



**University of  
Zurich<sup>UZH</sup>**

**Department of Informatics**

**MSc Thesis**

# **Application of Higher-Order Network Models to Representation Learning in Sequential Data**

**Michael Markus Studer**  
of Maschwanden ZH, Switzerland

Student-ID: 92-911-544  
[mstuder72@gmail.com](mailto:mstuder72@gmail.com)

November 18, 2020

supervised by Prof. Dr. Ingo Scholtes

# Application of Higher-Order Network Models to Representation Learning in Sequential Data

Master Thesis Proposal

Michael Studer (92-911-544)

May 18<sup>th</sup>-November 11<sup>th</sup>, 2020

Data Analytics Group, University of Zürich

Supervised by Prof. Dr. Ingo Scholtes

## Background

Unsupervised methods to learn low-dimensional representations of complex data are important tools to address the problem of feature extraction in machine learning [1]. Applying this concept to machine learning in graphs, researchers have recently developed graph representation learning or network embedded techniques[4, 8], which have been successfully applied to problems such as link prediction, node classification, clustering, or network visualization[5]. The common of those techniques is to learn a mapping function that maps the nodes of a graph to points in a low-dimensional vector space such that geometric distances between nodes reflect the topology of the network. Existing algorithms that approach this problem commonly utilize an edge-centric view that focuses on the **topology of edges**. However, apart from information on the topology of edges, i.e. who is connected to whom, rich time-stamped or sequential data on networks increasingly provides us with higher-order information like, e.g. in which chronological order biological interactions are activated, how users navigate hyperlinks in information networks, or along which paths information propagates in social networks. Previous works in the hosting group have shown that such sequential data contain higher-order, non-Markovian patterns that question edge-centric network analysis - and thus - graph embedding techniques. The lack of graph representation learning techniques that account for higher-order patterns in sequential data on networks has recently been highlighted by researchers in machine learning and network science [5, 6].

## Scope of this Project

Addressing the gap summarized above, the goals of this Master thesis are (i) to empirically investigate how higher-order interaction patterns in sequential data on networks influence representation learning techniques, and (ii) to use the higher-order models for paths in networks developed in the hosting group to develop a representation learning algorithm that accounts both for the topology of (dyadic) links in a network as well as patterns in the chronological ordering of links in sequential data. The concrete objectives of the project are as follows:

### Objectives

**O1** First, the candidate should perform a **literature survey of representation learning in networks**, with a particular focus on vector-space embeddings that use models for paths or walks in networks. The candidate should further **build a collection of implementations of representation learning algorithms** that minimally comprises node2vec[3], Deepwalk[9] and HONEM[11]. For this collection the candidate can use existing and publicly available implementations (e.g. of node2vec<sup>a</sup> or DeepWalk<sup>b</sup>), while for HONEM a custom implementation based on the description of the method in [11] may be necessary.

**O2** Evaluating the methods from O1, in a second step the candidate should **develop a probabilistic generative model for paths or walks in a network with a known embedding in a high-dimensional vector space**. The model should generate paths or walks such that the network embedding translates to non-Markovian properties for the generated walks or paths, i.e. for a sequences of nodes  $v_0, \dots, v_l$  traversed by a path of length  $l$  the sequence should exhibit “memory” that can be explained based on the known position of traversed nodes/edges in the vector space. This data should then be used to compare the performance of the representation learning techniques identified in O1. The expected outcome of this objective is that HONEM outperforms methods that do not account for sequential patterns in paths on networks.

**O3** The final objective of the thesis is to **apply higher-order models for paths in networks to improve representation learning in sequential network data**. For this, the candidate should investigate how the higher- and multi-order modelling frameworks developed in [2, 7, 10, 12, 13] can be applied to develop an algorithm that generates vector-space embeddings of nodes (or higher-order nodes). The developed algorithm should be evaluated both in the synthetic paths generated in O2 as well as in real data sets, e.g. performing a cross-validation experiment based on a similarity-based link prediction analogous to [11].

<sup>a</sup><https://github.com/aditya-grover/node2vec>

<sup>b</sup><https://github.com/phanein/deepwalk>

The thesis should motivate the problem, describe the current state of research in the field, propose and implement a solution based on prior work on higher-order models, and evaluate it by means of a validation in empirical data. The code of the implementation should be documented in line with typical software engineering standards and this documentation can be included in the appendix of the thesis. The hosting group offers regular meetings in which we will discuss the current progress of the work and clarify potential questions or issues.

## Bibliography

- [1] Yoshua Bengio, Aaron Courville, and Pascal Vincent. Representation learning: A review and new perspectives. *IEEE transactions on pattern analysis and machine intelligence*, 35(8):1798–1828, 2013.
- [2] Christoph Gute, Giona Casiraghi, Frank Schweitzer, and Ingo Scholtes. Mogen: A generative multi-order model to predict variable length paths in networks. 2020.
- [3] Aditya Grover and Jure Leskovec. node2vec: Scalable feature learning for networks. In *Proceedings of the 22nd ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 855–864, 2016.
- [4] Will Hamilton, Zhitao Ying, and Jure Leskovec. Inductive representation learning on large graphs. In *Advances in Neural Information Processing Systems*, pages 1024–1034, 2017.
- [5] William L Hamilton, Rex Ying, and Jure Leskovec. Representation learning on graphs: Methods and applications. *arXiv preprint arXiv:1709.05584*, 2017.
- [6] Renaud Lambiotte, Martin Rosvall, and Ingo Scholtes. From networks to optimal higher-order models of complex systems. *Nature physics*, page 1, 2019.
- [7] Timothy LaRock, Vahan Nanumyan, Ingo Scholtes, Giona Casiraghi, Tina Eliassi-Rad, and Frank Schweitzer. Hypa: Efficient detection of path anomalies in time series data on networks, 2019.
- [8] Thanh Tam Nguyen and Chi Thang Duong. A comparison of network embedding approaches. Technical report, School of Computer and Communication Sciences, EPFL, Technical Report, 2018.
- [9] Bryan Perozzi, Rami Al-Rfou, and Steven Skiena. Deepwalk: online learning of social representations. In Sofus A. Macskassy, Claudia Perlich, Jure Leskovec, Wei Wang, and Rayid Ghani, editors, *The 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD ’14, New York, NY, USA - August 24 - 27, 2014*, pages 701–710. ACM, 2014.
- [10] Martin Rosvall, Alcides V Esquivel, Andrea Lancichinetti, Jevin D West, and Renaud Lambiotte. Memory in network flows and its effects on spreading dynamics and community detection. *Nature communications*, 5:4630, 2014.
- [11] Mandana Saebi, Giovanni Luca Ciampaglia, Lance M Kaplan, and Nitesh V Chawla. Honem: Network embedding using higher-order patterns in sequential data. *arXiv preprint arXiv:1908.05387*, 2019.
- [12] Ingo Scholtes. When is a network a network?: Multi-order graphical model selection in pathways and temporal networks. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 1037–1046. ACM, 2017.
- [13] Ingo Scholtes, Nicolas Wider, René Pfitzner, Antonios Garas, Claudio J Tesone, and Frank Schweitzer. Causality-driven slow-down and speed-up of diffusion in non-markovian temporal networks. *Nature communications*, 5:5024, 2014.

## Abstract

Representation learning provides crucial input for machine learning algorithms. Learning these representations instead of manually engineering them accelerates the development of ML applications. Various such methods exist for networks, but they mostly rely on first-order Markov chains to generate random walks to explore the network. However, higher-order Markov chains are often better at modeling real-world spreading processes, and this model change may also affect community detection.

We review well-established methods and explore different approaches to upgrade them from first-order to higher-orders. We experiment with multi-class classification and visualization tasks to compare the original and upgraded methods, using an illustrative synthetic grid and real data on social interactions.

The Python source code of the methods and experiments is publicly available<sup>1</sup>.

## Keywords

node embedding, representation learning, random walk, directed graph, node classification, higher-order network, higher-order Markov model

---

<sup>1</sup><https://github.com/HONEembeddings/MScThesis>

## Acknowledgments

I want to thank my supervisors Prof. Dr. Ingo Scholtes and Vincenzo Perri, for the opportunity to investigate this exciting subject. Moreover, I am very grateful for the excellent preparation, valuable feedback, and pleasant cooperation.

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Contributions and Structure . . . . .	2
<b>2</b>	<b>Background</b>	<b>3</b>
2.1	Networks . . . . .	3
2.1.1	Random Walks on Networks . . . . .	3
2.1.2	First-Order Networks (FON) . . . . .	4
2.1.3	Higher-Order Networks (HON) . . . . .	5
2.2	Background on Embeddings for FON . . . . .	6
2.2.1	Skip-Gram Model . . . . .	6
2.3	Related FON Embeddings . . . . .	8
2.3.1	DeepWalk and Node2vec . . . . .	9
2.3.2	LINE . . . . .	10
2.3.3	GraRep . . . . .	11
2.3.4	NetMF and InfiniteWalk . . . . .	12
2.4	Related HON Embeddings . . . . .	13
2.4.1	Definition of HONEM . . . . .	14
2.4.2	Theoretical Analysis of HONEM . . . . .	14
<b>3</b>	<b>Methodology</b>	<b>17</b>
3.1	Approaches for HON embeddings . . . . .	17
3.1.1	Higher-order random walks . . . . .	17
3.1.2	Beyond pairwise interactions . . . . .	18
3.2	Higher-order network embeddings examined . . . . .	19
3.2.1	HON DeepWalk . . . . .	19
3.2.2	HON Node2vec . . . . .	19
3.2.3	HON NetMF . . . . .	20
3.2.4	HON GraRep . . . . .	20
3.2.5	HON Experimental Embedding . . . . .	21
3.2.6	HONEM . . . . .	23
3.3	Synthetic network (HON Lattice 2D) . . . . .	23
3.3.1	Definition and Properties of the HON Lattice 2D . . . . .	25
3.3.2	Measuring the Stretch Effect . . . . .	26
3.3.3	Visualizing Embeddings of HON Lattice 2D . . . . .	27

<b>4 Experimental results</b>	<b>31</b>
4.1 Visualization for a synthetic dataset	31
4.1.1 Network	32
4.1.2 Baselines	32
4.1.3 Results	32
4.1.4 Explanations	33
4.2 Multiclass Classification for SocioPatterns	38
4.2.1 Methodology	38
4.2.2 Networks	39
4.2.3 Baselines	40
4.2.4 Results	41
4.2.5 Explanations	42
<b>5 Conclusion and Discussions</b>	<b>49</b>
5.1 Summary and Reflection	49
5.2 Threats to Validity and Limitations	50
5.3 Future Work	50
<b>A Results for the Classification task</b>	<b>53</b>
<b>B Additional Experimental results</b>	<b>59</b>
B.1 Probability Prediction	59
B.1.1 Network	59
B.1.2 Baselines	60
B.1.3 Best Case: Complete Information	60
B.1.4 Methodology	60
B.1.5 Results	61
B.2 HON Penalized Transition Embedding	61
<b>C Source Code</b>	<b>63</b>
C.1 Introduction and Usage	63
C.2 HigherOrderPathGenerator.py	65
C.3 SyntheticNetworks.py	72
C.4 Datasets.py	74
C.5 Embedding.py	75
C.6 Visualizations.py	90
C.7 Lattice2D_sim.ipynb	99
C.8 SocioPatterns.ipynb	108
C.9 Classification_sim.ipynb	115
C.10 Plot_ExpClass.ipynb	131
C.11 ProbabilityPrediction.ipynb	135
<b>D Process Documentation</b>	<b>145</b>

# Chapter 1

## Introduction

Networks are essential to model complex systems. However, machine learning algorithms usually cannot use networks directly as input, requiring us to describe the network’s nodes through representations (or features) instead. Unless this task is automated, it becomes a bottleneck in an AI application. Hence, network embeddings<sup>1</sup> gained significant interest from researchers and practitioners since 2013 and ongoing.

Most of the embedding methods encode for each node the strengths of the association with the other nodes. Of which LINE [33] is the simplest example because it encodes the edges (i.e., direct neighbors), considering weights if applicable.

Computing an embedding from a network represents a loss of information, which is occasionally desirable if an embedding can generalize (similar to function approximation with neuronal networks). E.g., in social networks, the information about edges is incomplete, and we want to predict unobserved relationships rather than the ones already known. Hence, DeepWalk [23] characterizes a node not only by its direct neighbors but also by its indirect neighbors and therefore uses *random walks* to explore the neighborhood.

Most embeddings assume a classical network, referred to as a first-order network (FON) in this thesis.

However, some networks vary over time. E.g., social interactions (or disease spreads) depend on the chronological sequence of contacts. Such data may be modeled by first extracting paths respecting time constraints and fitting a higher-order Markov model to these paths afterward. We denote networks with random walks generated by a higher-order Markov chain as higher-order networks<sup>2</sup> (HON). According to Rosvall et al. [27], switching to a (more accurate) HON model affects spreading dynamics and community detection, changing the embedding in turn.

To our knowledge, HONEM [28] is the only embedding for higher-order networks. To close this research gap, we analyze the existing embeddings for FON and identify approaches to update them to HON. We compare our HON adaptations of embeddings with their FON counterparts for (multiclass) node

---

<sup>1</sup>We denote by ‘network embeddings’ methods that learn vector-space representations for the individual nodes (and possibly tuples of nodes, e.g., edges) of a network.

<sup>2</sup>The term higher-order network (HON) refers exclusively to networks with higher-order Markov chains in this thesis. Other extensions of the FON are mentioned in [17]; see § 2.1.3.

classification and visualization tasks.

## 1.1 Contributions and Structure

In summary, our contributions are:

- We adapt the FON embeddings DeepWalk [23], Node2vec [10], GraRep [3], and NetMF[25] to HON.
- We compare the performances of these embeddings in node classification tasks using real data.
- We analyze HONEM theoretically and also in the above node classification task.
- We introduce a synthetic network, HON Lattice 2D, and analyze its properties. To distinguish FON and HON embeddings, we changed the HON random walks' dynamic by tweaking the second-order transition probabilities while keeping the FON ones unchanged.

This thesis is organized as follows:

- § 2 explains the relevant background — including HON random walks and the skip-gram model — and introduces the necessary notation. We provide an overview of the existing embeddings and describe some of them. Finally, we present two concepts for later reference: clustering of embedded nodes by parity for lattices and tweaking higher-order probabilities to impact random walks' dynamic.
- § 3 contains our research: identifying different approaches (§ 3.1), our HON variants of embeddings (§ 3.2), and a synthetic network (§ 3.3).
- § 4 evaluates the embeddings in two experiments. § 4.1 applies the embeddings to the synthetic network to investigate theoretic properties and visualize the embeddings. § 4.2 classifies nodes based on social interaction data.
- § 5 summarizes the thesis.
- The appendix contains simulation results, additional experiments, the source code, and development-related communication.
- Finally, a nomenclature lists the definitions used.

# Chapter 2

## Background

### 2.1 Networks

A *graph*<sup>1</sup> describes objects (called ‘nodes’ or ‘vertices’) and their (*dyadic* i.e. pairwise) relations (called ‘edges’). The sets of nodes and edges are denoted by  $\mathcal{V}$  and  $\mathcal{E}$ , respectively. A *network* is a graph that allows for attributes of nodes or edges; examples for attributes are labeled nodes and weighted edges.

Networks are ubiquitous, as they describe social relationships, collaborations, web pages, shopping carts, protein interactions, and transportation of passengers, goods, and energy. To define a particular network, we only need to provide information for each node about the nodes in its neighborhood (i.e., those other nodes connected by an edge with the former node). Examining the network in its entirety, we can then rank nodes by importance, find the shortest paths between nodes, identify communities or bottlenecks, or simulate diseases’ transmission. It is incredible how complex systems can be specified by gathering only local information.

We distinguish between *undirected* and *directed networks*, where the former additionally assumes symmetric relations. In literature, networks are undirected by default. However, I will explicitly mention if a network is assumed to be undirected since this is a restriction compared to the general (i.e., directed) case. Moreover,  $(u, v) \in \mathcal{E}$  tests the existence of an edge, regardless of the network’s directionality.

#### 2.1.1 Random Walks on Networks

Another crucial concept is a *walk* on a network, essentially a sequence of nodes with the constraint that an edge must exist between two consecutive nodes in the sequence. As this is a fundamental concept for this thesis, let us introduce proper notation. For a tuple of nodes  $(s_1, \dots, s_L)$  in  $\mathcal{V}^L$ , use the abbreviation  $s_{1..L}$ . This tuple is a walk if  $(s_\ell, s_{\ell+1}) \in \mathcal{E}, \forall 1 \leq \ell < L$ , and we use the notation  $s_1 \rightarrow \dots \rightarrow s_L$  (shortened as  $\overrightarrow{s_{1..L}}$ ).

A *random walk*  $S_1 \rightarrow \dots \rightarrow S_L$  is a walk with random nodes  $(S_1, \dots, S_L)$ , and we denote its probability distribution by  $P_{\overrightarrow{S_{1..L}}}(\dots)$ .

---

<sup>1</sup>As in *graph theory*, see the Seven Bridges of Königsberg problem solved by L. Euler.

The reason for such a rigorous definition is that merely writing  $P(A \rightarrow B)$  for two nodes  $A$  and  $B$  conceals the ambiguity between  $P_{\overrightarrow{S_1..2}}(S_1 = A, S_2 = B)$  and e.g.  $P_{\overrightarrow{S_1..3}}(S_1 = A, S_2 = B)$ . This ambiguity is relevant as the following decomposition is fundamental for random walks:

$$P(A \rightarrow B \rightarrow C) = P(A) \cdot \underbrace{\frac{P(A \rightarrow B)}{P(A)}}_{P(B|A)} \cdot \underbrace{\frac{P(A \rightarrow B \rightarrow C)}{P(A \rightarrow B)}}_{P(C|A \rightarrow B)}$$

Hence, we have to assume that  $P_{\overrightarrow{S_1..L}}$  is **consistent** for different  $L$ ; this means that  $P_{\overrightarrow{S_1..L}}$  is the marginal distribution of  $P_{\overrightarrow{S_1..L_{\max}}}$  for some  $L_{\max}$ ,  $\forall L < L_{\max}$ . However, in practice, we use the right-hand side of Equation (2.1) to model walk probabilities consistent by construction.

$$P_{\overrightarrow{S_1..L}}(S_{1..L} = s_{1..L}) := P_{\overrightarrow{S_1}}(S_1 = s_1) \prod_{\ell=2}^L P_{\overrightarrow{S_1..L}}(S_\ell = s_\ell | S_{1..\ell-1} = s_{1..\ell-1}) \quad (2.1)$$

$$= P_{\overrightarrow{S_1..L-1}}(S_{1..L-1} = s_{1..L-1}) \cdot P_{\overrightarrow{S_1..L}}(S_L = s_L | S_{1..L-1} = s_{1..L-1}) \quad (2.2)$$

The simulation of random walks is straight forward by using Equation (2.2) iteratively.

### 2.1.2 First-Order Networks (FON)

An important assumption, known as *first-order Markov property*, states that the edges' distribution depends only on the current node (and not on previously visited ones):

$$P_{\overrightarrow{S_1..L}}(S_L = s_L | S_{1..L-1} = s_{1..L-1}) = P_{\overrightarrow{S_1..L}}(S_L = s_L | S_{L-1} = s_{L-1})$$

Furthermore, assume the transition probabilities are *time-homogeneous* or *stationary* since otherwise, they would also depend on the current node's index in the sequence:

$$P_{\overrightarrow{S_1..L}}(S_L = s_L | S_{L-1} = s_{L-1}) = P_{\overrightarrow{S_1..2}}(S_2 = s_L | S_1 = s_{L-1})$$

Define a *first-order network (FON)* as a network with random walks having stationary, first-order Markov transition probabilities, and introduce the notation

$$T_{s \rightarrow t} := P_{\overrightarrow{S_1..2}}(S_2 = t | S_1 = s), \quad \forall s, t \in \mathcal{V}$$

which allows simplifying the walk probability of Equation (2.1) to

$$P_{\overrightarrow{S_1..L}}(S_{1..L} = s_{1..L}) = P_{\overrightarrow{S_1}}(S_1 = s_1) \prod_{\ell=2}^L T_{s_{\ell-1} \rightarrow s_\ell}$$

Hence, the random walks in a FON are defined by  $P_{\overrightarrow{S_1}}$  and  $T_{s \rightarrow t}$ . An outstanding choice for the former is the *stationary distribution*  $\pi$ , which fulfills  $\pi(t) = \sum_{s \in \mathcal{V}} \pi(s) \cdot T_{s \rightarrow t}$ . While a FON has stationary transition probabilities, we do not require stationarity for the first node of a walk, i.e.,  $P_{\overrightarrow{S_1}} = \pi$ .

**Other Definitions** I have chosen to introduce FON as a particular case of random walks on a network, but let us also examine the traditional definitions.

Networks allow for edges having attributes, from which we may derive the transition probabilities. By convention, this is the ‘weight’ attribute, and therefore *weighted networks* are defined as networks, where each edge has a ‘weight’ attribute. Depending on the context, these edge weights may also store other quantities like strengths of relationship, transport capacities, distances, and time differences. For *unweighted networks*, use an edge weight of one instead.

Accordingly, define the adjacency matrix  $\mathbf{A}$  with elements  $[\mathbf{A}]_{s,t}$  equal to the weight of the edge  $s \rightarrow t$  (or zero if this edge does not exist)<sup>2</sup>. The degree matrix  $\mathbf{D}$  is diagonal and The transition matrix  $\mathbf{T}$  contains the transition probabilities, chosen proportional to the edge weights,  $\mathbf{T} = \mathbf{D}^{-1}\mathbf{A}$ .

The advantage of using  $\mathbf{A}$  over  $\mathbf{T}$  for undirected networks is that we can quickly determine the stationary distribution  $\pi$  from  $\mathbf{A}$  (using an equilibrium argument) compared to calculating an Eigenvector of  $\mathbf{T}$  otherwise.

### 2.1.3 Higher-Order Networks (HON)

Every model has its limitations. Hence, different approaches for higher-order network models have been considered, which relax the FON’s limitations. These involve allowing relations between more than two nodes (simplex, hyperedge), different types of edges, and non-Markovian random walks, see Lambiotte et al. [17].

Instead of abandoning the Markov property altogether, introduce *memory*. More precisely, the  $K^{\text{th}}\text{-order Markov property}$  (for some integer  $K > 1$ ) assumes that the transition probabilities  $P_{\overrightarrow{S_1..L}}(S_L | S_{1..L-1})$  are allowed to depend on the  $K$  most recently visited nodes  $(s_{L-K}, \dots, s_{L-1})$  in the walk sequence.

Assuming random walks with a higher-order Markov property is particularly interesting since this may change the random walk’s dynamic, see Rosvall et al. [27] and Scholtes et al. [30]. Let us examine the change in dynamic due to adding memory with a transportation network: passengers would have to decide at each station where to go next if no memory (first-order Markov) is assumed; introducing memory allows avoiding returning immediately to the previous station. Hence, first-order random walks are an inefficient way of utilizing some networks.

Stationarity is only required for the (maximal)  $K^{\text{th}}$ -order transition probabilities, allowing us to rewrite Equation (2.1) for  $L > K$ :

$$P_{\overrightarrow{S_1..L}}(S_{1..L} = s_{1..L}) := P_{\overrightarrow{S_1}}(S_1 = s_1) \underbrace{\prod_{\ell=1}^{K-1} T_{s_{1..L-\ell} \rightarrow s_{\ell+1}}^{(\ell)}}_{P_{\overrightarrow{S_1..K}}(S_{1..K} = s_{1..K})} \prod_{\ell=K}^{L-1} T_{s_{\ell+1-K..L-\ell} \rightarrow s_{\ell+1}}^{(K)} \quad (2.3)$$

using the following notation:

$$T_{s_{1..k} \rightarrow t}^{(k)} := P_{\overrightarrow{S_1..k+1}}(S_{k+1} = t | S_{1..k} = s_{1..k}), \quad \forall 1 \leq k \leq K$$

Observe that for a HON, the walk’s asymptotic dynamics depends entirely on the maximum order transition probabilities ( $T_{s_{1..K} \rightarrow t}^{(K)}$ ), while the lower order

---

<sup>2</sup>Due to the lack of universal convention regarding  $[\mathbf{A}]_{s,t}$  referring to the edges  $s \rightarrow t$  or  $t \rightarrow s$ , I decided for the unambiguous notation  $T_{s \rightarrow t}$ .

ones determine the first  $K$  nodes' distribution. Therefore, HON random walks are stationary only for particular choices of  $P_{\overrightarrow{S_1}}$  and  $T_{s_{1..k} \rightarrow t}^{(k)}$  for  $k < K$ .

**Throughout this thesis, a higher-order network model (HON) will denote a network satisfying the higher-order Markov property and stationarity (for the maximum order transition probabilities).**

**Estimating probabilities** Estimating the HON parameters from data is out of scope for this thesis, and we refer to BuildHON+ [38, 39], MON [29], and MoGEN [9] instead; see § 4.2 for an example using MON.

However, unless the random walks are stationary, inconsistencies quickly arise: on the one hand,  $T_{s \rightarrow t}$  determines the stationary distribution in a FON; on the other hand,  $T_{s_1 \rightarrow t}^{(1)}$  affects only the start of a random walk in a HON. Can we, therefore, convert a HON into a FON by merely dropping the higher-order transition probabilities, or should  $T_{s_1 \rightarrow t}^{(1)}$  differ from  $T_{s \rightarrow t}$ ? How about using the first-order stationary distribution  $\pi$  in conjunction with a HON?

Conversely, MoGEN uses a consistent model, including a stop symbol indicating a walk's end. However, this stop represents an absorbing state, leading to a degenerate stationary distribution. (Though, page rank is not affected.)

## 2.2 Background on Embeddings for FON

Despite the versatility of networks, some tasks require different representations of the information encoded in networks. In particular, networks are not well suited for machine learning tasks, which expect vector-valued data for each node. Hence, extracting suitable features (or representations) is a vital pre-processing step of many machine learning applications. Given this can be a labor-intensive task if done manually, it is no surprise that research about *representation learning* got quite some attention and is indeed an active field of research for a couple of years [2, 13, 40].

When visualizing a network, suitable locations (i.e., coordinates) for the nodes must be chosen. If we want to visualize communities and already have learned a representation suitable for classification, we could apply *t-SNE* to this representation for visualization. If we specify, which nodes should attract or repulse other ones, we can use a force-directed layout. Many popular dimension reduction techniques handle the particular case of pairwise distances (or dissimilarities) between nodes. In particular, PCA [6] and metric MDS [16] aim to find representations in a low-dimensional vector-space while preserving distances between points.

### 2.2.1 Skip-Gram Model

The *skip-gram model* [20] was introduced for natural language processing (NLP), aiming to predict the neighboring words (denoted as context) for any word. The first step in applying the skip-gram model is to obtain word–context pairs. For example, a sentence (interpreted as a list of words) generates word–context pairs by returning those pairs of words of the sentence, which neither are too far apart nor in the same position.

Let  $p(w, c)$  denote the sampling distribution of the word–context pairs. The skip-gram model aims to approximate the conditional distribution  $p(c|w)$  of the

context given the word  $w$ . However, the  $p(c|w)$  values are not stored directly due to the quadratic growth in storage complexity with the number of words and the inability to generalize to unobserved word–context pairs. Instead, embedding vectors  $\vec{w}$  and  $\vec{c}$  for  $w$  and  $c$ , respectively, are learned such that  $p(c|w)$  is approximated by  $\hat{p}(c|w)$ , which is the softmax (for fixed  $w$  and varying  $c$ ) of the scalar product  $\langle \vec{w}, \vec{c} \rangle$  between the embeddings of the word  $w$  and the context  $c$ .

$$\hat{p}(c|w) = \frac{\exp(\langle \vec{w}, \vec{c} \rangle)}{\sum_{v \in \mathcal{V}} \exp(\langle \vec{w}, \vec{v} \rangle)} \quad (2.4)$$

Note that the embedding usually depends on the role (word or context), so even if  $w$  and  $c$  are equal, their embeddings  $\vec{w}$  and  $\vec{c}$  might differ.

Hence, two embeddings  $\vec{w}_1$  and  $\vec{w}_2$  for two words  $w_1$  and  $w_2$  have higher cosine similarity, the closer their conditional distributions  $p(c|w_1)$  and  $p(c|w_2)$  are; see also § 2.3.2.

Training the skip-gram model aims for minimizing the KL-divergence from  $\hat{p}(c|w)$  to  $p(c|w)$ , where the different words use, e.g.,  $p(w) = \sum_c p(w, c)$  as weights for joint minimization. However, minimizing with stochastic gradient descent (SGD) is not straight forward due to the softmax’s denominator (normalization) depending on all context embeddings. Therefore, each gradient descent step must update all of them, which is generally too expensive. There are two approximations of the model provided by Mikolov et al. [21]:

- *Hierarchical softmax* uses a tree structure for the context embeddings, resulting in update costs growing only logarithmically with the number of contexts.
- *Negative sampling*<sup>3</sup> replaces the denominator by a subtrahend consisting of only  $N$  samples of contexts with distribution  $p_N(c)$ . This simplifies updating considerably while still avoiding the unbounded expansion of values, previously achieved by the denominator.

`Word2vec`<sup>4</sup> implements both.

Analytically determining the minimum of the objective of the skip-gram model with negative sampling (SGNS) — see [8, 19] — leads to (see § 3.2.5)

$$\langle \vec{w}, \vec{c} \rangle = \log(PMI(w, c)/N) \quad (2.5)$$

where  $PMI(w, c) = p(c|w)/p_N(c)$ , known as pointwise mutual information (in case  $p_N$  were the marginal distribution  $\sum_w p(w, c)$ ). Eq (2.5) is solved by *matrix factorization* (*MF*) (i.e. singular value decomposition, SVD).

However, matrix factorization can only approximate matrices of rank up to the dimension of the embedding. Hence, both SGNS and matrix factorization define a low-rank approximation of the same matrix, but there is a difference in using weighted or unweighted projections. Another difference arises because matrix factorization substitutes small values of PMI before calculating the logarithm.

In terms of scalability, SGD based methods have an advantage over matrix factorization. Moreover, SGD is easily modified to include penalties or

---

<sup>3</sup>Negative sampling approximates noise contrastive estimation (NCE), which has a stronger theoretical foundation. Nevertheless, this distinction is not central to this thesis.

<sup>4</sup><https://code.google.com/archive/p/word2vec/>

constraints. Conversely, matrix factorization is excellent for the theoretical understanding, and there are no parameters (number of iterations and learning rates) affecting convergence.

## 2.3 Related FON Embeddings

According to the classification of FON embeddings by Hamilton et al. [13], network embeddings correspond to a pair of encoders and decoders. The encoder maps the node of a network to its embedding vector. Training some embedding involves minimizing a loss function, which depends on the decoder. Following this classification, there are three groups of approaches:

- Shallow embeddings
- Autoencoders
- Convolution encoders

**Shallow embeddings** Shallow embeddings encode nodes by look-up, which requires allocating memory for each embedding vector, but besides this, they are stateless. Decoding generates a scalar value from a pair of embedding vectors.

Many embeddings fall into this category — particularly the ones based on the skip-gram model, which are summarized separately in table 2.1<sup>5</sup>.

- Classic embeddings (PCA, MDS, ...)
- Graph Factorization (GF) [1] factorizes (symmetric) adjacency matrices. The distributed calculation allows for massive graphs. However, for embeddings based on adjacency, LINE is preferred.
- GLEE [34] is based on the geometric properties of the Laplacian and scores with its solid theoretical foundation. Nevertheless, requiring an *undirected* network prevents its adaptation to HON.

**Autoencoders** Autoencoders do not store the individual embedding vectors but calculate them instead. The input is a vector of size  $|\mathcal{V}|$  for each node representing the similarities between this node and all network nodes. Function approximation via neuronal networks encodes this vector into a lower dimension before decoding it back to the original size. The state of the autoencoder consists of the parameters of the neuronal network.

This approach comprises two embeddings:

- DNGR [4] approximates the personalized page rank (PPR) similarity, similar to App/Verse.
- SDNE [37] approximates the adjacency vector, similarly to LINE.

---

<sup>5</sup>Khosla et al. [14] deserve credits for an excellent comparison study, in particular table 2.

Embedding	Word–Context	Train	Symmetric	# Embeddings
DeepWalk [23]	random walk	<b>HS</b>	yes	ignore 2nd
<u>LINE–1</u> [33]	adjacency	NS	yes	one
LINE–2 [33]	adjacency	NS	no	ignore 2nd
PTE [32]	adjacency	NS	no	ignore 2nd
GraRep [3]	multi-step adj.	NCE	no	ignore 2nd
<u>Node2vec</u> [10]	random walk	NS	yes	ignore 2nd
HOPE [24]	similarity	MF	no	two
App [41]	similarity	NS	no	two
Verse [35]	similarity	NCE	no	one
<u>NetMF</u> [25]	random walk	MF	yes / no	ignore 2nd
InfiniteWalk [5]	random walk	MF	yes	ignore 2nd

Table 2.1: The skip-gram model is the basis for many FON embeddings. (LINE–1 is the exception; see § 2.3.2.) The word–context probabilities refer to transitions to adjacent nodes, co–occurrences in random walks, or a similarity measure (e.g., PPR). The embeddings are trained by hierarchical softmax (HS), negative sampling (NS), noise contrastive estimation (NCE), and matrix factorization (MF). DeepWalk, Node2vec, and LINE–1 treat the word–context pairs symmetrically; InfiniteWalk and NetMF (for large window sizes) require undirected networks. Only Verse and LINE–1 use the same embeddings for word and context. The others calculate different embeddings, but only HOPE and App return both. The underlined embeddings will be pursued further.

**Convolution encoders** Both shallow embeddings and autoencoders are unable to generalize the embedding for new nodes. On the other hand, convolution encoders calculate the embedding while avoiding hardcoding the size of the network  $|\mathcal{V}|$ , and therefore can embed new nodes because they aggregate features from a neighborhood of each node. Random walks approximate the neighborhood, which might grow huge after just a few steps.

GraphSAGE [12] and various graph convolution networks (GCN) follow this approach.

**Considered embeddings** After reviewing the different embeddings concerning their potential for generalization from FON to HON (see the approaches in § 3.1), I decided to start with the simple ones. Hence, we will focus on the embedding marked by underlining. Note that GraRep and NetMF cover the second variant of LINE.

Alternatively, GraphSAGE looked promising, and the adaptation of its random walks to HON enjoyable, but its versatility might also require more time to find useful embeddings.

Below, we will explain some FON embeddings in detail.

### 2.3.1 DeepWalk and Node2vec

DeepWalk [23] was the first embedding for FON to adapt the skip-gram model. The word–context pairs have initially (in NLP) been generated from sentences and are replaced by random walks (of length  $L$ ) in DeepWalk.

Building on top of Word2vec, DeepWalk generates word–context pairs selecting all pairs of the random walk’s nodes within a window of size  $W$  (excluding the center of the window). Hence, word and context roles are indistinguishable, which suggests that DeepWalk better suited for *undirected networks*, although there is no hindrance in applying it to *directed networks*. DeepWalk implemented only *unweighted networks* and *hierarchical softmax*. Because DeepWalk adapted to weighted networks and negative sampling corresponds to a special case of Node2vec with parameters  $p = 1$  and  $q = 1$ , I decided to use the adapted variant of DeepWalk instead.

Interestingly, DeepWalk<sup>6</sup> implements a random walk with geometric restart probability (i.e., the personalized page rank, PPR), but this feature was disabled. So, App and Verse would independently suggest embedding the PPR similarity years later.

Node2vec [10] extends DeepWalk and parametrizes the random walk, turning the FON effectively into a 2<sup>nd</sup>-order network with

$$T_{s_1 \dots s_n \rightarrow t}^{(2)} \propto \begin{cases} \frac{1}{p} T_{s_2 \rightarrow t}, & \text{if } s_1 = t \text{ (identical)} \\ T_{s_2 \rightarrow t}, & \text{if } (s_1, t) \in \mathcal{E}, \text{ (adjacent)} \\ \frac{1}{q} T_{s_2 \rightarrow t}, & \text{else} \end{cases} \quad (2.6)$$

This adjustment (called bias) provides greater flexibility in adjusting the neighborhood exploration to the task at hand and is motivated by two strategies to traverse trees: breadth-first search (BFS) and depth-first search (DFS).

### 2.3.2 LINE

Tang et al. [33] did an excellent job explaining the two variants of LINE, and the difference is also crucial for the understanding of the skip-gram model. In particular, we will refer to the clustering of LINE-2 embedded lattices by parity later in this thesis.

For both variants of LINE, the word–context pairs correspond to edges (i.e., adjacent nodes), but they use different approximations to train the embeddings.

The first variant, LINE-1, uses a *different* though related approximation than the skip-gram model. Instead, it embeds the *first-order proximity*, which approximates  $p(w, c)$  with ( $\sigma$  is the *sigmoid* function)

$$\hat{p}(w, c) := \sigma(\langle \vec{w}, \vec{c} \rangle) = \frac{1}{1 + \exp(-\langle \vec{w}, \vec{c} \rangle)}$$

Only undirected networks can be embedded because word and context use the same embeddings. Hence, the approximated probability of an edge between two nodes  $w$  and  $c$  increases with the scalar product  $\langle \vec{w}, \vec{c} \rangle$  because  $\sigma$  is strictly monotonically increasing. Therefore, high values of  $\langle \vec{w}, \vec{c} \rangle$  should relate to adjacent pairs of nodes  $(w, c)$ .

However, we are more interested in distances than scalar products. If the embedding vectors had equal length, we could conclude that large scalar products correspond to small distances, because  $\|\vec{w} - \vec{c}\|^2 = \|\vec{w}\|^2 + \|\vec{c}\|^2 - 2\langle \vec{w}, \vec{c} \rangle$ . In figure 2.1, we apply LINE to a network with a two-dimensional lattice structure<sup>7</sup>

---

<sup>6</sup><https://github.com/phanein/deepwalk>

<sup>7</sup>no self-loops; also known as square grid graph

and display distances in the lattice, average lengths, and the angles between pairs of nodes ( $w_1, w_2$ ). Figure 2.2 displays the corresponding t-SNE visualizations. Admittedly, this required some tuning of t-SNE, but for LINE-1, this worked well.

In contrast, LINE-2 embeds the *second-order proximity*, which approximates  $p(c|w)$  with  $\hat{p}(c|w) = \exp(\langle \vec{w}, \vec{c} \rangle)/\text{const}(w)$ ; see § 2.2.1. So does this preserve the lattice structure too? No, according to figure 2.2, the embedded lattice structure consists of two separate clusters, and according to figure 2.1, that the scalar product  $\langle \vec{w}_1, \vec{w}_2 \rangle$  of a pair of nodes ( $w_1, w_2$ ) is generally bigger if the distance between  $w_1$  and  $w_2$  is two (compared to one or three). Pairs of nodes with (shortest path) distance two are also the only ones with overlapping neighborhoods.

Define a node’s **parity** (for a lattice) as the parity bit of the sum of its coordinates. Hence, pairs of nodes with overlapping neighborhoods must have the same parity. Intuitively, each embedded node  $\vec{w}$  is pulled towards all of its neighboring contexts’ embeddings  $\vec{c}$ , while the negative samples act as a counterforce. It appears that the nodes with the same parity are (indirectly) linked together, while there are no such connections otherwise.

Embeddings of lattices are the primary model for investigating the theoretical properties of embeddings throughout this thesis. Moreover, explaining the clustering by parity was the primary rationale for explaining LINE in detail.

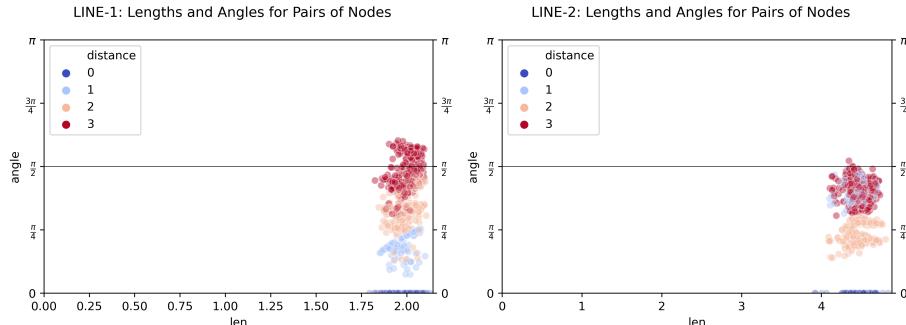


Figure 2.1: Examine the scalar product  $\langle \vec{w}_1, \vec{w}_2 \rangle = \text{len}^2 \cos(\text{angle})$  for a pair of nodes ( $w_1, w_2$ ) for LINE-1 (left) and LINE-2 (right). The ‘distance’ refers to the shortest path distance in the lattice. For distance zero, observe the concentration of the lengths  $\|\vec{w}\|$ , which confirms the analogy between scalar product  $\langle \vec{w}_1, \vec{w}_2 \rangle$  and distance  $\|\vec{w}_1 - \vec{w}_2\|$ . For LINE-1, the angle between pairs of nodes grows with their distance, indicating that the embedding might reflect the lattice structure. For LINE-2, the angles are considerably smaller for distance two (than one or three), indicating that these are the closest neighbors in the embedding space.

### 2.3.3 GraRep

GraRep [3] takes a middle ground between embedding the adjacency matrix (LINE-2) and a random walk (DeepWalk) by embedding multi-step transition probabilities separately for each step-size. However, it differs from the others by training the embedding with matrix factorization instead of SGD.

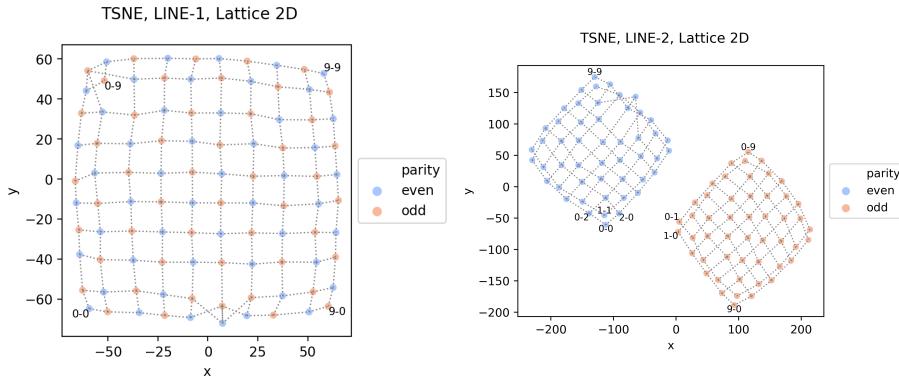


Figure 2.2: T-SNE visualizations of embeddings of a two-dimensional lattice with LINE-1 (left) and LINE-2 (right). The colors signify the parity bit of the sum of the coordinates of each node. LINE-1 mostly preserves the lattice structure, which is further highlighted by connecting adjacent nodes. The LINE-2 embeddings fall into two groups, the parity, and connecting adjacent edges would not be helpful. Instead, horizontal and vertical two-step neighbors are connected.

Hence, for each step-size  $k \in \{1, \dots, S\}$  define

$$PMI^{(k)}(w, c) = \frac{[\mathbf{T}^k]_{w,c}}{\frac{1}{|\mathcal{V}|} \sum_{u \in \mathcal{V}} [\mathbf{T}^k]_{u,c}}$$

and factor  $\langle \vec{w}, \vec{c} \rangle = \max(0, \log(PMI^{(k)}(w, c)/N))$  with truncated SVD to obtain  $\mathbf{U}_R^{(k)} \cdot \mathbf{\Lambda}_R^{(k)} \cdot \mathbf{V}_R^{(k)}$ , where  $R$  is the embedding dimension, and  $N$  is the number of negative samples<sup>8</sup>. Finally, concatenate the individual word embeddings  $\mathbf{U}_R^{(k)} \cdot \sqrt{\mathbf{\Lambda}_R^{(k)}}$  after scaling<sup>9</sup> them.

Note that when applying GraRep to a lattice, we observe the same clustering by parity issue as LINE-2. This happens due to not aggregating different step-sizes.

The authors of GraRep argue that aggregating over different step-sizes mixes up distinct features. While this argument initially sounds convincing, we should nevertheless remember that flexibility is a double-edged sword. Maybe a parameterized aggregation of different step-sizes before the matrix factorization would be useful — in the same way as Node2vec extends DeepWalk.

### 2.3.4 NetMF and InfiniteWalk

NetMF [25] adapts DeepWalk by using the connection between SGNS and matrix factorization discovered by Levy and Goldberg [19]. It assumes an *undirected, weighted network*, enabling the use of Eigendecomposition and a closed-form expression for the stationary distribution  $\pi$  proportional to the degree.

The distribution of word–context pairs corresponds to *co-occurrence probabilities of random walks* starting at the word–node and visiting  $W$  context–

<sup>8</sup><https://github.com/ShelsonCao/GraRep> has  $N = 1$  hardcoded.

<sup>9</sup>According to the source code. The paper does not mention this.

nodes. The length of the random walk  $W$  corresponds to the *window size* in DeepWalk and is therefore named the same.

$$p(c|w) = \frac{1}{W} \sum_{k=1}^W [\mathbf{T}^k]_{w,c} \quad (2.7)$$

NetMF implicitly assumes that the marginal distribution  $p(w)$  equals the stationary distribution  $\pi$ , justifying the distribution of negative samples  $p_N = \pi$  because the other marginal distribution  $p(c)$  is the same. Conversely, DeepWalk weights the start of random walks uniformly, and  $p(w)$  is somewhere ‘in between’ the uniform (for short  $L$ ) and the stationary (for large  $L$ ) distribution. Also, GraRep assumes a uniform distribution for  $p(w)$  because otherwise,  $PMI^{(k)}$  denominator would need to be a weighted average. Before NetMF, the authors of GraRep also introduced a similar method, E-SGNS, which differs from NetMF by the chosen  $p_N$ .

For *small window sizes*  $W$ , we directly calculate  $PMI(w, c) = p(c|w)/p_N(c)$  using Equation (2.7) and factor  $\langle \vec{w}, \vec{c} \rangle = \max(0, \log(PMI(w, c)/N))$  with truncated SVD of dimension  $R$  to obtain  $\mathbf{U}_R \cdot \mathbf{\Lambda}_R \cdot \mathbf{V}_R$ , where  $N$  is the number of negative samples. Finally, the resulting word embeddings are  $\mathbf{U}_R \cdot \sqrt{\mathbf{\Lambda}_R}$ .

**Approximation** For *large window sizes*  $W$ , Equation (2.7) is computationally challenging; hence we utilize the symmetry of the adjacency matrix  $\mathbf{A}$  of an *undirected* network. Define the *symmetrized transposition matrix*  $\tilde{\mathbf{T}} = \mathbf{D}^{-0.5} \mathbf{A} \mathbf{D}^{-0.5}$  and use the Eigendecomposition  $\tilde{\mathbf{T}} = \mathbf{Q} \mathbf{\Lambda} \mathbf{Q}^{-1}$ . This simplifies the calculation of powers of  $\mathbf{T} = \mathbf{D}^{-0.5} \mathbf{Q} \mathbf{\Lambda} \mathbf{Q}^{-1} \mathbf{D}^{+0.5}$ , and hence

$$p(c|w) = \mathbf{D}^{-0.5} \mathbf{Q} \left( \frac{1}{W} \sum_{k=1}^W \mathbf{\Lambda}^k \right) \mathbf{Q}^{-1} \mathbf{D}^{+0.5} \quad (2.8)$$

*Truncating the Eigendecomposition* (as in ‘truncated SVD’) further simplifies the calculation, and afterward proceed as above.

**InfiniteWalk** The limit of NetMF for  $W \rightarrow \infty$  is investigated by InfiniteWalk [5], which also assumes an undirected network to utilize Equation (2.8).

For this thesis, we have to assume directed networks, and NetMF with small window sizes is easily adapted and does not enforce symmetry (unlike DeepWalk). However, for NetMF with large window sizes and InfiniteWalk, this cannot be fixed. Therefore, we do not pursue InfiniteWalk further.

## 2.4 Related HON Embeddings

To our knowledge, HONEM [28] is the only published embedding for HON. Building on top of BuildHON+ [38, 39], it aggregates higher-order transition probabilities into a neighborhood matrix  $\mathbf{D}_{\text{neighbor}}$ , which is factored by SVD to obtain the embedding. However, given the vast number of FON embeddings, there is certainly still room for more competition.

### 2.4.1 Definition of HONEM

After fitting a  $K^{\text{th}}$ -order Markov model (including estimation of  $K$ ) with BuildHON+, we obtain transition probabilities ('rules')  $T_{s_{1..k} \rightarrow t}^{(k)}$  for any tuple  $s_{1..k}$  corresponding to a walk  $\overrightarrow{s_{1..k}}$  ( $k \leq K$ ). However, some tuples are unavailable due to the lack of data (i.e., walk  $\overrightarrow{s_{1..k}}$  has probability zero) or the removal of unnecessary transition probabilities ('rule pruning'). Hence,  $T_{s_{1..k} \rightarrow t}^{(k)}$  is only available for  $s \in \text{rule\_keys}$ , with

$$\text{rule\_keys} := \bigcup_{k=1}^K \{s_{1..k} \mid \exists T_{s_{1..k} \rightarrow t}^{(k)}\} \subseteq \bigcup_{k=1}^K \{s_{1..k} \in \mathcal{V}^k \mid \overrightarrow{s_{1..k}} \text{ is a walk}\}$$

The neighborhood matrix  $\mathbf{D}_{\text{neighbor}}$  is defined as

$$\mathbf{D}_{\text{neighbor}} := \sum_{k=1}^K \exp(1 - k) \cdot \mathbf{D}_{\text{neighbor}}^{(k)}$$

based on the  $k^{\text{th}}$ -order neighborhood matrix  $\mathbf{D}_{\text{neighbor}}^{(k)}$  — with  $\text{average}(\emptyset) = 0$

$$[\mathbf{D}_{\text{neighbor}}^{(k)}]_{i,j} := \text{average}(\{T_{s_{1..k} \rightarrow t}^{(k)} \mid s \in \text{rule\_keys}, s_1 = i, t = j\})$$

Note that  $\mathbf{D}_{\text{neighbor}}^{(1)}$  corresponds to the first-order transition matrix  $\mathbf{T}$ .

Applying truncated SVD with dimension  $R$  to  $\mathbf{D}_{\text{neighbor}}$  provides the factorization  $\mathbf{U}_R \cdot \mathbf{\Lambda}_R \cdot \mathbf{V}_R$ , and the embedding is defined as  $\mathbf{U}_R \cdot \sqrt{\mathbf{\Lambda}_R}$ .

