

Relationship Prediction in Dynamic Heterogeneous Information Networks

No Author Given

No Institute Given

Abstract. Most real-world information networks, such as social networks, are heterogeneous and as such, the relationships observed in these networks can be diverse and carry differing semantics. Therefore techniques for link prediction in homogeneous networks cannot be directly applied on heterogeneous ones. On the other hand, works that investigate link prediction in heterogeneous networks do not necessarily consider network dynamism in sequential time intervals. In this work we propose a technique that leverages a combination of latent and topological meta path-based features to predict a target relationship between two nodes of given types in a dynamic heterogeneous information network. **Ebrahim**
►Before saying "in our experiments" we need to enumerate the distinguishing aspects of our work and say why they are significant contributions.◀ Our experiment results on two real-world information network datasets show 10-40% increase in AUCROC, and 25-30% increase in prediction accuracy compared to the state of the art techniques.

1 Introduction

The goal of link prediction in a network [17] is to estimate the likelihood of a future relationship between two nodes based on the observed graph. Predicting such connections in a network can be applied in different contexts such as recommendation systems [4, 27, 19, 16, 12], network reconstruction [11], node classification [10], or biomedical applications such as predicting protein-protein interactions [14]. Traditional link prediction techniques, such as [17], consider networks to be homogeneous, i.e., graphs with only one type of nodes and edges. However, most real-world networks, such as social networks, scholar networks, patient networks [6] and knowledge graphs [33] are heterogeneous information networks (HINs) [26] and have multiple node and relation types. For example, in a bibliographic network, there are nodes of types authors, papers, and venues, and edges of types writes, cites and publishes.

In a HIN, relations between different entities carry different semantics. For instance the relationship between two authors are different in meaning when they are co-authors compared to the case when one cites another's paper. Thus techniques for homogeneous networks [17, 32, 18, 15, 1] cannot be directly applied on heterogeneous ones. A few works such as [28, 29] investigated the problem of link prediction in HINs, however, they do not consider the dynamism of networks

and overlook the potential benefits of analyzing a heterogeneous graph as a sequence of network snapshots. To this end, existing work has already shown that in homogenous networks incorporating temporal changes improves link prediction accuracy [38]. Previous work on temporal link prediction scarcely studied HINs and to the best of our knowledge, the problem of predicting relationships in dynamic heterogeneous networks has not been studied before. A dynamic heterogeneous information network (DHIN) is a HIN, whose links have associated timestamps.

In this work we study the problem of relationship prediction in a DHIN, which can be stated as: *Given a DHIN graph G at t consecutive time intervals, the objective is to predict the existence of a particular relationship between two given nodes at time $t + 1$.* The major challenge in relationship prediction in DHINs is how to effectively combine the HIN topology features and inferred latent features that incorporate temporal changes in order to exhibit the best performance. Also, the prediction technique should be computationally efficient for large-scale networks. To this end, the main contributions of our work include:

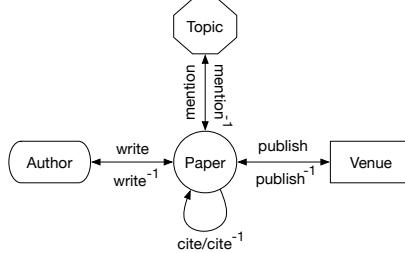
- We propose the problem of relationship prediction in a DHIN, and draw a contrast between this problem and existing link prediction techniques that have been proposed for dynamic and heterogeneous networks;
- We present a simple yet effective technique, called *MetaDynaMix*, that leverages topological meta path-based and latent features to predict a target relationship between two nodes in a DHIN;
- We empirically evaluate the accuracy of our proposed work on two real-world datasets, and the results show 10-40% increase in AUCROC, and 25-30% increase in prediction accuracy compared to the state of the art baselines.

2 Problem Statement

Our work is focused on heterogeneous information networks (graphs) that can change and evolve over time. As such, we first formally define the concept of *Dynamic Heterogeneous Information Networks*, as follows:

Definition 1 (Dynamic heterogeneous information network). *A dynamic heterogeneous information network (DHIN) is a directed graph $G = (V, E)$ with a node type mapping function $\phi : V \rightarrow \mathcal{A}$ and a link type mapping function $\psi : E \rightarrow \mathcal{R}$, where V , E , \mathcal{A} , and \mathcal{R} denote sets of nodes, links, node types, and relation types. Each node $v \in V$ belongs to a node type $\phi(v) \in \mathcal{A}$, each link $e \in E$ belongs to a relation $\psi(e) \in \mathcal{R}$, and $|\mathcal{A}| > 1$ and $|\mathcal{R}| > 1$. Also each edge $e = (u, v, t)$ is a temporal edge from a vertex u to a vertex v at time t . \square*

The DBLP bibliographic network is an example of a DHIN, containing different types of nodes such as papers, authors, topics, and publication venues, with publication links associated with a date. In the context of a heterogenous network, a *relation* can be in the form of a *direct link* or an *indirect link*, where an indirect link is a sequence of direct links in the network. Thus, two nodes

**Fig. 1.** Network schema for DBLP network.

might not be directly connected, however they might be considered to be indirectly connected through a set of intermediary links. In this work, we use the terms *relationship prediction* and *link prediction* interchangeably referring to predicting whether two nodes will be connected in the future via a *sequence of relations* in the graph, where the *length* of a sequence is greater than or equal to one. For instance in a bibliographic network, a direct link exists between an author and a paper she wrote, and an indirect link exists between her and her co-authors through the paper, which they wrote together. In order to better understand different types of nodes and their relation in a network, the concept of *network schema* [30] is used. A network schema is a meta graph structure that summarizes a HIN and is formally defined as follows:

Definition 2 (Network schema). For a heterogeneous network $G = (V, E)$, the network schema $S_G = (\mathcal{A}, \mathcal{R})$ is a directed meta graph where \mathcal{A} is the set of node types in V and \mathcal{R} is the set of relation types in E . \square

Figure 1 shows the network schema for the DBLP bibliographic network with $\mathcal{A} = \{\text{Author}, \text{Paper}, \text{Venue}, \text{Topic}\}$. In this paper, we refer to different types of nodes in the DBLP bibliographic network with abbreviations P for paper, A for author, T for topic, and V for venue.

