

Hyperedge2vec: Distributed Representations for Hyperedges

Ankit Sharma
University of Minnesota, USA
ankit@cs.umn.edu

Himanshu Kharakwal
LNM-IIT, India
himanshukharakwal765@gmail.com

Shafiq R. Joty
Nanyang Tech. University, Singapore
srjoty@ntu.edu.sg

Jaideep Srivastava
University of Minnesota, USA
srivasta@cs.umn.edu

ABSTRACT

Data structured in the form of overlapping or non-overlapping sets is found in a variety of domains, sometimes explicitly but often subtly. For example, teams, which are of prime importance in social science studies are “sets of individuals”; “item sets” in pattern mining are sets and for various types of analysis in language studies a sentence can be considered as a “set or bag of words”. Although building models and inference algorithms for structured data has been an important task in the fields of machine learning and statistics, research on “set-like” data still remains less explored. Relationships between pairs of elements can be modeled as edges in a graph. However, for modeling relationships that involve all members of a set, a hyperedge is a more natural representation. In this work, we focus on the problem of embedding hyperedges in a hypergraph (a network of overlapping sets) to a low dimensional vector space. We propose a probabilistic deep-learning based method as well as a tensor-based algebraic model, both of which capture the hypergraph structure in a principled manner without losing set-level information. Our central focus is to highlight the connection between hypergraphs (topology), tensors (algebra) and probabilistic models. We present a number of baselines, some of which adapt existing node-level embedding models to the hyperedge-level, as well as sequence based language techniques which are adapted for set structured hypergraph topology. The performance is evaluated with a network of social groups and a network of word phrases. Our experiments show that accuracy wise our methods perform similar to those of baselines which are not designed for hypergraphs. Moreover, our tensor based method is more efficient than deep-learning based auto-encoder method. We, therefore, argue that the proposed methods are more generic methods suitable for hypergraphs (and therefore also for graphs) that preserve accuracy and efficiency.

ACM Reference Format:

In *Proceedings of ACM Conference (Conference'18)*. ACM, London, UK, 11 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

1 INTRODUCTION

In group structured data we have multiple entities related by some form of group relationships. Such data is more abundantly found in the real world than has been usually studied [23]. In social networks domain, team data from massive online multi-player games [4] such as World of Warcraft, group communication tools such as Skype and Google Docs and research collaborations [1, 50]. There are

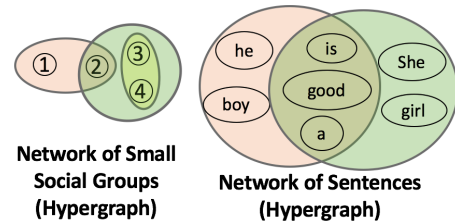


Figure 1: Example illustrating “set-like” hypergraph structure from two domains. Left is a collaboration network between four individuals and on right is a network resulting from two sentences: “He is a good boy” & “She is a good girl”; and eight word nodes.

other fields in which structural relationships between entities is important as well, and large datasets capturing them exist. Examples include Natural Language Processing [10], Biology [32, 33], e-commerce [17, 22] and Chemistry [7]. Figure 1 shows such examples for networks of groups and sentences.

Hypergraph [12], which is a generalization of graphs, is a popular model to naturally capture higher-order relationships between sets of objects (Figure 2) [23]. Within machine learning, algorithms guided by the structure of such higher order networks [61] have found applications in a variety of domains [25, 37, 45, 46, 56].

More recently, the interest of representation learning in NLP [38] has stimulated its application in graph embedding (learning low dimensional representation for graph nodes) [15, 27, 41, 52]. However, this new line of research is limited to simple graphs and we are unaware of any work that considers hyperedge embeddings (learning representations for hyperedges). In this paper, we propose methods which (1) learn hypergraph embeddings directly, (2) leverage the hypergraph topology and (3) not lose the hyperedge-level joint information. These learned embeddings can then be employed by a supervised or semi-supervised algorithm to perform various predictive tasks pertaining to hyperedges. For example, performance prediction of a team (set of individuals) engaged in a collaborative task, or sentiment analysis of a sentence (set of words).

Although, there have been a variety of attempts to learn *node* embeddings for hypergraphs by extending traditional graph embedding methods for hypergraph setting. These approaches differ in the extent they address the aims (2) and (3) as pointed previously. In the first category, there are a number of methods that incorporate the hypergraph topology using proxy graphs [32, 61], therefore, incurring loss of information. However, Agarwal et al.[2] do not agree that such representations can be learned by constructing graphs, which are proxies for the hypergraph structure. In the second category, there

are methods that take into account the hyperedge-level information, like the classical Neural Language Model which maintains the higher order sequence information [10]. But their model is designed for sequences and completely ignores the topology (sequences are linked by common words forming a network of sequences). However recently, there has been interest in modeling “set-like” structures within deep-learning community [43, 59]. But they do not consider the hypergraph structure between the sets and therefore, not model hypergraphs in a principled manner. The third category are those methods that do not ignore the topology completely and also retain the hyperedge-level information. Among this are k -way tensor based methods that explicitly do not work on hypergraph [14, 47], but the connection exists in the form that k -way tensor represents a k -uniform hypergraph [42]. But they are only restricted to uniform hypergraphs.

In this work, we propose two methods that (a) directly learn hyperedge embeddings for general hypergraphs; and (b) capture the hypergraph structure in a principled manner. The first method is a tensor-based algebraic method and the second method is an auto-encoder based deep-learning method. In addition, we also discuss the interesting connection between hypergraphs (topology), tensors (algebra) and probabilistic models. The proposed models are compared with a number of interesting baselines, some of which adapt existing node-level embedding models to hypergraph setting as well as adapting sequence based NLP techniques for set structured hypergraphs. The models are compared on two real world datasets: a social group dataset (network of teams as a hypergraph) from an multi-player online game for team (hyperedge) performance prediction as well as language networks (sentence/phrase hypergraphs) for phrase-level sentiment analysis. Our experiments show that accuracy wise our methods perform similar to those of baselines which are not designed for hypergraphs. Moreover, our tensor based method is more efficient than the deep-learning based auto-encoder method.

