

NETWORK DISTANCE PREDICTION

Dhiviyaa.S¹, Kanchana.V²

¹ Assistant professor, Department of BCA and MSC SS, Sri Krishna Arts and Science College, Coimbatore, India

² Student, Department of BCA and MSC SS, Sri Krishna Arts and Science College, Coimbatore, India

Abstract – Coordinates-based mechanisms have proven to be useful in a communication-to-communication architecture to predict Internet network route, even more so than the already existing ID Maps scheme. One such mechanism has been analyzed known as Global Network Positioning (GNP) to provide insight on a current state of the art technology in network distance prediction. Global Network Positioning instantiates a virtual geometric space by applying coordinates to nodes in a network, coordinates are computed by utilizing an objective function. A set of landmark hosts are first deployed into the geometric space to implement a set of reference points for any newly discovered host. Hosts also maintain their own coordinates, making it accessible to retrieve inter-communication network distance upon discovering each other by utilizing a route function. The variables associated with GNP must be tweaked to maximize efficiency. Through performed experiments using PlanetLab [1], a service that allows access to hosts at multiple research institutions across the globe, the factors that affect GNPs performance is analyzed.

Key Words: Predicting Network Distance, Global Network Positioning, Coordinates Based Network, Peer-To-Peer Network Optimization

1. INTRODUCTION

As communication-communication file sharing continues to grow in popularity ever since its debut over a ago, with such programs as Napster and Lime-wire, a need for predicting accurate network distance has emerged. A client’s ultimate goal within such applications is to seek for maximum available bandwidth between its self and its peers that contain the desired files to guarantee optimal transfer of data. Path optimization or other measurements of distance within a network is a somewhat impractical means of optimizing a clients bandwidth between peers; this approach becomes too costly in terms of speed and processing. Many efforts currently exist to effectively apply coordinates to a network or accurately predict network distance such as HTTP [8] or the Triangulated Heuristic. Some schemes are used more frequently than others while some may be just considered not feasible for real use. IDMaps is a state of the art system, where special HOPS servers are deployed, specifically utilized for the IDMaps scheme. HOPS servers hold a topological map of the net, a series of inter-communication measurements that are stored and may be retrieved by querying hosts for a prediction of distance between two nodes. Unfortunately, a few primary problems existed within the infrastructure of the scheme that needed

to be addressed. The distance between hosts x and y is defined by the distance between x and its nearest Tracer T_1 plus the distance between y and its nearest tracer T_2 , plus the shortest path between T_1 and T_2 [3][5]. This method presents a major flaw in accurately predicting inter-host distance. IDMaps constantly over-predicts the amount of distance between two hosts in a communication – communication network because it relies on Tracers to be close enough to a host so that a reasonably exactly prediction of network is possible. Like ID Maps other schemes carry similar problems of over prediction which was ultimately improved upon with GNP.

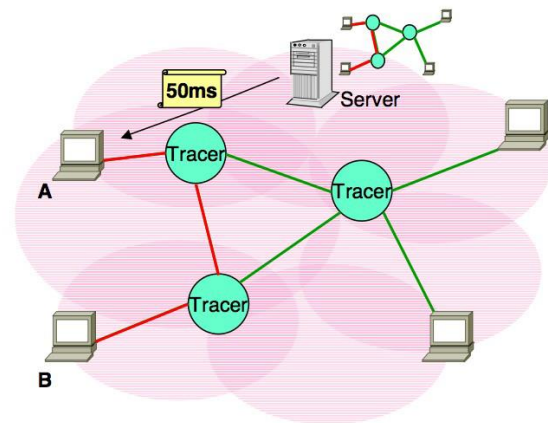


Fig. 1. An implementation of an IDMaps scheme

2. NETWORK DISTANCE PREDICTION BY MATRIX FACTORIZATION

This passage formulates the setback of consolidate distance of impossible feats by tricks abracadabra as matrix closing and describes its decree by matrix factorization. We furthermore provide a homogeneous look of antithetical approaches to absorb distance necromancy, the insights of which handle a unified optimization framework.□

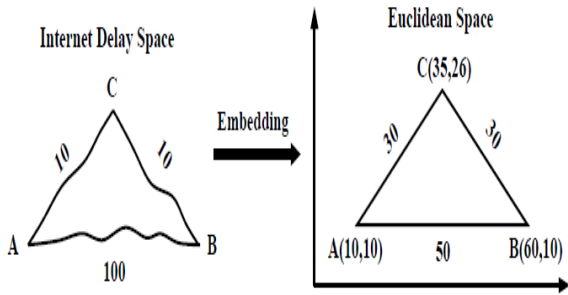


Fig. 1. Network distance prediction by Euclidean Embedding.

