

A Survey of Entity Similarity Measures on Heterogeneous Information Network

ZHU Zichen

Department of Computer Science
University of Hong Kong
Pokfulam Road, Hong Kong
zczhu2@cs.hku.hk

MA Chenhao

Department of Computer Science
University of Hong Kong
Pokfulam Road, Hong Kong
chma2@cs.hku.hk

ABSTRACT

According to recent studies, *heterogeneous information networks (HINs)* consisting of multiple types of entities and relations has shown its power in many disciplines, such as, computer science, social science, physics and so on. More and more researchers have noticed the importance of HIN analysis and many novel data mining tasks have been exploited in such networks, such as similarity search, clustering, and classification. Among those tasks, similarity measure on HINs, which is mainly to evaluate the similarity of entities, is the basis of many data mining tasks, such as clustering and classification. In this paper, we provide a survey of similarity measures on heterogeneous information networks. We will introduce basic concepts of heterogeneous information networks analysis, examine the recent developments on similarity measures and make a general evaluation of those similarity metrics.

1. INTRODUCTION

Recently, researchers use heterogeneous information networks (HINs) to model real world relationships in many applications, especially for real systems containing multi-typed interacting components. For instance, in bibliographic database, like DBLP¹ [10], papers are connected together via authors, venues and terms; and in Instagram², photos or videos are linked together via users, locations, hashtags and comments. Compared to widely-used homogeneous information networks [4, 5], which are extracted from real interacting systems by simply ignoring the heterogeneity of objects and links or only considering one type of relations among one type of objects, the heterogeneous information network can effectively fuse more information and contain rich and specific semantics in nodes and links [9].

Because of HIN's property of rich semantics and information, since the concept of heterogeneous information net-

¹<http://dblp.uni-trier.de>

²<http://instagram.com>

work and meta path proposed in 2009 [12] and 2011 [11], respectively, more and more researchers have noticed the importance of heterogeneous information network analysis and many novel data mining tasks have been developed in such networks, such as similarity search [8, 11], clustering [13]. In other words, HIN analysis has become a hot topic rapidly in the fields of data mining, database and information retrieval, involving similarity measure, clustering, classification, link prediction, ranking, recommendation and information fusion on HINs [9].

Among those analysis tasks, similarity measure is the fundamental problem of network analysis, because most high level tasks, such as clustering and classification, need to evaluate the similarity of objects or relations. What's more, most of the state-of-the-art homogeneous network similarity methods, which generally assume the networks don't carry semantics, do not generalize well in HINs, due to HINs' rich semantics. Thus, some similarity metrics based on meta-path, which represents semantics in HINs, have been designed to evaluate the similarity between entities or relations in HINs, such as PathSim [11] and RelSim [16]. **Maybe some figure of HIN and an example...**

In this paper, we attempts to clearly introduce basic concepts in heterogeneous network analysis and make a comprehensive investigation on contemporary research developments of similarity metrics on HINs. Then, we tentatively provide a general evaluation over some of these state-of-the-art algorithms to give some reference on similarity measure choosing in high level network analysis tasks, such as clustering and classification.

The following part is organized as follows. Section 2 introduces the basic concepts and examples about HIN. Section 3 presents recent designed similarity measures on HIN. Experiments and evaluation are conducted in Section 4. Finally, Section 5 summarizes and concludes this paper.

2. BASIC CONCEPTS AND DEFINITIONS

In this section, we introduce some basic concepts about HIN and give some HIN examples. We first define the information network and heterogeneous information network.

An information network represents an abstraction of the real world, focusing on the entities and the relations among these entities.

Definition 1. Information Network [10, 14]. An information network is defined as a directed graph $G = (V, E)$ with an entity type mapping function $\varphi : V \rightarrow \mathcal{A}$ and a relation type mapping function $\psi : E \rightarrow \mathcal{R}$. Each entity $v \in V$ belongs to one particular entity type in the entity type set $\mathcal{A} : \varphi(v) \in \mathcal{A}$, and each relation $e \in E$ belongs to a particular relation type in the relation type set $\mathcal{R} : \psi(e) \in \mathcal{R}$. If two relations belongs to the same relation type, the two links share the same starting entity type as well as the ending entity type.

Based on the definition of information network, we derive the definitions of heterogeneous/homogeneous information network.

Definition 2. Heterogeneous/homogeneous information network. The information network is called heterogeneous information network if the type of entities $|\mathcal{A}| > 1$ or the types of relations $|\mathcal{R}| > 1$; otherwise, it is a homogeneous information network.

For better understanding the entity types and relation types in a complex heterogeneous information network, the network schema provides a high-level description of a given heterogeneous information network.

