

# Subgraph Extraction for Trust Inference in Social Networks

Yuan Yao<sup>1,2</sup>, Hanghang Tong<sup>3</sup>, Feng Xu<sup>1,2</sup>, Jian Lu<sup>1,2</sup>

<sup>1</sup>State Key Laboratory for Novel Software Technology, Nanjing University, China

<sup>2</sup>Department of Computer Science and Technology, Nanjing University, China

<sup>3</sup>IBM T.J. Watson Research, Hawthorne NY, USA

yyao@smail.nju.edu.cn, htong@us.ibm.com, {xf, lj}@nju.edu.cn

**Abstract**—Trust inference is an essential task in many real world applications. Most of the existing inference algorithms suffer from the scalability issue, making themselves computationally costly, or even infeasible, for the graphs with more than thousands of nodes. In addition, the inference result, which is typically an abstract, numerical trustworthiness score, might be difficult for the end-user to interpret.

In this paper, we propose subgraph extraction to address these challenges. The core of the proposed method consists of two stages: *path selection* and *component induction*. The outputs of both stages can be used as an intermediate step to speed up a variety of existing trust inference algorithms. Our experimental evaluations on real graphs show that the proposed method can accelerate existing trust inference algorithms, while maintaining high accuracy. In addition, the extracted subgraph provides an intuitive way to interpret the resulting trustworthiness score.

## I. INTRODUCTION

Trust inference, which aims to infer a trustworthiness score from the trustor to the trustee in the underlying social network, is an essential task in many real world applications including e-commerce [1], peer-to-peer networks [2], mobile ad hoc networks [3], etc.

To date, many trust inference algorithms have been proposed, which can be categorized into two main classes (See Section VI for a review): (a) path-based inference [4], [5], [6], [7], [8] and (b) component-based inference [9], [10], [11], [12].

Despite their own success, most of the existing inference algorithms have two limitations. The first challenge lies in *scalability* - many existing algorithms become very time-consuming, or even computationally infeasible for the graphs with more than thousands of nodes. Additionally, some algorithms assume the existence of a subgraph while how to construct such a subgraph remains an open issue [8]. The second challenge is the *usability* of the inference results. Most, if not all, of the existing inference algorithms output an abstract numerical trustworthiness score. This gives a quantitative measure of *to what extent* the trustor should trust the trustee, but gives few cues on *how* the trustworthiness score is inferred. This usability/interpretation issue becomes more evident when the size of the underlying graph increases, since we cannot even display the entire graph to the end-users (See Fig. 9 for an example).

In this paper, we propose subgraph extraction to address these challenges. The core of our subgraph extraction consists

of two stages: *path selection* and *component induction*. In the first (path selection) stage, we extract a few, important paths from the trustor to the trustee. In the second (component induction) stage, we propose a novel evolutionary algorithm to generate a small subgraph based on the extracted paths. The outputs of these two stages are then used as an intermediate step to speed up the path-based inference and component-based inference algorithms, respectively. Our experimental evaluations on real graphs show that the proposed method can significantly accelerate existing trust inference algorithms (up to 2,400x speed-up), while maintaining high accuracy (P-error is less than 0.04). In addition, the extracted subgraph provides an intuitive way to interpret the resulting trustworthiness score by presenting a concise summarization on the relationship from the trustor to the trustee. To the best of our knowledge, we are the first to propose subgraph extraction for trust inference. We believe that our work can improve most of the existing trust inference algorithms by (1) scaling up as well as (2) delivering more usable (i.e., interpretation-friendly) inference results to the end-users.

The rest of the paper is organized as follows. Section II presents the detailed definition of the subgraph extraction problem for trust inference. Section III and Section IV describe the algorithms for the path selection stage and component induction stage, respectively. Section V presents our experimental setup and results. Section VI reviews the related work of subgraph extraction and trust inference. Section VII concludes the paper.

## II. PROBLEM DEFINITION

Following the standard notations in the existing trust inference algorithms, we model the trust relationships in social networks as a weighted directed graph [13], [14]. The nodes of the graph represent the participants in the network, and the weight on each edge indicates the local trust value derived from the historical interactions.

We then categorize the existing trust inference algorithms into two major classes: *path-based trust inference* and *component-based trust inference*.

*Definition 1:* Path-based trust inference.

Path-based trust inference includes the approaches, which are started by the trustor, to evaluating the trustworthiness of the

trustee, through a set of paths from the trustor to the trustee in the network.

**Definition 2:** Component-based trust inference.

Component-based trust inference includes the approaches, which are started by the trustor, to evaluating the trustworthiness of the trustee, through a connected component from the trustor to the trustee in the network.

Let us first explain the differences between these two classes. In path-based trust inference, trust is propagated along a path, and the propagated trust from multiple paths is combined to form a final trustworthiness score. In contrast, there is no explicit concept of paths in component-based trust inference. Instead, it takes the initial graph as input and treats trust as, for example, random walks on a Markov chain [15].

As to the similarity, both classes belong to the subjective trust metrics [11], where different trustors can form different opinions on the same trustee. Accordingly, path-based trust inference such as [4], [5], [6], [7], [8] and component-based inference such as [9], [10], [11], [12] all belong to trust inference algorithms. Although the main focus of this paper is on the subjective metrics, our proposed subgraph extraction can also be applied to the objective trust metrics.

Despite the success of most existing inference algorithms, they share the scalability and usability limitations. To address these issues, we propose subgraph extraction for trust inference. The core of our subgraph extraction consists of two stages. The first stage, which serves for path-based trust inference, selects a set of paths from the trustor to the trustee. The second stage aims to produce a connected component between the trustor and the trustee for component-based trust inference. In addition, the second stage of our subgraph extraction produces a relatively small subgraph which can be clearly displayed and help the end-user better understand the inference result.