Similar to the notion of network schema that provides a meta structure for the network, a *meta path* [30] provides a meta structure for paths between different node types in the network.

Definition 3 (Meta path). A meta path \mathcal{P} is a path in a network schema graph $S_G = (\mathcal{A}, \mathcal{R})$, denoted by $\mathcal{P}(A_1, A_{n+1}) = A_1 \xrightarrow{R_1} A_2 \dots \xrightarrow{R_n} A_{n+1}$, as a sequence of links between node types defining a composite relationship between a node of type A_1 and one of type A_{n+1} , where $A_i \in \mathcal{A}$ and $R_i \in \mathcal{R}$. \square

The *length* of a meta path is the number of relations in it. Note that given two node types A_i and A_j , there may exist multiple meta paths of different lengths between them. We call a path $p = (a_1 a_2 \dots a_{n+1})$ a *path instance* of a meta path $\mathcal{P} = A_1 - A_2 \dots - A_{n+1}$ if p follows \mathcal{P} in the corresponding HIN, i.e., for each node a_i in p , we have $\phi(a_i) = A_i$. The co-author relationship in DBLP can be described with the meta path $A \xrightarrow{\text{write}} P \xrightarrow{\text{write}^{-1}} A$ or in short $A-P-A$.

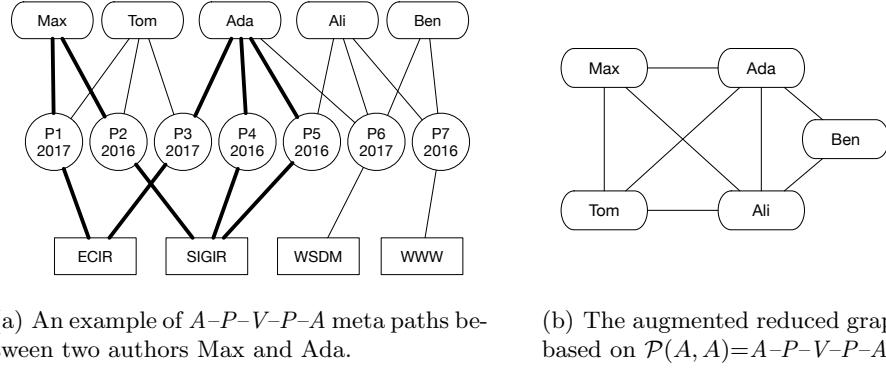


Fig. 2. An example of a publications network. Link formation time is presented as a year below a paper id rather than on links to make it easier to read.

Paths in thick solid lines in Figure 2(a) correspond to $A\text{-}P\text{-}V\text{-}P\text{-}A$ meta paths between *Max* and *Ada*, indicating they published in the same venue, such as *Max*-*P1*-*ECIR*-*P3*-*Ada*. Each meta path carries different semantics and defines a unique topology representing a special relation.

Meta Path-based Similarity Measures. Given a meta path $\mathcal{P} = (A_i, A_j)$ and a pair of nodes a and b such that $\phi(a) = A_i$ and $\phi(b) = A_j$, several *similarity measures* can be defined between a and b based on the path instances of \mathcal{P} . Examples of such similarity or proximity measures in a HIN are *path count* [30, 28], *PathSim* [30] or *normalized path count* [28], *random walk* [28], *HeteSim* [25], and *KnowSim* [34]. Without loss of generality, in this work we use *Path Count* (PC) as the default similarity measure. For example given the meta path $A\text{-}P\text{-}V\text{-}P\text{-}A$ and the HIN in Figure 2(a), $PC(\text{Max}, \text{Ada})=3$ and $PC(\text{Tom}, \text{Ada})=4$. We now formally define the problem that we target in this work as follows:

Definition 4 (Relationship prediction problem). *Given a DHIN graph G at time t , and a target relation meta path $\mathcal{P}(A_i, A_j)$ between nodes of type A_i and A_j , we aim to predict the existence of a path instance of \mathcal{P} between two given nodes of types A_i and A_j at time $t+1$. \square*

3 Proposed Relationship Prediction Approach

Given a DHIN graph $G = (V, E)$, we decompose G into a sequence of t HIN graphs G_1, \dots, G_t based on links with associated timestamps and then predict relationships in G_{t+1} . As mentioned in Definition 4, we intend to predict existence of a given type of relationship (target meta path) between two given nodes. Thus we define a new type of graph, called *augmented reduced graph*, that is generated according to a given heterogeneous network and a target relation meta path.