### 2.4.2 Theoretical Analysis of HONEM

Before addressing the weaknesses of HONEM, let us keep in mind that BuildHON+ decided to prune the rules to reduce memory requirements. Similarly, the HONEM neighborhood matrix's construction is computationally cheap, making HONEM suitable for large graphs. (As long as truncated SVD is feasible.)

First, observe that HONEM does not introduce any non-linearity (i.e., log) before factorizing the neighborhood matrix. However, InfiniteWalk [5] claims that embeddings using a low-dimensional matrix factorization without non-linearity perform poorly compared to skip-gram methods, which factorize the logarithm of (positive, shifted) PMI.

From a statistical modeling point of view, however, the construction of the neighborhood matrix — mainly the higher-order terms — is unsatisfactory, as it averages the transition probabilities of the last-step over walks of a given length. (Moreover, are walks with  $T_{s_{1..k} \rightarrow t}^{(k)} = 0$  or  $\approx 0$  included in the average?) The sum of the conditional (on the first node only!) probabilities of any walk ending at a given point would be better interpretable. However, if we also follow the advice of GraRep [3] to avoid aggregating multi-step transition probabilities of different orders, we would end up with HON GraRep; see § 3.2.4.

The rule pruning of BuildHON+ particularly has a detrimental effect on the HONEM neighborhood matrix: While the neighborhood matrix depends continuously on the output of BuildHON+, the pruning mechanism introduces a discontinuity, allowing the design of adversary examples — which is fun and informative. Inspired by [30], we construct a HON where the higher-order

probabilities are adjustable without affecting the first-order terms. Moreover, this is an excellent opportunity to introduce an essential technique for this thesis.

Figure 2.3 illustrates the effect of pruning on a simple network consisting of three nodes  $A$ ,  $B$ , and  $C$ . The first-order transition probabilities are  $T_{A \rightarrow B} = 1$ ,  $T_{B \rightarrow A} = 0.5 = T_{B \rightarrow C}$ , and  $T_{C \rightarrow B} = 1$ . For  $\omega \in (-0.5, 0.5)$ , define  $T_{AB \rightarrow C}^{(2)} = 1 + \omega = T_{CB \rightarrow A}^{(2)}$  and  $T_{AB \rightarrow A}^{(2)} = 1 - \omega = T_{CB \rightarrow C}^{(2)}$ . For  $\omega > 0$ , the walks  $A \rightarrow B \rightarrow C$  and  $C \rightarrow B \rightarrow A$  become more likely, which we denote as *speed-up*. Conversely,  $\omega < 0$  results in a *slow-down*. Note that the first-order stationary distribution  $\pi$  is not affected by this change. ( $\pi(B) = 0.5$  remains unchanged because a walk has to visit node  $B$  every second step, and  $\pi(A) = \pi(C)$  follows from symmetry.)

The  $(A, C)$ -entry of the second-order neighborhood matrix,  $[\mathbf{D}_{\text{neighbor}}^{(2)}]_{A,C}$  equals to  $T_{AB \rightarrow C}^{(2)}$ , unless this rule is pruned by BuildHON+ (i.e.,  $\omega \approx 0$  or equivalently  $T_{AB \rightarrow C}^{(2)} \approx T_{B \rightarrow C}$ ) resulting in  $[\mathbf{D}_{\text{neighbor}}^{(2)}]_{A,C} = 0$ . Note, that both *slow-down* ( $\omega < 0$ ) and *speed-up* ( $\omega > 0$ ) result in  $[\mathbf{D}_{\text{neighbor}}^{(2)}]_{A,C} > 0$ . This is unfortunate because tricks to save memory should have (almost) no impact on the embedding. Consequently, throughout this thesis, **HONEM is calculated without rule pruning**.

However, accommodating for the rule pruning would make HONEM's computation more expensive and give up its niche.

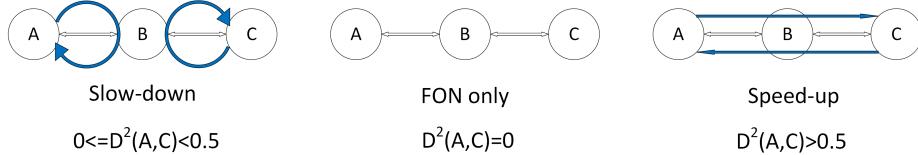


Figure 2.3: Effect of rule pruning on the second-order neighborhood matrix  $\mathbf{D}_{\text{neighbor}}^{(2)}$ . The blue arrows indicate second-order transition probabilities exceeding corresponding first-order ones. The only walk from  $A$  to  $C$  of length 2 is  $A \rightarrow B \rightarrow C$ , which has probability  $\pi(A) \cdot T_{A \rightarrow B} \cdot T_{AB \rightarrow C}^{(2)}$  and therefore is *slowed-down* or *speeded-up* if  $T_{AB \rightarrow C}^{(2)} < 0.5$  or  $> 0.5$ , respectively. By definition,  $[\mathbf{D}_{\text{neighbor}}^{(2)}]_{A,C} = T_{AB \rightarrow C}^{(2)}$ , but if  $T_{AB \rightarrow C}^{(2)} \approx T_{B \rightarrow C}$ , BuildHON+ might prune this rule resulting in  $[\mathbf{D}_{\text{neighbor}}^{(2)}]_{A,C} = 0$ . This *discontinuity* is undesirable.



# Chapter 3

## Methodology

### 3.1 Approaches for HON embeddings

After gaining an overview of existing embeddings for first-order networks, let us see how to adapt these to higher-order networks. Here, we lay out the ideas, while § 3.2 describes the embeddings considered.

#### 3.1.1 Higher-order random walks

Given that random walks make the difference between first- and higher-order networks, adapting the random walks is the most natural approach to generalizing FON embeddings to HON. Conveniently, no changes to the down-stream task are necessary.

Hence, let us identify the FON embeddings of § 2.3 depending on random walks, starting with embeddings that *sample* random walks:

- DeepWalk generates word-context pairs via random walks.
- Node2vec generates word-context pairs via biased random walks.
- App and Verse use the personalized page rank (PPR) similarity, defined as a random walk with geometric reset probability.
- GraphSAGE uses random walks to sample the neighborhood of each node.

Other embeddings calculate the *distribution of nodes visited* by random walks:

- GraRep embeds  $k$ -step transition probabilities, which are the distribution of the endpoint of random walks truncated to length  $k$ .
- NetMF is the matrix factorization variant of DeepWalk.
- DNGR encodes the personalized page rank (PPR) similarity.

These embeddings all lead to viable candidates for HON embeddings, and four will be pursued further in § 3.2.

However, while the adaptation of random walks from FON to HON is straightforward mathematically, it nevertheless poses computational challenges. A symmetric FON adjacency matrix  $\mathbf{A}$  enables calculating the distribution of nodes quickly using approximation (2.8). However, this trick does not work for directed and higher-order networks; see algorithm 3.1 instead.

### 3.1.2 Beyond pairwise interactions

All shallow embeddings in § 2.3 do model only *pairwise interactions*. We are concerned that this might be *too restrictive* because we cannot recover a multivariate distribution (e.g., the stationary distribution of the  $K^{\text{th}}$ -order Markov model) from bivariate marginal distributions. However, the multivariate distribution also has higher complexity than its marginals.

Moreover, the down-stream task may prevent embedding higher interactions, as it happens for this thesis.

Embedding higher-order interactions is straight forward: Each  $K^{\text{th}}$ -order Markov chain on a network corresponds to a first-order Markov chain with  $K$ -tuples of nodes<sup>1</sup> as state-space. Hence, we may use tuples of nodes (or walks) as words and contexts (instead of single nodes). To avoid redundancy, we could also assign walks as word and nodes as context.

This leaves us with three options to model interactions:

**Pairs of Nodes** Both words and contexts are single nodes. To utilize the HON, we combine this with HON random walks; see § 3.1.1. Inevitably, this cannot capture the full complexity of a HON model.

**Mixed** The word is a walk, and the context is a single node. Unlike the previous case, we can store the entire HON transition probabilities  $T_{w_{1..k} \rightarrow c}^{(k)}$  for  $k \leq K$  in the conditional word-context distribution  $p(c|w)$ . Similarly, we can distinguish different starts (e.g.,  $(s_1), (s_1, s_2)$ ) of random walks. This flexibility also comes at a (storage) cost.

**Pairs of Walks** Both words and contexts are walks. This is even more flexible than the ‘mixed’ case but also includes redundant information. (Not investigated further.)

To decrease the cost of embedding walks in the ‘mixed’ case, I considered two approaches:

- Exploiting the start-tuples’ hierarchical structure  $(s_K) \subset s_{K-1..K} \subset \dots \subset s_{1..K}$  and introducing a penalty on  $\sum_{\ell=1}^{K-1} \|\overrightarrow{s_{\ell+1..K}} - \overrightarrow{s_{\ell..K}}\|^2$  allows reducing the number of samples for training; see § 3.2.5.
- Instead of embedding each walk independently, combine the walk’s individual nodes’ embeddings with a Hadamard product. This reduces the embeddings’ storage, and leads to *Tensor decomposition* [15], brought to my attention by HOSGNS [24]. (Not pursued further.)

The HON embeddings in § 3.2 implement both ‘pairs of nodes’ (four times) and ‘mixed’ (three times) interactions.

However, the down-stream task ultimately decides whether higher-order interactions are beneficial. For the classification task in § 4.2, we learn labels for nodes, therefore learning embeddings for tuples of nodes is pointless, and we use only ‘pairs of nodes’ interactions.

---

<sup>1</sup>De Bruijn graph

## 3.2 Higher-order network embeddings examined

This section describes our HON embeddings, which all base on the skip-gram model (unlike HONEM). We simulate HON random walks iteratively using Equation (2.3) or equivalently:

$$P_{\overrightarrow{S_1..L+1}}(S_{L+1} = s_{L+1} | S_{1..L} = s_{1..L}) = \begin{cases} P_{\overrightarrow{S_1}}(S_1 = s_1) & \text{if } L = 0 \\ T_{s_{1..L} \rightarrow s_{L+1}}^{(L)} & \text{if } 1 \leq L \leq K \\ T_{s_{L+1-K..L} \rightarrow s_{L+1}}^{(K)} & \text{else} \end{cases}$$

Note that by tweaking the network, we can add functionality, such as FON variants of embeddings, cross-validation, or the biased random walk of Node2vec.

All embeddings implement a `decode`-method for calculating  $\hat{p}(c|w)$  from the embedding. We may improve this approximation by consulting the network topology (`use_neighborhood`), excluding self-loops (`no_self_loops`), and normalizing to one (`normalize`).

### 3.2.1 HON DeepWalk

HON DeepWalk combines HON random walks with Word2vec and therefore inherits the restriction that the word-context pairs are treated symmetrically. Consequently, it only supports ‘pair of nodes’ interactions.

The parameters for generating the HON random walks are `num_walks` (number of random walks per node) and `walk_length`  $L$  (length/number of steps of each random walk). The parameters for Word2vec are `dimension`  $R$  (dimension or size of the embedding), `window_size`  $W$ , `hs` (use hierarchical softmax), `negative`  $N$  (number of negative samples), `num_iter` (number of iterations), plus others. The `random_seed` ensures that random walks and embedding are reproducible.

As mentioned in § 2.3.1, I prefer negative sampling, even though FON DeepWalk uses hierarchical softmax.

### 3.2.2 HON Node2vec

Node2vec introduced a bias (with parameters  $p$  and  $q$ ) affecting all random steps except the first one. This bias perfectly fits into the HON framework and is implemented by tweaking the HON transition probabilities.

For HON Node2vec, first, ensure that the network contains at least the second-order probabilities  $T_{s_{1..2} \rightarrow t}^{(2)}$ ,  $\forall (s_1, s_2) \in \mathcal{E}$ , and define  $T_{s_{1..2} \rightarrow t}^{(2)} := T_{s_2 \rightarrow t}$  otherwise. Second, replace  $T_{s_{1..k} \rightarrow t}^{(k)}$  with  $T_{s_{1..k} \rightarrow t}^{(k),bias}$  for  $k > 1$  with

$$T_{s_{1..k} \rightarrow t}^{(k),bias} \propto \begin{cases} \frac{1}{p} T_{s_{1..k} \rightarrow t}^{(k)}, & \text{if } s_{k-1} = t \text{ (identical)} \\ T_{s_{1..k} \rightarrow t}^{(k)}, & \text{if } (s_{k-1}, t) \in \mathcal{E}, (\text{adjacent}) \\ \frac{1}{q} T_{s_{1..k} \rightarrow t}^{(k)}, & \text{else} \end{cases}$$

Finally, embed the modified network with HON DeepWalk.

### 3.2.3 HON NetMF

HON NetMF generates word–context pairs corresponding to HON random walks and embeds them with matrix factorization. It implements both ‘pairs of nodes’ and ‘mixed’ interactions.

For the distribution of the negative samples  $p_N$ , we use the stationary distribution  $\pi$ . For FON, this is natural because it follows from assuming  $p(w) = \pi(w)$ . For HON, however, this would require stronger assumptions than just stationarity of the maximal order  $T_{s_{1..K} \rightarrow t}^{(K)}$ ; see the remarks in § 2.1.3 on estimating probabilities.

The parameters are `dimension R` (dimension of the embedding), `negative N` (number of negative samples), `pairwise` (use ‘pairs of nodes’ instead of ‘mixed’ interactions, default=`False`), and `window_size W` (number of steps of random walk).

**Optimization** While  $p(c|w)$  is the distribution of nodes appearing in a random walk starting with  $w$ , we should not enumerate all these walks to calculate the probabilities. Because under the weak assumptions that each node has an out-degree of at least two, the number of walks *grows exponentially* with the number of steps  $W$ .

Note that FON NetMF (for small window sizes) calculates  $\sum_{k=1}^W \mathbf{T}^k$  sequentially, using  $\mathbf{T}^k = \mathbf{T}^{k-1} \cdot \mathbf{T}$ , and therefore in linear time. How about HON? For a HON of order  $K$ , simulating another step for a random walk requires only knowledge of the  $K$  most recently visited nodes. To calculate  $p(c|w)$  for a HON in *linear time*,  $O(W)$ , and *constant space*, we calculate the  $k$ -step distributions sequentially and still enumerate walks but have to *truncate* them after each step. The finite number of states of the HON bounds the number of different walks required to calculate the next step; see algorithm 3.1.

The algorithm also works for variable-order Markov models (e.g., Build-HON+), provided the keys  $s_{1..k}$  fulfill:

$$s_{1..k} \in \text{rules} \implies s_{2..k} \in \text{rules}$$

Credits go to the minibatch algorithm of GraphSAGE [12, § A], which inspired algorithm 3.1.

### 3.2.4 HON GraRep

The HON GraRep embedding calculates the multi-step transition probabilities congruent with HON random walks, embeds them with matrix factorization, and utilizes a similar optimization as HON NetMF.

The parameters are `num_steps S` (number of steps), `dimension R` (dimension of the embedding of a single step; the total dimension is  $R \cdot S$ ), `negative N` (number of negative samples), `pairwise` (use ‘pair of nodes’ instead of ‘mixed’ interactions, default `True`), `neg_stationary` (use stationary distribution instead of uniform for  $p_N$ , default `False`), `normalize` (scale the individual embeddings before concatenating them, default `False`).

I introduced the parameters `neg_stationary` and `normalize` for comparison with HON NetMF and had to turn the latter off due to poor performance together with LogisticRegression in § 4.2. Whether FON GraRep indeed uses `normalize=True` is dubious; see § 2.3.3.

```

1  # transition probabilities are stored in a dictionary
2  rules[s1..k] = { t :  $T_{s_{1..k} \rightarrow t}^{(k)}$  }
3
4  def find_rule_key(start): # start is tuple of nodes
5      # Find longest key such that start ends with key.
6      # Length of this key is at most the max. order K.
7      return sorted([key for key in rules.keys()
8                      if key==start[-len(key):]])
9          ,key=len)[-1]
10
11 def probability_of_context_given_word(start, num_steps):
12     pc_w = defaultdict(float)
13     walk_probs = {start: 1}
14     for step in range(num_steps):
15         walk_probs_new = defaultdict(float)
16         for w, pw in walk_probs.items():
17             for t, pt in rules[find_rule_key(w)].items():
18                 p = pw * pt
19                 pc_w[t] += p / num_steps
20                 w_new = (*w, t) # append t to w
21                 w_new = find_rule_key(w_new) # truncate
22                 walk_probs_new[w_new] += p
23         walk_probs = walk_probs_new
24     return pc_w

```

Algorithm 3.1: The optimized calculation for HON NetMF determines the frequency  $p(c|w)$  of nodes  $c$  visited in a random walk when starting with a provided tuple of nodes  $w$ . The complexity of enumerating all walks generally grows *exponentially* with the number of steps. Utilizing the maximal order  $K$  of the Markov model allows truncating the walks, resulting in *linear* complexity.

### 3.2.5 HON Experimental Embedding

Finally, this experimental embedding implements a HON LINE-2 with ‘mixed’ interactions and a penalty on the walk hierarchy. Its training uses SGD and negative sampling without relying on Word2vec (due to its limitations). The ‘mixed’ interactions are required because a HON LINE-2 with ‘pairs of nodes’ interactions would only capture the first-order probabilities. The start nodes are shuffled, and the negative samples’ distribution  $p_N$  is either uniform or the stationary distribution  $\pi$ . The embeddings of word and context are separate. **This embedding primarily serves to test ideas**, such as the penalty; see below. Without the penalty, it allows comparing SGD against NetMF’s matrix factorization (for  $W = 1$ ). It is principally a starting point to implement HON App/Verse, but its performance limits its use to a proof-of-concept. (The performance tuning of gensim’s Word2vec implementation is explained in a blog, see [26].)

The parameters are **dimension  $R$**  (dimension of the embedding), **steps** (number of iterations per start path), **negative  $N$**  (number of negative sam-

ples), `penalty`  $\tau$  (penalty parameter), `neg_stationary` (use stationary distribution instead of uniform for pN, default True), `learning_rate_start` and `learning_rate_end` (learning rates are linear interpolated between start and end), `seed` (random seed), `max_start_len` (train only short start tuples).

**Hierarchical structure of walks** Observe that the start–tuples of a walk follow a hierarchical structure  $(s_K) \subset s_{K-1..K} \subset \dots \subset s_{1..K}$ , and a larger tuple indicates more information about the next node. If the transition probabilities were *stationary*,  $T_{s_{1..k} \rightarrow t}^{(k)} = P_{\overrightarrow{S_{\ell+1..k+k+1}}} (S_{\ell+k+1} = t | S_{\ell+1..k+k} = s_{1..k})$ ,  $\forall \ell \geq 0$ , for all orders  $k \leq K$ ,  $T_{s_{2..k} \rightarrow t}^{(k-1)}$  would be a weighted average of  $T_{s_{1..k} \rightarrow t}^{(k)}$ , because

$$\underbrace{P_{\overrightarrow{S_{2..k+1}}} (S_k + 1 = t | S_{2..k} = s_{2..k})}_{T_{s_{2..k} \rightarrow t}^{(k-1)}, \text{ if stationary}} = \sum_{s_1} T_{s_{1..k} \rightarrow t}^{(k)} \cdot P_{\overrightarrow{S_{1..k+1}}} (S_1 = s_1 | S_{2..k} = s_{2..k})$$

Which is also confirmed by figure 4.5.

Since ‘mixed’ interactions require training embeddings for more words, this also needs more word–context pairs. As a countermeasure, we add a penalty on  $\overrightarrow{s_{k-1..K}} - \overrightarrow{s_{k..K}}$  for  $k \in \{2, \dots, K\}$  to SGNS’s loss function to utilize the dependency between  $T_{s_{2..k} \rightarrow t}^{(k-1)}$  and  $T_{s_{1..k} \rightarrow t}^{(k)}$ .

In § 4.2.5, paragraph ‘entropy of transition probabilities’, we expand on the ideas presented here; see also figure 4.12.

**Gradient of Loss** To calculate the gradients for minimization by SGD, we follow the derivation of Levy and Goldberg [19]. They define the loss of a single observed word–context pair  $(w, c)$  as (where  $\sigma$  is the *sigmoid function*)

$$-\left( \log \sigma(\langle \vec{w}, \vec{c} \rangle) + N \sum_{c_N} p_N(c_N) \log \sigma(-\langle \vec{w}, \vec{c}_N \rangle) \right)$$

Since we want to add a penalty term to the loss function, we average rather than sum these individual losses to make the loss  $\mathcal{L}$  independent from the total number of pairs. Hence, we replace their number of word–context pairs,  $\#(w, c)$ , with probabilities,  $p(w, c)$ .

$$\begin{aligned} \mathcal{L} &:= \mathcal{L}_{\text{penalty}} + \sum_w \sum_c \mathcal{L}(w, c) \\ \mathcal{L}(w, c) &:= -(p(w, c) \log \sigma(\langle \vec{w}, \vec{c} \rangle) + N \cdot p(w) \cdot p_N(c) \log \sigma(-\langle \vec{w}, \vec{c} \rangle)) \\ \mathcal{L}_{\text{penalty}} &:= \tau \sum_{(w_1, w_2) \in \mathcal{H}} \|\vec{w}_1 - \vec{w}_2\|^2 \end{aligned}$$

For the penalty term  $\mathcal{L}_{\text{penalty}}$ , we introduced a penalty parameter  $\tau \geq 0$  and defined  $\mathcal{H} = \bigcup_{k=2}^K \{(s_{1..k}, s_{2..k}) | s_1 \rightarrow \dots \rightarrow s_k \text{ is a walk}\}$ , the set of pairs of walks with a child–parent relation according to the walk hierarchy.

To calculate the gradients, set  $x = \langle \vec{w}, \vec{c} \rangle$  and use  $\partial \log \sigma(x)/\partial x = \sigma(-x)$

$$\begin{aligned}\nabla_{\vec{w}} \mathcal{L} &= \nabla_{\vec{w}} \mathcal{L}_{penalty} + \sum_c \nabla_{\vec{w}} \mathcal{L}(w, c) \\ \nabla_{\langle \vec{w}, \vec{c} \rangle} \mathcal{L}(w, c) &= -(p(w, c) \cdot \sigma(-\langle \vec{w}, \vec{c} \rangle)) - N \cdot p(w) \cdot p_N(c) \cdot \sigma(+\langle \vec{w}, \vec{c} \rangle)) \\ \nabla_{\vec{w}} \mathcal{L}(w, c) &= \nabla_{\langle \vec{w}, \vec{c} \rangle} \mathcal{L}(w, c) \cdot \vec{c} \\ \nabla_{\vec{w}} \mathcal{L}_{penalty} &= \tau \sum_{(w_1, w_2) \in \mathcal{H}} (\vec{w}_1 - \vec{w}_2) \cdot \begin{cases} 2 & \text{if } w = w_1 \\ -2 & \text{if } w = w_2 \\ 0 & \text{else} \end{cases} \\ \nabla_{\vec{c}} \mathcal{L} &= \sum_w \nabla_{\vec{v}} \mathcal{L}(w, c) \\ \nabla_{\vec{c}} \mathcal{L}(w, c) &= \nabla_{\langle \vec{w}, \vec{c} \rangle} \mathcal{L}(w, c) \cdot \vec{w}\end{aligned}$$

While we are at it, let us quickly verify Equation (2.5) for  $\tau = 0$  (as promised in § 2.2.1). The minimum condition is (using  $\sigma(x) = \exp(x)\sigma(-x)$ )

$$0 \stackrel{!}{=} \frac{\nabla_{\langle \vec{w}, \vec{c} \rangle} \mathcal{L}(w, c)}{\sigma(-\langle \vec{w}, \vec{c} \rangle)} = -(p(w, c) - N \cdot p(w) \cdot p_N(c) \cdot \exp(\langle \vec{w}, \vec{c} \rangle))$$

which leads to the known result

$$\exp(\langle \vec{w}, \vec{c} \rangle) = \frac{p(w, c)}{N \cdot p(w) \cdot p_N(c)} = \frac{PMI(w, c)}{N}$$

### 3.2.6 HONEM

Being the only HON embedding published, this is the benchmark to compare other embeddings with. This implementation calculates the neighborhood matrix  $\mathbf{D}_{neighbor}$  from HONEM and generates an embedding using matrix factorization.

While the original HONEM [28] used BuildHON+ (a variable order model) to estimate the transition probabilities, we favored using the same probabilities for all embeddings to facilitate comparison. So our implementation of HONEM is a fixed order model then, though technically it supports both fixed and variable order transition probabilities. However, this is no disadvantage as the pruning of rules in BuildHON+ causes inconsistencies in HONEM; see § 2.4.2.

## 3.3 Synthetic network (HON Lattice 2D)

Let us introduce the HON Lattice 2D, a synthetic network for the comparison of embeddings.

In § 2.3.2, we already embedded a two-dimensional lattice with both LINE variants to understand the embeddings' structure. We will see in § 4.1 that other embeddings also preserve the lattice structure. Combine this with the tweaking of higher-order terms to slow-down or speed-up walks from § 2.4.2. By construction, HON aware embeddings are affected by this kind of dynamic change, while the FON embeddings (see § 2.3) are not.

The **HON Lattice 2D** is a higher-order network with a parameter  $\omega \in (0, 1)$  controlling the higher-order dynamic, which speeds up horizontal movements and slows down vertical movements, as illustrated in figure 3.1. (The first-order lattice is equivalent to  $\omega = 0$ .) Speed-up means moving twice in the same direction has increased probability at the expense of moving in the opposite direction, while other movements are unchanged. E.g., for  $\omega = 0.5$  after moving (horizontally) from node '2-2' to '3-2', the transition probabilities are  $3/8$  to '4-2' (same direction),  $1/8$  to '2-2' (opposite direction), and  $1/4$  for '3-1' and '3-3' (otherwise). Slow-down is similar, exchanging 'same direction' and 'opposite direction'.

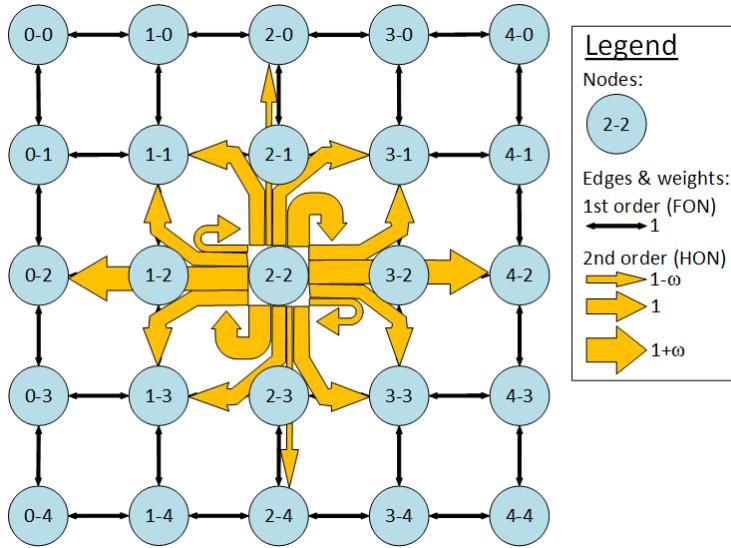


Figure 3.1: *HON Lattice 2D* is a synthetic network extending a first-order two-dimensional lattice. With transition probabilities proportional to weights. We display second-order weights for the node '2-2'. For horizontal transitions, moving twice in the same direction is amplified at the expense of U-turns. For vertical transitions, it is the other way round.

The HON lattice 2D comes with several explaining attributes, which do not affect the embeddings but are useful for visualization:

**key\_len** length of the source tuple (1 or 2).

**x\_orig**  $x$ -coordinate of the node. If  $key\_len = 2$  the last node is used.

**x\_orig**  $y$ -coordinate of the node. If  $key\_len = 2$  the last node is used.

**parity** parity ('even' or 'odd') of the sum of the coordinates; see § 2.3.2.

**direction** 'none' if  $key\_len = 1$  and 'left', 'right', 'up', or 'down' otherwise.

E.g., the transition probability  $P(N_3 | N_1 \rightarrow N_2)$  has a key  $(N_1, N_2)$  with key\_length two. The direction of this key is 'left' iff  $N_2$  is the left neighbor of  $N_1$ .

There is an interesting symmetry in the walk's distribution, as **a walk has the same probability as its reverse**. Recall that DeepWalk/Node2vec treat walks as bidirectional (due to relying on the implementation of Word2vec), while NetMF does not. Hence, this difference will be masked by that symmetry.

Moreover, the higher-order terms will not affect the (first-order) stationary distribution. For a complete definition (including boundary cases) and proofs, refer to § 3.3.1. We expect HON embeddings of the HON Lattice 2D to exhibit a stretch/compression due to the higher-order dynamic. Measuring this stretch effect is explained in § 3.3.2, and as a by-product, we improve the visualizations in § 3.3.3.

### 3.3.1 Definition and Properties of the HON Lattice 2D

For a formal definition, we have to establish some vocabulary: Equipping the lattice with cartesian coordinates, we denote the axes as horizontal and vertical. Each node  $n$  in the lattice has at most four neighbors, and we write  $L(n)$ ,  $R(n)$ ,  $U(n)$ , and  $D(n)$  to refer to the left, right, upper, and downward neighbors, respectively. To verify valid neighbors, use the following convention:

$$n \text{ has no such neighbor} \iff f(n) = \emptyset, \quad \forall f \in \{L, R, U, D\}$$

**Definition** Next, introduce weights, and we will choose the transition probabilities proportional to them below. For the FON, we use an undirected, unweighted lattice; hence each edge has weight one. More formally, we define the first order weights  $W_{n \rightarrow f(n)}^{(1)}$  as ( $\forall f \in \{L, R, U, D\}$ ):

$$W_{n \rightarrow f(n)}^{(1)} = \begin{cases} 0 & f(n) = \emptyset \quad \text{no such neighbor} \\ 1 & \text{else} \end{cases}$$

Since the node degree  $d_n = \sum_f W_{n \rightarrow f(n)}^{(1)}$ , we define  $T_{n \rightarrow f(n)} := (W_{n \rightarrow f(n)}^{(1)}) / d_n$ .

The second-order weights use a parameter  $\omega \in (0, 1)$  to assign modified weights corresponding to speed-up ( $1 + \omega$ ) or slow-down ( $1 - \omega$ ). The goal is to speed-up horizontal and slow-down vertical diffusion. Hence the second-order weights are defined as ( $\forall f_1, f_2 \in \{L, R, U, D\}$ )

$$W_{n \rightarrow f_1(n) \rightarrow f_2(f_1(n))}^{(2)} = \begin{cases} 0 & \text{if } f_1(n) = \emptyset \vee f_2(f_1(n)) = \emptyset \quad \text{no such neighbor} \\ 1 & \text{if } f_1(n) = \emptyset \quad \text{cannot move twice along } f_1 \\ 1 + \omega & \text{if } (f_1, f_2) \in \{(L, L), (R, R), (U, D), (D, U)\} \quad \text{speed-up} \\ 1 - \omega & \text{if } (f_1, f_2) \in \{(L, R), (R, L), (U, U), (D, D)\} \quad \text{slow-down} \\ 1 & \text{else} \end{cases}$$

Observe that speed-up and slow-down occur only if  $f_1 = f_2$  or  $f_2(f_1(n)) = n$  (i.e.,  $f_2$  is the inverse of  $f_1$ ;  $f_2 = f_1^{-1}$ ). To ensure that for fixed  $n$  and  $f_1$ , choosing both  $f_2 := f_1$  and  $f_2 := f_1^{-1}$  compensate any weight adjustments in total (i.e.,  $W_{n \rightarrow f_1(n) \rightarrow f_1(f_1(n))}^{(2)} + W_{n \rightarrow f_1(n) \rightarrow f_1^{-1}(f_1(n))}^{(2)} = 2$ ), we had to exclude the

case where moving twice in the first direction is impossible (i.e.,  $f_1(f_1(n)) = \emptyset$ ). This proves that  $(\forall f_1)$

$$\sum_{f_2} W_{n \rightarrow f_1(n) \rightarrow f_2(f_1(n))}^{(2)} = d_{f_1(n)}$$

Therefore define  $T_{n, f_1(n)_1..2 \rightarrow f_2(f_1(n))}^{(2)} := (W_{n \rightarrow f_1(n) \rightarrow f_2(f_1(n))}^{(2)}) / d_{f_1(n)}$ .

**Symmetry** Next, we will demonstrate that for any valid random walk (according to FON topology)  $N_1 \rightarrow \dots \rightarrow N_L$ , the **reverse walk**  $N_L \rightarrow \dots \rightarrow N_1$  has the **same probability** — provided the start node already follows the stationary distribution (i.e.,  $N_1 \sim \pi$  with  $\pi(N_1) = D_{N_1} / \sum_n d_n$ ). Because all edges have equal weight, all (valid) random walks  $N_1 \rightarrow N_2$  have the same probability. For  $L > 2$ , the probability of a random walk is as follows:

$$P(N_1 \rightarrow \dots \rightarrow N_L) = \underbrace{\frac{d_{N_1}}{\sum_n d_n}}_{const} \cdot \frac{W_{N_1 \rightarrow N_2}^{(1)}}{d_{N_1}} \cdot \prod_{i=2}^{L-1} \frac{W_{N_{i-1} \rightarrow N_i \rightarrow N_{i+1}}^{(2)}}{d_{N_i}}$$

which is non-zero, because we assume a valid walk. It remains to show that

$$W_{N_{i-1} \rightarrow N_i \rightarrow N_{i+1}}^{(2)} = W_{N_{i+1} \rightarrow N_i \rightarrow N_{i-1}}^{(2)}$$

Choose  $f_1$  and  $f_2$  such that  $f_1(N_{i-1}) = N_i$  and  $f_2(N_i) = N_{i+1}$ . Since two conditions result in weight one, we show equality only for weights different from one. The case  $f_2 = f_1^{-1}$  is trivial because, for  $N_{i-1} = N_{i+1}$ , the weights are equal. Then, the only possibility for speed-up or slow-down is  $f_1 = f_2$ , i.e., moving twice in the same direction, which is invariant under reversing the walk, and  $f_1(f_1(N_{i-1})) = N_{i+1} \neq \emptyset$  excludes another exception. This concludes the proof.

**Marginal distribution** Finally, establish that **second-order transition probabilities do not change the stationary distribution** (if starting with  $\pi$ ). Since any walk (of arbitrary length) and its reverse have the same probability, the first and the last node of this walk have identical marginal distribution, thus preserving the stationary distribution.

### 3.3.2 Measuring the Stretch Effect

Intuitively, the higher-order dynamic affects co-occurrence probabilities in random walks, where speed-up (slow-down) result in higher (lower) probabilities. In § 2.3.2, we noticed that higher probabilities are related to shorter distances in the embedding space. Hence, for a HON embedding based on the skip-gram model and random walks (e.g., HON DeepWalk), we expect the horizontal speed-up and vertical slow-downs to cause a contraction and a stretch in the respective directions of the embedding space.

To measure what we denote as the **stretch effect**, we first define *rightward edges*  $\mathcal{E}_R = \{(n, R(n)) \mid n \in \mathcal{V} \wedge R(n) \neq \emptyset\}$  and *upward edges*  $\mathcal{E}_U = \{(n, U(n)) \mid n \in \mathcal{V} \wedge U(n) \neq \emptyset\}$ . If the edges' directionality is irrelevant, instead refer to *horizontal and vertical edges* for  $\mathcal{E}_R$  and  $\mathcal{E}_U$ . We calculate the difference

in the embedding space and average its lengths among vertical and horizontal edges. Then we define the *stretch ratio* as:

$$\text{stretch} = \frac{\sum_{(u,v) \in \mathcal{E}_U} \|\vec{u} - \vec{v}\| / |\mathcal{E}_U|}{\sum_{(u,v) \in \mathcal{E}_R} \|\vec{u} - \vec{v}\| / |\mathcal{E}_R|} \quad (3.1)$$

Recall that for LINE-2, the embedding of a lattice consists of two separate clusters, each with common *parity*; see § 2.3.2. The same applies to GraRep, which does not aggregate multi-step transition probabilities over different step sizes. Moreover, LINE-2 is equivalent to DeepWalk, Node2vec, and NetMF with window size  $W = 1$  (i.e., a random walk of length one) — which are borderline cases, however. In these cases, the stretch ratio (eq. 3.1) is unsuitable for assessing the stretch effect, and we define the *two-step stretch ratio* as

$$\text{stretch}_2 = \frac{\sum_{(u,v) \in \mathcal{E}_{2U}} \|\vec{u} - \vec{v}\| / |\mathcal{E}_{2U}|}{\sum_{(u,v) \in \mathcal{E}_{2R}} \|\vec{u} - \vec{v}\| / |\mathcal{E}_{2R}|} \quad (3.2)$$

with *rightward two-step edges*  $\mathcal{E}_{2R} = \{(n, R(R(n))) \mid n \in \mathcal{V} \wedge R(R(n)) \neq \emptyset\}$  and *upward two-step edges*  $\mathcal{E}_{2U} = \{(n, U(U(n))) \mid n \in \mathcal{V} \wedge U(U(n)) \neq \emptyset\}$ .

### 3.3.3 Visualizing Embeddings of HON Lattice 2D

**Alignment** We also visualize the embeddings to confirm the stretch effect, though measuring the stretch effect in the embedding space is more meaningful. Also, align the t-SNE visualizations by an *orthogonal transformation* to correctly display the orientation ( $x$ -axis: left to right;  $y$ -axis: down to up) of a plot. (Both the skip-gram model and t-SNE rely only on scalar products or relative distances, respectively, which are unaffected by orthogonal transformations.) We calculate the *average edge direction* (i.e., differences of node pairs) for  $\mathcal{E}_R$  or  $\mathcal{E}_U$  in the visualization space to fit the orthogonal transformation. Moreover, we connect points corresponding to adjacent nodes ( $\mathcal{E}_R$  and  $\mathcal{E}_U$ ) — or their two-step counterpart ( $\mathcal{E}_{2R}$  and  $\mathcal{E}_{2U}$ ) if better suited; see figure 3.2.

**Projection** Due to the *laborious tuning of t-SNE*, therefore, the *lack of objectivity*, I initially considered using a *projection* instead. This projection assumed that the embedding is contained in a two-dimensional manifold, which is hopefully well-behaved to be approximated by its tangential space. The latter is approximated by the average differences between the embeddings of rightward/upward pairs of nodes.

While the visualizations generated by these projections (figure 3.3) looked interesting, I was finally discouraged when investigating lengths and angles (between individual embedding differences and corresponding averages), as shown in figure 3.4. The individual embedding differences turned out to be close to orthogonal to their corresponding (rightward or upward) averages. Under these circumstances, the projection seems unreliable. Hence, relying on intuition gained in low dimensions when tackling high-dimensional problems can be misleading.

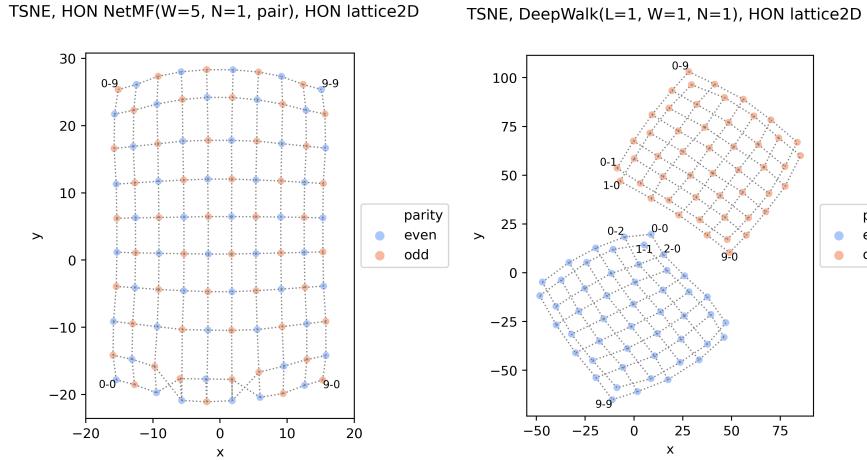


Figure 3.2: Embeddings of the HON Lattice 2D with  $\omega = 0.5$ . The left plot displays HON NetMF, exhibiting that on average horizontal edges  $\mathcal{E}_R$  are shorter than vertical edges  $\mathcal{E}_U$  due to the acceleration of horizontal movement. (These edges are both displayed as dotted lines.) The right plot displays a length-one DeepWalk (therefore without stretch effect), and the embedding splits into two clusters, as already seen in figure 2.2 for LINE-2. Hence, the dotted lines display two-step edges ( $\mathcal{E}_{2R}$  and  $\mathcal{E}_{2U}$ ) instead. For orientation, color the nodes by parity and annotate some. However, t-SNE reproduces the lattice structure only after tuning.

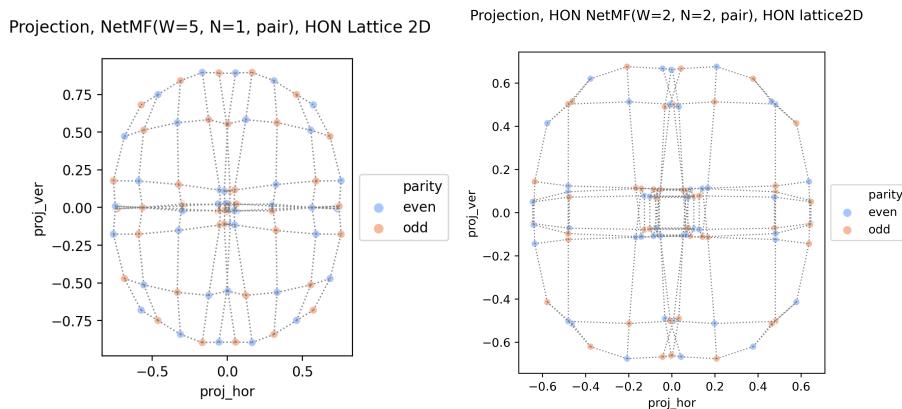


Figure 3.3: Visualization by projection of HON NetMF embeddings of the HON Lattice 2D. The stretch effect results in an asymmetry, but it is hard to tell in which direction for the right one. The same embeddings are visualized by t-SNE in figures 3.2 (left) and 4.2 (top, right).

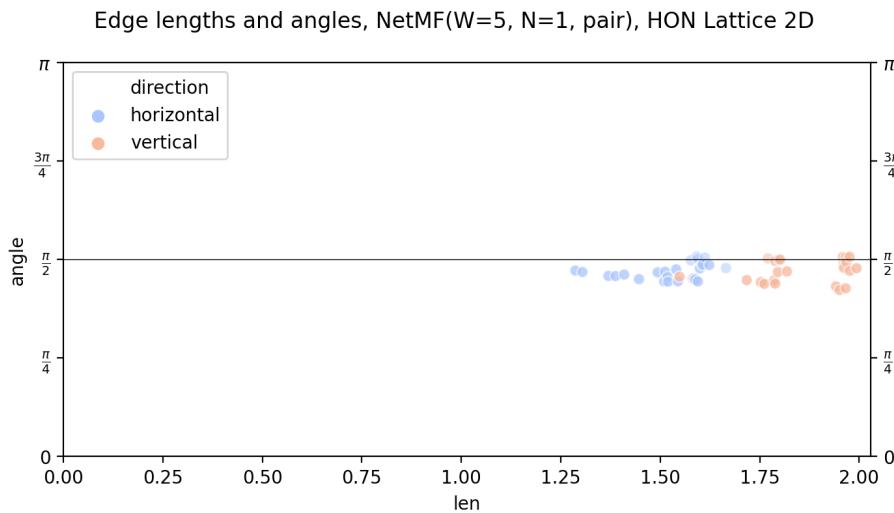


Figure 3.4: Distribution of lengths and angles (towards their average) for horizontal edges  $\mathcal{E}_R$  and vertical edges  $\mathcal{E}_U$ . This nicely illustrates the stretch of the lattice with the average edge lengths for  $\mathcal{E}_R$  (1.5) and  $\mathcal{E}_U$  (1.8) resulting in a stretch ratio of 1.2. However, the angles between the individual embedding edge-differences and their averages (per direction), indicate that we should not imagine the embedded lattice's manifold as essentially a linear subspace.



## Chapter 4

# Experimental results

A review of the embeddings mentioned in § 2 reveals the following applications of embeddings ranked by popularity:

1. Classification: HONEM, DeepWalk, Node2vec, NetMF, GraRep, Verse, LINE, InfiniteWalk, GraphSAGE
2. Link prediction: HONEM, Node2vec, Hope, App, Verse, Glee
3. Graph reconstruction: HONEM, Hope, Verse, Glee
4. Visualization: HONEM, GraRep, LINE
5. Clustering: GraRep, Verse
6. Recommendation: Hope, App

Deciding for a **multiclass classification** task is inevitable, allowing to compare all FON baselines against their HON extensions on a terrain of their preference and a real dataset. Furthermore, the **visualization** task is used to preview the classification complexity and explain a synthetic dataset’s theoretical properties. We also attempted to predict the probabilities (rather than the existence) of links without success; see § B.1.

### 4.1 Visualization for a synthetic dataset

In § 3.3, we introduced a two-dimensional lattice with higher-order dynamics. By design, only HON embeddings could notice the difference in diffusion speed between horizontal and vertical movements, resulting in a stretch/compression of the lattice in the embedding space, measured by stretch ratio (3.1) or two-step stretch ratio (3.2) for HON GraRep, and assessed visually.

This experiment will demonstrate that the different HON embeddings considered all exhibit a stretch ratio bigger than one. Importantly, this works for a broad range of parameters, as we are more interested in consistency than determining the ‘best’ embedding. By comparison, the stretch ratio for FON embeddings deviates from one only due to stochastic variability because FON embeddings do not utilize the higher-order probabilities.

We start with investigating embeddings with ‘pair of nodes’ interactions before looking at issues specific to embeddings with ‘mixed’ interactions.

### 4.1.1 Network

This experiment utilizes the HON Lattice 2D from § 3.3 of size  $10 \times 10$  with parameter  $\omega = 0.5$ . For  $\omega = 0$ , the resulting network is equivalent (w.r.t. the random walks generated) to a FON, which we use for comparison.

The symmetry concerning a reversal in the path’s distribution masks differences in directed networks between HON NetMF and HON DeepWalk.

### 4.1.2 Baselines

HONEM, as well as the HON variants of DeepWalk, Node2vec, NetMF, and GraRep, are all applied to this network, and we evaluate the stretch ratio, see Equations (3.1) and (3.2), on a grid of parameters.

- `dimension R`  $\in \{16, 32, 64, 100\}$
- `negative N`  $\in \{1, \dots, 20\}$
- `window_size W`  $\in \{2, \dots, 20\}$ ; for GraRep, the parameter `num_steps` plays a role similar to `window_size`. Note that random walks only utilize second-order probabilities if  $W \geq 2$ .
- `factor=` $\frac{\text{walk\_length}}{\text{window\_size}}$   $\in \{1, 2, 4\}$ ; only for DeepWalk and Node2vec.

To investigate ‘mixed’ interactions for the other experiment, we use HON NetMF with `pairwise=False` and parameters decided by the ‘pairs of nodes’ case.

### 4.1.3 Results

**Calculate the stretch** Figure 4.1 displays the grid search outcome over the parameters and visualizes it by fixing a parameter and aggregating the others to get the average, minimum, and maximum. For HON DeepWalk and HON GraRep, the lower bounds for the (one- or two-step) stretch ratios are larger than one. For NetMF, this holds for a small number of negative samples  $N$  only, as explained below. HONEM has only one parameter, the `dimension R`, and the stretch ratio is larger than one for all with a maximum at  $R = 16$  [not shown]. Therefore, we conclude that the stretch effect holds for a wide range of parameters, as predicted theoretically.

**Display the stretch** Next, we investigate the stretch effect visually and would like to compare HONEM with the other HON and FON embeddings. We will use `dimension R = 16` and `window_size W = 2`, as indicated by comparing the concentration around the start node of NetMF’s PMI (figure 4.3) and HONEM’s neighborhood (figure 4.6, left).

For the HON embeddings, we use their optimal parameters. However, the stretch ratios for FON embeddings deviate from one only due to blind chance, so we apply the corresponding HON parameters instead. Hence, it is fair to use the smallest `window_size` possible ( $W = 2$ ) because figure 4.1 indicates that the stretch ratio increases with  $W$  (for small values).

Figure 4.2 displays the t-SNE visualizations of NetMF, GraRep, and DeepWalk for  $W = 2$  — each in a FON and a HON variant. Moreover, figure 4.6

(right) displays the t-SNE visualization of HONEM, including the stretch. Note that FON and HON use the same parameters — optimized for HON. These figures confirm that the stretch effect exists for HON while absent for FON.

**Mixed case** The stretch effect also exists for ‘mixed’ interactions; see figure 4.5. Moreover, in § 3.2.5, we discussed the *hierarchical structure of walks* and hypothesized about  $T_{s_{2..k} \rightarrow t}^{(k-1)}$  being a weighted average of  $T_{s_{1..k} \rightarrow t}^{(k)}$  if stationary (which the HON Lattice 2D is). Despite the transformations by the embedding and t-SNE, the relationships due to the hierarchical structure are remarkably well represented.

#### 4.1.4 Explanations

**Negative** Noise contrastive estimation (NCE) — and therefore, negative sampling — compares the probability of a feature to a null-model, the network’s (first-order) stationary distribution. Increasing the number of negative samples, therefore, accentuates outstanding features (i.e.,  $p(c|w)/p_N(c)$  large). For small networks,  $p_N(c)$  is already relatively large, requiring a small  $N$ . Since matrix factorization operates on the positive  $\log(PMI/N)$  values, it amplifies the accentuation of outstanding features even more.

Unlike NetMF, GraRep does not average over different step sizes resulting in a higher concentration of PMI, resulting in an optimal  $N = 3$  for HON GraRep compared to  $N = 1$  for HON NetMF or HON DeepWalk. (See the upper bounds in figure 4.1; for  $W = 2$  in figure 4.2,  $N = 2$  is optimal for HON NetMF and HON DeepWalk.)

**Window size** The parameter `window_size` is the random walk length for NetMF and controls the concentration of the visitation probabilities (and PMI) around the starting node.

Figure 4.3 illustrates the concentration of PMI around a starting point. For small window sizes, the PMI is concentrated around the starting node, and the horizontal speed-up causes the region with high PMI (e.g.,  $> 1$ ) to grow faster horizontally than vertically. However, vertical growth can later catch up since asymptotically  $p(c|w)$  converges to the stationary distribution and  $PMI \rightarrow 1$ ; see figure 4.4.

The asymmetry of the concentration of PMI around the starting node causes the stretch in the embedding. Combining this explains the stretch ratio dependency on the `window_size`; see figures 4.1 and 4.4.

