

## Fast link prediction for large networks using spectral embedding

BENJAMIN PACHEV AND BENJAMIN WEBB<sup>†</sup>

*Department of Mathematics, Brigham Young University, Provo, UT, USA*

<sup>†</sup>Corresponding author. Email: bwebb@mathematics.byu.edu

Edited by: Ernesto Estrada

[Received on 23 March 2017; editorial decision on 1 June 2017; accepted on 12 June 2017]

Many link prediction algorithms require the computation of a similarity metric on each vertex pair, which is quadratic in the number of vertices and infeasible for large networks. We develop a class of link prediction algorithms based on a spectral embedding and the  $k$  closest pairs algorithm that are scalable to very large networks. We compare the prediction accuracy and runtime of these methods to existing algorithms on several large link prediction tasks. Our methods achieve comparable accuracy to standard algorithms but are significantly faster.

**Keywords:** link prediction; graph embedding; commute time; resistance distance; closest pairs.

### 1. Introduction

The study of networks has become increasingly relevant in our understanding of the technological, natural, and social sciences. This is owing to the fact that many important systems in these areas can be described in terms of networks [1], where vertices represent the system's individual components, for example computer routers, neurons, individuals, etc. and where edges represent interactions or relationships between these components.

An essential feature of the large majority of these networks is that they have a dynamic topology, that is a structure of interactions that evolves over time [2]. The structure of social networks, for instance, change over time as relationships are formed and dissolved. In information networks such as the WWW the network's structure changes as information is created, updated, and linked.

Although understanding the mechanisms that govern this structural evolution is fundamental to network science, these mechanisms are still poorly understood. Consequently, predicting a network's eventual structure, function, or whether the network is likely to fail at some point are all currently out of reach for even simple networks.

In an attempt to determine which processes cause changes in a network's structure we are lead to the following link prediction problem: Given a network, which of the links, that is edges between existing vertices, are likely to form in the near future. Here we adopt the standard convention that links are to be predicted solely on the basis of the network's current topology (see, for instance, [3]).

Importantly, the link prediction problem can be used to study more than just which edges will appear in a network. It can also be used to predict which of the non-network edges are, in fact, in the network but currently undetected. Similarly, it can be used to detect which of the current network edges have been falsely determined to be a part of the network.

This notion of link prediction is of central importance in numerous applications. Companies such as Facebook, Twitter and Google need to know the current state and efficiently predict the future structure

of the networks they use to accurately sort and organize data [4]. Biologists need to know whether biochemical reactions are caused by specific sets of enzymes to infer causality and so on [5].

The barrier in determining whether network links truly exist in these and other settings, is that testing and discovering interactions in a network requires significant experimental effort in the laboratory or in the field [6]. Similarly, determining experimentally when and where a new link will form may also be impractical, especially if the precise mechanism for link formation is unknown. For these reasons it is important to develop models for link prediction.

At present, there is an ever increasing number of proposed methods for predicting network links [7]. Not surprisingly, certain methods more accurately predict the formation of links in certain networks when compared with others. Additionally, each of these methods has a runtime that scales differently with the size of the network. In our experiments, we discover that a number of link predictors have a runtime that is so high that it effectively prohibits their use on moderately large networks.

Here we propose a class of link predicting algorithms that scale to large networks. This method, which we refer to as the *approximate resistance distance predictor*, integrates a spectral embedding of the network with a known algorithm for efficiently finding the  $k$  closest pairs of points in Euclidean space. The spectral embedding aspect of the algorithm is derived as a low-rank approximation of the effective resistance between network vertices, as in [8]. The  $k$  closest pairs component of the algorithm is taken from [9] and can be used to predict links based on this embedding.

Here we compare the prediction accuracy and runtime of this method against several well-known algorithms on a number of coauthorship networks and a social network consisting of a small subset of Facebook users. We find that our method achieves the best accuracy on some networks and scales to networks that many other link predictors cannot.

The article is structured as follows. We begin with a review of related work. In Section 3, we describe the link prediction problem and outline a number of standard link prediction algorithms. In Section 4, we introduce the method of resistance distance embedding and prove that it is optimal as a low rank approximation of effective resistance (see Proposition 4.1). In Section 5, we describe the experimental setup. Section 6 numerical results comparing the performance of the resistance distance embedding algorithm to other algorithms are given. Section 7 concludes with some closing remarks including a number of open questions for future work.

## 2. Related work

There are a vast number of proposed link prediction algorithms. A good survey is provided by Wang *et al.* [10]. We discuss a few representative works. Liben-Nowell and Kleinberg [3] studied link prediction in social networks, in which they examine a number of unsupervised methods based on network topology. Our experimental methodology and formulation of link prediction are patterned after theirs. Lichenwalter *et al.* [11] found that link prediction with supervised learning significantly outperforms unsupervised methods on certain networks. In the same work they also found that supervised methods for link prediction benefit from using the output of unsupervised methods as features. The implication is that developing good unsupervised link prediction algorithms can help improve the state-of-the-art of supervised methods for link prediction. In this work, our goal is to provide a scalable, performant class of unsupervised methods for link predictions.

Additionally, Fouss *et al.* [8] examined a low-rank approximation of the resistance distance as a link predictor. They found it performed well on several large collaborative recommendation tasks. Our contribution is an algorithm to scale this link predictor, and other link predictors of a similar form, to

large networks. This includes, for instance, the *amplified commute distance* introduced by von Luxburg *et al.* [12] (see Section 6).