The main contributions of this work are as follows:

- We propose the novel concept of a *dual tensor*, corresponding to the hypergraph dual that allows us to get a hyperedge embedding directly.
- We propose a general *hypergraph tensor decomposition* method designed for general hypergraphs (containing different cardinality hyperedges) unlike simple uniform hypergraph tensor decomposition, which is restricted to fixed cardinality hyperedges (i.e. uniform hypergraph). We are unaware of any such works or applications employing this approach.
- We propose the novel idea of using de-noising auto-encoders in a hypergraph setting. Moreover, we also develop techniques for creating noise using random-walks over hasse diagram topology, which is original and unique.
- We highlight and argue that embeddings from tensor based methods have a natural hypergraph theoretic interpretation unlike the deep learning based black-box approaches.
- Our proposed methods learn hypergraph embeddings based on the unique approach of leveraging the *existing structure* present in network data as context, in contrast to context generation based existing graph embedding techniques.

Following is the outline for the rest of the paper. In section 2, we describe the problem definition and statement followed by Section 3,

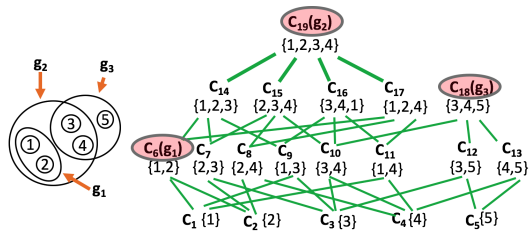


Figure 2: Example of a hypergraph (left) and the hasse diagram (right) corresponding to this hypergraph’s simplicial complex.

where we describe in detail the various methods proposed in this paper. Section 4 describes the datasets, experimental tasks & settings for the models. Section 5 provides an overview of the related literature followed by conclusion and appendix.

2 PRELIMINARIES

In this paper we consider the scenario where we have a collection of *elements*. These elements can represent individual actors in case of social groups or words in sentences or items in item-sets within a transaction database. In other words a social group or a sentence or an item-set are *sets* which contain these *elements*. Let $V = \{v_1, v_2, \dots, v_n\}$ represents n elements and we have m different sets defined over these elements, denoted by $G = \{g_1, g_2, \dots, g_m\}$, where $g_i \subseteq V$ represents the i^{th} set. The cardinality $|g_i|$ represents the number of elements in the set. Also each set $g_i \in G$ has an occurrence number $R(g_i)$, which denotes the number of times it has occurred. Such overlapping or non-overlapping sets can be modeled as a *hypergraph* [12], where the nodes and hyperedges, represent the elements and sets, respectively. This hypergraph is represented as $N_g = (V, G)$ with G as the collection of hyperedges over the nodes V . The incidence matrix $\mathbf{H} \in \{0, 1\}^{|G| \times |V|}$ for N_g represents the presence of nodes in different hyperedges with $\mathbf{H}(g_i, v) = 1$ if $v \in g_i$ else 0. We also define degree $d(v)$ of a vertex v as the number of hyperedges incident on this vertex i.e. $d(v) = \sum_{g_i \in G} \mathbf{H}(g_i, v)$. A specialization of hypergraph model is the *simplicial complex* [40] (Figure 2), in which additionally each hyperedge has the subset closure property, i.e., each subset of hyperedge is also a valid hyperedge.

Problem Statement: Given this setting, our goal is to learn the mapping $\phi : G \rightarrow \mathbb{R}^d$ from hyperedges to feature representations (i.e., embeddings) that can be used to build predictive models involving sets. Here d is a parameter specifying the number of dimensions of the embedding vector. Equivalently, ϕ can be thought of as a look-up matrix of size $|G| \times d$, where $|G|$ is the total number of sets or hyperedges.

3 METHODOLOGY

There are several methods to learn representation of nodes in a graph. Given that a hyperedge is a set of nodes, a natural question is if we can combine node level embeddings (learned using existing methods) within a given hyperedge to find a suitable representation of the hyperedge? Since a large number of ways of combining node embeddings is possible, it reduces to the problem of finding the best combination, with no guarantee that the best one is still found since this approach is greedy. Therefore, we propose two methods that

learn the embeddings for hyperedges directly in a more principled manner. First, a *tensor based method* is focused on retaining the set-level information intact while harnessing the hypergraph network structure. Second an *auto-encoder based method* harnesses the non-linearity of representation offered by deep learning to generate crisp and powerful embeddings. The tensor based method has a natural *hasse diagram* based topological interpretation, and can generate both node and hyperedge embeddings, whereas the deep learning method is more black-box, and can only generate hyperedge embeddings.

3.1 Hyperedge2vec using Hypergraph Tensor Decomposition

In this section we develop tensor (higher-order matrix) based linear algebraic methods that learn node as well as hyperedge embedding by taking into account the joint probability over a hyperedge. The idea behind using tensors is that they retain the set-level information contained in a hypergraph, unlike the proxy graphs (corresponding to hypergraphs) based techniques (used as baselines in our experiments), which approximate hyperedge or set-level information with dyadic edge-level information. This idea is not new, and was proposed earlier by Shashua et. al. [47], and more recently by [26]. The following proposition puts this argument more formally.

PROPOSITION 1. *Given a set of random variables X_1, \dots, X_c ($c \geq 2$), and $H(\cdot)$ as the information entropy, we have ([18, p. 34]):*

$$H(X_1, \dots, X_c) \leq \left(\frac{1}{c-1}\right) \sum_{(i,j) \subseteq [c]} H(X_i, X_j) \quad (1)$$

Therefore, the joint probability distribution over c cardinality hyperedges is more informative (lower entropy) than the sum total of information attained from probability distributions over each of the $\binom{c}{2}$ dyadic edges.