The stretch ratio for HON NetMF seems to prefer even `window_sizes`; see the zig-zag line in the upper-right of figure 4.1. In § 2.3.2, we already noticed that a pair of adjacent nodes have opposite parity. Hence, the nodes visited in a walk (excluding the start node) have balanced parity for even walk lengths and unbalanced for odd ones. Using the two-step stretch ratio (3.2) instead, which is suitable for embeddings splitting into two groups (e.g., GraRep), confirms that parity caused the zig-zag line [not shown].

**Node2vec** The simulation of Node2vec included bias parameters  $p = 1/q \in \{.5, 1, 2, 4\}$ . For  $p = 2$  and  $p = 4$ , diffusion is boosted in all directions, resulting

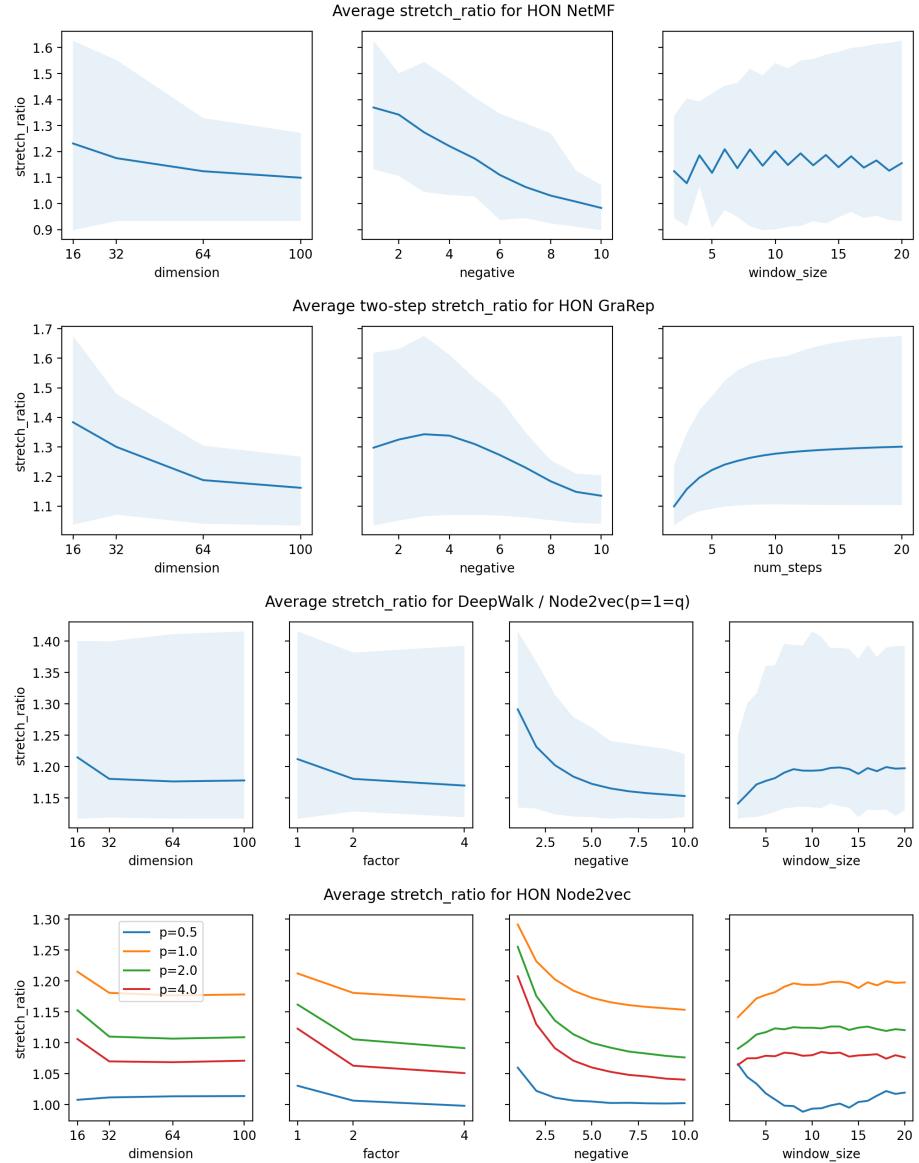


Figure 4.1: The stretch ratio for embeddings of HON Lattice 2D is calculated on a grid of parameters. Note that the minimal `window_size` or `num_steps` must be two to utilize the second-order dynamics. For visualization purposes, we fix one parameter and aggregate over the others. The individual plots display the mean as well as the range (except for Node2vec). Using a small number of negative samples,  $N$  seems crucial to get a large stretch ratio. Moreover, the biased random walks of HON Node2vec do not increase the stretch ratio compared to DeepWalk.

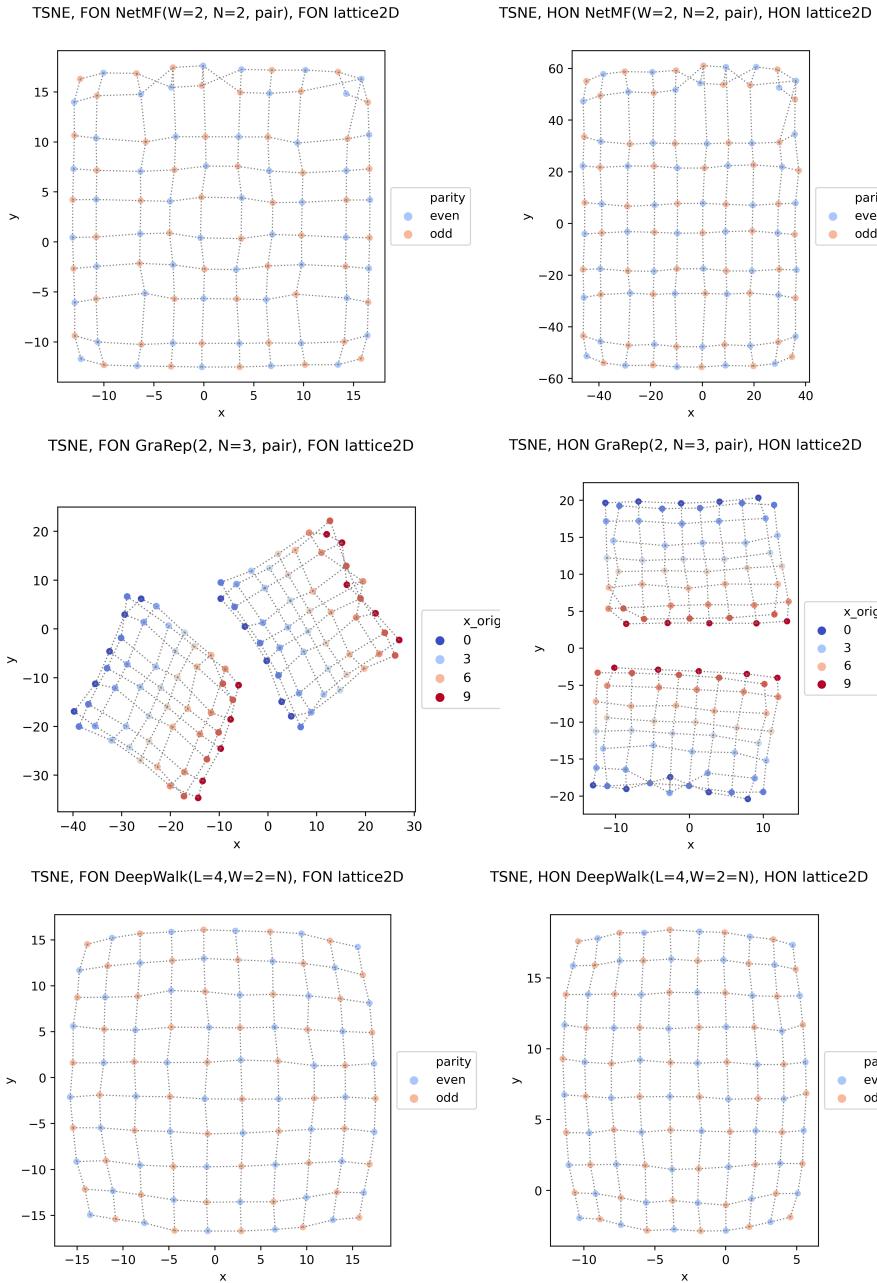


Figure 4.2: Compare of FON vs. HON embeddings. For the FON embeddings, the stretch ratio is one. The HON embeddings have stretch ratios of 1.34 (NetMF), 1.24 (GraRep), or 1.15 (DeepWalk). The `window_size` ( $W = 2$ ) was chosen for its similarity to HONEM (with a stretch ratio of 1.33, see figure 4.6). According to figure 4.1, the stretch effect is more pronounced for larger `window_sizes`. Recall that the stretch ratio is directly calculated on the embedding, while the t-SNE visualization required some manual parameter tuning.

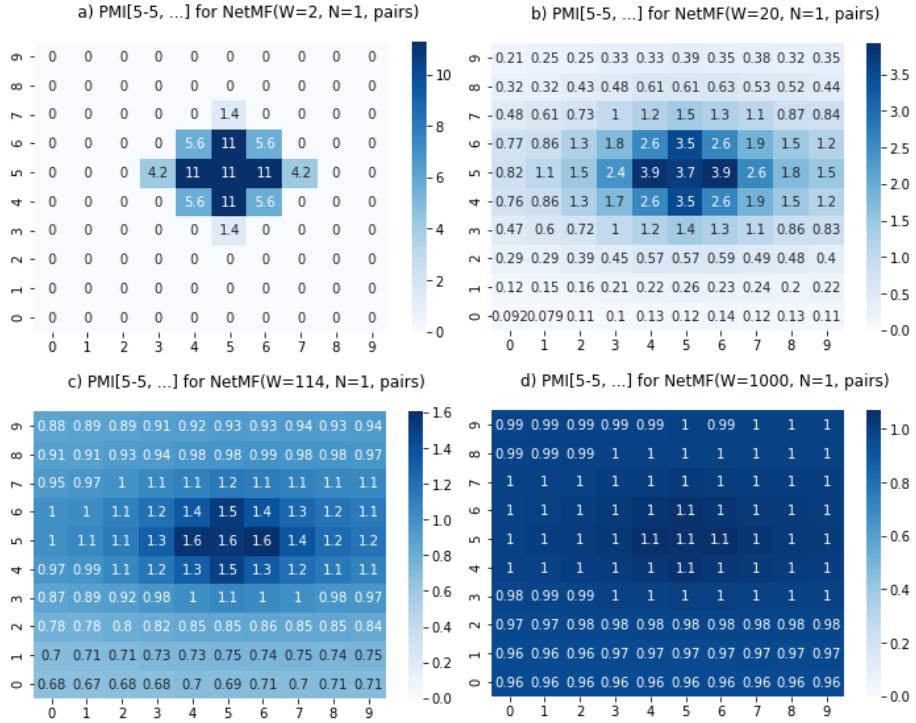


Figure 4.3: Display the concentration around the starting node ‘5-5’ of the pointwise mutual information (PMI) for different `window_sizes`  $W \in \{2, 20, 114, 1000\}$  of (pairwise) HON NetMF with one negative sample ( $N$ ). HON NetMF factorizes  $\log(PMI/N)^+$ , and PMI is defined as the random walk visitation probability divided by the stationary distribution. The second-order dynamic of the HON Lattice 2D ( $size = 10 \times 10, \omega = 0.5$ ) was constructed to speed-up horizontal and slow-down vertical transitions. Consequently, horizontal neighbors’ embeddings are, on average, closer than vertical ones which is measured by the stretch ratio, Eq. (3.1). This effect attains its maximum for window size  $W = 114$  and vanishes asymptotically for  $PMI \rightarrow 1$ .

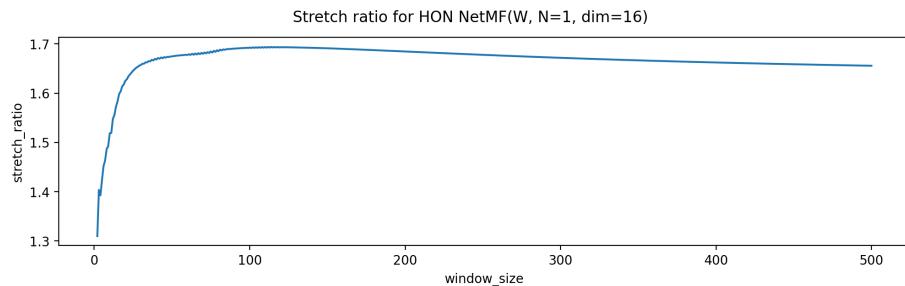


Figure 4.4: HON NetMF: stretch ratios for varying window sizes. The maximum 1.69 is at  $W = 114$ .

TSNE, HON NetMF(W=2, N=2, mix), HON lattice2D

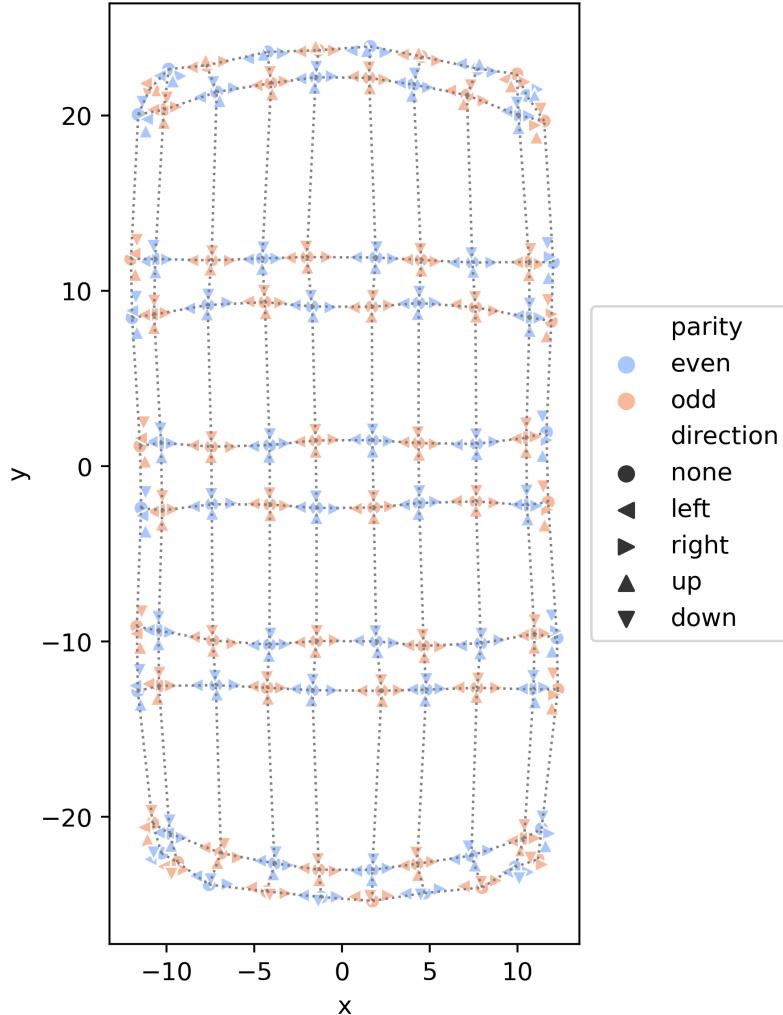


Figure 4.5: A HON NetMF embedding with ‘mixed interactions’ visualized by t-SNE displays the stretch effect (stretch ratio 1.23) — the grouping in pairs of rows disappears for higher dimensions. Several tuples (to be embedded) *end* with the same node and therefore have identical neighborhoods and similar embeddings. Since horizontal transitions are speeded up, a start tuple of nodes corresponding to a rightward movement ( $\blacktriangleright$ ) has a higher probability of moving rightward again and is therefore embedded towards its right neighbor. Conversely, starting with an upward movement ( $\blacktriangleup$ ) increases the probability of moving downwards, resulting in an embedding closer to its downward neighbor. For leftward ( $\blacktriangleleft$ ) and downward ( $\blacktriangledown$ ) movement, similar patterns emerge. Finally, starting with a single node ( $\bullet$ ) results in a mix of the above.

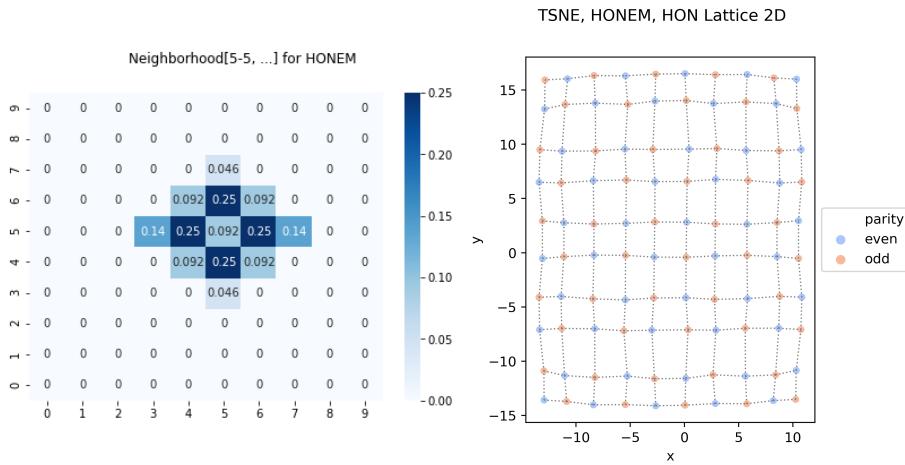


Figure 4.6: HONEM: neighborhood for the node ‘5-5’ and t-SNE visualization of the embedding (with dimension 16). The stretch ratio is 1.33.

in lower stretching compared to  $p = 1$  (i.e., DeepWalk). For  $p = 0.5$ , the stretching effect no longer appears consistently for different parameters. Hence, this specific application does not benefit from the additional tuning parameters for the bias of Node2vec compared to DeepWalk.

**Factor** For DeepWalk, the word–context pairs are generated from a random walk of length `factor*window_size`. While NetMF generally approximates DeepWalk, there are minor differences: first, nodes near the ends of a walk contribute to fewer word–context pairs, and second, the distribution of the words  $p(w)$  varies from almost uniform (small `walk_length`) to stationary (large `walk_length`); see also [25, Remark 1].

I decided to allow `walk_length` different from `window_size` in the simulation but will not investigate this further, as the HON Lattice 2D is too well behaved to exploit these differences.

## 4.2 Multiclass Classification for SocioPatterns

Apply the network embeddings to a real dataset for multi–class classification (see § 4.2.1). We will focus on comparing HON embeddings against their FON equivalents. Hence, use default settings for the classifiers and watch for a consensus instead of perfectly fitting a single classifier.

### 4.2.1 Methodology

Multi–class classification aims to predict the labels for the nodes of a network from its embedding.

In practice, some nodes are labeled, and we have to estimate the labels for the others. However, for a simulation study, all nodes are labeled, and we split the nodes into training and test sets to predict the labels of the nodes in the test

set based on the ones in the training set. Note that we still use all the nodes to calculate the embedding.

In  $k$ -fold cross-validation, we randomly split the nodes into  $k = 5$  sets  $\mathcal{V}_1, \dots, \mathcal{V}_k$  of roughly equal size and repeatedly ( $i = 1 \dots k$ ) estimate the labels for the test set of nodes  $\mathcal{V}_i$  based on the other nodes ( $\mathcal{V} - \mathcal{V}_i$ ) as the training set. Eventually, we have collected the predicted label and the probabilities for each label, from which we could calculate the *confusion matrix*, the *F1-scores*, or the *receiver operating characteristic* (ROC) curve. The confusion matrix counts the number of combinations of true and predicted label pairs.

The F1-score for binary classification is the *harmonic mean of precision and recall*,  $(\text{F1-score})^{-1} = \frac{1}{2}(\text{precision}^{-1} + \text{recall}^{-1})$ , with values between zero and one (bigger is better). Recall (a.k.a. true-positive rate, TPR) is a conditional probability of predicted given the true label. The precision is similar, with true and predicted labels exchanged. For multi-class classification, the difference between F1-micro and F1-macro is, whether we aggregate before or after the harmonic mean.

The ROC curve graphically depicts a fundamental trade-off in hypothesis testing between type 1 and type 2 errors. It uses the vocabulary **true-positive rate (TPR)** and **false-positive rate (FPR)**, which correspond to  $1 - P(\text{type 1 error})$  and  $P(\text{type 2 error})$ , respectively. The area under the ROC curve describes the toughness of a classification task. Generalizing from binary to multi-class classification, we can calculate the ROC curve for one class against either another class (one vs. one) or any other class (one vs. rest) and decided for the latter.

Moreover, we may repeat the  $k$ -fold cross-validation to reduce variability (due to random split) or get error bounds.

### 4.2.2 Networks

For the classification task, we need datasets with labels and a HON structure. SocioPatterns.org provides temporal network data, and I decided to use the following three ones.

**Primary school temporal network data [7, 31]** The labels name teachers or classes (for children).

**Contacts in a workplace [11]** The labels contain departments.

**Hospital ward dynamic contact network [36]** The labels identify patients and different kinds of staff.

All these record interactions with a 20 seconds resolution. Hence the next step is extracting<sup>1</sup> causal paths (named ‘walks’ in this thesis). There is a tuning parameter  $\Delta$  for the path extraction, which specifies the maximal time-difference allowed for two temporal edges to be included in a causal path.

Next, we fit a multi-order network (MON, [29]) and estimate the higher-order Markov model’s optimal order  $K$ . Finally, export the resulting models using the same data format as BuildHON+ [39] for later use in our software.

---

<sup>1</sup>Similar to task 1.4 of <https://ingoscholtes.github.io/kdd2018-tutorial/>

The maximal Delta reflects my computational resources, and I will continue with the following<sup>2</sup>:

- PrimarySchool,  $\Delta = 20$  sec,  $K = 2$
- Workplace,  $\Delta \in \{1 \text{ min}, 5 \text{ min}, 10 \text{ min}, 20 \text{ min}, 30 \text{ min}\}$ ,  $K = 2$
- HospitalWard,  $\Delta \in \{20 \text{ sec}, 40 \text{ sec}, 1 \text{ min}\}$ ,  $K = 3$  (or 2 for  $\Delta=20$  sec)

### 4.2.3 Baselines

Given the task of classifying nodes from their embeddings, we consider only HON embeddings from § 3.2 with ‘pairs of nodes’ interactions. (There is simply no use for embeddings of tuples of nodes longer than one.)

Each of the following embeddings consists of a HON and a FON variant, distinguished by their respective prefixes. Both variants use the same implementation, but the the FON variant is applied to a network converted from HON to FON (by discarding higher-order probabilities).

**DeepWalk** generates 100 walks per node, each of length  $L = 80$ . The `window_size` is  $W = 10$ ; Optimization uses negative sampling with one iteration.

**GraRep(S=4)** generates 4 embeddings for the different step sizes and concatenates them together. Therefore, each individual embedding uses only dimension 16, i.e., a fourth of the dimension.

**NetMF(W=10)** generates an embedding with `window_size`  $W = 10$ .

**NetMF(W=2)** generates an embedding with `window_size`  $W = 2$  because both PrimarySchool, Workplace, and HospitalWard ( $\Delta = 20$ ) have maximal order  $K = 2$ .

**Node2vec(p=0.5)** uses the bias parameters  $p = 0.5, q = 2$  and generates the embedding exactly as DeepWalk.

**Node2vec(p=2)** uses the bias parameters  $p = 2, q = 0.5$  and generates the embedding exactly as DeepWalk.

**HONEM** HONEM has neither parameters (besides the `dimension`) nor a pair of FON/HON variants.

Each embedding has `dimension`<sup>3</sup>  $R = 64$ , uses one negative sample ( $N = 1$ ), and uses only a ‘pair of nodes’ interaction.

The simulation uses the following classifiers with default settings (and for reproducibility `random_state=1`):

**LR** Logistic regression (`sklearn.linear_model.LogisticRegression`)

**RF** Random forest (`sklearn.ensemble.RandomForestClassifier`)

**SVC** Support vector classification (`sklearn.svm.SVC`)

---

<sup>2</sup>Due to the larger size than a dataset from a similar domain, I considered both ‘Highschool dynamic’ and ‘Workplace 2nd’ as a lower priority. Other datasets were too complicated or lacked metadata.

<sup>3</sup>I compared dimensions 32 and 64 earlier without noticing differences and did not reassess my decision since. Instead, I later focused on comparing different  $\Delta$  values.

	f1-macro			f1-micro		
	FON	HON	HONEM	FON	HON	HONEM
Hospital 20"	(20) 52%	(17) 53%	(38) 32%	(20) 61%	(18) 63%	(37) 47%
Hospital 40"	(18) 46%	(19) 45%	(37) 30%	(18) 55%	(21) 52%	(29) 48%
Hospital 1'	(16) 47%	(21) 42%	(35) 31%	(19) 55%	(21) 53%	(23) 50%
PrimarySchool	(26) 77%	(11) 83%	(38) 41%	(26) 82%	(11) 87%	(38) 45%
Workplace 1'	(14) 77%	(23) 73%	(38) 57%	(15) 88%	(22) 87%	(38) 74%
Workplace 5'	(18) 74%	(19) 73%	(38) 56%	(17) 86%	(20) 85%	(38) 72%
Workplace 10'	(14) 78%	(23) 73%	(38) 55%	(18) 85%	(19) 85%	(38) 73%
Workplace 20'	(15) 79%	(22) 74%	(38) 52%	(21) 85%	(16) 86%	(38) 68%
Workplace 30'	(10) 85%	(27) 73%	(38) 55%	(11) 90%	(26) 87%	(38) 73%

Table 4.1: Comparing *average* F1-scores (and their ranks in parentheses) among different embeddings and classifiers, with the embeddings grouped by FON, HON, and HONEM. For each combination of a dataset,  $\Delta$ , and type of score (macro or micro), we underline the values of the best group, as well as the best  $\Delta$  per dataset. An undisputed winner exists for PrimarySchool and Workplace ( $\Delta=30'$ ), while Hospital[Ward] ( $\Delta=20''$ ) is close and also has lower F1-scores.

#### 4.2.4 Results

We calculate the F1-scores (macro and micro) of every combination of embedding and classifier, display them as a scatter plot in figure 4.7, and provide a summary in table 4.1; see appendix A for the raw results.

**PrimarySchool** For  $\Delta=20$  sec (the only one available for this dataset), HON embeddings are superior to their FON counterparts, and an average rank of 11 among 18 items is exceptional — no FON embedding made it into the top ten. Moreover, the two F1-scores are highly correlated. HONEM performs worst by a huge margin.

**Workplace** For  $\Delta=30$  min, the FON embeddings outperform their HON counterparts in an equally (or even slightly more) resounding success as the other way around for PrimarySchool.

Things get closer for other  $\Delta$  values, and the two F1-scores agree less, with FON still winning for F1-macro and a tie between FON and HON for F1-micro. HONEM performs worst by a large margin.

**HospitalWard** The lower F1-scores and larger discrepancies between them indicate that this dataset is harder to classify. FON wins for  $\Delta=40$  sec and 1 min, but for  $\Delta=20$  sec, HON achieves the highest average F1-scores (among all  $\Delta$  values).

HONEM performs still worst for the F1-macro score (though without margin) and better for the F1-micro score. Surprisingly, while both HONEM and our HON embeddings utilize the higher-order terms, HONEM performs best when the FON embeddings have an advantage over the HON ones. (However, this is only anecdotal evidence.)

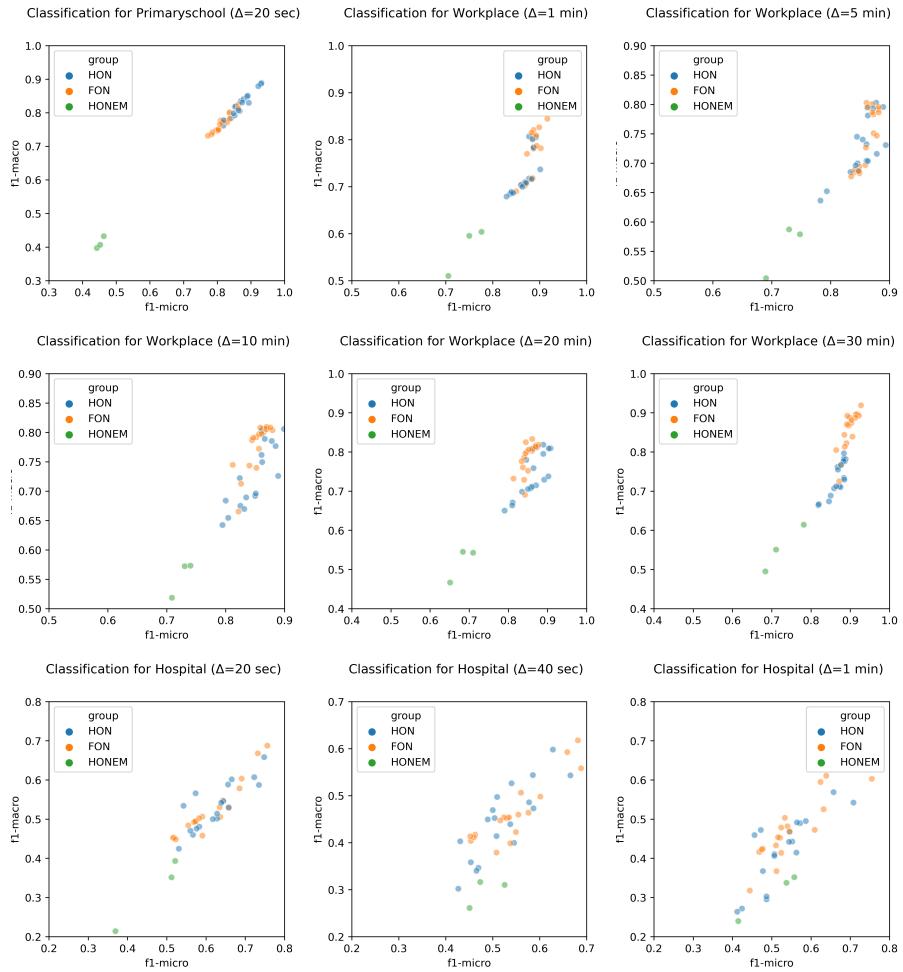


Figure 4.7: F1-scores (macro and micro) for pairs of embeddings and classifiers; see also table 4.1.

**Best embeddings** Next, we compare the classifications of a FON vs. a HON embedding for each dataset. For each, we display a t-SNE visualization of the embedding, a confusion matrix, and (one-vs-rest) ROC curves — both calculated from 10 repetitions. Hence, Figures 4.8, 4.9, and 4.10 display the best combinations of FON/HON embedding and classifier for PrimarySchool ( $\Delta=20$  sec), Workplace ( $\Delta=30$  min), and HospitalWard ( $\Delta=20$  sec), respectively; where ‘best’ means the highest sum of F1-micro and F1-macro.

#### 4.2.5 Explanations

Rosvall et al. [27] describe the effect of using HON on spreading dynamics and community detection and mention algorithms to investigate related phenomena.

**Modularity** Comparing Figures 4.8 and 4.10, we notice the classes’ cluster structure in the PrimarySchool network’s embeddings, which is not the case for

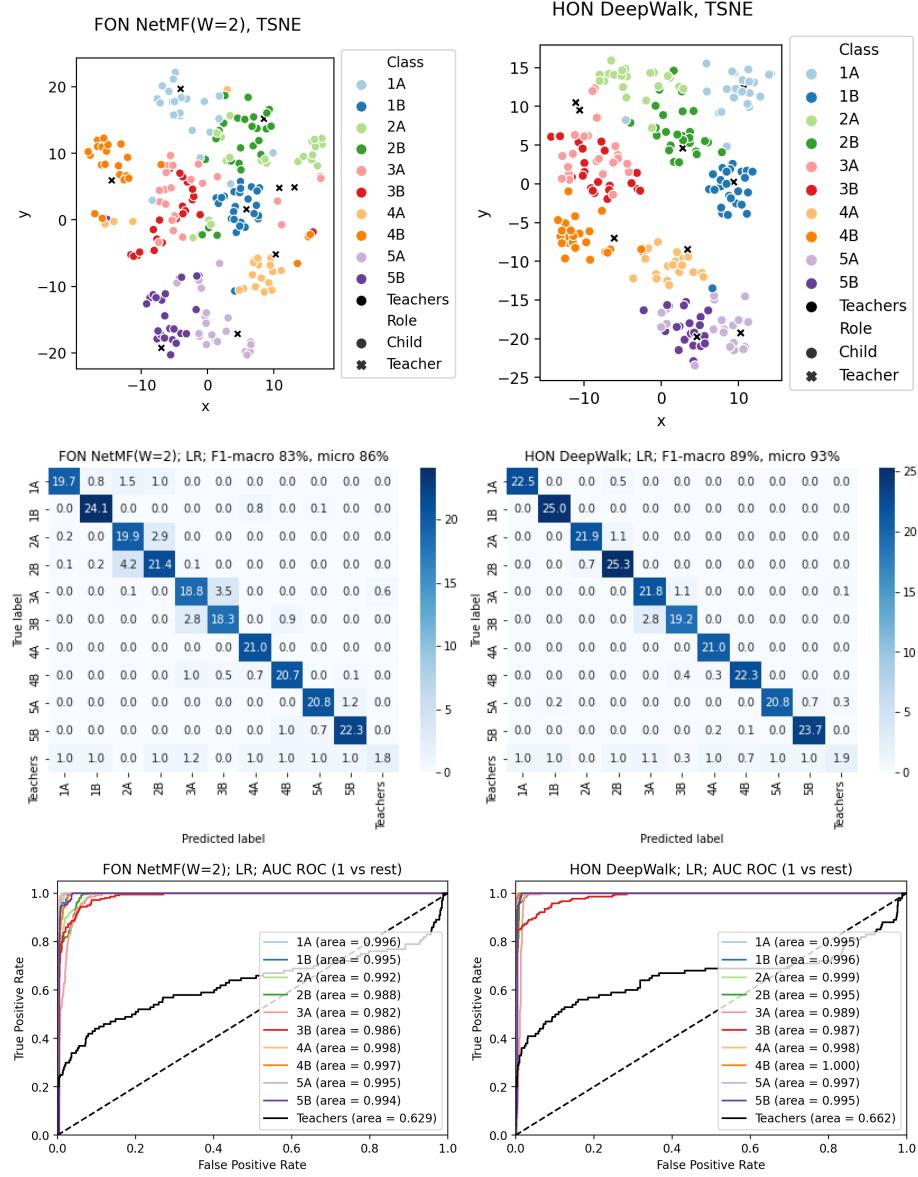


Figure 4.8: For the PrimarySchool network, combine several classifiers and embeddings and select the best pair (w.r.t. F1-micro + F1-macro) containing a FON (left) or a HON (right) embedding, respectively. These were ‘FON NetMF(W=2)’ and ‘HON DeepWalk’, both combined with logistic regression. T-SNE visualization (top) displays the clustering of children into classes, which is more pronounced for HON, confirmed by the confusion matrix (center). Both embeddings are unsuitable for classifying teachers. The area under ROC (bottom) illustrates the separability of any class from the rest.

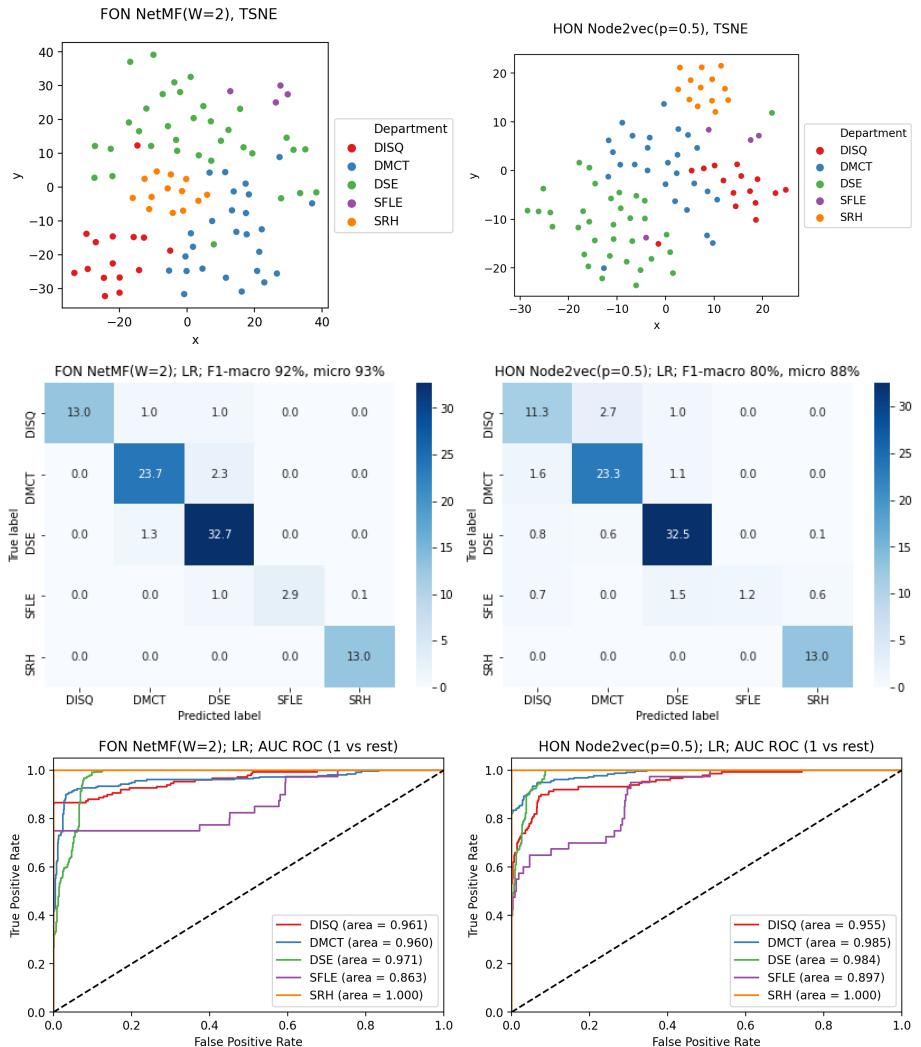


Figure 4.9: For the Workplace network with  $\Delta=30$  min, combine several classifiers and embeddings and select the best pair (w.r.t. F1-micro + F1-macro) containing a FON (left) or a HON (right) embedding, respectively. These were ‘FON NetMF(W=2)’ and ‘HON Node2vec( $p=0.5$ )’, both combined with logistic regression. T-SNE visualization (top) displays the clustering of individuals into departments. It is hard to decide visually, which one separates the clusters better, but the confusion matrix (center) reveals the differences, with FON performing better. The area under ROC (bottom) illustrates the separability of any class from the rest.

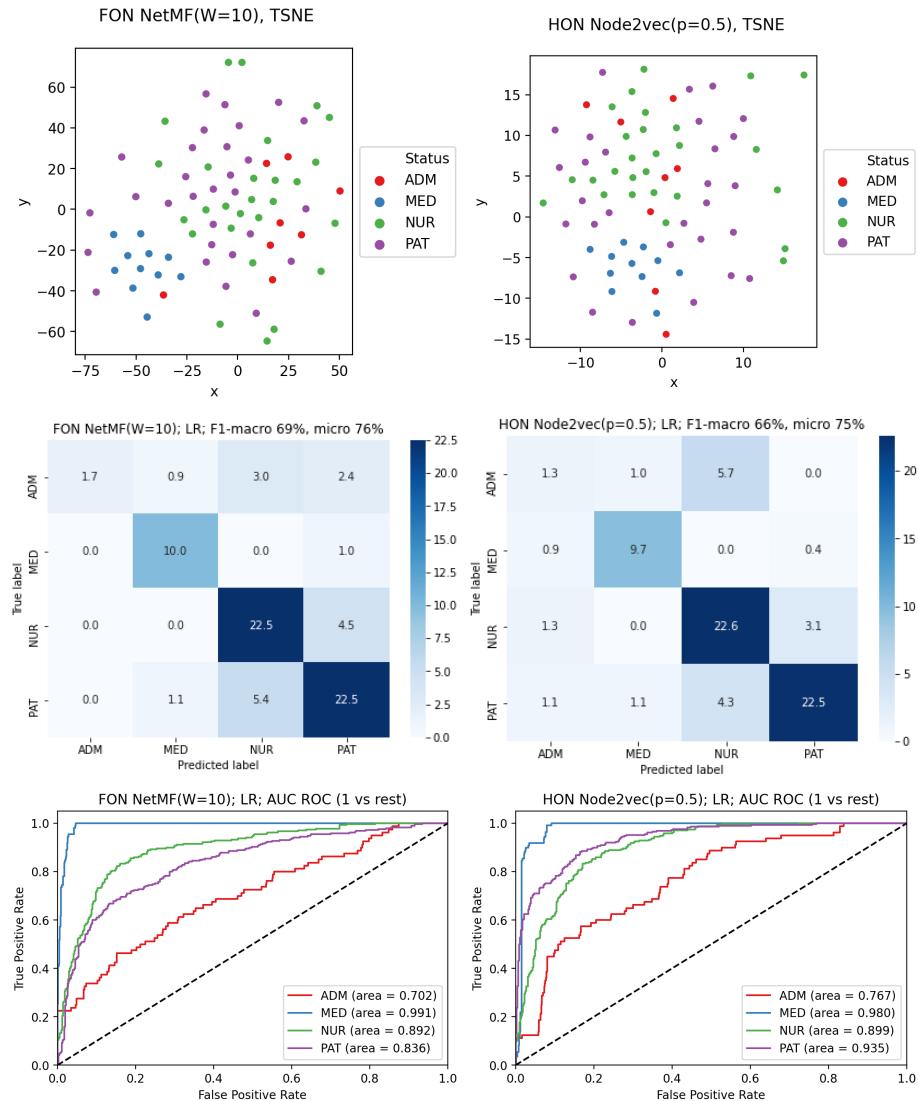


Figure 4.10: For the HospitalWard network with  $\Delta=20$  sec, combine several classifiers and embeddings and select the best pair (w.r.t. F1-micro + F1-macro) containing a FON (left) or a HON (right) embedding, respectively. These were ‘FON NetMF(W=10)’ and ‘HON Node2vec( $p=0.5$ )’, both combined with logistic regression. T-SNE visualization (top) displays the absence of clustering. FON performs slightly better than HON, see the confusion matrix (center), despite HON embeddings performing slightly better on average. The area under ROC (bottom) illustrates the separability of any class from the rest.

HospitalWard.

Literature defines the modularity for unweighted and undirected networks. The generalization to weighted [22] and directed [18, Eq. (3)] is straight forward.

Using higher-order De Bruijn graphs allows us to calculate the modularity for higher orders. Note that the adjacency matrix becomes increasingly sparse for higher dimensions, and it is essential to avoid iterating over all pairs of nodes<sup>4</sup>.

The modularities are 44% (PrimarySchool), 47% (Workplace), and 25% (HospitalWard), independent of the order. This investigation confirms that HospitalWard is more challenging to classify.

**Temporal aspects for PrimarySchool** Why do HON embeddings consistently outperform their FON counterparts for PrimarySchool? It turns out that the classification task performs poorly for teachers since they interact primarily with children rather than their peers. Examining the visualizations reveals that some teachers belong to a specific class (they are teaching), while others are separate. Hence, a teacher interacting with children of more than one class may function in a FON as a communication hub between children from different classes, which is impossible in HON, as the interaction takes place separately in time. Figure 4.11 confirms this kind of mechanism, although it is infrequent.

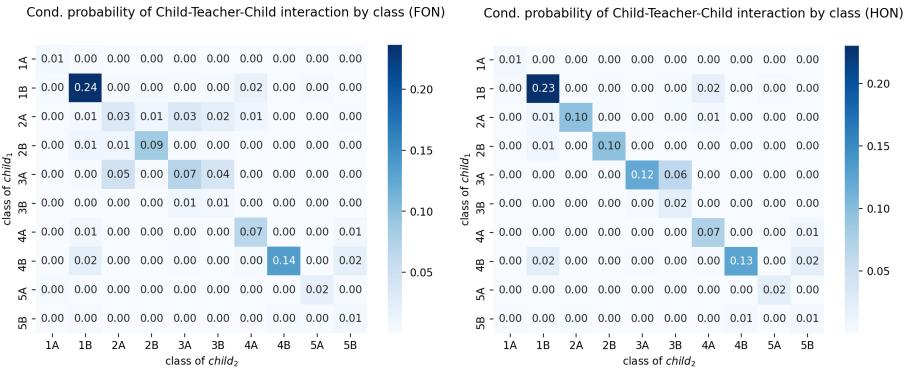


Figure 4.11: Examining interaction sequences from child to teacher and then to another child. FON ignores the temporal aspect when linking two children via the teacher, creating faulty connections between children of different classes sharing the same teacher. Observing that the diagonal contains more probability mass for HON than for FON confirms the hypothesis. The heatmaps display probabilities of classes of a pair of children. The children's pair follows the marginal distribution of random walks with the pattern  $child_1 \rightarrow teacher \rightarrow child_2$  (generated by rejection sampling from the stationary distribution). However, these walks are rare (< 0.4%).

**Entropy of Transition Probabilities** In § 3.2.5, we described the hierarchical structure of the start-tuples of a walk and demonstrated that  $T_{s_{2..k} \rightarrow t}^{(k-1)}$  is a weighted average of  $T_{s_{1..k} \rightarrow t}^{(k)}$  if the transition probabilities were *stationary*.

<sup>4</sup>Hint:  $\sum_{i,j} k_i^{in} k_j^{out} \delta_{c_i, c_j} = \sum_c (\sum_i k_i^{in} \delta_{c_i, c}) \cdot (\sum_j k_j^{out} \delta_{c_j, c})$  in [18, Eq. (3)]

Consequently,  $T_{s_{2..k} \rightarrow t}^{(k-1)}$  would tend to be less concentrated than  $T_{s_{1..k} \rightarrow t}^{(k)}$ , which we measure by their respective entropies.

Figure 4.12 confirms that higher-order probabilities are more concentrated (less entropy) than lower ones; see also [27, table 2]. This is in line with the previous observation that HON transitions are more restrictive than FON ones.

Examining transition probabilities' entropies in figure 4.12 also reveals that the PrimarySchool has higher average entropy than Workplace ( $\Delta=30$  min). Is the Workplace ( $\Delta=30$  min) insufficiently connected for a HON embedding? Unfortunately, I did not notice a pattern when comparing the Workplace network's entropies for different  $\Delta$  values [not shown].

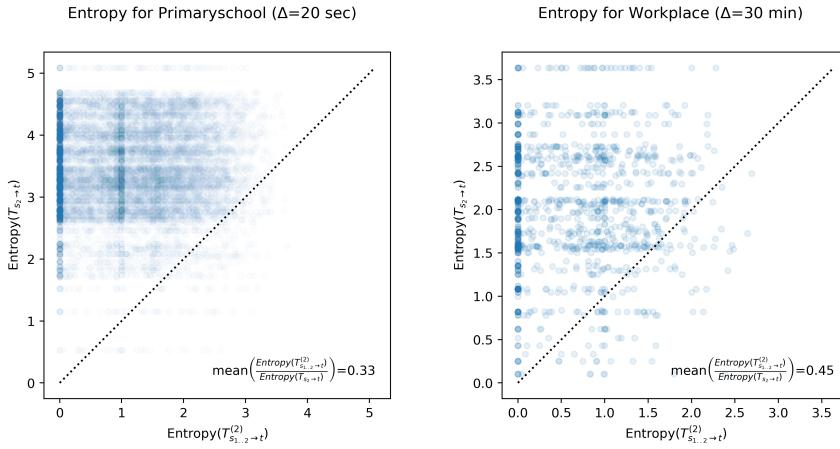


Figure 4.12: Comparing the entropy of pairs of transition probabilities  $T_{s_{1..2} \rightarrow t}^{(2)}$  and  $T_{s_2 \rightarrow t}$  for PrimarySchool ( $\Delta=20$  sec, left) and Workplace ( $\Delta=30$  min, right). The second-order transitions' entropy for *both networks* is considerably *smaller* than the first-order ones (on average). Note that PrimarySchool has higher average entropy than Workplace: 3.5 vs. 2.3 (first-order) and 1.1 vs. 0.7 (second-order), respectively.



# Chapter 5

# Conclusion and Discussions

## 5.1 Summary and Reflection

Representation learning for first-order networks (FON) is already an established research area that produced various embedding methods, which provide great starting points for higher-order network (HON) adaptations. For clarification, HON refers to a network with random walks following a higher-order Markov model.

An important early assumption, which shaped this thesis was that ‘pairwise interactions are not enough’; see § 3.1.2. This originates from objective two, which demanded generating paths from a known embedding. There is no way to decode higher-order transition probabilities from pairwise interactions (e.g.,  $\langle \vec{w}, \vec{c} \rangle$ ) in the skip-gram model with nodes  $w$  and  $c$ ) without losing information. Hence any approximation of  $T_{s_1..k \rightarrow s_{k+1}}^{(k)}$  must formally depend on all nodes of the tuple  $s_{1..k+1}$ .

In the end, it turned out to be sufficient to replace random walks by their HON counterparts (see § 3.1.1), and there was a lack of applications for the advanced approaches.

After setting the direction initially and analyzing the problem, we had plenty of options, a luxury which complicated decisions. The thesis reflects this by explaining the choices made. We have to finish with time running out and summarize some highlights, although there are still more interesting questions to explore.

**HON Lattice 2D** This synthetic network is useful for understanding the skip-gram model and explaining the impact of second-order dynamics. When implementing an embedding, it is valuable for debugging to have a theoretical understanding of the expected outcome when embedding a lattice.

**Classification** To publish a new method, one has to demonstrate the new method’s usefulness by surpassing existing methods in a comparison experiment. Fortunately, the HON embeddings outperform their FON counterparts for the PrimarySchool dataset, although the result for the Workplace dataset spoils the good impression.

However, the fundamental differences between FON and HON random walks are underlying circumstances determining whether HON or FON embeddings are a better tool. Ideally, we could determine the differences resulting from choosing a HON or FON random walk and assess them by subject-specific criteria.

**Probability Prediction** This experiment was the only opportunity to evaluate the approaches based on the above mentioned ‘pairwise is not enough’ assumption, and is therefore mentioned despite its failure. Moreover, it raises an important question: should an embedding generalize or reproduce a network?

**HONEM** Theoretical considerations reveal weaknesses of HONEM, but we acknowledge that fixing these weaknesses would result in an entirely different embedding, which is computationally more expensive. In our classification experiment, HONEM performed poorly.

## 5.2 Threats to Validity and Limitations

**Implications of using HON** We noticed in § 4.2.5 that HON might prohibit walks, which are valid in FON. This may remove unwanted connections between groups or break up clusters. Hence, we cannot conclude that HON embeddings are better than FON ones — even if the HON differs significantly from the FON, which is therefore incorrect. Criteria to assist this decision would be useful.

**Stationary Distribution** For a FON random walk, setting the start distribution  $P_{\overline{S_1}}$  equal to the stationary distribution  $\pi$  is sufficient for a stationary random walk. Therefore, it is natural for NetMF to set  $p_N = \pi$ .