An important special case of link prediction is triangle closure prediction. Leskovec *et al.* [13] examined a variety of local methods for triangle closure prediction in the context of a complete model of network growth. Estrada and Arrigo [14] proposed a model for triangle closure that is based on communicability distance and incorporates global network information. They validated this model on several large real-world networks.

### 3. The link prediction problem

The link prediction problem can be stated as follows. Given a connected graph  $G = (V, E)$ , and  $k$ , the number of predicted non-adjacent links, we seek  $k$  pairs of vertices which are most likely to become connected. While the choice of  $k$  depends on the application, we adopt the convention that  $1 \leq k \leq |E|$ .

The general paradigm for link prediction is to compute a similarity metric  $score(x, y)$  on each vertex pair  $(x, y)$ . The predicted links are then the  $k$   $(x, y) \in V \times V - E$  for which  $score(x, y)$  is maximal. By constructing a matrix from the scores, we obtain a *graph kernel*. We can also go in the other direction. Any real  $n \times n$  matrix, where  $n = |V|$ , defines a score function on pairs of vertices, and can be used for link prediction.

We now give a sampling of existing link prediction algorithms.

#### 3.1 Local methods

A *local method* for link prediction is an algorithm that uses vertex neighbourhoods to compute similarity.

**Common neighbours:** Common neighbours simply assign

$$score(x, y) = |\Gamma(x) \cap \Gamma(y)|, \quad (3.1)$$

where  $\Gamma(x)$  is the neighbour set for  $x \in V$ .

**Jaccard's coefficient:** Jaccard's coefficient is a normalized version of common neighbours that takes into account the total number of neighbours for both vertices. It is given by

$$score(x, y) = \frac{|\Gamma(x) \cap \Gamma(y)|}{|\Gamma(x) \cup \Gamma(y)|}. \quad (3.2)$$

**Preferential attachment:** Preferential attachment is based on the idea that highly connected nodes are more likely to form links, an observed pattern in coauthorship networks [15]. This leads to

$$score(x, y) = |\Gamma(x)| |\Gamma(y)|. \quad (3.3)$$

**Adamic-adar:**

$$score(x, y) = \sum_{z \in \Gamma(x) \cap \Gamma(y)} \frac{1}{\log |\Gamma(z)|}. \quad (3.4)$$

**Resource allocation:**

$$score(x, y) = \sum_{z \in \Gamma(x) \cap \Gamma(y)} \frac{1}{|\Gamma(z)|}. \quad (3.5)$$

### 3.2 Path-based methods

*Path-based methods* consider all or a subset of the paths between two vertices to compute similarity. Unlike local similarity measures, they can capture global information about the network.

**Shortest path:** This link predictor defines  $score(x, y)$  as the negated length of the shortest path from  $x$  to  $y$ .

**Katz:** The Katz metric counts all paths between two nodes, and discounts the longer paths exponentially. Define  $path_{x,y}^\ell$  to be the set of all paths of length  $\ell$  from  $x$  to  $y$ . Then given a weight  $0 < \beta < 1$ ,

$$score(x, y) = \sum_{\ell=1}^{\infty} \beta^\ell |path_{x,y}^\ell|. \quad (3.6)$$

A closed form for the associated graph kernel is given by  $(I - \beta A)^{-1} - I = \sum_{\ell=1}^{\infty} (\beta A)^\ell$ , where  $A$  is the adjacency matrix of  $G$ .

### 3.3 Random walks

A *random walk* on  $G$  starts at some node  $x$  and iteratively moves to new nodes with uniform probability. There are a multitude of link predictors based on random walks. These are some of the fundamental ones.

**Hitting and commute time:** The *hitting time*  $H_{x,y}$  is the expected number of steps required to reach  $y$  in a random walk starting at  $x$ . Commute time is defined as  $C_{x,y} = H_{x,y} + H_{y,x}$ . Negated hitting time can be used as a link predictor, but the hitting time is asymmetric in general, so we use instead the negated commute time, which is symmetric.

The commute time and its variants will be discussed further in Section 4.

**Rooted Page Rank:** A problem with hitting and commute time is that random walks can become lost exploring distant portions of the graph. Rooted Page Rank deals with this problem by introducing random resets. Given a root node  $x$ , we consider a random walk starting at  $x$ . At each step, with probability  $\alpha$  the walk returns back to  $x$ . With probability  $1 - \alpha$  the walk proceeds to a random neighbour. Given a root node  $x$ , for each other node  $y$ ,  $score(x, y)$  is defined as the stationary probability of  $y$  under the random walk rooted at  $x$ . The corresponding graph kernel is given by  $(1 - \alpha)(I - \alpha D^{-1}A)^{-1}$ , where  $D$  is the degree matrix and  $A$  is the adjacency matrix.

### 3.4 Scaling link predictors to large networks

Many link predictors, such as Katz, require the computation of a matrix inverse. This is heinously expensive for large networks, as it is cubic in the number of vertices. One way to circumvent such problems is via a low-rank approximation of the score matrix. We investigate such a low-rank approximation for the commute-time or resistance distance kernel in the next section.