Although tensors can retain higher order information, most research to date has focused on uniform hypergraphs using symmetric tensors. We propose an approach which is principally suited for *general* hypergraph structured data using higher-order tensors. For a given hypergraph we can extract a sub-hypergraph that only consists of the hyperedges with cardinality k . This sub-hypergraph is a k -uniform hypergraph or k -graph [20]. Corresponding to this k -uniform hypergraph, we can define a k^{th} order n -dimensional symmetric tensor [42] $\mathcal{A}_{\text{hyp}}^k = (a_{p_1, p_2, \dots, p_k}) \in \mathbb{R}^{[k, n]}$ whose elements are as follows:

$$a_{p_1, p_2, \dots, p_k} = R(g_i) \quad (2)$$

where $\{v_{p_1}, v_{p_2}, \dots, v_{p_k}\} \in g_i$ and $|g_i| = k, \forall i \in \{1, \dots, m\}$. Note that symmetry here implies that the value of element a_{p_1, p_2, \dots, p_k} is invariant under any permutation of its indices (p_1, p_2, \dots, p_k) . Rest of the elements in the tensor are zeros. We also define the lexicographically ordered index set for hyperedges:

$$\mathcal{P}^k = \{\mathbf{p} | \mathbf{p} = (p_1, p_2, \dots, p_k) \text{ where } \{v_{p_1}, v_{p_2}, \dots, v_{p_k}\} \in g_i, \forall g_i \in G \text{ s.t. } |g_i| = k \text{ and } p_1 < p_2 < \dots < p_k\}, \quad (3)$$

and we have different sets $\{\mathcal{P}^k\} \forall k \in \{c_{min}, \dots, c_{max}\}$ (c_{min} and c_{max} are the maximum and the minimum hyperedge cardinality in the given hypergraph). Notice that \mathcal{P}^k contains unique (non-repetitive) indexes as there is only a single \mathbf{p} corresponding to

each of the different hyperedges $g_i \in G$. Consequently, we have $|\mathcal{P}^k| = |\{g_i : |g_i| = k\}|$.

In a similar manner we can also define a *dual tensor*, corresponding to *hypergraph dual* where the roles of nodes and hyperedges are interchanged. We consider all the hyperedges in the hypergraph dual that are of cardinality k . This basically corresponds to all the vertices in the original hypergraph which have a degree of k , i.e., they are part of exactly k hyperedges in the original hypergraph. Corresponding to this k -uniform hypergraph dual, we can define a k^{th} order m -dimensional symmetric dual tensor $\mathcal{A}_{\text{dual}}^k = (a_{q_1, q_2, \dots, q_k}) \in \mathbb{R}^{[k, m]}$ whose elements are initialized as follows:

$$a_{q_1, q_2, \dots, q_k} = 1 \quad (4)$$

where $\{g_{q_1}, g_{q_2}, \dots, g_{q_k}\} \ni v_j$ and $d(v_j) = k, \forall j \in \{1, \dots, n\}$. Note that this tensor is also symmetric and rest all the elements in the tensor are zeros. Again, we define the lexicographically ordered index set for *dual* hyperedges (vertices in the *original* hypergraph):

$$\mathcal{Q}^k = \{\mathbf{q} | \mathbf{q} = (q_1, q_2, \dots, q_k) \text{ where } v_j \in \{g_{q_1}, g_{q_2}, \dots, g_{q_k}\}, \forall v_j \in V \text{ s.t. } |d(v_j)| = k \text{ and } q_1 < q_2 < \dots < q_k\}, \quad (5)$$

and we have different sets $\{\mathcal{Q}^k\} \forall k \in \{d_{min}, \dots, d_{max}\}$ (d_{min} and d_{max} are the maximum and the minimum vertex degree in the original hypergraph). Again, notice that \mathcal{Q}^k contains unique (non-repetitive) indexes as there is only a single \mathbf{q} corresponding to each of the different dual hyperedge (vertex in the original hypergraph) i.e. $v_i \in V$. Consequently, we have $|\mathcal{Q}^k| = |\{v_i : d(v_i) = k\}|$.

To realize our aim of learning node or hyperedge embeddings we perform Symmetric Tensor Decomposition [19], in a manner similar to [34], but jointly across different cardinality hypergraph tensors (for node embeddings) or dual tensors (for hyperedge embeddings). Specifically, for the hyperedge embeddings we consider the following optimization formulation:

$$f(\lambda, \mathbf{Z}) = \sum_{k=\alpha_1}^{\alpha_2} D_{\text{KL}}(\mathcal{M}^k \| \mathcal{A}_{\text{dual}}^k) = \sum_{k=\alpha_1}^{\alpha_2} \sum_{\mathbf{q} \in \mathcal{Q}^k} (\mathcal{M}_{\mathbf{q}}^k - a_{\mathbf{q}}^k \log \mathcal{M}_{\mathbf{q}}^k) \quad (6)$$

where,

$$\mathcal{M}^k = \sum_{r=1}^d \lambda_r \underbrace{(\mathbf{z}_r \otimes \mathbf{z}_r \otimes \dots \otimes \mathbf{z}_r)}_{k \text{ times}} \equiv \sum_{r=1}^d \lambda_r \mathbf{z}_r^{\otimes k} \quad (7)$$

with \otimes is the Kronecker product (generalized outer product, ref. [5]), $\mathbf{z}_r \in \mathbb{R}^m$, $\lambda_r \in \mathbb{R}$, $\mathbf{Z} \in \mathbb{R}^{m \times d}$ with $\mathbf{Z}(:, r) = \mathbf{z}_r$ and $a_{\mathbf{q}}^k$ is the a_{q_1, q_2, \dots, q_k} entry of $\mathcal{A}_{\text{dual}}^k$. Given the dual tensors $\mathcal{A}_{\text{dual}}^k, \forall k \in \{\alpha_1, \dots, \alpha_2\}$ (α_1 and α_2 are the minimum and maximum cardinalities considered for the dual hyperedges), Equation 6 aims to learn symmetric decomposition \mathcal{M}^k with \mathbf{Z} as the matrix containing d -dimensional embeddings for the m hyperedges. Notice that the embeddings in \mathbf{Z} are shared across the different order (k) decompositions \mathcal{M}^k . Moreover, we employ KL-divergence based loss as our hypergraph data is a count data which is more appropriately modeled using Poisson distribution [16]. Furthermore, although there are $k!$ entries for each k cardinality hyperedge in the dual tensor $\mathcal{A}_{\text{dual}}^k$, however, the summation in Equation 6 is performed over only the

unique set of indexes in Ω^k . This helps us escape the actual construction of the complete dual tensors ($\mathcal{A}_{\text{dual}}^k$) which otherwise may result in exponential space explosion. This idea of using only unique index sets (Ω^k) has been leveraged by both Kolda et. al. [34] [6] as well as by Shashua et. al. [47] in form of semi-norms.