Definition 5 (Augmented reduced graph). *Given a HIN graph $G = (V, E)$ and a target meta path $\mathcal{P}(A_i, A_j)$ between nodes of type A_i and A_j , an augmented*

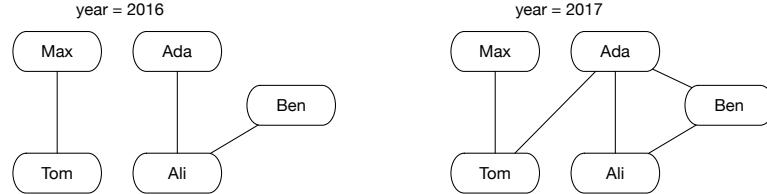


Fig. 3. Augmented reduced graphs for the network in Figure 2(a) with respect to the target meta path $A-P-A$ (co-authorship) in 2016 and 2017.

reduced graph $G^P = (V^P, E^P)$ is a graph, where $V^P \subseteq V$ and nodes in V^P are of type A_i and A_j , and edges in E^P indicate relationships of type P in G . \square

For example, an augmented reduced graph for the network in Figure 2(a) and target meta path $P(A, A)=A-P-V-P-A$ is a graph shown in Figure 2(b) whose nodes are of type *Author* and whose edges represent *publishing in the same venue*.

3.1 Homogenized Link Prediction

Once the given DHIN graph $G = (V, E)$ is decomposed to t HIN graphs G_1, \dots, G_t , one solution to the relationship prediction problem (Definition 4) is to build an augmented reduced graph G_i^P for each G_i with respect to the given target meta path P and then predict a link in G_i^P instead of a path in G_i . In other words, we generate a homogenized version of a graph snapshot and apply a link prediction method. Figure 3 shows examples of such graphs at different time intervals. The intuition behind considering different snapshots, i.e., a dynamic network, rather than a single snapshot for link prediction is that we can incorporate network evolution patterns to increase prediction accuracy. Our hypothesis is that the estimated graph \hat{G}_{i+1}^P depends on \hat{G}_i^P .

Recently researchers have focused on inferring latent space of networks for link prediction [38, 35, 23, 7, 21] based on the assumption that the probability of a link between two nodes depends on their positions in their latent space. Each dimension of the latent space characterizes an attribute, and the more two nodes share such attributes the more likely they connect (also known as homophily). Amongst such graph embedding methods, a few [7, 38] considered dynamic networks. Inspired by Zhu et al. [38], we formulate our problem as follows: Given a sequence of augmented reduced graphs G_1^P, \dots, G_t^P , we aim to infer a low rank k -dimensional latent space matrix Z_i for each adjacency matrix G_i^P at time i by minimizing

$$\begin{aligned} & \operatorname{argmin}_{Z_1, \dots, Z_t} \sum_{i=1}^t \left(\|G_i^P - Z_i Z_i^T\|_F^2 + \lambda \sum_{x \in V^P} (1 - Z_i(x) Z_{i-1}(x)^T) \right) \\ & \text{subject to : } \forall x \in V^P, i, Z_i \geq 0, Z_i(x) Z_i(x)^T = 1 \end{aligned} \quad (1)$$

Algorithm. 1 Homogenized Link Prediction

Input: A DHIN graph G , the number of snapshots t , a target meta path $\mathcal{P}(A, B)$, the latent space dimension k , the link to predict (a, b) at $t + 1$

Output: The probability of existence of link (a, b) in $G_{t+1}^{\mathcal{P}}$

```

1:  $\{G_1, \dots, G_t\} \leftarrow DecomposeGraph(G, t)$ 
2: for each graph  $G_i = (V_i, E_i)$  do
3:   for each node  $x \in V_i$  that  $\phi(x) = A$  do
4:     Follow  $\mathcal{P}$  to reach a node  $y \in V_i$  that  $\phi(y) = B$ 
5:     Add nodes  $x$  and  $y$ , and edge  $(x, y)$  to the augmented reduced graph  $G_i^{\mathcal{P}}$ 
6:   end for
7: end for
8:  $\{Z_1, \dots, Z_t\} \leftarrow MatrixFactorization(G_1^{\mathcal{P}}, \dots, G_t^{\mathcal{P}}, k)$ 
9: Return  $Pr((a, b) \in E_{t+1}^{\mathcal{P}}) \leftarrow \sum_{i=1}^k Z_t(a, i)Z_t(b, i)$ 

```

where $Z_i(x)$ is a temporal latent vector for node x at time i , λ is a regularization parameter, and $1 - Z_i(x)Z_{i-1}(x)^T$ penalizes sudden changes for x in the latent space. This optimization problem can be solved using gradient descent. The intuition behind the above formulation is two fold: (1) node with similar latent space representation are more likely to connect with each other, and (2) nodes typically evolve slowly over time and abrupt changes in their connection network are less likely to happen [37]. The matrix $G_{t+1}^{\mathcal{P}}$ can be estimated by $\Phi(f(Z_1, \dots, Z_t))$, where Φ and f are link and temporal functions, or simply by $Z_t Z_t^T$. Note that Z_i depends on Z_{i-1} as used in the temporal regularization term in Equation (1).

Algorithm 1 presents a concrete incorporation of Equation 1 for relation prediction. It takes as input a DHIN graph G , the number of graph snapshots t , a target relation meta path $\mathcal{P}(A, B)$, the latent space dimension k , and the link to predict (a, b) at $t+1$. It first decomposes G into a sequence of t graphs G_1, \dots, G_t by considering the associated timestamps on edges (line 1). Next from each graph G_i , a corresponding augmented reduced graph $G_i^{\mathcal{P}}$ is generated (lines 2-7) for which nodes are of type a and b (beginning and end of target meta path \mathcal{P}). For example given $\mathcal{P}(A, A)=A-P-A$, each $G_i^{\mathcal{P}}$ represents the co-authorship graph at time i . Finally by optimizing Equation (1) it infers latent spaces Z_1, \dots, Z_t (line 8) and estimates $G_{t+1}^{\mathcal{P}}$ using $Z_t Z_t^T$ (line 9).