Even the simpler local predictors such as common neighbours or preferential attachment face difficulties at scale. This is because for sufficiently large networks, it is not possible to compute scores for each pair of vertices and then find the maximal ones. Instead, efficient search techniques must be employed to search only a small subset of the potential links in order to find those of maximal score. In Section 4, we will demonstrate how a class of graph embedding based predictors can efficiently find the  $k$  links of maximal score.

#### 4. Spectral embedding

We begin by deriving the *approximate resistance distance link predictor* as a best low-rank approximation to commute time and show how to evaluate its link prediction scores with a spectral embedding. We then show that this link predictor is part of a family of graph embedding based link predictors that use the  $k$  closest pairs algorithm to efficiently find the links of maximal score. Finally, we discuss efficient ways to compute the spectral embedding upon which the approximate resistance distance predictor relies.

##### 4.1 Approximating commute time

Let  $L = D - A$  be the Laplacian matrix of a graph  $G = (V, E)$ , and let  $n = |V|$ . Let  $L^\dagger$  be the Moore–Penrose inverse of  $L$ . Then the commute time is given by

$$C_{x,y} = |E|(L_{x,x}^\dagger + L_{y,y}^\dagger - 2L_{x,y}^\dagger), \quad (4.1)$$

where the quantity  $r_{x,y} = (L_{x,x}^\dagger + L_{y,y}^\dagger - 2L_{x,y}^\dagger)$  is known as the *effective resistance* or the *resistance distance*. The effective resistance  $r_{x,y}$  is defined as the electrical resistance between nodes  $x$  and  $y$  where the network represents an electrical circuit and each edge a resistor. Many important properties connecting effective resistance and random walks can be found in [16], which describes the theory of electrical networks and traces the origin of this topic as far back as the nineteenth century (see references in [16]). Connections between effective resistance and the heat equation are described in [17].

More recently, it has been shown in [18] that effective resistance is, in fact, a distance on a graph, which is the reason it is also referred to as *resistance distance*. Other important properties of effective resistance are described in [19] including how network resistance can be allocated to minimize total effective resistance. Equation (4.1) states that resistance distance differs from commute-time by a (network-dependant) constant scaling factor. This property is shown in [20, 21], which allows us to use effective resistance and commute-time interchangeably for link prediction.

For many networks,  $G$  is too large to compute  $L^\dagger$  exactly, so an approximation must be used. A natural choice is a best rank- $d$  approximation to  $L^\dagger$  for some fixed dimension  $d$ . The resulting approximation of the resistance distances is closely related to distances between points in Euclidean space.

**PROPOSITION 4.1** Let  $d$  be a positive integer and let  $G = (V, E)$  be a connected, undirected graph. Then  $\exists$  a best rank- $d$  approximation  $S$  of  $L^\dagger$ , and a map  $f : V \rightarrow R^d$  so that  $\forall x, y \in V$ ,  $S_{x,x} + S_{y,y} - 2S_{x,y} = \|f(x) - f(y)\|_2^2$ . We call this map the *resistance distance embedding*.

*Proof.* For a connected graph, the Laplacian matrix is positive semidefinite, with eigenvalues  $0 = \lambda_1 < \lambda_2 \leq \dots \leq \lambda_n$  and corresponding eigenvectors  $v_1, v_2, v_3, \dots, v_n$ . Then we have the spectral decompositions

$$L = \sum_{i=2}^n \lambda_i v_i v_i^T$$

and

$$L^\dagger = \sum_{i=2}^n \frac{1}{\lambda_i} v_i v_i^T.$$

Hence,  $S = \sum_{i=2}^{d+1} \frac{1}{\lambda_i} v_i v_i^T$  is a best rank- $d$  approximation to  $L^\dagger$  in the 2-norm. Then note

$$\begin{aligned} S_{x,x} + S_{y,y} - 2S_{x,y} &= (e_x - e_y)^T S (e_x - e_y) \\ &= \sum_{i=2}^{d+1} \frac{1}{\lambda_i} (e_x - e_y)^T v_i v_i^T (e_x - e_y) \\ &= \sum_{i=2}^{d+1} \frac{1}{\lambda_i} (v_{i,x} - v_{i,y})^2 \\ &= \|f(x) - f(y)\|_2^2, \end{aligned}$$

where

$$f(x) = \left[ \frac{v_{2,x}}{\sqrt{\lambda_2}}, \frac{v_{3,x}}{\sqrt{\lambda_3}}, \dots, \frac{v_{d+1,x}}{\sqrt{\lambda_{d+1}}} \right]^T \in \mathbb{R}^d \quad (4.2)$$

□

We define the *approximate resistance distance link predictor* of dimension  $d$  by setting

$$\text{score}(x, y) = -(S_{x,x} + S_{y,y} - 2S_{x,y}) = -\|f(x) - f(y)\|_2^2, \quad (4.3)$$

where  $S$  and  $f$  are defined as in Proposition 4.1.

In the next section, we will see that the approximate resistance distance link predictor is part of a class of link predictors that avoid brute-force search when predicting links.

#### 4.2 Link prediction with graph embeddings

The resistance distance embedding is a special case of a *graph embedding*, which is a map  $f$  from  $V$  to  $\mathbb{R}^d$ ,  $d$  a positive integer. We can use graph embeddings to create link predictors. A natural choice is to set  $\text{score}(x, y) = -\|f(x) - f(y)\|_2$ , (so maximizing score corresponds to minimizing distance). We refer to this score function as the *Euclidean score*.