We highlight that $(\mathbf{z}_r^{\otimes k})_{\mathbf{q}} = \mathbf{z}_{q_1 r} \mathbf{z}_{q_2 r} \cdots \mathbf{z}_{q_k r}$ with $\mathbf{z}_{q_j r} = \mathbf{Z}(q_j, r) \in \mathbb{R}$. Then the gradients for Equation 6 are given by:

$$\frac{\partial f(\lambda, \mathbf{Z})}{\partial \lambda_r} = \sum_{k=\alpha_1}^{\alpha_2} \sum_{\mathbf{q} \in \Omega^k} \left\{ (\lambda_r \mathbf{z}_r^{\otimes k})_{\mathbf{q}} - \frac{a_{\mathbf{q}}^k (\lambda_r \mathbf{z}_r^{\otimes k})_{\mathbf{q}}}{\sum_{r=1}^d \lambda_r \mathbf{z}_r^{\otimes k}} \right\} \quad (8)$$

$$\frac{\partial f(\lambda, \mathbf{Z})}{\partial \mathbf{z}_{j r}} = \sum_{k=\alpha_1}^{\alpha_2} \sum_{\mathbf{q} \in \Omega^k} \left\{ (\lambda_r \mathbf{z}_{q_1 r} \mathbf{z}_{q_2 r} \cdots \mathbf{z}_{q_{j-1} r} \cdots \mathbf{z}_{q_k r}) - \frac{a_{\mathbf{q}}^k (\lambda_r \mathbf{z}_{q_1 r} \mathbf{z}_{q_2 r} \cdots \mathbf{z}_{q_{j-1} r} \cdots \mathbf{z}_{q_k r})}{\sum_{r=1}^d \lambda_r \mathbf{z}_r^{\otimes k}} \right\} \quad (9)$$

Furthermore, following [34] we add the following penalty to address the scaling ambiguity by enforcing the following penalty and gradient:

$$p_{\gamma}(\mathbf{Z}) = \gamma \sum_{r=1}^d (\mathbf{z}_r^{\top} \mathbf{z}_r - 1)^2, \quad \frac{\partial p_{\gamma}}{\partial \mathbf{z}_r} = 4\gamma (\mathbf{z}_r^{\top} \mathbf{z}_r - 1) \mathbf{z}_r \quad (10)$$

Another important issue that requires attention is that the optimization function $f(\lambda, \mathbf{Z})$ only considers dual hyperedges of cardinalities within $\{\alpha_1, \dots, \alpha_2\}$. As we will discuss in section 4, due to scalability reasons most likely $\{\alpha_1, \dots, \alpha_2\} \subset \{c_{min}, \dots, c_{max}\}$. Given this setting, it is therefore, possible that some hyperedges (in the original hypergraph) might not contain any vertex which has degree in the range $\{\alpha_1, \dots, \alpha_2\}$. Such hyperedges will then suffer from the cold start issue and to remedy that we leverage the auxiliary information in the form of the hypergraph structure. We therefore, introduce the following penalty and gradient:

$$p_{\eta}(\mathbf{Z}) = \eta \sum_{r=1}^d \mathbf{z}_r^{\top} \mathbf{L}_{\text{dual}} \mathbf{z}_r, \quad \frac{\partial p_{\eta}}{\partial \mathbf{z}_r} = 2\eta \mathbf{L}_{\text{dual}} \mathbf{z}_r \quad (11)$$

where \mathbf{L}_{dual} is the dual hypergraph laplacian (see detail in Appendix A.2). This topology smoothing constraint allows the diffusion of latent embeddings from the *not* cold-start vertices to nearby cold-start vertices, in order for the later to achieve meaningful non-zero embeddings. Notice that the above laplacian contains the entire hypergraph dual structure, therefore, allowing information exchange (1) across different cardinality hyperedges and (2) also between hyperedges that are not suffering from cold start. It would be interesting to develop other laplacians that (a) only affect the embeddings of the cold-start hyperedges, (b) only allow hyperedges of same cardinality to affect each others embeddings and (c) employ hasse diagram based laplacian [46] which explicitly models the cardinality hierarchy.

Putting it all together, the complete optimization objective function is:

$$f_{\text{tot}}(\lambda, \mathbf{Z}) = f(\lambda, \mathbf{Z}) + p_{\gamma}(\mathbf{Z}) + p_{\eta}(\mathbf{Z}) \quad (12)$$

where the choice of γ is the weight of the penalty on the norm of the columns of \mathbf{Z} and the choice of η determines the penalty on the extent of smoothness (similarity) between the embeddings of

hyperedges within neighborhood of each other. We call the algorithm that optimizes the above objective (Equation 12) as **Hypergraph-Tensor-Decomposition (HTD)** (Algorithm 1).

Algorithm 1 HTD ($d, m, \alpha_1, \alpha_2, \gamma, \eta, \mathcal{A}_{\text{dual}}^k \forall k \in \{\alpha_1, \dots, \alpha_2\}, \mathbf{L}_{\text{dual}}$)

- 1: randomly initialize $\mathbf{Z} \in \mathbb{R}^{m \times d}$ and $\lambda \in \mathbb{R}^d$
 - 2: **repeat**
 - 3: **for** $r = 1$ to d **do**
 - 4: Update λ_r using Equation 8
 - 5: **for** $j = 1$ to m **do**
 - 6: Update $\mathbf{z}_{j r}$ using Equation 9
 - 7: **end for**
 - 8: Update \mathbf{z}_r using Equation 10
 - 9: Update \mathbf{z}_r using Equation 11
 - 10: **end for**
 - 11: **until** fit criteria achieved or max. # of iterations exceeded
 - 12: **return** \mathbf{Z}
-

Notice that till now we have described the process of learning hyperedge embeddings. The same process can also be used for vertex embedding same algorithm 1 can be used to get the vertex embeddings, by calling **HTD** ($d, n, \beta_1, \beta_2, \gamma, \eta, \mathcal{A}_{\text{hyp}}^k \forall k \in \{\beta_1, \dots, \beta_2\}, \mathbf{L}_{\text{hyp}}$). Parameters β_1, β_2 are the limits for the hyperedge cardinalities considered in Equation 6 and \mathbf{L}_{hyp} is the hypergraph laplacian (see Appendix A.2). We shall jointly refer to the embeddings achieved for nodes and hyperedges via the above tensor decomposition techniques as **t2v**. Lastly, we would like to highlight that the tensors that we have employed are super-symmetric and hence able to capture distribution over sets rather than sequence. But in general we can employ a k -way tensor which is not symmetric to even capture sequence. In this sense tensors are more general purpose.