We now formally define the subgraph extraction problem for trust inference. In accordance to the corresponding two stages, the problem is divided into two subproblems: *path selection problem* and *component induction problem*.

**Definition 3:** Path Selection Problem.

Given: a weighted directed graph  $G(V, E)$ , two nodes  $s, t \in V$ , and an integer  $K$ ;

Find: a set  $C$  with  $K$  paths from  $s$  to  $t$  that minimizes the error function  $f(C)$ .

**Definition 4:** Component Induction Problem.

Given: a set  $C$  of paths from  $s$  to  $t$ , and an integer  $N$ ;

Find: an induced component  $H(V', E')$  with at most  $N$  edges that minimizes the error function  $g(H)$ , where  $V' \subseteq \{v | (u, v) \in P \text{ or } (v, u) \in P, P \in C\}$  and  $E' \subseteq \{e | e \in P, P \in C\}$ .

We next discuss the error function in the definitions. The error function  $f(C)$  in Definition 3 indicates the goodness of the extracted paths, and  $f(C)$  reaches its minimum value when  $C$  contains all the possible paths from  $s$  to  $t$ . Similarly, the error function  $g(H)$  in Definition 4 reaches its minimum value if  $H = G$ . In this paper, we use *P-error*, which is defined as

**Algorithm 1** KS algorithm. (See the appendix for the details)

**Input:** Weighted directed graph  $G(V, E)$ , two nodes  $s, t \in V$ , and a parameter  $K$  of path number

**Output:** Set  $C$  with  $K$  paths from  $s$  to  $t$

1:  $C = \text{k-shortest}(G, s, t, K)$   
2: **return**  $C$

follows, as the error function for both subproblems, i.e.,  $f = g = \text{P-error}$ .

**Definition 5:** P-error.

For a given trustor-trustee pair, the error function *P-error* is defined as

$$P - \text{error} = |p_{\text{sub}} - p_{\text{whole}}|,$$

where  $p_{\text{sub}}$  is the trustworthiness score inferred from the subgraph, and  $p_{\text{whole}}$  which serves as a ground truth, is the trustworthiness score inferred from the whole graph.

### III. PATH SELECTION

In the path selection stage, we aim to extract a few paths from the trustor to the trustee as an intermediate step to speed up path-based trust inference algorithms. These extracted paths will also serve as the input for the component induction stage.

There are two preprocessing steps in our extraction method. First of all, trust is interpreted as the probability by which the trustor expects that the trustee will perform a given action. This interpretation of trust is adopted by many existing trust inference algorithms, and it allows trust to be multiplicatively propagated along a path [16]. Second, we transform probability into weight by negative logarithm. Namely, the local trust value on the edge  $e$  is interpreted as probability  $p(e)$ , and the probability  $p(e)$  is transformed to weight  $w(e) = -\log(p(e))$ . Based on these two steps, the weight of a path  $P$  can be presented as

$$w(P) = \sum_{e \in P} -\log(p(e)) = -\log\left(\prod_{e \in P} p(e)\right) = -\log(\text{Pr}(P)).$$

As a result, finding a path of high trustworthiness in the original network is equivalent to finding a short path in the transformed network. We will use this transformed weighted graph  $G(V, E)$  as the input of our method.

Then, the path selection problem becomes to extract top- $k$  short paths from the trustor to the trustee in the transformed graph  $G(V, E)$ . Many existing algorithms can be plugged into this stage, such as Yen's k-shortest loopless paths (KS) [17], path sampling (PS) [18], etc. In our experiments, we found that KS algorithm performs best even if the multiplicative property of the interpretation does not hold, and we therefore recommend KS in this stage. A brief skeleton of the KS algorithm is shown in Algorithm 1, and the detailed algorithms for KS and PS are presented in the appendix for completeness.

#### A. Algorithm Analysis

The worst-case time complexity of KS is  $O(K|V|(|E| + |V|\log|V|))$ , which is known as the best result to ensure that k-shortest loopless paths can be found in a directed graph [19].

---

**Algorithm 2** EVO algorithm.

---

**Input:** Set  $C$  of paths from  $s$  to  $t$  and the directly induced component  $G^c(V^c, E^c)$ , as well as a constraint  $N$  of the edge number

**Output:** Induced component  $H(V', E')$  with at most  $N$  edges

```
1: define 0/1 vector  $B$  of size  $|E^c|$  where each element in  $B$ 
   stands for the existence of a corresponding edge in  $G^c$ 
2: initialize  $m$  vectors  $S \leftarrow \{B_1, B_2, \dots, B_m\}$ , with at most
    $N$  1-bits for each vector
3: while  $iter > 0$  do
4:   for each vector  $B_i$  in  $S$  do
5:     repeat
6:       mutate  $B_i$  to  $B_{i+m}$  with mutation probability
          $1/|E^c|$ 
7:     until the number of 1-bits in  $B_{i+m} \leq N$ 
8:   end for
9:   compute P-error results for the  $2m$  vectors
      $\{B_1, B_2, \dots, B_{2m}\}$ 
10:   $S \leftarrow$  the best  $m$  vectors from the  $2m$  ones
11:   $iter \leftarrow iter - 1$ 
12: end while
13:  $B_{final} \leftarrow$  the best vector in  $S$ 
14: return the corresponding component  $H(V', E')$  of
      $B_{final}$ 
```

---

However, the actual wall-clock time of KS on many real graphs is often much better than such worst case scenario [20]. In fact, based on our experiments, we find that it empirically scales near linearly wrt the graph size  $|V|$  in the chosen datasets.