If  $f$  is the resistance distance embedding, then link prediction with the Euclidean score is equivalent to the approximate resistance distance predictor. Recall that the approximate resistance distance score function is  $-\|f(x) - f(y)\|_2^2$ . The  $k$  predicted links of maximal score correspond to the  $k$  non-adjacent pairs of vertices  $(x, y)$  for which  $-\|f(x) - f(y)\|_2^2$  is maximal. These are precisely the  $k$  links for which  $\|f(x) - f(y)\|_2$  is minimal and are predicted with the Euclidean score.

Link prediction with the Euclidean score is related to the *k closest pairs problem*. The closest pairs problem is as follows. Given a set of vectors  $\{x_1, x_2, \dots, x_n\} \subset \mathbb{R}^d$  we seek the  $k$  unordered pairs  $(x_i, x_j)$ ,  $i \neq j$  of minimal distance (here we use the Euclidean norm but any  $L_p$  norm can be used,  $1 \leq p \leq \infty$ ). There is an algorithm to solve this problem in

$$O\left(d \left( n \log n + k \log n \log \left( \frac{n^2}{k} \right) \right)\right) \quad (4.4)$$

[9].

We can think of the link prediction problem as the closest pairs problem applied to the set of vectors  $\{f(y), y \in V\}$ , with the additional constraint that the best pairs must correspond to non-edges in  $G$ . The extra constraint can be handled by finding the  $|E| + k$  closest pairs, then selecting the best  $k$  which are non-edges. As there can be no more than  $|E|$  edges, this approach is sure to work. We then have the worst-case complexity bound of

$$O\left(d\left(n \log n + (|E| + k) \log n \log\left(\frac{n^2}{|E| + k}\right)\right)\right). \quad (4.5)$$

Recalling that we require  $1 \leq k \leq |E|$ , and assuming that  $G$  is connected so  $|E| \geq n - 1$ , this complexity bound can be simplified to

$$O(d |E| \log^2 n). \quad (4.6)$$

For large, sparse networks,  $|E| \ll n^2$ , and this is a tremendous speedup over the  $O(n^2)$  brute-force approach.

**Cosine similarity score:** Another link prediction score function that can be derived from a graph embedding is the cosine similarity score, defined by

$$\text{score}(x, y) = \frac{\langle f(x), f(y) \rangle}{\|f(x)\| \|f(y)\|}. \quad (4.7)$$

If the cosine similarity score is used, the link prediction problem can still be solved without brute-force search. It is equivalent to the link prediction problem with Euclidean score on a modified graph embedding. The modified embedding is obtained from the original by normalizing the embedding vectors as follows.

**PROPOSITION 4.2** Given a graph embedding  $f : V \rightarrow R^d$ , the link prediction problem using

$$\text{score}(x, y) = \frac{\langle f(x), f(y) \rangle}{\|f(x)\| \|f(y)\|} = \cos \theta$$

is equivalent to the link prediction problem with the Euclidean score function on the modified embedding given by  $g(y) = \frac{f(y)}{\|f(y)\|}$ .

*Proof.* Let  $x, y \in V$ . Note

$$\langle g(x), g(y) \rangle = \cos \theta = \text{score}(x, y).$$

We have

$$\begin{aligned} \|g(x) - g(y)\|_2^2 &= \|g(x)\|_2^2 + \|g(y)\|_2^2 - 2 \langle g(x), g(y) \rangle \\ &= 2 - 2 \cos \theta = 2 - 2 \text{score}(x, y). \end{aligned}$$

This shows that minimizing Euclidean distance for the modified embedding is the same as maximizing cosine similarity score on the original, so link prediction with Euclidean score on the modified embedding is equivalent to link prediction with the cosine similarity score on the original.  $\square$

This section introduced a class of link predictors that avoid a brute-force search when predicting links. These link predictors rely on a precomputed graph embedding. The graph embedding needs to be efficiently computable in order for the overall prediction algorithm to be fast. We are concerned with link predictors that rely on the resistance distance embedding. Consequently, rapid computation of this particular graph embedding is the subject of the next section.

#### 4.3 Computing the resistance distance embedding

Computing the resistance distance embedding of dimension  $d$  requires finding the smallest  $d$  non-zero eigenvalues and associated eigenvectors of the Laplacian matrix  $L$ . Fortunately, specialized, efficient algorithms exist for this problem which exploit the positive semi-definiteness and sparsity of  $L$ . These include TRACEMIN-Fiedler [22] and a multilevel solver MC73\_FIEDLER [23]. TRACEMIN-Fiedler is simpler to implement, and is also parallelizable, so we use it in our experiments.

### 5. Experimental setup

In this section, we compare the performance of our link prediction algorithm to others on several large social networks. In a social network, nodes correspond to persons or entities. Edges correspond to an interaction between nodes, such as coauthoring an article or becoming linked on a social media website.

#### 5.1 The networks

**Arxiv High Energy Physics Theory (hep-th):** This network is a coauthorship network obtained from the Konect network collection. [24, 25].

**Arxiv High Energy Physics Phenomenology (hep-ph):** This is another coauthorship network from the Konect network collection [24, 26].