3.2 Dynamic Meta Path-based Relationship Prediction

The above homogenized approach does not consider different semantics of meta paths between the source and destination nodes and assumes that the probability of a link between nodes depends only on their latent features. For instance, as depicted in Figure 3, *Tom* and *Ada* became co-authors in 2017 that can be due to publishing at the same venue in 2016, i.e., having two paths between them that passes through *SIGIR*, as shown in Figure 2. Similarly *Ben* and *Ada* who published with a same author, *Ali* in 2016, became co-authors in 2017.

We would like to further hypothesize that combining latent and topological features can increase prediction accuracy as we can learn latent features that fit the residual of meta path-based features. One way to combine these features is

to incorporate meta path measures in Equation (1) by changing the loss function and regularization term as:

$$\begin{aligned} & \underset{\boldsymbol{\theta}_i, Z_i}{\operatorname{argmin}} \sum_{i=1}^t \left\| G_i^P - \left(Z_i Z_i^T + \sum_{i=1}^n \theta_{i-1} \mathcal{F}_{i-1}^{P_i} \right) \right\|_F^2 + \\ & \lambda \sum_{i=1}^t \left(\sum_{x \in V^P} (1 - Z_i(x) Z_{i-1}(x)^T) + \sum_{i=1}^n \theta_i^2 \right) \end{aligned} \quad (2)$$

where n is the number of meta path-based features, \mathcal{F}^{P_i} is the i^{th} meta path-based feature matrix defined on G_i , and θ_i is the weight for feature f_i . Although we can use a fast block-coordinate gradient descent [38] to infer Z_i s, it cannot be efficiently applied to the above changed loss function. This is because it requires computing meta paths for all possible pairs of nodes in \mathcal{F}^{P_i} for all snapshots, which is not scalable, as calculating similarity measures, such as Path Count or PathSim, can be very costly. For example computing path counts for $A-P-V-P-A$ meta path can be done by multiply adjacency matrices $AP \times PV \times VP \times PA$.

As an alternative solution, we build a predictive model that considers a linear combination of topological and latent features. These features, however, can be interpolated in different ways that is beyond the scope of this work. Given the training pairs of nodes and their corresponding meta path-based and latent features, we apply logistic regression to learn the weights associated with these features. We define the probability of forming a *new link* in time $t+1$ from node a to b as $Pr(\text{label} = 1 | a, b; \boldsymbol{\theta}) = \frac{1}{e^{-z} + 1}$, where $z = \sum_{i=1}^n \theta_i f_t^{P_i}(a, b) + \sum_{j=1}^k \theta_{n+j} Z_t(a, j) Z_t(b, j)$, and $\theta_1, \theta_2, \dots, \theta_n$ and $\theta_{n+1}, \theta_{n+2}, \dots, \theta_{n+k}$ are associated weights for meta path-based features and latent features at time t between a and b . Given a training dataset with l instance-label pairs, we use logistic regression with L_2 regularization to estimate the optimal $\boldsymbol{\theta}$ as:

$$\hat{\boldsymbol{\theta}} = \underset{\boldsymbol{\theta}}{\operatorname{argmin}} \sum_{i=1}^l -\log Pr(\text{label} | a_i, b_i; \boldsymbol{\theta}) + \lambda \sum_{j=1}^{n+k} \theta_j^2 \quad (3)$$

We preferred to combine features in this learning framework since G_i is very sparse and thus the number of newly formed links are much less compared to all possible links. Consequently calculating meta path-based features for the training dataset is scalable compared to the matrix factorization technique. Moreover, similar to [28], in order to avoid excessive computation of meta path-based measures between nodes that might not be related, we confine samples to pairs that are located in a nearby neighborhood. More specifically, for each source node x in G_i^P , we choose target nodes that are within two hops of x but not in 1-hop, i.e., are not connected to x in G_i^P . We first find all target nodes that make a new relationship with x in G_{i+1}^P and label respective samples as positive. Next we sample an equal number of negative pairs, i.e., those targets that do not make new connection, in order to balance our training set. Once the dataset is built, we perform logistic regression to learn the model and then apply the predictive

Algorithm. 2 Dynamic Meta path-based Relationship Prediction

Input: A DHIN graph G , the number of snapshots t , a network schema S , a target meta path $\mathcal{P}(A, B)$, the maximum length of a meta path l , the latent space dimension k , the link to predict (a, b) at $t + 1$

Output: The probability of existence of link (a, b) in $G_{t+1}^{\mathcal{P}}$

- 1: $\{G_1, \dots, G_t\} \leftarrow DecomposeGraph(G, t)$
- 2: Generate target augmented reduced graphs $G_1^{\mathcal{P}}, \dots, G_t^{\mathcal{P}}$ following Algorithm 1 lines 2-7
- 3: $\{\mathcal{P}_1, \dots, \mathcal{P}_n\} \leftarrow GenerateMetaPaths(S, \mathcal{P}(A, B), l)$
- 4: $\{Z_1, \dots, Z_t\} \leftarrow MatrixFactorization(G_1^{\mathcal{P}}, \dots, G_t^{\mathcal{P}}, k)$
- 5: **for** each pair (x, y) , where $x \in V_{t-1}$ and $y \in N(x)$ is a nearby neighbor of x in $G_{t-1}^{\mathcal{P}}$ **do**
- 6: Add feature vector $\langle f_{t-1}^{\mathcal{P}_i}(x, y) \text{ for } i = 1..n, Z_{t-1}(x, j)Z_{t-1}(y, j) \text{ for } j = 1..k \rangle$ to the training set T with $label=1$ if (x, y) is a new link in $E_t^{\mathcal{P}}$ otherwise $label=0$.
- 7: **end for**
- 8: $model \leftarrow Train(T)$
- 9: Return $Pr((a, b) \in E_{t+1}^{\mathcal{P}}) \leftarrow Test(model, \langle f_t^{\mathcal{P}_i}(a, b) \text{ for } i = 1..n, Z_t(a, j)Z_t(b, j) \text{ for } j = 1..k \rangle)$

model to the feature vector for the target link. The output probability can be later interpreted as a binary value based on a cut-off threshold.