The maximal order transition probabilities  $(T_{s_{1..K} \rightarrow t}^{(K)})$  solely determine the stationary distribution of HON random walks. For a stationary distribution of the walk, we need to choose all other probabilities (start distribution and transition probabilities of order  $k$  smaller than the maximal order  $K$ ) accordingly; see § 2.1.3, estimating probabilities. Should we somehow ignore the lower order probabilities when generating random walks? Furthermore, should we replace the FON stationary distribution by a marginal of the HON one, given the maximal order terms are the hardest to estimate?

No big issue; I have just decided to stick with the FON stationary distribution without considering alternatives.

## 5.3 Future Work

**Fast SGD-based Implementations** For a thesis, matrix factorization based embeddings are rewarding, as we do not need to worry about convergence. My experimental SGD implementation (§ 3.2.5) lacks performance, putting any thoughts to include HON App/Verse in the evaluations to an end; see also [26] about tuning gensim’s Word2vec for Python. However, for large datasets, fast SGD implementations are essential. A HON App/Verse implementation is currently missing, and HON DeepWalk/Node2vec relies on Word2vec (which treats word and context symmetrically).

**Applications for Embeddings Beyond Pairs of Nodes** Is there any real-world application for embeddings with ‘mixed’ (or ‘pairs of walks’) interactions? We only attempted probability prediction (§ B.1) without success.



## Appendix A

# Results for the Classification task

The following tables contain the detailed results corresponding to figure 4.7 and table 4.1

	f1-macro			f1-micro		
	LR	RF	SVC	LR	RF	SVC
FON DeepWalk	0.772	0.742	0.749	0.815	0.786	0.797
FON GraRep(S=4)	0.777	0.771	0.746	0.824	0.829	0.801
FON NetMF(W=10)	0.799	0.768	0.777	0.836	0.817	0.810
FON NetMF(W=2)	0.825	0.813	0.801	0.863	0.862	0.837
FON Node2vec( $p=0.5$ )	0.782	0.751	0.748	0.837	0.803	0.803
FON Node2vec( $p=2$ )	0.766	0.732	0.736	0.808	0.772	0.783
HON DeepWalk	0.889	0.799	0.848	0.931	0.848	0.888
HON GraRep(S=4)	0.784	0.792	0.762	0.841	0.852	0.818
HON NetMF(W=10)	0.835	0.805	0.818	0.870	0.867	0.852
HON NetMF(W=2)	0.841	0.830	0.819	0.879	0.894	0.855
HON Node2vec( $p=0.5$ )	0.887	0.808	0.851	0.930	0.863	0.890
HON Node2vec( $p=2$ )	0.879	0.778	0.831	0.922	0.819	0.874
HONEM	0.397	0.407	0.432	0.442	0.452	0.463

Table A.1: F1-scores for the PrimarySchool dataset. Note that the best performing non-HON methods are ‘FON NetMF(W=2)’ with logistic regression and F1-macro 0.825 and ‘FON GraRep(S=4)’ with random forest and F1-micro 0.873, both have rank 12. HONEM performs worst by large margin.

	f1-macro			f1-micro		
	LR	RF	SVC	LR	RF	SVC
FON DeepWalk	0.782	0.809	0.770	0.902	0.886	0.873
FON GraRep(S=4)	0.718	0.707	0.718	0.884	0.872	0.885
FON NetMF(W=10)	0.781	0.702	0.705	0.887	0.868	0.867
FON NetMF(W=2)	0.803	0.690	0.705	0.885	0.850	0.862
FON Node2vec(p=0.5)	0.845	0.821	0.826	0.916	0.887	0.899
FON Node2vec(p=2)	0.787	0.815	0.809	0.893	0.883	0.892
HON DeepWalk	0.806	0.686	0.716	0.892	0.843	0.883
HON GraRep(S=4)	0.717	0.706	0.710	0.877	0.867	0.868
HON NetMF(W=10)	0.782	0.686	0.684	0.887	0.838	0.837
HON NetMF(W=2)	0.784	0.689	0.679	0.887	0.840	0.829
HON Node2vec(p=0.5)	0.807	0.703	0.737	0.877	0.860	0.901
HON Node2vec(p=2)	0.801	0.700	0.708	0.885	0.863	0.871
HONEM	0.510	0.604	0.595	0.705	0.776	0.750

Table A.2: F1-scores for the Workplace dataset with  $\Delta=1$  min. FON performs best. HONEM performs worst by large margin.

	f1-macro			f1-micro		
	LR	RF	SVC	LR	RF	SVC
FON DeepWalk	0.747	0.803	0.792	0.878	0.861	0.880
FON GraRep(S=4)	0.677	0.697	0.686	0.835	0.859	0.848
FON NetMF(W=10)	0.751	0.684	0.683	0.873	0.839	0.849
FON NetMF(W=2)	0.783	0.689	0.694	0.873	0.841	0.849
FON Node2vec(p=0.5)	0.786	0.801	0.796	0.882	0.871	0.882
FON Node2vec(p=2)	0.727	0.795	0.786	0.860	0.863	0.871
HON DeepWalk	0.803	0.740	0.716	0.877	0.854	0.878
HON GraRep(S=4)	0.696	0.704	0.687	0.842	0.861	0.847
HON NetMF(W=10)	0.781	0.700	0.652	0.863	0.846	0.793
HON NetMF(W=2)	0.745	0.685	0.636	0.845	0.834	0.783
HON Node2vec(p=0.5)	0.796	0.732	0.731	0.889	0.861	0.893
HON Node2vec(p=2)	0.793	0.795	0.703	0.872	0.862	0.863
HONEM	0.504	0.579	0.587	0.690	0.748	0.729

Table A.3: F1-scores for the Workplace dataset with  $\Delta=5$  min. FON performs best. HONEM performs worst by large margin.

	f1-macro			f1-micro		
	LR	RF	SVC	LR	RF	SVC
FON DeepWalk	0.797	0.787	0.806	0.861	0.845	0.867
FON GraRep(S=4)	0.713	0.772	0.665	0.826	0.857	0.822
FON NetMF(W=10)	0.800	0.745	0.792	0.864	0.812	0.852
FON NetMF(W=2)	0.797	0.744	0.740	0.857	0.840	0.852
FON Node2vec(p=0.5)	0.805	0.808	0.809	0.870	0.859	0.870
FON Node2vec(p=2)	0.804	0.791	0.808	0.879	0.847	0.876
HON DeepWalk	0.785	0.762	0.692	0.878	0.861	0.850
HON GraRep(S=4)	0.675	0.689	0.670	0.825	0.835	0.832
HON NetMF(W=10)	0.806	0.722	0.655	0.863	0.824	0.804
HON NetMF(W=2)	0.802	0.684	0.642	0.861	0.800	0.795
HON Node2vec(p=0.5)	0.806	0.789	0.726	0.899	0.866	0.889
HON Node2vec(p=2)	0.777	0.749	0.696	0.885	0.862	0.851
HONEM	0.519	0.573	0.572	0.709	0.740	0.730

Table A.4: F1-scores for the Workplace dataset with  $\Delta=10$  min. FON performs better for F1-macro while FON and HON are on par for F1-micro. HONEM performs worst by large margin.

	f1-macro			f1-micro		
	LR	RF	SVC	LR	RF	SVC
FON DeepWalk	0.811	0.786	0.818	0.871	0.840	0.877
FON GraRep(S=4)	0.732	0.752	0.691	0.813	0.850	0.842
FON NetMF(W=10)	0.801	0.817	0.803	0.848	0.870	0.861
FON NetMF(W=2)	0.833	0.761	0.729	0.861	0.837	0.840
FON Node2vec(p=0.5)	0.808	0.796	0.814	0.862	0.843	0.870
FON Node2vec(p=2)	0.776	0.825	0.806	0.834	0.845	0.852
HON DeepWalk	0.819	0.780	0.711	0.889	0.846	0.860
HON GraRep(S=4)	0.708	0.729	0.715	0.858	0.891	0.871
HON NetMF(W=10)	0.813	0.664	0.671	0.874	0.810	0.811
HON NetMF(W=2)	0.808	0.699	0.650	0.861	0.835	0.790
HON Node2vec(p=0.5)	0.809	0.795	0.738	0.908	0.889	0.902
HON Node2vec(p=2)	0.809	0.759	0.705	0.903	0.864	0.850
HONEM	0.467	0.543	0.545	0.651	0.710	0.684

Table A.5: F1-scores for the Workplace dataset with  $\Delta=20$  min. FON performs better for F1-macro while FON and HON are on par for F1-micro. HONEM performs worst by large margin.

	f1-macro			f1-micro		
	LR	RF	SVC	LR	RF	SVC
FON DeepWalk	0.896	0.873	0.887	0.914	0.902	0.913
FON GraRep(S=4)	0.805	0.768	0.725	0.864	0.876	0.872
FON NetMF(W=10)	0.897	0.839	0.893	0.918	0.905	0.922
FON NetMF(W=2)	0.919	0.813	0.822	0.927	0.886	0.890
FON Node2vec(p=0.5)	0.888	0.868	0.881	0.899	0.896	0.903
FON Node2vec(p=2)	0.844	0.892	0.871	0.885	0.892	0.891
HON DeepWalk	0.782	0.783	0.734	0.888	0.884	0.884
HON GraRep(S=4)	0.712	0.710	0.689	0.872	0.875	0.850
HON NetMF(W=10)	0.776	0.664	0.707	0.886	0.818	0.859
HON NetMF(W=2)	0.762	0.674	0.667	0.867	0.846	0.820
HON Node2vec(p=0.5)	0.797	0.755	0.730	0.884	0.868	0.885
HON Node2vec(p=2)	0.777	0.767	0.712	0.883	0.876	0.864
HONEM	0.495	0.614	0.551	0.684	0.782	0.711

Table A.6: F1-scores for the Workplace dataset with  $\Delta=30$  min. FON performs best. HONEM performs worst by large margin.

	f1-macro			f1-micro		
	LR	RF	SVC	LR	RF	SVC
FON DeepWalk	0.453	0.506	0.494	0.519	0.591	0.572
FON GraRep(S=4)	0.603	0.579	0.548	0.691	0.685	0.643
FON NetMF(W=10)	0.688	0.528	0.531	0.756	0.659	0.635
FON NetMF(W=2)	0.668	0.458	0.506	0.732	0.591	0.636
FON Node2vec(p=0.5)	0.452	0.494	0.492	0.516	0.573	0.569
FON Node2vec(p=2)	0.448	0.502	0.484	0.523	0.583	0.555
HON DeepWalk	0.588	0.546	0.530	0.656	0.644	0.657
HON GraRep(S=4)	0.602	0.607	0.471	0.665	0.723	0.560
HON NetMF(W=10)	0.566	0.476	0.481	0.573	0.576	0.583
HON NetMF(W=2)	0.534	0.424	0.460	0.543	0.531	0.567
HON Node2vec(p=0.5)	0.658	0.501	0.587	0.748	0.628	0.735
HON Node2vec(p=2)	0.542	0.500	0.514	0.639	0.617	0.628
HONEM	0.214	0.351	0.393	0.369	0.512	0.521

Table A.7: F1-scores for the HospitalWard dataset with  $\Delta=20$  sec. HON performs best. The scores are lower compared to the other datasets, indicating that classification does not work well on this dataset.

	f1-macro			f1-micro		
	LR	RF	SVC	LR	RF	SVC
FON DeepWalk	0.452	0.506	0.454	0.531	0.560	0.535
FON GraRep(S=4)	0.498	0.558	0.460	0.601	0.688	0.555
FON NetMF(W=10)	0.618	0.423	0.464	0.681	0.549	0.576
FON NetMF(W=2)	0.593	0.379	0.399	0.659	0.508	0.537
FON Node2vec(p=0.5)	0.454	0.417	0.447	0.524	0.463	0.516
FON Node2vec(p=2)	0.404	0.414	0.411	0.453	0.452	0.460
HON DeepWalk	0.526	0.452	0.439	0.540	0.504	0.537
HON GraRep(S=4)	0.486	0.543	0.414	0.577	0.665	0.508
HON NetMF(W=10)	0.449	0.400	0.347	0.489	0.545	0.469
HON NetMF(W=2)	0.403	0.358	0.302	0.431	0.453	0.427
HON Node2vec(p=0.5)	0.598	0.544	0.473	0.628	0.585	0.587
HON Node2vec(p=2)	0.497	0.469	0.340	0.509	0.500	0.465
HONEM	0.261	0.316	0.310	0.451	0.473	0.525

Table A.8: F1-scores for the HospitalWard dataset with  $\Delta=40$  sec. FON performs best. The scores are lower compared to the other datasets, indicating that classification does not work well on this dataset.

	f1-macro			f1-micro		
	LR	RF	SVC	LR	RF	SVC
FON DeepWalk	0.424	0.482	0.453	0.477	0.540	0.521
FON GraRep(S=4)	0.525	0.603	0.473	0.632	0.755	0.609
FON NetMF(W=10)	0.611	0.433	0.414	0.639	0.511	0.524
FON NetMF(W=2)	0.595	0.318	0.367	0.624	0.444	0.512
FON Node2vec(p=0.5)	0.454	0.479	0.469	0.516	0.524	0.545
FON Node2vec(p=2)	0.416	0.503	0.423	0.468	0.533	0.475
HON DeepWalk	0.490	0.411	0.443	0.572	0.507	0.552
HON GraRep(S=4)	0.569	0.542	0.415	0.657	0.708	0.563
HON NetMF(W=10)	0.459	0.295	0.272	0.456	0.487	0.424
HON NetMF(W=2)	0.472	0.303	0.264	0.472	0.487	0.412
HON Node2vec(p=0.5)	0.467	0.367	0.406	0.545	0.477	0.507
HON Node2vec(p=2)	0.495	0.492	0.442	0.587	0.564	0.544
HONEM	0.240	0.352	0.338	0.415	0.557	0.537

Table A.9: F1-scores for the HospitalWard dataset with  $\Delta=1$  min. FON performs best. The scores are lower compared to the other datasets, indicating that classification does not work well on this dataset.



## Appendix B

# Additional Experimental results

The following experiments are relegated to the appendix because the results are less relevant.

### B.1 Probability Prediction

According to the review in § 4, the second-most popular experiment is link prediction. The point about link prediction is assuming an incompletely observed network and guessing (unknown!) edges from the embedding. Experiments mimic this with cross-validation.

However, the existence of a link depends only on the network topology, a first-order concept. Moreover, ‘mixed’ interaction embeddings embed the start-tuple  $s_{1..k}$  of  $T_{s_{1..k} \rightarrow t}^{(k)}$ , which already reveals a few valid edges.

Hence, let us predict transition probabilities  $T_{s_{1..k} \rightarrow t}^{(k)}$  instead. Conveniently, we need an embedding with ‘mixed’ interactions for  $k > 1$ , which we had no opportunity to evaluate yet.

We predict the transition probabilities from the embedding by modifying formula (2.4) to consider the network topology:

$$\hat{p}(t|s_{1..k}) = \begin{cases} \frac{\exp(\langle \vec{s}_{1..k}, \vec{t} \rangle)}{\text{const}(s_{1..k})} & \text{if } (s_k, t) \in \mathcal{E} \\ 0 & \text{else} \end{cases}$$

This corresponds to the `decode`-method mentioned in § 3.2 with the parameters `use_neighborhood=True` and `normalize=True`.

#### B.1.1 Network

A network, whose probabilities an embedding manages to reproduce perfectly, is ill-suited because generalization is critical when dealing with incomplete data. The network’s structure should be resilient to removing or adding a few edges. A power grid or the internet looks promising due to the redundancy.

However, this is a topic for future research, and we simply use the PrimarySchool network from § 4.2.2.

### B.1.2 Baselines

The approaching deadline for submitting the thesis allows only for a small simulation study. We compare HON and FON variants of NetMF with the following parameters:

- **dimension**  $R \in \{16, 32, 64, 128\}$ ; smaller values increase the approximation error
- **negative**  $N \in \{1, 5\}$ ; larger values reduce the accuracy of the approximation
- **window\_size**  $W \in \{1, 2, 3, 5\}$ ; larger values blur the probabilities

### B.1.3 Best Case: Complete Information

Let us check first if probability prediction works at all, i.e., with complete information about the network.

NetMF(W=1, N=1, R=128) and NetMF(W=5, N=5, R=16)<sup>1</sup> represent two opposite poles: The former aims for good prediction, while the latter aims for stability.

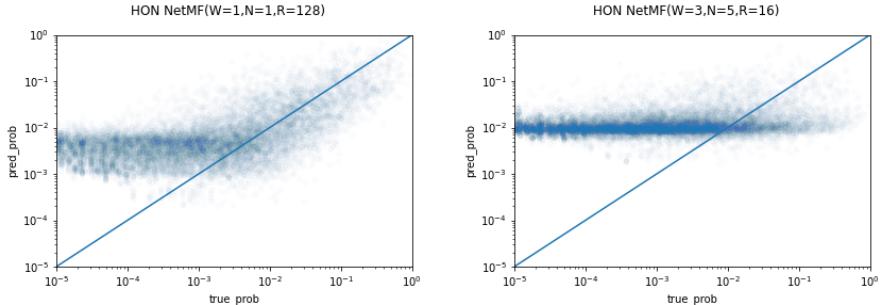


Figure B.1: As expected, HON NetMF(W=1, N=1, R=128) is better at predicting probabilities than NetMF(W=3, N=5, R=16)

### B.1.4 Methodology

The goal of cross-validation is to verify transition probabilities encoded by the embedding while making sure this information has not been used to train the embedding.

**Cross-validation for Weighted Optimization** Comparing matrix factorization via SVD and SGD, we notice that the latter is a weighted optimization problem while the former is an unweighted one. Setting the weight to zero is an easy way to encode missing information in a weighted optimization problem, which does not work for the unweighted one. Consequently, this approach works only for embeddings based on SGD.

Instead of splitting the data into training and test, we select a set of pair of nodes (FON edges)  $\mathcal{X}$  and make sure the `update(u, v, ...)` step in learning

---

<sup>1</sup>We use NetMF(W=3, N=5, R=16) in the figures instead

the embeddings is skipped if the pair  $(w, c) \in \mathcal{X}$ . If  $w$  or  $c$  are higher-order De Bruijn nodes, we use the pair  $(w[-1], w[-1])$  instead, which consists of the last FON nodes of  $w$  and  $c$  only. (This should also apply to the negative samples.) This effectively means replacing  $p(w, c)$  by the corresponding conditional distribution on  $(w, c) \notin \mathcal{X}$ .

Finally, compare those pairs excluded by  $\mathcal{X}$  the transition probabilities with their counterparts derived from the embedding and measure the approximation's quality by a quadratic loss. This approach has the disadvantage that we need control of the SGD optimization, which is complicated when using an existing implementation such as gensim's Word2vec. (Not to mention that SGD needs more fine-tuning than SVD.)

**Non-informative transition probabilities** Another idea is to tweak the probabilities, removing the information about the original value in order to use matrix factorization. Let  $s$  define  $\mathcal{X}_w = \{c | (w, c) \in X\}$  ('pair of nodes' case) or  $\mathcal{X}_w = \{c | (w[-1], c) \in X\}$  ('mixed' case).

Because higher-order transition probabilities already involve walks implies that the FON edges are already known. In order to hide any information about  $p(c|w)$  for  $c$  in  $\mathcal{X}_w$ , we replace it with the uniform distribution  $(1/d_w^{out}, \text{ where } d_w^{out} \text{ is the outdegree of } w)$ . This us our definition of *non-informative*.

$$\tilde{p}(c|w) = \begin{cases} 1/d_w^{out} & \text{if } c \in \mathcal{X}_w \\ p(c|w) \cdot \frac{1 - |\mathcal{X}_w|/d_w^{out}}{1 - \sum_{x \in \mathcal{X}_w} p(x|w)} & \text{else} \end{cases}$$

This is basically the idea for this experient: Use cross-validation to generate  $\mathcal{X}$ , create an embedding based on  $\tilde{p}(c|w)$ , and compare true and predicted probabilities.

### B.1.5 Results

Using the SSE criterion (sum of squared errors between true and predicted probability) — see the end of § C.11 — those estimators perform best, which perform worst in the best case (figure B.1). Unfortunately, this is not just a trade-off between good fit and the ability to generalize — it is just a shrinkage effect: a constant is a better approximation of the true probabilities than an uncorrelated prediction.

## B.2 HON Penalized Transition Embedding

We never displayed the experimental estimator (§ 3.2.5), because we did not find an applications for 'mixed' interactions, and figure 4.5 uses HON NetMF. Let us remedy this omission, see figure B.3.

Visualizing the loss during training (figure B.4), we observe a ledge, which is because the two parts (positive and negative) of the loss ( $\mathcal{L}(w, c)$ , see § 3.2.5) counteract each other:

$$\mathcal{L}(w, c) := - \left( \underbrace{p(w, c) \log \sigma(\langle \vec{w}, \vec{c} \rangle)}_{\text{positive}} + \underbrace{N \cdot p(w) \cdot p_N(c) \log \sigma(-\langle \vec{w}, \vec{c} \rangle)}_{\text{negative}} \right)$$

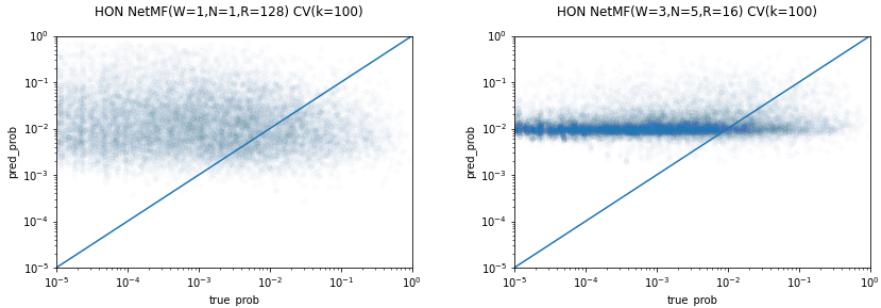


Figure B.2: Comparing true vs. predicted probabilities using cross-validation, using an non-informative distribution. Unfortunately, the probability prediction breaks down when using a non-informative distribution with 1% of the FON edges excluded.

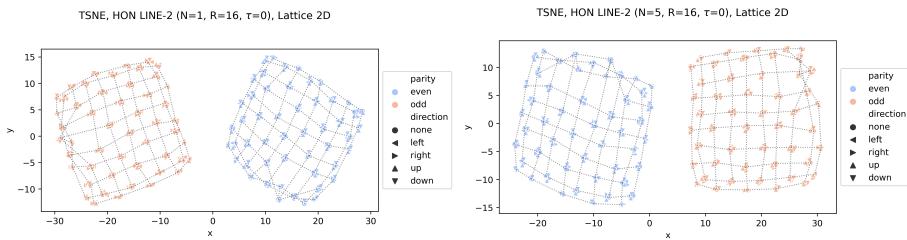


Figure B.3: Comparing the experimental estimator without penalty (i.e., HON LINE-2) for one (left) and five (right) negative samples  $N$ . The result is as expected when considering figures 4.5 and 2.2 (right).

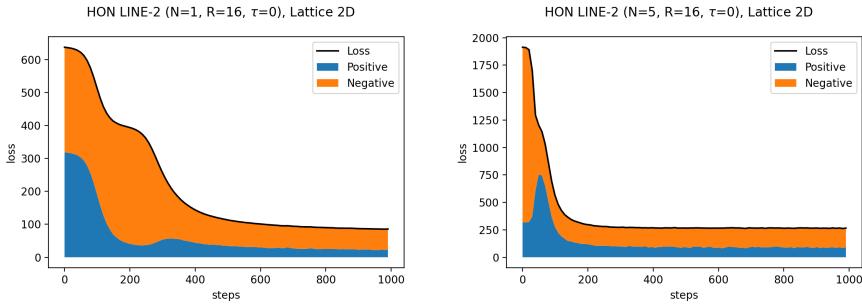


Figure B.4: Splitting the loss during training of the experimental estimator into positive and negative part explains the ledge in the sum. Note that for few negative samples ( $N = 1$ ) the positive part is optimized first, while for larger ones ( $N = 5$ ) it is the other way round.

# Appendix C

## Source Code

The codebase is divided between infrastructure and applications (e.g., code for an experiment or a specific figure). For the former, we provide information about its usage and display the source code in full length (§ C.2 to § C.6). For the latter, we display only the jupyter notebooks.

There is also a change in terminology: throughout the thesis, I avoided the term ‘path’ and used ‘walk’ instead because it is unclear whether a path must have distinct edges<sup>1</sup>. However, I noticed this late and did not refactor the code afterward.

### C.1 Introduction and Usage

**HigherOrderPathGenerator** First, we have to create an instance of the class `HigherOrderPathGenerator` and load the rules, configuration, and metadata; see `Plots_ExpClass.ipynb`. The rules are the transition probabilities are stored in a dictionary ( $\text{rules}[s_{1..k}] = \{t : T_{s_{1..k} \rightarrow t}\}$ ). The metadata contains information about the nodes, used for visualization and classification. The configuration identifies the rules (or transition probabilities) and will be exported together with the figure in `Visualizations.save_describe()`; see below.

```
1 from HigherOrderPathGenerator import HigherOrderPathGenerator
2 # read data/models/primaryschool_1.config into dictionary model_config
3 with open('data/models/primaryschool_1.config', 'r') as f:
4     model_config = { items[0]: 'u'.join(items[1:]) for items in [
5         line.strip('\n').split('\t') for line in f.readlines() ] }
6 gen = HigherOrderPathGenerator(node_sort_key=int, id='primaryschool_1',
7     config=model_config)
8 gen.load_BuildHON_rules('data/models/primaryschool_1.csv')
9 # read data/models/metadata_primaryschool.csv into dictionary metadata
10 with open('data/models/metadata_primaryschool.csv', 'r') as f:
11     metadata = {int(i): Ci for i,Ci [line.split() for line in f]}
12 gen.add_metadata('Class', metadata, use_last=True)
13 gen.add_metadata('Role', { i:'Teacher' if c=='Teachers' else 'Child' for i,c
14     in metadata.items()})
15 # checks
16 print(list(gen.check_transition_probs(tol=1e-14)))
17 print(list(gen.verify_stationarity(1, tol=1e-15)))
```

To avoid duplicate code, we may also use:

---

<sup>1</sup>[https://en.wikipedia.org/wiki/Path\\_\(graph\\_theory\)](https://en.wikipedia.org/wiki/Path_(graph_theory))

```

1 from Datasets import init_generator
2 gen = init_generator('primaryschool_1.csv')

```

Alternatively, a `HigherOrderPathGenerator` corresponding to the HON Lattice 2D (§ 3.3) is instantiated with:

```

1 from SyntheticNetworks import create_lattice_2nd_order_dynamic
2 gen = create_lattice_2nd_order_dynamic(size, omega, lattice_sep='-')
3 # type(gen) = HigherOrderPathGenerator

```

We can also convert a generator to a FON with `gen.to_FON()`, which discards all the higher-order transition probabilities.

We may also specify a `node.sort_key`, which defines the order of nodes and paths (which are first sorted by length and then by the individual nodes from start to end). This order is used when exporting the embedding as pandas DataFrame or decoding the probabilities as pandas Series.

**Embedding** We distinguish between *symmetric* and *asymmetric* embeddings. Asymmetric embeddings (e.g., NetMF and GraRep) calculate two embeddings `source_embedding` and `target_embedding`. Because they support both ‘pairs of nodes’ and ‘mixed’ interactions (see § 3.1.2), we assume that the source embeds paths (of length one in case of ‘pairs of nodes’), and the target embeds nodes. Because pandas works better with an index of strings rather than tuples, we convert the paths (`path2str`) and nodes (`node2str`) - both are configurable with `repr` as default. As this gets confusing, we will use an `EmbeddingView`, which unifies both variants, and refers to both single nodes and paths as `key`.

Symmetric embeddings (e.g., DeepWalk and Node2vec) are restricted to ‘pairs of nodes’ interactions and calculate a single `embedding`, and there is no need to distinguish between `key` and `node`.

To calculate an embedding, instantiate an embedding class and call `train()`.

```

1 from Embedding import HON_NetMF_EMBEDDING
2 emb = HON_NetMF_EMBEDDING(gen, dimension=128, pairwise=True)
3 %time emb.train(window_size=10, negative=1, optimized=True)

```

The `Generic.SkipGram_EMBEDDING` imports the output of, e.g., LINE and Node2vec instead of calculating an embedding; see `data/LINE/LINE.ipynb`.

### EmbeddingView and Visualization

```

1 from Visualizations import Visualization, EmbeddingView
2 ev = EmbeddingView(emb, use_source=True)
3 vis = ev.visualize_TSNE(random_state=1, n_iter=1000)
4 vis.plot1(hue='Class', figsize=(5,4), dpi=200)
5 #vis.annotate_node('xyz', 'XYZ', 'center', 'center')
6 vis.save_describe('tmp/plot.png', comment='...')

```

In the case of a HON Lattice 2D, we use `Lattice2D_EMBEDDING` instead of `EmbeddingView`, as the former provides additional functionality related to the grid structure.

The `EmbeddingView` is currently initialized after the embedding, but I better had the embedding provide these instances as properties `source` and `target` and tidied up the public interface. The `HigherOrderPathGenerator` would act as a factory and decide whether it represents a lattice.

## C.2 HigherOrderPathGenerator.py

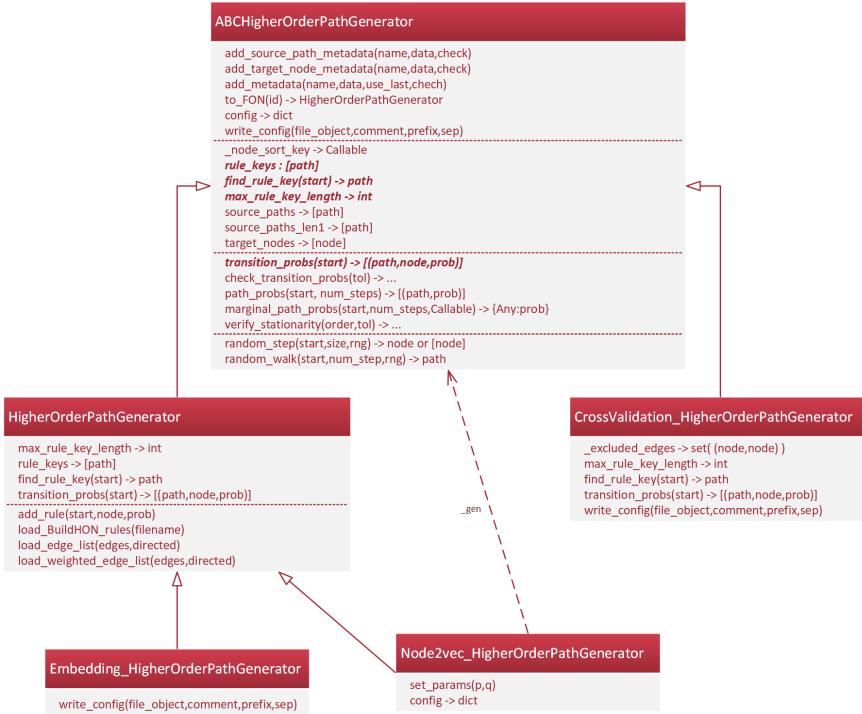


Figure C.1: UML class diagram for HigherOrderPathGenerator.py

These classes store the transition probabilities  $T_{s_{1..k} \rightarrow t}^{(k)}$  and enumerate or simulate random walks.

```

1  from abc import ABC, abstractmethod
2  from typing import Any, Callable, Dict, List, Iterator, Tuple, Optional #
3      https://mypy.readthedocs.io/en/stable/cheat_sheet_py3.html
4  from collections import defaultdict, Counter
5  import math
6  import numpy as np
7  import pandas as pd
8  #import random
8  import bisect
9
10 class ABCHigherOrderPathGenerator(ABC):
11     "Generator for higher-order paths (abstract) - storage of rules is not implemented here"
12     def __init__(self, node_sort_key=lambda node: node, id: Optional[str] = None,
13                  config: Dict[str, Any] = dict()):
14         self._rules_cumulated = None
15         self._node_sort_key = node_sort_key
16         self._id = type(self).__name__ if id is None else id
17         self._source_path_metadata = dict() # Dict[str, Dict[Tuple[Any, ...], Any]]
18         self._target_node_metadata = dict() # Dict[str, Dict[Any, Any]]
19         self._config = config
20
21     # metadata will be used in the visualizations
22     def add_metadata(self, name, data: Dict[Any, Any], use_last: bool = True,
23                      check: bool = True):
24         """
25             Add metadata describing nodes for both target_nodes and source_paths.
26             For the latter use either the first or last node of the path.
27         """
  
```

```

25      """
26      self.add_target_node_metadata(name, data, check)
27      index = -1 if use_last else 0
28      paths_data = { key: data.get(key[index], None) for key in self.
29                     source_paths }
30      self.add_source_path_metadata(name, paths_data, check=False)
31
32  def add_source_path_metadata(self, name: str, data: Dict[Tuple[Any,...],Any
33 ], check: bool = True):
34     """Add metadata describing the source paths"""
35     if check:
36         assert len(set(data.keys()).difference(self.source_paths))==0, '
37             metadata "%s" contains invalid keys' % name
38         self._source_path_metadata[name] = data
39
40  def add_target_node_metadata(self, name: str, data: Dict[Any,Any], check:
41     bool = True):
42     """Add metadata describing the target nodes"""
43     if check:
44         assert len(set(data.keys()).difference(self.target_nodes))==0, '
45             metadata "%s" contains invalid keys' % name
46         self._target_node_metadata[name] = data
47
48 @property
49 @abstractmethod
50 def rule_keys(self) -> Iterator[Tuple[Any,...]]:
51     pass
52
53 @abstractmethod
54 def max_rule_key_length(self) -> int:
55     pass
56
57 @property
58 def source_paths(self) -> List[Tuple[Any,...]]:
59     sort_key = lambda path:(len(path), *map(self._node_sort_key, path))
60     "returns a sorted list of source paths of the transition probabilities(
61         empty key is removed)"
62     return sorted({ key for key in self.rule_keys if len(key) >= 1 }, key=
63                 sort_key)
64
65 @property
66 def source_paths_leni(self) -> List[Tuple[Any,...]]:
67     "returns a sorted list of source paths of length 1 of the transition
68         probabilities"
69     sort_key = lambda path:self._node_sort_key(path[0])
70     #return sorted({ (v,) for key in self.rule_keys for v in key }, key=
71     #    sort_key) # safe & slower
72     return sorted({ key for key in self.rule_keys if len(key) == 1 }, key=
73                 sort_key)
74
75 @property
76 def target_nodes(self) -> List[Tuple[Any,...]]:
77     "returns a sorted list of target nodes of the transition probabilities"
78     #return sorted({ next_node for key in self.rule_keys for _, next_node,
79     #    in self.transition_probs(key) }) # safe & slower
80     return sorted({ key[0] for key in self.rule_keys if len(key) == 1 }, key=
81                 self._node_sort_key) # this should also work, as otherwise random
82     walks are broken.
83
84 @abstractmethod
85 def transition_probs(self, start: Tuple[Any,...]) -> Iterator[Tuple[Tuple[
86     Any,...],Any,float]]:
87     pass
88
89 def check_transition_probs(self, tol=1e-14) -> Iterator[Tuple[Tuple[Any
90     ,...],float]]:
91     "Verify that the transition probabilities sum up to one (with some
92         numerical tolerance)."
93     for key in self.rule_keys:

```

```

83     total_prob = sum(prob for _, prob in self.transition_probs(key))
84     if abs(total_prob - 1) > tol:
85         yield (key, total_prob)
86
87     def path_probs(self, start: Tuple[Any, ...], num_steps: int = 1) -> Iterator
88         [Tuple[Tuple[Any, ...], float]]:
89             "Conditional probabilities of all paths starting with some given sequence
90             of nodes."
91             if num_steps <= 0:
92                 yield (start, 1.0)
93             else: # elif num_steps >= 1:
94                 for outer_path, outer_prob in self.path_probs(start, num_steps - 1):
95                     for _, next_node, prob in self.transition_probs(outer_path):
96                         yield ((*outer_path, next_node), outer_prob * prob)
97
98     def marginal_path_probs(self, start: Tuple[Any, ...], num_steps: int = 1,
99                             projection: Callable[[Tuple[Any, ...]], Any] = lambda x: x[-1]) -> Dict[
100                                Any, float]:
101         "Calculates the marginal probabilities of path_probs. May be used to
102             calculate k-step transition probabilities."
103         probs = defaultdict(float)
104         for path, prob in self.path_probs(start, num_steps):
105             key = projection(path)
106             probs[key] += prob
107
108     return probs
109
110     def verify_stationarity(self, order=2, tol=1e-14) -> Iterator[Tuple[Any,
111                               float, float, float]]:
112         """
113             Verify that the stationary distribution is indeed stationary with respect
114             to the first order transition probabilities.
115
116             The (1st order) stationary distribution is stored as the transition
117             probabilities of an empty start path.
118         """
119         mprobs = self.marginal_path_probs(start=(), num_steps=order + 1, projection
120                                         =lambda x: x[-1]) # node-distribution after a few steps
121         for _, node, prob in self.transition_probs(start=()): # stationary
122             distribution
123                 mprob = mprobs[node]
124                 if abs(prob - mprob) > tol:
125                     yield (node, prob, mprob, prob - mprob)
126
127     @property
128     def rules_cumulated(self):
129         "Cumulates the transition probabilities for random.choices"
130         if self._rules_cumulated is None:
131             rules2 = dict()
132             for start in self.rule_keys:
133                 next_nodes = list()
134                 cum_weights = list()
135                 cum_prob = 0.0
136                 for _, next_node, prob in self.transition_probs(start):
137                     cum_prob += prob
138                     next_nodes.append(next_node)
139                     cum_weights.append(cum_prob)
140                 rules2[start] = (next_nodes, cum_weights)
141             self._rules_cumulated = rules2
142
143     return self._rules_cumulated
144
145     def random_step(self, start: Tuple[Any, ...], size=None, rng: Optional[np.
146                                     random.Generator] = None) -> Any:
147         if rng is None:
148             rng = np.random.default_rng()
149             rule_key = self.find_rule_key(start)
150             if rule_key == () and start != (): # or assert start is a tuple
151                 print('could not find a key matching %s' % start)
152             next_nodes, cum_weights = self.rules_cumulated[rule_key]
153             if size is None: # return scalar
154                 x = rng.random(1)
155                 # return random.choices(next_nodes, cum_weights=cum_weights, k=1)[0] #
156                 # fine, but wrong RNG
157             return next_nodes[bisect.bisect_left(cum_weights, x[0], hi=len(
158                 cum_weights) - 1)]

```

```

144     else: # return list
145         x = rng.random(size)
146         f = lambda y:next_nodes[bisect.bisect_left(cum_weights, y, hi=len(
147             cum_weights)-1)]
148         return list(map(f, x))
149
150     def random_walk(self, start: Tuple[Any,...], num_steps:int = 1, rng:
151         Optional[np.random.Generator] = None) -> Tuple[Any,...]:
152         if rng is None:
153             rng = np.random.default_rng()
154         walk = start
155         for _ in range(num_steps):
156             next_node = self.random_step(walk, size=None, rng=rng)
157             walk = (*walk,next_node)
158         return walk
159
160     def to_FON(self, id: Optional[str] =None) -> "HigherOrderPathGenerator":
161         "Return a HigherOrderPathGenerator containing only first-order rules (and
162             the stationary distribution)"
163         if id is None:
164             id = self._id + '_FON'
165         fon = HigherOrderPathGenerator(node_sort_key=self._node_sort_key, id=id,
166             config=self._config)
167         if hasattr(self, 'creator'): # see Lattice2D_2nd_order_dynamic & required
168             for Lattice2D_EMBEDDINGVIEW
169                 fon.creator = self.creator
170
171         for key in self.rule_keys:
172             if len(key)>1:
173                 continue
174             for start, next_node, prob in self.transition_probs(key):
175                 fon.add_rule(start, next_node, prob)
176         for name, data in self._source_path_metadata.items():
177             fon.add_source_path_metadata(name, {key: value for key,value in data.
178                 items() if len(key)==1 }, check=True)
179         fon._target_node_metadata = dict(self._target_node_metadata)
180
181     @property
182     def config(self):
183         "Get configuration"
184         cfg=dict(init_class=self.__class__.__name__, init_id=self._id)
185         cfg.update(self._config)
186         return cfg
187
188     def write_config(self, file_object, comment: str = '', prefix: str = '',
189                     sep: str = '\t'):
190         if comment != '':
191             file_object.write(prefix + comment + '\n' + prefix + '\n')
192             file_object.write(prefix + 'HigherOrderPathGenerator:\n')
193             for k,v in self.config.items():
194                 file_object.write(prefix + '%s%s\n' % (k,sep,v))
195
196     class HigherOrderPathGenerator(ABCHigherOrderPathGenerator):
197         "Generator for higher-order paths"
198         def __init__(self, node_sort_key=lambda node: node, id: Optional[str] =
199             None, config: Dict[str,Any] = dict()):
200             super().__init__(node_sort_key, id, config)
201             self.rules = defaultdict(dict)
202             self._max_rule_key_length = 1 # find_rule_key relies on
203                                         _max_rule_key_length > 0
204
205         @property
206         def max_rule_key_length(self) -> int:
207             return self._max_rule_key_length
208
209         def clear_rules(self):
210             "Clear the rules dictionary (incl. cleanup)"
211             self._max_rule_key_length = 1
212             self.rules.clear()
213             self._rules_cumulated = None
214
215         def add_rule(self, start: Tuple[Any,...], next_node: Any, prob: float):
216             "Add entries to rules. (load... methods must not manipulate the rules
217               directly.)"

```

```

208     assert type(start)==tuple, 'argument "start" expected a tuple, but got %r
209     % start
210     self._max_rule_key_length = max(len(start), self._max_rule_length)
211     self.rules[start][next_node] = prob
212
213     def load_BuildHON_rules(self, filename: str):
214         "Load output of BuildHON"
215         with open(filename, 'r') as f:
216             for line in f:
217                 items = line.split(' ')
218                 assert (len(items)>=3) and (items[-3]==='=>'), 'invalid line "%s" %
219                 line
220                 key = tuple(map(int,items[:-3]))
221                 self.add_rule(key, int(items[-2]), float(items[-1]))
222                 print('%d rules read' % sum([len(v) for v in self.rules.values()]))
223
224     def load_edge_list(self, edges : Iterator[Tuple[Any,Any]], directed=True):
225         def add_weights(edges_: Iterator[Tuple[Any,Any]]) -> Iterator[Tuple[Any,
226             Any,Any]]:
227             for u,v in edges_:
228                 yield u,v,1 # add weight
229             self.load_weighted_edge_list(add_weights(edges_), directed)
230
231     def load_weighted_edge_list(self, edges : Iterator[Tuple[Any,Any,Any]],
232         directed=True):
233         if not directed:
234             def to_directed(edges_):
235                 for u,v,w in edges_:
236                     yield u,v,w
237                     yield v,u,w
238             edges = to_directed(edges)
239             edges_list = list(edges) # iterating twice over edges does not work
240             out_weights = defaultdict(float)
241             for u,,,w in edges_list:
242                 out_weights[u] += w
243             for u,v,w in edges_list:
244                 self.add_rule((u,), v, w/out_weights[u])
245
246         @property
247         def rule_keys(self) -> Iterator[Tuple[Any,...]]:
248             return self.rules.keys()
249
250         def find_rule_key(self, start: Tuple[Any,...]) -> Tuple[Any,...]:
251             if len(start) > self._max_rule_key_length: # and (self.
252                 _max_rule_key_length > 0):
253                 start = start[-self._max_rule_key_length:]
254             while len(start)>0:
255                 if start in self.rules:
256                     return start
257                 else:
258                     start = start[1:]
259             return tuple()
260
261         def transition_probs(self, start: Tuple[Any,...]) -> Iterator[Tuple[ Tuple[
262             Any,...],Any,float]]:
263             rule_key = self.find_rule_key(start)
264             probs = self.rules[rule_key]
265             for next_node,prob in probs.items():
266                 yield(rule_key, next_node, prob)
267
268         class Node2vec_HigherOrderPathGenerator(HigherOrderPathGenerator):
269             """
270                 Adjusts the transition probabilities for the search bias of Node2vec.
271                 Note, that Node2vec corresponds to DeepWalk with biased random walks.
272             """
273             def __init__(self, gen: ABCHigherOrderPathGenerator, p: float = 1, q: float
274                         = 1, id_format: Optional[str] = None):
275                 # example: id_format = 'Node2vec(p={0}, q={1})'
276                 super().__init__(gen._node_sort_key, id=gen._id, config=gen._config)
277                 self._gen = gen
278                 self._source_path_metadata = dict(gen._source_path_metadata)
279                 self._target_node_metadata = dict(gen._target_node_metadata)
280                 if hasattr(gen, 'creator'): # see Lattice2D_2nd_order_dynamic & required
281                     for Lattice2D_EMBEDDINGVIEW

```

```

274     self.creator = gen.creator
275     self._id_format = id_format # if id_format is provided, _id will be
276         overwritten
277     self.set_params(p, q)
278
279     def set_params(self, p: float = 1, q: float = 1):
280         "Sets the parameters for the search bias"
281         self._p = p
282         self._q = q
283         if self._id_format is not None:
284             self._id = self._id_format.format(p,q)
285         self.clear_rules()
286         rules_tmp = defaultdict(list)
287         # copy all rules
288         for key in self._gen.rule_keys:
289             rules_tmp[key] = list( (next_node, prob) for _, next_node, prob in self
290                 ._gen.transition_probs(key) )
291         # generate rules for all direct neighbors if they do not yet exist
292         direct_neighbors = set()
293         for key1 in self._gen.source_paths_len1: # key1 is a tuple containing one
294             node
295             for _, key2, _ in self._gen.transition_probs(key1): # key2 is a node
296                 key = (key1[0], key2)
297                 direct_neighbors.add(key)
298                 if key not in rules_tmp:
299                     rules_tmp[key] = list( (next_node, prob) for _, next_node, prob in
300                         self._gen.transition_probs((key2,)) )
301         # add modified rules
302         for key, probs in rules_tmp.items():
303             if len(key) < 2:
304                 for next_node, prob in probs:
305                     self.add_rule(key, next_node, prob)
306             else:
307                 prev_node = key[-2]
308                 sum_probs = 0
309                 new_probs = list()
310                 for next_node, prob in probs:
311                     # search bias
312                     if next_node == prev_node:
313                         prob /= p
314                     elif (prev_node, next_node) not in direct_neighbors:
315                         prob /= q
316                     new_probs.append((next_node, prob))
317                     sum_probs += prob
318                 for next_node, prob in new_probs:
319                     self.add_rule(key, next_node, prob / sum_probs)
320
321     @property
322     def config(self):
323         "get configuration"
324         cfg=dict(init_class=self.__class__.__name__, init_id=self._id, p=self._p,
325                  q=self._q)
326         cfg.update(self._config)
327         return cfg
328
329     class CrossValidation_HigherOrderPathGenerator(ABCHigherOrderPathGenerator):
330         """
331             Adjusts the transition probabilities for cross validation
332         """
333         def __init__(self, gen: ABCHigherOrderPathGenerator, excluded_edges :
334             Iterator[Tuple[Any,Any]], id: Optional[str] = None):
335             self._gen = gen
336             self._excluded_edges = set(excluded_edges)
337             super().__init__(gen._node_sort_key, id, gen._config)
338             self._source_path_metadata = dict(gen._source_path_metadata)
339             self._target_node_metadata = dict(gen._target_node_metadata)
340             # init rules
341             self.rules = dict()
342             for start in gen.rule_keys:
343                 if len(start)==0:
344                     # todo: the stationary distribution must also be adjusted
345                     self.rules[start]= { next_node: prob for _,next_node,prob in gen.
346                         transition_probs(start)}
347                 continue

```

```

341     probs = list((next_node,prob) for _,next_node,prob in gen.
342                   transition_probs(start))
343     d = len(probs) # out degree
344     excluded = { next_node for next_node, _ in probs if (start[-1],next_node)
345                   ) in self._excluded_edges }
346     excluded_prob = sum(prob for next_node,prob in probs if next_node in
347                           excluded)
348     self.rules[start] = { next_node: 1/d if next_node in excluded else prob
349                           * (1-len(excluded)/d) / (1-excluded_prob)
350                           for next_node,prob in probs}
351
352     @property
353     def rule_keys(self) -> Iterator[Tuple[Any,...]]:
354         return self._gen.rule_keys
355
356     def find_rule_key(self, start: Tuple[Any,...]) -> Tuple[Any,...]:
357         return self._gen.find_rule_key(start)
358
359     @property
360     def max_rule_key_length(self) -> int:
361         return self._gen.max_rule_key_length
362
363     def transition_probs(self, start: Tuple[Any,...]) -> Iterator[Tuple[Tuple[
364         Any,...],Any,float]]:
365         rule_key = self._gen.find_rule_key(start)
366         probs = self.rules[rule_key]
367         for next_node,prob in probs.items():
368             yield(rule_key, next_node, prob)
369
370     def write_config(self, file_object, comment: str = '', prefix: str = '',
371                     sep: str = '\t'):
372         self._gen.write_config(file_object=file_object, comment=comment, prefix=
373                               prefix, sep=sep)
374         file_object.write(prefix + '\n' + prefix + 'CrossValidation-
375                           HigherOrderPathGenerator:\n')
376         for k,v in self.config.items():
377             file_object.write(prefix + '%s%s%s\n' % (k,sep,v))
378
379     class Embedding_HigherOrderPathGenerator(HigherOrderPathGenerator):
380         "Generator based on the decode-method of an asymmetric embedding"
381         def __init__(self, emb: 'ABCAsymmetricEmbedding', use_neighborhood: bool =
382                      True, no_self_loops: bool = False, id: Optional[str] = None, **kwargs):
383             :
384             config = dict(init_use_neighborhood=use_neighborhood, init_no_self_loops=
385                           no_self_loops, init_emb=emb._id, **kwargs)
386             super().__init__(emb._gen._node_sort_key, id=id, config=config)
387             gen = emb._gen
388             self._emb = emb
389             self._source_path_metadata = dict(gen._source_path_metadata)
390             self._target_node_metadata = dict(gen._target_node_metadata)
391             if hasattr(gen, 'creator'): # see Lattice2D_2nd_order_dynamic & required
392                 for Lattice2D_EMBEDDINGVIEW
393                     self.creator = gen.creator
394             str2nodes = { emb.node2str(node):node for node in emb.target_nodes }
395             for start in emb.source_paths:
396                 for next_node_str,prob in emb.decode_path(start=start, use_neighborhood
397                     =use_neighborhood, no_self_loops=no_self_loops, normalize=True, **
398                     kwargs).items():
399                     if prob>0:
400                         self.add_rule(start, str2nodes[next_node_str], prob)
401             # todo: calculate stationary distribution from the first-order rules
402             for _,next_node,prob in gen.transition_probs(start=()): # copy stationary
403                 distribution from original generator
404                 self.add_rule((), next_node, prob)
405
406             def write_config(self, file_object, comment: str = '', prefix: str = '',
407                             sep: str = '\t'):
408                 self._emb.write_config(file_object=file_object, comment=comment, prefix=
409                               prefix, sep=sep)
410                 file_object.write(prefix + '\n' + prefix + 'Embedding-
411                               HigherOrderPathGenerator:\n')
412                 for k,v in self.config.items():
413                     file_object.write(prefix + '%s%s%s\n' % (k,sep,v))

