^+ = (L - \frac{ee^T}{n})^{-1} + \frac{ee^T}{n}, \tag{33}$$

where e is a column vector consisting of 1's. The square root of Eq. (32) is called Euclidean commute time distance (ECTD) [51], so smaller value of this equation will represent higher similarity. The final expression representing this similarity index is thus, given by the squared reciprocal of Eq. (32) and by ignoring the constant term $|E|$.

$$S(x, y) = \frac{1}{l_{xx}^+ + l_{yy}^+ - 2l_{xy}^+}. \tag{34}$$

- (vii) Normalized Average Commute Time (NACT) [53] This is a variant of ACT that takes into account node degrees. For a high degree node (hub) y , $m(x, y)$ is usually small regardless of x , the similarity measure is normalized with stationary distribution π of the Markov chain describing random walker on the graph. This normalized measure can be computed with the following equation

$$S(x, y) = \frac{1}{(m(x, y)\pi_y + m(y, x)\pi_x)} \tag{35}$$

It is easy to show that, for a connected graph $\pi(x) = \frac{k_x}{\sum_y k_y}$.

- (viii) Matrix Forest Index (MF) [58] This index is based on the concept of spanning tree which is defined as the subgraph that spans total nodes without forming any cycle. The spanning tree may contain total or less number of links as compared to the original graph. Chebotarev and Shamis proposed a theorem called matrix-forest theorem [58] which states that the number of spanning tree in a graph is equal to the cofactor of any entry of Laplacian matrix of the graph. Here, the term forest represents the union of all rooted disjoint spanning trees. The similarity between two nodes x and y can be computed with Eq. (36) given below.

$$S = (I + L)^{-1}, \tag{36}$$

where $(I + L)_{(x,y)}$ is the number of spanning rooted forests (x as root) consisting of both the nodes x and y . Moreover, this quantity is equal to the cofactor of $(I + L)_{(x,y)}$.

- (ix) SimRank [59] SimRank is a measure of structural context similarity and shows object-to-object relationships. It is not domain-specific and recommends to apply in directed or mixed networks. The basic assumption of this measure is that two objects are similar if they are related to similar objects. SimRank computes how soon two random walkers meet each other, starting from two different positions. The measure is computed recursively using Eq. (37).

$$S(x, y) = \begin{cases} \frac{\alpha}{k_x k_y} \sum_{i=1}^{k_x} \sum_{j=1}^{k_y} S(\Gamma_i(x), \Gamma_j(y)) & x \neq y \\ 1 & x = y \end{cases} \tag{37}$$

where, $\alpha \in (0, 1)$ is a constant. $\Gamma_i(x)$ and $\Gamma_j(y)$ are the i th and j th elements in the neighborhood sets $\Gamma(x)$ and $\Gamma(y)$ respectively. Initially, $S(x, y) = A(x, y)$, i.e., $S(x, x) = 1$ and $S(x, y) = 0$ for $x \neq y$. This measure can be represented in matrix form as

$$S(x, y) = \alpha W^T S W + (1 - \alpha)I, \tag{38}$$

where W is the transformation matrix and computed by normalizing each column of adjacency matrix A as $W_{ij} = \frac{a_{ij}}{\sum_{k=1}^n a_{kj}}$.

The computational complexity of this measure is high for a large network, and to reduce its time, the authors [59] suggest pruning recursive branches after radius l . The time required to compute this score for each pair is $O(k^{2l+2})$, and total time is $O(n^2 k^{2l+2})$.

- (x) Rooted Pagerank (RPR) The idea of PageRank [60] was originally proposed to rank the web pages based on the importance of those pages. The algorithm is based on the assumption that a random walker randomly goes to a web page with probability α and follows hyper-link embedded in the page with probability $(1 - \alpha)$. Chung et al. [61] used this concept incorporated with a random walk in link prediction framework. The importance of web pages, in a random walk, can be replaced by stationary distribution. The similarity between two vertices x and y can be measured by the stationary probability of y from x in a random walk where the walker moves to an arbitrary neighboring vertex with probability α and returns to x with probability $(1 - \alpha)$. Mathematically, this score can be computed for all pair of vertices as

$$RPR = (1 - \alpha)(I - \alpha \hat{N})^{-1}, \tag{39}$$

where $\hat{N} = D^{-1}A$ is the normalized adjacency matrix with the diagonal degree matrix $D[i, i] = \sum_j A[i, j]$.

Table 1
Comparison of similarity-based approaches.

Properties	Local indices	Global indices	Quasi-local indices
Nature	Simple	Complex	Moderate
Features employed	Local neighborhood	Entire network	More than local neighborhood
Computational complexity	Low	high	Moderate
Parallelization	Easy	More complex	Moderate
Implementation	Feasible for large networks	Feasible for small networks	Feasible for large networks

2.1.3. Quasi-local indices

Quasi-local indices have been introduced as a trade-off between local and global approaches or performance and complexity, as shown in Table 1. These metrics are as efficient to compute as local indices. Some of these indices extract the entire topological information of the network. The time complexities of these indices are still below compared to the global approaches. Examples of such indices include local path index, local random walk index [53], local directed path (LDP) [62], etc.

- (i) Local Path Index (LP) With the intent to furnish a good trade-off between accuracy and computational complexity, the local path-based metric is considered [63]. The metric is expressed mathematically as

$$S^{LP} = A^2 + \varepsilon A^3, \tag{40}$$

where ε represents a free parameter. Clearly, the measurement converges to common neighbor when $\varepsilon = 0$. If there is no direct connection between x and y , $(A^3)_{xy}$ is equated to the total different paths of length 3 between x and y . The index can also be expanded to generalized form

$$S^{LP} = A^2 + \varepsilon A^3 + \varepsilon^2 A^4 + \dots + \varepsilon^{(n-2)} A^n, \tag{41}$$

where n is the maximal order. Computing this index becomes more complicated with the increasing value of n . The LP index [63] outperforms the proximity-based indices, such as RA, AA, and CN.

- (ii) Path of Length 3 (L3) [64] Georg Simmel, a German sociologist, first coined the concept “triadic closure” and made popular by Mark Granovetter in his work [65] “The Strength of Weak Ties”. The authors [64] proposed a similarity index in protein–protein interaction (PPI) network, called path of length 3 (or L3) published in the Nature Communication. They experimentally show that the triadic closure principle (TCP) does not work well with PPI networks. They showed the paradoxical behavior of the TCP (i.e., the path of length 2), which does not follow the structural and evolutionary mechanism that governs protein interaction. The TCP predicts well to the interaction of self-interaction proteins (SIPs), which are very small (4%) in PPI networks and fails in prediction between SIP and non SIP that amounts to 96%. They showed that the L3 index performs well in such conditions and give mathematical expression to compute this index as

$$S(x, y) = \sum_{u,v} \frac{a_{x,u} \cdot a_{u,v} \cdot a_{v,y}}{k_u \cdot k_v}. \tag{42}$$

Recently, Pech et al. [66] in Physica A, proposed a work that models the link prediction as a linear optimization problem. They introduced a theoretical explanation of how direct count of paths of length 3 significantly improves the prediction accuracy. Meanwhile, some more studies [67,68] focusing the length of path have been proposed in the literature. Muscoloni et al. [67] incorporate the concept of local community paradigm (LCP) with 2 and 3 length paths and introduced new similarity indices viz., Cannistraci–Hebb indices $CH2 - L2$ and $CH2 - L3$ corresponding to them. These indices are based on the common neighbor’s rewards to internal local community links (iLCL) and penalization to external local community links (eLCL) [67]. The mathematical expression to compute these two similarity indices are as follows.

$$S^{CH2-L2}(x, y) = \sum_{z \in \Gamma(x) \cap \Gamma(y)} \frac{1 + c_z}{1 + o_z}, \tag{43}$$

where c_z are total number of neighbors of z which are also members of $(\Gamma(x) \cap \Gamma(y))$ and o_z are those neighbors counting that are not in $(\Gamma(x) \cap \Gamma(y))$, also not x or y .

$$S^{CH2-L3}(x, y) = \sum_{z1 \in \Gamma(x), z2 \in \Gamma(y)} \frac{a_{z1,z2} \sqrt{(1 + \tilde{c}_{z1})(1 + \tilde{c}_{z2})}}{\sqrt{(1 + \tilde{o}_{z1})(1 + \tilde{o}_{z2})}}. \tag{44}$$

Here, $a_{z1,z2}$ is 1 when there is link between $z1$ and $z2$, 0, otherwise. \tilde{c}_{z1} is the number of links between $z1$ and the set of intermediate nodes on all 3 length paths between x and y . Similarly, \tilde{o}_{z1} is the number of links between $z1$ and nodes that are not in the set of intermediate nodes of any 3 length path between x and y , also not x or y .

- (iii) Similarity based on Local Random Walk and Superposed Random Walk (LRW and SRW) Liu and Lü [53] proposed new similarity measures by exploiting the random walk concept on graphs with limited walk steps. They defined node similarity based on random walks of lower computational complexity compared to the other random walk based methods [47,53]. Given a random walker, starting from the node x , the probability of reaching the random walker to the node y in t steps is

$$\vec{\pi}_x(t) = P^T \vec{\pi}_x(t-1), \quad (45)$$

where $\vec{\pi}_x(0)$ is a column vector with x th element as 1 while others are 0's and P^T is the transpose of the transition probability matrix P . P_{xy} entry of this matrix defines the probability of a random walker at node x will move to the next node y . It is expressed as $P_{xy} = \frac{a_{xy}}{k_x}$, where a_{xy} is 1 when there is a link between x and y and 0, otherwise. The authors computed the similarity score (LRW) between two nodes based on the above concept as

$$S^{LRW}(x, y) = \frac{k_x}{2|E|} \pi_{xy}(t) + \frac{k_y}{2|E|} \pi_{yx}(t). \quad (46)$$

This similarity measure focus on only few steps covered by the random walker (hence quasi-local) and not the stationary state compared to other approaches [47,53].

Random walk based methods suffer from the situation where a random walker moves far away with a certain probability from the target node whether the target node is closer or not. This is an obvious problem in social networks that show a high clustering index i.e., clustering property of the social networks. This degrades the similarity score between the two nodes and results in low prediction accuracy. One way to counter this problem is that continuously release the walkers at the starting point, which results in a higher similarity between the target node and the nearby nodes. By superposing the contribution of each walker (walkers move independently), SRW is expressed as

$$S^{SRW}(x, y)(t) = \sum_{l=1}^t S^{LRW}(l), \quad (47)$$

Remarks. Similarity-based approaches mostly focus on the structural properties of the networks to compute the similarity score. Local approaches consider, in general, neighborhood information (direct neighbors or neighbors of neighbor), which take less time for computation. This is the property that makes the local approaches feasible for massive real-world network datasets. Global approaches consider the entire structural information of the network; that is why time required to capture this information is more than local and quasi-local approaches. Also, sometimes, entire topological information may not be available at the time of computation, especially in a decentralized environment. So, parallelization over the global approaches may not possible or very complex compared to the local and quasi-local approaches. The performance or prediction accuracy of these approaches (i.e., global approaches) is better compared to local and quasi-local, as shown by the results in Tables 5–8. Quasi-local approaches extract more structural information than local and somehow less information compared to the global. Table 1 shows a simple comparison among similarity-based approaches to link prediction.

2.2. Probabilistic and maximum likelihood models

For a given network $G(V, E)$, the probabilistic model optimizes an objective function to set up a model that is composed of several parameters. Observed data of the given network can be estimated by this model nicely. At that point, the likelihood of the presence of a non-existing link (i, j) is evaluated using conditional probability $P(A_{ij} = 1|\theta)$. Several probabilistic models [69–71] and maximum likelihood models [72,73] have been proposed in the literature to infer missing links in the networks. The probabilistic models normally require more information like node or edge attribute knowledge in addition to structural information. Extracting these attribute information is not easy; moreover, the parameter tuning is also a big deal in such models that limit their applicability. Maximum likelihood methods are complex and time-consuming, so these models are not suitable for real large networks. Some seminal probabilistic and maximum likelihood models are tabulated in the Table 2 [74].

2.2.1. Local probabilistic model for link prediction

Wang et al. [69] proposed a local probabilistic model for link prediction in an undirected network. They employed three different types of features viz., topological, semantic, and co-occurrence probability features extracted from different sources of information. They presented an idea of a central neighborhood set derived from the local topology of the considered node-pair, which is relevant information for the estimation of a link between them. They computed non-derivable frequent itemsets (i.e., those itemsets whose occurrence statistics cannot be derived from other itemset patterns) from the network events log data, which is further used as training data for the model. An event corresponds to a publication of a paper (i.e., authors' interactions in the paper is an event, and a set of such events is the event

Table 2
Probabilistic and maximum likelihood models for link prediction.

Model	Network types	Characteristics	References
Hierarchical structure model (HSM)	Hierarchical networks	High accuracy for HSM and low for non-HSM structure	Clauset et al. [72]
Stochastic block model (SBM)	Noisy networks	Good at predicting spurious and missing links	Guimera et al. [73], Natalie Stanley et al. [75], Toni Valles-Catala et al. [76]
Parametric model	Dynamic networks	Extracts only topological features and performs better than structural methods	Kashima and Abe [77]
Non-parametric model	Dynamic networks	Explicitly clusters links instead of nodes	Sinead A. Williamson [78]
Local probabilistic model	Coauthorship networks	Combines co-occurrence features with topological and semantic features	Wang et al. [69]
Factor graph model	Heterogeneous social networks	Link prediction with aggregate statistics problem	Kuo et al. [79]
Affiliation model	Information and Social networks	soft-block assignment of each node	Jaewon Yang et al. [80]

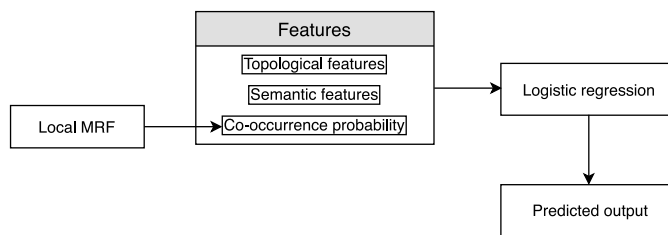


Fig. 4. Local probabilistic model for link prediction [74].

log) in the Coauthorship network. The event log consists of transactional.¹ data upon which frequent itemset mining approaches [81–85] are applied. The model [69] is shown in Fig. 4, which considers the following approach given below.

First, the central neighborhood set between x and y is calculated based on local event log data. One of the usual ways to find the central neighborhood set is to find the shortest path between two vertices of specified length, and the vertices are lying on this path can be included in the required set. There can be more than one shortest path between two vertices, so more neighborhood sets can be possible. Neighborhood sets of shorter lengths and more frequent (frequency score is used when more shortest paths of the same length are available) are chosen for the central neighborhood set. The authors considered the shortest path up to length 4 since the nodes lying on the shorter length path are more relevant.

In the second step, for a given central neighborhood set, non-derivable frequent itemsets are used to learn the local probabilistic model. Calders et al. [86] proposed a depth-first search method to calculate non-derivable itemsets and the same algorithm used by the authors [69]. [Why non-derivable frequent itemsets? Pavlov et al. [87] first introduced the concept of frequent itemset to construct an MRF [88]. They argued that a \mathcal{H} -itemset and its support represents a \mathcal{H} -way statistics, which can be viewed as a constraint on the true underlying distribution that generates the data. Given a set of itemset constraints, a maximum entropy distribution satisfying all these constraints is selected as the estimate for the true underlying distribution. This maximum entropy distribution is equivalent to an MRF. Since the number formed links are very few compared to all possible links in a sparse network, the authors [69] used a support threshold of one to extract all frequent itemsets. These extracted itemsets are large in number that results in expensive learning for the MRF. To reduce this cost, only non-derivable itemsets are extracted]. They find all such itemsets that lie entirely within the central neighborhood set. Using these itemsets [89], a Markov Random Field is learned.

In the last step, the iterative scaling algorithm [69] is used to learn a local MRF for the given central neighborhood set. This process continues overall itemset constraints and continuously updates the model until the model converges. Once the model learning process is over, one can infer the co-occurrence probability by computing the marginal probability over the constructed model. The Junction tree inference algorithm [90] is used to infer co-occurrence probability. The algorithm to induce co-occurrence probability feature for a pair of vertices can be found in [69].

¹ Typically, social networks are the results of evolution of chronological sets of events (e.g., authors participation in the Coauthorship networks). A transaction dataset consists of such events as described by [69].

2.2.2. Probabilistic relational model for link prediction (PRM)

Existing works show that node attributes play a significant role to improve the link prediction accuracy. However, no generic framework is available to incorporate node and link attributes and hence, not applicable to all scenarios. To this end, the probabilistic model is a good and concrete solution that provides a systematic approach to incorporate both node and link attributes in the link prediction framework. Pioneering works on PRM include Getoor et al. [14] study on directed networks, Taskar et al. [91] study on undirected networks, Jennifer Neville work on [70] for both networks, etc. [14] published in JMLR is based on Relational Bayesian network (RBN) where relation links are directed and [91] published in NIPS is based on Relational Markov network (RMN) where relational links are undirected.

PRM was originally designed for attribute prediction in relational data, and it later extended to link prediction task [14,70,91]. The authors employed the attribute prediction framework to link prediction. This casting can be understood with the following example [27]. Consider the problem of link prediction in a coauthorship network. Non-relational frameworks of link prediction consider only one entity type “person” as node and one relationship; however, relational frameworks (PRMs) include more entity types like article, conference venue, institution, etc. Each entity can have attributes like a person (attributes: name, affiliation institute, status (student, professor)), article (attributes: publication year, type (regular, review)), etc. Several relational links may possible among these entities like advisor–advisee/research scholar relation between two persons, author relationship between person and paper entities, and paper can be related to the conference venue with publish relationship. Moreover, relationships (links) among these entities can also have attributes viz., exists (if there is a link between the two involved entities), or not-exist (no link between the involved entities). This way, the link prediction can be reduced to an attribute prediction framework/model.

During the model training, a single link graph is constructed that incorporates above heterogeneous entities and relationships among them. Model parameters are estimated discriminately to maximize the probability of the link existence and other parameters with the given graph attribute information. The learned model is then applied using probabilistic inference to predict missing links. More details can be explored in [14,70,91].

2.2.3. Hierarchical structure model (HSM) [72]

These models are based on the assumption that the structures of many real networks are hierarchically organized, where nodes are divided into groups, which are further subdivided into subgroups and so forth over multiple scales. Some representative work [72] systematically encodes such structures from network data to build a model that estimates model parameters using statistical methods. These parameters are then used in estimating the link formation probability of unobserved links.

Some studies suggest that many real networks, like biochemical networks (protein interaction networks, metabolic networks, or genetic regulatory networks), Internet domains, etc. are hierarchically structured. In hierarchical networks, vertices are divided into groups, which are further sub-divided into subgroups and so forth over multiple scales [92]. Clauset et al. [72] proposed a probabilistic model that takes a hierarchical structure of the network into account. The model infers hierarchical information from the network data and further applies it to predict missing links.

The hierarchical structures are represented using a tree (binary), or dendrogram, where, the leaves (i.e., n) represent the number of total vertices in the network and each internal vertex out of $(n - 1)$ corresponds to the group of vertices descended from it. Each internal vertex r is associated with a probability p_r , then the existing edge probability p_{xy} between two vertices x and y is given by $p_{xy} = p_r$ where, r is their lowest common ancestor. The hierarchical random graph is then, represented by the dendrogram D^* with the set of probability $\{p_r\}$ as $(D^*, \{p_r\})$. Now the learning task is to find the hierarchical random graph(s) that best estimates the observed real-world network data. Assuming all possible dendrograms to be equally likely, Bayes theorem says that the probability of the dendrogram $(D^*, \{p_r\})$ that best estimates the data is proportional to the posterior probability or likelihood, L from which the model generates the observed network and our goal is to maximize L . The likelihood of a hierarchical random graph $(D^*, \{p_r\})$ is computed using the following equation

$$L(D^*, \{p_r\}) = \prod_{r \in D^*} p_r^{E_r} (1 - p_r)^{L_r R_r - E_r}, \quad (48)$$

where L_r and R_r are the left and right subtree rooted at r , and E_r is the number of links in the network whose endpoints have r as their lowest common ancestor in D^* . The above equation assumes the convention $0^0 = 1$. For a given dendrogram D^* , it is easy to compute the probability \bar{p}_r that maximizes $L(D^*, \{p_r\})$ i.e.

$$\bar{p}_r = \frac{E_r}{L_r R_r}. \quad (49)$$

This can be understood with the following example illustrated in Fig. 5 Now, this model can be used to estimate the missing links of the network as follows. Sample a large number of dendrograms with probability proportional to their likelihood. Then, compute the mean connecting probability \bar{p}_{xy} of each non-existing pair (x, y) by averaging the corresponding probability p_{xy} overall sampled dendrograms. Sort these vertices pairs scores in descending order and selects top- l links to be predicted.

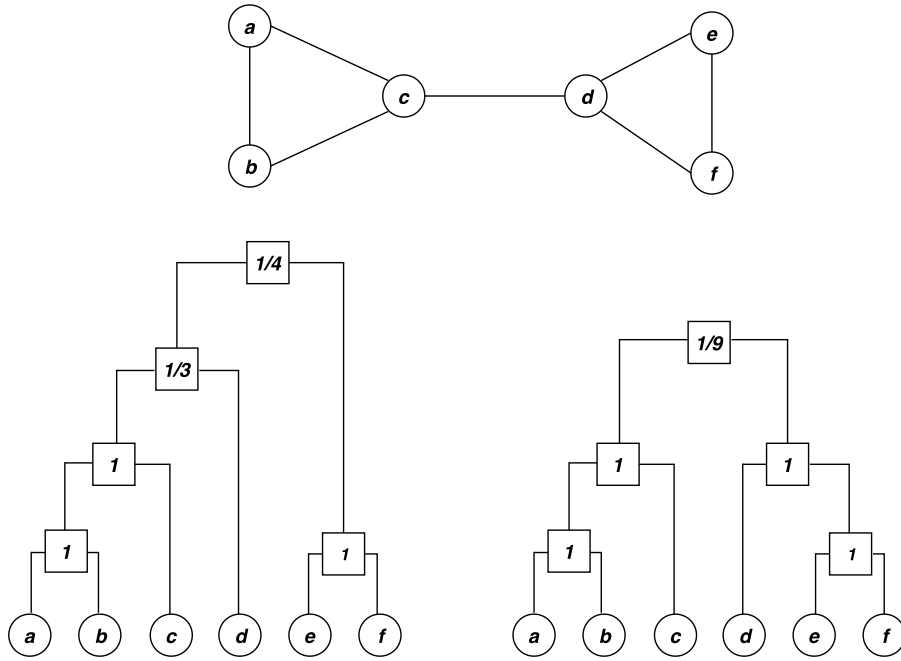


Fig. 5. An illustrating example of HSM for a graph of 6 nodes and its two possible dendrograms as described in the paper [72]. The internal nodes of each dendrogram are labeled as the maximum likelihood probability \bar{p}_r , defined by Eq. (49). The likelihoods of the left and the right dendrograms are $L(D_1) = (1/3)(2/3)^2 \cdot (1/4)^2(3/4)^6 = 0.00165$, and $L(D_2) = (1/9)(8/9)^8 = 0.0433$. Thus, the second (i.e., right) dendrogram is most probable as it divides the network in a balanced one at the first level.

2.2.4. Stochastic block model (SBM) [73]

Hierarchical structures may not represent most networks. A more general approach to represent these networks is block model [93,94] where vertices are distributed (partitioned) into blocks or communities and the connecting probability between two vertices depends on blocks they belong to. Guimerà et al. [73] presented a novel framework where stochastic block model representation of a network is employed to find missing and spurious links. The authors compute the reliability of the existence of links given an observed network that is further used to find missing links (non-existing links with higher reliabilities) and spurious links (existing links with lower probabilities).

The link reliability R_{xy} between the two vertices x and y is [73]

$$R_{xy} = p_{BM}(A_{xy} = 1|A^0).$$

i.e. probability that the link truly exists given the observed network A^0 , the block model BM .

Generally, complex networks are outcomes of combination of mechanisms, including modularity, role structure, and other factors. In SBM, partitioning vertices of network based on these mechanisms may result in different block models that capture different correlations (patterns) of the network. Assume that no prior knowledge of suitable models, the reliability is expressed as

$$R_{xy} = \frac{1}{Z} \sum_{P \in P^*} \left(\frac{l_{\sigma_x \sigma_y}^0 + 1}{r_{\sigma_x \sigma_y}^0 + 2} \right) \exp[-H(P)], \tag{50}$$

where the sum is over all possible partitions P^* of the network into groups, σ_x and σ_y are vertices x and y groups in partition P respectively. Moreover, $l_{\sigma_\alpha \sigma_\beta}^0$ and $r_{\sigma_\alpha \sigma_\beta}^0$ are the number of links and maximum possible links in the observed network between groups α and β . The function $H(P)$ is

$$H(P) = \sum_{\alpha < \beta} [\ln(r_{\alpha\beta}) + \ln \binom{r_{\alpha\beta}}{l_{\alpha\beta}^0}], \tag{51}$$

and $Z = \sum_{P \in P^*} \exp[-H(P)]$. Practically, solving equation (50), i.e., summing over all possible partitions is too complex even for a small network. However, the Metropolis algorithm [95] can be used to correctly sample the relevant partitions and obtain link reliability estimates.