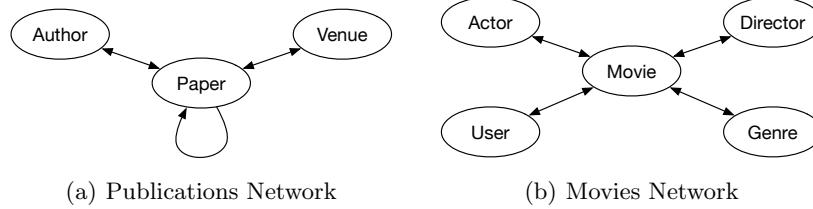
We describe steps for building and applying our predictive model, called *MetaDynaMix*, in Algorithm 2. The algorithm takes as input a DHIN graph G , the number of graph snapshots t , a network schema S , a target relation meta path $\mathcal{P}(A, B)$, the maximum length of a meta path l , the latent space dimension k , and the link to predict (a, b) at $t + 1$. Similar to Algorithm 1, it decomposes G into a sequence of graphs (line 1). Next it generates augmented reduced graphs $G_i^{\mathcal{P}}$ s from G_i s based on \mathcal{P} for nodes which are of type A and B (beginning and end of meta path \mathcal{P}) (line 2) as explained in Algorithm 1. It then produces the set of all meta paths between nodes of type A and type B defined in $\mathcal{P}(A, B)$ (line 3). This is done by traversing the network schema S (for instance through BFS traversal) and generating meta paths with the maximum length of l . It then applies matrix factorization to find latent space matrices Z_i (line 4). Next it creates a training dataset for sample pairs (x, y) with feature set containing meta path-based measures $f_t^{\mathcal{P}_i}(x, y)$ for each meta path \mathcal{P}_i , and latent features $Z_t(a, j)Z_t(b, j)$ for $j = 1..k$ at time t , and $label=1$ if (x, y) is a new link in $G_{t+1}^{\mathcal{P}}$ otherwise $label=0$ (lines 5-7). Subsequently the algorithm trains the predictive model (line 8), generates features for the given pair (a, b) , and tests it using the trained model (line 9).

4 Experiments

4.1 Experiment Setup

Dataset. We conduct our experiments on two real-world network datasets that have different characteristics and evolution behaviour.

Publications dataset: The AMiner citation dataset [31] version 8 (2016-07-14) is extracted from DBLP, ACM, and other sources. It contains 3,272,991 papers and 8,466,859 citation relationships for 1,752,443 authors, who published in 10,436 venues, from 1936 to 2016. Each paper is associated with an abstract,

**Fig. 4.** The simplified network schema used for our experiments.**Table 1.** Meta paths for publications dataset with $V=\{\text{Author}, \text{Paper}, \text{Venue}\}$ and movies dataset with $V=\{\text{User}, \text{Movie}, \text{Actor}, \text{Director}, \text{Genre}\}$.

Network	Meta path	Meaning
Publications	$A-P-A$	[The target relation] Authors are coauthors
	$A-P-V-P-A$	Authors publish in the same venue
	$A-P-A-P-A$	Authors have the same co-author
	$A-P-P-P-A$	Authors cite the same papers
Movies	$U-M$	[The target relation] A user watches a movie
	$U-M-A-M$	A user watches a movie with the same actor
	$U-M-D-M$	A user watches a movie with the same director
	$U-M-G-M$	A user watches a movie of the same genre
	$U-M-U-M$	A user watches a movie that another user

authors, year, venue, and title. We confined our experiments on papers published since 1996, which includes 2,935,679 papers. Similar to [28], we considered only authors with at least 5 papers.

Movies dataset: The RecSys HetRec movie dataset [3] is an extension of MovieLens10M dataset, published by the GroupLens research group that links the movies of MovieLens dataset with their corresponding web pages on IMDB and Rotten Tomatoes. It contains information of 2,113 users, 10,197 movies, 20 movie genres (avg. 2.04 genres per movie), 4,060 directors, 95,321 actors (avg. 22.78 actors per movie), 72 countries, and 855,598 ratings (avg. 404.92 ratings per user, and avg. 84.64 ratings per movie).

Experiment Settings. We describe meta paths and target relationships, baseline methods, and different parameter settings.

Meta Paths and Target Relationships. Figure 4 depicts network schemas for the two datasets. Note that we consider a simplified version and ignore nodes such as topic for papers or tag for movies. Table 1 presents a number of meta paths that we consider in our experiments, where target meta path relations are *co-authorship* and *watching*.