#### IV. COMPONENT INDUCTION

In the component induction, we take the output of path selection stage (i.e., a set of  $K$  paths) as input, and output a small connected component from the trustor to the trustee. The output of the component induction stage not only acts as an intermediate step to speed up component-based trust inference algorithms, but also helps to improve the usability of trust inference by interpreting the inference results for the end-users. Notice that although our upcoming proposed algorithm EVO could also be applied on the whole graph, we do not recommend it in practice for the following two reasons: (1) most trustworthy paths have already been captured by the path selection stage (i.e., KS, etc), and (2) applying EVO on the whole graph would cost more memory and time to achieve high accuracy. We will present more detailed experimental evaluations to validate this in the next section.

In general, our proposed EVO algorithm (shown in Algorithm 2) belongs to the so-called evolutionary methods [21]. It aims to minimize P-error under the constraint of edge number. The input component  $G^c(V^c, E^c)$  is directly induced from the set  $C$  of paths from  $s$  to  $t$ , where  $V^c = \{v | (u, v) \in P \text{ or } (v, u) \in P, P \in C\}$  and  $E^c = \{e | e \in P, P \in C\}$ . There are two implicit parameters in the algorithm, i.e., the initial vector number  $m$  and iteration number  $iter$ .

We now explain EVO in detail. The first step of EVO is to establish a one-to-one correspondence between the edges in  $G^c$  and the elements in vector  $B$ . Each element of  $B$  is a 0/1 bit where 1 indicates that the corresponding edge exists and 0 indicates otherwise. The vector has exactly  $|E^c|$  bits where  $|E^c|$  is the edge number of  $G^c$ . In the second step, the algorithm generates  $m$  vectors  $B_1, B_2, \dots, B_m$ , and each of them has at most  $N$  1-bits. In our implementation, we apply a constant-time search in  $C$  to find a subset of paths with minimized P-error. In the following steps, EVO adopts *mutation* on each of these vectors to separately generate  $m$  new vectors  $B_{m+1}, B_{m+2}, \dots, B_{2m}$ . In the mutation from  $B_i$  to  $B_{i+m}$ , each bit of  $B_i$  is changed with probability  $1/|E^c|$ . If the resulting vector has more than  $N$  1-bits, the mutation operation is redone. The error function, which is P-error in our case, is then computed on each of these  $2m$  vectors, and the  $m$  vectors with smallest P-error are kept to the next iteration. For efficiency, the P-error computation on vector  $B$  herein means computing the P-error between  $G^c(V^c, E^c)$  and the component corresponding to the vector  $B$ . Namely, we use the input component  $G^c(V^c, E^c)$  as an approximation of the ground truth in this stage.

##### A. Algorithm Analysis

The time complexity of EVO is summarized in the following lemma, which basically says that the expected time complexity of EVO scales linearly wrt both initial vector number  $m$  and iteration number  $iter$ .

*Lemma 1: The average-case time complexity of EVO is  $O(iter \cdot m(|E^c|/N + \theta))$ , where  $\theta$  is the time complexity of the error function computation.*

**Proof:** In the mutation step of EVO, with mutation probability  $1/|E^c|$ , the expected number of bit changes is 1. This step is expected to be redone only when the number of 1-bits is  $N$  and the bit change is from 0 to 1. Under this condition, the probability of bit change from 0 to 1 is  $(|E^c| - N)/|E^c|$ . Therefore, the expected iteration number of the mutation step is  $|E^c|/N$ . Therefore, the whole expected time complexity of EVO is  $O(iter \cdot m \cdot |E^c|/N + m\theta) = O(iter \cdot m(|E^c|/N + \theta))$ , which completes the proof.  $\square$

#### V. EXPERIMENTS

##### A. Experimental Setup

Before presenting the experimental results, we first describe the datasets and the representatives of path-based and component-based trust inference algorithms. All algorithms are implemented in Java, and have been run on a T400 ThinkPad with 1280m jvm heap space. Few other activities are done during the experiments.

1) *Datasets Description:* We use the advogato<sup>1</sup> datasets in our experiments, because advogato is a trust-based social network and it contains multilevel trust assertions. There are four levels of trust assertions in the network, i.e., ‘Observer’, ‘Apprentice’, ‘Journeyer’, and ‘Master’. These assertions can

<sup>1</sup>[http://www.trustlet.org/wiki/Advogato\\_dataset](http://www.trustlet.org/wiki/Advogato_dataset).

TABLE I  
HIGH LEVEL STATISTICS OF ADVOGATO DATASETS.

Graph	Nodes	Edges	Avg. degree	Avg. clustering [22]	Avg. diameter [23]	Date
advogato-1	279	2,109	15.1	0.45	4.62	2000-02-05
advogato-2	1,261	12,176	19.3	0.36	4.71	2000-07-18
advogato-3	2,443	22,486	18.4	0.31	4.67	2001-03-06
advogato-4	3,279	32,743	20.0	0.33	4.74	2002-01-14
advogato-5	4,158	41,308	19.9	0.33	4.83	2003-03-04
advogato-6	5,428	51,493	19.0	0.31	4.82	2011-06-23

be mapped into real numbers in  $[0,1]$ . In our experiments, we map ‘Observer’, ‘Apprentice’, ‘Journeyer’, and ‘Master’ to 0.1, 0.4, 0.7, and 0.9, respectively. The statistics of the datasets is shown in Table I.

2) *Trust Inference Representatives*: To evaluate our subgraph extraction method, we need to apply trust inference algorithms on the whole graph and on our extracted subgraph to compare their effectiveness and efficiency. We chose *CertProp* [7] as the representative of path-based inference algorithms, and *Appleseed* [11] as the representative of component-based inference algorithms.