**Facebook Friendship (facebook):** This social network consists of a small subset of facebook users, where edges represent friendships [27, 28].

**Arxiv Condensed Matter Physics (cond-mat):** This dataset was obtained from Mark Newman’s website [29], and is also a coauthorship network. Unlike the other datasets, the edges are not timestamped.

#### 5.2 Creating training graphs

In order to perform link prediction, we partition edges into a training set and a test set. Edges in the training set occur before those in the test set and are used to construct a training graph. We run link prediction algorithms on the training graph to predict the contents of the test set. In most cases, edges have timestamps, and we can choose a cutoff time to partition the edges.

For one network (cond-mat) the edges are not timestamped. However, there are two versions of the cond-mat network available. One contains all collaborations up to 2003. The second is an updated network with all collaborations up to 2005. We use the first network as the training graph. The test set consists of all edges in the second network for which both nodes are in the earlier network.



TABLE 1 *Training network statistics*

Network	Nodes	Edges	Average degree
cond-mat	15,803	60,989	7.7187
cond-mat train	13,861	44,619	6.4381
facebook	63,731	817,035	12.8201
facebook train	59,416	731,929	24.6374
hep-ph	28,093	3,148,447	112.0723
hep-ph train	26,738	2,114,734	158.1819
hep-th	22,908	2,444,798	106.7225
hep-th train	21,178	1,787,157	168.7749

Choosing the cutoff between the training and test edges is somewhat arbitrary. If too few edges are used for training, link predictors will struggle. If too few are left for testing, then results may be statistically insignificant. See Table 1 for a comparison of the training networks and original networks.

Our spectral embedding based link prediction algorithms require a connected graph. To solve this problem, we reduce each training graph to its largest connected component. For each network we consider, the largest component contains the vast majority of the vertices.

### 5.3 The predictors

We perform experiments with two spectral embedding based predictors. Each uses the resistance distance embedding of dimension  $d$ , with  $d$  a parameter to be varied. The first uses the Euclidean score function and is equivalent to the approximate resistance distance predictor of dimension  $d$ . The second uses the cosine similarity score. We refer to these link predictors as `spec_euclid` and `spec_cosine`, respectively (spec for spectral). In tables, the dimension of the embedding is indicated by a number after the predictor name. For example, `spec_euclid8` refers to the `spec_euclid` predictor using an eight-dimensional resistance distance embedding.

The other link prediction algorithms used in our experiments are preferential attachment, common neighbours, Adamic Adar, Rooted Page Rank and Katz (with  $\beta = .01$ ). Some networks are too large for certain algorithms to handle, so not every algorithm is run on each network. For example, the facebook training graph has 59,416 nodes. Computing the Katz score on this graph requires finding the inverse of a  $59,416 \times 59,416$  matrix, and is very expensive in time and space, so we do not use the Katz algorithm for the facebook graph.

All experiments were performed on the same 4 core machine. The common neighbours, preferential attachment, and Adamic Adar algorithms were implemented in Python and were not parallelized. Our spectral link predictors, Katz, and Rooted Page Rank use the Python library Numpy to parallelize linear algebra operations. All code that was used in the experiments in this paper can be found at the git repository [bitbucket.org/thorfa/spectral\\_research](https://bitbucket.org/thorfa/spectral_research).

For each network, we fix the number of links to be predicted. With the exception of hep-th, this number is equal to 10% of the maximum possible number of correct predictions (i.e. the number of new links in the test set). For the hep-th network we discovered that the `spec_euclid` and `spec_cosine` predictors achieve nearly perfect accuracy when predicting 1000 links. As this is not the case for any other network we considered, we report this unusual phenomenon.

TABLE 2 *Link prediction task setup*

Network	Links	Random accuracy (%)
cond-mat	1190	0.012
facebook	7858	0.004
hep-ph	101466	0.286
reduced hep-ph	1988	0.661
hep-th	1000	0.296
reduced hep-th	135	0.084

TABLE 3 *Performance of link predictors on the cond-mat network*

Predictor	Correct (%)	Time (s)
katz	5.97	62.96
commonNeighbors	5.97	1.55
prefattach	1.93	0.35
spec_euclid1	1.51	2.99
spec_cosine1	0.25	3.35
spec_euclid2	1.51	3.65
spec_cosine2	1.18	3.80
spec_euclid4	1.76	10.54
spec_cosine4	1.34	10.89
spec_euclid8	1.68	11.41
spec_cosine8	1.34	10.73
spec_euclid16	1.68	29.91
spec_cosine16	1.43	32.31

For all of the networks we consider, the probability of randomly predicting a correct link is very low. Most of the algorithms we consider do much better than the random baseline, but have low raw accuracy since there are few new links compared to the number of possible links. See Table 2 for a summary of the number of links predicted and baseline probability of randomly predicting a correct link.

## 6. Results

On the cond-mat and facebook networks, both the spec\_euclid and spec\_cosine predictors performed worse than the simple common neighbours predictor. In addition to the full networks, we also compared predictor accuracy on reduced versions of the hep-th and hep-ph networks, because the full networks are too large for methods like Katz, common neighbours, and Rooted Page Rank to complete in a reasonable amount of time. On our reduced version of the hep-th network, our embedding-based predictors did better than common neighbours but not as well as the Rooted Page Rank predictor. On the reduced hep-ph network, the spec\_euclid predictor performed significantly better than all other competitors, including our other embedding-based predictor, spec\_cosine.