3.2 Hyperedge2Vec Using Hasse De-noising Autoencoder

An autoencoder [9] takes an input vector $\mathbf{x} \in [0, 1]^n$ and maps it to a latent representation $\mathbf{y} \in [0, 1]^d$. This is typically done using an affine mapping followed by a non-linearity (more so when the input, like in our case, is binary [58]): $\mathbf{y} = f_{\theta}(\mathbf{x}) = \sigma(\mathbf{W}\mathbf{x} + \mathbf{b})$, with parameters $\theta = \{\mathbf{W}, \mathbf{b}\}$. Here, σ is a sigmoid function defined as $\sigma(x) = 1/(1 + e^{-x})$, \mathbf{W} is a $\mathbb{R}^{n \times d}$ weight matrix and \mathbf{b} is the offset. This latent representation is then used to reconstruct a vector $\mathbf{z} = g_{\theta'}(\mathbf{y}) = \sigma(\mathbf{W}'\mathbf{y} + \mathbf{b}')$, in the input space, $\mathbf{z} \in [0, 1]^n$ with parameters $\theta' = \{\mathbf{W}', \mathbf{b}'\}$. The mappings f_{θ} and $g_{\theta'}$ are referred to as the **encoder** and **decoder**, respectively. The representation \mathbf{y} is learned by minimizing the following reconstruction error:

$$\theta^*, \theta'^* = \arg \min_{\theta^*, \theta'^*} \frac{1}{m} \sum_{i=1}^m L(\mathbf{x}_i, \mathbf{z}_i) = \arg \min_{\theta^*, \theta'^*} \frac{1}{m} \sum_{i=1}^m L(\mathbf{x}_i, g_{\theta'}(f_{\theta}(\mathbf{x}_i))) \quad (13)$$

where L is a loss function, which in case of binary or bit probabilities is often chosen as the *cross-entropy loss*:

$$L(\mathbf{x}, \mathbf{z}) = \sum_{j=1}^n [\mathbf{x}(j) \log \mathbf{z}(j) + (1 - \mathbf{x}(j)) \log(1 - \mathbf{z}(j))] \quad (14)$$

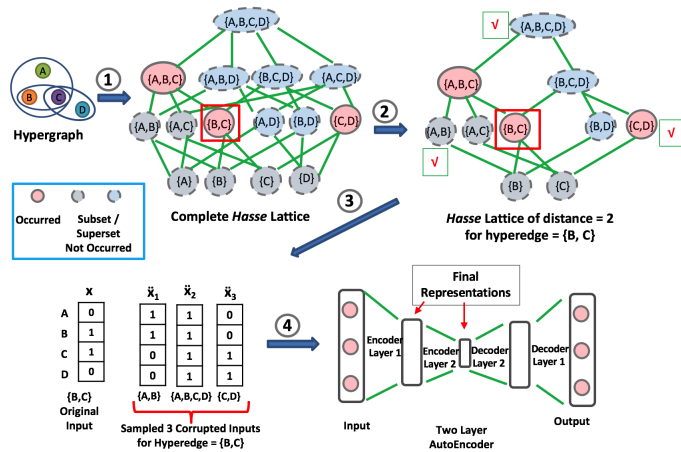


Figure 3: For the given hypergraph between four nodes (A,B,C,D) we consider the complete *hasse* lattice. For a given hyperedge {B,C} (square box) we then construct the sub-lattice made of hyperedges with distance $h = 2$ from {B,C}. We perform random walk (with $\tau = 0.2$) starting from the node corresponding to hyperedge {B,C} and sample $p = 3$ hyperedges (nodes visited by the random walk; shown with a check-mark ✓). Finally, we train the autoencoder to reconstruct the original hyperedge from these p noisy hyperedges.

In their paper, Vincent et al. [58] have shown that minimizing reconstruction amounts to maximizing the lower bound on the mutual information between input x and the representation y . However, they have further argued ([57]) that y retaining information about input x is insufficient. They further, propose the idea that the learned representation should be able to recover (*denoising*) the original input even after being trained with corrupted input (*adding noise*). They generate the corrupted input (\tilde{x}), using a stochastic mapping $q(\tilde{x}|x)$. Choice of noise is usually either *Gaussian* for real inputs and *Salt-and-pepper noise* for discrete inputs. The *denoising autoencoder* then learns the representation for each input, x , same as Equation 13, but with the following modified loss function: $L(x(i), g_{\theta}(f_{\theta}(\tilde{x}(i))))$.

We leverage the denoising autoencoder for learning representation for j^{th} hyperedge, by treating each hyperedge as an input, $x = H(j, \cdot)$. The size of this input vector for each hyperedge is n , which is the number of vertices in the hypergraph. In most natural hypergraphs, specially social networks, n can be quiet high ranging from thousands to millions or even billions (like Facebook for example). Therefore, randomly using a discrete noise like *salt-and-pepper*, might not be reasonable, as there are large number of possible permutations (as size n is large) and not all of them are related. Random addition of 1s or deletion of existing 1s from x , amounts to randomly adding or deleting vertices to the hyperedge corresponding to x . This might end up in new hyperedges that are completely unrelated to the given hyperedge (x). For example, users (nodes) in a social network from completely different regions of the network suddenly form a group (hyperedge). Such anomalous scenarios rarely happen in practice and social groups evolve in a gradual fashion via simple processes [45, 46].