Amin ► For the publications network we consider meta paths $A-P-V-P-A$, $A-P-A-P-A$, and $A-P-P-P-A$, as the study in [28] shows that shared co-authors, shared venues, and co-cited papers for two authors significantly contribute to their future collaborations. There

are two major differences between target relation types in our datasets. First, unlike a new co-authorship relation that happens at a particular time, users can watch/rate a movie once it is released. In other words each paper is published once and a new co-authorship is made at that time whereas users create new watching relations to an existing movie. Second, the target relation for the publications dataset, i.e., A-P-A, has the same node type at both ends, while the target meta path for the movie dataset, i.e., U-M, considers two different node types. Note that G_i^P 's in Equation 1 are square adjacency matrices. For the case of having target relations with two types of nodes at ends, we consider 0 value for the relationships of the same type in case no such relation actually exists in the network. \blacktriangleleft

Baseline Methods. Sun et al. [28] proposed a supervised learning framework for link prediction in HINs, called PathPredict, that learns coefficients associated with meta path-based features by maximizing the likelihood of new relationship formation. Their model is learned based on one past interval and does not consider temporal changes in different intervals. We perform comparative analysis of our work, denoted as MetaDynaMix, with four techniques: (1) The original PathPredict that considers only 3 intervals, (2) PathPredict applied on different time intervals, denoted as PathPredict+, (3) homogenized link prediction (Section 3.1), denoted as HLP, and (4) logistic regression on HLP latent features, denoted as LRHLP. The authors in [28] showed that PathPredict outperforms traditional link prediction approaches that use topological features defined in homogeneous networks such as common neighbors or Katz β , and thus we do not include these techniques in our experiments.

Parameters. We set the number of snapshots $t=3, 5$, and 7 to evaluate the effect of dynamic analysis of different time intervals. Note that $t=3$ refers to the default case for many link prediction algorithms that learn based on one interval and test based on another. More specifically in the training phase features are extracted based on T1 and labels are determined based on T2, and for the testing phase features are calculated based on T2 and labels are derived from T3. In our experiments we did not observe a considerable change in prediction performance by setting the number of latent features k to 5, 10, and 20, and thus all presented results are based on setting k to 20.

Implementation. We use the implementation of matrix factorization for inferring temporal latent spaces of a sequence of graph snapshots presented in [38]. We use all the default settings such as the number of latent features k to be 20, and the optimization algorithm to be the local block-coordinate gradient descent. For the classification part, we use the efficient LIBLINEAR [8] package and set the type of solver to L2-regularized logistic regression (primal).

Evaluation Metrics. To asses link prediction performance, we use Area Under Curves (AUC) for Receiver Operating Characteristic (ROC) [5] and accuracy (ACC). We also perform the McNemar's test [20] to assess the statistical significance of the difference between classification techniques.

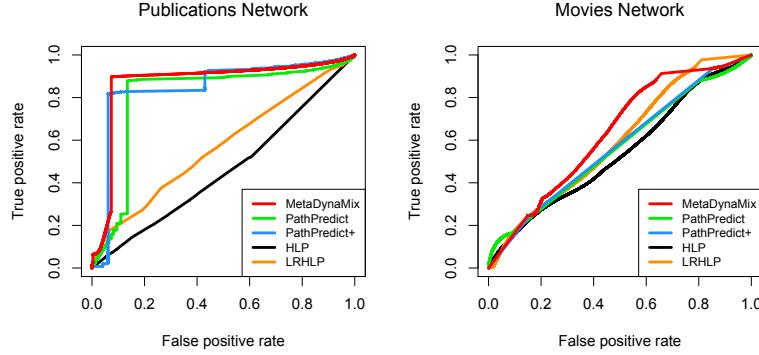


Fig. 5. The ROC curves for different methods and datasets.

4.2 Results and Findings

Link Prediction Accuracy. In this part we compare the prediction accuracy of different methods. The results shown in Figure 5 are based on setting the number of time intervals t to 7 for dynamic methods and 3 intervals for PathPredict. Table 2 shows more details considering different intervals. These results indicate the statistically significant improvement provided by the proposed MetaDynaMix prediction method compared to the baselines. The authors in [21, 38] showed that latent features are more predictive compared to unsupervised scoring techniques such as Katz, or Adamic. In our experiments we observed that combining latent with meta path-based features (MetaDynaMix) can increase prediction accuracy. However, if latent features learn similar structure as topological features do, then mixing them may not be beneficial. In such cases feature engineering techniques could be applied.

We also observe that PathPredict+ performs better than LRHLP in predicting links for the publications network but LRHLP predicts more accurate on the movies network. This implies that unlike the publications network, our meta path-based features for the movies network are not as predictive as latent features. However, in both cases combining the two set of features give better performance than either model individually.

Ebrahim ► Why there is this difference between accuracy on the two dataset? Sparsity? Sampling and neighborhood size? ◀

Significance of Performance Improvement. McNemar’s test, also called within-subjects χ^2 test, is used to compare statistically significant difference between the accuracy of two predictive models based on the contingency table of their predictions. The null hypothesis indicates that the performances of the two models are equal. We compare MetaDynaMix with the other four baselines and the test results show a p -value < 0.0001 for all cases and hence we reject the null hypothesis.

The Effect of Time Intervals. We set the number of time intervals t to 3, 5, and 7 and assess its impact on prediction performance. As presented in Table 2,

Table 2. Relationship prediction accuracy comparison. Bold values are determined to be statistically significant compared to the baselines based on McNemar’s test.