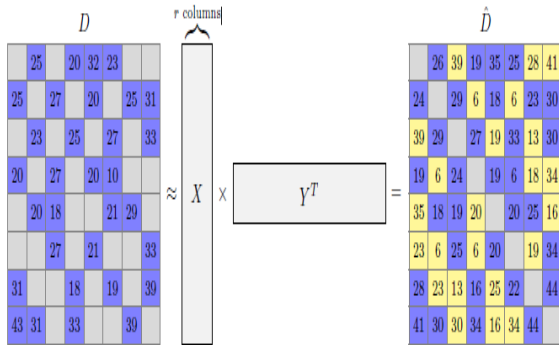


Fig. 2. Network distance prediction by matrix factorization. Note that the diagonal entries of D and D-hat are empty.

A. Problem Formulation

Assuming n nodes in the join, a n_n eclipse matrix is constructed mutually sprinkling distances during nodes measured and the others unmeasured. Denote D the measured top matrix by the whole of dij the measured eclipse from node I to node j and D-hat the predicted outstrip matrix by the whole of d-hat did the predicted transcend computed from some function.

Given the after notations, absorb distance illusion can be viewed as a matrix cessation problem that estimates the missing entries in D from a tiny number of met with entries. Its resolution generally amounts to minimizing a loss function of the following form

$$L(D, \hat{D}, W) = \sum_{i,j} w_{ij} l(d_{ij}; \hat{d}_{ij}); \quad (1)$$

where W is a weight matrix with w_{ij} taking values between 0 and 1. In a simple case, w_{ij} = 1 if d_{ij} is measured and 0 otherwise. Note that if the eclipse measurements are RTTs, once d_{ji} = d_{ij} as RTTs are necessarily symmetric. Consequently, w_{ji} = w_{ij} as d_{ji} and d_{ij} is, in turn, both met with or both unknown.

l is a loss function that penalizes the difference between an estimate and its desired or true value. The most commonly used loss function is the L2 or square loss function,

$$l(d; \hat{d}) = (d - \hat{d})^2; \quad (2)$$

We will discuss other loss functions in Section V.

B. Low-Rank Approximation and Matrix Factorization

Additional constraints are needed to solve the matrix completion problem in Eq. 1. A common approach is to constrain the rank of the approximate matrix D-hat so that

$$\text{Rank}(\hat{D}) = r; \quad (3)$$

where r_n for D of size n_n

The assumption in this low-rank approximation is that the entries of D are largely correlated, which causes D to have a low effective rank. To show that it holds for our problem, Figure 3 plots the singular values of two RTT matrices. It can be seen that the singular values of both matrices decrease fast as the 10th singular values are 5:7% and 2:9% of the largest ones respectively, indicating strong correlations in them. The low-rank nature of many other RTT datasets has been previously reported in.

Directly finding D-hat by minimizing Eq. 1 subject to Eq. 3 is considerably difficult due to the rank constraint. However, as D-hat is of low rank, we can factorize it into the product of two smaller matrices, i.e.,

$$\hat{D} = XY^T; \quad (4)$$

where X and Y are of size n_r. Therefore, we can get rid of the rank constraint by replacing D-hat by XY^T in eq. 1, and then look for X and Y instead by minimizing

$$L(D; X; Y; W) = \sum_{i,j} w_{ij} l(d_{ij}; x_{ij}y_{ij}^T); \quad (5)$$

where x_i is the ith row of X, y_i is the ith row of Y, and x_iy_i^T = d-hat_{ij} is the estimate of d_{ij}. Note that the factorization in Eq. 4 has no unique solution as

$$\hat{D} = XY^T = XGG^T Y^T; \quad (6)$$

where G is an arbitrary r_r invertible matrix. Thus, replacing X by XG and Y^T by G⁻¹Y^T results in the same D-hat