```

### C.3 SyntheticNetworks.py

This creates a `HigherOrderPathGenerator` containing a HON Lattice 2D (§ 3.3).

```

1  from HigherOrderPathGenerator import HigherOrderPathGenerator
2
3  from collections import Counter
4  import numpy as np
5  import pandas as pd
6  import math
7
8  def create_lattice_2nd_order_dynamic(size: int=10, omega: float=0,
9      lattice_sep: str='-', check: bool=False):
10     creator = Lattice2D_2nd_order_dynamic(size, lattice_sep)
11     return creator.create_generator(omega, check)
12
13 class Lattice2D_2nd_order_dynamic(object):
14     def __init__(self, size: int=10, lattice_sep: str='-'):
15         self.size=size
16         self.lattice_sep = lattice_sep
17         # '0-0' corresponds to the lower-left corner of the lattice. Format is (x
18         ,y).
18         self.neighbor_funcs = { (-1,0): self.left, (1,0): self.right, (0,1): self
19         .up, (0,-1): self.down }
19         self.coord2nodes = { (x,y): '%i%s%i' % (x, lattice_sep, y) for x in range
20         (size) for y in range(size) }
20         self.node2coords = { n:c for c,n in self.coord2nodes.items() }
21
22     def left(self, node:str)->str:
23         x,y = node.split(self.lattice_sep)
24         x = str(max(0, int(x)-1))
25         return x + self.lattice_sep + y
26     def right(self, node:str)->str:
27         x,y = node.split(self.lattice_sep)
28         x = str(min(self.size-1, int(x)+1))
29         return x + self.lattice_sep + y
30     def up(self, node:str)->str:
31         x,y = node.split(self.lattice_sep)
32         y = str(min(self.size-1, int(y)+1))
33         return x + self.lattice_sep + y
34     def down(self, node:str)->str:
35         x,y = node.split(self.lattice_sep)
36         y = str(max(0, int(y)-1))
37         return x + self.lattice_sep + y
38
39     @property
40     def horizontal_edges1(self):
41         return list( (u,v) for (u,v) in [(u,self.right(u)) for u in self.
42             coord2nodes.values()] if u!=v )
43
44     @property
45     def vertical_edges1(self):
46         return list( (u,v) for (u,v) in [(u,self.up(u)) for u in self.coord2nodes
47             .values()] if u!=v )
48
49     @property
50     def horizontal_edges2(self):
51         return list( ((u,v),(v,w)) for (u,v,w) in [(self.left(u),u,self.right(u))
52             for u in self.coord2nodes.values()] if u!=v and v!=w )
53
54     @property
55     def vertical_edges2(self):
56         return list( ((u,v),(v,w)) for (u,v,w) in [(self.down(u),u,self.up(u))
57             for u in self.coord2nodes.values()] if u!=v and v!=w )
58
59     def create_generator(self, omega: float=0, check: bool=False):
60         assert abs(omega)<=1
61         config = dict(code=__file__, init=self.__class__.__name__ + '.'
62             'create_generator',
63             size=self.size, lattice_sep=repr(self.lattice_sep), omega=omega)
64         # create list of nodes and (undirected) edges
65         # see also pathpy.generators.lattice_network(start=0, stop=size, dims=2)
66         nodes = list(self.coord2nodes.values())
67         undirected_edges = self.horizontal_edges1 + self.vertical_edges1

```

```

62     edges = undirected_edges + [(v,u) for u,v in undirected_edges]
63     degrees = Counter(u for u,v in edges)
64     sum_degrees = sum(degrees.values())
65
66     #node_sort_key = lambda x: tuple(map(int,x.split('-')))
67     node_sort_key = self.node2coords.__getitem__
68     latgen = HigherOrderPathGenerator(node_sort_key, 'Lattice2D(%d,%fomega=%f)
69         % (self.size, omega), config)
70     latgen.creator = self
71     latgen.load_edge_list(edges, directed=True)
72     # add 2nd order rules for (u,v) -> w
73     for u in nodes:
74         latgen.add_rule((), u, degrees[u]/sum_degrees) # prob of stationary
75             distribution
76         for (d1,f1) in self.neighbor_funcs.items():
77             v = f1(u)
78             if u == v: # cannot move / no self-loops
79                 continue
80             for (d2,f2) in self.neighbor_funcs.items():
81                 w = f2(v)
82                 if v == w: # cannot move / no self-loops
83                     continue
84                 if v == f1(v): # if we cannot move twice in direction f1 from u, do
85                     not change dynamic
86                 d = 0.0
87             else:
88                 d = d1[0]*d2[0] - d1[1]*d2[1] # horizontal speed-up and vertical
89                     slow-down
90             latgen.add_rule((u, v), w, (1 + d*omega)/degrees[v])
91
92     latgen.add_source_path_metadata('key_len', {n: len(n) for n in latgen.
93         source_paths}, check=False)
94     latgen.add_metadata('x_orig', {n: c[0] for c,n in self.coord2nodes.items
95         ()}, use_last=True, check=False)
96     latgen.add_metadata('y_orig', {n: c[1] for c,n in self.coord2nodes.items
97         ()}, use_last=True, check=False)
98     latgen.add_metadata('parity', {n: 'even' if sum(c)%2==0 else 'odd' for c
99         ,n in self.coord2nodes.items()}, use_last=True, check=False)
100    direction = dict()
101    for key in latgen.source_paths:
102        if len(key) == 1:
103            direction[key]='none'
104        else:
105            c0 = self.node2coords[key[0]]
106            c1 = self.node2coords[key[-1]]
107            delta = (c1[0]-c0[0], c1[1]-c0[1])
108            direction[key]=self.neighbor_funcs[delta].__name__
109    latgen.add_source_path_metadata('direction', direction, check=False)
110    if check:
111        # check whether probabilities sum up to one
112        print(list(latgen.check_transition_probs(tol=1e-14)))
113        # verify that the 2nd order rules did not modify the stationary
114            distribution
115        print(list(latgen.verify_stationarity(order=2, tol=1e-17)))
116    return latgen

```

## C.4 Datasets.py

The SocioPattern networks (§ 4.2.2) are loaded with `init_generator`.

```

1  from HigherOrderPathGenerator import HigherOrderPathGenerator
2  import os
3
4  def init_generator(filename, tol1=1e-14, tol2=1e-15, verbose=False):
5      filepath = 'data\\models\\' + filename
6      with open(os.path.splitext(filepath)[0] + '.config', 'r') as f:
7          model_config = {items[0]: ''.join(items[1:]) for items in [line.strip('\\n').split('\\t') for line in f.readlines()]}
8      if verbose:
9          print('Configuration:')
10         for k, v in model_config.items():
11             print('%s: %s' % (k, v))
12     model_config['filename'] = filename
13
14     gen = HigherOrderPathGenerator(node_sort_key=int, id=os.path.splitext(
15         filename)[0], config=model_config)
16     gen.load_BuildHON_rules('data\\models\\' + filename)
17     # add metadata
18     metadata_filename = 'metadata_' + (filename.split('.')[0].split('_')[0]) + '.csv'
19     metadata = dict()
20     with open('data\\models\\' + metadata_filename, 'r') as f:
21         for line in f:
22             i, Ci = line.split()
23             metadata[int(i)] = Ci
24     if metadata_filename == 'metadata_workplace.csv':
25         gen.add_metadata('Department', metadata, use_last=True)
26     elif metadata_filename == 'metadata_primaryschool.csv':
27         gen.add_metadata('Class', metadata, use_last=True)
28         gen.add_metadata('Role', {i: 'Teacher' if c == 'Teachers' else 'Child' for
29             i, c in metadata.items()})
30     elif metadata_filename == 'metadata_hospital.csv':
31         gen.add_metadata('Status', metadata, use_last=True)
32     elif metadata_filename == 'metadata_temporal-clusters.csv' or
33         metadata_filename == 'metadata_shuffled-temporal-clusters.csv':
34         gen.add_metadata('Color', metadata, use_last=True)
35     else:
36         raise Exception('unknown metadata_filename')
37     print(list(gen.check_transition_probs(tol1)))
38     print(list(gen.verify_stationarity(1, tol2)))
39     #print(list(gen.verify_stationarity(2, tol2))) # this does not hold
40
41     return gen

```

## C.5 Embedding.py

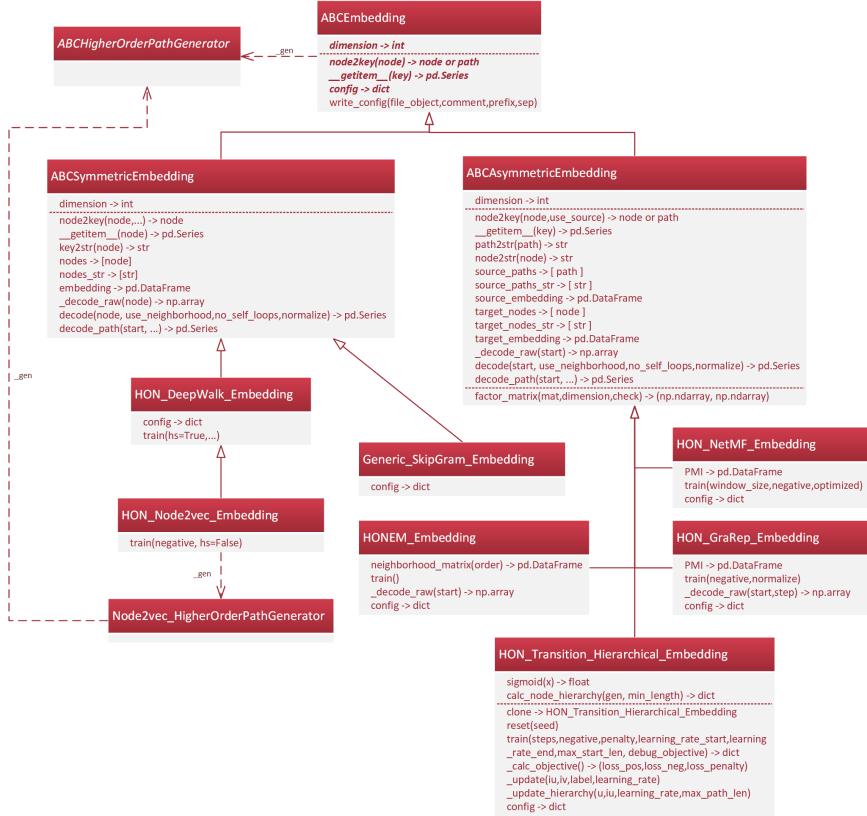


Figure C.2: UML class diagram for Embedding.py

All the embeddings mentioned in § 3.2 are implemented here.

```

1 from abc import ABC, abstractmethod
2 from typing import Any, Callable, Dict, Iterator, List, Tuple, Optional,
3     Union
4 import warnings
5 import numpy as np
6 import pandas as pd
7 import math
8 import collections
9 import os
10
11 from HigherOrderPathGenerator import ABCHigherOrderPathGenerator,
12     HigherOrderPathGenerator, Node2vec_HigherOrderPathGenerator
13 from gensim.models import Word2Vec
14 from scipy.linalg import svd
15 from sklearn.decomposition import TruncatedSVD
16 from zlib import adler32, crc32
17 import sys
18
19 # While the public interface of the embedding classes presents the
20 # embedding matrices as pandas DataFrames, they are internally stored
21 # in numpy ndarrays together with dictionaries for the row and column
22 # names. Besides better performance, the representation of paths as
23 # tuples did not fit well with pandas - where the paths had to be converted
# into strings. (And converting paths into str or int is basically the
# same in terms of readability of the code.)

```

```

24
25     class ABCEmbedding(ABC):
26         def __init__(self, gen: ABCHigherOrderPathGenerator, dimension: int,
27                      symmetric: bool=True):
28             self._gen = gen
29             self._symmetric = symmetric
30             self._dimension = dimension
31             self._id = type(self).__name__
32
33         @property
34         @abstractmethod
35         def dimension(self):
36             "effective dimension of embedding"
37             pass
38
39         @abstractmethod
40         def node2key(self, node: Any, use_source: bool=False):
41             pass
42
43         @abstractmethod
44         def __getitem__(self, key) -> pd.Series:
45             "Returns a row of the embedding matrix."
46             pass
47
48         @abstractmethod
49         def decode_path(self, start: Tuple[Any, ...], **kwargs) -> pd.Series:
50             pass
51
52         @property
53         @abstractmethod
54         def config(self):
55             "get configuration"
56             pass
57
58         def write_config(self, file_object, comment: str = '', prefix: str = '',
59                         sep: str = '\t'):
60             self._gen.write_config(file_object=file_object, comment=comment, prefix=
61                                   prefix, sep=sep)
62             file_object.write(prefix + '\n' + prefix + 'Embedding:\n')
63             for k, v in self.config.items():
64                 file_object.write(prefix + '%s%s\n' % (k, sep, v))
65
66     class ABCSymmetricEmbedding(ABCEmbedding):
67         def __init__(self, gen: ABCHigherOrderPathGenerator, dimension: int, nodes:
68                      Iterator[Any]):
69             super().__init__(gen, dimension, symmetric=True)
70             self._nodes = {n: i for i, n in enumerate(nodes)}
71             self._embedding = np.zeros(shape=(len(self._nodes), dimension))
72             self.key2str = repr
73
74         def __getitem__(self, key) -> pd.Series:
75             idx = self._nodes[key]
76             return pd.Series(self._embedding[idx, :], name=self.key2str(key))
77
78         def node2key(self, node: Any, use_source: bool=False):
79             return node
80
81         @property
82         def nodes(self) -> Iterator[Any]:
83             return self._nodes.keys()
84
85         @property
86         def nodes_str(self) -> List[str]:
87             return list(map(self.key2str, self.nodes))
88
89         @property
90         def embedding(self) -> pd.DataFrame:
91             return pd.DataFrame(self._embedding, index=self.nodes_str)
92
93

```

```

94     def _decode_raw(self, node: Any, step = None) -> np.array:
95         """Calculate the transition probabilities for start from the embedding.
96         Default implementation assumes the Skip-gram model.
97         (The parameter step is not used.)
98         """
99         return np.exp(self._embedding[self._nodes[node]] @ self._embedding.T)
100
101    def decode(self, node: Any, use_neighborhood: bool = False, no_self_loops:
102              bool = False, normalize: bool = True, **kwargs) -> pd.Series:
103        """Calculate probabilities from the embedding.
104        The parameter node must contain a value from self.nodes.
105
106        Parameters
107        -----
108        node : Any
109            key of the embedding.
110
111        use_neighborhood : bool, optional (default = False)
112            restrict the probabilities to the neighborhood (using knowledge about
113            network topology).
114
115        no_self_loops : bool, optional (default = False)
116            set the probability of a self loop to zero.
117
118        normalize : bool, optional (default = True)
119            normalize the probabilities, such that they sum up to one.
120
121        **kwargs are passed to _decode_raw(node)
122        """
123
124        data = self._decode_raw(node=node, **kwargs)
125        if use_neighborhood:
126            data_n = np.zeros(len(data))
127            for _, next_node, _ in self._gen.transition_probs(start=(node,)):
128                i = self._nodes[next_node]
129                data_n[i] = data[i]
130            data = data_n
131
132        if no_self_loops:
133            data[self._nodes[node]] = 0
134
135        if normalize:
136            data = data / max(1e-16, data.sum())
137
138        return pd.Series(data, index=self.nodes_str, name=self.key2str(node))
139
140    def decode_path(self, start: Tuple[Any,...], **kwargs) -> pd.Series:
141        """Returns decode() for the last node of the path.
142        """
143        return self.decode(start[-1], **kwargs)
144
145    class ABCAsymmetricEmbedding(ABCEmbedding):
146        def __init__(self, gen: ABCHigherOrderPathGenerator, dimension: int,
147                     source_paths: Iterator[Tuple[Any,...]], target_nodes: Iterator[Any]):
148            super().__init__(gen, dimension, symmetric=False)
149            self._source_paths = { n:i for i,n in enumerate(source_paths) }
150            self._target_nodes = { n:i for i,n in enumerate(target_nodes) }
151            self._source_embedding = np.zeros(shape=(len(self._source_paths),
152                                                    dimension))
153            self._target_embedding = np.zeros(shape=(len(self._target_nodes),
154                                                    dimension))
155            self.path2str = repr
156            self.node2str = repr
157
158        def __getitem__(self, key) -> pd.Series:
159            """If the key contains a node, a row of the target_embedding is returned.
160            And if the key contains a tuple of nodes, a row of the source_embedding
161            matrix is returned instead.
162            """
163            idx = self._target_nodes.get(key, None)
164            if idx is not None:
165                return pd.Series(self._target_embedding[idx,:], name=self.node2str(key))
166

```

```

161     idx = self._source_paths.get(key, None)
162     if idx is not None:
163         return pd.Series(self._source_embedding[idx, :], name=self.path2str(key))
164     else:
165         raise IndexError('Unknown key %s' % key, key=key)
166
167     def node2key(self, node: Any, use_source: bool=False):
168         return (node,) if use_source else node
169
170     @property
171     def source_paths(self) -> Iterator[Tuple[Any, ...]]:
172         return self._source_paths.keys()
173
174     @property
175     def source_paths_str(self) -> List[str]:
176         return list(map(self.path2str, self.source_paths))
177
178     @property
179     def target_nodes(self) -> Iterator[Any]:
180         return self._target_nodes.keys()
181
182     @property
183     def target_nodes_str(self) -> List[str]:
184         return list(map(self.node2str, self.target_nodes))
185
186     @property
187     def source_embedding(self) -> pd.DataFrame:
188         return pd.DataFrame(self._source_embedding, index=self.source_paths_str)
189
190     @property
191     def target_embedding(self) -> pd.DataFrame:
192         return pd.DataFrame(self._target_embedding, index=self.target_nodes_str)
193
194     @property
195     def dimension(self):
196         "effective dimension of embedding"
197         return self._source_embedding.shape[1]
198
199     def _decode_raw(self, start: Tuple[Any, ...], step = None) -> np.array:
200         """Calculate the transition probabilities for start from the embedding.
201             Default implementation assumes the Skip-gram model.
202             (The parameter step is not used.)
203         """
204         return np.exp(self._source_embedding[self._source_paths[start]] @ self._target_embedding.T)
205
206     def decode(self, start: Tuple[Any, ...], use_neighborhood: bool = False,
207               no_self_loops: bool = False, normalize: bool = True, **kwargs) -> pd.Series:
208         """Calculate probabilities from the embedding.
209             The parameter start must contain a value from self.source_paths.
210
211             Parameters
212             -----
213
214             start : Tuple[Any, ...]
215
216                 key of the source embedding.
217
218             use_neighborhood : bool, optional (default = False)
219
220                 restrict the probabilities to the neighborhood (using knowledge about
221                 the first-order network topology).
222
223             no_self_loops : bool, optional (default = False)
224
225                 set the probability of a self loop to zero.
226
227             normalize : bool, optional (default = True)
228
229                 normalize the probabilities, such that they sum up to one.
230
231             **kwargs are passed to _decode_raw(start)

```

```

230 """
231     data = self._decode_raw(start=start, **kwargs)
232     if use_neighborhood:
233         data_n = np.zeros(len(data))
234         for _, next_node, _ in self._gen.transition_probs(start=start[-1:]):
235             i = self._target_nodes[next_node]
236             data_n[i] = data[i]
237         data = data_n
238     if no_self_loops:
239         data[self._target_nodes[start[-1]]] = 0
240     if normalize:
241         data = data / max(1e-16, data.sum())
242     return pd.Series(data, index=self.target_nodes_str, name=self.path2str(
243         start))
244
245     def decode_path(self, start: Tuple[Any, ...], **kwargs) -> pd.Series:
246         """Finds the matching source path and returns decode()
247         if start in self._source_paths:
248             return self.decode(start, **kwargs)
249         if (len(start) > self._gen.max_rule_key_length): # and (self._gen.
250             max_rule_key_length > 0):
251             start = start[-self._gen.max_rule_key_length:]
252         while len(start)>0:
253             if start in self._source_paths:
254                 return self.decode(start, **kwargs)
255             else:
256                 start = start[1:]
257         return None
258
259     @staticmethod
260     def factor_matrix(mat: np.ndarray, dimension: int=None, check=False)->
261         Tuple[np.ndarray, np.ndarray]:
262         """Factor matrix using truncatedSVD
263         if (dimension is None) or (dimension >= min(mat.shape)): # SVD
264             dimension = min(mat.shape)
265             U,S,Vh = svd(mat)
266             sqrtS = S**0.5
267             if mat.shape[0] > mat.shape[1]:
268                 U = U[:, :dimension:] # U[:,1:] drops 2nd dimension (squeezing),
269                 while U[:, 1:] does not
270             elif mat.shape[0] < mat.shape[1]:
271                 Vh = Vh[:dimension:, :]
272                 source_embedding = U * (sqrtS[np.newaxis,:])
273                 target_embedding = np.transpose((sqrtS[:,np.newaxis]) * Vh)
274             else:
275                 tsvd = TruncatedSVD(n_components=dimension, random_state=1)
276                 tsvd.fit(mat) # returns svd
277                 W = tsvd.fit_transform(mat) # = U * Sigma, shape=(mat.shape[0],
278                                         dimension)
279                 H = tsvd.components_ # = V', shape=(dimension, mat.shape[1])
280                 # checking factorization
281                 #print('diff %g' % abs(mat - W@H).max())
282                 S = np.linalg.norm(W, axis=0) # len(S) == dimension
283                 S = S.clip(min=1e-20) # avoid division by zero
284                 source_embedding = W / ((S**0.5)[np.newaxis,:,:]) # = W @ np.diag(S
285                                         **-0.5)
286                 target_embedding = np.transpose(H * ((S**0.5)[:,np.newaxis])) # = (np.
287                                         diag(S**0.5) @ H).T
288                 if check: # checking factorization
289                     error = np.linalg.norm(mat - source_embedding @ target_embedding.
290                                         transpose())
291                     print('Approximation error %g' % (error/np.linalg.norm(mat)))
292             return (source_embedding, target_embedding)
293
294     class Generic_SkipGram_Embedding(ABCSymmetricEmbedding):
295         """Wrapper for embeddings based on the skip-gram model calculated externally
296         (e.g., LINE, Node2vec)"""
297         def __init__(self, gen: ABCHigherOrderPathGenerator, emb_path : str, id:
298             str = None, parse_node=lambda x:x, binary: bool=False, config=dict()):
299             assert binary==False, 'not implemented'
300             with open(emb_path, 'r') as f:
301                 dims = f.readline().split(' ')
302                 emb = pd.read_csv(f, sep=' ', header=None, comment='%',
303                                 index_col=0)
304             if all(emb[emb.columns[-1]].isna()):

```

```

294     # The output of LINE has a training space resulting in a columns with
295     # NaN; Node2vec does not.
296     emb.drop(emb.columns[-1], axis=1, inplace=True)
297     super().__init__(gen, dimension=int(dims[1]), nodes=gen.target_nodes)
298     self._id = os.path.splitext(os.path.split(emb_path)[-1])[0] if id is None
299     else id
300     self._emb_path = emb_path
301     #self._binary = binary
302     assert self._embedding.shape == emb.values.shape, f'shapes do not match:{self._embedding.shape}!={emb.values.shape}'
303     # verify row names and sort rows
304     emb.index = emb.index.map(parse_node)
305     keys_emb = set(emb.index)
306     keys_gen = set(gen.target_nodes)
307     if len(keys_emb.symmetric_difference(keys_gen))>0:
308         print('keys do not match (did you specify parse_node?)')
309         print('only in embedding:', keys_emb-keys_gen)
310         print('only in generator:', keys_gen-keys_emb)
311         assert False
312     ##emb.sort_index(axis=0, inplace=True, key=gen._node_sort_key) # requires
313     # pandas 1.1 and _node_sort_key must be vectorized
314     emb_sort = pd.Series(data=emb.index, index=emb.index).map(gen.
315     _node_sort_key).sort_values().index
316     emb = emb.loc[emb_sort]
317     self._embedding = emb.values
318     self._config = config
319
320     def train(self):
321         pass
322
323     @property
324     def config(self):
325         """get configuration"""
326         cfg = dict(init_class=self.__class__.__name__, init_gen=self._gen._id,
327         init_emb_path = self._emb_path, init_id=self._id), init_binary=self
328         ._binary)
329         cfg.update(self._config)
330         return cfg
331
332     class HON_DeepWalk_EMBEDDING(ABCSymmetricEmbedding):
333         """
334             Adapts [DeepWalk] to random walks in higher order models, see
335             [DeepWalk] Perozzi B., Al-Rfou R., and Skiena S. (2014)
336             'Deepwalk: Online learning of social representations',
337             https://doi.org/10.1145/2623330.2623732
338
339             Uses gensim.models.Word2Vec for the embedding, which treats the random
340             walks as bidirectional. Hence, the embedding is symmetric.
341             """
342
343         def __init__(self, gen: ABCHigherOrderPathGenerator, dimension: int,
344             reuse_walks: bool = False):
345             super().__init__(gen, dimension, gen.target_nodes)
346             self._id = 'DeepWalk'
347             self._reuse_walks = reuse_walks
348             self._walks = None
349             self._training_config = {'num_walks': 0, 'walk_length': 0, 'random_seed': 0}
350
351         # https://stackoverflow.com/questions/34831551
352         # Word2vec uses the hash function to initialize the vector for each word,
353         # no need for a cryptographic hash function
354         if sys.hash_info.width == 64:
355             @staticmethod
356             def str_hash(data:str) -> int:
357                 "deterministic hash function for Word2Vec (64 bit)"
358                 b = data.encode()
359                 return adler32(b) + (crc32(b) << 32)
360
361         else:
362             @staticmethod
363             def str_hash(data:str) -> int:
364                 "deterministic hash function for Word2Vec (32 bit)"
365                 return crc32(data.encode())
366
367         def train(self, num_walks:int=100, walk_length:int=80, window_size=10,

```

```

358     num_iter:int=1, min_count:int=0,
359     hs:bool = True, negative:int=5, workers:int=1, random_seed=None,
360     replace_hash:bool=True, **kwargs):
361     # reuse_walks allows for efficiently evaluating different params
362     # of Word2vec (e.g., negative, num_iter).
363     # walks depend only on num_walks, walk_length, and random_seed.
364     tc = self._training_config # config of previous training
365     can_reuse_walks = tc['num_walks']==num_walks and tc['walk_length']==
366     walk_length and tc['random_seed']==random_seed
367     self._training_config = dict(num_walks=num_walks, walk_length=walk_length
368     , window_size=window_size,
369     num_iter=num_iter, min_count=min_count, hs=hs, negative=negative,
370     workers=workers,
371     random_seed=random_seed, replace_hash=replace_hash, **kwargs)
372     if (not can_reuse_walks) or (self._walks is None):
373         walks = list()
374         rng = np.random.default_rng(random_seed) # for shuffle and random walks
375         for _ in range(num_walks):
376             nodes_s = list(self.nodes)
377             rng.shuffle(nodes_s)
378             for node in nodes_s:
379                 walks.append(self._gen.random_walk(start=(node,), num_steps=
380                 walk_length, rng=rng))
381         walks = [list(map(str, walk)) for walk in walks]
382         if self._reuse_walks:
383             self._walks = walks
384         else:
385             walks = self._walks
386         if replace_hash:
387             if 'hashfxn' in kwargs:
388                 print('The parameter hashfxn is ignored because of replace_hash=True')
389             kwargs['hashfxn'] = self.str_hash
390         model = Word2Vec(walks, size=self._dimension, window=window_size, sg=1,
391             min_count=min_count, iter=num_iter,
392             hs=1 if hs else 0, negative=0 if hs else negative, workers=workers,
393             seed=random_seed, **kwargs)
394         self.model = model # debug
395         for v,iv in self._nodes.items():
396             self._embedding[iv,:] = model.wv[str(v)]
397
398     @property
399     def config(self):
400         """get configuration
401         cfg = dict(init_class=self.__class__.__name__,
402                     init_gen=self._gen._id,
403                     init_dimension=self._dimension, init_id=self._id)
404         cfg.update(self._training_config)
405         return cfg
406
407     class HON_Node2vec_EMBEDDING(HON_DeepWalk_EMBEDDING):
408         """
409         Adapts [Node2vec] to random walks in higher order models, see
410         [Node2vec] Grover A. and Leskovec J. (2016)
411         'node2vec: Scalable Feature Learning for Networks', https://doi.org/10.1145/2939672.2939754
412
413         The differences between Node2vec and DeepWalk are:
414         - Node2vec uses biased random walks, implemented in
415             Node2vec_HigherOrderPathGenerator
416         - Both rely on gensim.models.Word2Vec, but Node2vec uses negative sampling
417             and DeepWalk uses hierarchical softmax.
418         """
419         def __init__(self, gen: ABCHigherOrderPathGenerator, dimension: int = 128,
420             p: float = 1, q: float = 1, reuse_walks:bool=False):
421             gen_n2v = Node2vec_HigherOrderPathGenerator(gen, p, q)
422             super().__init__(gen_n2v, dimension, reuse_walks)
423             self._id = 'Node2vec(p={0},q={1})'.format(p,q)
424
425         def set_params(self, p: float = 1, q: float = 1):
426             self.gen.set_params(p, q) # do not call this directly
427             self._id = 'Node2vec(p={0},q={1})'.format(p,q)
428             self._walks = None
429             #self._embedding = np.zeros(self._embedding.shape) # reset embedding

```

```

418     def train(self, negative:int=5, hs:bool = False, **kwargs): # Node2Vec uses
419         negative sampling by default
420         super().train(hs=hs, negative=negative, **kwargs)
421
422     class HONEM_EMBEDDING(ABCAsymmetricEmbedding):
423         """
424             Calculates the embedding according to
425             [HONEM] Saebi M., Ciampaglia G., Kaplan L., and Chawla N. (2019)
426             'HONEM: Network Embedding Using Higher-Order Patterns in Sequential Data',
427             arXiv:1908.05387
428
429             The instance of class HigherOrderPathGenerator is assumed to contain the
430                 rules detected by BuildHON+, see
431             [BuildHON+] Xu J., Saebi M., Ribeiro B., Kaplan L., and Chawla N. (2017)
432             'Detecting Anomalies in Sequential Data with Higher-order Networks', arXiv
433                 :1712.09658
434             """
435     def neighborhood_matrix(self, order: int=1, sort=True) -> pd.DataFrame:
436         """Calculates the 'v-th order neighborhood matrix' defined in the [HONEM] paper"""
437         keys = [key for key in self._gen.source_paths if len(key)==order]
438         data = [(self.path2str((key[0],)), self.node2str(next_node), prob)
439                 for key in keys for _, next_node, prob in self._gen.transition_probs(
440                     key)]
441         data_df = pd.DataFrame(data, columns=['src','trg','prob'])
442         res = data_df.groupby(['src','trg']).mean().unstack(fill_value=0)
443         res.columns = [c[1] for c in res.columns] # c[0] == 'prob'
444
445         # ensure that neighborhood matrices of different orders have identical
446             indices and columns
447         all_source_nodes = { self.path2str((key[0],)) for key in self._gen.
448             source_paths_len1 }
449         for v in all_source_nodes.difference(res.index):
450             res.loc[v]=0
451         all_target_nodes = { self.node2str(node) for node in self._gen.
452             target_nodes }
453         for v in all_target_nodes.difference(res.columns):
454             res[v]=0
455
456         if sort:
457             res = res.loc[self._gen.source_paths_str]
458             res = res[self._gen.target_nodes_str]
459
460     return res
461
462     def __init__(self, gen: HigherOrderPathGenerator, dimension: int):
463         super().__init__(gen, dimension, gen.source_paths_len1, gen.target_nodes)
464         self._id = 'HONEM'
465
466     def train(self):
467         # calculate neighborhood matrix
468         neighborhood = self.neighborhood_matrix()
469         if self._gen.max_rule_key_length > 1:
470             for order in range(2, self._gen.max_rule_key_length+1):
471                 neighborhood = neighborhood + math.exp(1-order) * self.
472                     neighborhood_matrix(order, sort=False)
473
474         # sort
475         neighborhood = neighborhood.loc[self._gen.source_paths_str]
476         neighborhood = neighborhood[self._gen.target_nodes_str]
477
478         source_embedding, target_embedding = self.factor_matrix(neighborhood.
479             values, self._dimension)
480         self._source_embedding = source_embedding
481         self._target_embedding = target_embedding
482
483     def _decode_raw(self, start: Tuple[Any,...], step = None) -> np.array:
484         """Calculate the transition probabilities for start from the embedding. (HONEM does not use the skip-gram model)"""
485         assert type(start) is tuple and len(start)==1, 'start must be a tuple of length 1'
486         return self._source_embedding[self._source_paths[start]] @ self.
487             _target_embedding.T
488
489     @property
490     def config(self):
491

```

```

478     "get_configuration"
479     return dict(init_class=self.__class__.__name__, init_gen=self._gen._id,
480                 init_dimension=self._dimension, init_id=self._id)
480
481 class HON_NetMF_EMBEDDING(ABCAsymmetricEmbedding):
482     """
483         Adapts [NetMF] to random walks in a higher-order model. (Additionally, the
484             assumption of undirected graphs is dropped.)
485
486         [NetMF] Qiu, Jiezhong and Dong, Yuxiao and Ma, Hao and Li, Jian and Wang,
487             Kuansan and Tang, Jie (2018)
488         'Network Embedding as Matrix Factorization: Unifying DeepWalk, LINE, PTE,
489             and Node2vec', https://doi.org/10.1145/3159652.3159706
490     """
491     def __init__(self, gen: ABCHigherOrderPathGenerator, dimension: int,
492                  pairwise=False):
493         source_paths = gen.source_paths_len1 if pairwise else gen.source_paths
494         super().__init__(gen, dimension, source_paths, gen.target_nodes)
495         self._id = 'NetMF_pairs' if pairwise else 'NetMF_paths'
496         self._pairwise = pairwise
497         self._training_config = {}
498         # caching the PMI calculation
499         self._window_size = 0
500         self._PMI = None
501
502     def train(self, window_size:int, negative:int=1, optimized:bool=True):
503         assert window_size > 0, 'window_size must be a positive integer'
504         self._training_config = dict(window_size=window_size, negative=negative,
505                                     optimized=optimized)
506         if window_size != self._window_size:
507             # calculate pointwise mutual information (PMI) - this takes some time
508             PMI = np.zeros(shape=(len(self.source_paths), len(self.target_nodes)))
509             idx = {v:i for i,v in enumerate(self.target_nodes)}
510             # stationary distribution (for target nodes)
511             sd = {v:p for _,v,p in self._gen.transition_probs(start=())}
512             if optimized:
513                 # We want to enumerate all paths and aggregate the visiting
514                 # probabilities of the individual nodes. However, enumerating the
515                 # paths in a depth-first style is generally expensive (exponential
516                 # growth in the length of number of paths). Instead use a breadth-
517                 # first enumeration of the paths. Because of the higher-order Markov
518                 # property of the transition probabilities, we need only to keep track
519                 # of the last k=max_rule_key_length nodes visited. This avoids the
520                 # exponential growth.
521                 # There are max. (#nodes ** max_rule_key_length) different subpaths.
522                 max_len = self._gen.max_rule_key_length
523                 for iu,u in enumerate(self.source_paths):
524                     paths = {u: 1/window_size}
525                     for _ in range(window_size):
526                         paths_new = collections.defaultdict(float)
527                         for source,source_prob in paths.items():
528                             for _,next_node,prob in self._gen.transition_probs(source):
529                                 prob_new = source_prob * prob
530                                 PMI[iu,idx[next_node]] += prob_new / sd[next_node]
531                                 source_new = self._gen.find_rule_key((source, next_node))
532                                 paths_new[source_new] += prob_new
533                         paths = paths_new
534                 else:
535                     # slower implementation enumerating all paths (debugging only)
536                     for iu,u in enumerate(self.source_paths):
537                         for path,prob in self._gen.path_probs(u, num_steps = window_size):
538                             prob_div = prob / window_size # = probability of this path
539                             conditional to start=u, divided by window_size
540                             for v in path[len(u)::]: # skip the start from the path
541                                 iv = idx[v]
542                                 PMI[iu,iv] += prob_div / sd[v]
543             # cache calculation
544             self._window_size = window_size
545             self._PMI = PMI
546             self._negative = negative
547             M = self._PMI / negative
548             mat = np.log(M.clip(1))
549             source_embedding, target_embedding = self.factor_matrix(mat, self.
550                                         _dimension)

```

```

544     self._source_embedding = source_embedding
545     self._target_embedding = target_embedding
546
547     @property
548     def PMI(self) -> pd.DataFrame:
549         return pd.DataFrame(self._PMI, index=self.source_paths_str, columns=self.
550                             target_nodes_str)
551
552     @property
553     def config(self):
554         "get configuration"
555         cfg = dict(init_class=self.__class__.__name__, init_gen=self._gen_id,
556                     init_dimension=self._dimension, init_id=self._id, init_pairwise=self.
557                     _pairwise)
558         cfg.update(self._training_config)
559         return cfg
560
561     class HON_GraRep_EMBEDDING(ABCAsymmetricEmbedding):
562         """
563             Adapts [GraRep] to random walks in a higher-order model.
564
565             [GraRep] Cao, Shaosheng and Lu, Wei and Xu, Qiongkai (2015)
566             'GraRep: Learning Graph Representations with Global Structural Information
567             ', https://doi.org/10.1145/2806416.2806512
568         """
569         def __init__(self, gen: ABCHigherOrderPathGenerator, dimension: int,
570                      num_steps: int = 4, pairwise: bool=True, neg_stationary: bool=False):
571             """
572                 Init
573
574                 Parameters
575                 -----
576                 dimension : int, >0
577                     The dimension of the embedding for each individual step.
578                     Effectively, the dimension is num_steps * dimension.
579
580                 num_steps : int, >0 (default = 4)
581
582                 pairwise : bool (default = True)
583                     Indicates whether nodes (True) or paths are embedded.
584
585                 neg_stationary : bool (default = False)
586                     The PMI is divided by the corresponding averages (False, as in GraRep)
587                     or by the stationary distribution (True, as in HON NetMF).
588             """
589             source_paths = gen.source_paths_len1 if pairwise else gen.source_paths
590             self._dimension_per_step = dimension = min(dimension, len(gen.
591                                         target_nodes)) # dimension returned by factor_matrix()
592             super().__init__(gen, self._dimension_per_step * num_steps, source_paths,
593                             gen.target_nodes)
594             self._dimension = dimension # specified dimension
595             self._id = 'GraRep_pairs' if pairwise else 'GraRep_paths'
596             self._pairwise = pairwise
597             self._num_steps = num_steps
598             self._neg_stationary = neg_stationary # distribution for negative samples
599             : stationary (True) or averages (False)
600             self._training_config = {}
601             # caching the PMI calculation
602             self._PMI = None
603
604         def train(self, negative: int = 1, normalize: bool = False):
605             """
606                 # optimized PMI calculation from HON_NetMF_EMBEDDING (with small
607                 # modifications).
608             # HON_NetMF_EMBEDDING._PMI matches self._PMI.mean(axis=2) if window_size
609             # equals num_steps and neg_stationary=True
610             self._training_config = dict(negative=negative, normalize=normalize)
611             if self._PMI is None:
612                 idx = {v:i for i,v in enumerate(self.target_nodes) }
613                 PMI = np.zeros(shape=(len(self.source_paths), len(self.target_nodes),
614                                 self._num_steps))
615                 max_len = self._gen.max_rule_key_length
616                 for iu,u in enumerate(self.source_paths):
617                     paths = {u: 1} # 1 instead of 1/window_size in NetMF

```

```

607     for step0 in range(self._num_steps): # step0 = step - 1
608         paths_new = collections.defaultdict(float)
609         for source,source_prob in paths.items():
610             for _,next_node,prob in self._gen.transition_probs(source):
611                 prob_new = source_prob * prob
612                 PMI[iu, idx[next_node], step0] += prob_new
613                 source_new = self._gen.find_rule_key((source, next_node))
614                 paths_new[source_new] += prob_new
615         paths = paths_new
616     # The PMI values calculated above correspond to  $p(c|w)$ , and we still
617     # have to divide it by  $p(c)$  - or similar. GraRep divides it by
618     #  $\sum_w p(c|w) / |V|$ , see Table 1 in [GraRep]. (The division by the
619     # number of nodes  $|V|$  is borrowed from the parameter beta.)
620     # neg_stationary=True instead divides by the stationary distribution
621     # for comparison with NetMF.
622     averages = np.zeros(PMI.shape[1])
623     if self._neg_stationary:
624         for _,v,p in self._gen.transition_probs(start=()): # stationary
625             distribution
626             averages[idx[v]] += p
627
628     for step0 in range(self._num_steps): # step0 = step - 1
629         if not self._neg_stationary:
630             averages = PMI[:, :, step0].mean(axis=0) # divide by mean over column
631             PMI[:, :, step0] /= averages[np.newaxis, :]
632
633     self._PMI = PMI
634     self._negative = negative
635     for step0 in range(self._num_steps):
636         M = self._PMI[:, :, step0] / negative
637         mat = np.log(M.clip(1))
638         source_embedding, target_embedding = self.factor_matrix(mat, self.
639             _dimension_per_step)
640         indices = slice(step0 * self._dimension_per_step, (step0 + 1) * self.
641             _dimension_per_step)
642         # Table 1 in [GraRep] does not mention normalizing the individual
643         # embeddings before concatenating them together.
644         # However, https://github.com/ShelsonCao/GraRep scales each embedding
645         # by its L2 norm.
646         if normalize:
647             source_embedding = source_embedding / np.linalg.norm(source_embedding
648             )
649             target_embedding = target_embedding / np.linalg.norm(target_embedding
650             )
651             self._source_embedding[:, indices] = source_embedding
652             self._target_embedding[:, indices] = target_embedding
653
654     def PMI(self, step: int = 1) -> pd.DataFrame:
655         return pd.DataFrame(self._PMI[:, :, step - 1], index=self.
656             source_paths_str, columns=self.target_nodes_str)
657
658     def _decode_raw(self, start: Tuple[Any,...], step: int = 1) -> pd.Series:
659         indices = slice((step - 1) * self._dimension_per_step, step * self.
660             _dimension_per_step)
661         return np.exp(self._source_embedding[self._source_paths[start], indices]
662             @ self._target_embedding[:, indices].T)
663
664     @property
665     def config(self):
666         "get configuration"
667         cfg = dict(init_class=self.__class__.__name__, init_gen=self._gen._id,
668             init_dimension=self._dimension, init_id=self._id,
669             init_num_steps=self._num_steps, init_pairwise=self._pairwise,
670             init_neg_stationary=self._neg_stationary)
671         cfg.update(self._training_config)
672         return cfg
673
674     class HON_Transition_Hierarchical_Embedding(ABCAsymmetricEmbedding):
675         """
676         Extends HON_NetMF_Embedding with window_size=1 and pairwise=False by adding
677         a penalty term which ties the
678         embeddings together along the specified hierarchy.
679         Due to the penalty, we cannot use singular value decomposition (SVD) and
680         had to rely on stochastic gradient

```

```

667     descent (SGD) instead - similar to gensim.models.Word2Vec.
668     We use skip-gram with negative sampling (SGNS) as optimization criterion
669         and added a penalty proportional to
670         sum( norm(u-pu)**2 for u,pu in node_hierarchy.items() ).
671
672     The intuition for introducing the penalty stems from the observation, that
673         the probability of a path
674     (e.g. P[A->B->C]) can be calculated from the sums of probabilities of
675         longer paths (e.g. sum_x P[x->A->B->C])
676     - assuming stationarity or introducing a 'begin of path'-symbol.
677     Hence, the transition probabilities for some start path correspond to
678         averages of transition probabilities
679     for longer paths.
680     While for short paths, there is enough data (i.e. high support for the
681         estimated transition probabilities),
682     longer paths may benefit from being tied towards their parent path (i.e.
683         path[1:]).  

684     The notion of "short" and "longer" is controlled by the parameter
685         min_length of calc_node_hierarchy.  

686
687     Note that, gensim.models.Word2Vec tries very hard to speed calculation up,
688         as explained by a blog article of
689     its author. Unfortunately, this implementation will be much slower.
690     """
691     @staticmethod
692     def calc_node_hierarchy(gen: ABCHigherOrderPathGenerator, min_length:int=1)
693         -> Dict[Tuple[Any,...],Tuple[Any,...]]:
694         """
695         Calculates the default node_hierarchy, which is used in __init__(...,
696             node_hierarchy='calc').  

697
698         min_length specifies the minimal length of those paths, which are
699             affected by the penalty.
700         """
701         assert min_length > 0, 'min_length must be > 1'
702         return {key: key[1:] for key in gen.rule_keys if len(key) >= min_length + 1}  

703
704     def __init__(self, gen: ABCHigherOrderPathGenerator, dimension: int=128,
705         node_hierarchy:Union[str,Dict[Tuple[Any,...],Tuple[Any,...]]]='calc',
706         seed=None, neg_stationary:bool=True):
707         self._neg_stationary = neg_stationary # distribution for negative samples
708             : stationary (True) or uniform (False)
709         source_paths = gen.source_paths
710         target_nodes = gen.target_nodes
711         super().__init__(gen, dimension, source_paths, target_nodes)
712         self._id = 'Hierarchical_GF' # todo
713         self.reset(seed)
714         if node_hierarchy == 'calc':
715             node_hierarchy = self.calc_node_hierarchy(gen)
716             assert isinstance(node_hierarchy, collections.abc.Mapping), "
717                 node_hierarchy must be either the string 'calc' or a dictionary with
718                 keys containing paths and values containing the parent of each path
719                 , i.e. key[1:]"
720             self._parent = node_hierarchy
721             children = collections.defaultdict(set)
722             for u,pu in node_hierarchy.items():
723                 children[pu].add(u)
724             self._children = children
725
726     def clone(self):
727         "Clones the instance"
728         res = HON_Transition_Hierarchical_EMBEDDING(self._gen, self._dimension,
729             self._parent, None, self._neg_stationary)
730         res._total_steps = self._total_steps
731         res._source_embedding = self._source_embedding.copy()
732         res._target_embedding = self._target_embedding.copy()
733         res._rng.bit_generator.state = self._rng.bit_generator.state
734         res._training_history = list(self._training_history)
735         return res
736
737     def reset(self, seed=None):
738         "Resets the embedding (matrices are initialized with random values)."
739         dimension = self._dimension

```

```

723     self._seed = seed
724     self._rng = np.random.default_rng(seed)
725     self._total_steps = 0
726     self._source_embedding = self._rng.normal(0, 1/dimension, (len(self.
727         _source_paths),dimension))
727     self._target_embedding = self._rng.normal(0, 1/dimension, (len(self.
728         _target_nodes),dimension))
729     self._training_history = list()
730
731     @staticmethod
732     def sigmoid(x: float) -> float: # better use a lookup table
733         # https://stackoverflow.com/questions/3985619/how-to-calculate-a-logistic
734         # -sigmoid-function-in-python ... or use lookup table
735         # sigmoid(-710) raises a math range error
736         return 1 / (1 + math.exp(-max(x,-700)))
737
738     def _calc_objective(self) -> Tuple[float,float,float]:
739         "calculates the loss function (SGNS and penalty are returned separately)"
740         loss = 0
741         loss_neg = 0
742         # neg_probs = distribution of negative samples
743         if self._neg_stationary: # stationary distribution
744             neg_probs = { next_node:prob for _,next_node,prob in self._gen.
745                 transition_probs(start=()) }
746         else: # uniform distribution
747             n_target = len(self._target_nodes)
748             neg_probs = { next_node:1/n_target for next_node in self._target_nodes
749                 }
750
751         for u,iu in self._source_paths.items(): # uniform sampling
752             probs = { next_node:prob for _,next_node,prob in self._gen.
753                 transition_probs(u) }
754             for v,iv in self._target_nodes.items():
755                 score = self._source_embedding[iu,:].dot(self._target_embedding[iv,:])
756                 sigmoid_value = self.sigmoid(score)
757                 loss -= math.log(sigmoid_value) * probs.get(v,0)
758                 loss_neg -= math.log(1-sigmoid_value) * neg_probs[v]
759         loss_penalty = 0
760         for u,pu in self._parent.items():
761             iu = self._source_paths[u]
762             ipu = self._source_paths[pu]
763             loss_penalty += np.linalg.norm(self._source_embedding[iu,:] - self.
764                 _source_embedding[ipu,:])**2
765         return (loss, loss_neg, loss_penalty) # objective is loss + negative *
766             loss_neg + penalty * loss_penalty
767
768     def _update(self, iu, iv, label, learning_rate):
769         "SGD update"
770         eu = self._source_embedding[iu,:]
771         ev = self._target_embedding[iv,:]
772         score = eu @ ev
773         grad = learning_rate * (self.sigmoid(score) - label)
774         self._source_embedding[iu,:] -= grad * ev
775         self._target_embedding[iv,:] -= grad * eu
776
777     def _update_hierarchy(self, u, iu, learning_rate, penalty, max_start_len:
778         int):
779         "Calculates the gradient of the penalty and updates the source embedding
780         accordingly."
781         factor = min(2 * learning_rate * penalty, 1) # avoid overshooting, which
782             could result in a "math domain error" in the sigmoid calculation
783         if u in self._children:
784             if len(u) == max_start_len:
785                 return
786             ichildren = list(self._source_paths[cu] for cu in self._children[u])
787             ec = self._source_embedding[ichildren,:].mean(axis=0) # average
788                 embedding for all children
789             self._source_embedding[iu,:] += factor * (ec - self._source_embedding[
790                 iu,:])
791
792         pu = self._parent.get(u, None)
793         if not pu is None:
794             ipu = self._source_paths[pu]
795             self._source_embedding[iu,:] += factor * (self._source_embedding[ipu,:]
796                 - self._source_embedding[iu,:])