The authors employed the link reliability concept to find missing links and to identify the spurious link in the networks with the following procedure. (i) Generate the observed network A^0 by removing/adding some random links (for finding

missing/spurious links) from/to the true network A^t . (ii) Compute the link reliability for non-observed links (i.e. non-existing + missing/spurious links). (iii) Arrange these links with their reliability score in decreasing order and decide the top- l links as desired ones (i.e., missing/spurious links).

Probabilistic and maximum likelihood methods extract useful features and valuable correlation among the data using hierarchical and stochastic block models, which result in significant improvements in prediction results as compared to some similarity-based methods. However, these are quite complex and time-consuming even on small datasets that limit their applicability on large scale real-world network datasets.

2.2.5. Exponential random graph model (ERGM) or P-star model

Exponential random graphs were first studied by Holland and Leinhardt [96], further explored by [88], and practically used by several works [97–99]. ERGM is an ensemble model where one defines it as consisting of a set of all simple undirected graphs and specifies a probability corresponding to each graph in the ensemble. Properties of the ERGM is computed by averaging over the ensemble [98]. Pan et al. [99] also proposed a similar probabilistic framework (ERGM) to find missing and spurious links in the network. They employed predefined structural Hamiltonian for the score computation. The Hamiltonian is selected based on some organizing principle such that the observed network can have lower Hamiltonian than its randomized one. They defined the structure Hamiltonian by generalizing the 3-order loop to higher-order as

$$H(A) = - \sum_{l=3}^{\infty} \beta_l \ln(\text{Tr}(A^l)), \quad (52)$$

where A is the adjacency matrix of the network, β_l is temperature parameter. Here, the number of loops of length l starting and ending at the node i is $[A^l]_{ii}$. For undirected network, loops are counted several times when counting occurs for each involved node of the loop, also, for a given node it is counted twice (clockwise and anti-clockwise). Therefore, $\text{Tr}(A^l)$ counts approximated to $2l$ times the number of loops of length l that can be taken care of by the parameter β_l [99].

For large value of l , increment in $\text{Tr}(A^l)$ reaches to the leading eigen value λ_1 and small world phenomenon of a social network ensures to have l to a lower cut-off l_c .

$$H(A) = - \sum_{l=3}^{l_c} \beta_l \ln\left(\sum_{i=1}^n \lambda_i^l\right) \quad (53)$$

Note that the above equation is result of diagonalization of the adjacency matrix A^l as follows

$$\begin{aligned} \text{Tr}(A^l) &= \text{Tr}(U^T A^l U) \\ &= \text{Tr}(A^l U U^T) = \text{Tr}(A^l) = \sum_{i=1}^n \lambda_i^l \end{aligned}$$

Once, the structural Hamiltonian is defined to capture different parameters (higher order loops here), the probability of the appearance of the observed network $A^O = A - A^P$ in an ensemble \mathcal{M} is

$$p(A^O) = \frac{1}{Z} \exp[-H(A^O)], \quad (54)$$

where, $Z = \sum_{p \in \mathcal{M}} \exp[-H(A^p)]$ is the partition function. The parameters β_l are chosen to maximize the probability expressed in the above equation.

Now, the score of non-observed links can be computed by the conditional probability of the appearance of link (x, y)

$$S(x, y) = \frac{1}{Z_{xy}} \exp[-H(\tilde{A}(x, y))], \quad (55)$$

where $\tilde{A}(x, y)$ is the observed network by adding the link (x, y) , and Z_{xy} is a normalization factor defined as follows [99]

$$Z_{xy} = \exp[-H(\tilde{A}(x, y))] + \exp[-H(A^O)].$$

Here, the prediction is based on the assumption that there is no significant change in the topological structure after adding the link (x, y) to the observed network and the parameter β_l for $\tilde{A}(x, y)$ is almost similar to that of the observed network A^O .

2.3. Link prediction using dimensionality reduction

The curse of dimensionality is a well-known problem in machine learning. Some researchers [100,101] employ dimension reduction techniques to tackle the above problem and apply it in the link prediction scenario. Recently, many authors are working on network embedding and matrix decomposition techniques, which are also considered as dimension reduction techniques.

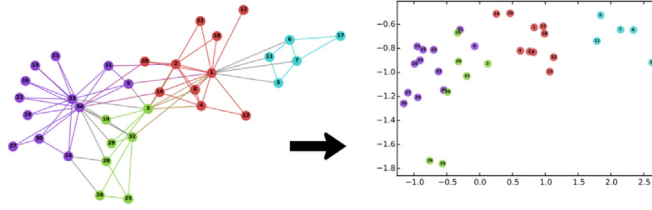


Fig. 6. The Karate club network (left) and its representation in the embedding space with the DeepWalk [102] algorithm.

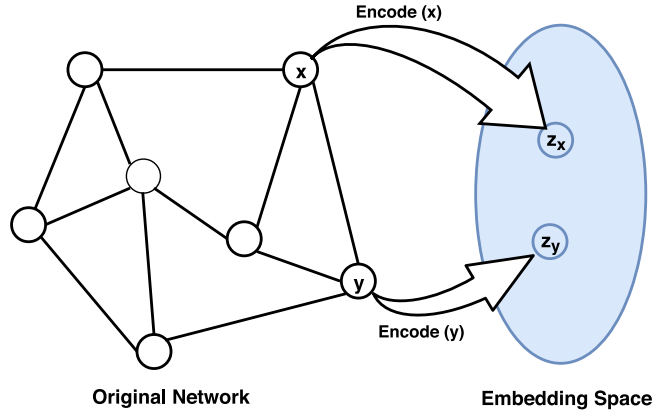


Fig. 7. Embedding of nodes x and y to the embedding space.

2.3.1. Embedding-based link prediction

The network embedding is considered as a dimensionality reduction technique in which higher D dimensional nodes (vertices) in the graphs are mapped to a lower d ($d \ll D$) dimensional representation (embedding) space by preserving the node neighborhood structures. In other words, find the embedding of nodes to a lower d -dimensions such that similar nodes (in the original network) have similar embedding (in the representation space). Fig. 6 shows the structure of Zachary Karate club social network (left) and the representation of nodes in the embedding space using DeepWalk [102] (right). The nodes are colored based on the membership of their communities (See Fig. 6).

The main component of the network embedding is the encoding function or encoder f_{en} that map each node to the embedding space as shown in Fig. 7.

$$f_{en}(x) = z_x, \tag{56}$$

where z_x is the d -dimensional embedding of the node x . The embedding matrix is $Z \in \mathbb{R}^{d \times |V|}$, each column of which represents an embedding vector of a node. Now, a similarity function is $S(x, y)$ is defined that specifies how to model the vector (embedding) space relationships equivalent to the relationships in the original network, i.e.,

$$S(x, y) \approx z_x^T z_y. \tag{57}$$

Here, $S(x, y)$ is the function that reconstructs pairwise similarity values from the generated embedding. The term $S(x, y)$ is the one that differ according to the function used in different factorization-based embedding approaches. For example, graph factorization [103] directly employ adjacency matrix A i.e. ($S(x, y) \triangleq A_{(x,y)}$) to capture first order proximity, GraRep [104] selects ($S(x, y) \triangleq A_{(x,y)}^2$) and HOPE [105] uses other similarity measures (e.g. Jaccard neighborhood overlap). Most embedding methods realize the reconstruction objective by minimizing the loss function, L

$$L = \sum_{(x,y) \in \{V \times V\}} l(z_x^T z_y, S(x, y)). \tag{58}$$

Once Eq. (58) is converged (i.e. trained), one can use the trained encoder to generate nodes embedding, which can further be employed to infer missing link and other downstream machine learning tasks.

Recently, some network embedding techniques [102,106–109] have been proposed and applied successfully in link prediction problem. The Laplacian eigenmaps [106], Logically linear embedding (LLE) [109], and Isomap [110,111] are examples based on the simple notion of embedding. such embedding techniques are having quite complex in nature and face scalability issues. To tackle the scalability issue, graph embedding techniques have leveraged the sparsity of real-world networks. For example, DeepWalk [102] extracts local information of truncated random walk and embeds the nodes in

Table 3

Deep learning models for embedding based link prediction.

	Model	Proximity preserved	Embedding type	Scalability	Learning	Reference
With random walk	DeepWalk	Higher order	Shallow	Yes	Unsupervised	[102]
	Node2vec	Higher order	Shallow	Yes	Semi-supervised	[107]
	HARP	Higher order	Shallow	Yes	Supervised	[117]
	Walklets	Higher order	Shallow	Yes	Unsupervised	[118]
Without random walk	LINE	First and second order	Shallow	Yes	Supervised	[119]
	SDNE	First and second order	Deep	No	Semi-supervised	[120]
	DNGR	Higher order	Deep	Yes	Unsupervised	[121]
	GCN	Higher order	Deep	Yes	Semi-supervised	[122]
	VGAE	Higher order	Deep	No	Unsupervised	[123]
	SEAL	First and second order	Deep	Yes	Supervised	[124]
	ARGA	Higher order	Deep	No	Unsupervised	[125]

representation space by considering the walk as a sentence in the language model [112,113]. It preserves higher order proximity by maximizing the probability of co-occurrence of random walk of length $2k + 1$ (previous and next k nodes centered at a given node). Node2vec [107] also uses a random walk to preserve higher order proximity but it is biased which is a trade-off between the breadth-first search (BFS) and depth-first search (DFS). The experimental results show that the Node2vec performs better than the Deepwalk. In next step, Trouillon et al. [114] introduced complex embedding in which simple matrix and tensor factorization have been used for link prediction that uses a vector with complex values. Such composition of complex embedding includes all possible binary relations especially symmetric and anti-symmetric relations. Recently, some more studies have been published in link prediction using embedding, for example, Cao et al. subgraph embedding [115], Li et al. deep dynamic network embedding [116], Kazemi et al. [108], etc. some seminal works in network embedding are listed in Table 3.

2.3.2. Matrix factorization/decomposition-based link prediction

From last decade, matrix factorization has been used in lots of papers based on link prediction [126–133] and recommendation systems [134]. Typically, the latent features are extracted and using these features, each vertex is represented in latent space, and such representations are used in a supervised or unsupervised framework for link prediction. To further improve the prediction results, some additional node/link or other attribute information can be used. In most of the works, non-negative matrix factorization has been used. Some authors also applied the singular value decomposition technique [135].

Let the input data matrix is represented by $X = (x_1, x_2, \dots, x_n)$ that contains n data vectors as columns. Now, factorization of this matrix can be expressed as

$$X \approx FG^T, \quad (59)$$

where $X \in \mathbb{R}^{p \times n}$, $F \in \mathbb{R}^{p \times k}$, and $G \in \mathbb{R}^{n \times k}$. Here, F contains the bases of the latent space and is called the basis matrix. G contains combination of coefficients of the bases for reconstructing the matrix X , and is called the coefficient matrix. k is the dimension of latent space ($k < n$). Several well-known matrix factorizations are expressed based on some constraints on either of the three matrices, for example, [136],

SVD:

$$X_{\pm} \approx F_{\pm} G_{\pm}^T. \quad (60)$$

NMF:

$$X_{+} \approx F_{+} G_{+}^T. \quad (61)$$

Semi-NMF:

$$X_{\pm} \approx F_{\pm} G_{+}^T. \quad (62)$$

Convex-NMF:

$$X_{\pm} \approx X_{\pm} W_{+} G_{\pm}^T. \quad (63)$$

In the above four equations, Z_{\pm} represents the nature of the entries in the matrix Z , i.e. both positive and negative entries allowed in the matrix Z . In the last equation, $F = XW$ represents the convex combinations of the columns of F . Generally, such a factorization problem can be modeled as the following Frobenius norm optimization problem

$$\begin{aligned} \min_{f,g} \quad & \|X - FG^T\|_{fro}^2 \\ \text{subject to} \quad & F \geq 0, G \geq 0. \end{aligned} \quad (64)$$

Here, $\|Z\|_{fro}^2$ is the Frobenius norm of Z and the constraints represent NMF factorization. However, any of the above four constraints can be used depending on the requirement of the problem underlying.

After solving the above optimization problem, the similarity between a non-existing pair (x, y) can be computed by the similarity of the x th and y th row vectors in the coefficient matrix G .

Acar et al. [126] expressed temporal link prediction as a matrix completion problem and solve it through the matrix and tensor factorization. They proposed a weighted method to collapsed the temporal data in a single matrix and factorize it using CANDECOMP/PARAFAC (CP) [137,138] tensor decomposition method. Ma et al. [127] also applied matrix factorization to temporal networks where features of each network are extracted using graph communicability and then collapsed into a single feature matrix using weighted collapsing tensor (WCT) [128]. They showed the equivalence between eigen decomposition of Katz matrix and non-negative matrix factorization (NMF) of the communicability matrix that serves as the foundation of their framework. Further, a notable work by Menon et al. [129] is proposed for structural link prediction. Here, the problem is modeled as matrix completion problem [139], and matrix factorization are used to solve it. They introduced a supervised matrix decomposition framework that learns latent (unobserved) structural features of the graph and incorporates it with additional node/link explicit feature information to make a better prediction. Additionally, they allowed the factorization model to solve class imbalance problem [140] by optimizing ranking loss. Chen et al. [130] proposed somehow similar to work [129], where the authors extracted topological matrix and attribute matrix and factorized these matrices using non-negative matrix factorization. The final score matrix is obtained by integrating these two matrices in the latent space. Recently some more works [131–133] have been published in this area.

2.4. Other approaches

2.4.1. Learning-based frameworks for link prediction

Earlier described approaches (e.g., similarity and probabilistic methods) deal with the computing a score of each non-observed link either by a similarity or a probabilistic function. However, the link prediction problem can also be modeled as a learning-based model to exploit graph topological features and attribute information. The problem is cast as a supervised classification model where a point (i.e., training data) corresponds to a vertex-pair in the network, and the label of the point represents the presence or absence of an edge (link) between the pair. In other words, consider a vertex-pair (x, y) in the graph $G(V, E)$ and the label of the corresponding data point in the classification model is $l_{(x,y)}$. Then,

$$l_{(x,y)} = \begin{cases} +1 & \text{if } (x, y) \in E, \\ -1 & \text{if } (x, y) \notin E. \end{cases} \quad (65)$$

This is typically a binary classification task where several classifiers (e.g., decision tree, naive Bayes, support vector machine, etc.) can be employed to predict the label of unknown data points (corresponding to missing links in the network).

One of the major challenges of this model (i.e., machine learning) is the selection of appropriate feature set [27]. Majority of the existing research works [2,9,77] extract feature sets from the network topology (i.e., topological information of the network). These features are generic and domain-independent that are applicable to any network. Such features are typical, neighborhood, and path-based features. Some other works [9,141] concentrate on extracting node and edge features that play a crucial role to improve the performance of link prediction. Hasan et al. [9] employed vertex attribute viz., the degree of overlap among research keywords incorporated with other features in the coauthorship network, and showed that the author-pairs having higher values of these features are top rankers in the list. The cost of extraction of such features is cheap and easy, while the main disadvantage is the domain-specific nature of them.

2.4.2. Information theory-based link prediction

Several complex networks have utilized the concept of information theory to compute their complexity on different scales [142,143]. They defined several correlation measures and modeled some networks (e.g., star, tree, lattice, ER graph, etc.). They also showed that the real networks spanned noise entropy space. Bauer et al. [144] used the maximum entropy principle to assign a statistical weight to any graph and introduced random graph construction with arbitrary degree distribution.

Tan et al. [145] posed the link prediction problem in the framework of information theory. They mainly focus on local assortativity to capture local structural properties of the network and showed that mutual information (MI) method performs well on both low and highly correlated networks. Motivated by [145], Zhu, B. and Xia [146] added more local features (i.e., links information of neighbors of the seed nodes as well as their common neighbors) in their framework and called it as neighbor set information (NSI) index. Thus, they showed that the different features could be combined in an information-theoretic model to improve the link prediction accuracy.

Xu et al. [147] considered path entropy as a similarity metric for the link prediction problem. The authors assumed that there is no correlation among the degrees of the nodes in the network. Consider the following notations based on

the paper [147]: L_{xy}^0 shows no link exists between two vertices x and y , and the corresponding existence is represented by L_{xy}^1 . Probability of existence of a link between the above two vertices is given as

$$P(L_{xy}^1) = 1 - P(L_{xy}^0) = 1 - \frac{C_M^{k_y - k_x}}{C_M^{k_y}}, \quad (66)$$

where $C_M^{k_y}$ represents the number of candidate link sets for the vertex y with all links incident with y and $C_M^{k_y - k_x}$ denotes the number of candidate link sets for the vertex y with all links incident with y but none of them is incident with x .

$$I(L_{xy}^1) = -\log P(L_{xy}^1) = -\log\left(1 - \frac{C_M^{k_y - k_x}}{C_M^{k_y}}\right) \quad (67)$$

They show that the likelihood of occurrence of a path having no loops equates to multiplication of the occurrence probabilities of the links involved in that path. i.e., given a simple path $D = v_0, v_1, v_2, \dots, v_\gamma$ of length γ , the co occurrence probability of path D is evaluated to

$$P(D) \approx \prod_{i=0}^{\gamma-1} P(L_{v_i v_{i+1}}^1) \quad (68)$$

and, the sum of links entropies involved in a path equals to the entropy of the path.

$$I(D) \approx \sum_{i=0}^{\gamma-1} I(L_{v_i v_{i+1}}^1). \quad (69)$$

Further, they calculated similarity based on entropy of the path which is the negative of conditional entropy

$$S_{xy}^{PE} = -I(L_{xy}^1 | \cup_{i=2}^{\max len} D_{xy}^i), \quad (70)$$

where D_{xy}^i represents the set consisting of all simple paths of length i between the two vertices and $\max len$ is the maximum length of simple path of the network. Outcome results on several networks demonstrate that the similarity index based on path entropy performs better than other indices in terms of prediction accuracy and precision. Xu et al. [148] extend the previous work [147] to the weighted network by considering the weight of the paths. Recently, some more efforts have been applied in this direction based on different features of the networks like influential nodes [149], combining node attributes with structural similarity [150], local likelihood [151], and maximal entropy random walk [152].

2.4.3. Clustering-based link prediction

This paragraph gives an overview of the clustering-based link prediction. Huang [153] presented a paper on graph topology-based link prediction where a generalized clustering coefficient is used as a predictive parameter. The author introduces a cycle formation model that shows the relationship between link occurrence probability and its ability to form different length cycles. This model suggests that the occurrence probability of a particular link depends on the number of different lengths cycles formed by adding this link. The model is based on the assumption of the stationary property of the degree of clustering of the network [154]. This model captures longer cycles by extending the higher-order clustering coefficients [155] and defines the generalized clustering coefficient $C(k)$ as

$$C(k) = \frac{\text{number of } k\text{-length cycles}}{\text{number of } k\text{-length paths}}, \quad (71)$$

where k is the degree of the cycle formation model.

The author treats the link occurrence probability as governed by t link generation mechanisms $g(1), g(2), \dots, g(k)$ of cycle formation model, each described by a single parameter c_1, c_2, \dots, c_k . The above mentioned link generation mechanism can be understood with the help of Fig. 8. Consider a cycle formation model ($CF(k)$) of degree ($k = 3$). The Seed link (x, y), here, can be generated by the following three mechanisms; the random link occurrence $g(1)$, length-2 cycle generation $g(2)$ i.e. ($x - a - y$ and $x - c - y$), and length-4 cycle generation $g(3)$ i.e. ($x - b - d - y$). The main issue is to combine several generation mechanisms to compute total link occurrence probability. The author [153] posits a method to combine both path and cycle (of different lengths) generation mechanism in the framework. The expected general clustering coefficient of degree k for this model can be estimated as [153]

$$\begin{aligned} E[C(k)] &= f(c_1, c_2, \dots, c_k) \\ &= \sum_i |G_i| p(G_i) p((e_{l,k+1}) \in E|G_i), \end{aligned} \quad (72)$$

where $|G_i|$ is the number of subgraph possible corresponding to the graph pattern G_i , listed in Table 1 of the paper [153], $p(G_i)$ is the probability of occurrence of one of such graphs G_i , and $p(e_{l,k+1})$ is the probability of edge $e_{l,k+1}$ to occur given

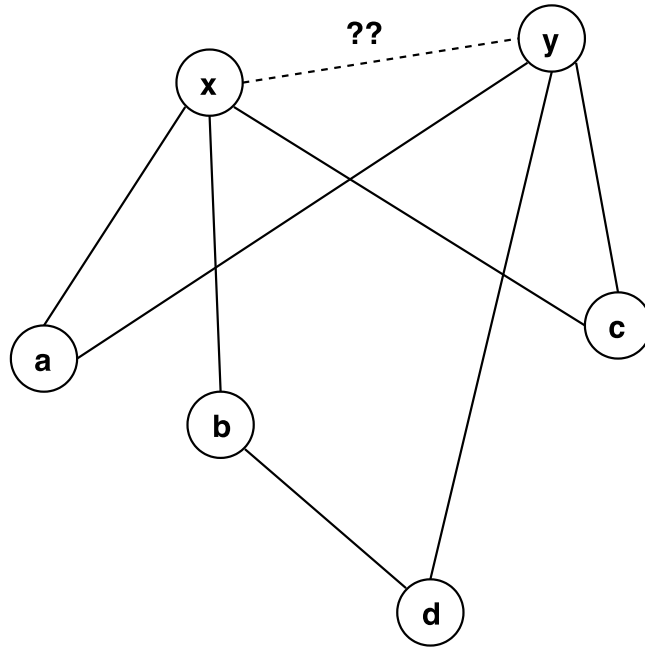


Fig. 8. An example illustrating the cycle formation link probability model [153], where the probability of the missing link $(x - y)$ is generated by the following three mechanisms; random link occurrence $g(1)$, length-2 cycle generation $g(2)$ i.e. $(x - a - y, x - c - y)$, and length-4 cycle generation $g(3)$ i.e. $(x - b - d - y)$.

the pattern G_i . Finally, given the coefficients, the probability of existence of link is

$$p_{x,y}(c_1, \dots, c_k) = \frac{c_1 \prod_{i=2}^k c_i^{|path_{x,y}^i|}}{c_1 \prod_{i=2}^k c_i^{|path_{x,y}^i|} + (1 - c_1) \prod_{i=2}^k (1 - c_i)^{|path_{x,y}^i|}} \tag{73}$$

Liu et al. [156] proposed degree related clustering coefficient to quantify the clustering ability of nodes. They applied the same to paths of shorter lengths and introduced a new index Degree related Clustering ability Path (DCP). They performed the degree of robustness (DR) test for their index and showed that missing links have a small effect on the index. Recently Wu et al. [42] extracted triangle structure information in the form of node clustering coefficient of common neighbors. Their experiments on several real datasets show comparable results to the CAR index in [38]. The same concept of the clustering coefficient also introduced in the work presented by Wu et al. [43]. Authors introduce both node and link clustering information in their work [43]. Their experiments on different network datasets showed better performance results against existing methods, especially on middle and large network datasets. Kumar et al. [157] explored the concept of node clustering coefficient to the next level (level-2) that captures more clustering information of a network. The comprehensive results on several real-world datasets show better performance compared to local methods and comparable to the node embedding method Node2vec [107]. Meanwhile, Benson et al. [158] studied simplicial closure events to capture higher-order structures in several temporal networks. The simplicial closure events are the process of closure of timestamped simplices (simplicial complexes² are set of nodes with different sizes) available in a dataset. These structures are common in several real-time complex systems, for example, communication in a group, collaboration of authors for a paper, etc. To assess these higher-order structures, the authors study the simplicial closure events on triples of nodes (for simplicity) and suggest that the open triangles or triples of nodes with strong ties are more likely to close in the future.

2.4.4. Structural perturbation method (SPM) [159]

Lü et al. introduced a new framework of computing predictability of links in the networks. They coined a structural consistency index to quantify the link predictability. This index is based on the assumption that links in a network are highly predictable if no significant changes occur in the structural feature after the addition or deletion of a small fraction of the link. Based on this index, they proposed a new similarity index, namely structural perturbation method (SPM). The experimental results show the outstanding performance compared to the state-of-the-art in their paper.

² https://en.wikipedia.org/wiki/Simplicial_complex.

		Actual Class	
		Link Available	Link Not Available
Prediction Class	Predicted	True Positive (TP)	False Positive (FP)
	Not Predicted	False Negative (FN)	True Negative (TN)

Fig. 9. Confusion matrix.

3. Experimental setup and results analysis

3.1. Evaluation metrics

The link prediction problem is treated as a binary classification task [9], so most of the evaluation metrics of any binary classification task can be used in link prediction evaluation. The evaluation of a binary classification task having two classes can be represented as a confusion matrix [160], as given in Fig. 9.

In the confusion matrix,

- True Positive (TP): The positive data item (Link Available) predicted as positive (Predicted).
- True Negative (TN): The negative data item (Link Not Available) predicted as negative (Not Predicted).
- False Positive (FP): The negative data item (Link Not Available) predicted as positive (Predicted).
- False Negative (FN): The positive data item (Link Available) predicted as negative (Not Predicted).