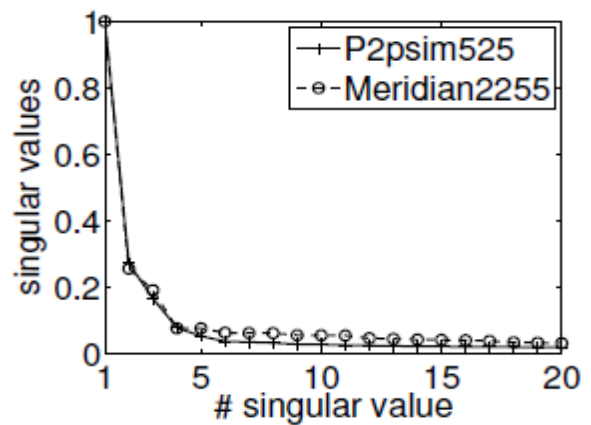


Fig. 3. The singular values of an RTT matrix of 2255_2255, extracted from the Meridian dataset [30] and called "Meridian2255", and of an RTT matrix of 525_525,

extracted from the P2psim dataset [30] and called "P2psim525". The singular values are normalized so that the largest singular values of both matrices are equal to 1. Generally, the class of techniques to solve the low-rank approximation is matrix factorization. When D is complete, analytic solutions can be found by using singular value decomposition (SVD) [31]. With missing entries, the factorization is usually done by iterative optimization methods such as Gradient Descent or Newton algorithms [32]. Note that additional constraints can be imposed in eq. 5. For instance, the entries of X and Y can be required to be non-negative in order to recover a nonnegative matrix, leading to the nonnegative matrix factorization (NMF) [33].

C. Incorporation of the Regularization

Matrix completion by matrix factorization suffers from a well-known problem called overfitting in the field of machine learning [34]. In words, directly optimizing eq. 5 often leads to a "perfect" model with no or small errors on the training data while having large errors on the unseen data which are not used in learning. The problem is more severe when D is sparse or when r is large.

A common way to avoid overfitting is through regularization that penalizes the norms of the solutions, resulting in the following regularized loss function,

$$L(D;X;Y;W; \lambda) = \sum_{i,j=1}^n w_{ij} l(d_{ij}; x_i^T y_j) + \lambda \sum_{i=1}^n \|x_i\|^2 + \lambda \sum_{i=1}^n \|y_i\|^2 \tag{7}$$

where λ is the regularization coefficient that controls the extent of regularization? Besides avoiding overfitting, the regularization also helps overcome the drift of the solutions due to the non-uniqueness of the factorization (see eq. 6), which often leads to the overflows of the solutions. Among the infinite number of pairs

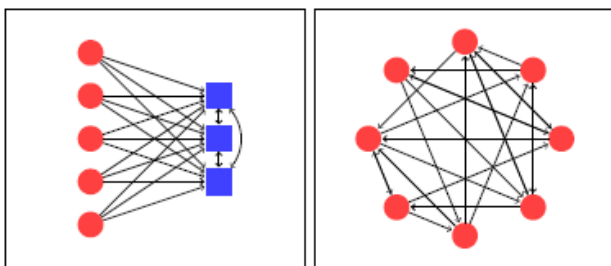


Fig. 4. Design of landmark-based, the left and right decentralized systems for distance vector prediction. The squares are landmarks and circles are ordinary node. The directed path from node I to node J means that node I probe node J and therefore $w_{IJ} = 1$. Of X and Y which produce the same \hat{D} , the incorporation of the regularization will force to choose the pair with the smallest norm.

D. A Unified View of Approaches to Network Distance Prediction Although near to one heart approaches to became

lost in eclipse illusion vary by adopting disparate models including Euclidean embedding and matrix factorization and by adopting disparate architectures of as a choice landmark-based or landmark-less and by means of this decentralized, these seemingly antithetical approaches bodily optimize the same field in Eq. 1 nonetheless differs unattended in the stage set of w_{ij} and in the associated distance functions to speculate \hat{d}_{ij} . Setting of w_{ij} : For landmark-based methods, as bodily paths, mid landmarks are measured and deformed nodes seize only the landmarks,

$$w_{ij} = \begin{cases} 1 & \text{if node } j \text{ is a landmark} \\ 0 & \text{otherwise} \end{cases}$$

For decentralized methods, as each node equally probes a number of nodes,

$$w_{ij} = \begin{cases} 1 & \text{if node } I \text{ probe node } j \\ 0 & \text{otherwise} \end{cases}$$

Figure 4 illustrates the architectures of landmark-based and decentralized systems. Distance functions to calculate \hat{d}_{ij} : For matrix factorization, as described above,

$$\hat{d}_{ij} = \|x_i - y_j\| \tag{8}$$

For Euclidean embedding, the Euclidean distance is defined as

$$\hat{d}_{ij} = \sqrt{(x_i - x_j)^T (x_i - x_j)} \tag{9}$$

where x_i and x_j are the Euclidean coordinates of node I and node j. The above insights suggest a unified framework to treat and to solve equally network distance prediction under different models and different architectures. For instance, the decentralized matrix factorization algorithms proposed in the following sections can be used to solve both Euclidean embedding and landmark-based systems with little modification.