```

```

783
784     def train(self, steps: int=1, negative: int=1, penalty: float=0,
785               learning_rate_start: float=0.0025, learning_rate_end: Optional[float]=
786               None,
787               max_start_len: Optional[int]=None, debug_objective: Optional[int]=None):
788         """
789         Trains the embeddings using SGD.
790
791         Parameters
792         -----
793
794         steps : int, >0
795             The number of SGD updates (not counting negative samples) per each path
796             in source_paths.
797
798         negative : int, >= 0
799             Number of negative samples.
800
801         penalty: float, >= 0
802             Factor for the penalty terms.
803
804         learning_rate_start: float, >0
805             Learning rate. Linearly decreasing learning rates are specified with
806             learning_rate_end != None.
807
808         learning_rate_end: float, optional (default=None)
809             Allows for linearly falling learning rates.
810
811         max_start_len: int, optional (default = None)
812             Allows for training source embeddings corresponding to short
813             source_paths only.
814             At the end, the remaining source embeddings are replaced by their
815             corresponding ancestors.
816             (max_start_len = None disables this feature.)
817             Source embeddings for a path longer than max_start_len are replaced by
818             the ones corresponding
819             to shorter ancestors (i.e. path[-max_start_len:]).
```

820

```

821         debug_objective: int, optional (default = None)
822             If not None, this method returns the objective function evaluated in
823             periodic intervals.
824             The length of the interval is specified by the parameter
825             debug_objective.
```

826

```

827         self._training_history.append(dict(steps=steps, negative=negative,
828                                           penalty=penalty,
829                                           learning_rate_start=learning_rate_start, learning_rate_end=
830                                           learning_rate_end,
831                                           max_start_len=max_start_len)) # skip debug_objective
832         list_of_objectives = []
833         list_of_objectives_pos = []
834         list_of_objectives_neg = []
835         list_of_objectives_penalty = []
836         learning_rate_delta = 0 if learning_rate_end is None else
837             learning_rate_end - learning_rate_start
838         if max_start_len is None:
839             source_paths_s = list((u,iu) for u,iu in self._source_paths.items())
840             max_start_len_eff = max(len(key) for key in source_paths_s)
841         else:
842             source_paths_s = list((u,iu) for u,iu in self._source_paths.items() if
843                 len(u) <= max_start_len)
844             max_start_len_eff = max_start_len
845     def copy_long_key_embeddings():
846         if max_start_len is not None:
847             # copy embeddings corresponding to keys longer than max_start_len
848             for u,iu in self._source_paths.items():
849
```

```

844     if len(u) > max_start_len:
845         pu = u[(-max_start_len_eff):] # FIXME: hardcoded assumption about
846         # node_hierarchy!
847         ipu = self._source_paths[pu]
848         self._source_embedding[iu,:] = self._source_embedding[ipu,:]
849
850         n_to = len(self._target_nodes)
851         self._rng.shuffle(source_paths_s)
852         i = 0
853         imax = steps * len(source_paths_s) - 1
854         for step in range(steps):
855             if debug_objective is not None:
856                 copy_long_key_embeddings()
857                 if step % int(debug_objective) == 0: # calculate objective every {
858                     debug_objective} step
859                     loss, loss_neg, penalty_loss = self._calc_objective()
860                     list_of_objectives.append(loss + negative * loss_neg + penalty *
861                     penalty_loss)
862                     list_of_objectives_pos.append(loss)
863                     list_of_objectives_neg.append(loss_neg)
864                     list_of_objectives_penalty.append(penalty_loss)
865                     for u,iu in source_paths_s: # uniform sampling
866                         learning_rate = learning_rate_start + learning_rate_delta * (i/imax)
867                         v = self._gen.random_step(start=u, rng=self._rng)
868                         iv = self._target_nodes[v]
869                         self._update(iu, iv, 1, learning_rate) # 1-sigma(x) = sigma(-x)
870                         if self._neg_stationary: # stationary distribution
871                             for v_neg in self._gen.random_step(start=(), size=negative, rng=
872                             self._rng):
873                                 iv_neg = self._target_nodes[v_neg]
874                                 self._update(iu, iv_neg, 0, learning_rate) # 0-sigma(x) = -sigma(
875                                 x)
876                         else: # uniform distribution
877                             for iv_neg in self._rng.integers(n_to, size=negative):
878                                 self._update(iu, iv_neg, 0, learning_rate) # 0-sigma(x) = -sigma(
879                                 x)
880                         if penalty > 0:
881                             self._update_hierarchy(u, iu, learning_rate, penalty, max_start_len
882                             )
883                             i+=1
884                         copy_long_key_embeddings()
885                         if debug_objective is not None:
886                             # finally, evaluate the objective once more.
887                             loss, loss_neg, penalty_loss = self._calc_objective()
888                             list_of_objectives.append(loss + negative * loss_neg + penalty *
889                             penalty_loss)
890                             list_of_objectives_pos.append(loss)
891                             list_of_objectives_neg.append(loss_neg)
892                             list_of_objectives_penalty.append(penalty_loss)
893                             self._total_steps += steps
894                         if debug_objective is not None:
895                             return dict(total_steps=self._total_steps, steps=steps, penalty=penalty
896                             ,
897                             negative=negative, neg_stationary=self._neg_stationary,
898                             learning_rate_start=learning_rate_start, learning_rate_end=
899                             learning_rate_start+learning_rate_delta,
900                             dimension=self.dimension, max_start_len=max_start_len, objectives=np.
901                             array(list_of_objectives),
902                             objectives_pos=np.array(list_of_objectives_pos), objectives_neg=np.
903                             array(list_of_objectives_neg),
904                             objectives_penalty=np.array(list_of_objectives_penalty))
905
906             @property
907             def config(self):
908                 "get configuration"
909                 cfg = dict(init_class=self.__class__.__name__, init_gen=self._gen._id,
910                 init_dimension=self._dimension, init_id=self._id,
911                 init_seed=self._seed, init_neg_stationary=self._neg_stationary)
912                 # node_hierarchy is missing
913                 for i,th in enumerate(self._training_history):
914                     cfg.update({f'train{i}_{k}':v for k,v in th.items()})
915
916             return cfg

```

## C.6 Visualizations.py

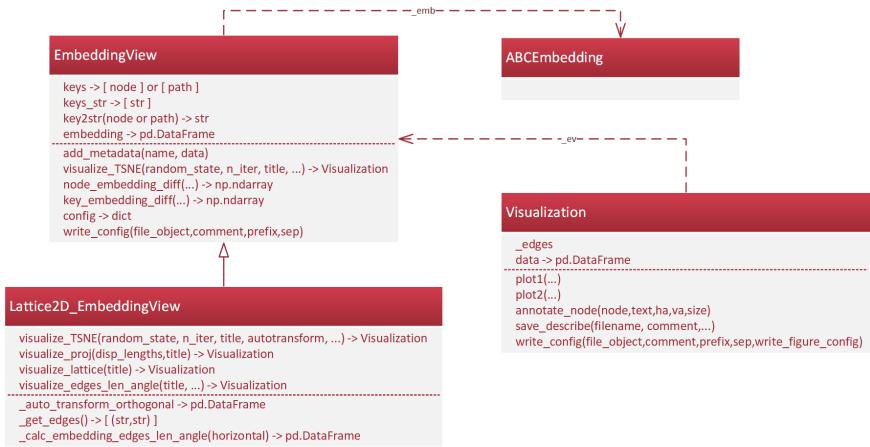


Figure C.3: UML class diagram for Visualizations.py

Due to the decision to implement also ‘mixed’ interactions (see § 3.1.2), a word or context may refer to a single node or a tuple of nodes, for which we use the umbrella term *key* — as already mentioned in § C.1, paragraph embedding. All embedding implementations provide either one or two pandas DataFrames whose rows contain the individual embeddings. How do we locate the embedding corresponding to a given node? First, determine the corresponding key (either the node itself or a tuple containing only the node). Second, we convert the key into a string. This process is simply

$$\text{node} \xrightarrow{\text{node2key}} \text{key} \xrightarrow{\text{key2str}} \text{key\_str},$$

but it is tedious to distinguish, e.g., whether the list of keys equals nodes, source\_paths, or target\_nodes. To obtain a common interface for the different cases, we use the `EmbeddingView` as an adapter, whose primary task is to provide data (i.e., transformation by t-SNE) for the `Visualization`. The latter provides two plot methods, although one was basically obsoleted by adding grid lines to the plot. For a HON Lattice 2D, we use the derived class `Lattice2D_EMBEDDINGVIEW` instead.

```

1  from typing import Any, Callable, Dict, List, Iterator, Tuple, Optional
2  #from SyntheticNetworks import *
3  from Embedding import ABCEmbedding, ABCSymmetricEmbedding
4  from sklearn.manifold import TSNE
5  import matplotlib.pyplot as plt
6  from matplotlib.collections import LineCollection
7  import seaborn as sns
8  import numpy as np
9  import pandas as pd
10 import math
11 import os
12
13 class Visualization(object):
14     def __init__(self, ev: "EmbeddingView", title, data, cols, config=dict(),
15                  edges=None):
16         self._ev = ev
17         self._title = title

```

```

17     self._data = data
18     self._cols = cols
19     self.config = config
20     self._edges = edges
21     self._edge_args = dict(color='gray', linewidth=1, linestyle=':')
22     self._figure = None
23     self._figure_config = {}
24
25     def __repr__(self):
26         othercols = sorted(list(set(self._data.columns).difference(self._cols)))
27         return f"""paths: {self.paths_id}
28 embedding: {self.emb_id}
29 dim: {self._data.shape}
30 embedding columns: {str(self._cols)}
31 other columns: {str(othercols)}"""
32
33     @property
34     def paths_id(self) -> str:
35         "ID of the paths (transition probabilities)"
36         return self._ev._emb._gen._id
37
38     @property
39     def emb_id(self) -> str:
40         "ID of the embedding"
41         return self._ev._emb._id + ('_source' if self._ev.use_source else '')
42
43     @property
44     def data(self) -> pd.DataFrame:
45         "returns a DataFrame containing the (visualization of the) embedding
46         together with explaining attributes."
47         return self._data
48
49     def _add_edges(self, ax, data, **kwargs):
50         if self._edges is None:
51             return
52         data_xy = data[self._cols]
53         lines = list([data_xy.loc[n1], data_xy.loc[n2]] for n1,n2 in self._edges)
54         lc = LineCollection(lines, **kwargs)
55         ax.add_collection(lc)
56
57     def plot1(self, figsize: Tuple[float,float]=(6,4), dpi: int=200, figureargs
58             =dict(),
59             filter_col: Optional[str]=None, filter_values=set(), rotate:float=0,
60             return_figure=False, **kwargs):
61         """
62             Displays a visualization by a seaborn.scatterplot
63
64             Parameters
65             -----
66             figsize : (float,float)
67                 Width, height in inches.
68
69             dpi : int, default=200
70                 Resolution of the figure.
71
72             figureargs : dict
73                 These are passed to pyplot.figure
74
75             filter_col : str, optional
76                 If specified, the data is filtered by this column, see filter_values.
77
78             filter_values : set
79                 If filter_column is specified, consider only records where the column
80                 filter_col has values in filter_values.
81
82             kwargs :
83                 These are passed to seaborn.scatterplot (e.g. hue='col1', style='col2',
84                 alpha=0.5, palette='coolwarm').
85
86             data = self._data
87             x_name, y_name = self._cols

```

```

87     if rotate != 0: # rotate before edges
88         data = data.copy()
89         c,s = math.cos(rotate), math.sin(rotate)
90         x = c*data[x_name] + s*data[y_name]
91         y = -s*data[x_name] + c*data[y_name]
92         data[x_name] = x
93         data[y_name] = y
94     fig = plt.figure(figsize=figsize, dpi=dpi, **figureargs)
95     if self._title != '':
96         fig.suptitle(self._title)
97     self._add_edges(fig.gca(), data=data, **self._edge_args)
98     palette = kwargs.pop('palette', 'coolwarm')
99     if filter_col is not None: # apply filter after edges
100        data = data[data[filter_col].isin(filter_values)]
101    ax = sns.scatterplot(data=data, x=x_name, y=y_name, palette=palette, ax=
102        fig.gca(), **kwargs)
103    ax.set_aspect(1)
104    # caution: by default there nothing to display in legend; should check
105    # for hue, style, ... in kwargs.
106    ax.legend(loc='center_left', bbox_to_anchor=(1, 0.5), ncol=1)
107    self._figure = fig
108    self._figure_config = dict(figsize=figsize, dpi=dpi, figureargs=repr(
109        figureargs),
110        filter_col=filter_col, filter_values=repr(filter_values), rotate=
111        rotate, **kwargs)
112    if return_figure:
113        return fig
114
115    def plot2(self, figsize: Tuple[float,float]=(15,10), dpi: int=200,
116             figureargs=dict(),
117             filter_col: Optional[str]=None, filter_values=set(),
118             return_figure=False, **kwargs):
119        """
120        Displays a visualization by two seaborn.scatterplots.
121        The columns x_orig and y_orig added by Lattice2D_EMBEDDINGVIEW.
122        _add_lattice_coord() are displayed as hue.
123
124        Parameters
125        -----
126        figsize : (float,float)
127            Width, height in inches.
128        dpi : int, default=200
129            Resolution of the figure.
130
131        figureargs : dict
132
133        filter_col : str, optional
134            If specified, the data is filtered by this column, see filter_values.
135
136        filter_values : set
137            If filter_column is specified, consider only records where the column
138            filter_col has values in filter_values.
139
140        return_figure: bool, default=False
141            If True, the figure is returned.
142
143        kwargs :
144            These are passed to seaborn.scatterplot (e.g. style='key_len', alpha
145            =0.5, palette='coolwarm').
146
147        data = self._data
148        x_name, y_name = self._cols
149
150        constrained_layout=figureargs.pop('constrained_layout',True)
151        fig, axs = plt.subplots(1,2, figsize=figsize, dpi=dpi, **figureargs,
152            constrained_layout=constrained_layout)
153        if self._title != '':
154            fig.suptitle(self._title)
155        self._add_edges(axs[0], data, **self._edge_args)
156        self._add_edges(axs[1], data, **self._edge_args)
157
158        if filter_col is not None: # apply filter after edges
159            data = data[data[filter_col].isin(filter_values)]
```

```

152
153     palette = kwargs.pop('palette', 'coolwarm')
154     scatterplot_args = dict(data=data, x=x_name, y=y_name, palette=palette,
155                             **kwargs)
156
157     ax = sns.scatterplot(hue='x_orig', ax=axs[0], **scatterplot_args)
158     ax.set_aspect(1)
159     ax.legend(loc='center_left', bbox_to_anchor=(1, 0.5), ncol=1)
160
161     ax = sns.scatterplot(hue='y_orig', ax=axs[1], **scatterplot_args)
162     ax.set_aspect(1)
163     ax.legend(loc='center_left', bbox_to_anchor=(1, 0.5), ncol=1)
164     self._figure = fig
165     self._figure_config = dict(figsize=figsize, dpi=dpi, figureargs=repr(
166                               figureargs),
167                               filter_col=filter_col, filter_values=repr(filter_values), **kwargs)
168     if return_figure:
169         return fig
170
171
172     def annotate_node(self, node, text=None, ha='center', va='center', size=8):
173         # ha in { 'left', 'center', 'right' }, va in { 'top', 'center', 'bottom' }
174         ev = self._ev
175         coord = self.data.loc[ev.key2str(ev.node2key(node)), self._cols]
176         if text is None:
177             text = str(node)
178         for ax in self._figure.axes:
179             ax.annotate(text, coord, horizontalalignment=ha, verticalalignment=va,
180                         size=size)
181
182     def write_config(self, file_object, comment: str = '', prefix: str = '',
183                     sep: str = '\t', write_figure_config: bool = True):
184         self._ev.write_config(file_object=file_object, comment=comment, prefix=
185                               prefix, sep=sep)
186         file_object.write(prefix + '\n' + prefix + 'Visualization:\n')
187         for k,v in self.config.items():
188             file_object.write(prefix + '%s%s%s\n' % (k,sep,v))
189         if write_figure_config:
190             file_object.write(prefix + '\n' + prefix + 'Plot:\n')
191             for k,v in self._figure_config.items():
192                 file_object.write(prefix + '%s%s%s\n' % (k,sep,v))
193
194     def save_describe(self, filename, comment : str = '', bbox_inches='tight',
195                       **kwargs):
196         "saves the most recent figure created by this instance and exports the
197          config to a text file."
198         self._figure.savefig(filename, bbox_inches=bbox_inches, **kwargs)
199         ev = self._ev
200         emb = ev._emb
201         gen = emb._gen
202         with open(filename + '.txt','w') as f:
203             self.write_config(f, comment=comment, prefix='', sep='\t',
204                              write_figure_config=True)
205
206     class EmbeddingView(object):
207         "Helper class for visualizing embeddings"
208         def __init__(self, emb : ABCEmbedding, use_source: bool = True):
209             self._emb = emb
210             self._series = dict()
211             if isinstance(emb, ABCSymmetricEmbedding):
212                 self._keys = list(emb.nodes)
213                 self._keys_str = emb.nodes_str
214                 self.key2str = emb.key2str
215                 use_source = False
216             elif use_source:
217                 self._keys = list(emb.source_paths)
218                 self._keys_str = emb.source_paths_str
219                 self.key2str = emb.path2str
220             else:
221                 self._keys = list(emb.target_nodes)
222                 self._keys_str = emb.target_nodes_str
223                 self.key2str = emb.node2str
224             self.use_source = use_source # True if key is a tuple of nodes
225             if use_source:

```

```

217     for name, data in emb._gen._source_path_metadata.items():
218         self.add_metadata(name, data)
219     else:
220         for name, data in emb._gen._target_node_metadata.items():
221             self.add_metadata(name, data)
222
223     def add_metadata(self, name, data):
224         # data maps the keys to some grouping criterion
225         self._series[name] = pd.Series({self.key2str(key): data[key] for key in
226                                         self.keys})
227     return self
228
229     def __getitem__(self, series_name) -> pd.Series:
230         return self._series[series_name]
231
232     def _append_metadata(self, data: pd.DataFrame, copy=False) -> pd.DataFrame:
233         res = data.copy() if copy else data
234         for name, s in self._series.items():
235             res[name] = s
236         return res
237
238     def node2key(self, node:Any):
239         "converts a node into a key"
240         #return self._emb.node2key(node, self.use_source) # todo: remove
241         #ABCEmbedding.node2key?
242         return (node,) if self.use_source else node
243
244     @property
245     def keys(self): # -> List[Any]:
246         "Keys used for the embedding"
247         return self._keys
248
249     @property
250     def keys_str(self): # -> List[str]:
251         "String representation (i.e self.key2str) of the keys. (Equals embedding.
252         index)"
253         return self._keys_str
254
255     @property
256     def embedding(self) -> pd.DataFrame:
257         "Embedding, index corresponds to keys_str and columns are the dimensions
258         of the embedding space."
259         if isinstance(self._emb, ABCSymmetricEmbedding):
260             return self._emb.embedding
261         elif self.use_source:
262             return self._emb.source_embedding
263         else:
264             return self._emb.target_embedding
265
266     def node_embedding_diff(self, node_pairs) -> np.ndarray:
267         "For a given list of first order edges (pairs of nodes), determine the
268         difference in the embedding space between them"
269         key_pairs = list((self.node2key(v), self.node2key(w)) for v,w in
270                           node_pairs)
271         return self.key_embedding_diff(key_pairs)
272
273     def key_embedding_diff(self, key_pairs) -> np.ndarray:
274         "For a given list of pairs of keys, determine the difference in the
275         embedding space between them"
276         tmp = np.zeros(shape=(len(key_pairs), self._emb.dimension))
277         i=0
278         for start,end in key_pairs:
279             e_start = self._emb[start] # = self.embedding.loc[key2str(start)]
280             e_end = self._emb[end] # = self.embedding.loc[key2str(end)]
281             tmp[i,:] = e_end - e_start
282             i+=1
283         return tmp
284
285     def visualize_TSNE(self, random_state=1, n_iter=1000, title: Optional[str]=
286         None, **kwargs) -> Visualization:
287         "Returns TSNE visualization"
288         config = dict(creator='visualize_TSNE', random_state=random_state, n_iter
289                     =n_iter, title=title, **kwargs)
290         embedding = self.embedding

```

```

282     cols = ['x', 'y']
283     data = TSNE(n_components=2, random_state=random_state, n_iter=n_iter, **
284                 kwargs).fit_transform(embedding)
285     data = pd.DataFrame(data, index=embedding.index, columns=cols)
286     if title is None:
287         title = 'TSNE(random_state=%d)' % random_state
288     return Visualization(self, title, self._append_metadata(data), cols,
289                          config)
290
291     @property
292     def config(self):
293         "get configuration"
294         return dict(init_class=self.__class__.__name__, use_source=self.
295                     use_source)
296
297     def write_config(self, file_object, comment: str = '', prefix: str = '',
298                      sep: str = '\t'):
299         self._emb.write_config(file_object=file_object, comment=comment, prefix=
299                               prefix, sep=sep)
300         file_object.write(prefix + '\n' + prefix + 'EmbeddingView:\n')
301         for k,v in self.config.items():
302             file_object.write(prefix + '%s%s%s\n' % (k,sep,v))
303
304     class Lattice2D_EMBEDDINGVIEW(EmbeddingView):
305         "EmbeddingView for the synthetic network 'lattice2D_2nd_order_dynamic'"
306         def __init__(self, emb : ABCEmbedding, use_source: bool = True,
307                      edge_distance: Optional[int] = 1):
308             super().__init__(emb, use_source)
309             self.creator = emb._gen.creator # todo: check has_attr and type
310             self._edge_distance = edge_distance
311             if edge_distance is None:
312                 self._hor_edge_pairs = []
313                 self._ver_edge_pairs = []
314             elif edge_distance == 1:
315                 self._hor_edge_pairs = self.creator.horizontal_edges1
316                 self._ver_edge_pairs = self.creator.vertical_edges1
317             else:
318                 self._hor_edge_pairs = list((k1[0],k2[1]) for k1,k2 in self.creator.
319                                             horizontal_edges2)
320                 self._ver_edge_pairs = list((k1[0],k2[1]) for k1,k2 in self.creator.
321                                             vertical_edges2)
322
323         def visualize_proj(self, disp_lengths: bool=True, title: Optional[str]=None
324                           ) -> Visualization:
325             """
326             Returns a projection visualization.
327             The projection is determined utilizing knowledge about lattice structure
328             (i.e. 'right' and 'down' neighbors).
329             """
330             config = dict(creator='visualize_proj', disp_lengths=disp_lengths, title=
331                           title)
332             # Since TSNE may distort the data, find average directions of horizontal
333             # and vertical edges
334             max_key_len = max(len(key) for key in self._keys) if self.use_source else
335                         1
336             if max_key_len == 1:
337                 hor_edge_diffs = self.node_embedding_diff(self.creator.
338                                              horizontal_edges1)
339                 ver_edge_diffs = self.node_embedding_diff(self.creator.vertical_edges1)
340             else:
341                 hor_edge_diffs = self.key_embedding_diff(self.creator.horizontal_edges2)
342                 ver_edge_diffs = self.key_embedding_diff(self.creator.vertical_edges2)
343             if disp_lengths:
344                 hor_edge_lengths = np.linalg.norm(hor_edge_diffs, axis=1)
345                 ver_edge_lengths = np.linalg.norm(ver_edge_diffs, axis=1)
346                 print('max_key_len', max_key_len)
347                 print('average_horizontal_edge_length', hor_edge_lengths.mean())
348                 print('average_vertical_edge_length', ver_edge_lengths.mean())
349                 print('ratio', hor_edge_lengths.mean()/ver_edge_lengths.mean())
350                 mean_hor_edge = hor_edge_diffs.mean(axis=0)
351                 mean_ver_edge = ver_edge_diffs.mean(axis=0)
352                 embedding = self.embedding
353                 cols = ['proj_hor', 'proj_ver']

```

```

341     emb_proj_hor = embedding @ (mean_hor_edge / np.linalg.norm(mean_hor_edge))
342     emb_proj_ver = embedding @ (mean_ver_edge / np.linalg.norm(mean_ver_edge))
343     data = pd.DataFrame(np.array([emb_proj_hor, emb_proj_ver]).T, index=
344         embedding.index, columns=cols)
345     if title is None:
346         title = 'Projection\u20a6along\u20a6average\u20a6horizontal\u20a6&\u20a6vertical\u20a6edge'
347     edges = self._get_edges()
348     return Visualization(self, title, self._append_metadata(data), cols,
349         config, edges)
350
351     def visualize_lattice(self, title: Optional[str]=None) -> Visualization:
352         """Displays the coordinates of the embeddings"""
353         config = dict(creator='visualize_lattice', title=title)
354         data = pd.DataFrame({'x': self._series['x_orig'], 'y': self._series['
355             y_orig']})
356         cols = data.columns.tolist()
357         edges = self._get_edges()
358         return Visualization(self, title, self._append_metadata(data), cols,
359             config, edges)
360
361     def _auto_transform_orthogonal(self, vis_data, verbose=False) -> pd.
362         DataFrame:
363         """
364             Visualizations are often defined in terms of relative distances,
365             which are not affected by orthogonal transformations.
366             Given the ground truth is known, let us try to align the visualization.
367             """
368             # We would like to use horizontal/vertical edges to determine the
369             # direction of the x and y axis.
370             # The edges should be oriented from left to right and from down to up.
371             # However, some embeddings separate the nodes into two clusters (even and
372             # odd).
373             # Therefore, take horizontal and vertical edges of distance 2 unless
374             # specified otherwise.
375             def vis_diff(vis_data, node_pairs):
376                 tmp = np.zeros(shape=(len(node_pairs), vis_data.shape[1]))
377                 i=0
378                 for start,end in node_pairs:
379                     v_start = vis_data.loc[self.key2str(self.node2key(start))]
380                     v_end = vis_data.loc[self.key2str(self.node2key(end))]
381                     tmp[i,:] = v_end - v_start
382                     i += 1
383                 return tmp
384             hor_edge_diffs = vis_diff(vis_data, self._hor_edge_pairs)
385             ver_edge_diffs = vis_diff(vis_data, self._ver_edge_pairs)
386             hor_edge_mean = hor_edge_diffs.mean(axis=0)
387             ver_edge_mean = ver_edge_diffs.mean(axis=0)
388             det = np.linalg.det(np.matrix([hor_edge_mean, ver_edge_mean]))
389             if verbose:
390                 print(hor_edge_mean, ver_edge_mean, det)
391             if det < 0:
392                 y_col = vis_data.columns[1]
393                 vis_data = vis_data.copy()
394                 vis_data[y_col] = -vis_data[y_col]
395                 hor_edge_mean[1] *= -1
396                 ver_edge_mean[1] *= -1
397                 # determine the angle by which hor_edge_mean should be rotated to point
398                 # towards [1,0].
399                 angle1 = -math.atan2(hor_edge_mean[1], hor_edge_mean[0])
400                 c, s = math.cos(angle1), math.sin(angle1)
401                 rot = np.matrix([[c,-s],[s,c]])
402                 if verbose:
403                     print('1st\u20a6rotation\u20a6by', angle1, 'results\u20a6in', (rot @ hor_edge_mean).A1
404                         , (rot @ ver_edge_mean).A1)
405
406                 # Determine the angle by which ver_edge_mean should be rotated to point
407                 # towards [0,1].
408                 # To avoid calculations modulo 2*pi estimate the angle2 after applying
409                 # the rotation by angle1.
410                 # Checking the determinant already ensures that angle2 is in [-pi/2, pi
411                 # /2], and I expect angle2 to be small.
412                 rot_ver_edge_mean = (rot @ ver_edge_mean).A1

```

```

400     angle2 = math.atan2(rot_ver_edge_mean[0], rot_ver_edge_mean[1])
401     if verbose:
402         print('angle2', angle2)
403     angle = angle1 + angle2/2 # take average
404     c,s = math.cos(angle),math.sin(angle)
405     rot = np.matrix([[c,-s],[s,c]])
406     if verbose:
407         print('2nd_rotation_by', angle, 'results_in', (rot @ hor_edge_mean).A1,
408               (rot @ ver_edge_mean).A1)
408     return pd.DataFrame(vis_data.values @ rot.T, index=vis_data.index,
409                          columns=vis_data.columns)
409
410 def _get_edges(self):
411     if self._edge_distance is None:
412         return None
413     edges = self._hor_edge_pairs + self._ver_edge_pairs
414     # convert pairs of nodes into pairs of string representations of keys
415     return list((self.key2str(self.node2key(n1)),self.key2str(self.node2key(
416         n2))) for n1,n2 in edges)
416
417 def visualize_TSNE(self, random_state=1, n_iter=1000, title: Optional[str]=
418     None,
419     autotransform: bool = True, autotransform_verbose:bool = False, **kwargs) -> Visualization:
420     "Returns a TSNE visualization"
421     config = dict(creator='visualize_TSNE', random_state=random_state, n_iter=
422     n_iter, title=title,
423     autotransform=autotransform, **kwargs)
424     embedding = self.embedding
425     cols = ['x','y']
426     data = TSNE(n_components=2, random_state=random_state, n_iter=n_iter, **
427     kwargs).fit_transform(embedding)
428     data = pd.DataFrame(data, index=embedding.index, columns=cols)
429     if title is None:
430         title = 'TSNE(random_state=%d)' % random_state
431     if autotransform:
432         data = self._auto_transform_orthogonal(data, verbose=
433             autotransform_verbose)
434     edges = self._get_edges()
435     return Visualization(self, title, self._append_metadata(data), cols,
436     config, edges)
437
438 def _calc_embedding_edges_len_angle(self, horizontal: bool):
439     "Calculates lengths and angles (to their average) of embedded edges"
440     edges = self._hor_edge_pairs if horizontal else self._ver_edge_pairs
441     keys = list(f'{n1}→{n2}' for n1,n2 in edges)
442     emb_edges = self.node_embedding_diff(edges)
443     mean_edge = emb_edges.mean(axis=0)
444     mean_edge1 = mean_edge / np.linalg.norm(mean_edge)
445     edges_len = np.linalg.norm(emb_edges, axis=1)
446     edges_spr = (emb_edges @ mean_edge1) / edges_len
447     edges_angle = np.arccos( edges_spr.clip(-1,1) )
448     return pd.DataFrame({'len': edges_len, 'angle': edges_angle, 'angle360':
449         edges_angle *180/math.pi,
450         'direction': ('horizontal' if horizontal else 'vertical')}, index=
451         keys)
452
453 def visualize_edges_len_angle(self, title: Optional[str]=None,
454     figsize: Tuple[float,float]=(6,4), dpi: int=200, figureargs=dict(), **kwargs):
455     config=dict(creator='visualize_edges_len_angle')
456     hor_stats = self._calc_embedding_edges_len_angle(True)
457     ver_stats = self._calc_embedding_edges_len_angle(False)
458     stats = hor_stats.append(ver_stats)
459     stats.sort_values('len', inplace=True) # in case of overplotting, treat
460         # horizontal and vertical edges equally
461     cols = ['len','angle']
462     if title is None:
463         title = 'Edge lengths and angles'
464     # This plot serves only to illustrate conceptual limitations of
465         # visualize_proj(),
466     # which does not justify changes in class Visualization.
467     # As a workaround, return an instance of Visualization, which has already

```

```

    created a figure.
460 vis = Visualization(self, title, stats, cols, config, edges=None)
461 vis.plot1(figsize=figsize, dpi=dpi, hue='direction', figureargs=
462     figureargs, **kwargs)
463 ax = vis._figure.gca()
464 ax.set_aspect('auto')
465 ax.set_xlim(0, None)
466 ax.set_ylim(0, math.pi)
467 ax.tick_params(labelright=True, right=True)
468 ax.set_yticks(math.pi * np.linspace(0,1,num=5))
469 ax.set_yticklabels(['0',r'$\frac{\pi}{4}$',r'$\frac{\pi}{2}$',r'$\frac{3\pi}{4}$',r'$\pi$'])
470 ax.hlines(math.pi/2, *ax.get_xlim(), linewidths=0.5)
471 ax.legend(loc='upper left')
472 vis._figure_config['workaround'] = 'see visualize_edges_len_angle' # let
473     save_describe() know about this workaround
474 return vis
475
476 @property
477 def config(self):
478     "get configuration"
479     cfg = super().config
480     cfg['edge_distance'] = self._edge_distance
481     if self._edge_distance is not None:
482         cfg['mean_horizontal_edge_length'] = self.
483             _calc_embedding_edges_len_angle(True)['len'].mean()
484         cfg['mean_vertical_edge_length'] = self._calc_embedding_edges_len_angle
485             (False)['len'].mean()
486         cfg['stretch'] = cfg['mean_vertical_edge_length'] / cfg['mean_
487             horizontal_edge_length']
488     return cfg
489
490 ## Usage:
491 # from Visualizations import Lattice2D_EMBEDDINGVIEW, EmbeddingView,
492 #     Visualization
493 # from SyntheticNetworks import create_lattice_2nd_order_dynamic
494 # from Embedding import HON_DeepWalk_EMBEDDING
495 # latgen = create_lattice_2nd_order_dynamic(size=10, omega=0.5, lattice_sep
496 #     ='-')
497 # emb = HON_DeepWalk_EMBEDDING(latgen, 128)
498 # %time emb.train(random_seed=1)
499 # ev = Lattice2D_EMBEDDINGVIEW(emb)
# vis = ev.visualize_TSNE(random_state=1, n_iter=1000)
# vis.plot1(hue='even_odd', style='direction', rotate=2.4)
# vis.plot2(style='direction', alpha=0.5)
#
# only_len_1 = dict(filter_col='key_len', filter_values={1})
# only_len_2 = dict(filter_col='key_len', filter_values={2})
# vis.plot2(style='direction', alpha=0.5, **only_len_2)

```

## C.7 Lattice2D\_sim.ipynb

For the visualization experiment in § 4.1, we performed a grid search over the parameter space (figures 4.1 and 4.4). Below, we display a shortened variant of the corresponding jupyter notebook, `Lattice2D_sim (short).ipynb`. Other plots were created by `Plots_ExpVis_synth.ipynb` and `Plots_ExpVis_Word2vec.ipynb`, which we omit here, because the underlying functionality is part of the infrastructure — with the exception of the concentration analysis of the PMI (figure 4.3) and the neighborhood matrix (figure 4.6, left).

11/18/2020

Lattice2D\_sim (short)

In [1]:

```
comment = 'calculated in Lattice2D_sim (short).ipynb'
verbose = True
```

In [2]:

```
%matplotlib inline
#%%matplotlib notebook
import numpy as np
import pandas as pd
import math
import os
import matplotlib.pyplot as plt
#from tqdm import tqdm
from tqdm.notebook import tqdm
```

In [3]:

```
from HigherOrderPathGenerator import CrossValidation_HigherOrderPathGenerator
from Embedding import HON_DeepWalk_EMBEDDING, HON_Node2vec_EMBEDDING, HONEM_EMBEDDING,
HON_NetMF_EMBEDDING, HON_GraRep_EMBEDDING, HON_Transition_Hierarchical_EMBEDDING
from SyntheticNetworks import create_lattice_2nd_order_dynamic
from Visualizations import Visualization, EmbeddingView, Lattice2D_EMBEDDINGVIEW
```

In [4]:

```
size = 10
omega = 0.5
gen = create_lattice_2nd_order_dynamic(size, omega, lattice_sep='-', check=False)
```

In [5]:

```
if not os.path.exists('tmp'):
    os.makedirs('tmp')
```

In [6]:

```
def iter_cross_prod(**kwargs):
    #kwargs = { k:list(v) for k,v in kwargs.items() } # generators to list
    total_len = np.prod([len(v) for v in kwargs.values()])
    for i in tqdm(range(total_len)):
        current = dict()
        #for k,v in kwargs.items(): # simpler, but wrong order
        #    current[k] = v[i % len(v)]
        #    i = i // len(v)
        j = total_len
        for k,v in kwargs.items():
            j = j // len(v) # = product of lengths over remaining items
            current[k] = v[i // j]
            i = i % j
        yield current
```

11/18/2020

Lattice2D\_sim (short)

In [7]:

```
# verifying order of iteration
for current in iter_cross_prod(a=[1,2], b=[3,4,5], c=[6,7]):
    print(current)
```

```
{'a': 1, 'b': 3, 'c': 6}
{'a': 1, 'b': 3, 'c': 7}
{'a': 1, 'b': 4, 'c': 6}
{'a': 1, 'b': 4, 'c': 7}
{'a': 1, 'b': 5, 'c': 6}
{'a': 1, 'b': 5, 'c': 7}
{'a': 2, 'b': 3, 'c': 6}
{'a': 2, 'b': 3, 'c': 7}
{'a': 2, 'b': 4, 'c': 6}
{'a': 2, 'b': 4, 'c': 7}
{'a': 2, 'b': 5, 'c': 6}
{'a': 2, 'b': 5, 'c': 7}
```

In [8]:

```
def get_stretch_ratio(emb, edge_distance = 1):
    ev = Lattice2D_EMBEDDINGVIEW(emb, use_source=True, edge_distance=edge_distance)
    hor_edges_len = np.linalg.norm(ev.node_embedding_diff(ev._hor_edge_pairs), axis=1)
    ver_edges_len = np.linalg.norm(ev.node_embedding_diff(ev._ver_edge_pairs), axis=1)
    return ver_edges_len.mean() / hor_edges_len.mean() # = ev.config['stretch']
```

In [9]:

```
def disp_stretch_W(df, emb=None, title=None, figsize=(12,3), dpi=200, bbox_inches='tight',
                   filename=None):
    df = df.sort_values('window_size')
    fig = plt.figure(figsize=figsize, dpi=dpi)
    if title is not None:
        fig.suptitle(title)
    ax = fig.gca()
    ax.plot(df['window_size'], df['stretch_ratio'])
    ax.set_xlabel('window_size')
    ax.set_ylabel('stretch_ratio')
    if filename is not None:
        fig.savefig(filename, bbox_inches=bbox_inches)
        with open(filename + '.txt', 'w') as f:
            if comment != '':
                f.write(comment + '\n\n')
            if emb is not None:
                emb.write_config(f, comment=comment, prefix='', sep='\t')
            f.write('\nData:\n')
            f.write(df[['window_size', 'stretch_ratio']].sort_values('stretch_ratio', ascending=False).to_csv(index=False, line_terminator='\n'))
```

11/18/2020

Lattice2D\_sim (short)

In [10]:

```

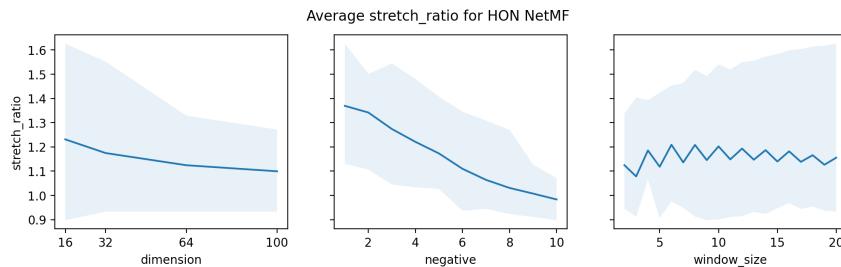
def disp_marginals(df, emb=None, title=None, figsize=(12,3), dpi=200, bbox_inches='tight',
t', filename=None, exclude_FON=True):
    "average stretch_ratio grouped by the values of each parameter"
    cols = sorted(set(df.columns) - {'stretch_ratio', 'name', 'q', 'p', 'num_walks'})
    p_vals = sorted(df['p'].unique()) if 'p' in df.columns else [None]
    if exclude_FON:
        if 'window_size' in df.columns:
            df = df[df['window_size']>1]
        if 'num_steps' in df.columns:
            df = df[df['num_steps']>1]
    fig,axs = plt.subplots(nrows=1, ncols=len(cols), sharey=True, figsize=figsize, dpi=dpi)
    if title is not None:
        fig.suptitle(title)
    for c,ax in zip(cols, axs):
        for p in p_vals:
            df_p = df if p is None else df[df['p']==p]
            avg = df_p.groupby(c)[['stretch_ratio']].mean().sort_index()
            x_vals = avg.index
            y_vals = avg.values
            label = None if p is None else 'p=' + str(p)
            ax.plot(x_vals, y_vals, label=label)
            if p is None:
                y_min = df_p.groupby(c)[['stretch_ratio']].min().sort_index().values
                y_max = df_p.groupby(c)[['stretch_ratio']].max().sort_index().values
                ax.fill_between(x_vals, y_min, y_max, alpha=0.1)
            ax.set_xlabel(c)
            if c == 'dimension' or c == 'factor':
                ax.set_xticks(x_vals)
    axs[0].set_ylabel('stretch_ratio')
    if len(p_vals) > 1:
        axs[0].legend()
    if filename is not None:
        fig.savefig(filename, bbox_inches=bbox_inches)
        with open(filename + '.txt', 'w') as f:
            if comment != '':
                f.write(comment + '\n\n')
            f.write('Params:\n')
            for c in cols:
                vals = sorted(df[c].unique())
                f.write(f'{c}\t{vals}\n')
            f.write('\nMarginals:\n')
            for c in cols:
                f.write(str(df.groupby(c)[['stretch_ratio']].agg(['min', 'mean', 'max'])) + '\n')
            if emb is not None:
                emb.write_config(f, comment=comment, prefix='', sep='\t')
        f.write('\nTop 10:\n')
        df = df.sort_values('stretch_ratio', ascending=False)
        f.write(str(df.head(10)) + '\n')
        f.write('\nBottom 10:\n')
        f.write(str(df.tail(10)) + '\n')

```

## NetMF

11/18/2020 Lattice2D\_sim (short)

```
In [11]:
# NetMF
res_NetMF = []
emb = None
# Loop over dimension and window_size first to allow NetMF to reuse its PMI calculation.
for current in iter_cross_prod(dimension=[16, 32, 64, 100], window_size=range(1,21), negative=range(1,11)):
    if (emb is None) or (emb._dimension!=current['dimension']):
        emb = HON_NetMF_EMBEDDING(gen, dimension=current['dimension'], pairwise=True)
        emb.train(window_size=current['window_size'], negative=current['negative'])
        stretch = get_stretch_ratio(emb, edge_distance = 1)
        config = dict(stretch_ratio=stretch, name='NetMF')
        config.update(current)
        res_NetMF.append(config)
res_NetMF = pd.DataFrame(res_NetMF).sort_values('stretch_ratio', ascending=False)
disp_marginals(res_NetMF, emb, title='Average stretch_ratio for HON NetMF', filename='tmp/stretch_sim_NetMF.png')
#res_NetMF.head(10)
```



In [12]:

```
res_NetMF[res_NetMF['window_size']==2].head()
```

Out[12]:

	stretch_ratio	name	dimension	window_size	negative
11	1.337149	NetMF	16	2	2
10	1.309763	NetMF	16	2	1
12	1.302841	NetMF	16	2	3
211	1.275984	NetMF	32	2	2
210	1.258308	NetMF	32	2	1

## HONEM

11/18/2020

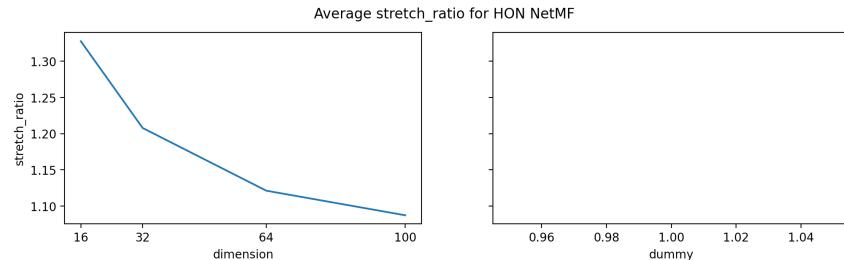
Lattice2D\_sim (short)

In [13]:

```
# HONEM
res_HONEM = []
emb = None
# disp_marginals requires at least two parameters
for current in iter_cross_prod(dimension=[16, 32, 64, 100], dummy=[1]):
    emb = HONEM_Embbedding(gen, dimension=current['dimension'])
    emb.train()
    stretch = get_stretch_ratio(emb)
    config = dict(stretch_ratio=stretch, name='HONEM')
    config.update(current)
    res_HONEM.append(config)
res_HONEM = pd.DataFrame(res_HONEM).sort_values('stretch_ratio', ascending=False)
disp_marginals(res_HONEM, emb, title='Average stretch_ratio for HON NetMF', filename='tmp/stretch_sim_HONEM.png')
res_HONEM
```

Out[13]:

	stretch_ratio	name	dimension	dummy
0	1.327574	HONEM	16	1
1	1.207816	HONEM	32	1
2	1.121258	HONEM	64	1
3	1.087279	HONEM	100	1



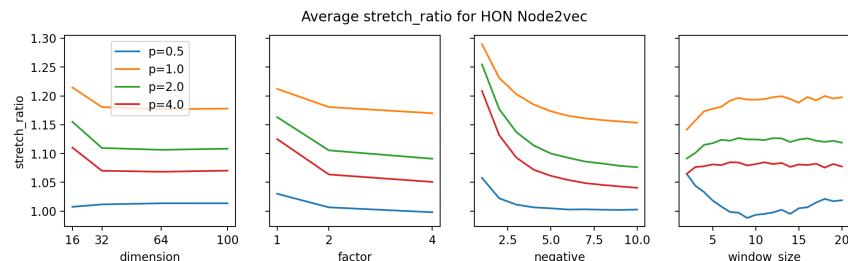
## DeepWalk & Node2vec

11/18/2020

Lattice2D\_sim (short)

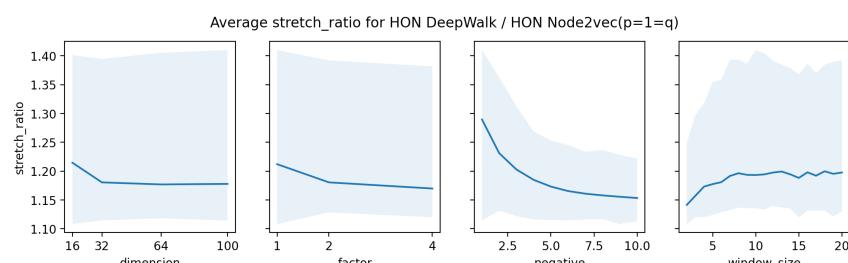
In [14]:

```
# Node2vec (7 hrs)
res_N2V = []
n2v_train_param = dict(num_iter=5, random_seed=100)
emb = None
for current in iter_cross_prod(p=[0.5,1,2,4], dimension=[16, 32, 64, 100], window_size=range(1,21), factor=[1,2,4], negative=range(1,11)):
    if (emb is None) or (emb._dimension!=current['dimension']) or (emb._gen._p!=current['p']):
        emb = HON_Node2vec_EMBEDDING(gen, current['dimension'], p=current['p'], q=1/current['p'], reuse_walks=True)
        window_size = current['window_size']
        negative = current['negative']
        factor = current['factor']
        emb.train(num_walks=1000/window_size, walk_length=factor*window_size, window_size=window_size, hs=False, negative=negative, **n2v_train_param)
        stretch = get_stretch_ratio(emb)
        config = dict(stretch_ratio=stretch, name='Node2vec', q=1/current['p'])
        config.update(current)
        res_N2V.append(config)
res_N2V = pd.DataFrame(res_N2V).sort_values('stretch_ratio', ascending=False)
disp_marginals(res_N2V, emb, title='Average stretch_ratio for HON Node2vec', filename='tmp/stretch_sim_N2V.png')
#res_N2V.head(10)
```



In [15]:

```
# DeepWalk -> see Node2vec(p=1=q)
res_DW = res_N2V[res_N2V['p']==1].drop(['p','q'], axis=1)
res_DW['name'] = 'DeepWalk'
disp_marginals(res_DW, title='Average stretch_ratio for HON DeepWalk / HON Node2vec(p=1=q)', filename='tmp/stretch_sim_DW.png')
```



11/18/2020

Lattice2D\_sim (short)

In [16]:

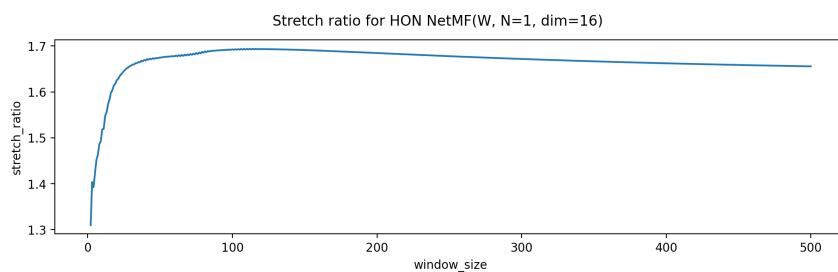
```
res_DW[res_DW['window_size']==2].head()
```

Out[16]:

	stretch_ratio	name	dimension	window_size	factor	negative
3640	1.248872	DeepWalk	64	2	2	1
2440	1.248437	DeepWalk	16	2	2	1
4240	1.248239	DeepWalk	100	2	2	1
3040	1.242888	DeepWalk	32	2	2	1
4241	1.188573	DeepWalk	100	2	2	2

In [17]:

```
# NetMF (getting the optimal window size)
if True:
    res_NetMF2 = []
    emb = None
    # Loop over dimension and window_size first to allow NetMF to reuse its PMI calculation.
    for current in iter_cross_prod(dimension=[16], window_size=range(2,501), negative=[1]):
        if (emb is None) or (emb._dimension!=current['dimension']):
            emb = HON_NetMF_EMBEDDING(gen, dimension=current['dimension'], pairwise=True)
        emb.train(window_size=current['window_size'], negative=current['negative'])
        stretch = get_stretch_ratio(emb)
        config = dict(stretch_ratio=stretch, name='NetMF')
        config.update(current)
        res_NetMF2.append(config)
    res_NetMF2 = pd.DataFrame(res_NetMF2).sort_values('stretch_ratio', ascending=False)
    disp_stretch_W(res_NetMF2, emb, title='Stretch ratio for HON NetMF(W, N=1, dim=16)'
, filename='tmp/stretch_sim_NetMF_by_W.png')
    res_NetMF2.head(10)
```



## GraRep

11/18/2020

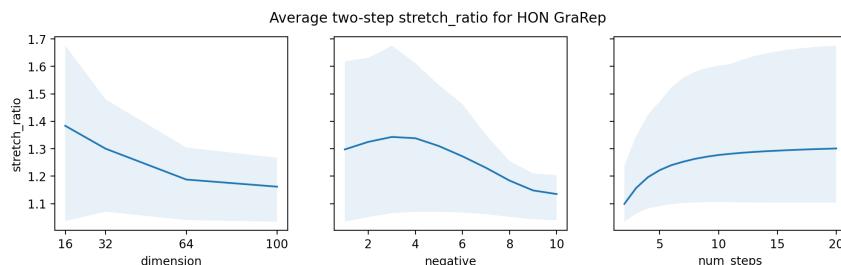
Lattice2D\_sim (short)

In [18]:

```

import warnings
# GraRep
res_GraRep = []
emb = None
# Loop over dimension and window_size first to allow NetMF to reuse its PMI calculation.
for current in iter_cross_prod(dimension=[16, 32, 64, 100], num_steps=range(1,21), negative=range(1,11)):
    if (emb is None) or (emb._dimension!=current['dimension']) or (emb._num_steps!=current['num_steps']):
        emb = HON_GraRep_EMBEDDING(gen, dimension=current['dimension'], num_steps=current['num_steps'], pairwise=True)
    with warnings.catch_warnings():
        warnings.simplefilter("ignore")
        emb.train(negative=current['negative'])
    stretch = get_stretch_ratio(emb, edge_distance=2)
    config = dict(stretch_ratio=stretch, name='GraRep')
    config.update(current)
    res_GraRep.append(config)
res_GraRep = pd.DataFrame(res_GraRep).sort_values('stretch_ratio', ascending=False)
disp_marginals(res_GraRep, emb, title='Average two-step stretch_ratio for HON GraRep',
filename='tmp/stretch_sim_GraRep.png')
#res_GraRep.head(10)

```



In [19]:

res\_GraRep[res\_GraRep['num\_steps']==2].head()

Out[19]:

	stretch_ratio	name	dimension	num_steps	negative
13	1.237977	GraRep	16	2	4
14	1.235931	GraRep	16	2	5
12	1.235618	GraRep	16	2	3
15	1.228402	GraRep	16	2	6
11	1.222282	GraRep	16	2	2

## C.8 SocioPatterns.ipynb

As discussed in § 4.2.2, we had to estimate the transition probabilities from the original data using `Python/data/models/SocioPatterns.ipynb`. Note that this requires `pathpy`<sup>2</sup>.