Rather, we take advantage of the hypergraph structure to systematically guide us in generating this noise. A hypergraph can be defined by its corresponding *hasse lattice* [46]. For a given hyperedge (x), we consider the sub-lattice consisting of only those hyperedges that are at distance h from it in the complete lattice. On this sub-lattice we sample p hyperedges (nodes in sub-lattice) by performing random walk starting at the given hyperedge’s node (see Figure 3). We assume that all the nodes in the sub-lattice which correspond to previously occurred hyperedge g_i have weight $R(g_i)$ and rest all nodes have a constant weight τ . During random walk from a given node we choose a neighboring node (as the next node) in proportion to this neighboring node’s weight as compared to the other neighboring nodes. Our stochastic mapping $q(\tilde{x}|x)$ is therefore, a random walk on the sub-lattice of hyperedge (x) containing hyperedges at distance h from it. Intuitively, the hyperedges coming within a reasonable distance will affect each others representations and will have more similar representations. We will refer to the hyperedge representations learned by the above autoencoder technique, as **h2v-auto**.

4 EXPERIMENTS

4.1 Dataset Description

As the first dataset, we use group interaction log-data of the Sony’s Online multi-player game EverQuest II (**EQ II**) (everquest2.com) for a time period of nine months. In this game, several players work together as a team to perform various tasks. Each team earn points after completion of each task, and as the teams progress by earning points, they are escalated to different levels of game play. The interestingness of the game increases with each level. The points earned by the teams are treated as a measure of group performance. Each set of players who played together is treated as a hyperedge. We treat the number of times same set of players play together again as hyperedge occurrence number ($R(g_i)$). Players can participate in several teams over time, therefore, resulting in a hypergraph with overlapping hyperedges. We consider hyperedges of cardinality $\in [2, 6]$ as almost 90% of our hyperedges lie within this range. The resulting dataset contains a total of 5964 hyperedges (teams) among 6536 nodes (players).

Second dataset, is the fully labeled and publicly available sentiment analysis corpus of *Stanford Sentiment Treebank* (**LangNet**) [49]. This dataset is based on the reviews from a movie review website (rottentomatoes.com) and contains 215,154 unique phrases. Each of the phrases are labeled with a sentiment score (a real number $\in [0, 1]$, larger value indicates positive sentiment) by human annotators. Each phrase is a set or hyperedge of words. As there is no occurrence information for a phrase hyperedge we consider $R(g_i) = 1, \forall i \in \{1, \dots, m\}$. Given that words are shared across various phrases, these common words connect the phrase hyperedges, resulting in a phrase hypergraph with overlapping phrase hyperedges. Again, we only consider hyperedges of cardinality $\in [2, 6]$. After applying this cardinality filter we are left with 141,410 hyperedges (phrases) and 21,122 nodes (words).

4.2 Evaluation Methodology and Experimental Setup

4.2.1 Methods Compared. As mentioned before we refer to our proposed methods: tensor based *hypergraph tensor decomposition* and deep auto-encoder based *hypergraph auto-encoder*, as **t2v** and **h2v-auto**, respectively. We compare our proposed methods against six baselines each of which generates hyperedge as well as node embeddings. (1) **h2v-DM** (refer section A.1) (2) **h2v-DBOW** (refer section A.1) (3) **h2v-inv** (refer section A.2) (4) **h2v-dual** (refer section A.2) (5) **e2v** (refer section A.2) (6) **e2v-hyp** (refer section A.2). Methods (1-2) are adapted from sequence based language models for set structured hypergraph data. Methods (3-6) are various kind of dyadic graph based embedding methods adapted for hypergraph setting.

Except for **h2v-auto**, each of the baselines as well as **t2v** outputs both node as well as hyperedge embeddings of dimension $d = 128$. We further combine the node and hyperedge embedding using five different strategies: (i) node embedding summation (dimension $d = 128$), (ii) node embedding summation and concatenation with hyperedge embedding (dimension $2 \times d = 256$), (iii) node embedding averaging (dimension $d = 128$), (iv) node embedding averaging and concatenation with hyperedge embedding (dimension $2 \times d = 256$), and (v) only hyperedge embedding (dimension $d = 128$). **h2v-auto** only produces hyperedge embeddings of dimension $d = 128$. But it builds the embeddings using three different scenarios as mentioned in next section. Therefore, in total we have 38 ($= 35 + 3$) different scenarios each resulting in a different hyperedge embedding.

4.2.2 Evaluation Tasks and Setup. We perform two regression based tasks for the two datasets. In **EQII** dataset each team (hyperedge) has a team performance score associated with it. This team performance score is a real number, equal to the number of points earned by the team while performing one or more tasks within a gaming session. We treat the embedding learned for a given team (hyperedge) as its feature vector which is associated with a real number (team performance). We therefore, perform on regression over all the hyperedges (teams) with team performance as the dependent variable.

Similarly, in **LangNet** dataset each phrase (hyperedge) has a sentiment score associated with it, which again is a real number. Similar to the team dataset above, we treat the embedding learned for a given phrase (hyperedge) as its feature vector which is associated with a real number (sentiment score). We therefore, treat this as a regression task with sentiment score as the dependent variable and perform regression using the feature matrix containing embeddings of all the phrases.

For both the tasks we just described, we perform several evaluation runs. In each run we randomly choose 30% of hyperedges (teams or phrases) as the test set and learn ridge regression parameters using the remaining 70% training hyperedges for each of the 38 different embedding scenarios. Root mean squared error (RMSE) was chosen as the evaluation metric (the lower, the better). RMSE was calculated for each of the 38 scenarios and for each run. Final RMSE score was taken as the average RMSE score across five runs. Ridge regression's hyper-parameter was chosen by 5-fold cross-validation.

For **t2v** we found $\gamma = 0.1$ and $\eta = 0.2$ as the best parameters using grid search. For the auto-encoder method (**h2v-auto**) we consider three scenarios: (1) single hidden layer ($L1$) of $d = 128$; (2) two hidden layers ($L1$ & $L2$) with size of $L1 : d = 96$ and of $L2 :$

$d = 32$. We concatenate these embedding to get a single $d = 128$ size embedding; and (3) two hidden layers ($L1$ & $L2$) with size of $L1 : d = 512$ and of $L2 : d = 128$. We use the output of $L2$, which is of dimension $d = 128$, as the embedding. For sampling, we use the distance parameter $h = 2$ for generating the sub-lattice and $\tau = 0.2$ for both datasets. Also, $p = 10$ & $p = 5$ number of hyperedges are sampled (corresponding to each hyperedge) from **EQII** and **LangNet**, respectively.