P-error computation in *CertProp* needs to first compute the ground truth  $p_{whole}$  by finding all paths from the trustor to the trustee in the whole graph. This computation, however, easily causes the overflow of the jvm heap space even on the advogato-1 graph. Following the suggestions in the original *CertProp* [7], we apply the fixed search strategy and search all paths whose length is not longer than seven as an approximation of the ground truth. For *CertProp*, we define *collapsed samples* as the trustor-trustee pairs of which the P-error computation either exceeds the range of `Java.lang.Double` or runs out of the jvm heap space. We randomly select 100 node pairs out of 122 samples, where the rest 22 of them are collapsed samples. Our experimental results are all based on the average of these 100 samples. Notice that, as discussed in the path selection section, the multiplicative property of the probability interpretation does not hold in *CertProp*. As to *Appleseed*, we apply linear normalization on the outputs, since the algorithm can produce arbitrary trustworthiness scores.

## B. Experimental Results

We now present the experimental results of our subgraph extraction method. In our experiments, the effectiveness, efficiency comparisons, and interpretation results are all based on the advogato-1 graph, as we found *CertProp* on the whole graph becomes computationally infeasible on all the other larger datasets. We evaluate the scalability of our method using all the datasets (i.e., advogato-1 to advogato-6). As for EVO, we set  $m = 5$  and  $iter = 10$  unless otherwise specified. The edge constraint  $N$  is set as  $K/2$ .

1) *Effectiveness*: For effectiveness, we first study how *CertProp* and *Appleseed* perform on the KS subgraph (the output of path selection stage) and EVO subgraph (the output of component induction stage), respectively. The results are shown in Fig. 1. We can observe that all the P-error values of *CertProp* and *Appleseed* are less than 0.04, indicating that our extracted subgraphs, which are based on a small set of

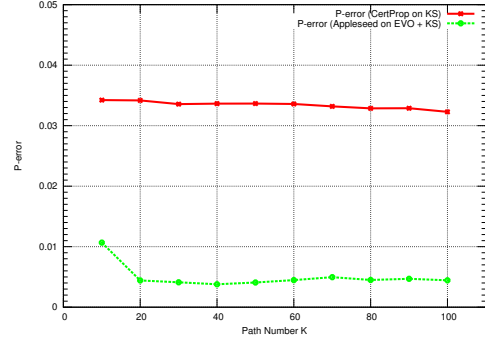


Fig. 1. Effectiveness of our subgraph extraction method with edge number constraint  $N = K/2$ . In all cases, the P-error is less than 0.04.

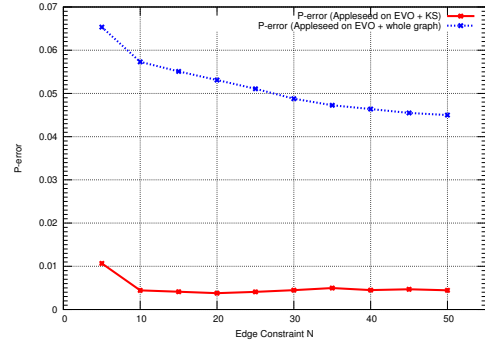


Fig. 2. Comparison of EVO on KS vs. EVO on the whole graph with edge number constraint  $N = K/2$ . EVO on KS outperforms EVO on the whole graph.

carefully selected paths and an evolutionary strategy, provide high accuracy for the trust inference algorithms.

Remember that the proposed EVO is always applied on the output of the path selection stage (referred to as ‘EVO+KS’). Here, for comparison purpose, we also apply EVO on the entire graph (referred to as ‘EVO+whole graph’). With the same parameter setting, the results are shown in Fig. 2. It can be seen that EVO on KS outperforms EVO on the whole graph. The reason is as follows. As an evolutionary algorithm, EVO (either on KS or on the entire graph) finds a local minima. By restricting the search space to those highly trustworthy paths (i.e., the output of KS), it converges to a better local minima in terms of P-error.

Finally, to compare EVO with existing component induction algorithms, we implement the *Monte Carlo pruning (MC)*

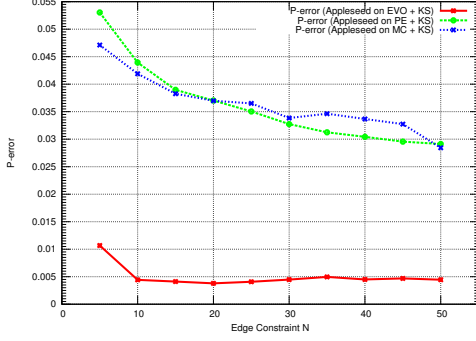


Fig. 3. Comparison of different component induction algorithms with edge number constraint  $N = K/2$ . EVO outperforms the existing component induction algorithms.

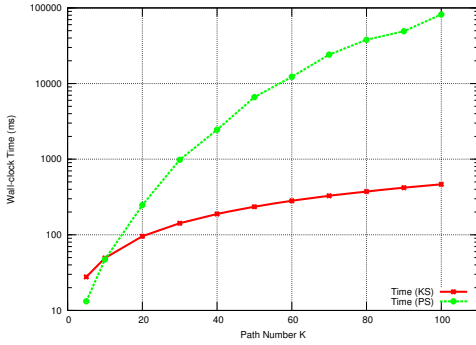


Fig. 4. The average wall-clock time of KS and PS. The average wall-clock time of KS is much faster than that of PS when  $K$  is greater than 30.

method [24] and the *proximity extraction (PE)* method [25], and plot the results in Fig. 3. Again, we can see that EVO outperforms both MC and PE wrt P-error. In fact, MC induces a component by successively deleting edges (edge-level component induction), while PE only selects a smaller set of paths (path-level component induction). Our EVO algorithm combines these two levels of component induction by searching a smaller set of paths in the initial step and then evolving the resulting component on the edge level.