TABLE 4 *Performance of link predictors on the facebook network*

Predictor	Correct (%)	Time (s)
commonNeighbors	5.29	151.76
prefattach	0.41	7.00
spec_euclid1	0.42	9.86
spec_cosine1	0.00	47.20
spec_euclid2	0.50	11.52
spec_cosine2	0.42	12.96
spec_euclid4	1.40	24.05
spec_cosine4	1.02	25.96
spec_euclid8	1.95	61.21
spec_cosine8	2.58	62.27

TABLE 5 *Performance of link predictors on the hep-th network*

Predictor	Correct (%)	Time (s)
prefattach	0.00	9.47
spec_euclid1	94.50	10.88
spec_cosine1	1.50	17.17
spec_euclid2	98.70	17.88
spec_cosine2	99.60	15.97
spec_cosine4	100.00	15.67
spec_euclid4	99.90	18.53
spec_euclid8	100.00	29.81
spec_cosine8	100.00	29.55
spec_euclid16	99.90	96.47
spec_cosine16	99.90	100.93

As Table 3 shows, the best predictors for the cond-mat network were Katz and common neighbours. Note that for both spec\_euclid and spec\_cosine, the accuracy increases with the dimension of the embedding.

As previously mentioned, the facebook graph was too large to run the Katz predictor on it in a reasonable amount of time. As with the cond-mat network, the simple common neighbours predictor performs best (see Table 4).

Our spectral embedding link predictors performed significantly better on the hep-th and hep-ph networks, as Tables 5 and 6 show.

The common neighbours algorithm did not scale to the hep-th and hep-ph networks, unlike the facebook network. Although the facebook network had more nodes, it has a lower average node degree and fewer distance two pairs. The common neighbours algorithm computes intersections of neighbour sets for each distance two pair. Because the average node degree is higher for the hep-th and hep-ph networks, these intersections are more expensive to compute and there are more distance two pairs for which intersections must be computed.

TABLE 6 *Performance of link predictors on the hep-ph network*

Predictor	Correct (%)	Time (s)
prefattach	0.00	16.85
spec_euclid1	3.93	18.25
spec_cosine1	0.14	32.37
spec_euclid2	9.16	21.98
spec_cosine2	3.44	23.46
spec_euclid4	19.25	27.04
spec_cosine4	13.65	26.09
spec_euclid8	22.90	46.69
spec_cosine8	21.12	49.62
spec_euclid16	24.62	148.51
spec_cosine16	23.97	135.83

TABLE 7 *Performance of link predictors on the reduced hep-ph network*

Predictor	Correct (%)	Time (s)
prefattach	0.00	1.52
katz	1.16	2.75
commonNeighbors	9.36	23.96
pageRank	11.87	2.94
adamicAdar	8.85	754.40
spec_euclid1	1.81	2.80
spec_cosine1	0.10	6.21
spec_euclid2	4.73	3.73
spec_cosine2	2.57	6.86
spec_euclid4	13.13	5.80
spec_cosine4	11.12	9.01
spec_euclid8	16.40	7.39
spec_cosine8	9.31	7.44
spec_euclid16	14.13	17.32
spec_cosine16	4.93	14.90

In order to compare the performance of our spectral predictors to other predictors on the hep-ph and hep-th network data, we conducted another experiment using downsampled versions of these networks. To downsample, we used only the top 10% highest degree nodes. Our spectral predictors performed the best on the reduced hep-ph network (see Table 7), while the Rooted Page Rank algorithm was best for the reduced hep-th network (see Table 8).

It is worth mentioning that for large random geometric graphs the effective resistance  $r_{x,y}$  is approximately given by

$$r_{x,y} \approx \frac{1}{d_x} + \frac{1}{d_y}, \quad (6.1)$$

TABLE 8 *Performance of link predictors on the reduced hep-th network*

Predictor	Correct (%)	Time (s)
prefattach	0.00	1.28
katz	0.00	1.61
commonNeighbors	2.22	16.70
pageRank	11.11	1.97
adamicAdar	2.22	788.11
spec_euclid1	0.00	2.02
spec_cosine1	0.00	3.77
spec_euclid2	0.74	4.45
spec_cosine2	0.00	4.61
spec_euclid4	0.74	6.84
spec_cosine4	0.00	6.25
spec_euclid8	2.22	10.59
spec_cosine8	1.48	11.02
spec_euclid16	8.89	26.82
spec_cosine16	5.93	24.92

TABLE 9 *Comparison of the inverse degree link predictor to spectral embedding*

Network	Spectral embedding (%)	Inverse degree (%)	Common (%)
cond-mat	1.70	2.00	0.00
facebook	1.70	0.30	0.00
hep-ph	99.70	0.00	0.00
reduced hep-ph	18.10	0.00	0.00
hep-th	100.00	0.00	0.00
reduced hep-th	0.90	0.00	0.20

where  $d_x$  and  $d_y$  are the degree of node  $x$  and node  $y$ , respectively [12]. Hence, for such graphs effective resistance reduces to a local method of link prediction (see Section 3). The implication is that effective resistance, and by extension *approximate resistance distance*, may lose their effectiveness as a tool for link prediction for this important class of networks.

