

Date of publication xxxx 00, 0000, date of current version xxxx 00, 0000.

Digital Object Identifier 10.1109/ACCESS.2017.DOI

## Learning Heterogeneous Network Embedding from Text and Links

YUNFEI LONG<sup>1\*</sup>, RONG XIANG<sup>1\*</sup>, QIN LU<sup>1</sup>, DAN XIONG<sup>1</sup>, CHU-REN HUANG<sup>2</sup>, CHENLIN BI<sup>3</sup>, and MINGLEI LI<sup>1</sup>

<sup>\*</sup>Both authors contribute to this paper equally

<sup>1</sup>Department of Computing, The Hong Kong Polytechnic University, Hong Kong (e-mail: yunfei.long@connect.polyu.hk, csrxiang, csuqin, csdxiong@comp.polyu.edu.hk, minglei.li@connect.polyu.hk)

<sup>2</sup>Department of Chinese and Bilingual Studies, The Hong Kong Polytechnic University, Hong Kong. (e-mail: churen.huang@polyu.edu.hk)

<sup>3</sup>Advanced Micro Devices (Shanghai), Shanghai, China

Corresponding author: Minglei Li (e-mail: minglei.li@connect.polyu.hk).

The work is partially supported by the research grants from Hong Kong Polytechnic University (PolyU RTVU) , GRF grant (CERG PolyU 15211/14E, PolyU 152006/16E).

The work is partially supported by the research grants from Hong Kong Polytechnic University (PolyU RTVU) , GRF grant (CERG PolyU 15211/14E, PolyU 152006/16E).

**ABSTRACT** Finding methods to represent multiple types of nodes in heterogeneous networks is both challenging and rewarding as there is much less work in this area compared to that of homogeneous networks. In this paper, we propose a novel approach to learn node embedding for heterogeneous networks through a joint learning framework of both network links and text associated with nodes. A novel attention mechanism is also used to make good use of text extended through links to obtain much larger network context. Link embedding is first learned through a random walk based method to process multiple types of links. Text embedding is separately learned at both sentence level and document level to capture salient semantic information more comprehensively. Then, both types of embeddings are jointly fed into a hierarchical neural network model to learn node representation through mutual enhancement. The attention mechanism follows linked edges to obtain context of adjacent nodes to extend context for node representation. The evaluation on a link prediction task in a heterogeneous network dataset shows that our method outperforms the current state-of-the-art method by 2.5% to 5.0% in AUC values with *p-value* less than  $10^{-9}$ , indicating very significant improvement.

**INDEX TERMS** Network embedding, Heterogeneous network, Attention mechanism, Text processing.

### I. INTRODUCTION

Nowadays, networks are ubiquitous and many applications need to mine information within these networks. Network applications include DNA networks in biology [1]–[3], friendship/follower networks in social sciences [4]–[7], Internet of Things [8], [9], and word co-occurrence networks in linguistics [10]–[12], etc. With the wide use of networks in modeling and applications, network embedding, a method to use fixed dimension vectors to represent nodes in a network, becomes a hot research topic [6], [13]–[15] [16]–[19].

According to the unitarity of node characters or features, there are two types of networks: homogeneous networks and heterogeneous networks. A homogeneous network only consists of one type of nodes and behaviors. A heterogeneous network contains different types of nodes. In many heterogeneous networks, especially user related networks, users (the subjects of behaviors) and products (the objects of behaviors) are commonly regarded as two different types of nodes.

Previous tasks in network embedding mostly focus on homogeneous networks. Different methods are proposed such as matrix factorization [20], random walk [14], and neural networks [21]. For homogeneous networks, the use of text information as context data is also common. Due to the complexity of this issues, few studies work on network embedding of heterogeneous networks. In heterogeneous networks, we consider it particularly important to leverage on both link information and other types of information such as text for product description and comments written by friends in social networks. Link information is considered naturally structured as a graph. Text information, on the other hand, is often free-structured or semi-structured [22].

To integrate link structure and text in the same network, two main issues need to be addressed: the first issue is how to learn node representation by integrating link information and text content coherently; the second issue is how to distinguish different types of nodes in the representation framework. This

Long, Yunfei, Rong Xiang, Qin Lu, Dan Xiong, Chu-Ren Huang, Chenlin Bi, Minglei Li. 2018. Learning Heterogeneous Network Embedding from Text and Links. *IEEE Access*. 6.1.1-11.

Published Online 2018-10-02.

To Appear: 2018-12

DOI: [10.1109/ACCESS.2018.2873044](https://doi.org/10.1109/ACCESS.2018.2873044)

• Early Access:

<https://ieeexplore.ieee.org/document/8478654/>

• IEEE Access<sup>®</sup> is a multidisciplinary, applications-oriented, all-electronic archival journal that continuously presents the results of original research or development across all of IEEE's fields of interest.

• JCR IF: 3.557 (2017)

• SJR Q1 in Computer Science (Misc.) and Engineering (Misc.)