Based on the confusion matrix, several metrics can be derived as follows [160].

True Positive Rate (TPR)/Recall/Sensitivity

$$TPR = \frac{\#TP}{\#TP + \#FN}. \quad (74)$$

False Positive Rate (FPR)

$$FPR = \frac{\#FP}{\#FP + \#TN}. \quad (75)$$

True Negative Rate (TNR)/Specificity

$$TNR = \frac{\#TN}{\#TN + \#FP}. \quad (76)$$

$$Precision = \frac{\#TP}{\#TP + \#FP}. \quad (77)$$

In the above equations # represents 'the number of'.

Our approach is evaluated on four metrics viz., Area under the ROC curve (AUROC) [161,162], Area under the precision–recall curve (AUPR) [163], Average precision [160] and Recall@k [160].

3.1.1. Area under the receiver operating characteristics curve (AUROC)

A roc curve is a plot between the true positive rate (sensitivity) on Y-axis and the false positive rate (1-specificity) on the X-axis. The true positive rate and false positive rate can be evaluated using Eqs. (74) and (75) respectively. Sensitivity is a performance of the whole positive part, and specificity is a performance of the whole negative part of a dataset. The area under the roc curve [162] is a single point summary statistics between 0 and 1 that can be computed using the trapezoidal rule which sums all the trapezoids under the curve. The value of the auroc of a predictor should be greater than 0.5, which is the value of a random predictor, i.e., higher the value of auroc better the performance of the predictor.

3.1.2. Area under the precision–recall curve (AUPR)

AUPR is also a single point summary statistics used to evaluate the performance of a binary classifier (predictor). This value is computed based on the precision–recall curve, which is a plot between the precision values on Y-axis and the recall values on X-axis. The precision and recall values can be computed using Eq. (77) and Eq. (74) respectively. The precision–recall curve is more useful and informative when applied to binary classification on imbalanced datasets [164]. A higher value of aupr of a model represents the better model.

Table 4
Topological information of real-world network datasets.

Datasets	$ V $	$ E $	$\langle D \rangle$	$\langle K \rangle$	$\langle C \rangle$
Karate	34	78	2.337	4.588	0.570
Dolphin	62	159	3.302	5.129	0.258
Macaque	91	1401	1.658	30.791	0.742
Football	115	613	2.486	10.661	0.403
Jazz	198	2742	2.235	27.697	0.620
C. Elegans	297	2148	2.447	14.456	0.308
USAir97	332	2126	2.738	12.807	0.749
Netscience	1589	2742	5.823	3.451	0.878

3.1.3. Average precision

This metric is also a single point summary value computed based on varying threshold³ values of recall. The average precision value is equal to the precision averaged over all values of recall between 0 and 1, i.e.,

$$\text{Average precision} = \int_{r=0}^1 p(r)dr,$$

where p is the precision at different threshold value of recall r .

Practically, integral is approximated to sum over the precisions at each threshold value, multiplied by the change in the recall, i.e.,

$$\text{Average precision} = \sum_{k=1}^R p(k)\Delta r(k), \quad (78)$$

where R is the set of different threshold values.

3.1.4. Recall@k

This metric⁴ is almost same as the metric given in Eq. (74), but it considers only top- k data items instead.

3.2. Datasets

This work used 8 network datasets from various fields to study the performance of our approach. Karate⁵ [165]: A friendship network of 34 members of a Karate club at a US university. Dolphin³ [166]: A social network of dolphins living in Doubtful Sound in New Zealand. Macaque⁶ [167]: is a biological network of cerebral cortex of Rhesus macaque. Football³ [168]: American football games network played between Division IA colleges during regular season Fall 2000. Jazz⁷ [169]: A collaboration network of 115 jazz musician where a link between two musicians denotes music played by both in a band. C. Elegans³ [44]: A neural network of C. Elegans compiled by D. Watts and S. Strogatz in which each node refers a neuron and, a link joins two neurons if they are connected by either a synapse or a gap junction. USAir97⁸ is an airline network of US where a node represents an airport, and a link shows the connectivity between two airports. Netscience³ [170] is a Coauthorship network of researchers in the network theory domain where a node is denoted by a researcher, and an edge denotes coauthorship of at least one paper between two researchers.

Table 4 shows some basic topological properties of the considered network datasets. $|V|$ and $|E|$ are the total numbers of nodes and links of the networks, respectively. $\langle D \rangle$ represents the average shortest distance, $\langle K \rangle$, the average degree, and $\langle C \rangle$, the average clustering coefficient of the network.

3.3. Accuracy results

Four accuracy measures have been used to evaluate each similarity-based algorithm and some other representative methods. We report these results in Tables 5–8 for similarity-based approaches and Tables 9–12 for other representative methods. In the tables of other representative methods, the first method (i.e., HSM) is the maximum likelihood-based method followed by the next three embedding-based methods followed by the three clustering methods. The last method of the table belongs to other category. The best results are highlighted in the table on each dataset. These results are generated with the help of code implemented by Gregorio Alanis-Lobato.⁹

³ <https://sanchom.wordpress.com/tag/average-precision/>.

⁴ https://ils.unc.edu/courses/2013_spring/inls509_001/lectures/10-EvaluationMetrics.pdf.

⁵ <http://www-personal.umich.edu/~mejn/netdata/>.

⁶ <https://neurodata.io/project/connectomes/>.

⁷ <http://deim.urv.cat/~alexandre.arenas/data/welcome.htm>.

⁸ <http://vlado.fmf.uni-lj.si/pub/networks/data/>.

⁹ <https://github.com/galanis/LinkPrediction>.

Table 5
Recall results.

Methods	Karate	Macaque	C. Elegans	Football	Jazz	USAir97	Dolphin	Netscience
CN	0.11363	0.30152	0.08918	0.23800	0.50078	0.40944	0.12500	0.52000
JC	0	0.01908	0.02540	0.31800	0.52007	0.07722	0.08750	0.60322
AA	0.10000	0.28320	0.09945	0.22600	0.52125	0.40444	0.11250	0.67419
RA	0.02500	0.27328	0.09513	0.23400	0.52795	0.45833	0.13125	0.70709
PA	0.05000	0.33053	0.05459	0	0.10984	0.32611	0.03125	0.00129
SALTON	0.17500	0.27709	0.09891	0.23600	0.51181	0.39000	0.13750	0.52129
SORENSEN	0.15000	0.27633	0.08648	0.25200	0.50039	0.37333	0.13125	0.53870
CAR	0.20000	0.27251	0.09243	0.26200	0.51850	0.38333	0.12500	0.54129
CAA	0.15000	0.28015	0.10594	0.33600	0.52362	0.38611	0.13750	0.58000
CRA	0.07500	0.27709	0.11459	0.31800	0.56732	0.43888	0.13125	0.61806
CPA	0.11111	0.29313	0.10108	0.20400	0.48858	0.38555	0.08750	0.33225
HPI	0.10000	0.28167	0.07783	0.23800	0.51259	0.40444	0.10625	0.50451
HDI	0.15000	0.28473	0.09405	0.24600	0.48897	0.38277	0.18750	0.52129
NLC	0.05000	0.30763	0.07351	0.08800	0.44803	0.39444	0.05000	0.00516
LNBCN	0.10000	0.00305	0.08972	0.26200	0.38897	0.40555	0.11250	0.58322
LHNL	0.12500	0.27251	0.09297	0.24400	0.49409	0.41222	0.09375	0.52451
CCLP	0.05000	0.28015	0.09675	0.29000	0.52244	0.40555	0.10625	0.60000
KATZ	0.05000	0.34122	0.08162	0.20600	0.44212	0.39722	0.10625	0.43741
RWR	0.10000	0.38855	0.10216	0.24600	0.33937	0.08222	0.06000	0.30925
Shortest Path	0	0.09160	0.02702	0.03200	0.02007	0.02111	0.02000	0.13868
LHNG	0	0	0	0.36200	0.10669	0.00388	0.01250	0.05185
ACT	0.02500	0.30076	0.04972	0.03600	0.15748	0.33444	0.02000	0.20740
NACT	0	0	0.00540	0.32800	0.33740	0.00888	0	0.34024
Cos ⁺	0.02500	0.20610	0.04540	0.30400	0.13464	0.01888	0.02000	0.07037
MF	0.05000	0.19923	0.04324	0.30200	0.15590	0.04111	0.10000	0.41642
SPM	0.10000	0.51297	0.16216	0.28000	0.65000	0.47111	0.15000	0.63161
L3	0.05000	0.38549	0.11189	0.20200	0.34409	0.37777	0.07500	0.34645
LP	0.15000	0.38931	0.10594	0.23000	0.36692	0.40000	0.13125	0.37677

3.3.1. Recall@k

The recall results for each similarity-based method on all the datasets have been shown in Table 5. This measure represents the ability to find all positive/relevant samples by a classifier. We observe that the SPM outperforms against the existing methods on four datasets (Macaque, C. Elegans, Jazz, and USAir97). CAR method best performs on Karate and HDI on dolphins. The global version of LHN (i.e., LHNG) works best on football dataset, and RA is the best performing approach on Netscience. Local similarity methods extract relevant documents more precisely on 3 datasets, and the global methods retrieve more accurate on 5 datasets. Quasi-local approaches and CAR-based indices lie in top-5 ranked algorithms. The quasi-local methods have average good performance compared to the global approaches.

Table 9 shows the recall results for other representative methods where SPM outperforms on C. Elegans, Jazz, USAir97, and Netscience. HSM is a good indicator for Dolphin, Node2vec for Macaque, and CCLP2 for Karate and Football networks. On Karate, both the CCLP2 and The SPM show equally good performance.

3.3.2. Area under the precision–recall curve (AUPR)

Area under the precision–recall curve (AUPR) is proved to be more informative for imbalanced datasets. The real-world networks are highly imbalanced as the number of positive samples is very less than the negative samples. Table 6 shows the AUPR results on eight datasets. Here, We observe that the aupr results resemble the recall@k, i.e., SPM best performs on five datasets as that of recall results and CAR, HDI, LHNG, and RA outperform on karate, dolphin, football and netscience datasets respectively.

AUPR results corresponding to other representative methods are tabulated in Table 10. Here, The Node2vec shows the best performance on all datasets except the Macaque, where SPM performs well. We also observe that the Laplacian eigenmaps (Leig) and Isomap are the worst performers on all datasets.

3.3.3. Area under the receiver operating characteristics curve (AUROC)

The AUROC (or AUC) results have been reported in Table 7. Here, we observe that the global approaches perform best on Macaque, C. Elegans, Football, jazz, and dolphin, while the RWR is the best-ranking algorithm on C. Elegans and Dolphin, the SPM is the best ranker on the Macaque and Jazz networks and Cos⁺ is best on Football. The table shows the local approaches (i.e., RA and LNBCN) best result on usair97 and Netscience. The best performance of the RA index on usair97 is that this network is highly heterogeneous with a higher clustering coefficient and absence of a strongly assortative linking pattern. One more thing to note that the PA index works well on networks that follow rich-club phenomenon but, here we observe that the auroc results (please see Table 7) on netscience dataset (having a member of rich-club phenomenon), the PA and its CAR version (i.e., CPA) are having lowest values compared to all other methods. The reason is that this network is disconnected (consists of many connected components), and hence many nodes are isolated and

Table 6
AUPR results.

Methods	Karate	Macaque	C. Elegans	Football	Jazz	USAir97	Dolphin	Netscience
CN	0.07030	0.28027	0.04234	0.17455	0.51238	0.39009	0.11986	0.58330
JC	0.01397	0.03590	0.01626	0.23470	0.51956	0.04744	0.05666	0.48951
AA	0.05710	0.27442	0.05181	0.16056	0.54088	0.39757	0.08579	0.74908
RA	0.04075	0.26368	0.04786	0.16674	0.56630	0.43537	0.10066	0.76599
PA	0.02765	0.34323	0.02085	0.00506	0.06746	0.28537	0.02237	0.00276
SALTON	0.12384	0.26720	0.04554	0.15954	0.52720	0.36980	0.10740	0.58826
SORENSEN	0.06751	0.25750	0.04205	0.17329	0.51526	0.36370	0.10011	0.60575
CAR	0.21733	0.25305	0.04368	0.18561	0.53572	0.36416	0.10058	0.59813
CAA	0.16841	0.27093	0.04799	0.25907	0.55630	0.35679	0.06648	0.63863
CRA	0.10457	0.26834	0.05257	0.22227	0.60831	0.42311	0.08433	0.66808
CPA	0.09365	0.28103	0.04094	0.13220	0.50689	0.36225	0.05561	0.29581
HPI	0.03955	0.26042	0.04161	0.17890	0.52255	0.38245	0.09205	0.57262
HDI	0.10237	0.27818	0.04596	0.16474	0.50829	0.37266	0.17651	0.48951
NLC	0.07886	0.29955	0.03999	0.05210	0.41044	0.35726	0.04364	0.00123
LNBCN	0.07839	0.02712	0.03942	0.17575	0.41752	0.39835	0.05266	0.65339
LHNL	0.07935	0.24733	0.04664	0.17493	0.51164	0.39038	0.09565	0.59308
CCLP	0.05137	0.26904	0.04862	0.21183	0.55269	0.39442	0.07150	0.68109
KATZ	0.07354	0.33081	0.04047	0.16959	0.43602	0.38977	0.08983	0.51171
RWR	0.07874	0.38532	0.06197	0.21580	0.25769	0.09175	0.04901	0.20418
Shortest Path	0.01768	0.06612	0.01401	0.02278	0.02293	0.01120	0.02404	0.09914
LHNG	0.01393	0.02578	0.00781	0.30867	0.08835	0.00593	0.02896	0.04370
ACT	0.02262	0.31290	0.01823	0.01571	0.10025	0.29573	0.03442	0.17097
NACT	0.01379	0.02903	0.00725	0.25823	0.22698	0.00594	0.01488	0.19203
Cos ⁺	0.03392	0.19253	0.02463	0.23921	0.11242	0.02406	0.02476	0.02591
MF	0.06056	0.18374	0.02706	0.22821	0.14706	0.04203	0.04985	0.37889
SPM	0.08345	0.54004	0.08662	0.20540	0.68789	0.44342	0.08450	0.67826
L3	0.07943	0.38737	0.04389	0.13276	0.29754	0.35378	0.05865	0.32287
LP	0.06738	0.38600	0.04183	0.17927	0.31623	0.37741	0.07059	0.38762

Table 7
AUROC results.

Methods	Karate	Macaque	C. Elegans	Football	Jazz	USAir97	Dolphin	Netscience
CN	0.66139	0.78749	0.87663	0.86928	0.95854	0.96328	0.81040	0.99832
JC	0.60817	0.40191	0.81608	0.85834	0.96612	0.93311	0.77303	0.99945
AA	0.65683	0.78724	0.88189	0.84932	0.96488	0.97402	0.74429	0.99932
RA	0.72318	0.79259	0.88736	0.85685	0.97466	0.97728	0.78629	0.99953
PA	0.67080	0.91731	0.76145	0.25727	0.76620	0.91754	0.66843	0.74877
SALTON	0.71231	0.78275	0.86753	0.85063	0.96092	0.96565	0.76276	0.99933
SORENSEN	0.66563	0.77183	0.86141	0.85060	0.95879	0.96409	0.75858	0.99919
CAR	0.50134	0.78093	0.84277	0.84644	0.96118	0.95233	0.68292	0.95230
CAA	0.45998	0.79080	0.83838	0.84661	0.96122	0.95184	0.66276	0.94446
CRA	0.51796	0.78359	0.84652	0.83926	0.96944	0.96218	0.66596	0.94907
CPA	0.62594	0.81236	0.77510	0.67895	0.94639	0.91978	0.57540	0.76972
HPI	0.65191	0.78251	0.86798	0.88764	0.96274	0.96466	0.75740	0.99957
HDI	0.74994	0.78996	0.87213	0.86098	0.96041	0.96735	0.80747	0.99887
NLC	0.70911	0.84349	0.86321	0.80929	0.95393	0.90702	0.73127	0.59213
LNBCN	0.41790	0.21923	0.75164	0.78726	0.87271	0.96541	0.60227	0.99964
LHNL	0.70088	0.77854	0.86965	0.85957	0.95841	0.96452	0.75054	0.99937
CCLP	0.67934	0.79209	0.88046	0.86181	0.96510	0.96933	0.80289	0.99849
KATZ	0.73788	0.86200	0.86363	0.86480	0.94763	0.96276	0.80212	0.99934
RWR	0.84177	0.92849	0.90660	0.90110	0.95999	0.97113	0.88020	0.99347
Shortest Path	0.61283	0.61649	0.79575	0.75712	0.68467	0.83586	0.85132	0.95818
LHNG	0.64006	0.13196	0.72447	0.89802	0.90150	0.73105	0.77995	0.98541
ACT	0.51413	0.86292	0.74735	0.56118	0.80132	0.92371	0.81672	0.94354
NACT	0.67774	0.23142	0.65767	0.90085	0.91021	0.69823	0.78073	0.94845
Cos ⁺	0.80465	0.70811	0.86812	0.90329	0.91351	0.91748	0.82919	0.95781
MF	0.75465	0.71741	0.87700	0.89923	0.92916	0.95041	0.84830	0.95504
SPM	0.74565	0.95551	0.90499	0.84371	0.97807	0.95203	0.75200	0.99148
L3	0.84751	0.91431	0.84857	0.84796	0.91036	0.93500	0.77610	0.97264
LP	0.76661	0.91639	0.84942	0.87907	0.91572	0.94884	0.78649	0.99915

lower degree values. The Quasi-local method viz., the path of length 3 (i.e., $L3$), is the best performer on the karate network. On these datasets, only quasi-local and global approaches lie in top-5. The quasi-local methods lie among top-5 on almost all datasets considered here.

Table 8

Average precision results.

Methods	Karate	Macaque	C. Elegans	Football	Jazz	USAir97	Dolphin	Netscience
CN	0.01821	0.11240	0.01319	0.03005	0.06329	0.01728	0.02566	0.00113
JC	0.01164	0.03315	0.00988	0.03158	0.06432	0.01179	0.02240	0.00112
AA	0.01585	0.11140	0.01381	0.02882	0.06464	0.01764	0.02238	0.00117
RA	0.01775	0.11130	0.01380	0.02942	0.06649	0.01850	0.02538	0.00118
PA	0.01457	0.13441	0.01004	0.00373	0.03275	0.01556	0.01523	0.00032
SALTON	0.01938	0.11038	0.01321	0.02895	0.06384	0.01709	0.02419	0.00112
SORENSEN	0.01658	0.10809	0.01288	0.02963	0.06347	0.01705	0.02333	0.00112
CAR	0.01882	0.10872	0.01277	0.02994	0.06429	0.01688	0.02151	0.00104
CAA	0.01644	0.11117	0.01291	0.03153	0.06465	0.01692	0.01889	0.00104
CRA	0.01556	0.11031	0.01335	0.03064	0.06685	0.01782	0.01967	0.00104
CPA	0.01932	0.11484	0.01164	0.02437	0.06280	0.01653	0.01514	0.00064
HPI	0.01461	0.10972	0.01294	0.03078	0.06397	0.01714	0.02288	0.00112
HDI	0.02094	0.11184	0.01334	0.02948	0.06339	0.01722	0.02800	0.00112
NLC	0.01725	0.12147	0.01303	0.02145	0.06099	0.01615	0.02017	0.00020
LNBCN	0.01141	0.01766	0.01185	0.02665	0.05660	0.01760	0.01864	0.00114
LHNL	0.01830	0.10794	0.01317	0.02979	0.06336	0.01733	0.02323	0.00112
CCLP	0.01529	0.11138	0.01371	0.03056	0.06488	0.01754	0.02424	0.00115
KATZ	0.01915	0.12693	0.01284	0.02957	0.06022	0.01723	0.02383	0.00108
RWR	0.02122	0.13976	0.01436	0.03425	0.05601	0.01427	0.01647	0.00511
Shortest Path	0.01252	0.06152	0.00914	0.01551	0.02105	0.00731	0.01347	0.00142
LHNG	0.01212	0.02209	0.00687	0.03401	0.04081	0.00520	0.01889	0.00391
ACT	0.01155	0.12567	0.00950	0.01200	0.03723	0.01554	0.01414	0.00354
NACT	0.01244	0.02610	0.00650	0.03337	0.05158	0.00519	0.01091	0.00359
Cos ⁺	0.01698	0.09097	0.01154	0.03260	0.04411	0.00970	0.01290	0.00324
MF	0.01748	0.09138	0.01193	0.03243	0.04785	0.01174	0.01561	0.00166
SPM	0.01960	0.15628	0.01565	0.03046	0.06972	0.01798	0.02276	0.00115
L3	0.02324	0.13855	0.01256	0.02746	0.05304	0.01644	0.02174	0.00095
LP	0.01951	0.13865	0.01258	0.03025	0.05424	0.01689	0.02297	0.00103

Table 9

Recall results for other representative methods.

Methods	Karate	Macaque	C. Elegans	Football	Jazz	USAir97	Dolphin	Netscience
HSM	0.07500	0.34885	0.07405	0.24400	0.29606	0.22666	0.17000	0.17407
Leig	0.02500	0.07022	0.01405	0.06000	0.12007	0.01777	0.02000	0.08333
Isomap	0	0.01221	0.01081	0.04200	0.10433	0.01444	0.02000	0.19814
Node2vec	0.03898	0.67234	0.02167	0.09887	0.09738	0.02741	0.01947	0.08102
CCLP	0.05000	0.28015	0.09675	0.29000	0.52244	0.40555	0.10625	0.60000
CCLP2	0.10000	0.40305	0.09675	0.32600	0.41850	0.38555	0.12500	0.41419
NLC	0.05000	0.30763	0.07351	0.08800	0.44803	0.39444	0.05000	0.00516
SPM	0.10000	0.51297	0.16216	0.28000	0.65000	0.47111	0.15000	0.63161

The auoc results of other representative methods are shown in Table 11, where SPM is the best performer on Macaque, C. Elegans, and Jazz networks. HSM performs best on Football and Dolphin networks, CCLP is the best method on USAir97 and Netscience networks. On Karate, Isomap is the best performing similarity index.

3.3.4. Average precision

Table 8 shows the average precision results of similarity-based methods on eight datasets. The global approaches here, also are the best performer on all datasets except Karate and usair97, and dolphin where the quasi-local index (L3), and local indices RA and HDI respectively are the best. Here, SPM performs overall best on Macaque, C. Elegans, and jazz networks. The Resource allocation and HDI methods are top rankers on usair97 and dolphin networks, respectively.

Table 12 represents the average precision results of the other representative methods. From the table, it is observed that the Node2vec shows the highest average precision values against all networks except Macaque and Dolphin, where SPM and CCLP respectively show the best results.

Parameters settings. We conduct 10-fold cross-validation to evaluate each method on four different evaluation metrics described in the earlier subsection. The disadvantage with the global approaches is parameters tuning that needs to be done carefully to obtain good results. The dumping parameter β of the Katz index is set to 0.01, the return probability $(1 - c) = 0.3$ in the random walk with restart method. The ϕ of the global version of Leicht–Holme–Newman, i.e., LHNG, is set to 0.5 that equally balances both self and neighborhood similarity terms. The free parameter $\varepsilon = 0.5$ and path up to the length 3 is considered in the local path index.

Table 10
AUPR results for other representative methods.

Methods	Karate	Macaque	C. Elegans	Football	Jazz	USAir97	Dolphin	Netscience
HSM	0.06145	0.33552	0.03442	0.18720	0.23787	0.15064	0.12695	0.14478
Leig	0.03166	0.05281	0.00944	0.04602	0.07086	0.01360	0.01759	0.04326
Isomap	0.03072	0.03123	0.01084	0.03013	0.06972	0.01223	0.01948	0.11020
Node2vec	0.90000	0.08058	0.72930	0.79032	0.91563	0.84788	0.65000	0.85818
CCLP	0.05137	0.26904	0.04862	0.21183	0.55269	0.39442	0.07150	0.68109
CCLP2	0.11750	0.40541	0.04602	0.27284	0.44192	0.36988	0.08118	0.41091
NLC	0.07886	0.29955	0.03999	0.05210	0.41044	0.35726	0.04364	0.00123
SPM	0.08345	0.54004	0.08662	0.20540	0.68789	0.44342	0.08450	0.67826

Table 11
AUROC results for other representative methods.

Methods	Karate	Macaque	C. Elegans	Football	Jazz	USAir97	Dolphin	Netscience
HSM	0.78390	0.91606	0.84022	0.88921	0.87704	0.92754	0.85570	0.97286
Leig	0.72189	0.50369	0.70811	0.81513	0.81454	0.81281	0.77543	0.91977
Isomap	0.80300	0.32430	0.74469	0.79262	0.85134	0.81146	0.81074	0.96592
Node2vec	0.76850	0.63184	0.80230	0.85278	0.87941	0.85538	0.71474	0.89241
CCLP	0.67934	0.79209	0.88046	0.86181	0.96510	0.96933	0.80289	0.99849
CCLP2	0.74751	0.91921	0.84180	0.87319	0.93687	0.94741	0.77198	0.97351
NLC	0.70911	0.84349	0.86321	0.80929	0.95393	0.90702	0.73127	0.59213
SPM	0.74565	0.95551	0.90499	0.84371	0.97807	0.95203	0.75200	0.99148

Table 12
Average precision results for other representative methods.