---

<sup>2</sup>We use version 2.2.0 installed from <https://pypi.org/project/pathpy2/>

18.11.2020

SocioPatterns

## Data preparation

[SocioPatterns.org/Datasets \(<http://www.sociopatterns.org/datasets/>\)](http://www.sociopatterns.org/datasets/)

In [ ]:

```
import os
base_dir = os.path.expanduser(r'~\Desktop\MScThesis\SocioPatterns')
work_contacts_dir = os.path.join(base_dir, r'Contacts in a workplace')
primary_school_dir = os.path.join(base_dir, r'Primary school temporal')
hospital_ward_dir = os.path.join(base_dir, r'Hospital ward')
out_dir = os.path.expanduser(r'~\Desktop\MScThesis\Python\data\models')
```

In [ ]:

```
# this workbook needs pathpy2 installed
%conda list pathpy
```

In [ ]:

```
import pathpy as pp
import numpy as np
import scipy.sparse.linalg as sla
from collections import defaultdict
```

In [ ]:

```
def iter_transitions(hon, include_subpaths=bool=True):
    "iterate over transition matrix of HigherOrderNetwork"
    mat = hon.transition_matrix(include_subpaths)
    if hon.order == 1:
        nodes = list((n,) for n in hon.nodes.keys())
    else:
        nodes = list(tuple(n.split(hon.separator)) for n in hon.nodes.keys())
    for r in range(mat.shape[0]):
        for ind in range(mat.indptr[r], mat.indptr[r+1]):
            yield nodes[mat.indices[ind]], nodes[r][-1], mat.data[ind] # yield (start,
next_node, prob)
```

18.11.2020

SocioPatterns

In [ ]:

```
# orig: Calc_q.ipynb
def calc_q(net, C, weighted=False, **kwargs):
    nodes = list(net.nodes.keys())
    mat = net.adjacency_matrix(weighted=weighted, **kwargs) # scipy.sparse.csc.csc_matrix
    mat_sum = 0
    mat_sum_by_row_C = defaultdict(float)
    mat_sum_by_col_C = defaultdict(float)
    q = 0
    for c in range(mat.shape[1]):
        c_C = C[nodes[c]] # class of current column's node
        for ind in range(mat.indptr[c], mat.indptr[c+1]):
            r = mat.indices[ind]
            v = mat.data[ind]
            # assert mat[r,c]==v
            r_C = C[nodes[r]] # class of current rows's node
            mat_sum += v
            mat_sum_by_row_C[r_C] += v
            mat_sum_by_col_C[c_C] += v
            if c_C == r_C:
                q += v
    q_exp = sum( v*mat_sum_by_col_C[c] for c,v in mat_sum_by_row_C.items() ) / (mat_sum**2)
    q = q/mat_sum - q_exp
    q_max = 1 - q_exp
    return (q, q_max)
```

```

18.11.2020           SocioPatterns

In [ ]:

def export_rules(filename, paths, config, metadata, max_order:int=1, replace_space=False,
                 include_subpaths=True):
    "Export the transition probabilities (rules) up to the estimated order"
    mog = pp.MultiOrderModel(paths, max_order)
    estimated_order = mog.estimate_order(paths)
    print('Estimated order:', estimated_order)
    stats = dict()
    with open(filename, 'w') as f:
        for order in range(1,estimated_order+1):
            hon = mog.layers[order]
            if order == 1:
                print('Exporting stationary distribution')
                probs = hon.transition_matrix(include_subpaths)
                _,ev = sla.eigs(probs, k=1, which='LM')
                ev = np.abs(ev).flatten()
                ev_sum = ev.sum()
                stat_dist = {n:v/ev_sum for n,v in zip(hon.nodes.keys(), ev)}
                for n,v in sorted(list(stat_dist.items())):
                    f.write('=> %s %r\n' % (n, v))
            print('Exporting rules for order',order)
            for start,next_node,prob in sorted(list(iter_transitions(hon, include_subpaths))):
                if replace_space:
                    next_node = next_node.replace(' ', '_')
                    start = tuple(n.replace(' ', '_') for n in start)
                for n in start:
                    f.write(n + ' ')
                f.write('=>')
                f.write(' %s %r\n' % (next_node,prob))
            # calc q, q_max
            node2cat = {n:metadata[n.split(',')[-1]] for n in hon.nodes} # last node
            is relevant for category
            q, q_max = calc_q(hon, node2cat, weighted=True, include_subpaths=include_subpaths)
            print(f'order {hon.order}: q={q}, q_max={q_max}, q/q_max={q/q_max}')
            stats.update({f'q[{hon.order}]': q, f'q_max[{hon.order}]': q_max, f'q/q_max[{hon.order}]': q/q_max })
            config_filename = os.path.splitext(filename)[0]+'.config'
            with open(config_filename,'w') as f:
                f.write('code\t%s\n' % os.path.join(os.getcwd(), 'SocioPatterns.ipynb'))
                for k,v in config.items():
                    f.write('%s\t%s\n' % (k,v))
                delta_min,delta_sec = config['delta'] // 3, config['delta'] % 3 *20
                delta_time = f'{delta_min}:{delta_sec:02d}'
                delta_text = '%d min' % delta_min if delta_sec==0 else '%d sec' % delta_sec if
                delta_min==0 else delta_time
                config_loc = dict(max_order=max_order, estimated_order=estimated_order, replace_
                _space=replace_space, include_subpaths=include_subpaths,
                                  delta_time=delta_time, delta_text=delta_text)
                for k,v in config_loc.items():
                    f.write('%s\t%s\n' % (k,v))
                for k,v in stats.items():
                    f.write('%s\t%s\n' % (k,v))

```

18.11.2020

SocioPatterns

## Contacts in workplace

delta=90 finished in less than 6 hrs; delta=180 did not finish within 24 hrs.

In [ ]:

```
filename = os.path.join(work_contacts_dir, 'tij_InVS.dat')
c_workplace=dict(source=filename)
t_workplace = pp.TemporalNetwork()
with open(filename,'r') as f:
    for line in f:
        t,i,j = line.split()
        t_workplace.add_edge(i, j, int(t)//20)
        t_workplace.add_edge(j, i, int(t)//20)
print(t_workplace.summary())
```

In [ ]:

```
# copy metadata
m_workplace = dict()
with open(os.path.join(work_contacts_dir, 'metadata_InVS13.txt'), 'r') as f:
    with open(os.path.join(out_dir, 'metadata_workplace.csv'), 'w') as g:
        #g.write(f.read())
        for line in f:
            i,Ci = line.split()
            m_workplace[i]=Ci
            g.write(line)
```

In [ ]:

```
#extraction_param = dict(delta=90)
extraction_param = dict(delta=30)
c_workplace.update(extraction_param)
p_workplace = pp.path_extraction.paths_from_temporal_network_dag(t_workplace, **extraction_param)
```

In [ ]:

```
%time export_rules(os.path.join(out_dir, 'workplace_%d.csv' % extraction_param['delta']), p_workplace, c_workplace, m_workplace, max_order=4) # estimated_order=2
```

In [ ]:

```
del p_workplace
```

## Primary school

Unable to extract paths for delta=2 (on a PC with 16GB RAM)

18.11.2020

SocioPatterns

In [ ]:

```

filename = os.path.join(primary_school_dir, 'primaryschool.csv')
c_primaryschool = dict(source=filename)
t_primaryschool = pp.TemporalNetwork()
m_primaryschool = dict() # metadata
with open(filename,'r') as f:
    for line in f:
        t,i,j,Ci,Cj = line.split()
        t_primaryschool.add_edge(i, j, int(t)//20)
        t_primaryschool.add_edge(j, i, int(t)//20)
        m_primaryschool[i]=Ci
        m_primaryschool[j]=Cj
print(t_primaryschool.summary())

```

In [ ]:

```

# export metadata
with open(os.path.join(out_dir, 'metadata_primaryschool.csv'),'w') as g:
    for i,c in sorted(list(m_primaryschool.items())):
        g.write('%s\t%s\n' % (i,c))

```

In [ ]:

```

# delta=1 took 2 hours
# delta=2 uses > 50GB RAM
extraction_param = dict(delta=1)
c_primaryschool.update(extraction_param)
p_primaryschool = pp.path_extraction.paths_from_temporal_network_dag(t_primaryschool, *extraction_param)

```

In [ ]:

```
%time export_rules(os.path.join(out_dir, 'primaryschool_%d.csv')) % extraction_param['delta'],
p_primaryschool, c_primaryschool, m_primaryschool, max_order=3)
```

In [ ]:

```
del p_primaryschool
```

## Hospital ward

delta=3 took 12 hrs to generate paths and 5 hrs to export rules

18.11.2020

SocioPatterns

In [ ]:

```
filename = os.path.join(hospital_ward_dir, 'detailed_list_of_contacts_Hospital.dat')
c_hospitalward = dict(source=filename)
t_hospitalward = pp.TemporalNetwork()
m_hospitalward = dict()
with open(filename,'r') as f:
    for line in f:
        t,i,j,Si,Sj = line.split()
        t_hospitalward.add_edge(i, j, int(t)//20)
        t_hospitalward.add_edge(j, i, int(t)//20)
        m_hospitalward[i]=Si
        m_hospitalward[j]=Sj
print(t_hospitalward.summary())
```

In [ ]:

```
# export metadata
with open(os.path.join(out_dir, 'metadata_hospital.csv'), 'w') as g:
    for i,c in sorted(list(m_hospitalward.items())):
        g.write('%s\t%s\n' % (i,c))
```

In [ ]:

```
#extraction_param = dict(delta=3)
extraction_param = dict(delta=2)
c_hospitalward.update(extraction_param)
p_hospitalward = pp.path_extraction.paths_from_temporal_network_dag(t_hospitalward, **extraction_param)
```

In [ ]:

```
%time export_rules(os.path.join(out_dir, 'hospital_%d.csv' % extraction_param['delta']), p_hospitalward, c_hospitalward, m_hospitalward, max_order=4)
```

In [ ]:

```
del p_hospitalward
```

## C.9 Classification\_sim.ipynb

This jupyter notebook calculates the tables and plots for § 4.2. We display a shorter version, which includes two flags in the first code block to control the size of the output.

18.11.2020

Classification\_sim (short)

## Init

In [1]:

```
comment = 'calculated in Classification_sim (short).ipynb'
verbose = False
tiny=True
```

In [2]:

```
%matplotlib inline
#%matplotlib notebook
import os
import pathpy as pp
import numpy as np
import pandas as pd
import math
from collections import defaultdict
import matplotlib.pyplot as plt
import seaborn as sns
```

In [3]:

```
from HigherOrderPathGenerator import HigherOrderPathGenerator, CrossValidation_HigherOrderPathGenerator
from Embedding import ABCEmbedding, HON_DeepWalk_EMBEDDING, HON_Node2vec_EMBEDDING, HON_EM_EMBEDDING, HON_NetMF_EMBEDDING, HON_GraRep_EMBEDDING, HON_Transition_Hierarchical_EMBEDDING
#from SyntheticNetworks import create_lattice_2nd_order_dynamic
from Visualizations import Visualization, EmbeddingView, Lattice2D_EMBEDDINGView
from Datasets import init_generator
```

In [4]:

```
#filename='workplace_3.csv'
#filename='workplace_15.csv'
#filename='workplace_30.csv'
#filename='workplace_60.csv'
#filename='workplace_90.csv'
#
filename='primaryschool_1.csv'
#filename='hospital_1.csv'
#filename='hospital_2.csv'
#filename='hospital_3.csv'
#filename='temporal-clusters.csv'
#filename='shuffled-temporal-clusters.csv'
```

```
18.11.2020 Classification_sim (short)

In [5]:
gen = init_generator(filename, verbose=verbose)
delta_text = gen._config['delta_text']
palette='Set1'
dim=64
metadata_name = list(gen._source_path_metadata.keys())[0]
metadata2_name = None
if filename=='primaryschool_1.csv':
    palette=sns.color_palette('Paired', n_colors = 11)
    palette[-1] = (0,0,0) # replace yellow by black
    metadata_name='Class'
    metadata2_name = 'Role'
elif filename=='temporal-clusters.csv' or filename=='shuffled-temporal-clusters.csv':
    dim=16
    palette = sns.color_palette(['#0000FF','#00CC00','#FF0000'])
filename_noext = f'tmp/{os.path.splitext(filename)[0]}-{dim}'

97223 rules read
[]
[]
```

In [6]:

```
gen_FON = gen.to_FON()
```

In [7]:

```
if not os.path.exists('tmp'):
    os.makedirs('tmp')
```

## Embeddings

In [8]:

```
embeddings_list = []
```

In [9]:

```
emb_N10 = HON_NetMF_EMBEDDING(gen, dim, pairwise=True)
%time emb_N10.train(window_size=10, negative=1, optimized=True)
emb_N10._id = 'HON NetMF(W=10)'
embeddings_list.append(emb_N10)
```

Wall time: 8min 8s

In [10]:

```
emb_N10_FON = HON_NetMF_EMBEDDING(gen_FON, dim, pairwise=True)
%time emb_N10_FON.train(window_size=10, negative=1, optimized=True)
emb_N10_FON._id = 'FON NetMF(W=10)'
embeddings_list.append(emb_N10_FON)
```

Wall time: 1min 24s

18.11.2020

Classification\_sim (short)

In [11]:

```
emb_N2 = HON_NetMF_EMBEDDING(gen, dim, pairwise=True)
%time emb_N2.train(window_size=2, negative=1, optimized=True)
emb_N2._id = 'HON NetMF(W=2)'
embeddings_list.append(emb_N2)
```

Wall time: 1.25 s

In [12]:

```
emb_N2_FON = HON_NetMF_EMBEDDING(gen_FON, dim, pairwise=True)
%time emb_N2_FON.train(window_size=2, negative=1, optimized=True)
emb_N2_FON._id = 'FON NetMF(W=2)'
embeddings_list.append(emb_N2_FON)
```

Wall time: 3.64 s

In [13]:

```
emb_G4 = HON_GraRep_EMBEDDING(gen, dim//4, num_steps=4, pairwise=True, neg_stationary=False)
%time emb_G4.train(negative=1, normalize=False) # LogisticRegression performs poorly if
normalize=True
emb_G4._id = 'HON GraRep(S=4)'
embeddings_list.append(emb_G4)
```

Wall time: 1min 17s

In [14]:

```
emb_G4_FON = HON_GraRep_EMBEDDING(gen_FON, dim//4, num_steps=4, pairwise=True, neg_stat
ionary=False)
%time emb_G4_FON.train(negative=1, normalize=False) # LogisticRegression performs poorl
y if normalize=True
emb_G4_FON._id = 'FON GraRep(S=4)'
embeddings_list.append(emb_G4_FON)
```

Wall time: 23.3 s

In [15]:

```
deepwalk_para = dict(num_walks=100, hs=False, negative=1, random_seed=1)
emb_D10 = HON_DeepWalk_EMBEDDING(gen, dim)
%time emb_D10.train(window_size=10, **deepwalk_para)
emb_D10._id = 'HON DeepWalk'
embeddings_list.append(emb_D10)
```

Wall time: 32.1 s

In [16]:

```
emb_D10_FON = HON_DeepWalk_EMBEDDING(gen_FON, dim)
%time emb_D10_FON.train(window_size=10, **deepwalk_para)
emb_D10_FON._id = 'FON DeepWalk'
embeddings_list.append(emb_D10_FON)
```

Wall time: 29.3 s

18.11.2020

Classification\_sim (short)

In [17]:

```
node2vec_para = dict(num_walks=100, hs=False, negative=1, random_seed=1)
emb_N2V_10a = HON_Node2vec_Embedding(gen, dim, p=0.5, q=2)
%time emb_N2V_10a.train(window_size=10, **node2vec_para)
emb_N2V_10a._id = 'HON Node2vec(p=0.5)'
embeddings_list.append(emb_N2V_10a)
```

Wall time: 32.1 s

In [18]:

```
emb_N2V_10a_FON = HON_Node2vec_Embedding(gen_FON, dim, p=0.5, q=2)
%time emb_N2V_10a_FON.train(window_size=10, **node2vec_para)
emb_N2V_10a_FON._id = 'FON Node2vec(p=0.5)'
embeddings_list.append(emb_N2V_10a_FON)
```

Wall time: 32.8 s

In [19]:

```
emb_N2V_10b = HON_Node2vec_Embedding(gen, dim, p=2, q=0.5)
%time emb_N2V_10b.train(window_size=10, **node2vec_para)
emb_N2V_10b._id = 'HON Node2vec(p=2)'
embeddings_list.append(emb_N2V_10b)
```

Wall time: 33.6 s

In [20]:

```
emb_N2V_10b_FON = HON_Node2vec_Embedding(gen_FON, dim, p=2, q=0.5)
%time emb_N2V_10b_FON.train(window_size=10, **node2vec_para)
emb_N2V_10b_FON._id = 'FON Node2vec(p=2)'
embeddings_list.append(emb_N2V_10b_FON)
```

Wall time: 33.6 s

In [21]:

```
emb_H = HONEM_Embedding(gen, dim)
%time emb_H.train()
emb_H._id = 'HONEM'
embeddings_list.append(emb_H)
```

Wall time: 358 ms

In [22]:

```
embeddings = { e._id:e for e in embeddings_list}
```

## Visualization

In [23]:

```
embedding_views = [ EmbeddingView(emb, use_source=True) for emb in embeddings_list ]
```

18.11.2020

Classification\_sim (short)

In [24]:

```
def plot_ev(ev, filename=None, random_state=1, comment=comment):
    vis = ev.visualize_TSNE(random_state=random_state, title=ev._emb._id + ', TSNE')
    metadata_classes = sorted(list(ev[metadata_name].unique()))
    plot_args = dict(hue=metadata_name, palette=palette, hue_order=metadata_classes)
    if metadata2_name is not None:
        plot_args['style'] = metadata2_name
    if tiny:
        plot_args['figsize']=(4.5,3)
        plot_args['dpi']=75
    vis.plot1(**plot_args)
    if filename is not None:
        vis.save_describe(filename, comment)
if verbose:
    for ev in embedding_views:
        plot_ev(ev, random_state=1, filename=None)
```

## Classification

In [25]:

```
from sklearn.multiclass import OneVsRestClassifier
from sklearn.svm import SVC
from sklearn.linear_model import LogisticRegression
from sklearn.ensemble import RandomForestClassifier
```

In [26]:

```
models = {
    'LR': LogisticRegression(random_state=1),
    'SVC': SVC(random_state=1, probability=True),
    'RF': RandomForestClassifier(random_state=1)
}
```

18.11.2020

Classification\_sim (short)

In [27]:

```
from sklearn.model_selection import KFold
from sklearn.metrics import f1_score
def get_stats(ev, model, n_splits=5, n_repeat=10, verbose=True, target_name=metadata_name):
    stats = defaultdict(list)
    X = ev.embedding.values
    y = ev[target_name].values
    #y_classes = sorted(np.unique(x))
    for random_state in range(n_repeat):
        if verbose: print('.', end='')
        y_true = []
        y_pred = []
        scores = []
        kf = KFold(n_splits, shuffle=True, random_state=random_state)
        for train_index, test_index in kf.split(X):
            X_train, X_test = X[train_index], X[test_index]
            y_train, y_test = y[train_index], y[test_index]
            model.fit(X_train, y_train)
            scores.append(model.score(X_test, y_test))
            y_true.extend(y_test)
            y_pred.extend(model.predict(X_test))
        stats['score'].append(np.mean(scores))
        stats['f1-micro'].append(f1_score(y_true, y_pred, average='micro'))
        stats['f1-macro'].append(f1_score(y_true, y_pred, average='macro'))
    if verbose: print()
    res = {k: np.mean(v) for k,v in stats.items()}
    res.update({k+' (min)': np.min(v) for k,v in stats.items()})
    res.update({k+' (max)': np.max(v) for k,v in stats.items()})
    res.update({k+' (std)': np.std(v) for k,v in stats.items()})
    return res
```

In [28]:

```
stats_list = []
for ev in embedding_views:
    if verbose: print('Embedding', ev._emb._id)
    for model_name, model in models.items():
        if verbose: print('- Model', model_name, end='')
        stats = dict(embedding=ev._emb._id, model=model_name)
        stats.update(get_stats(ev, model, verbose=verbose))
        stats_list.append(stats)
stats_df = pd.DataFrame(stats_list).set_index(['model', 'embedding'])
if verbose:
    stats_df
```

In [29]:

```
stats4tex = stats_df[['f1-macro', 'f1-micro']].reset_index().set_index(['embedding', 'model']).unstack()
if verbose:
    stats4tex
```

18.11.2020

Classification\_sim (short)

In [30]:

```
def summarize(col='f1-macro'):
    df = stats_df[[col]].sort_values(col, ascending=False).copy().reset_index(drop=False)
    df['group'] = df['embedding'].map(lambda x:x.split(' ')[0])
    df['rank'] = df.index+1
    df['score'] = col
    df['value'] = df[col]
    df['dataset'] = gen._id
    return df.groupby(['dataset', 'score', 'group'])[['value', 'rank']].mean()
summarize().append(summarize('f1-micro'))
```

Out[30]:

			value	rank
dataset	score	group		
primaryschool_1	f1-macro	FON	0.769721	25.611111
		HON	0.825248	11.388889
		HONEM	0.412202	38.000000
	f1-micro	FON	0.815680	26.166667
		HON	0.871786	10.833333
		HONEM	0.452342	38.000000

18.11.2020

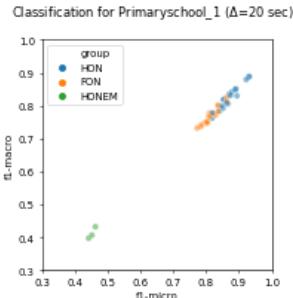
Classification\_sim (short)

In [31]:

```

with open(filename_noext + '.tex', 'w') as f:
    gen.write_config(f, comment=comment, prefix='%', sep='=')
    f.write('%\n% Embeddings:\n')
    for emb in embeddings.values():
        f.write(f'% {emb._id}={repr(emb.config)}\n')
    f.write('\n')
    stats4tex.to_latex(f, float_format=".3f", index_names=False)
# export summary
summarize('f1-macro').append(summarize('f1-micro')).to_csv(filename_noext + '.csv', sep='\t', encoding='utf-16')
# export plot
df = stats_df.copy().reset_index(drop=False)
df['group'] = df['embedding'].map(lambda x:x.split(' ')[0])
fig = plt.figure(figsize=(4, 4.3), dpi=50 if tiny else 400)
fig.suptitle('Classification for ' + gen_id.title() + ' ($\Delta$=' + delta_text + ')')
ax = sns.scatterplot(x='f1-micro', y='f1-macro', hue='group', data=df, ax=fig.gca(), alpha=0.5)
f1_min = min(df['f1-micro'].min(), df['f1-macro'].min())
f1_max = max(df['f1-micro'].max(), df['f1-macro'].max())
f1_min = np.floor(f1_min*10)/10
f1_max = np.ceil(f1_max*10)/10
ax.set_xlim(f1_min, f1_max)
ax.set_ylim(f1_min, f1_max)
ax.set_aspect(1)
ax.figure.savefig(filename_noext + '.png')

```



In [32]:

```

sorted_f1_micro = stats_df['f1-micro'].sort_values(ascending=False)
sorted_f1_micro.to_csv(filename_noext + '_f1-micro.csv', sep='\t', encoding='utf-16')
if verbose:
    sorted_f1_micro

```

In [33]:

```

sorted_f1_macro = stats_df['f1-macro'].sort_values(ascending=False)
sorted_f1_macro.to_csv(filename_noext + '_f1-macro.csv', sep='\t', encoding='utf-16')
if verbose:
    sorted_f1_macro

```

18.11.2020

Classification\_sim (short)

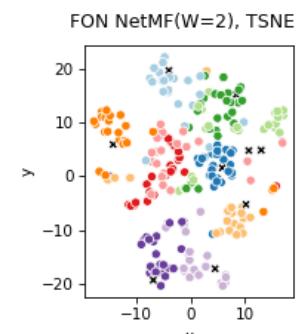
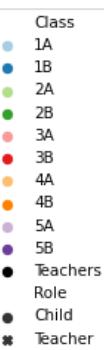
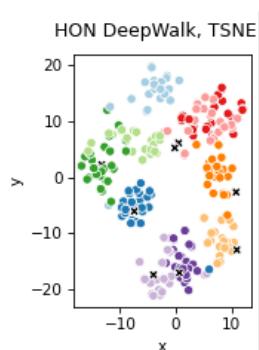
In [34]:

```
sorted_f1 = (stats_df['f1-micro'] + stats_df['f1-macro']).sort_values(ascending=False)
best_HON =sorted_f1[sorted_f1.index.map(lambda x:x[1].startswith('HON '))].index[0]
best_FON =sorted_f1[sorted_f1.index.map(lambda x:x[1].startswith('FON '))].index[0]
print('best HON model and embedding', best_HON)
print('best FON model and embedding', best_FON)
```

best HON model and embedding ('LR', 'HON DeepWalk')  
 best FON model and embedding ('LR', 'FON NetMF(W=2)')

In [35]:

```
# export tsne visualization of best embeddings
plot_ev(EmbeddingView(embeddings[best_HON[1]], use_source=True), random_state=2, filename=f'{filename_noext}_tsne_hon.png')
plot_ev(EmbeddingView(embeddings[best_FON[1]], use_source=True), random_state=1, filename=f'{filename_noext}_tsne_fon.png')
```



## Confusion matrix

18.11.2020

Classification\_sim (short)

In [36]:

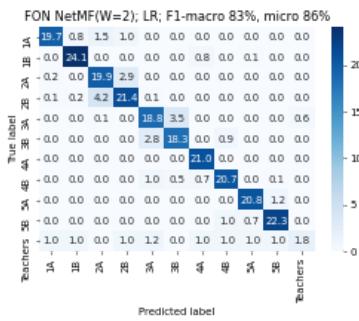
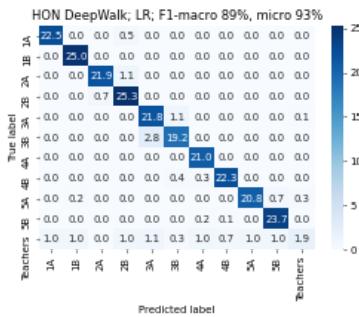
```
from sklearn.metrics import confusion_matrix
from sklearn.base import BaseEstimator
def create_confusion_matrix(ev: EmbeddingView, model: BaseEstimator, n_splits:int=5, n_
repeat=10, plot:bool=False, title=None, target_name=metadata_name):
    if verbose:
        print(ev._emb._id, model)
    X = ev.embedding.values
    y = ev[target_name].values
    classes = sorted(np.unique(y))
    y_true, y_pred=[], []
    for random_state in range(n_repeat):
        kf = KFold(n_splits, shuffle=True, random_state=random_state)
        for train_index, test_index in kf.split(X):
            X_train, X_test = X[train_index], X[test_index]
            y_train, y_test = y[train_index], y[test_index]
            model.fit(X_train, y_train)
            y_true.extend(y_test)
            y_pred.extend(model.predict(X_test))
    cm = confusion_matrix(y_true, y_pred, labels=classes) / n_repeat
    if plot:
        plt.figure(figsize=(6,4) if tiny else(7,5), dpi=50 if tiny else 100)
        if title is not None:
            plt.title(title)
        sns.heatmap(cm, annot=True, fmt='.'1f', cmap=plt.cm.Blues, xticklabels=classes,
        yticklabels=classes)
        plt.xlabel('Predicted label')
        plt.ylabel('True label')
    return pd.DataFrame(cm, index=pd.MultiIndex.from_product([[ 'True label'], classes
]), columns=pd.MultiIndex.from_product([[ 'Predicted label'], classe
s]))
```

18.11.2020

Classification\_sim (short)

In [37]:

```
def export_confusion_matrix(best, filename=None, comment=comment):
    title = f'{best[1]}; {best[0]}; F1-macro {sorted_f1_macro[best]:.0%}, micro {sorted_f1_micro[best]:.0%}'
    emb = embeddings[best[1]]
    cf_para = dict(model=models[best[0]], n_splits=5, n_repeat=10)
    data = create_confusion_matrix(EmbeddingView(emb, use_source=True), plot=True, title=title, **cf_para)
    if filename is not None:
        plt.gcf().savefig(filename, bbox_inches='tight')
        with open(filename + '.txt', 'w') as f:
            comment_loc = (comment + ('' if comment == '' else '\n\n')) + f'F1-macro\t{sorted_f1_macro[best]:f}\nF1-micro\t{sorted_f1_micro[best]:f}\n'
            emb.write_config(f, comment=comment_loc, prefix='', sep='\t')
            f.write('\nConfusion matrix:\n')
            for k, v in cf_para.items():
                f.write(f'{k}\t{v}\n')
            f.write('\n\nData:\n')
            f.write(data.to_csv(line_terminator='\n'))
    export_confusion_matrix(best_HON, filename=f'{filename_noext}_confusion_hon.png')
    export_confusion_matrix(best_FON, filename=f'{filename_noext}_confusion_fon.png')
```



## AUC / ROC

18.11.2020

Classification\_sim (short)

In [38]:

```
# https://stackoverflow.com/questions/45332410/sklearn-roc-for-multiclass-classification
n
from sklearn.metrics import roc_curve, auc
def create_auc_roc_old(ev: EmbeddingView, model: BaseEstimator, n_splits:int=5, n_repeat=10, title=None, plot:bool=False, target_name=metadata_name):
    X = ev.embedding.values
    #y = ev.target.values
    y = ev[target_name].values
    y_true, y_score = [], []
    for random_state in range(n_repeat):
        kf = KFold(n_splits, shuffle=True, random_state=random_state)
        for train_index, test_index in kf.split(X):
            X_train, X_test = X[train_index], X[test_index]
            y_train, y_test = y[train_index], y[test_index]
            model.fit(X_train, y_train)
            y_true.extend(y_test)
            y_score.append(model.decision_function(X_test)) # corresponding labels are
    in model.classes_
    y_score = np.concatenate(y_score, axis=0) # rbind...
    y_true = np.array(y_true)

    fpr, tpr, roc_auc = dict(), dict(), dict()
    for i,c in enumerate(model.classes_):
        # roc_curve for y_true == c vs. rest
        fpr_, tpr_, _ = roc_curve(y_true == c, y_score[:,i])
        fpr[c] = fpr_
        tpr[c] = tpr_
        roc_auc[c] = auc(fpr_, tpr_)
    if plot:
        colors = sns.color_palette(palette, len(model.classes_))
        plt.figure(dpi=200)
        for c in model.classes_: # or sorted(np.unique(y))
            plt.plot(fpr[c], tpr[c], label='%s (area = %0.2f)' % (c,roc_auc[c]), color=colors.pop(0))
        plt.plot([0, 1], [0, 1], 'k--')
        plt.xlim([0.0, 1.0])
        plt.ylim([0.0, 1.05])
        plt.xlabel('False Positive Rate')
        plt.ylabel('True Positive Rate')
        if title is None:
            plt.title('Receiver operating characteristic (1 vs rest) - %d times %d-fold CV' % (n_repeat, n_splits))
        else:
            plt.title(title)
        plt.legend(loc="lower right")
        #plt.show()

    return pd.Series(roc_auc, index=model.classes_, name='auc')#, fpr, tpr
```

18.11.2020

Classification\_sim (short)

In [39]:

```

from sklearn.metrics import roc_curve, auc, roc_auc_score
def create_auc_roc(ev: EmbeddingView, model: BaseEstimator, n_splits:int=5, n_repeat=10
, title=None, plot:bool=False, target_name=metadata_name, check=True):
    X = ev.embedding.values
    #y = ev.target.values
    y = ev[target_name].values
    y_true, y_prob = [], []
    for random_state in range(n_repeat):
        kf = KFold(n_splits, shuffle=True, random_state=random_state)
        for train_index, test_index in kf.split(X):
            X_train, X_test = X[train_index], X[test_index]
            y_train, y_test = y[train_index], y[test_index]
            model.fit(X_train, y_train)
            y_true.extend(y_test)
            y_prob.append(model.predict_proba(X_test)) # corresponding labels are in model.classes_
    y_prob = np.concatenate(y_prob, axis=0) # rbind(...)
    y_true = np.array(y_true)

    fpr, tpr, roc_auc = dict(), dict(), dict()
    for i,c in enumerate(model.classes_):
        # roc_curve for y_true == c vs. rest
        fpr_, tpr_, _ = roc_curve(y_true == c, y_prob[:,i])
        fpr[c] = fpr_
        tpr[c] = tpr_
        roc_auc[c] = auc(fpr_, tpr_)
    if check:
        # "Area under ROC for the multiclass problem" according to sklearn documentation
        # https://scikit-learn.org/stable/auto_examples/model_selection/plot_roc.html
        macro_roc_auc_ovr = roc_auc_score(y_true, y_prob, multi_class='ovr', average='macro')
        # getting the 'macro'-average of roc_auc, which is supposed to calculate the same
        auc_mean = np.mean(list(roc_auc.values()))
        if abs(auc_mean - macro_roc_auc_ovr) > 1e-16:
            print('The individual AUCs are inconsistent with roc_auc_score(ovr,macro)=%f'%macro_roc_auc_ovr)
    if plot:
        colors = sns.color_palette(palette, len(model.classes_))
        plt.figure(dpi=50 if tiny else 200)
        for c in model.classes_: # or sorted(np.unique(y))
            plt.plot(fpr[c], tpr[c], label='%s (area = %.3f)' % (c,roc_auc[c]), color=colors.pop(0))
        plt.plot([0, 1], [0, 1], 'k--')
        plt.xlim([0.0, 1.0])
        plt.ylim([0.0, 1.05])
        plt.xlabel('False Positive Rate')
        plt.ylabel('True Positive Rate')
        if title is None:
            plt.title('Receiver operating characteristic (1 vs rest) - %d times %d-fold CV' % (n_repeat, n_splits))
        else:
            plt.title(title)
        plt.legend(loc="lower right")
        #plt.show()

    return pd.Series(roc_auc, index=model.classes_, name='auc')#, fpr, tpr

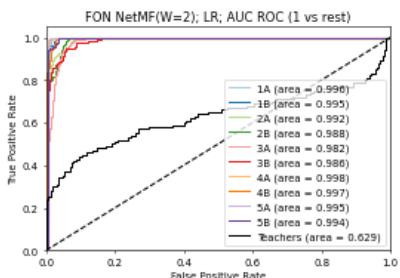
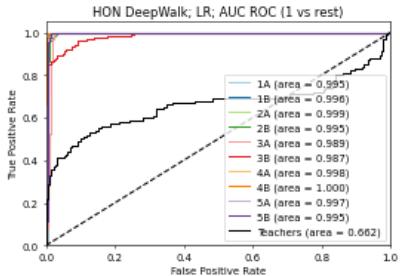
```

18.11.2020

Classification\_sim (short)

In [40]:

```
def export_auc_roc(best, filename=None, comment='comment'):
    title = f'{best[1]}; {best[0]}; AUC ROC (1 vs rest)'
    emb = embeddings[best[1]]
    auc_para = dict(model=models[best[0]], n_splits=5, n_repeat=10)
    data = create_auc_roc(EmbeddingView(emb, use_source=True), title=title, plot=True,
    **auc_para)
    if filename is not None:
        plt.gcf().savefig(filename, bbox_inches='tight')
        with open(filename + '.txt', 'w') as f:
            comment_loc = (comment + ('' if comment == '' else '\n\n')) +
                f'F1-macro\t{sorted_f1_macro[best]:f}\nF1-micro\t{sorted_f1_micro[bes
t]:f}\n'
            + f'macro_roc_auc_ovr\t{data.mean():f})'
        emb.write_config(f, comment=comment_loc, prefix='', sep='\t')
        f.write('\nAUC ROC:\n')
        for k,v in auc_para.items():
            f.write('%s\t%s\n' % (k,v))
        f.write('\n\nData:\n')
        f.write(data.to_csv(line_terminator='\n'))
    export_auc_roc(best_HON, filename=f'{filename_noext}_auc-roc_hon.png')
    export_auc_roc(best_FON, filename=f'{filename_noext}_auc-roc_fon.png')
```



## Entropy

In [41]:

```
entropies = dict()
for key in gen.rule_keys:
    tmp = 0
    for _, _, prob in gen.transition_probs(start=key):
        tmp += prob * math.log2(prob)
    entropies[key] = -tmp
```

18.11.2020

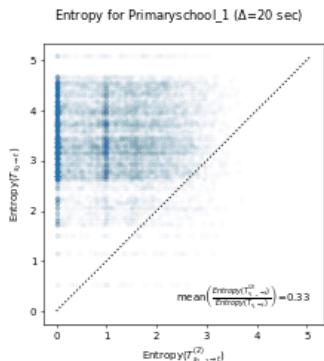
Classification\_sim (short)

In [42]:

```
entropy_data = [
    dict(start=str(start), entropy=entropy,
         parent_entropy=entropies[start[:-1]],
         start_len=len(start))
    )
for start,entropy in entropies.items() if len(start)>0
]
entropy_data = pd.DataFrame(entropy_data).set_index('start')
entropy_data['entropy_ratio'] = entropy_data['entropy']/entropy_data['parent_entropy']
```

In [43]:

```
entropy_data2 = entropy_data[entropy_data['start_len']==2]
avg_ratio_desc = 'mean$\left(\frac{\text{Entropy}(T^{(2)}_{s_1..2})}{\text{Entropy}(T_{s_2})}\right) = %.2f' % entropy_data["entropy_ratio"].mean()
fig = plt.figure(figsize=(5,5), dpi=50 if tiny else 400)
alpha = 0.01 if filename == 'primaryschool_1.csv' else 0.1
ax = entropy_data2.plot.scatter('entropy', 'parent_entropy', alpha=alpha, ax=plt.gca())
max_entropy = entropy_data2['parent_entropy'].max()
ax.set_aspect(1)
ax.set_xlabel('Entropy($T^{(2)}_{s_1..2}$)')
ax.set_ylabel('Entropy($T_{s_2}$)')
ax.plot([0, max_entropy], [0, max_entropy], 'k:')
ax.figure.suptitle('Entropy for ' + gen._id.title() + ' ($\Delta=' + delta_text + ')')
ax.annotate(avg_ratio_desc, (max_entropy,0), ha='right', va='bottom')
fig.savefig(f'tmp/{os.path.splitext(filename)[0]}_entropy.png')
```



In [44]:

```
entropy_data['entropy'].mean()
```

Out[44]:

1.1137836903958092

In [45]:

```
entropy_data['parent_entropy'].mean()
```

Out[45]:

3.501841638726696

## **C.10 Plot\_ExpClass.ipynb**

This jupyter notebook calculates figure 4.11.

17.11.2020

Plots\_ExClass

In [1]:

```
comment = 'calculated in Plots_ExClass.ipynb'
```

In [2]:

```
%matplotlib inline
#%matplotlib notebook
import os
import pandas as pd
from collections import defaultdict
import matplotlib.pyplot as plt
import seaborn as sns
```

In [3]:

```
from Datasets import init_generator # returns HigherOrderPathGenerator
```

In [4]:

```
gen_school = init_generator('primaryschool_1.csv')
```

```
97223 rules read
[]
[]
```

In [5]:

```
if not os.path.exists('tmp'):
    os.makedirs('tmp')
```

## Check child-teacher-child interactions

In [6]:

```
def is_teacher(node):
    return gen_school._target_node_metadata['Role'][node] == 'Teacher'
def is_child(node):
    return gen_school._target_node_metadata['Role'][node] == 'Child'
def get_class(node):
    return gen_school._target_node_metadata['Class'][node]
```

17.11.2020

Plots\_ExpClass

In [7]:

```

####ctc_paths = list() # child-teacher-child paths
def get_class_class_probs_via_teacher(HON=True, title=None, filename=None, transition_p
robs=False, figsize=(7,5), dpi=200, bbox_inches='tight', **kwargs):
    """
    returns a dictionary { (c1,c2) : prob}
    where (c1,c2) is a pair of classes
    and prob = P(class(n1) = c1, class(n3)=c2 | is_child(n1), is_teacher(n2), is_child
(n3))
    where (n1,n2,n3) is a random path of length 3 starting with n1 ~ stationary
    """
    gen = gen_school if HON else gen_school.to_FON()
    total_prob = 0
    cp = defaultdict(float) # prob of class
    ccp = defaultdict(float) # prob of pair of classes
    for (n1,n2,n3),prob in gen.path_probs(start=(), num_steps=3):
        if is_child(n1) and is_teacher(n2) and is_child(n3):
            c1 = get_class(n1)
            c3 = get_class(n3)
            total_prob += prob
            cp[c1] += prob
            ccp[(c1,c3)] += prob
    print('total probability', total_prob)
    if transition_probs:
        ccp = { cc:prob/cp[cc[0]] for cc,prob in ccp.items()} # transition probability
c1->c2
    else:
        ccp = { cc:prob/total_prob for cc,prob in ccp.items()} # probability of pairs
    res = pd.Series(ccp).unstack()
    if filename is not None:
        fig = plt.figure(figsize=figsize, dpi=dpi)
        if title is not None:
            fig.suptitle(title)
        classes = sorted(cp.keys())
        ax = fig.gca()
        sns.heatmap(res, annot=True, fmt=' .2f', cmap=plt.cm.Blues, xticklabels=classes,
yticklabels=classes, ax=ax)
        if transition_probs:
            ax.set_xlabel('target')
            ax.set_ylabel('source')
        else:
            ax.set_xlabel('class of $child_2$')
            ax.set_ylabel('class of $child_1$')
        fig.savefig(filename, bbox_inches=bbox_inches, **kwargs)
        with open(filename + '.txt','w') as f:
            gen.write_config(f, comment=comment, prefix='', sep='\t')
            f.write('\n\nProbability of a child->teacher->child path %f\n' % total_prob
)
        if transition_probs:
            f.write('\nTransition probabilities:\n')
        else:
            f.write('\nPair probabilities:\n')
            f.write(res.to_csv(line_terminator='\n'))
    return res

```

17.11.2020

Plots\_ExpClass

In [ ]:

```
title_HON = 'Cond. probability of Child-Teacher-Child interaction by class (HON)'
c2c_HON = get_class_class_probs_via_teacher(True, title_HON, filename='tmp/school_hon_i
nteraction_via_teacher.png')
```

In [ ]:

```
title_FON = 'Cond. probability of Child-Teacher-Child interaction by class (FON)'
c2c_FON = get_class_class_probs_via_teacher(False, title_FON, filename='tmp/school_fon_
interaction_via_teacher.png')
```

Get frequencies by class from stationary distribution

In [10]:

```
class_probs = defaultdict(float)
for _, n, prob in gen_school.transition_probs(start=()): # stationary distribution
    class_probs[get_class(n)] += prob
class_probs
```

Out[10]:

```
defaultdict(float,
    {'5B': 0.08779785547204058,
     '5A': 0.04456183510387105,
     '4A': 0.04474650478145715,
     'Teachers': 0.004609362345706892,
     '3B': 0.15727510236331413,
     '4B': 0.05062492826777912,
     '2A': 0.11235273025299394,
     '1B': 0.2166310296376714,
     '2B': 0.12849865639303715,
     '1A': 0.04554663582826103,
     '3A': 0.10735535955386742})
```

## C.11 ProbabilityPrediction.ipynb

This jupyter notebook calculates the simulation and plots for § B.1.

11/18/2020

ProbabilityPrediction\_sim

## Init

```
In [1]: verbose=False
```

```
In [2]: %matplotlib inline
#%matplotlib notebook
import os
import pathpy as pp
import numpy as np
import pandas as pd
import math
from collections import defaultdict
import matplotlib.pyplot as plt
#from tqdm import tqdm
from tqdm.notebook import tqdm
```

```
In [3]: from HigherOrderPathGenerator import HigherOrderPathGenerator, CrossValidation
HigherOrderPathGenerator, ABCHigherOrderPathGenerator
from Embedding import ABCEmbedding, HON_DeepWalk_EMBEDDING, HONEM_EMBEDDING, H
ON_NetMF_EMBEDDING, HON_GraRep_EMBEDDING, HON_Transition_Hierarchical_EMBEDDING
from Visualizations import Visualization, EmbeddingView, Lattice2D_EMBEDDINGVi
ew
from Datasets import init_generator
```

```
In [4]: gen_HON = init_generator('primaryschool_1.csv')
#gen_HON = init_generator('workplace_30.csv')
#gen_HON = init_generator('hospital_1.csv')
```

```
97223 rules read
[]
[]
```

```
In [5]: gen_FON = gen_HON.to_FON(gen_HON._id + ' (FON)')
```

## Embeddings

11/18/2020 ProbabilityPrediction\_sim

```
In [6]: class embedding_builder(object):
    """Instantiates an embedding and trains it."""
    def __init__(self, name:str, init_para=None, **train_para):
        self._name = name
        self._init_para = dict() if init_para is None else init_para
        self._train_para = train_para

    def build(self, gen: ABCHigherOrderPathGenerator, dimension:int = 128):
        if self._name == 'NetMF':
            emb = HON_NetMF_EMBEDDING(gen, dimension, **self._init_para)
            emb.train(**self._train_para)
        elif self._name == 'GraRep':
            emb = HON_GraRep_EMBEDDING(gen, dimension, **self._init_para)
            emb.train(**self._train_para)
        elif self._name == 'Experiment':
            #emb = HON_Transition_Hierarchical_EMBEDDING(gen, dimension, **self._init_para)
            emb = HON_CV_Transition_Hierarchical_EMBEDDING(gen, dimension, **self._init_para)
            emb.train(**self._train_para)
        elif self._name == 'DeepWalk':
            emb = HON_DeepWalk_EMBEDDING(gen, dimension)
        elif self._name == 'HONEM':
            emb = HONEM_EMBEDDING(gen, dimension)
            emb.train()
        else:
            assert False, 'Invalid name %s' % name
        return emb

    @staticmethod
    def get_true_probs(gen: ABCHigherOrderPathGenerator, excluded_edges):
        res = defaultdict(dict) # use nested dictionaries, because decode is expensive
        for source in gen.source_paths_len1:
            for _, next_node, prob in gen.transition_probs(source):
                if (source[-1], next_node) in excluded_edges:
                    res[source][next_node] = prob
        return res

    def evaluate_pairs(self, true_probs, gen_build: ABCHigherOrderPathGenerator,
                      dimension: int = 128, **kwargs):
        emb = self.build(gen_build, dimension)
        node2str = emb.key2str if emb._symmetric else emb.node2str
        source2str = emb.key2str if emb._symmetric else emb.path2str
        res = []
        for source, probs in true_probs.items():
            source_str = source2str(source)
            predicted_probs = emb.decode_path(source, **kwargs)
            for next_node, prob in probs.items():
                next_node_str = node2str(next_node)
                res.append(dict(source=source_str, target=next_node_str, true_prob=prob, pred_prob=predicted_probs[next_node_str]))
        return res

    def evaluate(self, true_probs, gen_build: ABCHigherOrderPathGenerator, dimension:int = 128, **kwargs):

```

localhost:8888/nbconvert/html/ProbabilityPrediction\_sim.ipynb?download=false 2/8

11/18/2020

ProbabilityPrediction\_sim

```

emb = self.build(gen_build, dimension)
node2str = emb.key2str if emb._symmetric else emb.node2str
SSE = 0
SSElog = 0
min_prob=1e-10
for source, probs in true_probs.items():
    predicted_probs = emb.decode_path(source, **kwargs)
    for next_node, prob in probs.items():
        SSE += (prob - predicted_probs[node2str(next_node)])**2
        SSElog += (math.log2(max(prob,min_prob)) - math.log2(max(predicted_probs[node2str(next_node)],min_prob)))**2
return dict(SSE=SSE,SSElog=SSElog) # include time? include SSE of log probabilities?