Method	Metric	Publications Network			Movies Network		
		<i>t</i> =3	<i>t</i> =5	<i>t</i> =7	<i>t</i> =3	<i>t</i> =5	<i>t</i> =7
PathPredict	ROC	0.78	—	—	0.56	—	—
	ACC	0.55	—	—	0.54	—	—
PathPredict+	ROC	0.78	0.80	0.83	0.56	0.57	0.57
	ACC	0.55	0.58	0.60	0.54	0.54	0.55
HLP	ROC	0.42	0.43	0.46	0.51	0.53	0.54
	ACC	0.50	0.50	0.50	0.51	0.52	0.53
LRHLP	ROC	0.49	0.50	0.52	0.52	0.56	0.59
	ACC	0.47	0.50	0.51	0.52	0.56	0.58
MetaDynaMix	ROC	0.85	0.87	0.87	0.57	0.59	0.63
	ACC	0.78	0.80	0.82	0.56	0.60	0.62

accuracy increases with the number of snapshots. The intuition is that shorter time interval results in less changes in the graph and thus leads to more reliable predictions. For example considering a meta path $A-P-V-P-A$, with smaller number of intervals, i.e., longer time intervals, we have more distinct authors who have published in a venue in different years and thus more similar path count values. However, by considering more intervals fewer authors will have such relations and more diverse path counts can contribute to more accurate prediction for the next time interval.

5 Related Work

The problem of link prediction in static and homogeneous networks has been extensively studied in the past [17, 32, 18, 15, 1, 2], for which the probability of forming a link between two nodes is generally considered as a function of their topological similarity. However, such techniques cannot be directly applied to heterogeneous networks. A few works such as [28, 29] investigated the problem of link prediction in HINs. Sun et al. [28] showed that *PathPredict* outperforms traditional link prediction approaches that use topological features defined on homogeneous networks such as common neighbors [22], preferential attachment [22], Jaccard’s coefficient [17], and Katz β [13]. Different from the original link prediction problem, Sun et al. [29] studied the problem of predicting the time of relationship building in HINs. These works, however, do not consider the dynamism of networks and overlook the potential benefits of analyzing a HIN as a sequence of network snapshots.

Research works on static latent space inference of networks [24, 21, 36, 23, 35] have assumed that the latent positions of nodes are fixed, and only few graph embedding methods [9, 7, 38] have considered dynamic networks. Dunlavy et al. [7] developed a tensor-based latent space modeling to predict temporal links. Zhu et al. [38] added a temporal-smoothing regularization term to a non-negative

matrix factorization objective to penalizes abrupt large changes in the latent positions and optimized it using a block-coordinate gradient descent algorithm. These works, however, do not consider heterogeneity of network structure.

6 Conclusions and Future Work

We studied the problem of relationship prediction in DHINs and proposed a supervised learning framework based on a combined set of latent and topological meta path-based features. Our results show that the proposed technique significantly improves prediction accuracy compared to the baseline methods. In this work we did not evaluate the running time and efficiency of our approach. Since our major computational bottleneck is calculating meta path-based measures, such as path count, we would like to investigate approximation techniques to make the prediction scalable. Furthermore we are interested in enhancing the matrix factorization technique based on a loss function that does not require full topological features matrix. In addition to model improvement, another interesting direction is to investigate the effectiveness of our proposed approach in other applications, such as predicting interests of users in a social media, that can be formulated as a link/relationship prediction problem.

References

1. Al Hasan, M., Chaoji, V., Salem, S., Zaki, M.: Link prediction using supervised learning. In: SDM06: workshop on link analysis, counter-terrorism and security (2006)
2. Al Hasan, M., Zaki, M.J.: A survey of link prediction in social networks. In: Social network data analytics, pp. 243–275. Springer (2011)
3. Cantador, I., Brusilovsky, P., Kuflik, T.: 2nd workshop on information heterogeneity and fusion in recommender systems (hetrec 2011). In: Proceedings of the 5th ACM conference on Recommender systems. RecSys 2011, ACM, New York, NY, USA (2011), <http://www.grouplens.org>
4. Chen, H., Li, X., Huang, Z.: Link prediction approach to collaborative filtering. In: Digital Libraries, 2005. JCDL'05. Proceedings of the 5th ACM/IEEE-CS Joint Conference on. pp. 141–142. IEEE (2005)
5. Davis, J., Goadrich, M.: The relationship between precision-recall and roc curves. In: Proceedings of the 23rd international conference on Machine learning. pp. 233–240. ACM (2006)
6. Denny, J.C.: Mining electronic health records in the genomics era. PLoS computational biology **8**(12), e1002823 (2012)
7. Dunlavy, D.M., Kolda, T.G., Acar, E.: Temporal link prediction using matrix and tensor factorizations. ACM Transactions on Knowledge Discovery from Data (TKDD) **5**(2), 10 (2011)
8. Fan, R.E., Chang, K.W., Hsieh, C.J., Wang, X.R., Lin, C.J.: Liblinear: A library for large linear classification. Journal of machine learning research **9**(Aug), 1871–1874 (2008), <https://github.com/cjlin1/liblinear>
9. Fu, W., Song, L., Xing, E.P.: Dynamic mixed membership blockmodel for evolving networks. In: Proceedings of the 26th annual international conference on machine learning. pp. 329–336. ACM (2009)