3. EXPERIMENTS AND EVALUATIONS

In this section, we evaluate DMF3 and compare it with two popular NCS algorithms: Vivaldi and IDes. The former is based on metric space embedding, while the latter is also based on matrix factorization but uses landmarks. All the experiments are performed on two typical disclosure sets collecting trustworthy Internet measurements: the P2psim data exist which contains the measured distances surrounded by 1740 Internet DNS servers, and the Meridian data apply which contains the measured distances during 2500 nodes. While DMF bounces in principle act with regard to asymmetric eclipse matrices, in our demonstrate, we took $d_{i,j} = d_{j,i}$ and bounded these distances as the half of the round-trip-time mid nodes i and j. The agnate assumption is

adopted in Vivaldi and has the biggest slice of the cake of to a great extent simplifying the implementation of the algorithm, as reflection one-way restrain is abstract in practice. In the simulations, we randomly selected a node and updated its coordinates at each step. An iteration of a simulation is defined by a fixed round of node updates. Since Vivaldi updates its coordinates with respect to only one neighbor in contrast to DMF that does it with respect to all neighbors, an iteration in Vivaldi is defined by $n \times k$ node updates whereas in DMF an iteration is n node updates, where n is the number of nodes and k is the number of neighbors. In doing so, we ensure that, on average, all nodes have a chance to update their coordinates with respect to all neighbors. Note that IDES is not an iterative method. The coordinates of the nodes are unchanged. We recognize them from that day forward classical criticism criteria. [2]

- Cumulative Distribution of Relative Estimation Error Relative Estimation

Error (REE) is most zoned as

$$REE = \frac{|\hat{d}_{i,j} - d_{i,j}|}{d_{i,j}}$$

- Stress deciding the complete fitness of the embedding is bounded as [2]

$$\text{stress} = \frac{\sum_{i,j} (d_{i,j} - \hat{d}_{i,j})^2}{\sum_{i,j} d_{i,j}^2}$$

- Median Absolute Estimation Error (MAEE) is defined as

$$MAEE = \text{median}_{i,j} (|d_{i,j} - \hat{d}_{i,j}|).$$

Our DMF algorithm utilizes only a small % of the route measurements in the datasets to estimate the coordinates of the nodes, but the evaluation of the above criteria is done using all distance measurements.

4. CONCLUSION

Recent studies symbolize that the fascination quality of existing of impossible feats by tricks abracadabra mechanisms boot be incapable from the debate perspective. This handout has uncovered that interim it might be used to enliven the foreboding quality over intelligent landmark letter from uncle sam, it is unclear at which point to engineering the levy procedure in sending up the river to corroborate good foreboding quality from one end to the other all has a jump on ranges. Although choosing nearby nodes as landmarks can result in higher prediction accuracy for short links, longer links may suffer significantly degraded prediction accuracy. [2]

In light of this problem, we have proposed a hierarchical approach for network distance prediction. The hierarchical prediction leverages multiple coordinates at multiple distance scales. The right scale is chosen for predicting the target distance. We study two hierarchical prediction schemes. The first scheme leverages a shared landmark hierarchy. The breath scheme allows malleable landmark assignment at companionless nodes and the "hierarchy" is

marked by constantly smaller top scales. Experiments by the whole of Internet intensity traces bring to light that the hierarchical act outperforms the hybrid landmark assignment scheme: swiftly links cut back be predicted by the whole of higher accuracy by all of the little violence on the mind reader or daydream links. [2]

4. REFERENCES

[1] "Network Distance Prediction Based on Decentralized Matrix Factorization", Yongjun Liao¹, Pierre Geurts^{2,3}, Guy Leduc¹

¹ Research Unit in Networking (RUN), University of Liège, Belgium, ² Systems and Modeling, University of Liège, Belgium, ³ Research associate, FRS-F.N.R.S. (Belgium)

[2] "DMFSGD: A Decentralized Matrix Factorization Algorithm for Network Distance Prediction", Yongjun Liao, Wei Du, Pierre Geurts and Guy Leduc

[3] "A Hierarchical Approach to Internet Distance Prediction", Rongmei Zhang¹, Y. Charlie Hu¹, Xiaojun Lin¹, and Sonia Fahmy², School of Electrical and Computer Engineering, Department of Computer Science Purdue University, West Lafayette [2]

[4] "Analysis of Network Distance Prediction with Global Network Positioning Mathieu Rodrigue UCONN BioGrid, REU Summer 2009 Department of Computer Science, University of Hartford, 200 Bloomfield Ave, West Hartford.