Methods	Karate	Macaque	C. Elegans	Football	Jazz	USAir97	Dolphin	Netscience
HSM	0.01995	0.13304	0.01201	0.03063	0.04953	0.01455	0.01862	0.00427
Leig	0.01585	0.05008	0.00751	0.02061	0.03569	0.00776	0.01124	0.00325
Isomap	0.01764	0.02514	0.00816	0.01789	0.03630	0.00747	0.01200	0.00382
Node2vec	0.04780	0.08355	0.02206	0.10330	0.09875	0.02816	0.02305	0.08175
CCLP	0.01529	0.11138	0.01371	0.03056	0.06488	0.01754	0.02424	0.00115
CCLP2	0.02145	0.14074	0.01259	0.03241	0.05924	0.01669	0.02390	0.00100
NLC	0.01725	0.12147	0.01303	0.02145	0.06099	0.01615	0.02017	0.00020
SPM	0.01960	0.15628	0.01565	0.03046	0.06972	0.01798	0.02276	0.00115

3.4. Efficiency

We have performed our experiment on a 64-bit core i7 Intel system having 8 GB internal memory and 3.60 GHz speed without a dedicated graphics card. To reduce the computational time, some optimization strategies can be applied (if possible) for example, the union and intersection of two sets of sizes m and n can be computed in $O(m + n)$ using hash tables. The computational complexity of the addition and the subtraction of two matrices are $O(n^2)$, however, these operations can be performed in $O(nt)$ in sparse networks where, $t \ll n$. The matrix multiplication of two dense matrices of sizes $m \times n$ and $m \times p$ are done in $O(mnp)$, while it is $O(mtp)$, where $t \ll n$. The matrix inversion typically takes $O(n^3)$ time for a square matrix of size $n \times n$ however, some improvements are available that reduce the time to $O(n^{2.81})$ or even less. The time complexity of the similarity-based methods have been tabulated in Table 13, in which most complexities are explained in [26]. In the table, computational complexity of each method are shown using big O notation where n , e , and K are the number of nodes, number of links, and average degree of the networks.

4. Variations of link prediction problem

As earlier mentioned that the techniques listed in this work mainly focus on a simple abstract graph (i.e., a graph with no vertex or edge attribute). The networks considered in this work simple undirected and unweighted. However, some modification needs to be done to apply on weighted and directed networks. In such networks, links are assigned with weights that represent the strengths of these links. Two types of link direction can be possible of a node (i.e., incoming and outgoing). So, a node x can have two types of neighbors (degrees) viz., in-neighbors $\Gamma_i(x)$ and out-neighbors $\Gamma_o(x)$. Based on these modifications, earlier similarity approaches can be redefined as given below.

In a directed network, the common neighbor method based on in-neighbors and out-neighbors are expressed as

$$S_i(x, y) = |\Gamma_i(x) \cap \Gamma_i(y)|, \tag{79}$$

and

$$S_o(x, y) = |\Gamma_o(x) \cap \Gamma_o(y)|. \tag{80}$$

Table 13

The computational Complexity of similarity-based methods and the corresponding references.

Method	Time complexity	Reference
Local similarity index		
CN	$O(nK^3)$	[2]
JC	$O(nK^3)$	[33]
AA	$O(nK^2)$	[29]
RA	$O(nK^3)$	[30]
PA	$O(nK^2)$	[28]
Salton	$O(nK^3)$	[35]
Sorenson	$O(nK^2)$	[36]
CAR	$O(nK^4)$	[38]
CAA	$O(nK^4)$	[38]
CRA	$O(nK^4)$	[38]
CPA	$O(nK^3)$	[38]
HPI	$O(nK^3)$	[39]
HDI	$O(nK^3)$	[39]
LNBCN	$O(n \cdot O(f(z) + nK^3))$	[40]
LHNL	$O(nK^3)$	[41]
CCLP	$O(n^2K^2)$	[42]
NLC	$O(nK^3)$	[43]
Global Similarity Index		
Katz	$O(nK^3)$	[46]
RWR	$O(cn^2K)$	[47]
Shortest Path	$O(n \log n)$	[2]
LHNG	$O(cn^2K)$	[41]
ACT	$O(n^3)$	[53]
NACT	$O(n^3)$	[53]
L+	$O(n^3)$	[51]
MF	$O(n^3)$	[58]
SPM	$O(n^3)$	[159]
Quasi-local similarity Index		
LP	$O(ln^2K)$	[63]
L3	$O(n^3)$	[64]

In weighted directed network, the expression are

$$S_i^{weight}(x, y) = \sum_{z \in (I_i^+(x) \cap I_i^+(y))} \frac{w(z, x) + w(z, y)}{2}, \quad (81)$$

and

$$S_o^{weight}(x, y) = \sum_{z \in (I_o^+(x) \cap I_o^+(y))} \frac{w(x, z) + w(y, z)}{2}. \quad (82)$$

In a similar way, other approaches can be modified for directed and weighted networks. The point to be noted here is that first, define several topological features (e.g., degree, path, clustering coefficient, etc.) in weighted directed networks and apply these features to implement several link prediction algorithms.

Mostly works on link prediction focus mainly on simple undirected networks due to simplicity. The cost for this simplicity is that it fails to extract rich information available in most real-world networks, which are not undirected in general. Some notable works on directed weighted networks are nicely presented in [171–181]. Lichtenwalter et al. [171] work, published in SIGKDD, extracts 12 topological features on a large directed weighed network of over 5 million nodes and performs ensembles of classification algorithms (C4.5, J48, Naive Bayes). The training over such a big network (millions of edges) is problematic; to mitigate this issue, the authors defined edge features of vertices of 2 and 4 hops only. They perform a quasi-local training to obtain the final model, and their results are outperforming compared to the state-of-the-art. Further, Bütün et al. [176] proposed a new topological similarity metric published in ASONAM that takes into account temporal and weighted information in directed networks which are useful for the improvement of the accuracy. They extract all possible triad pattern features and incorporate them with the weighted version of baseline topological similarity metrics (CN, JC, AA, RA, and PA). They employed a supervised learning framework using these metrics as features and predict missing links. Recently, Bütün et al. [178] introduced a supervised learning model for predicting the citation count of scientists (PCCS). They formulate the problem of PCCS as a link prediction problem and predict links with their weights in weighted temporal and directed (citation) networks. Their model incorporated both local and global topological features and claims the excellence of their proposed work.

Table 14
Link prediction in temporal networks.

Models	Network types	Characteristics	References
Learning-based models	Coauthorship networks	High computational cost	Vu et al. [184], Pujari et al. [185], Zeng et al. [186], He et al. [187], Bao et al. [188], Madadhain et al. [189], Bringmann et al. [190]
Heuristics-based models	Twitter, Collaboration and Coauthorship networks	Fast convergence and high precision	Catherine et al. [191], Sherkat et al. [192]
Probabilistic model	Nodes-attributed graphs	Characterize the stochastic and dynamic relations. Need prior link distribution so impractical for real networks	Hu et al. [193], Barbieri et al. [194], Gao et al. [183], Ji Liu et al. [195], Hanneke et al. [196]

The link prediction problem is being explored in several other types of networks, including temporal networks, signed social networks, heterogeneous networks, bipartite networks, etc. Some of these variations are studied in this section.

4.1. Link prediction in temporal networks

Today's scenario shows that the relationships among users in social networks are continuously changing; for example, each time in the Facebook network, some users join, and some others quit. It results in the networks to be highly complex. Here, time is an important parameter to consider for the evolution of networks. In temporal link prediction, time is considered as the third dimension and represented by a third-order tensor A .

$$A(i, j, T) = \begin{cases} 1 & \text{if node } i \text{ is connected with node } j \text{ at time } T, \\ 0 & \text{Otherwise} \end{cases} \quad (83)$$

Thus, for a given sequence of snapshots of a network at different time interval T_1, T_2, \dots, T_t , the link prediction finds links that evolves at the next time slot T_{t+1} .

Several efforts have been employed by the researchers in this direction in the last decade. Purnamrita et al. [182] introduced a nonparametric method for temporal network link prediction where the time dimension is partitioned into subsequences of snapshots of the graph. This approach predicts links based on topological features and local neighbors. Dunlavy et al. [25] employ matrix and tensor techniques in a framework where matrix part collapses sequence of snapshots of networks into a single matrix and computes link scores using truncated svd and extended Katz methods. The tensor part computes the scores using heuristics and temporal forecasting. The tensor part captures the temporal patterns effectively in the network, but it costs heavily also. Moreover, Gao et al. [183] proposed a model based on latent matrix factorization that employs content values with the structural information to capture the temporal patterns of links in the networks. Table 14 shows some more works in this direction.

4.2. Link prediction in bipartite networks

Till now, we have reviewed link prediction methods in unipartite networks in which links may present between any pair of vertices. Now, we review the link prediction problem in a specific graph where only two sets of vertices are present, and a link can be possible between a pair of vertices in which one vertex belongs to one set of vertices and the other vertex to another set of vertices. Such types of networks are called bipartite networks. Lots of social networks logically can be considered as bipartite such as Term-Document network [197], Scientists-Papers cooperation network [198], RNA-PI network [199], IMDb network, and many more.

Kunegis et al. [200] study the link prediction problem in bipartite networks and observed that most common neighbor-based approaches (e.g., Common Neighbors, Adamic/Adar, Resource Allocation, etc.) are not applicable to these networks. The reason is that adjacent nodes belong to different clusters and are connected with the path of odd lengths only. Also, common neighbor-based approaches are based on the path of length two. The authors give hyperbolic Sine and Von Neumann kernels of odd order to compute the similarity between vertices. Only the PA method is applicable to these networks in its natural form because it considers the degree of the neighbors. Some researchers [201–203] have implemented common neighbor-based methods (e.g., CN, AA, RA, PA, LCP-CN, etc.) in bipartite networks. Xia et al. [201] studied the link prediction problem by exploiting structural holes in bipartite networks. They proposed two implementations of structural holes viz., absent links (consisting of c-type and s-type links), and minimum description length [204,205].

Recently several methodologies of link prediction in bipartite networks have been addressed. Baltakiene et al. [206] implemented maximum entropy principle, an extension of the recent one [151]. They used probability of Bipartite Configuration Model [207] as the score function. Allali et al. [208] presented the term "internal link" based on which they proposed a new link prediction algorithm.

Table 15

Link prediction in heterogeneous networks.

Models	Network types	Characteristics	Reference
Supervised models	Youtube, Gene, Climate	Proposed both unsupervised and supervised approach to link prediction	Davis et al. [213]
	DBLP	Extracts meta path-based topological features and applies logistic regression as prediction model	Y. Sun et al. [215]
	Epinions, Slashdot, Wikivote, Twitter	Define social pattern-based features (social balance and microscopic mechanism), input to the inference model namely (transfer) factor graph models	Y. Dong et al. [216]
Collective LP models	MovieLens, Book-Crossing, Douban	Non-parametric Bayesian model that considers the similarity between tasks when leveraging all the link data together	B. Cao et al. [217]
	Flickr, DBLP	Distance feature extraction using both network and node features and for learning Multi-Task Structure Preserving Metric Learning (MTSPML) is used	S. Negi et al. [218]

4.3. Link prediction in heterogeneous networks

Most of the contemporary approaches of link prediction focus on homogeneous networks where the object and the link are of single (same) types such as author collaboration networks. These networks comprise less information, like which two authors have collaborated with a paper that causes less accuracy for the prediction task. In heterogeneous networks, the underlying assumption of a single type of object and links does not hold good. Such networks contain different types of objects as well as links that carry more information compared to homogeneous networks and hence more fruitful to link prediction, also called multi-relational link prediction (MRLP). Examples of such networks are DBLP bibliography¹⁰ and Flickr networks.¹¹ In the bibliography database, authors, papers, venue, terms are different types of objects/nodes, and relationships are paper–author, author–author, paper–term, paper–venue, and so on.

Sun et al. [209,210] coined the concept of heterogeneous information network (HIN) and subsequently meta path concept [211], since then it becomes popular among researchers. The key idea to multi-relational link prediction (MRLP) is to employ an appropriate weighting scheme to combine different link types. The authors predict the relationship building time between two objects by encoding the target relation and topological features to meta paths in a supervised framework. Moreover, Yang et al. [212] proposed a new topological feature, namely multi-relational influence propagation to capture the correlation between different types of links and further incorporate temporal features to improve link prediction accuracy. Davis et al. [213] proposed a novel probabilistic framework, a weighted extension of Adamic/Adar measure. Their approach is based on the idea that the non-existing node pair forms a partial triad with their common neighbor, and their probabilistic weight is based on such triad census. Then the prediction score is computed for each link type by adding such weights. Meanwhile, Sun et al. [214] a new supervised framework for HIN where meta path-based topological features (i.e., path count, random walk) are extracted, and then logistic regression is applied to build the relationship prediction model that learns the weight associated with these features. Table 15 lists some more works on link prediction in heterogeneous networks.

5. Link prediction applications

5.1. Network reconstruction

Guimerà et al. [73] proposed a framework that applies link prediction for network reconstruction. They reconstruct of the true network is done from the observed network based on missing links (removed one) and the spurious links (added links). Although it is not obvious because no one knows about the amount of missing and spurious links in the networks. For this, the authors describe the reliability of networks based on the reliability of both missing and spurious links by formulating the link prediction problem as a stochastic block model [94]. The reliability of the network A is [73]

$$R(A) = \prod_{A_{xy}=1, x < y} R_{xy} = \prod_{A_{xy}=1, x < y} L(A_{xy} = 1/A^0), \quad (84)$$

where R_{xy} is the reliability of the link (x, y) that is defined by the likelihood that the link (x, y) truly exists given the observed network A^0 . This equation can be solved by finding out the network A that maximizes the reliability given by (84).

¹⁰ <https://dblp.uni-trier.de>.

¹¹ <http://www.flickr.com/>.

The computational cost of the equation is high, so the authors [73] give a greedy algorithm to compute it. The algorithm starts with computing the link reliability of all pairs of vertices. At each step, the algorithm removes the least reliable link and adds the most reliable link (non-existing in the current network). This change in the network is accepted when the reliability of the network increases and rejected otherwise. In case of rejection, the next step selects the least reliable existing link and the highest reliable non-existing links for swapping. The algorithm stops when there are no five consecutive changes (swaps) in the network. The reliability of the network improves from the initial observed network, which is the reconstructed ones. Now, the authors compare the six (6) global properties of both the observed and the reconstructed networks and show that the reconstruction improves the estimates.

5.2. Recommender system

The recommender systems [5,6,134,219] (also called information filtering systems) have been widely applied in social media (like Facebook, Twitter) and online shopping websites (e.g., Flipkart, Amazon, etc.). Such systems recommend new friends, followers, and followers on social networking platforms and new products on online shopping portals based on users' previous browsing history (such as interests, preferences, ratings, etc.). Even though collaborative filtering (CF) is a successful recommendation paradigm that applies transaction (Transaction/purchase is essentially an implicit and coarse rating on preferring an item [220]), information to enrich user and item features for recommendation. Although they have been applied in many recommender systems, they are greatly limited by data sparsity problem [221]. The recommender system in bipartite networks can be mapped to link prediction problem as follows [222]. Consider U^* and O be the sets of users (first set of vertices) and objects/items (second set of vertices). Construct the user-item interaction graph $G = (V, E)$ from the available transactions T (purchasing patterns), where $V = U^* \cup O$ and $E = \{(u, o) : u \in U^*, o \in O, u \rightarrow o \in T\}$.

Huang et al. [5] and Li et al. [223] proposed approaches, where the recommender system (user-item recommendation) is represented as a bipartite graph, and employed basic link prediction approaches for the items recommendation. Sadilek et al. [224] proposed FLAP (Friendship + Location Analysis and Prediction) system in which both friendship and location prediction tasks are implemented. They employed users tweets, their locations, and their neighboring information as model features and inferred social ties and location using MRF. More related works can be found in [222,225,226].

5.3. Network completion problem

In general, the network representation of the real-world problem is incomplete or partially observed or incremental with both missing links and nodes such as wall posts on Facebook, tweets in the Tweeter, etc. The problem arises due to several reasons like security, data aggregation overhead, manual errors, etc. Predicting such nodes and links is the network completion problem in which some notable works like [227] in SIGKDD, [228] in ASONAM, and [139] in EPL. Filling missing entries of the adjacency matrix of a network is link prediction, which can be considered as a subset of network completion problem. Kim et al. [229] cast network completion problem to the Expectation-Maximization (EM) framework and proposed KronEM, an EM model based on Kronecker graphs. They, first, represent the network as Kronecker graph and estimate the model parameters as well as missing links using KronEM algorithm. The estimated network is then considered as the complete network and re-estimate the model, and this process is repeated until the convergence. Further, Pech et al. [139] employed the robust principle component analysis (Robust PCA) [230] [to recover both low rank and sparse components of a data adjacency matrix] in link prediction framework and introduce a novel global prediction method using both the components. They reconstruct the original network using the robust PCA where these components are extracted by minimizing the weighted combination of the nuclear norm and of the l_1 norm [230].

$$\min_{X^*, \mathcal{E}} \|X^*\|_* + \lambda \|\mathcal{E}\|_1, \quad (85)$$

where $\|\cdot\|_*$ is nuclear norm (i.e., sum of singular values) of the matrix and $\|\cdot\|_1$ is the l_1 norm, \mathcal{E} is error or noise matrix (sparse matrix containing spurious links as positive entries and missing links as negative entries) and λ is a positive weighing parameter that balances the contribution of both the components (low rank and sparse components). $X^* [= A^O - \mathcal{E}$, here, A^O is the observed network] is the set of patterns (links that are newly predicted and some links that are eliminated). From which only the newly appeared links are extracted and added to the observed network A^O to recover the original matrix (also known as reconstructed matrix). Once the reconstructed matrix is obtained, link prediction can be performed accordingly.

5.4. Spam mail detection

Spreading and receiving irrelevant emails is common in today's world that consumes network bandwidth, memory, etc. Many email service companies are trying to implement several filter mechanism to stop such emails known as spam mails. To implement spam filter mechanism, spam detection is a necessary task. In this context, Zan Huang and Daniel D. Zeng [231] proposed a model to detect spam emails using link prediction. They construct an email graph (directed and weighted) based on the email data, consisting of a sender, recipient, and timestamp of the communication as attributes. Many email communication links between the sender and the receiver are mapped to a weight of the link between them. Then, an anomaly score is computed for each distinct sender-recipient pair using the Adamic/Adar link prediction approach by making it adaptive based on the spreading activation algorithm. Some more work related to this can be found in [232,233].

5.5. Privacy control in social networks

Lots of users share personal posts, audios, videos, and other sensitive information to social networking websites. Trust is an important parameter to evaluate users' relationships on such media, i.e., the strength of a relationship between two users can be determined based on the trust in the form of link weights. Thus, it is important for companies to maintain the privacy of users from anomalous ones. Oufi et al. [234] proposed a framework implementing a capacity-based algorithm that employs Advogato trust metric [235,236] to compute the level of trust between users. This means that the framework identifies all possible trustworthy users of a seed user, which results in the privacy of that user in the network from anomalous users.

5.6. Identifying missing references in a publication

A research article may contain some irrelevant references and miss some relevant ones. Identifying such missing references in a research article is an important task to avoid plagiarism. It becomes more critical for the point of a novice researcher due to a lack of literature survey carried out by him. Kc et al. [237] proposed a machine learning approach to link prediction tackle this problem. They provide a framework for the generation of links between referenced and otherwise interlinked documents. The nodes of the graph represent documents, and the links between them show references available between them. They find new links/references of documents based on this graph using Probability Measure Graph Self-Organizing Map (PM-GraphSOM).

5.7. Routing in networks

In complex network theory, link prediction in social networks resembles link quality prediction in wireless sensor network [238]. The routing problem in a network finds the shortest path (optimal) between the sender and the receiver. The strength of signal frequently varies in mobile and ad hoc networks that results in frequent breaks in routes and degrades the performance result. Weiss et al. [239] and Yadav et al. [240] proposed some models to estimate the signal strength-based link availability prediction for optimal routing. Such link information is beneficial to estimate the link breakage time and hence, to repair the existing route or to discover a new route for the packets. This reduces end-to-end routing delay and packet drops, thereby improving the performance. Once broken links are identified earlier (using link breakage time), routing management protocol needs either to repair the broken link or to find an alternate route. Several works state that link prediction may play a crucial role in this scenario that results in low latency in packet delivery to the receiver and hence improves reliability. Hu and Hou [241] presented link prediction-based traffic prediction for the best routing of packets in a wireless network. Some more works in this area can be found in [242,243]. Recently, Zhao et al. [238] proposed a neighborhood-based NMF model to estimate the link quality in the wireless sensor network. They extend the link prediction model to the wireless sensor network, where they predict the quality of a link based on NMF associated with structural (neighborhood) information.

5.8. Incorporating user's influence in link prediction

Lots of works based on individual influence have been proposed in social network analysis, such as link prediction [244,245], information diffusion [246–248], influence maximization [249–254], community detection [255,256], etc. Particularly, the role of individual influence in link prediction provides a new perspective/insight into the problem. Influence maximization (IM) [249] is one of the fundamental problems in social network analysis where the goal is to find a set of users (seed set) that can be further utilized to maximize the expected influence spread (defined as the expected number of influenced users) among others. The influence (social influence here) is propagated through certain channels (i.e., intermediate nodes), that are captured by diffusion models [246]. IM and diffusion models are a cooperative and correlated task as for IM, several Diffusion models are used in the computing framework. Zhang et al. [257] proposed a new framework of link diffusion to predict more links in the microblogging networks. They find the triadic structure to be the crucial factor that affects the link diffusion process and hence, link prediction. Earlier, Cervantes et al. [258] proposed a supervised learning model to find an influential collaborative researcher in the collaboration network. They employ the model to the whole network and compare its result with those sub-networks generated each time when a distinct vertex is removed from the training set. Finally, results are ranked and examine the collaborative influential of each researcher based on the presence or absence of it in the network. Finding influential users (i.e., the seed set) is useful in many applications like viral marketing, where an influential user can be used to advertise the product to maximize the profit. Other application areas may be disease prevention using vaccinate to the most influential patient.

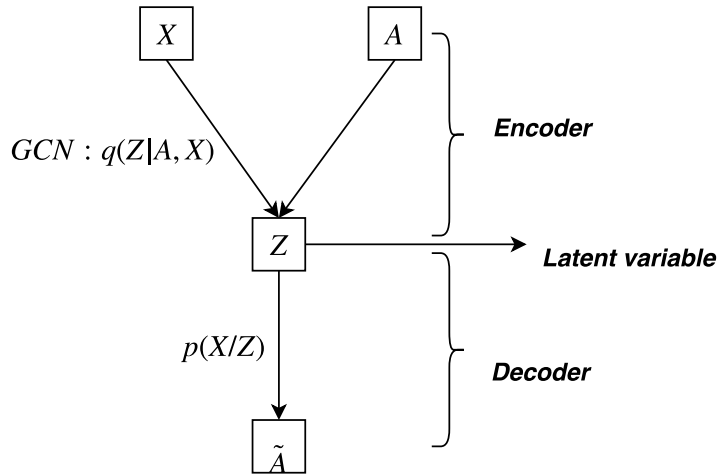


Fig. 10. A schematic diagram of variational graph autoencoder (VGAE) [123] where, $GCN : q(Z|A, X)$ is the probability of estimating the latent variable Z given the inputs feature matrix X and adjacency matrix A (here, GCN shows that the encoder part of the model is framed as graph convolutional network having two hidden layers.). $p(X/Z)$ is the probability that the input is reconstructed with the given latent representation Z .

6. Recent developments

6.1. Link prediction using deep learning

The deep learning methodology is a part of machine learning, which is based on data representation learning and has been used in various applications like image processing, computer vision, natural language processing, pattern recognition, etc. A very few works [259–261] i.e., deep learning-based link prediction in the literature are available. Li et al. [259] proposed a deep generative model viz., Conditional Temporal Restricted Boltzmann Machine (ctRBM) that captures complex transitional variance and local influences to find link structure in dynamic networks. This work and some other works (Wang et al. [120] and Wang et al. [262]) consider only structural features to performs link prediction in static/dynamic networks that limit their performance. To improve the performance, Wang et al. [260] devise a hierarchical Bayesian model that incorporates both structural, and node features to perform relational deep learning. Meanwhile, Schlichtkrul et al. [263] proposed a new model called the relational graph convolutional neural network (R-GCN) that employs GCN as a building block to model relational data (knowledge graph in particular). The authors represent the knowledge base (i.e., relational data) to a directed multigraph where both nodes (entities) and edges (relations) are labeled. The model encodes each relationship (both for incoming and outgoing edges) separately with including individual vertex feature also. In other words, the latent representation of a vertex depends on all neighbors (either incoming or outgoing) and the vertex itself. Once the encoding part is complete, the given framework can be applied in node classification and link prediction. They consider their model as an autoencoder for link prediction tasks where R-GCN (encoder part) is used for the latent representation of each entity and a tensor factorization model, namely, DistMult [264] as reconstruction.