```

In [7]:

```

class HON_CV_Transition_Hierarchical_Embedding(HON_Transition_Hierarchical_Emb
edding):
    """In a cross-validation setting, some edges should be excluded from train
    ing the embedding.
    While this is impossible for methods based on (unweighted) matrix factoriz
    ation using SVD,
    this is simple for methods based on weighted matrix factorization using SG
    D.

    The trick is to skip the update step if it happens to coincide with one of
    the excluded edges.
    """
    def __init__(self, gen, dimension: int=128, node_hierarchy='calc', seed=None,
                 neg_stationary:bool=True):
        super().__init__(gen, dimension, node_hierarchy, seed, neg_stationary)
        if type(gen) == CrossValidation_HigherOrderPathGenerator:
            self._excluded_edges = { (iu,iv) for u,iu in self._source_paths.it
ems() for v,iv in self._target_nodes.items()
                           if (u[-1],v) in gen._excluded_edges }
        else:
            self._excluded_edges = set() # warn?

    def _update(self, iu, iv, label, learning_rate):
        if (iu,iv) in self._excluded_edges:
            return # skip training (both positive and negative samples) for th
e excluded edges
        super()._update(iu, iv, label, learning_rate)

```

In [8]:

```
dimensions = [16, 32, 64, 128]
```

In [9]:

```
decode_args = dict(use_neighborhood=True, no_self_loops=True, normalize=True,
step=1) # decode ignores additional parameters
```

11/18/2020

ProbabilityPrediction\_sim

```
In [10]: netmf_pairs = dict(pairwise=True)
grarep_para = dict(num_steps=2, pairwise=True)
#deepwalk10_para = dict(num_walks=100, walk_length=100, window_size=10, negative=5)
embedding_builders = {
    'NetMF(W1,N1)': embedding_builder('NetMF', init_para=netmf_pairs, window_size=1, negative=1),
    'NetMF(W1,N5)': embedding_builder('NetMF', init_para=netmf_pairs, window_size=1, negative=5),
    'NetMF(W2,N1)': embedding_builder('NetMF', init_para=netmf_pairs, window_size=2, negative=1),
    'NetMF(W2,N5)': embedding_builder('NetMF', init_para=netmf_pairs, window_size=2, negative=5),
    'NetMF(W3,N1)': embedding_builder('NetMF', init_para=netmf_pairs, window_size=3, negative=1),
    'NetMF(W3,N5)': embedding_builder('NetMF', init_para=netmf_pairs, window_size=3, negative=5),
    'NetMF(W5,N1)': embedding_builder('NetMF', init_para=netmf_pairs, window_size=5, negative=1),
    'NetMF(W5,N5)': embedding_builder('NetMF', init_para=netmf_pairs, window_size=5, negative=5),
    # 'NetMF(W10,N1)': embedding_builder('NetMF', init_para=netmf_pairs, window_size=10, negative=1),
    # 'NetMF(W10,N5)': embedding_builder('NetMF', init_para=netmf_pairs, window_size=10, negative=5),
    #
    #'GraRep(2,N5)': embedding_builder('GraRep', init_para=grarep_para, negative=5), # GraRep behaves like NetMF(W1)
    #'HONEM': embedding_builder('HONEM'),
}
```

```
In [11]: FON_edges = list((start[-1],next_node) for start in gen_HON.source_paths_len1
for _, next_node, _ in gen_HON.transition_probs(start))
#FON_edges
```

## debug builder.evaluate\_pairs

Before running extensive simulations, examine the best case - no cross validation:

- true\_probs contains all transition probabilities for the excluded\_edges, but we want all of them (set excluded\_edges=FON\_edges)
- the embedding uses the original transition probabilities (gen\_HON)
- The embedding approximates the 1-step transitions [NetMF(window\_size=1) or Experiment (its SGD equivalent)] instead of using random walks
- After decoding the probabilities, keep only those connected by a FON link and re-normalize

**The outcome is worse than I had hoped for, rendering the whole analysis useless.**

11/18/2020

ProbabilityPrediction\_sim

```
In [12]: def debug_evaluate_pairs(builder, dimension=128, title=None, only_FON_edges=True,
    FON_edges=FON_edges, decode_args=decode_args, limit=(0.0001,1)):
    true_probs = builder.get_true_probs(gen_HON, excluded_edges=FON_edges) # true_probs for all edges
    gen_XE = gen_HON # no cross validation
    res = builder.evaluate_pairs(true_probs, gen_XE, dimension, **decode_args)
    df = pd.DataFrame(res).set_index(['source', 'target'])

    df.plot.scatter('true_prob', 'pred_prob', alpha=0.01, logx=True, logy=True
    , xlim=limit, ylim=limit)
    if title is not None: plt.suptitle(title)
    ax = plt.gca()
    ax.plot(limit,limit, '-')
```

```
In [13]: if verbose:
    %time debug_evaluate_pairs(embedding_builders['NetMF(W1,N1)'], 128, 'HON NetMF(W=1,N=1,R=128)') # 3.5 sec
    plt.savefig('tmp/prob_netmf-w1n1r128.png')
    %time debug_evaluate_pairs(embedding_builders['NetMF(W3,N5)'], 16, 'HON NetMF(W=3,N=5,R=16)') # 11.5 sec
    plt.savefig('tmp/prob_netmf-w3n5r16.png')
    #%%time debug_evaluate_pairs(embedding_builders['NetMF(W10,N5)'], 16, 'HON NetMF(W=10,N=5,R=16)') # 5 min
    #plt.savefig('tmp/prob_netmf-w310n5r16.png')

    #%%time debug_evaluate_pairs(embedding_builders['NetMF(W10,N5)'], 256, 'NetMF(W10,N5)')
    #E%%time debug_evaluate_pairs(embedding_builder('Experiment', dict(seed=1,
    neg_stationary=True), steps=1000, negative=1), 256, 'Experiment')
```

```
In [14]: #for name, builder in embedding_builders.items():
#    %%time debug_evaluate_pairs(builder, 16, name)
```

## Crossvalidation of probability prediction for single embedding

Compares true and predicted probabilities using cross validation.

Cross validation means, that for a set of edges, a new instance `gen_XE` of `CrossValidation_HigherOrderPathGenerator` is generated, where all information about the transition probabilities along these edges is hidden. An embedding is trained using `gen_XE` and the (predicted) transition probabilities are decoded from this embedding.

11/18/2020

ProbabilityPrediction\_sim

```
In [15]: from sklearn.model_selection import KFold
def evaluate_pairs(gen, builder, builder_id, dimension=128, n_splits=20, use_FON=False, plot=False, decode_args=decode_args, limit=(0.00001,1)):
    gen_build = gen
    if use_FON:
        gen_build = HigherOrderPathGenerator(node_sort_key=gen._node_sort_key,
                                              id=gen._id + ' (FON)')
    for key in gen.rule_keys:
        if len(key)>1:
            continue
        for start, next_node, prob in gen.transition_probs(key):
            gen_build.add_rule(start, next_node, prob)
    res = []
    kf = KFold(n_splits=n_splits, shuffle=True, random_state=0)
    progress = tqdm(range(kf.n_splits))
    progress_iter = iter(progress)
    for i_split, (_, test_index) in enumerate(kf.split(FON_edges)):
        next(progress_iter)
        #progress.set_postfix(dict(split=i_split))
        excluded_edges = list(e for i,e in enumerate(FON_edges) if i in test_index)
        true_probs = embedding_builder.get_true_probs(gen, excluded_edges)
        gen_XE = CrossValidation_HigherOrderPathGenerator(gen_build, excluded_edges, '%s exclude %d' % (gen_build._id,i_split))
        res.extend(builder.evaluate_pairs(true_probs, gen_XE, dimension, **decode_args))
    try:
        next(progress_iter)
    except StopIteration:
        pass
    res = pd.DataFrame(res).set_index(['source', 'target'])
    if plot:
        res.plot.scatter('true_prob', 'pred_prob', alpha=0.01, logx=True, logy=True, xlim=limit, ylim=limit)
        ax = plt.gca()
        ax.plot(limit, limit, '-')
        plt.gcf().suptitle((('FON' if use_FON else 'HON') + ' ' + builder_id + f' CV(k={n_splits})'))
    return res
```

```
In [16]: if verbose:
    df = evaluate_pairs(gen_HON, embedding_builders['NetMF(W1,N1)'], 'NetMF(W=1,N=1,R=128)', 128, n_splits=100, use_FON=False, plot=True)
    plt.savefig('tmp/cv-prob_netmf-w1n1r128.png')
```

```
In [17]: if verbose:
    df2 = evaluate_pairs(gen_HON, embedding_builders['NetMF(W2,N1)'], 'NetMF(W=2,N=1,R=16)', 16, n_splits=100, use_FON=False, plot=True)
    plt.savefig('tmp/cv-prob_netmf-w2n1r16.png')
```

```
In [18]: if verbose:
    df2 = evaluate_pairs(gen_HON, embedding_builders['NetMF(W3,N5)'], 'NetMF(W=3,N=5,R=16)', 16, n_splits=100, use_FON=False, plot=True)
    plt.savefig('tmp/cv-prob_netmf-w3n5r16.png')
```

11/18/2020

ProbabilityPrediction\_sim

## Compare embeddings

```
In [19]: from sklearn.model_selection import KFold
res = []
kf = KFold(n_splits=20, shuffle=True, random_state=0)
progress = tqdm(range(kf.n_splits * len(embedding_builders) * len(dimensions) * 2))
progress_iter = iter(progress)
for i_split, (_, test_index) in enumerate(kf.split(FON_edges)):
    ##excluded_edges = FON_edges[test_index]
    excluded_edges = list(e for i,e in enumerate(FON_edges) if i in test_index)
    true_probs = embedding_builder.get_true_probs(gen_HON, excluded_edges)
    for gen_name,gen in [('HON', gen_HON), ('FON', gen_FON)]:
        gen_XE = CrossValidation_HigherOrderPathGenerator(gen, excluded_edges,
'%'s exclude %d' % (gen_name,i_split))
        for builder_name, builder in embedding_builders.items():
            for dimension in dimensions:
                next(progress_iter)
                progress.set_postfix_str(f'{i_split}: {gen_name} {builder_name}{dimension}')
                #progress.set_postfix(dict(split=i_split, gen=gen_name, model=builder_name, dim=dimension))
                out = dict(model=gen_name, dimension=dimension, embedding=builder_name, split=i_split)
                out.update(builder.evaluate(true_probs, gen_XE, dimension, **decode_args))
                res.append(out)
try:
    next(progress_iter)
except StopIteration:
    pass
df = pd.DataFrame(res)
df.head()
```

Out[19]:

	model	dimension	embedding	split	SSE	SSElog
0	HON	16	NetMF(W1,N1)	0	3.811116	27535.852668
1	HON	32	NetMF(W1,N1)	0	5.826772	29355.270917
2	HON	64	NetMF(W1,N1)	0	7.298458	31640.753643
3	HON	128	NetMF(W1,N1)	0	8.006961	33005.995388
4	HON	16	NetMF(W1,N5)	0	4.222847	28841.761619

11/18/2020

ProbabilityPrediction\_sim

```
In [20]: df_sum = df.groupby(['model', 'dimension', 'embedding'])['SSE', 'SSElog'].sum()
df_sum.to_csv('tmp/cv-prob.csv', sep='\t', encoding='utf-16')
```

C:\Users\mstud\Anaconda3\envs\pathpy\lib\site-packages\ipykernel\_launcher.py:  
 1: FutureWarning: Indexing with multiple keys (implicitly converted to a tuple  
 of keys) will be deprecated, use a list instead.  
 """Entry point for launching an IPython kernel.

```
In [21]: df_sum['SSE'].sort_values()
```

```
Out[21]: model dimension embedding
FON    16      NetMF(W5,N5)    47.792190
          NetMF(W3,N5)    49.280430
          NetMF(W5,N1)    49.350664
HON    16      NetMF(W5,N5)    49.384262
          NetMF(W5,N1)    50.371449
          ...
          128      NetMF(W1,N5)    93.412453
          64       NetMF(W1,N1)    106.354252
FON    64       NetMF(W1,N1)    106.354252
HON    128      NetMF(W1,N1)    132.641770
FON    128      NetMF(W1,N1)    132.641770
Name: SSE, Length: 64, dtype: float64
```

```
In [22]: df_sum['SSElog'].sort_values()
```

```
Out[22]: model dimension embedding
HON    16      NetMF(W5,N1)    524060.605983
          32       NetMF(W5,N1)    526176.746898
          16       NetMF(W3,N1)    531328.282973
          32       NetMF(W3,N1)    538490.524515
          16       NetMF(W2,N1)    539406.553885
          ...
          128      NetMF(W1,N5)    636259.354630
FON    64       NetMF(W1,N1)    656262.486685
HON    64       NetMF(W1,N1)    656262.486685
          128      NetMF(W1,N1)    679515.663336
FON    128      NetMF(W1,N1)    679515.663336
Name: SSElog, Length: 64, dtype: float64
```

```
In [ ]:
```



# **Appendix D**

# **Process Documentation**

The following documents illustrate my plans for this thesis.

## MSc Thesis: first ideas

According to the encoder-decoder framework of [Hamilton,Ying,Leskovec\_2017b], the principal decision is whether a **shallow embedding** is learned, or whether a **neighbourhood aggregation** is used. At this early stage of the thesis, I cannot assess which approach works best. Hence, different ideas may be pursued in parallel for now.

### Adaptation of a generic shallow embedding

Many of the baseline methods for creating shallow embeddings ([node2vec], [DeepWalk], [NetMF], [Honem], [Hope], [Line-2]) aim at finding a representation suitable to reconstruct **pairwise** (kind of) proximities – which corresponded to combinations of one- or more-step transition probabilities.

How are we going to generalize these pairwise probabilities to higher order networks?

#### Reasons for embedding higher-order (De Bruijn) memory nodes

For the first order network, the embedding  $Z_j$  of a node  $X_j$  is used to reconstruct  $P(X_j|X_{j-1})$  from the scalar product  $\langle Z_j, Z_{j-1} \rangle$  (and some transformation, e.g. softmax). Therefore, the natural **generalization to higher-order networks** is IMHO reconstructing  $P(X_j|X_{j-1}, \dots, X_{j-k})$  via the scalar product of  $\langle Z_j, Z_{j-1}, \dots, Z_{j-k} \rangle$  where the tuple  $(X_{j-1}, \dots, X_{j-k})$  is embedded as  $Z_{j-1}, \dots, Z_{j-k}$ .

The transition probability  $P(X_j|X_{j-1}, \dots, X_{j-k})$  equals the corresponding edge probability in a  $k^{\text{th}}$  order De Bruijn Graph. Observing that the tuples of different lengths form a **hierarchy**, we evaluate the transition probabilities for different orders instead of using a fixed order  $k$ . Furthermore, all these transition probabilities form the basis for generating the null model for anomaly detection with [Hypa].

This hierarchy is also suitable to **regularize** embeddings for configurations lacking data, see [Hamilton,Ying,Leskovec\_2017b, § 2.5.2 tying node embeddings across layers]. It also allows for curriculum learning.

#### Why (not) stick to pairwise proximities

Note that the term “High-order proximity” is used to describe pairwise proximities in [Hope].

A generalization might base on a similarity measure between two points, which somehow is based on paths through these two points.

However, the whole issue about bothering to use higher order models is the assumption that some transition probabilities are sensitive to steps of the random walk previous to the last one. Therefore, I guess that link prediction – apparently the most prominent task – based on pairwise similarities is doomed to fail. (No need to use higher order models to conclude that points within a community are more likely to be connected...)

*Therefore, I had intended not to follow this approach.*

However, I like the idea in [Hope] of breaking the symmetry by using **different embeddings for source and target nodes** of a directed edge.

### Adaptation of GraphSAGE

[GraphSAGE] repeatedly aggregates features over a neighbourhood, but relies on random sampling instead of aggregating over the whole neighbourhood. This has the advantage of dealing with a fixed amount of data in each iteration.

However, the repeated aggregation inherently assumes a first order model (i.e. memoryless random walk). Hence, instead of  $k$  iterations of first order neighbourhoods, we should first aggregate over the neighbourhoods of each memory node in a  $k^{\text{th}}$  order De Bruijn graph, and repeat this for decreasing orders  $k$ . (Whether the second iteration uses  $(k-1)^{\text{th}}$  order only or the union with  $k^{\text{th}}$  order needs to be decided.)

### Memory

Since calculating the embedding builds on top of fitting a higher-order model to paths, we need to assume, that for the latter memory consumption is not prohibitive. However, I assume that the average degree of a higher-order (De Bruijn) memory node is much smaller than the dimension of the embedding (e.g. 128). Therefore, storing all the embedded values (assuming shallow embedding) could be much costlier (e.g. 50 times) than storing the transition probabilities in a sparse matrix.

The situation is similar for the higher-order adaption of GraphSAGE.

Therefore, there is no clear winner in terms of memory requirement.

14.11.2020 [https://idimail07.lotus.uzh.ch/mail/mistude.nsf/\(%24Sent\)/FFB0584FCC1AF3EEC12585AF004E495C/?OpenDocument&Form=h\\_PrintU...](https://idimail07.lotus.uzh.ch/mail/mistude.nsf/(%24Sent)/FFB0584FCC1AF3EEC12585AF004E495C/?OpenDocument&Form=h_PrintU...)

**Von:** Michael Markus Studer/at/UZH  
**An:** perri@ifi.uzh.ch  
**Kopie:** scholtes@ifi.uzh.ch

**Datum:** Freitag, 24. Juli 2020 20:03  
**Betreff:** MSc Thesis Michael

Dear Ingo  
 Dear Vincenzo

@Ingo: summary of the summary - embedding lattice with slow-down/speed-up successful;  
 [InfiniteWalk] is limit of [NetMF] (which approximates [DeepWalk]) for window size to infinity and creates embedding from Laplacian, which might be applied to higher-order De Bruijn graph; Progress: code & new papers.

In the meeting of July 16, Ingo mentioned, that I might discuss with you (Vincenzo) some topics of my thesis. Not sure, what Ingo already told you about and whether you indeed need know my thoughts explained below...

Since I have been reading into the topic for some time, I thought that I might give you some summary about my ideas. This will hopefully simplify the discussion.

---

You probably already know that I am investigating embeddings for higher-order networks (i.e. higher-order markovian).

As far as transition probabilities in higher-order networks are concerned, I will build on top of MON, MoGEN, or [BuildHON+].

\* Example network using non-markovian lattice

As mentioned in [Slow-down/speed-up], the higher-order terms might affect the diffusion speed, and therefore Ingo suggested investigating a two dimensional lattice where horizontal transitions are encouraged (e.g. moving right after moving right gets a higher probability at the expense of moving left after moving right) and vertical transitions are discouraged (e.g. after moving up, moving down is preferred over moving up again).

\* Skip-gram with negative sampling (SGNS, aka Word2Vec) model

So far, I have focused on SGNS and therefore adjusted the random walk in [DeepWalk] to reflect the higher-order dynamics. As expected, the embeddings of the above lattice network reveals, that the nodes connected by a horizontal edge are closer than those connected by a vertical edge.

[LevyGoldberg] demonstrate, that SGNS corresponds to a matrix factorization (via SVD) of the logarithmic shifted positive pointwise mutual information (PMI), which inspired another embedding [NetMF]. (Which indeed confirms the effects of speed-up and slow-down on the embeddings by DeepWalk.) But we have to keep in mind, that the result of [LevyGoldberg] holds only if the embedding allows for perfect reconstruction.

[InfiniteWalk] covers the case of DeepWalk with window size tending to infinity and its connection to the Laplacian.

\* Directed graphs

These embeddings all assume undirected graphs, which simplifies the calculation of multi-step probabilities but is generally not fulfilled. In a higher-order model, these multi-step probabilities

14.11.2020 [https://idimail07.lotus.uzh.ch/mail/mistude.nsf/\(%24Sent\)/FFB0584FCC1AF3ECC12585AF004E495C/?OpenDocument&Form=h\\_PrintU...](https://idimail07.lotus.uzh.ch/mail/mistude.nsf/(%24Sent)/FFB0584FCC1AF3ECC12585AF004E495C/?OpenDocument&Form=h_PrintU...)

anyway have to be replaced - no big deal. Unfortunately, Word2Vec translates paths into word-context pairs by treating paths essentially as bidirectional. Hence, I will have to move on from using simple black boxes (Word2Vec and SVD) to more general methods of optimization (SGD).  
 Concerning directed graphs: I intend to borrow from [HOPE] the idea of having different embeddings for source and target nodes.

#### \* Beyond pairwise interactions

For a first-order model, it is fine to consider embeddings, which optimize a loss function depending on pairwise interactions of nodes only. However, there is no way to reconstruct the distribution of paths of length 3 from two-dimensional marginal distributions. In order to reconstruct/approximate transition probabilities from an embedding, I will need estimates of  $P(S_k|S_0, \dots, S_{k-1})$  and not just  $P(S_k|S_0)$  for random walk  $[S_0, \dots, S_k]$ .

Speaking in terms of encoders and decoders [EncDec], I consider two approaches for approximating  $P(S_k|S_0, \dots, S_{k-1})$ :

- 1) DEC( ENC\_source( $S_0 \rightarrow \dots \rightarrow S_{k-1}$ ), ENC\_target( $S_k$ ) )
- 2) DEC( ENC\_0( $S_0$ ), ..., ENC\_{k-1}( $S_{k-1}$ ), ENC\_target( $S_k$ ) )

For Approach 1, the decoder combines the two encodings with a scalar product (and a soft-max), which might be generalized for approach 2 to a higher-order inner product (used in [HOSGNS]), which is solved by tensor factorization.

Approach 1 still fits into the PMI framework and can be solved by SVD. For fully connected graphs, the number of paths to embed individually grows exponentially with the length. (But so do MON & MoGEN...)

This might be addressed by the fact that the subpaths form a hierarchy and that  $P(S_k|S_1, \dots, S_{k-1}; S_0 \text{ exists}) = \sum_x P(S_k|S_0=x, S_1, \dots, S_{k-1}) * P(S_0=x|S_1, \dots, S_{k-1})$   
 i.e.  $P(S_k|S_1, \dots, S_{k-1}; S_0 \text{ exists})$  is a weighted average of  $P(S_k|S_0, S_1, \dots, S_{k-1})$ .  
 [EncDec, equation (18)] suggests regularizing penalties across layers. Keeping this in mind, I intend to get some experience with SGD too.

Related to approach 1, I could apply [InfiniteWalk] to the maximal-order De Bruijn graph of the higher-order model:

DEC( ENC\_source( $S_0 \rightarrow \dots \rightarrow S_{k-1}$ ), ENC\_target( $S_1 \rightarrow \dots \rightarrow S_k$ ) )

But I will have to see how to replace ENC\_target( $S_1 \rightarrow \dots \rightarrow S_k$ ) by ENC\_target( $S_k$ ) and I will also need to dig into the proofs to see, whether they hold for directed graphs too.

#### \* Experiments

So far, I did not spend much time thinking about experiments. I primarily had the reconstruction of the transition matrix in mind; maybe predicting the top 10 links, etc. However, having already trained embedding vectors for each subpath (in approach 1) means, that the first-order topology of the network is already known.

It is hard to tell, whether an embedding like ENC\_source( $S_0 \rightarrow \dots \rightarrow S_{k-1}$ ) is still useful for other tasks.

Most likely, I will seek your advice about experiments in future discussions.

Best,  
 Michael

References ([LevyGoldberg] and [EncDec] are key to understanding):

[BuildHON+] <https://arxiv.org/abs/1712.09658>

[Slow-down/speed-up] <https://www.nature.com/articles/ncomms6024/>

[DeepWalk] <https://dl.acm.org/doi/pdf/10.1145/2623330.2623732>

[LevyGoldberg] <http://papers.nips.cc/paper/5477-neural-word-embedding-as>

[NetMF] <https://dl.acm.org/doi/pdf/10.1145/3159652.3159706>

14.11.2020 [https://idimail07.lotus.uzh.ch/mail/mistude.nsf/\(%24Sent\)/FFB0584FCC1AF3EEC12585AF004E495C/?OpenDocument&Form=h\\_PrintU...](https://idimail07.lotus.uzh.ch/mail/mistude.nsf/(%24Sent)/FFB0584FCC1AF3EEC12585AF004E495C/?OpenDocument&Form=h_PrintU...)

[InfiniteWalk] <https://arxiv.org/abs/2006.00094>  
[HOPE] <https://dl.acm.org/doi/pdf/10.1145/2939672.2939751>  
[EncDec] <https://arxiv.org/abs/1709.05584>  
[HOSGNS] <https://arxiv.org/abs/2006.14330>

[https://idimail07.lotus.uzh.ch/mail/mistude.nsf/\(%24Sent\)/FFB0584FCC1AF3EEC12585AF004E495C/?OpenDocument&Form=h\\_PrintUI&PresetField...](https://idimail07.lotus.uzh.ch/mail/mistude.nsf/(%24Sent)/FFB0584FCC1AF3EEC12585AF004E495C/?OpenDocument&Form=h_PrintUI&PresetField...) 3/3

14.11.2020 [https://idimail07.lotus.uzh.ch/mail/mistude.nsf/\(%24Sent\)/8EDBAC1639915E00C12585D9006C4727/?OpenDocument&Form=h\\_PrintUI...](https://idimail07.lotus.uzh.ch/mail/mistude.nsf/(%24Sent)/8EDBAC1639915E00C12585D9006C4727/?OpenDocument&Form=h_PrintUI...)

**Von:** Michael Markus Studer/at/UZH  
**An:** scholtes@ifi.uzh.ch, perri@ifi.uzh.ch

---

**Datum:** Samstag, 05. September 2020 00:23  
**Betreff:** Status MSc Thesis Michael Studer

---

Dear Ingo  
 Dear Vincenzo

Since I managed to implement the analyses discussed with Vincenzo one week ago, I would like to provide some status update. Some PDFs displaying Jupyter notebooks are also added to provide detailed information.

Examining the various papers on embeddings revealed the following tasks (starting with the most important one):

- 1) Classification
- 2) Link prediction
- 3) Graph reconstruction
- 4) Visualization

...  
 While HONEM (my direct competitor) covers all these four tasks, there is an unanimously preference for **multi-label classification** by DeepWalk, Node2vec, and NetMF (my primary inspiration).

#### Classification task

Following Vincenzo's advice, I decided, that the following datasets from SocioPatterns are best suited for a **multi-class classification** task:

- Contacts in a workplace (2013)
- Primary school temporal network data
- Hospital ward dynamic contact network

Using pathpy2 and the kdd2018-tutorial, I extracted paths, estimated the order (which happens to be 2 for all datasets), and exported the transition probabilities of order 1 and 2, see the attached SocioPatterns.pdf.

Next follows calculating various embeddings for visualization and classification learning - the latter calculates also score, F1-micro, F1-macro, confusion matrix, and auc roc. The models for classification were Logistic Regression, SVC, and Random Forest (20 and 100 estimators) - basically what came to my mind first. The following embeddings were used: HON NetMF (window\_size 2 and 10), FON NetMF (window\_size 2 and 10), HON DeepWalk (window\_size 10), FON DeepWalk (window\_size 10), and HONEM. In order to investigate the effect of using a HON vs a FON, I choose to include both a FON and a HON version of each embedding. Consequently, I excluded Node2vec as this would require a HON variant too - but this is perfectly feasible by manipulating the transition probabilities before using DeepWalk.

I have to admit, that I did not perform an exhaustive search over many parameters, and the results are rather a PoC than what is expected for the thesis.

The results were encouraging for the primary school dataset (see attached Classification\_sim\_v1\_primaryschool.pdf), and I calculated the following:

- Visualizations (output 19) are ok
- F1-micro (output 25/26) and F1-macro (output 27/28) reveals that HON performs better than FON (averages over 10 simulations).

14.11.2020 [https://idimail07.lotus.uzh.ch/mail/mistude.nsf/\(%24Sent\)/8EDBAC1639915E00C12585D9006C4727/?OpenDocument&Form=h\\_PrintUI...](https://idimail07.lotus.uzh.ch/mail/mistude.nsf/(%24Sent)/8EDBAC1639915E00C12585D9006C4727/?OpenDocument&Form=h_PrintUI...)

- There is also a confusion matrix (output 31) - for 1 simulation only, hence not directly comparable to F1 scores. (I might do this later.)
- For AUC (output 33), I aggregated pairs of true labels and decision scores over 10 simulations with 5 splits each and calculated the roc curves, using a one-vs-rest "approximation" (i.e. I hope this is correct). However, for the workplace dataset, *HON performs worse* (see Classification\_sim\_v1\_workplace.pdf). Finally, for the hospital ward data, the F1 scores are generally lower than for the other two datasets and there is a tie between FON and HON (see Classification\_sim\_v1\_hospital.pdf).

#### Link prediction (or rather probability prediction) task

I also tried to predict transition probabilities using a **cross validation** approach, where I carefully removed any information (besides its existence) about the probabilities of following a link by tweaking the transition probabilities (in CrossValidation\_HigherOrderPathGenerator, which is not shown though).

I had intended to do a grid search to identify parameters of (HON or FON) NetMF, which are best suited for this task.

Before actually running the simulation, I stepped back and examined, how well the probability reconstruction from the embedding actually works in the **best case** (i.e. no cross validation and even utilizing knowledge about network topology), see end of page 5 of the attached ProbabilityPrediction\_sim.pdf. The scatter plots (page 7) were quite **disappointing**, causing me to consider this task a failure. May be datasets from the MON paper (when is a network a network?) - AIR, WIKI, TUBE - work better.

#### Conclusion

I think, the classification task is both relevant and convincing.

However, since the link prediction failed, I still have no example supporting my initial hypothesis, that for HON embeddings, pairwise decoding is not enough. (In fact, the findings using synthetic data do explain the benefit of going beyond pairwise.)

Unless instructed otherwise, I will focus on writing the thesis until our next meeting.

Best,  
Michael

Anhänge:

SocioPatterns.pdf	Classification_sim_v1-primaryschool.pdf	Classification_sim_v1-workplace.pdf	Classification_sim_v1-hospital.pdf
ProbabilityPrediction_sim.pdf			

15.10.2020

Calc\_q

## Modularity calculation

Calculating the modularity of a higher order network may take a long time - especially when looping over all the pairs of nodes. Observe, that for a network  $n$  first-order nodes (each with constant degrees  $d$ ), the number of  $k$ -th order nodes is  $n \cdot d^{order-1}$ . Therefore, the adjacency matrix is very sparse and its number of non-zero entries is only  $n \cdot d^{order}$ . Hence, it is crucial to avoid iterating over all  $n^2 \cdot d^{2order-2}$  pairs of nodes.

The modularity  $q$  compares the number of actual edges connecting nodes of the same class with its expected counterpart. To count the actual edges, we need only to iterate over the non-zero entries of the adjacency matrix. While the term containing the expected number of edges apparently needs a loop over all pairs of nodes, its special structure (outer product) enables a faster calculation:

$$\sum_{i,j} a[i] \cdot b[j] \cdot 1_{C[i]=C[j]} = \sum_c (\sum_i a[i] \cdot 1_{C[i]=c}) (\sum_j b[j] \cdot 1_{C[j]=c})$$

Furthermore, the implementation supports generalization of modularity to

- directed networks, as in Leicht and Newman (2007), "Community structure in directed networks"
- weighted networks, as in Newman (2004), "Analysis of weighted networks"

In [1]:

```
# this workbook needs pathpy2 installed
%conda list pathpy

# packages in environment at C:\Users\mstud\Anaconda3\envs\pathpy2:
#
# Name           Version      Build Channel
pathpy2          2.2.0       pypi_0   pypi
Note: you may need to restart the kernel to use updated packages.
```

In [2]:

```
import pathpy as pp
import numpy as np
from collections import defaultdict
```

15.10.2020

Calc\_q

In [3]:

```
def calc_q(net, C, weighted=False, **kwargs):
    nodes = list(net.nodes.keys())
    mat = net.adjacency_matrix(weighted=weighted, **kwargs) # scipy.sparse.csc.csc_matrix
    mat_sum = 0
    mat_sum_by_row_C = defaultdict(float)
    mat_sum_by_col_C = defaultdict(float)
    q = 0
    for c in range(mat.shape[1]):
        c_C = C[nodes[c]] # class of current column's node
        for ind in range(mat.indptr[c], mat.indptr[c+1]):
            r = mat.indices[ind]
            v = mat.data[ind]
            # assert mat[r,c]==v
            r_C = C[nodes[r]] # class of current rows's node
            mat_sum += v
            mat_sum_by_row_C[r_C] += v
            mat_sum_by_col_C[c_C] += v
            if c_C == r_C:
                q += v
    q_exp = sum( v*mat_sum_by_col_C[c] for c,v in mat_sum_by_row_C.items() ) / (mat_sum**2)
    q = q/mat_sum - q_exp
    q_max = 1 - q_exp
    return (q, q_max)
```

calc\_q() still uses for loops. Consider converting with .tocoos(), which represents the matrix with three arrays (row, col, data), and use numpy array operations.

## Some testdata

to compare with existing implementation

In [4]:

```
n = 100
net = pp.algorithms.random_graphs.erdos_renyi_gnp(n, 0.2, self_loops=False, directed=False)
for e in net.edges:
    net.edges[e]['weight'] = np.random.exponential()
net_D = pp.algorithms.random_graphs.erdos_renyi_gnp(n, 0.2, self_loops=False, directed=True)
for e in net_D.edges:
    net_D.edges[e]['weight'] = np.random.exponential()
```

In [5]:

```
classes = { str(i): 'group %d' % (i%5) for i in range(n) }
```

15.10.2020

Calc\_q

In [6]:

```
q = pp.algorithms.modularity.q(net, C=classes)
q_max = pp.algorithms.modularity.q_max(net, C=classes)
if q_max < 0:
    q_max+= 1
print(q,q_max)
```

-0.016743908659001307 0.7998893578740652

In [7]:

```
calc_q(net, classes, transposed=True) # weighted=False
```

Out[7]:

( -0.016743908659001383, 0.7998893578740648 )

Calculation works for all variants of adjacency matrices (weighted and/or directed):

In [8]:

```
calc_q(net, classes, weighted=True, transposed=True)
```

Out[8]:

( -0.01852957868856711, 0.798714644094694 )

In [9]:

```
calc_q(net_D, classes, weighted=True, transposed=True)
```

Out[9]:

( 0.010614981011942137, 0.8010128521279705 )

In [10]:

```
# calculation is invariant under transposition of the matrix
calc_q(net_D, classes, weighted=True, transposed=False)
```

Out[10]:

( 0.010614981011942137, 0.8010128521279705 )

In [11]:

```
# works also for HigherOrderNetwork.adjacency_matrix(include_subpaths, weighted, transposed)
# calc_q(hon, classes, weighted=True, include_subpaths=True)
```

In [ ]:



# Nomenclature

## Network

**A** Adjacency matrix

**D** Degree matrix

$d_u$  Degree of node  $u$

**D<sub>neighbor</sub><sup>k</sup>**  $k^{\text{th}}$ -order neighborhood matrix (HONEM)

$\mathcal{E}$  Set of edges

$\mathcal{E}_R, \mathcal{E}_U$  Rightward/horizontal and upward/vertical edges (lattice 2D)

$\mathcal{E}_{2R}, \mathcal{E}_{2U}$  Rightward and upward two-step edges (lattice 2D)

$K$  Max order of Markov model

$k$  Number of folds

**L** Laplacian matrix

$L, R, U, D$  Directions (HON Lattice 2D)

$\overrightarrow{S_{1..L}}$  Random walk (abbreviation for  $S_1 \rightarrow \dots \rightarrow S_L$ )

$T_{s \rightarrow t}$  First-order transition probability

$T_{s_{1..k} \rightarrow t}^{(k)}$   $k^{\text{th}}$ -order transition probability

$\mathcal{V}$  Set of nodes

$\pi$  Stationary distribution

$\omega$  Parameter for HON lattice 2D controlling the higher-order dynamic

## Skipgram model and related embeddings

$p, q$  Bias parameters (Node2vec)

$p(w, c)$  Distribution of word–context pairs (skip–gram model)

$p(w)$  Marginal distribution of  $p(w, c)$  (skip–gram model)

$p(c|w)$  Conditional distribution of context given the word (skip–gram model)

$\hat{p}(c|w)$  Approximation  $\langle \vec{w}, \vec{c} \rangle / \text{const}(w)$  of  $p(c|w)$  based on embeddings (skip-gram model)

$\vec{w}, \vec{c}$  Embedding of word  $w$  and context  $c$  (skip-gram model)

$p_N(c)$  Distribution of the context for negative samples (skip-gram model)

$PMI(w, c)$  Pointwise mutual information ( $p(c|w)/p_N(c)$ ) (skip-gram model)

$L$  Either last index or length of a walk.

$L_{\max}$  Maximal size of a tuple for the consistency of walk probabilities

$\mathcal{L}$  Loss function

$N$  Number of negative samples

$R$  Dimension of embedding space

$S$  Number of steps (GraRep)

$\tau$  Penalty parameter

$W$  Window size

### Indices

$\ell$   $\ell \leq L$  or  $\ell \leq k$

$w, c$  Word and context (skip-gram model)

$u, v$  nodes of a network (nodes of a nedge); also De Bruijn nodes

$s, t$  Source/start and target (transition probability)

# List of Figures

2.1	Lengths and angles for LINE embeddings of Lattice 2D . . . . .	11
2.2	T-SNE for LINE embeddings of Lattice 2D . . . . .	12
2.3	Effect of rule pruning on HONEM neighborhood matrix. . . . .	15
3.1	HON Lattice 2D weights . . . . .	24
3.2	Two Embeddings of HON Lattice 2D displaying stretch effect or clustering by parity . . . . .	28
3.3	Projection of HON NetMF embeddings of the HON Lattice 2D .	28
3.4	Distribution of embedded edge lengths and their angles towards average . . . . .	29
4.1	Stretch ratios for embeddings of HON Lattice 2D . . . . .	34
4.2	Comparing FON and HON embedding visually for stretch effect	35
4.3	Concentration of PMI for HON NetMF . . . . .	36
4.4	Stretch ratios by window size for HON NetMF . . . . .	36
4.5	Embedding of HON NetMF with ‘mixed’ interactions . . . . .	37
4.6	Neighborhood and t-SNE for HONEM and HON Lattice 2D . .	38
4.7	F1-scores for classification of SocioPatterns data . . . . .	42
4.8	Best HON vs. FON for classification of PrimarySchool with $\Delta=20$ sec . . . . .	43
4.9	Best HON vs. FON for classification of Workplace with $\Delta=30$ min.	44
4.10	Best HON vs. FON for classification of HospitalWard with $\Delta=20$ sec . . . . .	45
4.11	FON vs. HON for child–teacher–child interaction probabilities .	46
4.12	Comparing the entropy of pairs of transition probabilities . . . .	47
B.1	Comparing true vs. predicted probabilities (best case) . . . . .	60
B.2	Comparing true vs. predicted probabilities using cross-validation	62
B.3	t-SNE visualization of the experimental estimator . . . . .	62
B.4	Loss during training of the experimental estimator . . . . .	62
C.1	UML class diagram for HigherOrderPathGenerator.py . . . . .	65
C.2	UML class diagram for Embedding.py . . . . .	75
C.3	UML class diagram for Visualizations.py . . . . .	90



# List of Tables

2.1	FON embeddings based on the skip-gram model . . . . .	9
4.1	Average F1-scores for classification of SocioPatterns data . . . . .	41
A.1	F1-scores for the PrimarySchool dataset with $\Delta=20$ sec . . . . .	53
A.2	F1-scores for the Workplace dataset with $\Delta=1$ min . . . . .	54
A.3	F1-scores for the Workplace dataset with $\Delta=5$ min . . . . .	54
A.4	F1-scores for the Workplace dataset with $\Delta=10$ min . . . . .	55
A.5	F1-scores for the Workplace dataset with $\Delta=20$ min . . . . .	55
A.6	F1-scores for the Workplace dataset with $\Delta=30$ min . . . . .	56
A.7	F1-scores for the HospitalWard dataset with $\Delta=20$ sec . . . . .	56
A.8	F1-scores for the HospitalWard dataset with $\Delta=40$ sec . . . . .	57
A.9	F1-scores for the HospitalWard dataset with $\Delta=1$ min . . . . .	57



# List of Algorithms

3.1 Optimized calculation for HON NetMF . . . . .	21
---	----



# Bibliography

- [1] Amr Ahmed, Nino Shervashidze, Shravan Narayananmurthy, Vanja Josifovski, and Alexander J. Smola. [Distributed Large-Scale Natural Graph Factorization](#). In *Proceedings of the 22nd International Conference on World Wide Web*, WWW '13, page 37–48, New York, NY, USA, 2013. Association for Computing Machinery. ISBN 9781450320351. doi: 10.1145/2488388.2488393. URL <https://doi.org/10.1145/2488388.2488393>.
- [2] Y. Bengio, A. Courville, and P. Vincent. [Representation Learning: A Review and New Perspectives](#). *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 35(8):1798–1828, 2013.
- [3] Shaosheng Cao, Wei Lu, and Qiongkai Xu. [GraRep: Learning Graph Representations with Global Structural Information](#). In *Proceedings of the 24th ACM International on Conference on Information and Knowledge Management*, CIKM '15, page 891–900, New York, NY, USA, 2015. Association for Computing Machinery. ISBN 9781450337946. doi: 10.1145/2806416.2806512. URL <https://doi.org/10.1145/2806416.2806512>.
- [4] Shaosheng Cao, Wei Lu, and Qiongkai Xu. [Deep neural networks for learning graph representations](#). In *AAAI*, volume 16, pages 1145–1152, 2016.
- [5] Sudhanshu Chanpuriya and Cameron Musco. [InfiniteWalk: Deep Network Embeddings as Laplacian Embeddings with a Nonlinearity](#), 2020.
- [6] Karl Pearson F.R.S. [LIII. On lines and planes of closest fit to systems of points in space](#). *The London, Edinburgh, and Dublin Philosophical Magazine and Journal of Science*, 2(11):559–572, 1901. doi: 10.1080/14786440109462720. URL <https://doi.org/10.1080/14786440109462720>.
- [7] Valerio Gemmetto, Alain Barrat, and Ciro Cattuto. [Mitigation of infectious disease at school: targeted class closure vs school closure](#). *BMC infectious diseases*, 14(1):695, December 2014. ISSN 1471-2334. doi: 10.1186/PREACCEPT-6851518521414365. URL <http://www.biomedcentral.com/1471-2334/14/3841>.
- [8] Yoav Goldberg and Omer Levy. [word2vec Explained: deriving Mikolov et al.’s negative-sampling word-embedding method](#), 2014.
- [9] Christoph Gute, Giona Casiraghi, Frank Schweitzer, and Ingo Scholtes. [Mogen: A generative multi-order model to predict variable length paths in networks](#). KDD '20, 2020.

- [10] Aditya Grover and Jure Leskovec. [Node2vec: Scalable Feature Learning for Networks](#). In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD '16, page 855–864, New York, NY, USA, 2016. Association for Computing Machinery. ISBN 9781450342322. doi: 10.1145/2939672.2939754. URL <https://doi.org/10.1145/2939672.2939754>.
- [11] MATHIEU GÉNOIS, CHRISTIAN L. VESTERGAARD, JULIE FOURNET, ANDRÉ PANISSON, ISABELLE BONMARIN, and ALAIN BAR-RAT. [Data on face-to-face contacts in an office building suggest a low-cost vaccination strategy based on community linkers](#). *Network Science*, 3:326–347, 9 2015. ISSN 2050-1250. doi: 10.1017/nws.2015.10. URL [http://journals.cambridge.org/article\\_S2050124215000107](http://journals.cambridge.org/article_S2050124215000107).
- [12] Will Hamilton, Zhitao Ying, and Jure Leskovec. [Inductive representation learning on large graphs](#). In *Advances in neural information processing systems*, pages 1024–1034, 2017.
- [13] William L. Hamilton, Rex Ying, and Jure Leskovec. [Representation Learning on Graphs: Methods and Applications](#), 2017.
- [14] M. Khosla, V. Setty, and A. Anand. [A Comparative Study for Unsupervised Network Representation Learning](#). *IEEE Transactions on Knowledge and Data Engineering*, pages 1–1, 2019.
- [15] Tamara G. Kolda and Brett W. Bader. [Tensor Decompositions and Applications](#). *SIAM Review*, 51(3):455–500, 2009. doi: 10.1137/07070111X. URL <https://doi.org/10.1137/07070111X>.
- [16] Joseph B Kruskal. [Multidimensional scaling by optimizing goodness of fit to a nonmetric hypothesis](#). *Psychometrika*, 29(1):1–27, 1964.
- [17] Renaud Lambiotte, Martin Rosvall, and Ingo Scholtes. [From networks to optimal higher-order models of complex systems](#). *Nature physics*, 15(4):313–320, 2019.
- [18] E. A. Leicht and M. E. J. Newman. [Community Structure in Directed Networks](#). *Phys. Rev. Lett.*, 100:118703, Mar 2008. doi: 10.1103/PhysRevLett.100.118703. URL <https://link.aps.org/doi/10.1103/PhysRevLett.100.118703>.
- [19] Omer Levy and Yoav Goldberg. [Neural Word Embedding as Implicit Matrix Factorization](#). In Z. Ghahramani, M. Welling, C. Cortes, N. D. Lawrence, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 27*, pages 2177–2185. Curran Associates, Inc., 2014. URL <http://papers.nips.cc/paper/5477-neural-word-embedding-as-implicit-matrix-factorization.pdf>.
- [20] Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. [Efficient Estimation of Word Representations in Vector Space](#), 2013.

- [21] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. [Distributed Representations of Words and Phrases and their Compositionality](#). In C. J. C. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 26*, pages 3111–3119. Curran Associates, Inc., 2013. URL <http://papers.nips.cc/paper/5021-distributed-representations-of-words-and-phrases>.
- [22] M. E. J. Newman. [Analysis of weighted networks](#). *Phys. Rev. E*, 70:056131, Nov 2004. doi: 10.1103/PhysRevE.70.056131. URL <https://link.aps.org/doi/10.1103/PhysRevE.70.056131>.
- [23] Bryan Perozzi, Rami Al-Rfou, and Steven Skiena. [Deepwalk: Online learning of social representations](#). In *Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 701–710, 2014.
- [24] Simone Piaggesi and André Panisson. [Time-varying Graph Representation Learning via Higher-Order Skip-Gram with Negative Sampling](#), 2020.
- [25] Jiezhong Qiu, Yuxiao Dong, Hao Ma, Jian Li, Kuansan Wang, and Jie Tang. [Network Embedding as Matrix Factorization: Unifying DeepWalk, LINE, PTE, and Node2vec](#). In *Proceedings of the Eleventh ACM International Conference on Web Search and Data Mining*, WSDM ’18, page 459–467, New York, NY, USA, 2018. Association for Computing Machinery. ISBN 9781450355810. doi: 10.1145/3159652.3159706. URL <https://doi.org/10.1145/3159652.3159706>.
- [26] Radim Rehurek. [Deep learning with word2vec and gensim](#), 2013. URL <https://rare-technologies.com/deep-learning-with-word2vec/>. [Online; accessed 2020-11-15].
- [27] Martin Rosvall, Alcides V Esquivel, Andrea Lancichinetti, Jevin D West, and Renaud Lambiotte. [Memory in network flows and its effects on spreading dynamics and community detection](#). *Nature communications*, 5(1):1–13, 2014.
- [28] Mandana Saebi, Giovanni Luca Ciampaglia, Lance M Kaplan, and Nitesh V Chawla. [HONEM: Network Embedding Using Higher-Order Patterns in Sequential Data](#), 2019.
- [29] Ingo Scholtes. [When is a Network a Network? Multi-Order Graphical Model Selection in Pathways and Temporal Networks](#). In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD ’17, page 1037–1046, New York, NY, USA, 2017. Association for Computing Machinery. ISBN 9781450348874. doi: 10.1145/3097983.3098145. URL <https://doi.org/10.1145/3097983.3098145>.
- [30] Ingo Scholtes, Nicolas Wider, René Pfitzner, Antonios Garas, Claudio J Tessone, and Frank Schweitzer. [Causality-driven slow-down and speed-up of diffusion in non-Markovian temporal networks](#). *Nature communications*, 5(1):1–9, 2014.

- [31] Juliette Stehlé, Nicolas Voirin, Alain Barrat, Ciro Cattuto, Lorenzo Isella, Jean-François Pinton, Marco Quaggiotto, Wouter Van den Broeck, Corinne Régis, Bruno Lina, and Philippe Vanhems. [High-Resolution Measurements of Face-to-Face Contact Patterns in a Primary School](#). *PLOS ONE*, 6(8):e23176, 08 2011. doi: 10.1371/journal.pone.0023176. URL <http://dx.doi.org/10.1371/journal.pone.0023176>.
- [32] Jian Tang, Meng Qu, and Qiaozhu Mei. [PTE: Predictive Text Embedding through Large-Scale Heterogeneous Text Networks](#). In *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD '15, page 1165–1174, New York, NY, USA, 2015. Association for Computing Machinery. ISBN 9781450336642. doi: 10.1145/2783258.2783307. URL <https://doi.org/10.1145/2783258.2783307>.
- [33] Jian Tang, Meng Qu, Mingzhe Wang, Ming Zhang, Jun Yan, and Qiaozhu Mei. [LINE: Large-Scale Information Network Embedding](#). In *Proceedings of the 24th International Conference on World Wide Web*, WWW '15, page 1067–1077, Republic and Canton of Geneva, CHE, 2015. International World Wide Web Conferences Steering Committee. ISBN 9781450334693. doi: 10.1145/2736277.2741093. URL <https://doi.org/10.1145/2736277.2741093>.
- [34] Leo Torres, Kevin S Chan, and Tina Eliassi-Rad. [GLEE: Geometric Laplacian Eigenmap Embedding](#). *Journal of Complex Networks*, 8(2), 03 2020. ISSN 2051-1329. doi: 10.1093/comnet/cnaa007. URL <https://doi.org/10.1093/comnet/cnaa007>. cnaa007.
- [35] Anton Tsitsulin, Davide Mottin, Panagiotis Karras, and Emmanuel Müller. [VERSE: Versatile Graph Embeddings from Similarity Measures](#). In *Proceedings of the 2018 World Wide Web Conference*, WWW '18, page 539–548, Republic and Canton of Geneva, CHE, 2018. International World Wide Web Conferences Steering Committee. ISBN 9781450356398. doi: 10.1145/3178876.3186120. URL <https://doi.org/10.1145/3178876.3186120>.
- [36] Philippe Vanhems, Alain Barrat, Ciro Cattuto, Jean-François Pinton, Nagham Khanafer, Corinne Régis, Byeul-a Kim, Brigitte Comte, and Nicolas Voirin. [Estimating Potential Infection Transmission Routes in Hospital Wards Using Wearable Proximity Sensors](#). *PLoS ONE*, 8(9):e73970, 09 2013. doi: 10.1371/journal.pone.0073970. URL <http://dx.doi.org/10.1371%2Fjournal.pone.0073970>.
- [37] Daixin Wang, Peng Cui, and Wenwu Zhu. [Structural Deep Network Embedding](#). In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD '16, page 1225–1234, New York, NY, USA, 2016. Association for Computing Machinery. ISBN 9781450342322. doi: 10.1145/2939672.2939753. URL <https://doi.org/10.1145/2939672.2939753>.
- [38] Jian Xu, Thanuka L. Wickramarathne, and Nitesh V. Chawla. [Representing higher-order dependencies in networks](#). *Science Advances*, 2(5), 2016. doi: 10.1126/sciadv.1600028. URL <https://advances.sciencemag.org/content/2/5/e1600028>.

- [39] Jian Xu, Mandana Saebi, Bruno Ribeiro, Lance M. Kaplan, and Nitesh V. Chawla. [Detecting Anomalies in Sequential Data with Higher-order Networks](#), 2017.
- [40] D. Zhang, J. Yin, X. Zhu, and C. Zhang. [Network Representation Learning: A Survey](#). *IEEE Transactions on Big Data*, 6(1):3–28, 2020.
- [41] Chang Zhou, Yuqiong Liu, Xiaofei Liu, Zhongyi Liu, and Jun Gao. [Scalable graph embedding for asymmetric proximity](#). In *Thirty-First AAAI Conference on Artificial Intelligence*, 2017.