2) *Efficiency*: First, we compare the different algorithmic choices in the path selection stage. To this end, we compare the wall-clock time of KS with an alternative path selection algorithm *path sampling (PS)* [18]. The results are shown in Fig. 4. Note that the y-axis is of log scale. As we can see from the figure, although PS is slightly faster than KS when  $K = 5$ , the wall-clock time of PS is much longer than that of KS when  $K$  is greater than 30. For example, the wall-clock time of PS is more than 170x longer than that of KS when  $K = 100$ . Therefore, we recommend using KS for path selection.

Next, we study the computational savings by applying the proposed subgraph extraction as the intermediate steps for the existing trust inference algorithms. To this end, we report the wall-clock time of CertProp on the output of the path selection stage, and Applesed on the output of the component induction stage, respectively. The results are shown in Fig. 5

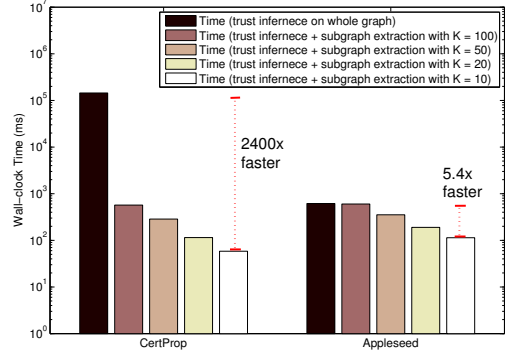


Fig. 5. The average wall-clock time of CertProp on KS and Applesed on KS+EVO. We achieve up to 2400x speed-up.

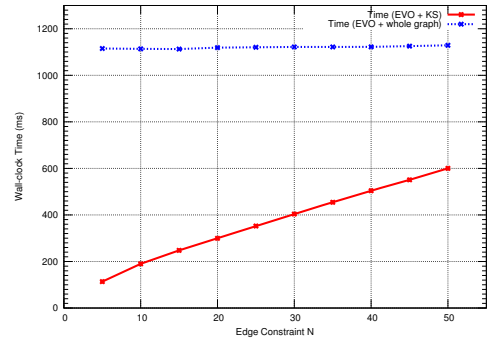


Fig. 6. The average wall-clock time of EVO on KS and EVO on the whole graph with edge number constraint  $N = K/2$ . EVO on KS is much faster.

where the y-axis is of log scale. Notice that the reported time includes the wall-clock time of both subgraph extraction and trust inference. In the figure, we also plot the wall-clock time of CertProp and Applesed on the entire graph for comparison. We can see that our subgraph extraction method saves the wall-clock time for both path-based trust inference and component-based trust inference, especially for the former one. For example, when  $K = 10$ , our subgraph extraction method achieves up to 2,400x and 5.4x speed-up for CertProp and Applesed, respectively. Even when  $K$  grows to more than 60, our method can still achieve 200-400x speed-up for CertProp.

Next, we compare the efficiency between applying EVO on KS and applying EVO on the whole graph. With  $N = K/2$ , the results are shown in Fig. 6. As we can see, the wall-clock time of EVO on KS (which includes the wall-clock time of both EVO and KS) is much faster than EVO on the whole graph. Together with the effectiveness results (Fig. 2), we recommend running EVO on the KS subgraph in practice.

Finally, we evaluate how the parameters  $m$  and  $iter$  in EVO affect the wall-clock time. In this experiment, we fix  $K = 20$  and  $N = 10$ , and the results are shown in Fig. 7. We can observe that the wall-clock time of EVO scales linearly wrt  $iter$  for any fixed  $m$ , which is consistent with the time complexity analysis shown before.

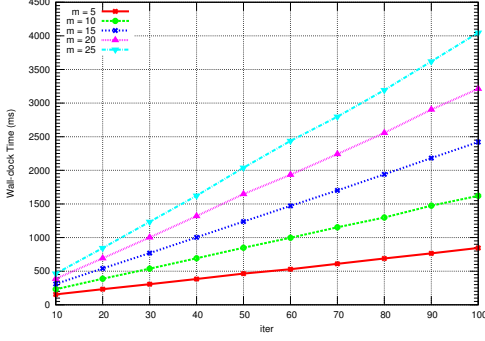


Fig. 7. The average wall-clock time of EVO with  $K = 20$  and  $N = 10$ . EVO scales linearly wrt  $iter$  for the fixed  $m$ .

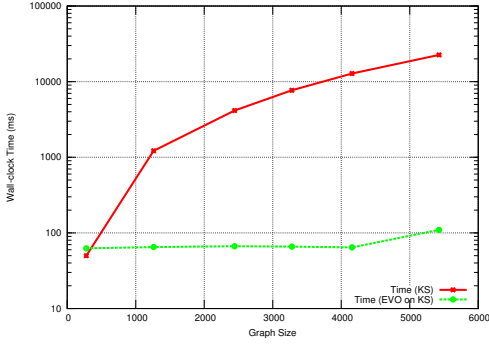


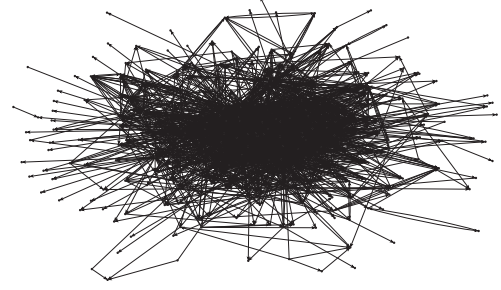
Fig. 8. The scalability of our subgraph extraction method. KS scales near linearly wrt the graph size, while the wall-clock time of EVO stays almost constant.