Kipf et al. [123] introduce a variational graph autoencoder (VGAE) framework that learns latent representation on graph structured data (Fig. 10). The model takes two parameters as inputs; the adjacency matrix $A_{|V| \times |V|}$ of the network and a feature matrix $X_{|V| \times D}$, where $|V|$ is the total number of vertices, and D is the number of input features. The input adjacency matrix is preprocessed by adding a self-loop of each vertex to include its feature, and the matrix is normalized by the diagonal node degree matrix to resolve the feature scaling issue. This normalization becomes more useful when it is symmetric as

$$A_{norm} = \hat{D}^{-1/2} \hat{A} \hat{D}^{-1/2}, \tag{86}$$

where $\tilde{A} = A + I$ (I : the identity matrix of A) to enforce each vertex to include own feature and \tilde{D} is diagonal node degree matrix. The encoder part of the VGAE acts as graph convolutional network (GCN) that use two hidden layer and employ the following update rule to optimize the parameters of the model

$$f(H^{(l)}, A) = \sigma(\hat{D}^{-1/2} \hat{A} \hat{D}^{-1/2} H^{(l)} W^{(l)}), \tag{87}$$

where $W^{(l)}$ is the weight matrix for the l th layer of neural network $H^{(l+1)} = f(H^{(l)}, A)$, and $\sigma(\cdot)$ is the non linear activation function (e.g. RELU is normally used in GCN as non linear activation function). Now, the decoder part (generative model) of the model employs simple inner product of the latent variables for reconstruction. The authors perform experiments with and without input features (X) in their two models (viz., VGAE, and GAE) and compare with well-known embedding methods spectral clustering (SC) [265] and DeepWalk (DW) [102]. The experiment results with/without the input feature

Table 16

Deep learning frameworks for link prediction.

Model name	Architecture	Features used	Reference
ctRBM	Restricted Boltzmann Machine	Nodes transitional patterns and local neighbors influence	Li et al. [259]
Multi-relational model	General Neural Machine	Relational embeddings	Yang et al. [264]
VGAE	Graph Autoencoder	Node and structural features	Kipf et al. [123]
Collaborative filtering model	Graph Autoencoder	Node and structural features	Berg et al. [266]
GraphGAN	Generative Adversarial Network	Structural information	Wang et al. [267]
Graphite	Graph Autoencoder	Structural neighborhood features	Grover et al. [268]
WLNLM	Weisfeiler–Lehman Neural Machine	Enclosing subgraph	Zhang et al. [269]
SEAL	Graph Convolutional Network	Structural roles	Zhang et al. [124]
Multi-relational model	Graph Autoencoder	Structural neighborhood features	Schlichtkrull et al. [263]

matrix show significantly higher/comparable to SC and DW. Recently, some notable works [124,266–269] based on deep learning in link prediction have been proposed (Refer to Table 16). Zhang et al. [269] introduced a novel framework viz., Weisfeiler–Lehman Neural Machine (WLNLM), based on Weisfeiler–Lehman algorithm that labels the nodes of a graph and determines the vertex ordering using the topology of the underlying graph (especially based on structural roles). For each non-existing link, WLNLM extracts subgraphs in the neighborhood and encodes it as an adjacency matrix. Finally, a neural network is trained on these matrices to build a predictive model. Further, Zhang et al. [124] proposed a new heuristic learning paradigm (SEAL framework: learning from Subgraphs, Embeddings and Attributes for Link prediction) that captures first, second, and higher-order structural informations in the form of local subgraphs similar to the previous work [269]. The SEAL framework unifies all three types of information (i.e., local subgraph, embedding, and attribute information) using the graph convolutional network. The experimental results show that the learning based on these three information outperforms several heuristics (individually based on heuristic methods and latent feature methods).

6.2. Fuzzy model-based link prediction

L. A. Zadeh [270] introduces the concept of fuzzy in his paper published in Information and Control, which became further a well-known modeling paradigm for various types of applications. Bastani et al. [271] used a fuzzy paradigm in link prediction problem. They proposed two approaches viz. Fuzzy Clustering Coefficient (FCC)-based and fuzzy Cluster Overlapping (FCO)-based link prediction. In addition, they presented a hybrid model and an inference engine developed by creating a synergy between these models. Bastani et al. Their work is considered as an initial attempt to employ the concept of fuzzy logic in link prediction tasks in a complex network. In [271], authors developed a fuzzy model based on Yager's Paradigm for Intelligent Social Network Analysis (PISNA) [272]. In PISNA, Yager observed the possibility of a fuzzy set method that bridges the gap between human thoughts and a formal model of the network. Fuzzy logic acts as an important concept for mapping human thinking and reasoning to mathematical constructs.

In a social network, Clustering Coefficient (CC) [44] is an essential idea that represents the inter-connectivity among a given node and its neighbors. In other words, CC delineates the cliquishness of a given subgraph in a network. Yager proposed a softer definition for the Clustering Coefficient or cliquishness of a node in a graph. If a clique in a graph is represented by S , then the term S can be defined by the following criteria: $C1$: "Most vertices belonging to S are closely connected". $C2$: "No vertex in S is too far from remaining vertices". $C3$: "any vertex outside the clique S must not be better connected to that inside the clique". The above three criteria consist of some linguistic terms that should be defined in terms of fuzzy. Examples of such terms are "Closely connected", "Most", "Far", "Not Far" etc. Yager focused on providing machine understanding of each criterion. The linguistic term "Closely connected" (refer to Fig. 11) can be defined as the ramp function $Q : \text{positive integer} \rightarrow [0, 1]$ such that $Q(k)$ is the degree of closeness with shortest path of at most k edges between two vertices. Yager [272] proposed a prototypical definition of "Close" as a ramp function in the Fig. 11 and expressed using the following equation.

$$Q(k) = \begin{cases} 1 & k \leq a \\ \frac{b-k}{b-a} & a \leq k < b \\ 0 & k > b. \end{cases} \quad (88)$$

In Fig. 11, shortest path of $\leq a$ represent close, not close for $> b$ and partial closeness for $a \leq k < b$. This is also expressed using Eq. (88). The small-world phenomenon suggests that any two nodes can have maximum separation of 6 hops between them. Thus, the closeness function [271] must decrease exponentially to satisfy the small world phenomenon [273].

$$\text{Close}(x_i, x_j) = \frac{q_{ij}}{2 \times 10^{q_{ij}^2}}, \quad (89)$$

where q_{ij} is the path length between x_i and x_j and if it is a shortest path then

$$\text{Close}(x_i, x_j) = Q(q_{ij}). \quad (90)$$

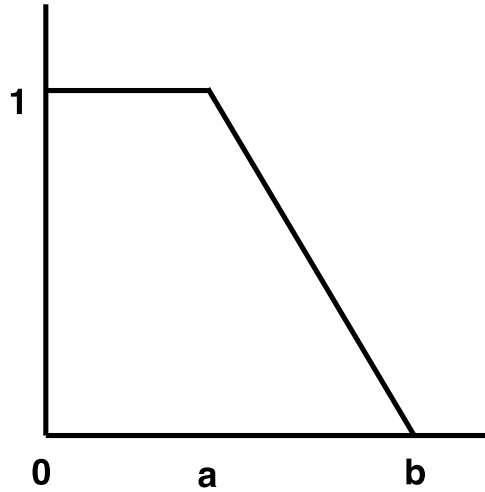


Fig. 11. Prototypical definition of Close.

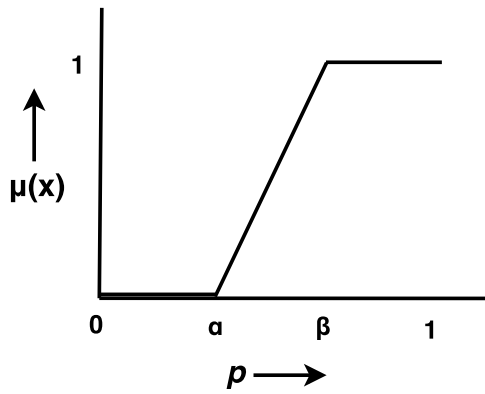


Fig. 12. The Most as a fuzzy set.

Here, $Q(q_{ij})$ shows the closeness between two vertices of path length q_{ij} . The above function Close is defined for an unweighted and undirected network, which can be extended for a weighted network described in [271].

$$Close(i, j) = \begin{cases} 1 & q_{ij} < 2 \\ \frac{w(i,k)+w(k,j)}{2*10q_{ij}^{-2}} & q_{ij} = 2 \\ \frac{w(i,k)+w(k,e)+w(e,j)}{2*10q_{ij}^{-2}} & q_{ij} = 3 \\ 0 & q_{ij} > 3. \end{cases} \tag{91}$$

Next, the linguistic term “Most” is also described in terms of fuzzy function $M(p)$ that expresses closeness of the considered node to all other nodes in the cluster [272]. Function $M(p)$ calculates the number of connected nodes that are as close as possible and that satisfy the following criteria [272]:

$$M(p) = \begin{cases} 0 & p = 0 \\ M(p_1) \geq M(p_2) & p_2 \geq p_1 \\ 1 & p = 1. \end{cases} \tag{92}$$

Above Eq. (92) can also be expressed as a ramp function [271] which is shown in Eq. (93) (represented by Fig. 12) and interpreted similar to Eq. (88).

$$M(p) = \begin{cases} 0 & p \leq \alpha \\ \frac{\beta-p}{\beta-\alpha} & (\alpha \leq p \leq \beta) \\ 1 & p \geq \beta. \end{cases} \tag{93}$$

where, the value p for a node in the network can be calculated as

$$p_{x_i} = \sum_{j=1, i \neq j}^{n_s} \frac{\text{Close}(x_i, x_j)}{n_s - 1} \quad (94)$$

Clearly, criteria C_1 can be computed using amount of $M(p_{x_i})$. "Far" is a function like "Close", and its negation is "Not Far". For each pair of vertices, Not Far can be computed as

$$\text{Not Far}(x, y) = \text{Max}_k^n [R^k(x, y) \wedge \neg F(k)], \quad (95)$$

where $R^k(x, y)$ is a path of length k between the two vertices x and y . Consider U_x as a set consisting of different elements pairs of the cluster of node x . For such pair $u \in U_x$, the second criterion can be calculated as [271]

$$C_2(x) = \text{Min}_{u \in U_x} [\text{NotFar}(u)]. \quad (96)$$

Here, maximally separated pair of vertices in S are found then the degree to which they are not far, is computed. The last criterion can be interpreted as every node, outside the cluster of a given node, must not be closed to most of the nodes within the cluster. Consider two nodes, one of which belongs to the cluster, and the other does not. If a node y is not in the cluster S , the extent to which node y is closer to most of the nodes in S can be computed as [271]

$$M(y/S) = \text{Most} \left(\frac{\sum_{j=1}^{n_s} \text{Close}(y, x_j)}{n_s} \right), \quad (97)$$

and for any node $x_i \in S$,

$$M(x_i/S) = \text{Most} \left(\frac{\sum_{j=1, j \neq i}^{n_s} \text{Close}(x_i, x_j)}{n_s - 1} \right). \quad (98)$$

Now, if $M(y/S) < M(x_i/S)$ for all nodes in the cluster then $C_3 = 1$, otherwise $C_3 = 0$.

Local clustering coefficient-based fuzzy link prediction (FCC). Previously, the Clustering Coefficient(CC) [44] idea is restricted just to the triangles related to a node. Here, the same idea of clustering coefficients has been used to develop a new similarity index, which considers more extensive clusters. In [271], a novel fuzzy quasi-local Clustering Coefficient(CC) model has been proposed for the link prediction task. In the model, the score of a link is calculated using the sum of clustering values of the corresponding nodes, i.e., for a given link $x_i - x_j$, the score can be calculated as $C(x_i) + C(x_j)$. And the CC for every node according to the above definition of clustering is the minimum satisfaction of criteria, i.e., for a given node x_i the CC of this node is computed as

$$CC(x_i) = \text{Min}_j [C_j(x_i)], \quad (99)$$

where C_j is the criteria j . For better prediction result, the value of α and β have been calculated by trial and error, and observed values are $\alpha = 0.3$ and $\beta = 0.7$. The FCC algorithm pseudo code is given in [271].

Cluster overlapping-based fuzzy link prediction (FCO). In the above model, overlapping between clusters has not been considered, which might be a possible condition when working with the local clustering coefficient-based fuzzy link prediction. In a social network, cluster overlapping between clusters of two nodes can be defined as those nodes belong to both clusters. Bastani et al. [271] have considered a path as cluster overlapping of two nodes, which crosses the common nodes of both clusters. A new index has been proposed by them to compute the overlapping between clusters of the nodes x_i and x_j as

$$S(x_i, x_j) = \frac{\sum_{z=1}^n \text{Close}(x_i, x_j)}{|\sum_{u,z \in S_{x_i}} W_{uz}| + |\sum_{u,z \in S_{x_j}} W_{uz}|}. \quad (100)$$

The above equation represents the ratio between the sum of the closeness of two nodes (considering every path between these two nodes) to the weighted sum of all the links inside the clusters. Here S_{x_i} represents the cluster of node x_i and W , the weight of the links. The pseudo code of the algorithm FCO is also given in [271]. The authors in [271] have created a synergy between models, called Intelligent Hybrid model shown in Fig. 13. For the creation of synergy, a two-step methodology has been used. The first step of which consists of the selection of high scored candidates based on proximity measures, and in the second step, an inference engine is generated to predict the links of the highest strength (see Figs. 13 and 14) The authors have given a pseudo code for the hybrid model in their paper. They tested their algorithm on a collaboration network among scientists, generated by Newman [274]. They demonstrated that the accuracy of the suggested hybrid model is better to that of other models. They have also demonstrated the FCO model as the second most accurate model in their implementation.

Some works [275–277] employ a fuzzy approach to link prediction framework and show their effectiveness on different network datasets. Bhawsar et al. [275] focus basically on several attributes (additional features) of nodes of a network and applied fuzzy soft set [278] and Markov model [54] to predict missing links. Moradabadi et al. [276] applied learning automata (See Fig. 14) in distributed manner (DLA) in the fuzzy social network for evaluating link scores. First, the authors

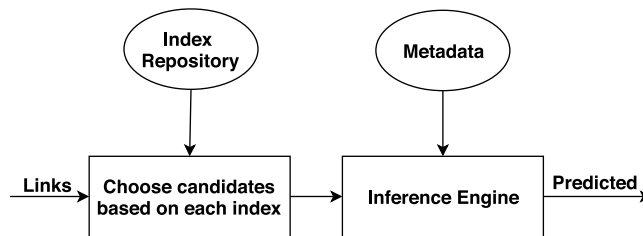


Fig. 13. The hybrid model of link prediction.

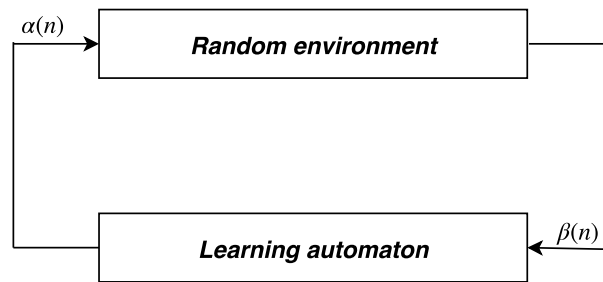


Fig. 14. The relationship between a learning automaton and its random environment [276]. Here, $\alpha(n)$ is the finite set of possible actions, and $\beta(n)$ is the set of reinforcement signals. A learning automaton is an adaptive decision-making system that improves its performance by choosing the optimal actions from a set $\alpha(n)$ through repeated interaction between learning the automaton and the random environment as shown in this figure.

convert given social network (coauthorship networks where the date of creation of links are available) to the fuzzy social network where nodes are represented by learning automata, and the links correspond to the L–R fuzzy numbers computed by creation date of links in the original network. Then the fuzzy strength of each link is calculated. Now DLA is employed on the fuzzy social network to compute the strength of each non-existing link (seed link) based on a path between the two endpoints of the seed link. More paths may exist between these two points; the proposed algorithm finds the path that minimizes the total penalties of the learning automata in the path and assigned this strength as the score of the seed link. Links are sorted based on their strengths, and top- l links are determined as predicted links.

7. Conclusion and future directions

In this survey, we have gone through several link prediction methods broadly classified into similarity-based, probabilistic models, dimensionality reduction-based, entropy-based, and clustering-based. We have also reviewed some recent approaches, including fuzzy models and link prediction in bipartite networks. The experiment of similarity-based approaches on seven network datasets has been conducted and evaluated on four well-known measures. We observed that local and quasi-local approaches perform well, usually. Global approaches are mostly based on exploring paths which are complex to compute and increase the noise in the networks. The running time of these methods with their algorithmic complexity in big 'O' notation has been reported in this survey.

Although several link prediction methods have been explored in the literature, it is still an open research problem. Several problems are yet to be explored, for example, which structural properties perform better on each technique, also how to deal with the large size of the network. Can we devise an approach to predict missing links where strengths/weights are changing with time? As outlier concept is useful to detect spam emails, so outliers detection may be another framework where link prediction approaches would make a fruitful contribution. Most real-world networks are highly sparse where the number of positive instances is very few compared to negative instances, so handling imbalanced datasets in the context of link prediction may be another aspect. Limited works on multiplex and multilayer networks are available in the literature; this can be more explored in the future. Today's communications comprise several recipients (more than 2) that shift our attention from a one-to-one communication to one-to-many and many-to-many. Link prediction can be useful in such scenarios.

Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

References

- [1] L. Lü, T. Zhou, Link prediction in complex networks: a survey, *Physica A* 390 (6) (2011) 1150–1170, <http://dx.doi.org/10.1016/j.physa.2010.11.027>, <http://www.sciencedirect.com/science/article/pii/S037843711000991X>.
- [2] D. Liben-Nowell, J. Kleinberg, The link prediction problem for social networks, in: Proceedings of the Twelfth International Conference on Information and Knowledge Management, in: CIKM '03, ACM, New York, NY, USA, 2003, pp. 556–559, <http://dx.doi.org/10.1145/956863.956972>, <http://doi.acm.org/10.1145/956863.956972>.
- [3] S.F. Adafre, M. de Rijke, Discovering missing links in wikipedia, in: Proceedings of the 3rd International Workshop on Link Discovery, LinkKDD '05, pp. 90–97.
- [4] J. Zhu, J. Hong, J.G. Hughes, Using Markov models for web site link prediction, in: Proceedings of the Thirteenth ACM Conference on Hypertext and Hypermedia, HYPERTEXT '02, pp. 169–170.
- [5] Z. Huang, X. Li, H. Chen, Link prediction approach to collaborative filtering, in: Proceedings of the 5th ACM/IEEE-CS Joint Conference on Digital Libraries, JCDL '05, pp. 141–142.
- [6] L. Lü, M. Medo, C.H. Yeung, Y.-C. Zhang, Z.-K. Zhang, T. Zhou, Recommender systems, *Phys. Rep.* 519 (1) (2012) 1–49, <http://dx.doi.org/10.1016/j.physrep.2012.02.006>, <http://www.sciencedirect.com/science/article/pii/S0370157312000828>.
- [7] E.M. Airoldi, D.M. Blei, S.E. Fienberg, E.P. Xing, T. Jaakkola, Mixed membership stochastic block models for relational data with application to protein-protein interactions, in: Proceedings of the International Biometrics Society Annual Meeting, vol. 15, 2006.
- [8] M.E.J. Newman, Clustering and preferential attachment in growing networks, *Phys. Rev. E* 64 (2001) 025102, <http://dx.doi.org/10.1103/PhysRevE.64.025102>.
- [9] M.A. Hasan, V. Chaoji, S. Salem, M. Zaki, Link prediction using supervised learning, in: Proc. of SDM 06 Workshop on Link Analysis, Counterterrorism and Security, 2006.
- [10] A. Popescul, L.H. Ungar, Statistical relational learning for link prediction, in: IJCAI Workshop on Learning Statistical Models from Relational Data, 2003.
- [11] A. Popescul, L.H. Ungar, Structural logistic regression for link analysis, Departmental Papers (CIS), 2003, p. 133.
- [12] B. Taskar, M.-F. Wong, P. Abbeel, D. Koller, Link prediction in relational data, in: Proceedings of the 16th International Conference on Neural Information Processing Systems, in: NIPS'03, MIT Press, Cambridge, MA, USA, 2003, pp. 659–666, <http://dl.acm.org/citation.cfm?id=2981345.2981428>.
- [13] R.R. Sarukkai, Link prediction and path analysis using Markov chains, *Comput. Netw.* 33 (1–6) (2000) 377–386.
- [14] L. Getoor, N. Friedman, D. Koller, B. Taskar, Learning probabilistic models of link structure, *J. Mach. Learn. Res.* 3 (2002) 679–707.
- [15] R.M. Nallapati, A. Ahmed, E.P. Xing, W.W. Cohen, Joint latent topic models for text and citations, in: Proceedings of the 14th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '08, pp. 542–550.
- [16] W. Fu, L. Song, E.P. Xing, Dynamic mixed membership blockmodel for evolving networks, in: Proceedings of the 26th Annual International Conference on Machine Learning, ICML '09, pp. 329–336.
- [17] Z. Xu, V. Tresp, S. Yu, K. Yu, Nonparametric relational learning for social network analysis, in: KDD'2008 Workshop on Social Network Mining and Analysis, 2008.
- [18] A.-L. Barabasi, R. Albert, Emergence of scaling in random networks, *Science* 286 (5439) (1999) 509–512, <http://dx.doi.org/10.1126/science.286.5439.509>, <http://science.sciencemag.org/content/286/5439/509>, [arXiv:1999.05464](http://arxiv.org/abs/1999.05464), <http://science.sciencemag.org/content/286/5439/509.full.pdf>.
- [19] J.M. Kleinberg, Navigation in a small world, *Nature* 406 (6798) (2000) 845.
- [20] J. Leskovec, J. Kleinberg, C. Faloutsos, Graphs over time: densification laws, shrinking diameters and possible explanations, in: Proceedings of the Eleventh ACM SIGKDD International Conference on Knowledge Discovery in Data Mining, in: KDD '05, ACM, New York, NY, USA, 2005, pp. 177–187, <http://dx.doi.org/10.1145/1081870.1081893>, <http://doi.acm.org/10.1145/1081870.1081893>.
- [21] W. Wang, Q. Zhang, T. Zhou, Evaluating network models: a likelihood analysis, *CoRR* abs/1112.4597 (2011) <http://arxiv.org/abs/1112.4597>, [arXiv:1112.4597](http://arxiv.org/abs/1112.4597).
- [22] Q.-M. Zhang, X.-K. Xu, Y.-X. Zhu, T. Zhou, Measuring multiple evolution mechanisms of complex networks, *Sci. Rep.* 5 (2015) 10350, <http://dx.doi.org/10.1038/srep10350>, [arXiv:1410.3519](http://arxiv.org/abs/1410.3519).
- [23] T. Tylenda, R. Angelova, S. Bedathur, Towards time-aware link prediction in evolving social networks, in: Proceedings of the 3rd Workshop on Social Network Mining and Analysis, SNA-KDD '09, pp. 9:1–9:10.
- [24] H.H. Song, T.W. Cho, V. Dave, Y. Zhang, L. Qiu, Scalable proximity estimation and link prediction in online social networks, in: Proceedings of the 9th ACM SIGCOMM Conference on Internet Measurement, in: IMC '09, ACM, New York, NY, USA, 2009, pp. 322–335, <http://dx.doi.org/10.1145/1644893.1644932>, <http://doi.acm.org/10.1145/1644893.1644932>.
- [25] D.M. Dunlavy, T.G. Kolda, E. Acar, Temporal link prediction using matrix and tensor factorizations, *ACM Trans. Knowl. Discov. Data* 5 (2) (2011) 10:1–10:27, <http://dx.doi.org/10.1145/1921632.1921636>, <http://doi.acm.org/10.1145/1921632.1921636>.
- [26] V. Martínez, F. Berzal, J.-C. Cubero, A survey of link prediction in complex networks, *ACM Comput. Surv.* 49 (4) (2016) 69:1–69:33, <http://dx.doi.org/10.1145/3012704>, <http://doi.acm.org/10.1145/3012704>.
- [27] M.A. Hasan, M.J. Zaki, A survey of link prediction in social networks, in: C.C. Aggarwal (Ed.), *Social Network Data Analytics*, Springer US, Boston, MA, 2011, pp. 243–275, http://dx.doi.org/10.1007/978-1-4419-8462-3_9.
- [28] A. Barabasi, H. Jeong, Z. Neda, E. Ravasz, A. Schubert, T. Vicsek, Evolution of the social network of scientific collaborations, *Physica A* 311 (2002) 590–614, [http://dx.doi.org/10.1016/S0378-4371\(02\)00736-7](http://dx.doi.org/10.1016/S0378-4371(02)00736-7).
- [29] A. Lada, E. Adar, Friends and neighbors on the web, *Social Networks* 25 (2003) 211–230, [http://dx.doi.org/10.1016/S0378-8733\(03\)00009-1](http://dx.doi.org/10.1016/S0378-8733(03)00009-1).
- [30] T. Zhou, L. Lü, Y.-C. Zhang, Predicting missing links via local information, *Eur. Phys. J. B* 71 (4) (2009) 623–630, <http://dx.doi.org/10.1140/epjb/e2009-00335-8>.
- [31] G. Kossinets, D.J. Watts, Origins of homophily in an evolving social network, *Am. J. Sociol.* 115 (2009) 405–450, <http://www.journals.uchicago.edu/doi/abs/10.1086/599247>.
- [32] G. Kossinets, D.J. Watts, Empirical analysis of an evolving social network, *Science* 311 (5757) (2006) 88–90, <http://dx.doi.org/10.1126/science.1116869>, <https://science.sciencemag.org/content/311/5757/88>, [arXiv:https://science.sciencemag.org/content/311/5757/88.full.pdf](http://arxiv.org/abs/https://science.sciencemag.org/content/311/5757/88.full.pdf).
- [33] P. Jaccard, Distribution de la flore alpine dans le bassin des Dranses et dans quelques régions voisines, *Bull. Soc. Vaud. Sci. Nat.* 37 (1901) 241–272.
- [34] Q. Ou, Y.-D. Jin, T. Zhou, B.-H. Wang, B.-Q. Yin, Power-law strength-degree correlation from resource-allocation dynamics on weighted networks, *Phys. Rev. E* 75 (2) (2007) <http://dx.doi.org/10.1103/physreve.75.021102>.
- [35] G. Salton, M.J. McGill, *Introduction to Modern Information Retrieval*, McGraw-Hill, Inc., New York, NY, USA, 1986.
- [36] T. Sørensen, A method of establishing groups of equal amplitude in plant sociology based on similarity of species and its application to analyses of the vegetation on Danish commons, *Biol. Skr.* 5 (1948) 1–34.
- [37] B. McCune, J.B. Grace, D.L. Urban, *Analysis of Ecological Communities*, MjM Software Design, Gleneden Beach, Oregon, 2002.