To test this we compare the accuracy of our method of approximate resistance distance against what we refer to as the *inverse degree* link predictor with score function given by  $score(x, y) = 1/d_x + 1/d_y$  based on equation (6.1). The results of this are shown in Table 9 where Spectral Embedding(%) gives the percent of  $k = 1000$  links predicted using the Euclidean score with  $d = 8$  for the networks described in Section 5.1. Inverse degree (%) refers to the percent correct of  $k = 1000$  links predicted by the inverse degree predictor for the same networks and Common (%) is the percent of these thousand links that were predicted by both methods.

As can be seen the inverse degree as a predictor is fairly ineffective in identifying links that will form in all but the cond-mat network when compared with the approximate resistance distance. Perhaps more importantly though, in each network the number of links predicted by both methods is essentially zero. This suggests that the methods of *approximate resistance distance* and *inverse degree* are quite different.

Reasons for this may include that the networks we consider are not random geometric graphs, possibly the size of the networks are not large enough, or the approximation given in (6.1) does not hold for our approximation of resistance distance. Although this is currently an open question, the fact that approximate resistance distance can accurately predict a number of link comparable to other well-known predictors in a way that also scales to large networks indicates the potential usefulness of this method in applications.

It is also worth mentioning that in [12] the authors suggest that the potential deficiency expressed in equation (6.1) for large random geometric graphs can be ‘corrected for’ by defining an *amplified commute distance*. This is given by

$$C_{x,y}^{amp} = \begin{cases} r_{x,y} - \frac{1}{d_x} - \frac{1}{d_y} + \frac{2a_{x,y}}{d_x d_y} + \frac{a_{x,x}}{d_x^2} + \frac{a_{y,y}}{d_y^2}, & \text{if } x \neq y \\ 0, & \text{otherwise} \end{cases}$$

where  $a_{x,y}$  is the  $x,y$ -entry of the network’s adjacency matrix  $A$ . They then show that the amplified commute distance has nice limiting properties as the size of the network tends to infinity and is a distance function on any graph.

For us the amplified commute distance is an important example since it can be shown that  $C_{x,y}^{amp} = (e_x - e_y)^T K (e_x - e_y)$  for some positive semidefinite matrix  $K$  (see page 82, [30]). Hence, it is possible to also approximate amplified commute distance via a spectral embedding in a way analogous to our approximation of effective resistance (see Proposition 4.1). However, the extent to which this approximation of amplified commute distance is useful as a link predictor is left for future work. We simply mention it as an example for how our methodology of spectral embedding can be used to potentially speed up other link prediction methods while preserving the method’s accuracy, which is the case for the approximation of effective resistance we consider in this article.

## 7. Conclusion

We present a link prediction framework that can scale to very large networks by avoiding the quadratic costs inherent in methods that exhaustively search all candidate pairs of non-adjacent nodes. We investigated the performance of a set of predictors based on this framework and the spectrum and eigenvectors of the graph’s Laplacian matrix. These methods achieved high levels of accuracy on certain real-world link prediction tasks, and scaled well to networks with tens of thousands of nodes and millions of edges.

We emphasize that there are many other possible graph embeddings to investigate. Virtually all the runtime of our spectral link predictors is spent computing the resistance distance embedding. The  $k$  closest pairs component of our algorithm is very fast in practice, with nearly linear temporal complexity in the number of edges. Replacing the resistance distance embedding with one that is cheaper to compute could potentially produce link predictors that can scale to much larger networks than the ones we consider in this article.

Our approximate resistance distance link predictor was derived as a low-rank approximation of resistance distance, an established link prediction score that is expensive to compute. Many other well-known predictors are expensive to compute, such as Katz and Rooted Page Rank. There is much room to explore low-rank approximations of such predictors and investigate whether they can be converted into accurate, scalable, graph embedding based, link predictors of the form we considered.

## Acknowledgements

The authors would like to thank the anonymous referee for their comments and suggestion.

## Funding

Defense Threat Reduction Agency (HDTRA1-15-1-0049).