4.3 Results and Discussion

Tables 1 & 2 show the RMSE scores of **t2v** (as compared with baselines) for the tasks of team performance prediction and sentiment score prediction, respectively. These tables contain scores for all the 35 different scenarios: columns represent 7 (6 baselines and the proposed *hypergraph tensor decomposition* (**t2v**)) different models while rows represent combination strategies. The scores for *hypergraph auto-encoder* (**h2v-auto**) are shown in the separate Table 4 as auto-encoder only generates hyperedge embeddings.

Accuracy & Run-times As we can observe that for our datasets and for both the regression tasks, almost all the embeddings (baselines and proposed) are performing very similarly in terms of accuracy. However, we can observe in Table 3 that tensor based method (**t2v**) have significant less time than the expensive auto-encoder technique (**h2v-auto**) (Table 4). Among the baselines **e2v-hyp** is the fastest. We also observe that sentence based techniques run faster on text-based **LangNet** dataset as compared to node2vec based methods which are not designed for text data and vice-versa.

Interpretability All the matrix or tensor based algebraic techniques: **e2v**, **e2v-hyp** and **t2v**, have a natural graph theoretic interpretation. **e2v** and **e2v-hyp** are both well studied spectral techniques with several eigen-value based interpretations. **t2v** has a hierarchical *hasse diagram* based interpretation. In contrast the sentence and node2vec based techniques can be understood only intuitively in terms of the cost function. Apart from the noise which has a random walk based interpretation, **h2v-auto** is a deep-learning based method which exploits multiple level of non-linearity and is over all a black-box approach.

Information Loss To reiterate one of the primary aims of this study is to design methods that retain the hyperedge-level joint information. The proposed tensor-based **t2v** principally capture the joint distribution over various cardinality hyperedges unlike conditional distribution like sentence embedding (**h2v-dm**) (which are more appropriate for sequences). In comparison the other method proposed **h2v-auto**, although is not directly designed to retain hyperedge-level joint distribution, but we hypothesize that the deep layered neural network should output highly informative non-linear representations. Several other baselines only retain pair-wise information. For example, the methods using node2vec (in Section A.2) are based on the skip-gram model, which learns embedding of nodes while maximizing the pair-wise conditional probability of a node given another node in a context. Similarly, the spectral methods (Section A.2) are inherently two dimensional as they are based on matrix. Same is the case with skip-gram based sentence embedding (**h2v-dbow**). However, sentence embedding based on the **DM** architecture (**h2v-dm**) maximizes the conditional probability of a word given the

Embed Combination	Baselines						Hypergraph
	Sentence Embed based		Node2Vec based		Spectral methods		Tensor Decomp.
	h2v-DM	h2v-DBOW	h2v-inv	h2v-dual	e2v	e2v-hyp	(t2v)
Node Embed Sum	0.79308	0.79567	0.80418	0.79956	0.81183	0.81405	0.81341
Node Embed Sum + Hyperedge Embed	0.79651	0.80241	0.81362	0.80636	0.81113	0.81652	0.81299
Node Embed Average	0.81584	0.81733	0.82407	0.82281	0.81234	0.81369	0.81303
Node Embed Avg + Hyperedge Embed	0.8182	0.82077	0.83378	0.82896	0.81223	0.81608	0.8127
Only Hyperedge Embed	0.81203	0.81522	0.82189	0.81984	0.81233	0.81608	0.81341

Table 1: RMSE Scores of (t2v) compared to baselines for EQ II Team Performance Analysis

Embed Combination	Baselines						Hypergraph
	Sentence Embed based		Node2Vec based		Spectral methods		Tensor Decomp.
	h2v-DM	h2v-DBOW	h2v-inv	h2v-dual	e2v	e2v-hyp	(t2v)
Node Embed Sum	0.14081	0.14029	N/A	N/A	0.14633	0.14854	0.14194
Node Embed Sum + Hyperedge Embed	0.14028	0.13883	N/A	N/A	0.14627	0.14845	0.14144
Node Embed Average	0.14245	0.14115	N/A	N/A	0.14665	0.14852	0.14381
Node Embed Avg + Hyperedge Embed	0.14178	0.14007	N/A	N/A	0.14661	0.14845	0.14333
Only Hyperedge Embed	0.14194	0.14147	N/A	N/A	0.14744	0.14844	0.1482

Table 2: RMSE Scores of (t2v) compared to baselines for LangNet Sentiment Analysis

Dataset	Baselines						Hypergraph
	Sentence Embed based		Node2Vec based		Spectral methods		Tensor Decomp.
	h2v-DM	h2v-DBOW	h2v-inv	h2v-dual	e2v	e2v-hyp	(t2v)
EQ2	455.84	103.47	90.05	93.41	128.03	12.01	213.37
LangNet	80.61	62.31	211.97*	207.86*	221.46	47.12	483.81

* these are average time taken for learning vertex embeddings only

Table 3: Average Runtime (seconds) of (t2v) compared to baselines across datasets

Layer Sizes	EQ II			LangNet		
	L1:128	L1:96/L2:32	L1:512/L2:128	L1:128	L1:96/L2:32	L1:512/L2:128
RMSE	0.81104	0.81512	0.81635	0.14568	0.14529	0.14784
Run Time	52 min	40 min	1 hr 20 min	2 hr 10 min	3 hr 20 min	6 hr

Table 4: RMSE Scores & Run-times of (h2v-auto)

previous set of words (*context*) without breaking this context (by concatenating or averaging the embedding of the previous words).

Leveraging Hypergraph Topology is another feature that we stress as a desirable property in a method which aims to build hyper-edge embedding. Hypergraph topology is an important auxiliary information, if left unused is wastage of resources at hand. We observe that except the sentence-based baselines, both proposed methods as well as other baselines, which are adapted for the hypergraph setting, leverage hypergraph topology using various hypergraph representations. Proposed **h2v-auto** generated noise using the *hasse lattice* and **t2v** directly models the hasse lattice by storing hyperedge level joint distribution along with hypergraph laplacian regularization. Node2vec based and spectral baselines use various kinds of matrices (that capture hypergraph topology up to varying degrees) as enlisted in (3-6) in Section 4.2.1.