Definition 3. Network schema [10, 14]. The network schema, denoted as $T_G = (\mathcal{A}, \mathcal{R})$, is a meta template for an information network $G = (V, E)$ with the entity type mapping $\varphi : V \rightarrow \mathcal{A}$ and the relation type mapping $\psi : E \rightarrow \mathcal{R}$, which is a directed graph defined over entity types \mathcal{A} , with edges as relations from \mathcal{R} .

The network schema of a HIN specifies type constraints on the sets of entities and relationships among those entities. An information network following a network schema is called a **network instance** of the network schema. For a relation type R connecting entity type S to entity type T , i.e., $S \xrightarrow{R} T$, S and T are the **source entity type** and **target entity type** of relation type R , which can be denoted as $R.S$ and $R.T$, respectively. The inverse relation R^{-1} holds naturally for $T \xrightarrow{R^{-1}} S$. Generally, R is not equal to R^{-1} , unless R is symmetric.

Another important concept, meta-path, is proposed to systematically define relations between entities at the schema level.

Definition 4. Meta-path [11]. A meta-path \mathcal{P} is a path defined on a schema $S = (\mathcal{A}, \mathcal{R})$, and is denoted in the form of $A_1 \xrightarrow{R_1} A_2 \xrightarrow{R_2} \dots \xrightarrow{R_l} A_{l+1}$, which defines a composite relation $R = R_1 \circ R_2 \circ \dots \circ R_l$ between entities A_1, A_2, \dots, A_{l+1} , where \circ denotes the composition operator on relations.

For simplicity, we can also use entity types to denote the meta-path if there are no multiple relation types between

the same pair of entity types: $\mathcal{P} = (A_1 A_2 \dots A_{l+1})$. We say a concrete path $p = (a_1 a_2 \dots a_{l+1})$ between entities a_1 and a_{l+1} in network G is a **path instance** of the relevance path \mathcal{P} , if $\forall a_i, \varphi(a_i) = A_i$ and $\forall e_i = \langle a_i, a_{i+1} \rangle, \psi(e_i) = R_i$ in \mathcal{P} . It can be denoted as $p \in \mathcal{P}$. A meta-path \mathcal{P} is a **symmetric path**, if the relation R defined by it is symmetric (i.e., \mathcal{P} is equal to \mathcal{P}^{-1}), such as APA . Two meta-paths $\mathcal{P}_1 = (A_1 A_2 \dots A_l)$ and $\mathcal{P}_2 = (B_1 B_2 \dots B_k)$ are **concatenable** if and only if A_l is equal to B_1 , and the concatenated path is written as $\mathcal{P} = (\mathcal{P}_1 \mathcal{P}_2)$, which equals to $(A_1 A_2 \dots A_l B_2 \dots B_k)$.

The commuting matrix is defined in [11] to compute the frequencies of all the paths related to a meta-path.

Definition 5. Commuting matrix. Given a network $G = (V, E)$ and its network schema T_G , a commuting matrix $M_{\mathcal{P}}$ for a meta-path $\mathcal{P} = (A_1, A_2, \dots, A_{l+1})$ is defined as $M_{\mathcal{P}} = W_{A_1 A_2} W_{A_2 A_3} \dots W_{A_l A_{l+1}}$, where $W_{A_i A_j}$ is the adjacency matrix between types A_i and A_j . $M_{\mathcal{P}}(i, j)$ represents the number of path instances between entities v_i and v_j , where $\varphi(v_i) = A_1$ and $\varphi(v_j) = A_{l+1}$, under meta-path \mathcal{P} .

Example 1. Example and figure about HIN, schema, meta-path

3. SIMILARITY MEASURES

Similarity measure is to evaluate the similarity of entities. It is the basis of many data mining tasks, such as classification, clustering, and recommendation system. Similarity measure has been well studied on different kinds of data types for a long time. These studies can be roughly categorized into two types: **feature based approaches** and **link based approaches**. The feature based approaches measure the similarity of entities based on their feature/attribute values, such as cosine similarity, Jaccard similarity and Euclidean distance. The link based approaches measure the similarity of entities based on their link structures in a network. For instance, Personalized PageRank [3] evaluates the probability starting from a source object to a target object by randomly walking with restart, and SimRank[2] evaluates the similarity of two objects by their neighbors' similarities.

In recent years, similarity measures on heterogeneous information networks begin to be noticed by more and more researchers. Apart from the structure similarity addressed by most homogeneous similarity metrics, similarity metrics on HIN also need to take the meta-path connecting these two objects into account. As we know, there are different meta-paths connecting two objects, and these meta paths contain different semantics meanings, which may lead to different similarities. So, the similarity measure on HIN is meta-path constraint [9]. **Example?** We present the recent state-of-the-art similarity metrics in the following part.

3.1 PathSim

PathSim [11] is the first meta-path based similarity measure to evaluate the similarity of same-typed entities based on symmetric meta-paths.