3) *Scalability*: We now evaluate the scalability of our method on datasets advogato-1 to advogato-6. Fig. 8 shows the results, where the y-axis is of log scale. In this experiment, we fix  $K = 10$  and  $N = 5$ .

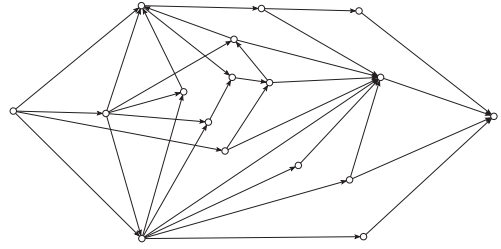
We can observe from the figure that even on the largest graph of 5,428 nodes and 51,293 edges, KS can help to infer the trustworthiness score within 25 seconds. In addition, KS scales near linearly wrt the underlying graph size. As to EVO, the wall-clock time stays stable in spite of the growth of the graph size. The reason is that  $|E^c|$  scales near linearly to  $K$  due to many overlapping edges, and  $N$  is set to  $K/2$ . Consequently,  $|E^c|/N$  is close to a constant, and the time complexity of EVO can be approximated to  $O(iter \cdot m \cdot \theta)$ .

4) *Usability/Interpretation*: Another important goal of the proposed EVO is to improve the usability in trust inference by interpreting the inferred trustworthiness score for end-users. An illustrative example is shown in Fig. 9. The whole graph and the induced KS subgraph by the path selection stage are also plotted for comparison.

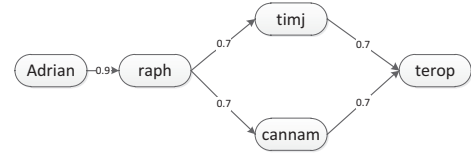
From the figures we can see that the whole graph is hard for interpretation. As to the KS subgraph, although the number of edges has significantly decreased compared with the original whole graph, there are still some redundant edges which might diverge end-users' attention. On the other hand, the EVO subgraph only presents the most important participants and



(a) The original whole graph



(b) KS subgraph with  $K = 20$ . The paths are from 'Adrian' (the leftmost node) to 'terop' (the rightmost node).



(c) EVO subgraph with  $N = 10$  on KS-20. The component is from 'Adrian' to 'terop'.

Fig. 9. The interpretation example of the whole graph, KS-20, and EVO-10 on KS-20.

their trust opinions, providing a much clearer explanation on how the trustworthiness score is inferred.

## VI. RELATED WORK

We review the related work in this section, which can be categorized into two parts: trust inference algorithms and subgraph extraction.

### A. Trust Inference

We categorize existing trust inference algorithms into two main classes: path-based trust inference and component-based trust inference.

In the first class of path-based inference, Wang and Singh [5], [26] as well as Hang et al. [7] propose operators to concatenate trust along a path and aggregate trust from multiple paths. Liu et al. [6] argue that not only trust values but social relationships and recommendation role are important for trust inference. However, these algorithms are only suitable for small networks due to their complexity. Some other path-based trust inference algorithms, such as [4], [8], assume the existence of an extracted subgraph while how to construct such a subgraph remains an open issue [8].

In the second class of component-based inference, EigenTrust [2] tries to compute an objective trustworthiness score

for each node in the graph. In contrast to EigenTrust, our main focus is to provide support for subjective trust metrics where different trustors can form different opinions on the same trustee. Existing subjective trust algorithms, including [9], [10], [11], [27], [28], take the initial graph as input and treat trust as random walks on a Markov chain or on a graph. Our subgraph extraction method not only can speed up many of these algorithms but also can provide interpretive result which is not considered by the existing algorithms.

Overall, our subgraph extraction is motivated to address the two common challenges (i.e., scalability and usability) shared by most of these existing trust inference algorithms.

### B. Subgraph Extraction

Several end-to-end subgraph extraction algorithms are developed to solve different problems.

In the field of graph mining, Faloutsos et al. [29] refer to the idea of electrical current where trust relationships are modeled as resistors, and try to find a connection subgraph that maximizes the current flowing from source to target. Later, Tong et al. [30] generalize the connection subgraph to directed graphs and use the subgraph to compute proximities between nodes. Similar to Tong et al., Koren et al. [25] also try to induce a subgraph for proximity computation. In addition, Koren et al. search the  $k$ -shortest paths to provide a basis for measuring the proximity.

Recently, several algorithms are proposed for reliable subgraph extraction. Among them, Monte Carlo pruning [24] measures the relevance of each edge by Monte Carlo simulations, and tends to remove the edge of lowest relevance one by one. The most related work is perhaps the randomized Path Covering algorithm [18] which also consists of two stages of path sampling and subgraph construction. However, both Monte Carlo pruning and Path Covering tend to find a subgraph with highest probability to be connected, while we aim to find a subgraph to address the scalability and usability issues in trust inference.

## VII. CONCLUSIONS

In this paper, we have proposed subgraph extraction to address the scalability and usability challenges of existing trust inference algorithms. The core of our subgraph extraction has two stages, and the outputs of both stages can be used as an intermediate step to speed up a variety of existing trust inference algorithms. Our experimental evaluations on real graphs show that the proposed method can significantly accelerate existing trust inference algorithms (up to 2,400x speed-up), while maintaining high accuracy (P-error is less than 0.04). In addition, the extracted subgraph provides an intuitive way to interpret the inferred trustworthiness score for end-users. Future work includes incorporating distrust in the subgraph extraction.