- [38] C.V. Cannistraci, G. Alanis-Lobato, T. Ravasi, From link-prediction in brain connectomes and protein interactomes to the local-community-paradigm in complex networks, *Sci. Rep.* 3 (2013) 1613, <http://dx.doi.org/10.1038/srep01613>.
- [39] E. Ravasz, A.L. Somera, D.A. Mongru, Z.N. Oltvai, A.-L. Barabási, Hierarchical organization of modularity in metabolic networks, *Science* 297 (5586) (2002) 1551–1555, <http://dx.doi.org/10.1126/science.1073374>, <http://science.sciencemag.org/content/297/5586/1551>, arXiv:<http://science.sciencemag.org/content/297/5586/1551.full.pdf>.
- [40] Z. Liu, Q.-M. Zhang, L. Lü, T. Zhou, Link prediction in complex networks: a local naïve Bayes model, *Europhys. Lett.* 96 (4) (2011) 48007, <http://stacks.iop.org/0295-5075/96/i=4/a=48007>.
- [41] E.A. Leicht, P. Holme, M.E.J. Newman, Vertex similarity in networks, *Phys. Rev. E* 73 (2006) 026120, <http://dx.doi.org/10.1103/PhysRevE.73.026120>, <https://link.aps.org/doi/10.1103/PhysRevE.73.026120>.
- [42] Z. Wu, Y. Lin, J. Wang, S. Gregory, Link prediction with node clustering coefficient, *Physica A* 452 (2016) 1–8, <http://dx.doi.org/10.1016/j.physa.2016.01.038>.
- [43] Z. Wu, Y. Lin, H. Wan, W. Jamil, Predicting top-L missing links with node and link clustering information in large-scale networks, *J. Stat. Mech. Theory Exp.* 8 (2016) 083202, <http://dx.doi.org/10.1088/1742-5468/2016/08/083202>.
- [44] D.J. Watts, S.H. Strogatz, Collective dynamics of ‘small-world’ networks, *Nature* 393 (6684) (1998) 440–442, <http://dx.doi.org/10.1038/30918>.
- [45] P.W. Holland, S. Leinhardt, Transitivity in structural models of small groups, *Comp. Group Stud.* 2 (2) (1971) 107–124, <http://dx.doi.org/10.1177/104649647100200201>, arXiv:<https://doi.org/10.1177/104649647100200201>.
- [46] L. Katz, A new status index derived from sociometric analysis, *Psychometrika* 18 (1) (1953) 39–43.
- [47] H. Tong, C. Faloutsos, J.-Y. Pan, Fast random walk with restart and its applications, in: Proceedings of the Sixth International Conference on Data Mining, in: ICDM '06, IEEE Computer Society, Washington, DC, USA, 2006, pp. 613–622, <http://dx.doi.org/10.1109/ICDM.2006.70>.
- [48] E.W. Dijkstra, A note on two problems in connexion with graphs, *Numer. Math.* 1 (1) (1959) 269–271, <http://dx.doi.org/10.1007/BF01386390>.
- [49] R. Bellman, On a routing problem, *Quart. Appl. Math.* 16 (1958) 87–90.
- [50] R.W. Floyd, Algorithm 97: shortest path, *Commun. ACM* 5 (6) (1962) 345, <http://dx.doi.org/10.1145/367766.368168>, <http://doi.acm.org/10.1145/367766.368168>.
- [51] F. Fouss, A. Pirotte, J.-M. Renders, M. Saerens, Random-walk computation of similarities between nodes of a graph with application to collaborative recommendation, *IEEE Trans. Knowl. Data Eng.* 19 (3) (2007) 355–369, <http://dx.doi.org/10.1109/TKDE.2007.46>.
- [52] D.A. Spielman, Spectral graph theory and its applications, in: 48th Annual IEEE Symposium on Foundations of Computer Science, FOCS'07, 2007, pp. 29–38, <http://dx.doi.org/10.1109/FOCS.2007.56>.
- [53] W. Liu, L. Lü, Link prediction based on local random walk, *Europhys. Lett.* 89 (5) (2010) 58007, <http://stacks.iop.org/0295-5075/89/i=5/a=58007>.
- [54] J.R. Norris, Markov Chains, in: Cambridge Series in Statistical and Probabilistic Mathematics, Cambridge University Press, 1998, <http://dx.doi.org/10.1017/CBO9780511810633>.
- [55] J.G. Kemeny, J.L. Snell, Finite Markov chains, New York Springer-Verlag, 1983, <http://openlibrary.org/books/OL3174668M>, Bibliography: p. 224. Originally published: Princeton, N.J. : Van Nostrand, 1960.
- [56] F. Göbel, A. Jagers, Random walks on graphs, *Stochastic Process. Appl.* 2 (4) (1974) 311–336, [http://dx.doi.org/10.1016/0304-4149\(74\)90001-5](http://dx.doi.org/10.1016/0304-4149(74)90001-5), <http://www.sciencedirect.com/science/article/pii/0304414974900015>.
- [57] D.J. Klein, M. Randić, Resistance distance, *J. Math. Chem.* 12 (1) (1993) 81–95, <http://dx.doi.org/10.1007/BF01164627>.
- [58] P. Chebotarev, E. Shamis, The matrix-forest theorem and measuring relations in small social groups, *CoRR abs/math/0602070* (2006) <http://arxiv.org/abs/math/0602070>, arXiv:<http://arxiv.org/abs/math/0602070>.
- [59] G. Jeh, J. Widom, Simrank: a measure of structural-context similarity, in: Proceedings of the Eighth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, in: KDD '02, ACM, New York, NY, USA, 2002, pp. 538–543, <http://dx.doi.org/10.1145/775047.775126>, <http://doi.acm.org/10.1145/775047.775126>.
- [60] S. Brin, L. Page, The anatomy of a large-scale hypertextual web search engine, *Comput. Netw. ISDN Syst.* 30 (1–7) (1998) 107–117, [http://dx.doi.org/10.1016/S0169-7552\(98\)00110-X](http://dx.doi.org/10.1016/S0169-7552(98)00110-X).
- [61] F. Chung, W. Zhao, Pagerank and random walks on graphs, in: G.O.H. Katona, A. Schrijver, T. Szőnyi, G. Sági (Eds.), Fete of Combinatorics and Computer Science, Springer Berlin Heidelberg, Berlin, Heidelberg, 2010, pp. 43–62, http://dx.doi.org/10.1007/978-3-642-13580-4_3.
- [62] X. Wang, X. Zhang, C. Zhao, Z. Xie, S. Zhang, D. Yi, Predicting link directions using local directed path, *Physica A* 419 (2015) 260–267, <http://dx.doi.org/10.1016/j.physa.2014.10.007>, <http://www.sciencedirect.com/science/article/pii/S0378437114408450>.
- [63] L. Lü, C.-H. Jin, T. Zhou, Similarity index based on local paths for link prediction of complex networks, *Phys. Rev. E* 80 (2009) 046122, <http://dx.doi.org/10.1103/PhysRevE.80.046122>.
- [64] I.A. Kovács, K. Luck, K. Spirohn, Y. Wang, C. Pollis, S. Schlabach, W. Bian, D.-K. Kim, N. Kishore, T. Hao, M.A. Calderwood, M. Vidal, A.-L. Barabási, Network-based prediction of protein interactions, *bioRxiv* (2018) <http://dx.doi.org/10.1101/275529>, <https://www.biorxiv.org/content/early/2018/03/02/275529>.
- [65] M.S. Granovetter, The strength of weak ties, *Am. J. Sociol.* 78 (6) (1973) 1360–1380, <http://www.jstor.org/stable/2776392>.
- [66] R. Pech, D. Hao, Y. Lee, Y. Yuan, T. Zhou, Link prediction via linear optimization, *CoRR abs/1804.00124* (2018) <http://arxiv.org/abs/1804.00124>, arXiv:<http://arxiv.org/abs/1804.00124>.
- [67] A. Muscoloni, I. Abdelhamid, C.V. Cannistraci, Local-community network automata modelling based on length-three-paths for prediction of complex network structures in protein interactomes, food webs and more, *bioRxiv* (2018) 346916, <http://dx.doi.org/10.1101/346916>, <https://www.biorxiv.org/content/early/2018/06/18/346916>.
- [68] T. Zhou, Y. Lee, G. Wang, Experimental analyses on 2-hop-based and 3-hop-based link prediction algorithms, *CoRR abs/1909.00174* (2019) <http://arxiv.org/abs/1909.00174>.
- [69] C. Wang, V. Satuluri, S. Parthasarathy, Local probabilistic models for link prediction, in: Proceedings of the 2007 Seventh IEEE International Conference on Data Mining, ICDM '07, pp. 322–331.
- [70] J. Neville, *Statistical models and analysis techniques for learning in relational data* (Ph.D. thesis), 2006.
- [71] K. Yu, W. Chu, S. Yu, V. Tresp, Z. Xu, Stochastic relational models for discriminative link prediction, in: Proceedings of the 19th International Conference on Neural Information Processing Systems, in: NIPS'06, MIT Press, Cambridge, MA, USA, 2006, pp. 1553–1560, <http://dl.acm.org/citation.cfm?id=2976456.2976651>.
- [72] A. Clauset, C. Moore, M.E.J. Newman, Hierarchical structure and the prediction of missing links in networks, *Nature* 453 (7191) (2008) 98–101, <http://dx.doi.org/10.1038/nature06830>.
- [73] R. Guimerà, M. Sales-Pardo, Missing and spurious interactions and the reconstruction of complex networks, *Proc. Natl. Acad. Sci.* 106 (52) (2009) 22073–22078, <http://dx.doi.org/10.1073/pnas.0908366106>, <https://www.pnas.org/content/106/52/22073>, arXiv:<https://www.pnas.org/content/106/52/22073.full.pdf>.
- [74] P. Wang, B. Xu, Y. Wu, X. Zhou, Link prediction in social networks: the state-of-the-art, *CoRR abs/1411.5118* (2014) <http://arxiv.org/abs/1411.5118>, arXiv:<http://arxiv.org/abs/1411.5118>.
- [75] N. Stanley, T. Bonacci, R. Kwitt, M. Niethammer, P.J. Mucha, Stochastic block models with multiple continuous attributes, *CoRR abs/1803.02726* (2018) <http://arxiv.org/abs/1803.02726>, arXiv:<http://arxiv.org/abs/1803.02726>.

- [76] T. Vallès-Català, T.P. Peixoto, M. Sales-Pardo, R. Guimerà, Consistencies and inconsistencies between model selection and link prediction in networks, *Phys. Rev. E* 97 (6) (2018) 062316, <http://dx.doi.org/10.1103/physreve.97.062316>, <https://app.dimensions.ai/details/publication/pub.1105229576>, <http://arxiv.org/pdf/1705.07967>, Exported from <https://app.dimensions.ai> on 2019/01/03.
- [77] H. Kashima, N. Abe, A parameterized probabilistic model of network evolution for supervised link prediction, in: Proceedings of the Sixth International Conference on Data Mining, in: ICDM'06, IEEE Computer Society, Washington, DC, USA, 2006, pp. 340–349, <http://dx.doi.org/10.1109/ICDM.2006.8>.
- [78] S.A. Williamson, Nonparametric network models for link prediction, *J. Mach. Learn. Res.* 17 (1) (2016) 7102–7121, <http://dl.acm.org/citation.cfm?id=2946645.3053484>.
- [79] T.-T. Kuo, R. Yan, Y.-Y. Huang, P.-H. Kung, S.-D. Lin, Unsupervised link prediction using aggregative statistics on heterogeneous social networks, in: Proceedings of the 19th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, in: KDD '13, ACM, New York, NY, USA, 2013, pp. 775–783, <http://dx.doi.org/10.1145/2487575.2487614>, <http://doi.acm.org/10.1145/2487575.2487614>.
- [80] J. Yang, J.J. McAuley, J. Leskovec, Community detection in networks with node attributes, in: 2013 IEEE 13th International Conference on Data Mining, Dallas, TX, USA, December 7–10, 2013, 2013, pp. 1151–1156 <http://dx.doi.org/10.1109/ICDM.2013.167>, <https://doi.org/10.1109/ICDM.2013.167>.
- [81] R. Agrawal, R. Srikant, Fast algorithms for mining association rules in large databases, in: Proceedings of the 20th International Conference on Very Large Data Bases, VLDB '94, pp. 487–499, <http://dl.acm.org/citation.cfm?id=645920.672836>.
- [82] J. Han, J. Pei, Y. Yin, Mining frequent patterns without candidate generation, in: Proceedings of the 2000 ACM SIGMOD International Conference on Management of Data, SIGMOD '00, pp. 1–12.
- [83] M.J. Zaki, S. Parthasarathy, M. Ogihara, W. Li, New algorithms for fast discovery of association rules, in: Proceedings of the Third International Conference on Knowledge Discovery and Data Mining, KDD'97, 1997, pp. 283–286, <http://dl.acm.org/citation.cfm?id=3001392.3001454>.
- [84] K. Singh, A. Kumar, S.S. Singh, H.K. Shakya, B. Biswas, EHNL: an efficient algorithm for mining high utility itemsets with negative utility value and length constraints, *Inform. Sci.* 484 (2019) 44–70, <http://dx.doi.org/10.1016/j.ins.2019.01.056>, <http://www.sciencedirect.com/science/article/pii/S0020025519300696>.
- [85] K. Singh, S.S. Singh, A. Kumar, B. Biswas, TKEH: an efficient algorithm for mining top-k high utility itemsets, *Appl. Intell.* 49 (3) (2019) 1078–1097, <http://dx.doi.org/10.1007/s10489-018-1316-x>.
- [86] T. Calders, B. Goethals, Depth-first non-derivable itemset mining, in: Proceedings of the 2005 SIAM International Conference on Data Mining, SDM 2005, Newport Beach, CA, USA, April 21–23, 2005, pp. 250–261, <http://dx.doi.org/10.1137/1.9781611972757.23>, <https://doi.org/10.1137/1.9781611972757.23>.
- [87] D.N. Pavlov, H. Mannila, P. Smyth, Beyond independence: probabilistic models for query approximation on binary transaction data, *IEEE Trans. Knowl. Data Eng.* 15 (6) (2003) 1409–1421, <http://dx.doi.org/10.1109/TKDE.2003.1245281>.
- [88] O. Frank, D. Strauss, Markov graphs, *J. Amer. Statist. Assoc.* 81 (395) (1986) 832–842, <http://www.jstor.org/stable/2289017>.
- [89] C. Wang, S. Parthasarathy, Summarizing itemset patterns using probabilistic models, in: Proceedings of the 12th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, in: KDD '06, ACM, New York, NY, USA, 2006, pp. 730–735, <http://dx.doi.org/10.1145/1150402.1150495>, <http://doi.acm.org/10.1145/1150402.1150495>.
- [90] S.L. Lauritzen, D.J. Spiegelhalter, Local computations with probabilities on graphical structures and their application to expert systems, *J. Royal Stat. Soc. B* 50 (1988) 157–224.
- [91] B. Taskar, M.F. Wong, P. Abbeel, D. Koller, Link prediction in relational data, in: Advances in Neural Information Processing Systems 16 [Neural Information Processing Systems, NIPS 2003, December 8–13, 2003, Vancouver and Whistler, British Columbia, Canada], 2003, pp. 659–666, <http://papers.nips.cc/paper/2465-link-prediction-in-relational-data>.
- [92] M. Sales-Pardo, R. Guimerà, A.A. Moreira, L.A.N. Amaral, Extracting the hierarchical organization of complex systems, *Proc. Natl. Acad. Sci.* 104 (39) (2007) 15224–15229, <http://dx.doi.org/10.1073/pnas.0703740104>, <https://app.dimensions.ai/details/publication/pub.1027402365>, <http://www.pnas.org/content/104/39/15224.full.pdf>, Exported from <https://app.dimensions.ai> on 2019/02/14.
- [93] H. White, S. Boorman, R. Breiger, Social structure from multiple networks. i. blockmodels of roles and positions, *Am. J. Sociol.* 81 (4) (1976) 730–780.
- [94] P.W. Holland, K.B. Laskey, S. Leinhardt, Stochastic blockmodels: first steps, *Social Networks* 5 (2) (1983) 109–137, [http://dx.doi.org/10.1016/0378-8733\(83\)90021-7](http://dx.doi.org/10.1016/0378-8733(83)90021-7), <http://www.sciencedirect.com/science/article/pii/0378873383900217>.
- [95] I. Beichl, F. Sullivan, The metropolis algorithm, *Comput. Sci. Eng.* 2 (1) (2000) 65–69, <http://dx.doi.org/10.1109/5992.814660>.
- [96] P.W. Holland, S. Leinhardt, An exponential family of probability distributions for directed graphs, *J. Amer. Statist. Assoc.* 76 (373) (1981) 33–50, <http://dx.doi.org/10.1080/01621459.1981.10477598>, <https://amstat.tandfonline.com/doi/abs/10.1080/01621459.1981.10477598>, arXiv:<https://arxiv.org/abs/10.1080/01621459.1981.10477598>, <https://amstat.tandfonline.com/doi/pdf/10.1080/01621459.1981.10477598>.
- [97] S. Wasserman, P. Pattison, Logit models and logistic regressions for social networks: I. An introduction to Markov graphs andp, *Psychometrika* 61 (3) (1996) 401–425, <http://dx.doi.org/10.1007/BF02294547>.
- [98] J. Park, M.E.J. Newman, Solution of the two-star model of a network, *Phys. Rev. E* 70 (6) (2004) <http://dx.doi.org/10.1103/physreve.70.066146>.
- [99] L. Pan, T. Zhou, L. Lü, C.-K. Hu, Predicting missing links and identifying spurious links via likelihood analysis, *Sci. Rep.* 6 (1) (2016) 22955, <http://dx.doi.org/10.1038/srep22955>, <https://app.dimensions.ai/details/publication/pub.1006608090>, <http://www.nature.com/articles/srep22955.pdf>, Exported from <https://app.dimensions.ai> on 2019/02/14.
- [100] A. Pecli, B. Giovanini, C. Pacheco, C. Moreira, F. Ferreira, F. Tosta, J. Tesolin, M. Dias, S. Filho, M.C. Cavalcanti, R. Goldschmidt, Dimensionality reduction for supervised learning in link prediction problems, in: ICEIS 2015 - 17th International Conference on Enterprise Information Systems, Proceedings, vol. 1, 2015, pp. 295–302.
- [101] K. Fukumizu, F.R. Bach, M.I. Jordan, Dimensionality reduction for supervised learning with reproducing kernel Hilbert spaces, *J. Mach. Learn. Res.* 5 (2004) 73–99, <http://dl.acm.org/citation.cfm?id=1005332.1005335>.
- [102] B. Perozzi, R. Al-Rfou, S. Skiena, Deepwalk: online learning of social representations, in: Proceedings of the 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, in: KDD '14, ACM, New York, NY, USA, 2014, pp. 701–710, <http://dx.doi.org/10.1145/2623330.2623732>, <http://doi.acm.org/10.1145/2623330.2623732>.
- [103] A. Ahmed, N. Shervashidze, S. Narayanamurthy, V. Josifovski, A.J. Smola, Distributed large-scale natural graph factorization, in: Proceedings of the 22Nd International Conference on World Wide Web, in: WWW '13, ACM, New York, NY, USA, 2013, pp. 37–48, <http://dx.doi.org/10.1145/2488388.2488393>, <http://doi.acm.org/10.1145/2488388.2488393>.
- [104] S. Cao, W. Lu, Q. Xu, GraRep: learning graph representations with global structural information, in: Proceedings of the 24th ACM International on Conference on Information and Knowledge Management, in: CIKM '15, ACM, New York, NY, USA, 2015, pp. 891–900, <http://dx.doi.org/10.1145/2806416.2806512>, <http://doi.acm.org/10.1145/2806416.2806512>.
- [105] M. Ou, P. Cui, J. Pei, Z. Zhang, W. Zhu, Asymmetric transitivity preserving graph embedding, in: Proceedings of the 22Nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, in: KDD '16, ACM, New York, NY, USA, 2016, pp. 1105–1114, <http://dx.doi.org/10.1145/2939672.2939751>, <http://doi.acm.org/10.1145/2939672.2939751>.