Definition 6. PathSim: Given a symmetric meta-path \mathcal{P} , PathSim between two entities u and v of the same entity type is:

$$\begin{aligned} \text{PathSim}(u, v) &= \frac{2 \times |\{p_{u \rightsquigarrow v} \in \mathcal{P}\}|}{|\{p_{u \rightsquigarrow u} \in \mathcal{P}\}| + |\{p_{v \rightsquigarrow v} \in \mathcal{P}\}|} \\ &= \frac{2 \cdot M_{\mathcal{P}}(u, v)}{M_{\mathcal{P}}(u, u) + M_{\mathcal{P}}(v, v)} \end{aligned} \quad (1)$$

3.2 Distant Meta-Path Similarity

Distant meta-path similarity [15] is designed to evaluate text-based meta-path similarity between two distant (relatively isolated) entities. Here, distant entities means those two entities can not connected by the given meta-path.

Definition 7. Distant meta-path similarity. The distant meta-path similarity between an entity pair describes the proximity of the pair's neighborhood entities. Neighborhood entities are defined as the entities kinked via meta-path(s) to the pair. Let $\{M_{\mathcal{P}}(u, w)\}_{w=1}^N$ denotes the meta-path instances between entity u and its neighborhood entities. The distant meta-path similarity between u and v is the decided by the proximity of $\{M_{\mathcal{P}}(u, w)\}_{w=1}^N$ and $\{M_{\mathcal{P}}(v, w)\}_{w=1}^N$. Entities u and v are called as distant neighbors to each other.

There are 53 similarity metrics, i.e., metrics to measure the similarity between $\{M_{\mathcal{P}}(u, w)\}_{w=1}^N$ and $\{M_{\mathcal{P}}(v, w)\}_{w=1}^N$, tested in [15] to find the best way to define a distant meta-path similarity. Experimental results in [15] show cosine similarity is consistent good for general use. Thus, we present cosine similarity based distant meta-path similarity here.

$$\begin{aligned} \text{DistantSim}(u, v) &= \frac{\sum_{m=1}^M \sum_{w=1}^N M_{\mathcal{P}_m}(u, w) M_{\mathcal{P}_m}(v, w)}{\sqrt{\sum_{m=1}^M \sum_{w=1}^N M_{\mathcal{P}_m}(u, w)^2} \sqrt{\sum_{m=1}^M \sum_{w=1}^N M_{\mathcal{P}_m}(v, w)^2}} \end{aligned} \quad (2)$$

3.3 HeteSim

The similarity of objects with different are needed in many applications, such as recommendation system [1] and medicine annotation analysis [6]. Thus, HeteSim [7] is proposed for evaluating the similarity of entities with different types.

Before giving the definition of HeteSim, we first introduce the decomposition of meta-path.

Definition 8. Decomposition of meta-path. An arbitrary meta-path $\mathcal{P} = (A_1, A_2, \dots, A_{l+1})$ can be decomposed into two equal-length path \mathcal{P}_L and \mathcal{P}_R , i.e., $\mathcal{P} = \mathcal{P}_L \mathcal{P}_R$, where $\mathcal{P}_L = (A_1, A_2, \dots, A_{\lfloor l/2 \rfloor}, B)$ and $\mathcal{P}_R = (B, A_{\lfloor l/2 \rfloor + 1}, \dots, A_{l+1})$. If l is even, $B = A_{\lfloor l/2 \rfloor}$. Otherwise, B is the middle type entity E between the atomic relation $A_{\lfloor l/2 \rfloor}$ and $A_{\lfloor l/2 \rfloor + 1}$. The new path becomes $\mathcal{P}' = (A_1, \dots, E, \dots, A_{l+1})$, so B is also the middle item of \mathcal{P}' .

Obviously, for a symmetric path $\mathcal{P} = \mathcal{P}_L \mathcal{P}_R$, \mathcal{P}_R^{-1} is equal to \mathcal{P}_L . After transforming the original meta-path, when its length is odd, the definition of HeteSim can be expressed:

Definition 9. HeteSim. Given a relevance path $\mathcal{P} = (A_1, A_2, \dots, A_{l+1})$, the HeteSim score between two entities u and v ($u \in A_1, v \in A_{l+1}$) is:

$$\text{HeteSim}(u, v) = \frac{\sum_{w=1}^N M_{\mathcal{P}_L}(u, w) \cdot M_{\mathcal{P}_R^{-1}}(v, w)}{\sqrt{\sum_{w=1}^N M_{\mathcal{P}_L}(u, w)^2} \sqrt{\sum_{w=1}^N M_{\mathcal{P}_R^{-1}}(v, w)^2}} \quad (3)$$

3.4 RelSim

RelSim [16] is a meta-path based relation similarity measure. It measures the similarity between two relation instances based on the latent semantic relation (LSR): two relation instances are more similar when sharing more important (heavily weighted) meta-paths.