10. Gallagher, B., Tong, H., Eliassi-Rad, T., Faloutsos, C.: Using ghost edges for classification in sparsely labeled networks. In: Proceedings of the 14th ACM SIGKDD international conference on Knowledge discovery and data mining. pp. 256–264. ACM (2008)
11. Guimerà, R., Sales-Pardo, M.: Missing and spurious interactions and the reconstruction of complex networks. *Proceedings of the National Academy of Sciences* **106**(52), 22073–22078 (2009)
12. Guy, I.: Social recommender systems. In: *Recommender Systems Handbook*, pp. 511–543. Springer (2015)
13. Katz, L.: A new status index derived from sociometric analysis. *Psychometrika* **18**(1), 39–43 (1953)
14. Lei, C., Ruan, J.: A novel link prediction algorithm for reconstructing protein–protein interaction networks by topological similarity. *Bioinformatics* **29**(3), 355–364 (2012)
15. Leroy, V., Cambazoglu, B.B., Bonchi, F.: Cold start link prediction. In: Proceedings of the 16th ACM SIGKDD international conference on Knowledge discovery and data mining. pp. 393–402. ACM (2010)
16. Li, X., Chen, H.: Recommendation as link prediction in bipartite graphs: A graph kernel-based machine learning approach. *Decision Support Systems* **54**(2), 880–890 (2013)
17. Liben-Nowell, D., Kleinberg, J.: The link-prediction problem for social networks. *Journal of the American society for information science and technology* **58**(7), 1019–1031 (2007)
18. Lichtenwalter, R.N., Lussier, J.T., Chawla, N.V.: New perspectives and methods in link prediction. In: Proceedings of the 16th ACM SIGKDD international conference on Knowledge discovery and data mining. pp. 243–252. ACM (2010)
19. Lü, L., Medo, M., Yeung, C.H., Zhang, Y.C., Zhang, Z.K., Zhou, T.: Recommender systems. *Physics Reports* **519**(1), 1–49 (2012)
20. McNemar, Q.: Note on the sampling error of the difference between correlated proportions or percentages. *Psychometrika* **12**(2), 153–157 (1947)
21. Menon, A.K., Elkan, C.: Link prediction via matrix factorization. In: Joint european conference on machine learning and knowledge discovery in databases. pp. 437–452. Springer (2011)
22. Newman, M.E.: Clustering and preferential attachment in growing networks. *Physical review E* **64**(2), 025102 (2001)
23. Qi, G.J., Aggarwal, C.C., Huang, T.: Link prediction across networks by biased cross-network sampling. In: *Data Engineering (ICDE), 2013 IEEE 29th International Conference on*. pp. 793–804. IEEE (2013)
24. Sarkar, P., Moore, A.W.: Dynamic social network analysis using latent space models. *ACM SIGKDD Explorations Newsletter* **7**(2), 31–40 (2005)
25. Shi, C., Kong, X., Huang, Y., Philip, S.Y., Wu, B.: Hetesim: A general framework for relevance measure in heterogeneous networks. *IEEE Trans. Knowl. Data Eng.* **26**(10), 2479–2492 (2014)
26. Shi, C., Li, Y., Zhang, J., Sun, Y., Philip, S.Y.: A survey of heterogeneous information network analysis. *IEEE Transactions on Knowledge and Data Engineering* **29**(1), 17–37 (2017)
27. Song, H.H., Cho, T.W., Dave, V., Zhang, Y., Qiu, L.: Scalable proximity estimation and link prediction in online social networks. In: Proceedings of the 9th ACM SIGCOMM conference on Internet measurement. pp. 322–335. ACM (2009)

28. Sun, Y., Barber, R., Gupta, M., Aggarwal, C.C., Han, J.: Co-author relationship prediction in heterogeneous bibliographic networks. In: Proceedings of the 2011 International Conference on Advances in Social Networks Analysis and Mining. pp. 121–128. ASONAM ’11, IEEE Computer Society (2011)
29. Sun, Y., Han, J., Aggarwal, C.C., Chawla, N.V.: When will it happen?: Relationship prediction in heterogeneous information networks. In: Proceedings of the Fifth ACM International Conference on Web Search and Data Mining. pp. 663–672. WSDM ’12, ACM, New York, NY, USA (2012)
30. Sun, Y., Han, J., Yan, X., Yu, P.S., Wu, T.: Pathsim: Meta path-based top-k similarity search in heterogeneous information networks. In: Proceedings of the VLDB Endowment 4 (11). pp. 992–1003. VLDB Endowment (2011)
31. Tang, J., Zhang, J., Yao, L., Li, J., Zhang, L., Su, Z.: Arnetminer: Extraction and mining of academic social networks. In: KDD’08. pp. 990–998 (2008), <https://aminer.org/citation>
32. Wang, C., Satuluri, V., Parthasarathy, S.: Local probabilistic models for link prediction. In: icdm. pp. 322–331. IEEE (2007)
33. Wang, C., Song, Y., El-Kishky, A., Roth, D., Zhang, M., Han, J.: Incorporating world knowledge to document clustering via heterogeneous information networks. In: Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining. pp. 1215–1224. ACM (2015)
34. Wang, C., Song, Y., Li, H., Zhang, M., Han, J.: Text classification with heterogeneous information network kernels. In: AAAI. pp. 2130–2136 (2016)
35. Ye, J., Cheng, H., Zhu, Z., Chen, M.: Predicting positive and negative links in signed social networks by transfer learning. In: Proceedings of the 22nd international conference on World Wide Web. pp. 1477–1488. ACM (2013)
36. Yin, J., Ho, Q., Xing, E.P.: A scalable approach to probabilistic latent space inference of large-scale networks. In: Advances in neural information processing systems. pp. 422–430 (2013)
37. Zhang, J., Wang, C., Wang, J., Yu, J.X.: Inferring continuous dynamic social influence and personal preference for temporal behavior prediction. Proceedings of the VLDB Endowment 8(3), 269–280 (2014)
38. Zhu, L., Guo, D., Yin, J., Steeg, G.V., Galstyan, A.: Scalable temporal latent space inference for link prediction in dynamic social networks. IEEE Transactions on Knowledge and Data Engineering (TKDE) 28(10), 2765–2777 (2016), <https://github.com/linhongseba/Temporal-Network-Embedding>