- [106] M. Belkin, P. Niyogi, Laplacian eigenmaps and spectral techniques for embedding and clustering, in: Proceedings of the 14th International Conference on Neural Information Processing Systems: Natural and Synthetic, in: NIPS'01, MIT Press, Cambridge, MA, USA, 2001, pp. 585–591, <http://dl.acm.org/citation.cfm?id=2980539.2980616>.
- [107] A. Grover, J. Leskovec, Node2vec: scalable feature learning for networks, in: Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, in: KDD '16, ACM, New York, NY, USA, 2016, pp. 855–864, <http://dx.doi.org/10.1145/2939672.2939754>.
- [108] S.M. Kazemi, D. Poole, Simple embedding for link prediction in knowledge graphs, in: Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, 3–8 December 2018, Montréal, Canada, 2018, pp. 4289–4300.
- [109] S.T. Roweis, L.K. Saul, Nonlinear dimensionality reduction by locally linear embedding, *Science* 290 (5500) (2000) 2323–2326, <http://dx.doi.org/10.1126/science.290.5500.2323>, <http://science.sciencemag.org/content/290/5500/2323>, arXiv:<http://science.sciencemag.org/content/290/5500/2323.full.pdf>.
- [110] J.B. Tenenbaum, V.d. Silva, J.C. Langford, A global geometric framework for nonlinear dimensionality reduction, *Science* 290 (5500) (2000) 2319–2323, <http://dx.doi.org/10.1126/science.290.5500.2319>, <https://science.sciencemag.org/content/290/5500/2319>, arXiv:<https://science.sciencemag.org/content/290/5500/2319.full.pdf>.
- [111] O. Kuchaiev, M. Kasajski, D.J. Higham, N. Przulj, Geometric de-noising of protein-protein interaction networks, *PLoS Comput. Biol.* 5 (8) (2009) <http://dx.doi.org/10.1371/journal.pcbi.1000454>.
- [112] T. Mikolov, K. Chen, G. Corrado, J. Dean, Efficient estimation of word representations in vector space, *CoRR abs/1301.3781* (2013) <http://arxiv.org/abs/1301.3781>, arXiv:1301.3781.
- [113] T. Mikolov, I. Sutskever, K. Chen, G. Corrado, J. Dean, Distributed representations of words and phrases and their compositionality, *CoRR abs/1301.4546* (2013) <http://arxiv.org/abs/1301.4546>, arXiv:1301.4546.
- [114] T. Trouillon, J. Welbl, S. Riedel, E. Gaussier, G. Bouchard, Complex embeddings for simple link prediction, in: Proceedings of the 33rd International Conference on International Conference on Machine Learning - Volume 48, in: ICML'16, JMLR.org, 2016, pp. 2071–2080, <http://dl.acm.org/citation.cfm?id=3045390.3045609>.
- [115] Z. Cao, L. Wang, G. de Melo, Link prediction via subgraph embedding-based convex matrix completion, in: Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence, (AAAI-18), the 30th innovative Applications of Artificial Intelligence (IAAI-18), and the 8th AAAI Symposium on Educational Advances in Artificial Intelligence (EAAI-18), New Orleans, Louisiana, USA, February 2-7, 2018, 2018, pp. 2803–2810 <https://www.aaai.org/ocs/index.php/AAAI/AAAI18/paper/view/16442>.
- [116] T. Li, J. Zhang, P.S. Yu, Y. Zhang, Y. Yan, Deep dynamic network embedding for link prediction, *IEEE Access* 6 (2018) 29219–29230, <http://dx.doi.org/10.1109/ACCESS.2018.2839770>, <https://doi.org/10.1109/ACCESS.2018.2839770>.
- [117] H. Chen, B. Perozzi, Y. Hu, S. Skiena, HARP: Hierarchical representation learning for networks, *CoRR abs/1706.07845* (2017) <http://arxiv.org/abs/1706.07845>, arXiv:1706.07845.
- [118] B. Perozzi, V. Kulkarni, H. Chen, S. Skiena, Don't walk, skip!: online learning of multi-scale network embeddings, in: Proceedings of the 2017 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining 2017, Sydney, Australia, July 31–August 03, 2017, 2017, pp. 258–265, <https://doi.org/10.1145/3110025.3110086>.
- [119] J. Tang, M. Qu, M. Wang, M. Zhang, J. Yan, Q. Mei, LINE: Large-scale information network embedding, *CoRR abs/1503.03578* (2015) <http://arxiv.org/abs/1503.03578>, arXiv:1503.03578.
- [120] D. Wang, P. Cui, W. Zhu, Structural deep network embedding, in: Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, in: KDD '16, ACM, New York, NY, USA, 2016, pp. 1225–1234, <http://dx.doi.org/10.1145/2939672.2939753>.
- [121] S. Cao, W. Lu, Q. Xu, Deep neural networks for learning graph representations, in: Proceedings of the Thirtieth AAAI Conference on Artificial Intelligence, in: AAAI'16, AAAI Press, 2016, pp. 1145–1152, <http://dl.acm.org/citation.cfm?id=3015812.3015982>.
- [122] T.N. Kipf, M. Welling, Semi-supervised classification with graph convolutional networks, *CoRR abs/1609.02907* (2016) <http://arxiv.org/abs/1609.02907>, arXiv:1609.02907.
- [123] T.N. Kipf, M. Welling, Variational graph auto-encoders, *CoRR abs/1611.07308* (2016) <http://arxiv.org/abs/1611.07308>, arXiv:1611.07308.
- [124] M. Zhang, Y. Chen, Link prediction based on graph neural networks, *CoRR abs/1802.09691* (2018) <http://arxiv.org/abs/1802.09691>, arXiv:1802.09691.
- [125] S. Pan, R. Hu, G. Long, J. Jiang, L. Yao, C. Zhang, Adversarially regularized graph autoencoder, *CoRR abs/1802.04407* (2018) <http://arxiv.org/abs/1802.04407>, arXiv:1802.04407.
- [126] E. Acar, D.M. Dunlavy, T.G. Kolda, Link prediction on evolving data using matrix and tensor factorizations, in: 2009 IEEE International Conference on Data Mining Workshops, 2009, pp. 262–269, <http://dx.doi.org/10.1109/ICDMW.2009.54>.
- [127] X. Ma, P. Sun, G. Qin, Nonnegative matrix factorization algorithms for link prediction in temporal networks using graph communication, *Pattern Recognit.* 71 (2017) 361–374, <http://dx.doi.org/10.1016/j.patcog.2017.06.025>, <http://www.sciencedirect.com/science/article/pii/S0031320317302480>.
- [128] U. Sharan, J. Neville, Temporal-relational classifiers for prediction in evolving domains, in: 2008 Eighth IEEE International Conference on Data Mining, 2008, pp. 540–549, <http://dx.doi.org/10.1109/ICDM.2008.125>.
- [129] A.K. Menon, C. Elkan, Link prediction via matrix factorization, in: Proceedings of the 2011 European Conference on Machine Learning and Knowledge Discovery in Databases - Volume Part II, in: ECML PKDD'11, Springer-Verlag, Berlin, Heidelberg, 2011, pp. 437–452, <http://dl.acm.org/citation.cfm?id=2034117.2034146>.
- [130] B. Chen, F. Li, S. Chen, R. Hu, L. Chen, Link prediction based on non-negative matrix factorization, *PLOS ONE* 12 (8) (2017) 1–18, <http://dx.doi.org/10.1371/journal.pone.0182968>.
- [131] W. Wang, F. Cai, P. Jiao, L. Pan, A perturbation-based framework for link prediction via non-negative matrix factorization, *Sci. Rep.* 6 (2016) 38938, <http://dx.doi.org/10.1038/srep38938>.
- [132] N.M. Ahmed, L. Chen, Y. Wang, B. Li, Y. Li, W. Liu, DeepEye: link prediction in dynamic networks based on non-negative matrix factorization, *Big Data Min. Anal.* 1 (1) (2018) 19–33, <http://dx.doi.org/10.26599/BDMA.2017.9020002>.
- [133] G. Liyuan, W. Zhiqiang, L. Jiye, Link prediction algorithm by matrix factorization based on importance of edges 31 (2) 150 <http://dx.doi.org/10.16451/j.cnki.issn1003-6059.201802006>, http://manu46.magtech.com.cn/jweb_prai/EN/abstract/article_11455.shtml.
- [134] Y. Koren, R. Bell, C. Volinsky, Matrix factorization techniques for recommender systems, *Computer* 42 (8) (2009) 30–37, <http://dx.doi.org/10.1109/MC.2009.263>.
- [135] Z. Wu, Y. Chen, Link prediction using matrix factorization with bagging, in: 2016 IEEE/ACIS 15th International Conference on Computer and Information Science, ICIS, 2016, pp. 1–6, <http://dx.doi.org/10.1109/ICIS.2016.7550942>.
- [136] C.H.Q. Ding, T. Li, M.I. Jordan, Convex and semi-nonnegative matrix factorizations, *IEEE Trans. Pattern Anal. Mach. Intell.* 32 (1) (2010) 45–55, <http://dx.doi.org/10.1109/TPAMI.2008.277>.
- [137] J.D. Carroll, J.-J. Chang, Analysis of individual differences in multidimensional scaling via an n-way generalization of “Eckart-Young” decomposition, *Psychometrika* 35 (3) (1970) 283–319, <http://dx.doi.org/10.1007/BF02310791>.

- [138] R.A. Harshman, Foundations of the PARAFAC procedure: models and conditions for an 'explanatory' multi-modal factor analysis, *UCLA Working Papers in Phonetics* 16, 1970, pp. 1–84.
- [139] R. Pech, D. Hao, L. Pan, H. Cheng, T. Zhou, Link prediction via matrix completion, *CoRR* abs/1606.06812 (2016) <http://arxiv.org/abs/1606.06812>, [arXiv:1606.06812](https://arxiv.org/abs/1606.06812).
- [140] R. Lichtenwalter, N.V. Chawla, Link prediction: fair and effective evaluation, in: *Proceedings of the 2012 International Conference on Advances in Social Networks Analysis and Mining (ASONAM 2012)*, in: *ASONAM '12*, IEEE Computer Society, Washington, DC, USA, 2012, pp. 376–383, <http://dx.doi.org/10.1109/ASONAM.2012.68>.
- [141] J.R. Doppa, J. Yu, P. Tadepalli, L. Getoor, Learning algorithms for link prediction based on chance constraints, in: *Proceedings of the 2010 European Conference on Machine Learning and Knowledge Discovery in Databases: Part I*, in: *ECML PKDD'10*, Springer-Verlag, Berlin, Heidelberg, 2010, pp. 344–360, <http://dl.acm.org/citation.cfm?id=1888258.1888288>.
- [142] K. Anand, G. Bianconi, Entropy measures for networks: toward an information theory of complex topologies, *Phys. Rev. E* 80 (2009) 045102, <http://dx.doi.org/10.1103/PhysRevE.80.045102>.
- [143] R.V. Sole, S. Valverde, Information theory of complex networks: on evolution and architectural constraints, in: E. Ben-Naim, H. Frauenfelder, Z. Toroczkai (Eds.), *Complex Networks*, 2004, pp. 189–207, http://dx.doi.org/10.1007/978-3-540-44485-5_9.
- [144] M. Bauer, D. Bernard, Maximal entropy random networks with given degree distribution, 2002, [arXiv:cond-mat/0206150](https://arxiv.org/abs/cond-mat/0206150).
- [145] F. Tan, Y. Xia, B. Zhu, Link prediction in complex networks: a mutual information perspective, *PLOS ONE* 9 (9) (2014) 1–8, <http://dx.doi.org/10.1371/journal.pone.0107056>.
- [146] B. Zhu, Y. Xia, An information-theoretic model for link prediction in complex networks, *Sci. Rep.* 5 (2015) 13707, <http://dx.doi.org/10.1038/srep13707>.
- [147] Z. Xu, C. Pu, J. Yang, Link prediction based on path entropy, *Physica A* 456 (2016) 294–301, <http://dx.doi.org/10.1016/j.physa.2016.03.091>.
- [148] Z. Xu, C. Pu, R.R. Sharafat, L. Li, J. Yang, Entropy-based link prediction in weighted networks, *Chin. Phys. B* 26 (1) (2017) 018902, <http://stacks.iop.org/1674-1056/26/j=1/a=018902>.
- [149] T. Wang, X.-S. He, M.-Y. Zhou, Z.-Q. Fu, Link prediction in evolving networks based on popularity of nodes, *Sci. Rep.* 7 (2017) 7147, <http://dx.doi.org/10.1038/s41598-017-07315-4>.
- [150] L. Yin, H. Zheng, T. Bian, Y. Deng, An evidential link prediction method and link predictability based on shannon entropy, *Physica A* 482 (2017) 699–712, <http://dx.doi.org/10.1016/j.physa.2017.04.106>, <http://www.sciencedirect.com/science/article/pii/S0378437117304302>.
- [151] F. Parisi, G. Caldarelli, T. Squartini, Entropy-based approach to missing-links prediction, *Appl. Netw. Sci.* 3 (1) (2018) 17, <http://dx.doi.org/10.1007/s41109-018-0073-4>.
- [152] Y.-W. Niu, H. Liu, G.-H. Wang, G.-Y. Yan, Maximal entropy random walk on heterogenous network for MIRNA-disease association prediction, *Math. Biosci.* 306 (2018) 1–9, <http://dx.doi.org/10.1016/j.mbs.2018.10.004>, <http://www.sciencedirect.com/science/article/pii/S0025556418304668>.
- [153] Z. Huang, Link prediction based on graph topology: the predictive value of generalized clustering coefficient, in: *Proceedings of the Workshop on Link Analysis*, 2006, <http://dx.doi.org/10.2139/ssrn.1634014>.
- [154] R. Albert, A.-L. Barabási, Statistical mechanics of complex networks, *Rev. Modern Phys.* 74 (2002) 47–97, <http://dx.doi.org/10.1103/RevModPhys.74.47>, <https://link.aps.org/doi/10.1103/RevModPhys.74.47>.
- [155] A. Fronczak, J.A. Holyst, M. Jedynak, J. Sienkiewicz, Higher order clustering coefficients in Barabási–Albert networks, *Physica A* 316 (1) (2002) 688–694, [http://dx.doi.org/10.1016/S0378-4371\(02\)01336-5](http://dx.doi.org/10.1016/S0378-4371(02)01336-5), <http://www.sciencedirect.com/science/article/pii/S0378437102013365>.
- [156] Y. Liu, C. Zhao, X. Wang, Q. Huang, X. Zhang, D. Yi, The degree-related clustering coefficient and its application to link prediction, *Physica A* 454 (2016) 24–33, <http://dx.doi.org/10.1016/j.physa.2016.02.014>.
- [157] A. Kumar, S.S. Singh, K. Singh, B. Biswas, Level-2 node clustering coefficient-based link prediction, *Appl. Intell.* 49 (7) (2019) 2762–2779, <http://dx.doi.org/10.1007/s10489-019-01413-8>.
- [158] A.R. Benson, R. Abebe, M.T. Schaub, A. Jadbabaie, J.M. Kleinberg, Simplicial closure and higher-order link prediction, *Proc. Natl. Acad. Sci. USA* 115 (48) (2018) E11221–E11230, <http://dx.doi.org/10.1073/pnas.1800683115>.
- [159] L. Lü, L. Pan, T. Zhou, Y.-C. Zhang, H.E. Stanley, Toward link predictability of complex networks, *Proc. Natl. Acad. Sci.* 112 (8) (2015) 2325–2330, <http://dx.doi.org/10.1073/pnas.1424644112>, <http://www.pnas.org/content/112/8/2325>, [arXiv:1408.2325](https://arxiv.org/abs/1408.2325), <http://www.pnas.org/content/112/8/2325.full.pdf>.
- [160] C.D. Manning, P. Raghavan, H. Schütze, *Introduction to Information Retrieval*, Cambridge University Press, New York, NY, USA, 2008.
- [161] J.A. Hanley, B.J. McNeil, The meaning and use of the area under a receiver operating characteristic (ROC) curve, *Radiology* 143 (1) (1982) 29–36, <http://dx.doi.org/10.1148/radiology.143.1.7063747>.
- [162] T. Fawcett, An introduction to ROC analysis, *Pattern Recognit. Lett.* 27 (8) (2006) 861–874, <http://dx.doi.org/10.1016/j.patrec.2005.10.010>.
- [163] J. Davis, M. Goadrich, The relationship between precision-recall and roc curves, in: *Proceedings of the 23rd International Conference on Machine Learning*, in: *ICML '06*, ACM, New York, NY, USA, 2006, pp. 233–240, <http://dx.doi.org/10.1145/1143844.1143874>, <http://doi.acm.org/10.1145/1143844.1143874>.
- [164] T. Saito, M. Rehmsmeier, The precision-recall plot is more informative than the ROC plot when evaluating binary classifiers on imbalanced datasets, *PLOS ONE* 10 (2015) 1–21, <http://dx.doi.org/10.1371/journal.pone.0118432>.
- [165] W. Zachary, An information flow model for conflict and fission in small groups, *J. Anthropol. Res.* 33 (1977) 452–473.
- [166] D. Lusseau, K. Schneider, O.J. Boisseau, P. Haase, E. Slooten, S.M. Dawson, The bottlenose dolphin community of doubtful sound features a large proportion of long-lasting associations, *Behav. Ecol. Sociobiol.* 54 (4) (2003) 396–405, <http://dx.doi.org/10.1007/s00265-003-0651-y>.
- [167] N.T. Markov, M.M. Ercsey-Ravasz, A.R. Ribeiro Gomes, C. Lamy, L. Magrou, J. Vezoli, P. Misery, A. Falchier, R. Quilodran, M.A. Gariel, J. Sallet, R. Gamanut, C. Huisoud, S. Clavagnier, P. Giroud, D. Sappey-Marinié, P. Barone, C. Dehay, Z. Toroczkai, K. Knoblauch, D.C. Van Essen, H. Kennedy, A weighted and directed interareal connectivity matrix for macaque cerebral cortex, *Cerebral Cortex* 24 (1) (2014) 17–36, <http://dx.doi.org/10.1093/cercor/bhs270>, [arXiv:oup/backfile/content_public/journal/cercor/24/1/10.1093_cercor_bhs270/1/bhs270.pdf](https://arxiv.org/abs/oup/backfile/content_public/journal/cercor/24/1/10.1093_cercor_bhs270/1/bhs270.pdf).
- [168] M. Girvan, M.E.J. Newman, Community structure in social and biological networks, 99 (12) (2002) 7821–7826, <http://dx.doi.org/10.1073/pnas.122653799>.
- [169] P.M. Gleiser, L. Danon, Community structure in jazz, *Adv. Complex Syst.* 6 (4) (2003) 565–574, <http://dx.doi.org/10.1142/S0219525903001067>.
- [170] M.E.J. Newman, Finding community structure in networks using the eigenvectors of matrices, *Phys. Rev. E* 74 (2006) 036104, <http://dx.doi.org/10.1103/PhysRevE.74.036104>.
- [171] R.N. Lichtenwalter, J.T. Lussier, N.V. Chawla, New perspectives and methods in link prediction, in: *Proceedings of the 16th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, in: *KDD '10*, ACM, New York, NY, USA, 2010, pp. 243–252, <http://dx.doi.org/10.1145/1835804.1835837>, <http://doi.acm.org/10.1145/1835804.1835837>.
- [172] A. Mantrach, L. Yen, J. Callut, K. Francoise, M. Shimbo, M. Saerens, The sum-over-paths covariance kernel: a novel covariance measure between nodes of a directed graph, *IEEE Trans. Pattern Anal. Mach. Intell.* 32 (6) (2010) 1112–1126, <http://dx.doi.org/10.1109/TPAMI.2009.78>.
- [173] J. Kunegis, J. Fliege, Predicting directed links using nondiagonal matrix decompositions, in: *2012 IEEE 12th International Conference on Data Mining*, 2012, pp. 948–953, <http://dx.doi.org/10.1109/ICDM.2012.16>.
- [174] Q.-M. Zhang, L. Lü, W.-Q. Wang, Yu-Xiao, T. Zhou, Potential theory for directed networks, *PLOS ONE* 8 (2) (2013) 1–8, <http://dx.doi.org/10.1371/journal.pone.0055437>.

- [175] X. Zhang, C. Zhao, X. Wang, D. Yi, Identifying missing and spurious interactions in directed networks, *Int. J. Distrib. Sen. Netw.* 2015 (2015) 27:27, <http://dx.doi.org/10.1155/2015/507386>.
- [176] E. Bütün, M. Kaya, R. Alhaji, A new topological metric for link prediction in directed, weighted and temporal networks, in: 2016 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining, ASONAM, 2016, pp. 954–959.
- [177] B. Chen, Y. Hua, Y. Yuan, Y. Jin, Link prediction on directed networks based on AUC optimization, *IEEE Access* 6 (2018) 28122–28136, <http://dx.doi.org/10.1109/ACCESS.2018.2838259>.
- [178] E. Bütün, M. Kaya, Predicting citation count of scientists as a link prediction problem, *IEEE Trans. Cybern.* (2019) 1–12, <http://dx.doi.org/10.1109/TCYB.2019.2900495>.
- [179] M. Lu, X. Wei, D. Ye, Y. Dai, A unified link prediction framework for predicting arbitrary relations in heterogeneous academic networks, *IEEE Access* 7 (2019) 124967–124987, <http://dx.doi.org/10.1109/ACCESS.2019.2939172>.
- [180] E. Bütün, M. Kaya, A pattern based supervised link prediction in directed complex networks, *Physica A* 525 (2019) 1136–1145, <http://dx.doi.org/10.1016/j.physa.2019.04.015>, <http://www.sciencedirect.com/science/article/pii/S0378437119303796>.
- [181] G. Salha, S. Limnios, R. Hennequin, V. Tran, M. Vazirgiannis, Gravity-inspired graph autoencoders for directed link prediction, *CoRR abs/1905.09570* (2019) <http://arxiv.org/abs/1905.09570>, arXiv:1905.09570.
- [182] P. Sarkar, D. Chakrabarti, M.I. Jordan, Nonparametric link prediction in dynamic networks, in: Proceedings of the 29th International Conference on Machine Learning, in: ICML'12, Omnipress, USA, 2012, pp. 1897–1904, <http://dl.acm.org/citation.cfm?id=3042573.3042815>.
- [183] S. Gao, L. Denoyer, P. Gallinari, Temporal link prediction by integrating content and structure information, in: Proceedings of the 20th ACM International Conference on Information and Knowledge Management, in: CIKM '11, ACM, New York, NY, USA, 2011, pp. 1169–1174, <http://dx.doi.org/10.1145/2063576.2063744>, <http://doi.acm.org/10.1145/2063576.2063744>.
- [184] D.Q. Vu, A.U. Asuncion, D.R. Hunter, P. Smyth, Continuous-time regression models for longitudinal networks, in: Proceedings of the 24th International Conference on Neural Information Processing Systems, in: NIPS'11, Curran Associates Inc., USA, 2011, pp. 2492–2500, <http://dl.acm.org/citation.cfm?id=2986459.2986737>.
- [185] M. Pujari, R. Kanawati, Supervised rank aggregation approach for link prediction in complex networks, in: Proceedings of the 21st International Conference on World Wide Web, in: WWW '12 Companion, ACM, New York, NY, USA, 2012, pp. 1189–1196, <http://dx.doi.org/10.1145/2187980.2188260>, <http://doi.acm.org/10.1145/2187980.2188260>.
- [186] Z. Zeng, K. Chen, S. Zhang, H. Zhang, A link prediction approach using semi-supervised learning in dynamic networks, in: 2013 Sixth International Conference on Advanced Computational Intelligence, ICACI, 2013, pp. 276–280, <http://dx.doi.org/10.1109/ICACI.2013.6748516>.
- [187] Y.-l. He, J.N. Liu, Y.-x. Hu, X.-z. Wang, Owa operator based link prediction ensemble for social network, *Expert Syst. Appl.* 42 (1) (2015) 21–50, <http://dx.doi.org/10.1016/j.eswa.2014.07.018>.
- [188] Z. Bao, Y. Zeng, Y.C. Tay, Sonlp: social network link prediction by principal component regression, in: Proceedings of the 2013 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining, in: ASONAM '13, ACM, New York, NY, USA, 2013, pp. 364–371, <http://dx.doi.org/10.1145/2492517.2492558>, <http://doi.acm.org/10.1145/2492517.2492558>.
- [189] J. O'Madadhain, J. Hutchins, P. Smyth, Prediction and ranking algorithms for event-based network data, *SIGKDD Explor. Newsl.* 7 (2) (2005) 23–30, <http://dx.doi.org/10.1145/1117454.1117458>, <http://doi.acm.org/10.1145/1117454.1117458>.
- [190] B. Bringmann, M. Berlingerio, F. Bonchi, A. Gionis, Learning and predicting the evolution of social networks, *IEEE Intell. Syst.* 25 (4) (2010) 26–35, <http://dx.doi.org/10.1109/MIS.2010.91>.
- [191] C.A. Bliss, M.R. Frank, C.M. Danforth, P.S. Dodds, An evolutionary algorithm approach to link prediction in dynamic social networks, *J. Comput. Sci.* 5 (5) (2014) 750–764, <http://dx.doi.org/10.1016/j.jocs.2014.01.003>, <http://www.sciencedirect.com/science/article/pii/S1877750314000040>.
- [192] Structural link prediction based on ant colony approach in social networks, *Physica A* 419 (2015) 80–94, <http://dx.doi.org/10.1016/j.physa.2014.10.011>, <http://www.sciencedirect.com/science/article/pii/S0378437114008498>.
- [193] F. Hu, H. San Wong, Labeling of human motion based on CBGA and probabilistic model, *Int. J. Smart Sens. Intell. Syst.* 6 (2013) 583–609, <http://dx.doi.org/10.21307/ijssis-2017-556>.
- [194] N. Barbieri, F. Bonchi, G. Manco, Who to follow and why: link prediction with explanations, in: Proceedings of the 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, in: KDD '14, ACM, New York, NY, USA, 2014, pp. 1266–1275, <http://dx.doi.org/10.1145/2623330.2623733>, <http://doi.acm.org/10.1145/2623330.2623733>.
- [195] J. Liu, G. Deng, Link prediction in a user-object network based on time-weighted resource allocation, *Physica A* 388 (17) (2009) 3643–3650, <http://dx.doi.org/10.1016/j.physa.2009.05.021>, <http://www.sciencedirect.com/science/article/pii/S0378437109004099>.
- [196] S. Hanneke, W. Fu, E.P. Xing, Discrete temporal models of social networks, *Electron. J. Stat.* 4 (2010) 585–605, <http://dx.doi.org/10.1214/09-EJS548>.
- [197] P. Klimek, A.S. Jovanovic, R. Eglöf, R. Schneider, Successful fish go with the flow: citation impact prediction based on centrality measures for term-document networks, *Scientometrics* 107 (2016) 1265–1282.
- [198] Y. Li, A. Wen, Q. Lin, R. Li, Z. Lu, Name disambiguation in scientific cooperation network by exploiting user feedback, *Artif. Intell. Rev.* 41 (4) (2014) 563–578, <http://dx.doi.org/10.1007/s10462-012-9323-5>.
- [199] M. Ge, A. Li, M. Wang, A bipartite network-based method for prediction of long non-coding RNA–protein interactions, *Genom. Proteom. Bioinf.* 14 (1) (2016) 62–71, <http://dx.doi.org/10.1016/j.gpb.2016.01.004>.
- [200] J. Kunegis, E.W. De Luca, S. Albayrak, The link prediction problem in bipartite networks, in: Proceedings of the Computational Intelligence for Knowledge-based Systems Design, and 13th International Conference on Information Processing and Management of Uncertainty, in: IPMU'10, Springer-Verlag, Berlin, Heidelberg, 2010, pp. 380–389, <http://dl.acm.org/citation.cfm?id=1876326.1876373>.
- [201] S. Xia, B.T. Dai, E. Lim, Y. Zhang, C. Xing, Link prediction for bipartite social networks: the role of structural holes, in: International Conference on Advances in Social Networks Analysis and Mining, ASONAM 2012, Istanbul, Turkey, 26–29 August 2012, 2012, pp. 153–157, <http://dx.doi.org/10.1109/ASONAM.2012.35>.
- [202] Y. Chang, H. Kao, Link prediction in a bipartite network using wikipedia revision information, in: 2012 Conference on Technologies and Applications of Artificial Intelligence, 2012, pp. 50–55, <http://dx.doi.org/10.1109/TAAL.2012.49>.
- [203] S. Daminelli, J.M. Thomas, C. Durán, C.V. Cannistraci, Common neighbours and the local-community-paradigm for topological link prediction in bipartite networks, *CoRR abs/1504.07011* (2015) <http://arxiv.org/abs/1504.07011>, arXiv:1504.07011.
- [204] J. Rissanen, Paper: modeling by shortest data description, *Automatica* 14 (5) (1978) 465–471, [http://dx.doi.org/10.1016/0005-1098\(78\)90005-5](http://dx.doi.org/10.1016/0005-1098(78)90005-5).
- [205] D. Chakrabarti, S. Papadimitriou, D.S. Modha, C. Faloutsos, Fully automatic cross-associations, in: Proceedings of the Tenth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, in: KDD '04, ACM, New York, NY, USA, 2004, pp. 79–88, <http://dx.doi.org/10.1145/1014052.1014064>, <http://doi.acm.org/10.1145/1014052.1014064>.
- [206] M. Baltakiene, K. Baltakys, D. Cardamone, F. Parisi, T. Radicioni, M. Torricelli, J.A. van Lidth de Jeude, F. Saracco, Maximum entropy approach to link prediction in bipartite networks, *CoRR abs/1805.04307* (2018) <http://arxiv.org/abs/1805.04307>, arXiv:1805.04307.
- [207] F. Saracco, R. di Clemente, A. Gabrielli, T. Squartini, Randomizing bipartite networks: the case of the world trade web, *Sci. Rep.* 5 (2015) 10595, <http://dx.doi.org/10.1038/srep10595>.