## VIII. ACKNOWLEDGEMENT

This work is supported by the National Natural Science Foundation of China (No. 61021062, 61073030), and the

National 973 Program of China (No. 2009CB320702). The second author was partly sponsored by the Army Research Laboratory and was accomplished under Cooperative Agreement Number W911NF-09-2-0053.

## REFERENCES

- [1] L. Xiong and L. Liu, "Peertrust: Supporting reputation-based trust for peer-to-peer electronic communities," *IEEE Transactions on Knowledge and Data Engineering*, vol. 16, no. 7, pp. 843–857, 2004.
- [2] S. D. Kamvar, M. T. Schlosser, and H. Garcia-Molina, "The Eigentrust algorithm for reputation management in p2p networks," in *WWW*. ACM, 2003, pp. 640–651.
- [3] S. Buchegger and J.-Y. Le Boudec, "A robust reputation system for mobile ad-hoc networks," KTH Royal Institute of Technology, Theoretical Computer Science Group, Tech. Rep., 2004.
- [4] L. Mui, M. Mohtashemi, and A. Halberstadt, "A computational model of trust and reputation," in *HICSS*. IEEE, 2002, pp. 2431–2439.
- [5] Y. Wang and M. P. Singh, "Trust representation and aggregation in a distributed agent system," in *AAAI*, 2006, pp. 1425–1430.
- [6] G. Liu, Y. Wang, and M. Orgun, "Trust inference in complex trust-oriented social networks," in *ICCSE*. IEEE, 2009, pp. 996–1001.
- [7] C.-W. Hang, Y. Wang, and M. P. Singh, "Operators for propagating trust and their evaluation in social networks," in *AAMAS*, 2009, pp. 1025–1032.
- [8] G. Wang and J. Wu, "Multi-dimensional evidence-based trust management with multi-trusted paths," *Future Generation Computer Systems*, vol. 27, no. 5, pp. 529–538, 2011.
- [9] R. Guha, R. Kumar, P. Raghavan, and A. Tomkins, "Propagation of trust and distrust," in *WWW*. ACM, 2004, pp. 403–412.
- [10] P. Massa and P. Avesani, "Controversial users demand local trust metrics: An experimental study on opinions. com community," in *AAAI*, 2005, pp. 121–126.
- [11] C. Ziegler and G. Lausen, "Propagation models for trust and distrust in social networks," *Information Systems Frontiers*, vol. 7, no. 4, pp. 337–358, 2005.
- [12] R. Zhou and K. Hwang, "Powertrust: A robust and scalable reputation system for trusted peer-to-peer computing," *IEEE Transactions on Parallel and Distributed Systems*, vol. 18, no. 4, pp. 460–473, 2007.
- [13] G. Barbian, "Assessing trust by disclosure in online social networks," in *ASONAM*, 2011, pp. 163–170.
- [14] Y. Yao, J. Zhou, L. Han, F. Xu, and J. Lü, "Comparing linkage graph and activity graph of online social networks," in *SocInfo*, 2011, pp. 84–97.
- [15] M. Richardson, R. Agrawal, and P. Domingos, "Trust management for the Semantic Web," in *ISWC*. Springer, 2003, pp. 351–368.
- [16] G. Liu, Y. Wang, and M. Orgun, "Optimal social trust path selection in complex social networks," in *AAAI*, 2010, pp. 1391–1398.
- [17] J. Yen, "Finding the  $k$  shortest loopless paths in a network," *Management Science*, pp. 712–716, 1971.
- [18] P. Hintsanen, H. Toivonen, and P. Sevon, "Fast discovery of reliable subnetworks," in *ASONAM*, 2010, pp. 104–111.
- [19] J. Hershberger, M. Maxel, and S. Suri, "Finding the  $k$  shortest simple paths: A new algorithm and its implementation," *ACM Transactions on Algorithms*, vol. 3, no. 4, p. 45, 2007.
- [20] E. Martins and M. Pascoal, "A new implementation of Yens ranking loopless paths algorithm," *4OR: A Quarterly Journal of Operations Research*, vol. 1, no. 2, pp. 121–133, 2003.
- [21] T. Bäck, *Evolutionary algorithms in theory and practice: evolution strategies, evolutionary programming, genetic algorithms*. Oxford University Press, USA, 1996.
- [22] D. Watts and S. Strogatz, "Collective dynamics of 'small-world' networks," *Nature*, vol. 393, no. 6684, pp. 440–442, 1998.
- [23] J. Leskovec, J. Kleinberg, and C. Faloutsos, "Graphs over time: densification laws, shrinking diameters and possible explanations," in *KDD*. ACM, 2005, pp. 177–187.
- [24] P. Hintsanen and H. Toivonen, "Finding reliable subgraphs from large probabilistic graphs," *Data Mining and Knowledge Discovery*, vol. 17, no. 1, pp. 3–23, 2008.
- [25] Y. Koren, S. North, and C. Volinsky, "Measuring and extracting proximity in networks," in *KDD*. ACM, 2006, pp. 245–255.
- [26] Y. Wang and M. P. Singh, "Formal trust model for multiagent systems," in *IJCAI*, 2007, pp. 1551–1556.



- [27] U. Kuter and J. Golbeck, “Sunny: A new algorithm for trust inference in social networks using probabilistic confidence models,” in *AAAI*, 2007, pp. 1377–1382.
- [28] K. Nordheimer, T. Schulze, and D. Veit, “Trustworthiness in networks: A simulation approach for approximating local trust and distrust values,” in *IFIPTM*, vol. 321. Springer-Verlag, 2010, pp. 157–171.
- [29] C. Faloutsos, K. S. McCurley, and A. Tomkins, “Fast discovery of connection subgraphs,” in *KDD*. ACM, 2004, pp. 118–127.
- [30] H. Tong, C. Faloutsos, and Y. Koren, “Fast direction-aware proximity for graph mining,” in *KDD*. ACM, 2007, pp. 747–756.
- [31] J. Golbeck and J. Hendler, “Inferring binary trust relationships in Web-based social networks,” *ACM Transaction on Internet Technology*, vol. 6, pp. 497–529, 2006.
- [32] G. Robins, P. Pattison, Y. Kalish, and D. Lusher, “An introduction to exponential random graph (p\*) models for social networks,” *Social Networks*, vol. 29, no. 2, pp. 173–191, 2007.