## REFERENCES

1. NEWMAN, M. (2010) *Networks: An Introduction*. Oxford, UK: Oxford University Press.
2. GROSS, T. & SAYAMA, H. (2000) *Adaptive Networks: Theory, Models and Applications*. Dordrecht, Heidelberg, London, New York: Springer Publishing Company.
3. LIBEN-NOWELL, D. & KLEINBERG, J. (2001) The link-prediction problem for social networks. *J. Amer. Soc. Inf. Sci. Tech.*, **58**, 1019–1031.
4. QUERCIA, D., ASKHAM, H. & CROWCROFT, J. (2012) TweetLDA: supervised topic classification and link prediction in Twitter. *Proceedings of the 4th Annual ACM Web Science Conference*. New York, NY, USA: ACM, pp. 247–250.
5. BARZEL, B. & BARABÁSI, A.-L. (2013) Network link prediction by global silencing of indirect correlations. *Nature Biotechnol.*, **31**, 720–725.
6. CLAUSET, A., MOORE, C. & NEWMAN, M. E. (2008) Hierarchical structure and the prediction of missing links in networks. *Nature*, **453**, 98–101.
7. SRINIVAS, V. & MITRA, P. (2016) *Link Prediction in Social Networks-Role of Power Law Distribution*. Switzerland: Springer.
8. FOUSS, F., PIROTTE, A., RENDERS, J.-M. & SAERENS, M. (2007) Random-walk computation of similarities between nodes of a graph with application to collaborative recommendation. *IEEE Trans. Knowl. Data Eng.*, **19**, 355–369.
9. LENHOF, H.-P. (1992) The k closest pairs problem. <http://people.scs.carleton.ca/michiel/k-closestnote.pdf>.
10. WANG, P., XU, B., WU, Y. & ZHOU, X. (2015) Link prediction in social networks: the state-of-the-art. *Sci. China Inf. Sci.*, **58**, 1–38.
11. LICHTENWALTER, R. N., LUSSIER, J. T. & CHAWLA, N. V. (2010) New perspectives and methods in link prediction. *Proceedings of the 16th ACM SIGKDD International conference on Knowledge Discovery and Data Mining*. New York, NY, USA: ACM, pp. 243–252.
12. VON LUXBURG, U., RADL, A. & HEIN, M. (2010) Getting lost in space: Large sample analysis of the resistance distance, *Advances in Neural Information Processing Systems 23: 24th Annual Conference on Neural Information Processing Systems 2010* (J. Lafferty, C. K. I. Williams, J. Shawe-Taylor, R. S. Zemel & A. Culotta eds), Norwich, UK: Curran, pp. 2622–2630.
13. LESKOVEC, J., BACKSTROM, L., KUMAR, R. & TOMKINS, A. (2008) Microscopic evolution of social networks. *Proceedings of the 14th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. New York, NY, USA: ACM, pp. 462–470.
14. ESTRADA, E. & ARRIGO, F. (2015) Predicting triadic closure in networks using communicability distance functions. *SIAM J. Appl. Math.*, **75**, 1725–1744.
15. NEWMAN, M. E. (2001) Clustering and preferential attachment in growing networks. *Phys. Rev. E*, **64**, 025102.
16. DOYLE, P. G. & SNELL, J. L. (1984) *Random Walks in Electrical Networks*. 1st edn, vol. 22, Mathematical Association of America, JSTOR, [www.jstor.org/stable/10.4169/j.ctt5hh804](http://www.jstor.org/stable/10.4169/j.ctt5hh804).
17. HERSH, R. & GRIEGO, R. (1969) Brownian motion and potential theory. *Sci. Amer.* **220**, 67–74.
18. KLEIN, D. J. & RANDIC, M. (1993) Resistance distance. *J. Math. Chem.*, **12**, 81–95.
19. GHOSH, A., BOYD, S. & AMIN, S. (2008) Minimizing effective resistance of a graph. *SIAM Rev.*, **50**, 37–66.
20. CHANDRA, A. K., RAGHAVAN, P., RUZZO, W. L., & SMOLENSKY, R. (1989) The electrical resistance of a graph captures its commute and cover times. *Proceeding STOC'89 Proceedings of the twenty-first annual ACM symposium on Theory of computing*, New York, NY, USA: ACM, pp. 574–586.
21. HERSH, A. K., RAGHAVAN, P., RUZZO, W. L., SMOLENSKY, R., & TIWARI, P. (1996) The electrical resistance of a graph captures its commute and cover times. *Comput. Complexity*, **6**, 312–340.

22. MANGUOGLU, M., COX, E., SAIED, F. & SAMEH, A. (2010) TRACEMIN-Fiedler: a parallel algorithm for computing the Fiedler vector. *International Conference on High Performance Computing for Computational Science – VECPAR 2010*. Lecture Notes in Computer Science, vol. 6449. Heidelberg, Berlin: Springer, pp. 449–455.
23. HU, Y. & SCOTT, J. (2003) *HSL MC73: A Fast Multilevel Fiedler and Profile Reduction Code*. Oxford, UK: Rutherford Appleton Laboratory.
24. LESKOVEC, J., KLEINBERG, J. & FALOUTSOS, C. (2007) Graph evolution: densification and shrinking diameters. *ACM Trans. Knowl. Discov. Data*, **1**, 1–40.
25. NETWORKS COLLECTION, K. (2016) arXiv hep-th network dataset—KONECT. <http://konect.uni-koblenz.de/networks/ca-cit-HepTh> (accessed on 1 March 2017).
26. NETWORKS COLLECTION, K. (2016) arXiv hep-ph network dataset—KONECT. <http://konect.uni-koblenz.de/networks/ca-cit-HepPh> (accessed on 1 March 2017).
27. NETWORKS COLLECTION, K. (2016) Facebook friendships network dataset—KONECT. <http://konect.uni-koblenz.de/networks/facebook-wosn-links> (accessed on 1 March 2017).
28. VISWANATH, B., MISLOVE, A., CHA, M., & GUMMADI, K. P. (2009) On the evolution of user interaction in Facebook. *Proceedings of the 2nd ACM workshop on Online social networks (WOSN'09)*, New York, NY, USA: ACM, pp. 37–42.
29. NEWMAN, M. E. (2001) The structure of scientific collaboration networks. *Proc. Nat. Acad. Sci.*, **98**, 404–409.
30. FOUSS, F., SAERENS, M. & SHIMBO, M. (2016) *Algorithms and Models for Network Data and Link Analysis*. New York, NY, USA: Cambridge University Press.