Definition 10. RelSim. Given an LSR (latent semantic relation), denoted as $\{w_m, \mathcal{P}_m\}_{m=1}^M$, RelSim between two relation instances $r = \langle v^{(1)}, v^{(2)} \rangle$ and $r' = \langle v^{(1)'}, v^{(2)'} \rangle$ is defined as:

$$\text{RelSim}(r, r') = \frac{2 \times \sum_m w_m \min(x_m, x'_m)}{\sum_m w_m x_m + \sum_m w_m x'_m} \quad (4)$$

where x_m is the number of path instances between $v^{(1)}$ and $v^{(2)}$ in relation r following meta-path \mathcal{P}_m , and x'_m is the number of path instances between $v^{(1)'}$ and $v^{(2)'}$ in relation r' following meta-path \mathcal{P}_m . We use a vector $\mathbf{x} = [x_1, \dots, x_m, \dots, x_M]$ to characterize a relation instance r , and a vector $\mathbf{w} = [w_1, \dots, w_m, \dots, w_M]$ to denote the corresponding weights. M is the number of meta-paths.

4. REFERENCES

- [1] M. Jamali and L. Lakshmanan. Heteromf: recommendation in heterogeneous information networks using context dependent factor models. In *Proceedings of the 22nd international conference on World Wide Web*, pages 643–654. ACM, 2013.
- [2] G. Jeh and J. Widom. Simrank: a measure of structural-context similarity. In *Proceedings of the eighth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 538–543. ACM, 2002.
- [3] G. Jeh and J. Widom. Scaling personalized web search. In *Proceedings of the 12th international conference on World Wide Web*, pages 271–279. ACM, 2003.
- [4] V. Leroy, B. B. Cambazoglu, and F. Bonchi. Cold start link prediction. In *Proceedings of the 16th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 393–402. ACM, 2010.
- [5] R. N. Lichtenwalter, J. T. Lussier, and N. V. Chawla. New perspectives and methods in link prediction. In *Proceedings of the 16th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 243–252. ACM, 2010.
- [6] G. Palma, M.-E. Vidal, E. Haag, L. Raschid, and A. Thor. Measuring relatedness between scientific

- entities in annotation datasets. In *Proceedings of the International Conference on Bioinformatics, Computational Biology and Biomedical Informatics*, page 367. ACM, 2013.
- [7] C. Shi, X. Kong, Y. Huang, S. Y. Philip, and B. Wu. Hetesim: A general framework for relevance measure in heterogeneous networks. *IEEE Transactions on Knowledge and Data Engineering*, 26(10):2479–2492, 2014.
- [8] C. Shi, X. Kong, P. S. Yu, S. Xie, and B. Wu. Relevance search in heterogeneous networks. In *Proceedings of the 15th International Conference on Extending Database Technology*, pages 180–191. ACM, 2012.
- [9] C. Shi, Y. Li, J. Zhang, Y. Sun, and S. Y. Philip. A survey of heterogeneous information network analysis. *IEEE Transactions on Knowledge and Data Engineering*, 29(1):17–37, 2017.
- [10] Y. Sun and J. Han. Mining heterogeneous information networks: a structural analysis approach. *ACM SIGKDD Explorations Newsletter*, 14(2):20–28, 2013.
- [11] Y. Sun, J. Han, X. Yan, P. S. Yu, and T. Wu. Pathsim: Meta path-based top-k similarity search in heterogeneous information networks. *Proceedings of the VLDB Endowment*, 4(11):992–1003, 2011.
- [12] Y. Sun, J. Han, P. Zhao, Z. Yin, H. Cheng, and T. Wu. Rankclus: integrating clustering with ranking for heterogeneous information network analysis. In *Proceedings of the 12th International Conference on Extending Database Technology: Advances in Database Technology*, pages 565–576. ACM, 2009.
- [13] Y. Sun, B. Norick, J. Han, X. Yan, P. S. Yu, and X. Yu. Pathselclus: Integrating meta-path selection with user-guided object clustering in heterogeneous information networks. *ACM Transactions on Knowledge Discovery from Data (TKDD)*, 7(3):11, 2013.
- [14] Y. Sun, Y. Yu, and J. Han. Ranking-based clustering of heterogeneous information networks with star network schema. In *Proceedings of the 15th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 797–806. ACM, 2009.
- [15] C. Wang, Y. Song, H. Li, Y. Sun, M. Zhang, and J. Han. Distant meta-path similarities for text-based heterogeneous information networks. 2017.
- [16] C. Wang, Y. Sun, Y. Song, J. Han, Y. Song, L. Wang, and M. Zhang. Relsim: relation similarity search in schema-rich heterogeneous information networks. In *Proceedings of the 2016 SIAM International Conference on Data Mining*, pages 621–629. SIAM, 2016.