#### APPENDIX

To find  $K$  short paths from graph  $G(V, E)$  in the path selection stage, many existing algorithms can be used. We consider two representative algorithms from the literature. Here, we present the detailed algorithm description for completeness.

The first algorithm is *Yen’s k-shortest loopless paths (KS)* algorithm [17], which is shown in Algorithm 3.

---

#### Algorithm 3 Detailed KS algorithm.

---

**Input:** Weighted directed graph  $G(V, E)$ , two nodes  $s, t \in V$ , and a parameter  $K$  of path number

**Output:** Set  $C$  with  $K$  paths from  $s$  to  $t$

```

1:  $X \leftarrow$  shortest path from  $s$  to  $t$ 
2:  $C \leftarrow$  shortest path from  $s$  to  $t$ 
3: while  $|C| < K$  and  $X \neq \emptyset$  do
4:    $P \leftarrow$  remove the shortest path in  $X$ 
5:    $d \leftarrow$  the deviation node of  $P$ 
6:   for each node  $v$  between  $d$  (inclusive) and trustee  $t$  (exclusive) in  $P$  do
7:      $pre \leftarrow$  subpath from trustor  $s$  to  $v$  in  $P$ 
8:      $post \leftarrow$  the deviated shortest path from  $v$  to  $t$ 
9:     combine  $pre$  and  $post$ , and add it to  $X$ 
10:  end for
11:   $C \leftarrow C +$  the shortest path in  $X$ 
12: end while
13: return  $C$ 
```

---

In the algorithm, we use Dijkstra’s algorithm for finding a shortest path. All the computed paths are loopless by temporarily removing visited nodes. The key idea of the KS algorithm is *deviation*. The *deviation node*  $d$  of path  $P$  is the node that makes  $P$  deviate from existing paths in the candidate set  $C$ . For each node  $v$  between  $d$  (inclusive) and trustee  $t$  (exclusive) in  $P$ , the *deviated shortest path* from node  $v$  to  $t$  is computed by temporarily removing the edge starting at  $v$  in  $P$ . The computed deviated shortest path  $post$  and the subpath  $pre$  (the path from  $s$  to  $v$  in  $P$ ) are combined to form a possible path candidate. For the nodes before  $d$ , possible shortest paths are already computed and included in  $X$ . Based on deviation, KS finds the  $K$  shortest paths from trustor  $s$  to trustee  $t$  one by one. Following Martins and Pascoal’s implementation [20], we compute the deviated shortest path from deviation node  $d$  to the trustee in a reverse order.

---

#### Algorithm 4 PS algorithm.

---

**Input:** Weighted directed graph  $G(V, E)$ , two nodes  $s, t \in V$ , and a parameter  $K$  of path number

**Output:** Set  $C$  with  $K$  paths from  $s$  to  $t$

```

1:  $C \leftarrow$  shortest path from  $s$  to  $t$ 
2: while  $|C| < K$  do
3:   re-decide all the edges in  $E$ 
4:   for each path  $P$  in  $C$  do
5:     if  $P$  is decided as true then
6:        $F \leftarrow F + P$ 
7:     end if
8:   end for
9:   while  $F \neq \emptyset$  do
10:    re-decide the most overlapped edge in  $F$  as failed
11:    remove failed paths from  $F$ , if there are any
12:  end while
13:   $P \leftarrow$  the shortest path among the non-failed edges from  $s$  to  $t$ 
14:  if  $P \neq \emptyset$  then
15:     $C \leftarrow C + P$ 
16:  end if
17: end while
18: return  $C$ 
```

---

The other algorithm is the randomized algorithm *path sampling (PS)* [18], which is proposed for the *most reliable subgraph problem* [24]. While PS is proposed for undirected graphs, trust relationships in social networks should be directed as trust is asymmetric in nature [31]. Therefore, we adapt PS (as shown in Algorithm 4) for a directed graph.

PS considers the input graph as a Bernoulli random graph [32], and the algorithm is based on the *edge decision* of this random graph. An edge is randomly decided as true with probability  $p(e)$ , and a path is decided as true if all the edges on the path are decided as true. At the beginning of each iteration, all the edges of the graph are re-decided, and these *graph decisions* provide opportunities for distrust information to be contained. Like KS, PS first adds a shortest path into candidate set  $C$ . PS then tries to find a graph decision based on which none of the paths in  $C$  is true. To avoid the situation when this graph decision is hardly found, PS stores the true paths in  $C$  to a temporary set  $F$ , and deliberately fails the most overlapping edges in  $F$  until none of the paths in  $F$  is true. Finally, based on the results of graph decision and edge failing, PS finds the shortest path  $P$  among the non-failed edges from trustor  $s$  to trustee  $t$ , and adds it to  $C$ . The algorithm ends until  $K$  paths are found.

PS allows some distrust information to be incorporated into the extracted subgraph, which could in turn lower the P-error based on our experiments. However, the time complexity of PS is difficult to estimate, since the wall-clock time depends on the graph density. In addition, as shown in our experiments, the wall-clock time of PS is especially long when  $K$  becomes sufficiently large. We conjecture that PS can be used in dense graphs where numerous paths exist between node pairs.