- [208] O. Allali, C. Magnien, M. Latapy, Link prediction in bipartite graphs using internal links and weighted projection, in: Third International Workshop on Network Science for Communication Networks (Netscom 2011), In Conjunction with IEEE Infocom 2011., IEEE, 2011, pp. 936–941, <http://dx.doi.org/10.1109/INFCOMW.2011.5928947>, <https://hal.archives-ouvertes.fr/hal-01286948>.
- [209] Y. Sun, Y. Yu, J. Han, Ranking-based clustering of heterogeneous information networks with star network schema, in: Proceedings of the 15th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, in: KDD '09, ACM, New York, NY, USA, 2009, pp. 797–806, <http://dx.doi.org/10.1145/1557019.1557107>, <http://doi.acm.org/10.1145/1557019.1557107>.
- [210] Y. Sun, J. Han, P. Zhao, Z. Yin, H. Cheng, T. Wu, Rankclus: integrating clustering with ranking for heterogeneous information network analysis, in: Proceedings of the 12th International Conference on Extending Database Technology: Advances in Database Technology, in: EDBT '09, ACM, New York, NY, USA, 2009, pp. 565–576, <http://dx.doi.org/10.1145/1516360.1516426>, <http://doi.acm.org/10.1145/1516360.1516426>.
- [211] Y. Sun, J. Han, X. Yan, P.S. Yu, T. Wu, Pathsim: meta path-based top-k similarity search in heterogeneous information networks, *PVLDB* 4 (11) (2011) 992–1003, <http://www.vldb.org/pvldb/vol4/p992-sun.pdf>.
- [212] Y. Yang, N. Chawla, Y. Sun, J. Han, Predicting links in multi-relational and heterogeneous networks, in: Proceedings of the 2012 IEEE 12th International Conference on Data Mining, in: ICDM '12, IEEE Computer Society, Washington, DC, USA, 2012, pp. 755–764, <http://dx.doi.org/10.1109/ICDM.2012.144>.
- [213] D. Davis, R. Lichtenwalter, N.V. Chawla, Multi-relational link prediction in heterogeneous information networks, in: Proceedings of the 2011 International Conference on Advances in Social Networks Analysis and Mining, in: ASONAM '11, IEEE Computer Society, Washington, DC, USA, 2011, pp. 281–288, <http://dx.doi.org/10.1109/ASONAM.2011.107>.
- [214] Y. Sun, R. Barber, M. Gupta, C.C. Aggarwal, J. Han, Co-author relationship prediction in heterogeneous bibliographic networks, in: Proceedings of the 2011 International Conference on Advances in Social Networks Analysis and Mining, in: ASONAM '11, IEEE Computer Society, Washington, DC, USA, 2011, pp. 121–128, <http://dx.doi.org/10.1109/ASONAM.2011.112>.
- [215] Y. Sun, J. Han, C.C. Aggarwal, N.V. Chawla, When will it happen?: relationship prediction in heterogeneous information networks, in: Proceedings of the Fifth ACM International Conference on Web Search and Data Mining, in: WSDM '12, ACM, New York, NY, USA, 2012, pp. 663–672, <http://dx.doi.org/10.1145/2124295.2124373>, <http://doi.acm.org/10.1145/2124295.2124373>.
- [216] Y. Dong, J. Tang, S. Wu, J. Tian, N.V. Chawla, J. Rao, H. Cao, Link prediction and recommendation across heterogeneous social networks, in: 2012 IEEE 12th International Conference on Data Mining, 2012, pp. 181–190, <http://dx.doi.org/10.1109/ICDM.2012.140>.
- [217] B. Cao, N.N. Liu, Q. Yang, Transfer learning for collective link prediction in multiple heterogeneous domains, in: Proceedings of the 27th International Conference on International Conference on Machine Learning, in: ICML'10, Omnipress, USA, 2010, pp. 159–166, <http://dl.acm.org/citation.cfm?id=3104322.3104344>.
- [218] S. Negi, S. Chaudhry, Link prediction in heterogeneous social networks, in: Proceedings of the 25th ACM International on Conference on Information and Knowledge Management, in: CIKM '16, ACM, New York, NY, USA, 2016, pp. 609–617, <http://dx.doi.org/10.1145/2983323.2983722>, <http://doi.acm.org/10.1145/2983323.2983722>.
- [219] I. Esslimani, A. Brun, A. Boyer, Densifying a behavioral recommender system by social networks link prediction methods, *Soc. Netw. Anal. Min.* 1 (3) (2011) 159–172, <http://dx.doi.org/10.1007/s13278-010-0004-6>.
- [220] Y. Hu, Y. Koren, C. Volinsky, Collaborative filtering for implicit feedback datasets, in: 2008 Eighth IEEE International Conference on Data Mining, 2008, pp. 263–272, <http://dx.doi.org/10.1109/ICDM.2008.22>.
- [221] B. Allison, D. Guthrie, L. Guthrie, Another look at the data sparsity problem, in: P. Sojka, I. Kopeček, K. Pala (Eds.), *Text, Speech and Dialogue*, Springer Berlin Heidelberg, Berlin, Heidelberg, 2006, pp. 327–334.
- [222] X. Li, H. Chen, Recommendation as link prediction: a graph kernel-based machine learning approach, in: Proceedings of the 9th ACM/IEEE-CS Joint Conference on Digital Libraries, in: JCDL '09, ACM, New York, NY, USA, 2009, pp. 213–216, <http://dx.doi.org/10.1145/1555400.1555433>, <http://doi.acm.org/10.1145/1555400.1555433>.
- [223] J. Li, L. Zhang, F. Meng, F. Li, Recommendation algorithm based on link prediction and domain knowledge in retail transactions, in: 2nd International Conference on Information Technology and Quantitative Management, ITQM 2014, *Procedia Comput. Sci.* 31 (2014) 875–881, <http://dx.doi.org/10.1016/j.procs.2014.05.339>, <http://www.sciencedirect.com/science/article/pii/S187705091400516X>.
- [224] A. Sadilek, H. Kautz, J.P. Bigham, Finding your friends and following them to where you are, in: Proceedings of the Fifth ACM International Conference on Web Search and Data Mining, in: WSDM '12, ACM, New York, NY, USA, 2012, pp. 723–732, <http://dx.doi.org/10.1145/2124295.2124380>, <http://doi.acm.org/10.1145/2124295.2124380>.
- [225] X. Li, H. Chen, Recommendation as link prediction in bipartite graphs: a graph kernel-based machine learning approach, *Decis. Support Syst.* 54 (2) (2013) 880–890, <http://dx.doi.org/10.1016/j.dss.2012.09.019>.
- [226] X. Li, H. Chen, Recommendation as link prediction in bipartite graphs, *Decis. Support Syst.* 54 (2) (2013) 880–890, <http://dx.doi.org/10.1016/j.dss.2012.09.019>.
- [227] J. Tang, S. Wu, J. Sun, H. Su, Cross-domain collaboration recommendation, in: Proceedings of the 18th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, in: KDD '12, ACM, New York, NY, USA, 2012, pp. 1285–1293, <http://dx.doi.org/10.1145/2339530.2339730>, <http://doi.acm.org/10.1145/2339530.2339730>.
- [228] F. Masrour, I. Barjesteh, R. Forsati, A. Esfahanian, H. Radha, Network completion with node similarity: A matrix completion approach with provable guarantees, in: 2015 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining (ASONAM), 2015, pp. 302–307, <http://dx.doi.org/10.1145/2808797.2809407>.
- [229] M. Kim, J. Leskovec, The network completion problem: inferring missing nodes and edges in networks, in: Proceedings of the Eleventh SIAM International Conference on Data Mining, SDM 2011, April 28–30, 2011, Mesa, Arizona, USA, 2011, pp. 47–58, <https://doi.org/10.1137/1.9781611972818.5>.
- [230] E.J. Candès, X. Li, Y. Ma, J. Wright, Robust principal component analysis?, *J. ACM* 58 (3) (2011) 11:1–11:37, <http://dx.doi.org/10.1145/1970392.1970395>, <http://doi.acm.org/10.1145/1970392.1970395>.
- [231] Z. Huang, D. Zeng, A link prediction approach to anomalous email detection, in: Conference Proceedings - IEEE International Conference on Systems, Man and Cybernetics, vol. 2, 2007, pp. 1131–1136, <http://dx.doi.org/10.1109/ICSMC.2006.384552>.
- [232] S.-d. Lin, H. Chalupsky, Unsupervised link discovery in multi-relational data via rarity analysis, in: Proceedings of the Third IEEE International Conference on Data Mining, in: ICDM '03, IEEE Computer Society, Washington, DC, USA, 2003, pp. 171–178, <http://dl.acm.org/citation.cfm?id=951949.952188>.
- [233] M.J. Rattigan, D. Jensen, The case for anomalous link detection, in: Proceedings of the 4th International Workshop on Multi-relational Mining, in: MRDM '05, ACM, New York, NY, USA, 2005, pp. 69–74, <http://dx.doi.org/10.1145/1090193.1090205>, <http://doi.acm.org/10.1145/1090193.1090205>.
- [234] S. Al-Oufi, H.-N. Kim, A. El Saddik, Controlling privacy with trust-aware link prediction in online social networks, in: Proceedings of the Third International Conference on Internet Multimedia Computing and Service, in: ICIMCS '11, ACM, New York, NY, USA, 2011, pp. 86–89, <http://dx.doi.org/10.1145/2043674.2043699>, <http://doi.acm.org/10.1145/2043674.2043699>.
- [235] R. Levien, Attack-resistant trust metrics, in: J. Golbeck (Ed.), *Computing with Social Trust*, Springer London, London, 2009, pp. 121–132, http://dx.doi.org/10.1007/978-1-84800-356-9_5.

- [236] P. Massa, K. Souren, Trustlet, open research on trust metrics, in: D. Flejter, S. Grzonkowski, T. Kaczmarek, M. Kowalkiewicz, T. Nagle, J. Parkes (Eds.) BIS 2008 Workshop Proceedings, 2008, pp. 31–43, <http://sunsite.informatik.rwth-aachen.de/Publications/CEUR-WV/Vol-333/saw3.pdf>.
- [237] M. Kc, R. Chau, M. Hagenbuchner, A.C. Tsoi, V. Lee, A machine learning approach to link prediction for interlinked documents, in: Proceedings of the Focused Retrieval and Evaluation, and 8th International Conference on Initiative for the Evaluation of XML Retrieval, in: INEX'09, Springer-Verlag, Berlin, Heidelberg, 2010, pp. 342–354, <http://dl.acm.org/citation.cfm?id=1881065.1881106>.
- [238] Y. Zhao, S. Li, J. Hou, Link quality prediction via a neighborhood-based nonnegative matrix factorization model for wireless sensor networks, *Int. J. Distrib. Sens. Netw.* 11 (10) (2015) 828493, <http://dx.doi.org/10.1155/2015/828493>, arXiv:<https://doi.org/10.1155/2015/828493>.
- [239] E. Weiss, K. Kurowski, S. Hischke, B. Xu, Avoiding route breakage in ad hoc networks using link prediction, in: Proceedings of the Eighth IEEE Symposium on Computers and Communications, vol. 1, pp. 57–62 <http://dx.doi.org/10.1109/ISCC.2003.1214101>.
- [240] A. Yadav, Y.N. Singh, R.R. Singh, Improving routing performance in AODV with link prediction in mobile adhoc networks, *Wirel. Pers. Commun.* 83 (1) (2015) 603–618, <http://dx.doi.org/10.1007/s11277-015-2411-5>.
- [241] C. Hu, J.C. Hou, A link-indexed statistical traffic prediction approach to improving ieee 802.11 psm, *Ad Hoc Netw.* 3 (5) (2005) 529–545, <http://dx.doi.org/10.1016/j.adhoc.2004.08.003>.
- [242] Q. Han, Y. Bai, L. Gong, W. Wu, Link availability prediction-based reliable routing for mobile ad hoc networks, *IET Commun.* 5 (16) (2011) 2291–2300, <http://dx.doi.org/10.1049/jiet-com.2010.0946>.
- [243] J. Chen, Y. Han, D. Li, J. Nie, Link prediction and route selection based on channel state detection in uasns, *Int. J. Distrib. Sens. Netw.* 7 (1) (2011) 939864, <http://dx.doi.org/10.1155/2011/939864>, arXiv:<https://doi.org/10.1155/2011/939864>.
- [244] Z. Huo, X. Huang, X. Hu, Link prediction with personalized social influence, in: Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence, (AAAI-18), the 30th innovative Applications of Artificial Intelligence (IAAI-18), and the 8th AAAI Symposium on Educational Advances in Artificial Intelligence (EAAI-18), New Orleans, Louisiana, USA, February 2–7, 2018, 2018, pp. 2289–2296, <https://www.aaai.org/ocs/index.php/AAAI/AAAI18/paper/view/16971>.
- [245] G. D'Angelo, L. Severini, Y. Velaj, Recommending links through influence maximization, *CoRR abs/1706.04368* (2017) <http://arxiv.org/abs/1706.04368>, arXiv:1706.04368.
- [246] E.M. Rogers, Diffusion of innovations, *J. Pharm. Sci.* 52 (6) (1963) 612, <http://dx.doi.org/10.1002/jps.2600520633>, <http://www.sciencedirect.com/science/article/pii/S0022354915340132>.
- [247] S.S. Singh, K. Singh, A. Kumar, H.K. Shakya, B. Biswas, A survey on information diffusion models in social networks, in: A.K. Luhach, D. Singh, P.-A. Hsiung, K.B.G. Hawari, P. Lingras, P.K. Singh (Eds.), *Advanced Informatics for Computing Research*, Springer Singapore, Singapore, 2019, pp. 426–439.
- [248] S.S. Singh, K. Singh, A. Kumar, B. Biswas, Influence maximization on social networks: a study, *Recent Patents on Computer Science* 12, <http://dx.doi.org/10.2174/2213275912666190417152547>.
- [249] D. Kempe, J. Kleinberg, E. Tardos, Maximizing the spread of influence through a social network, in: Proceedings of the Ninth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, in: KDD '03, ACM, New York, NY, USA, 2003, pp. 137–146, <http://dx.doi.org/10.1145/956750.956769>.
- [250] S.S. Singh, A. Kumar, K. Singh, B. Biswas, C2im: community based context-aware influence maximization in social networks, *Physica A* 514 (2019) 796–818, <http://dx.doi.org/10.1016/j.physa.2018.09.142>, <http://www.sciencedirect.com/science/article/pii/S0378437118312822>.
- [251] S.S. Singh, A. Kumar, K. Singh, B. Biswas, LAPSO-IM: a learning-based influence maximization approach for social networks, *Appl. Soft Comput.* (2019) 105554, <http://dx.doi.org/10.1016/j.asoc.2019.105554>.
- [252] S.S. Singh, K. Singh, A. Kumar, B. Biswas, Mim2: multiple influence maximization across multiple social networks, *Physica A* 526 (2019) 120902, <http://dx.doi.org/10.1016/j.physa.2019.04.138>.
- [253] S.S. Singh, A. Kumar, K. Singh, B. Biswas, IM-SSO: Maximizing influence in social networks using social spider optimization, *Concurr. Comput.: Pract. Exper.* e5421, e5421 cpe.5421, arXiv:<https://onlinelibrary.wiley.com/doi/pdf/10.1002/cpe.5421>, <http://dx.doi.org/10.1002/cpe.5421>, <https://onlinelibrary.wiley.com/doi/abs/10.1002/cpe.5421>.
- [254] S.S. Singh, K. Singh, A. Kumar, B. Biswas, ColM: community-based influence maximization in social networks, in: A.K. Luhach, D. Singh, P.-A. Hsiung, K.B.G. Hawari, P. Lingras, P.K. Singh (Eds.), *Advanced Informatics for Computing Research*, Springer Singapore, Singapore, 2019, pp. 440–453.
- [255] A. Biswas, B. Biswas, Community-based link prediction, *Multimedia Tools Appl.* 76 (18) (2017) 18619–18639, <http://dx.doi.org/10.1007/s11042-016-4270-9>.
- [256] A. Biswas, B. Biswas, Investigating community structure in perspective of ego network, *Expert Syst. Appl.* 42 (20) (2015) 6913–6934, <http://dx.doi.org/10.1016/j.eswa.2015.05.009>.
- [257] J. Zhang, Z. Fang, W. Chen, J. Tang, Diffusion of “following” links in microblogging networks, *IEEE Trans. Knowl. Data Eng.* 27 (8) (2015) 2093–2106, <http://dx.doi.org/10.1109/TKDE.2015.2407351>.
- [258] E. Perez-Cervantes, J.P. Mena-Chalco, M.C.F.D. Oliveira, R.M. Cesar, Using link prediction to estimate the collaborative influence of researchers, in: 2013 IEEE 9th International Conference on e-Science, 2013, pp. 293–300, <http://dx.doi.org/10.1109/eScience.2013.32>.
- [259] X. Li, N. Du, H. Li, K. Li, J. Gao, A. Zhang, A deep learning approach to link prediction in dynamic networks, in: Proceedings of the 2014 SIAM International Conference on Data Mining, Philadelphia, Pennsylvania, USA, April 24–26, 2014, 2014, pp. 289–297, <https://doi.org/10.1137/1.9781611973440.33>.
- [260] H. Wang, X. Shi, D. Yeung, Relational deep learning: a deep latent variable model for link prediction, in: Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence, February 4–9, 2017, San Francisco, California, USA, 2017, pp. 2688–2694, <http://aaai.org/ocs/index.php/AAAI/AAAI17/paper/view/14346>.
- [261] A. Dadu, A. Kumar, H.K. Shakya, S.K. Arjaria, B. Biswas, A study of link prediction using deep learning, in: A.K. Luhach, D. Singh, P.-A. Hsiung, K.B.G. Hawari, P. Lingras, P.K. Singh (Eds.), *Advanced Informatics for Computing Research*, Springer Singapore, Singapore, 2019, pp. 377–385.
- [262] X.-W. Wang, Y. Chen, Y.-Y. Liu, Link prediction through deep learning, *bioRxiv* (2018) <http://dx.doi.org/10.1101/247577>, <https://www.biorxiv.org/content/early/2018/11/27/247577>, arXiv:<https://www.biorxiv.org/content/early/2018/11/27/247577.full.pdf>.
- [263] M. Schlichtkrull, T.N. Kipf, P. Bloem, R. van den Berg, I. Titov, M. Welling, Modeling relational data with graph convolutional networks, in: A. Gangemi, R. Navigli, M.-E. Vidal, P. Hitzler, R. Troncy, L. Hollink, A. Tordai, M. Alam (Eds.), *The Semantic Web*, Springer International Publishing, Cham, 2018, pp. 593–607.
- [264] B. Yang, W. Yih, X. He, J. Gao, L. Deng, Embedding entities and relations for learning and inference in knowledge bases, *CoRR abs/1412.6575* (2014) <http://arxiv.org/abs/1412.6575>, arXiv:1412.6575.
- [265] L. Tang, H. Liu, Leveraging social media networks for classification, *Data Min. Knowl. Discov.* 23 (3) (2011) 447–478, <http://dx.doi.org/10.1007/s10618-010-0210-x>.
- [266] R. van den Berg, T.N. Kipf, M. Welling, Graph convolutional matrix completion, *CoRR abs/1706.02263* (2017) <http://arxiv.org/abs/1706.02263>, arXiv:1706.02263.
- [267] H. Wang, J. Wang, J. Wang, M. Zhao, W. Zhang, F. Zhang, X. Xie, M. Guo, Graphgan: graph representation learning with generative adversarial nets, *CoRR abs/1711.08267* (2017) <http://arxiv.org/abs/1711.08267>, arXiv:1711.08267.

- [268] A. Grover, A. Zweig, S. Ermon, Graphite: iterative generative modeling of graphs, CoRR abs/1803.10459 (2018) <http://arxiv.org/abs/1803.10459>, arXiv:1803.10459.
- [269] M. Zhang, Y. Chen, Weisfeiler-Lehman neural machine for link prediction, in: Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, Halifax, NS, Canada, August 13–17, 2017, 2017, pp. 575–583, <https://doi.org/10.1145/3097983.3097996>.
- [270] L. Zadeh, Fuzzy sets, Inf. Control 8 (3) (1965) 338–353, [http://dx.doi.org/10.1016/S0019-9958\(65\)90241-X](http://dx.doi.org/10.1016/S0019-9958(65)90241-X).
- [271] S. Bastani, A.K. Jafarabad, M.H.F. Zarandi, Fuzzy models for link prediction in social networks, Int. J. Intell. Syst. 28 (8) (2013) 768–786.
- [272] R.R. Yager, Intelligent social network analysis using granular computing, Int. J. Intell. Syst. 23 (11) (2008) 1197–1219.
- [273] S. Milgram, The small world problem, Psychol. Today 2 (1967) 60–67.
- [274] M.E.J. Newman, Finding community structure in networks using the eigenvectors of matrices, Phys. Rev. E 74 (2006) 036104, <http://dx.doi.org/10.1103/PhysRevE.74.036104>.
- [275] Y. Bhawsar, G. Thakur, Performance evaluation of link prediction techniques based on fuzzy soft set and markov model, Fuzzy Inf. Eng. 8 (1) (2016) 113–126, <http://dx.doi.org/10.1016/j.fiae.2016.03.007>, arXiv:<https://doi.org/10.1016/j.fiae.2016.03.007>.
- [276] B. Moradabadi, M.R. Meybodi, Link prediction in fuzzy social networks using distributed learning automata, Appl. Intell. 47 (3) (2017) 837–849, <http://dx.doi.org/10.1007/s10489-017-0933-0>.
- [277] B. Moradabadi, M.R. Meybodi, Link prediction in stochastic social networks: learning automata approach, J. Comput. Sci. 24 (2018) 313–328, <http://dx.doi.org/10.1016/j.jocs.2017.08.007>, <http://www.sciencedirect.com/science/article/pii/S187750317300534>.
- [278] P. Maji, R. Biswas, A. Roy, Fuzzy soft sets, J. Fuzzy Math. 9 (3) (2001) 589–602.



Ajay Kumar completed his master of technology in Computer Science & Engineering from Samrat Ashok Technological Institute Vidisha (M.P.) and Bachelor of Technology in Computer Science & Engineering from R.K.D.F Institute of Science and Technology Bhopal (M.P.). He is pursuing a Ph.D. in Computer Science and Engineering from Indian Institute of Technology (BHU), Varanasi. His research interests include Link Prediction and Influence Maximization in social/complex networks.



Shashank Sheshar Singh received M.Tech. Degree in Computer Science and Engineering from Indian Institute of Technology, Roorkee (IITR). He received B.Tech. Degree in Computer Science and Engineering from Kali Charan Nigam Institute of Technology (KCNT), Banda affiliated to GBTU University, Lucknow. He is working toward the Ph.D. in Computer Science and Engineering from Indian Institute of Technology (BHU), Varanasi. His research interests include Data Mining, Influence Maximization, Link Prediction, and Social Network Analysis.



Kuldeep Singh is currently pursuing his Ph.D. in Computer science & Engineering from Indian Institute of Technology (BHU) Varanasi. His research interest includes High utility itemsets mining, social network analysis, and data mining. He received his M.Tech degree in Computer science and Engineering from Guru Jambheshwar University of Science and Technology, Hisar (Haryana). He has 8 years of teaching and research experience. His research interests include Data Mining and Social Network Analysis.



Bhaskar Biswas received a Ph.D. in Computer Science and Engineering from Indian Institute of Technology (BHU), Varanasi. He received the B.Tech. Degree in Computer Science and Engineering from Birla Institute of Technology, Mesra. He is working as an Associate Professor at Indian Institute of Technology (BHU), Varanasi in the Computer Science and Engineering department. His research interests include Data Mining, Text Analysis, Machine Learning, Influence Maximization, Link Prediction, Social Network Analysis.