Ability to generate hyperedge embeddings directly is another critical aim as was highlighted in the introduction. We emphasize that none of the baselines are designed to give hyperedge embeddings in a principled manner. However, both tensor based (**t2v** using $\mathcal{A}_{\text{dual}}^k$) as well as auto-encoder (**h2v-auto**) methods proposed in this paper are specifically designed to give hyperedge representations directly.

Choice of Method Both tensor based (**t2v** using $\mathcal{A}_{\text{dual}}^k$) as well as auto-encoder (**h2v-auto**) methods proposed in this paper are able to generate embeddings for all hyperedges of various cardinalities present in data. However, the auto-encoder method *cannot* generate

node embeddings unlike tensor method (using $\mathcal{A}_{\text{hyp}}^k$). On contrary, the main advantage of auto-encoders is the crisp embedding achieved via multiple levels of non-linearity offered by deep neural networks. Although this non-linearity improves the accuracy of auto-encoder methods very slightly (Table 4) as compared to tensor method (bottom right corner of (Table 1 & 2), the computation cost of tensor is far less than the auto-encoder method (Table 3 & 4). Intuitively, it seems that depending upon the task at hand, in some tasks the retention of joint hyperedge-level information is more important while modeling and data has lesser non-linear structure; or vice versa. But the low computational cost (as reflected in our tasks) makes tensor method definitely more lucrative. Together with the natural graph theoretic interpretation, tensor methods in comparison to black-box approach of deep learning based auto-encoder, makes them even more viable.

Highlight Another distinguishing feature of our work is that we leverage the existing structure (group structure) within the data. Recent attempts along these lines in network literature [27, 41, 52] argue that language models have ready-made context in the form of sentences or paragraphs to train the model, which are not available in networks, and therefore, they propose different ways to generate this context. In contrast, we focus on networks where the context is already present, e.g. collaboration networks where collaborative teams are hyperedges, or language hyper-networks where sentences are hyperedges.

4.4 Scalability

Both tensor based (**t2v** using $\mathcal{A}_{\text{dual}}^k$) as well as auto-encoder (**h2v-auto**) methods proposed in this paper are able to generate embeddings for *all hyperedges* of various cardinalities present in data.

Scalability for t2v The scalability issue in tensor based approach arise on two fronts: (1) the increase in number of hyperedges (m) or number of vertices (n) and (2) increase in parameter α (maximum cardinality hyperedge considered in the cost function) or β (maximum degree vertex considered in the cost function). Also there are two kinds of cost: (1) enumeration cost associated while building hypergraph and dual tensors and (2) cost for hypergraph tensor decomposition.

In case of hyperedge embeddings, for each vertex v of degree (i.e. number of hyperedges it is part of) $d(v)$ there are $d(v)!$ entries (permutations of hyperedges incident on v) to be enumerated to fill the tensor $\mathcal{A}_{\text{dual}}^{d(v)}$. This amounts to a worst case enumeration cost of $\sum_{v \in V} d(v)!$. However, $d(v)!$ can grow prohibitively large for vertices with large vertex degree $d(v)$. We therefore, propose to restrict ourselves to vertices of degree $d(v) \leq \alpha$. In fact we propose the following augmented initialization scheme for the dual tensors. For degree $k (\leq \alpha)$ vertex, $v_j \in \{g_{q_1}, g_{q_2}, \dots, g_{q_k}\}$ and $d(v_j) = k \leq \alpha, \forall j \in \{1, \dots, n\}$, we initialize the following elements of $\mathcal{A}_{\text{dual}}^k = (a_{q_1, q_2, \dots, q_k}) \in \mathbb{R}^{[k, m]}$, same as Equation 4. For degree $k (> \alpha)$ vertex, $v_j \in \{g_{q_1}, g_{q_2}, \dots, g_{q_k}\}, \forall j \in \{1, \dots, n\}$, we initialize

elements of $\mathcal{A}_{\text{dual}}^\alpha$ as,

$$\begin{aligned} a_{q_1, q_2, \dots, q_\alpha} &= 1 \\ a_{q_{(\alpha+1)}, q_{(\alpha+2)}, \dots, q_{2\alpha}} &= 1 \\ &\dots\dots \\ &\dots\dots \end{aligned} \quad (15)$$

$$a_{q_{(\lfloor \frac{k}{\alpha} \rfloor - 1)\alpha + 1}, q_{(\lfloor \frac{k}{\alpha} \rfloor - 1)\alpha + 2}, \dots, q_{(\lfloor \frac{k}{\alpha} \rfloor)\alpha}} = 1$$

and if $\delta = (k - (\lfloor \frac{k}{\alpha} \rfloor)\alpha) \neq 0$, then we also initialize the following element of $\mathcal{A}_{\text{dual}}^\delta$ as:

$$a_{q_{(\lfloor \frac{k}{\alpha} \rfloor)\alpha + 1}, q_{(\lfloor \frac{k}{\alpha} \rfloor)\alpha + 2}, \dots, q_k} = 1. \quad (16)$$

The above initialization basically partitions the k hyperedges incident on a vertex into $(\lfloor \frac{k}{\alpha} \rfloor + 1)$ (or $(\lfloor \frac{k}{\alpha} \rfloor)$, if α is a divisor of k) partitions. There are several techniques to obtain such a partition, but in this paper we simply sort the incident hyperedges by their cardinality and then sequentially make sets of α elements each (except the last partition which is of size $(k - (\lfloor \frac{k}{\alpha} \rfloor)\alpha)$ is not an integer) as shown in Equations(15- 16). This augmentation shall decrease the enumeration cost to: $\Delta_{\text{dual}} = \sum_{v \in V_a} d(v)! + \sum_{v \in V_b} \left(\lfloor \frac{d(v)}{\alpha} \rfloor \alpha! + (d(v) - \lfloor \frac{d(v)}{\alpha} \rfloor \alpha)! \right)$, where $V_a = \{v | v \in V, d(v) \leq \alpha\}$ and $V_b = (V - V_a)$. Depending upon the availability of resources (single machine or a distributed computing resources) we can choose the cut-off degree (α) accordingly. Notice that although we perform this degree thresholding we still get embedding for all the hyperedges. However, what we sacrifice is the higher order information that a given vertex connects a set of incident hyperedges jointly. But for higher order vertices we argue that as the vertex degree get too high our belief in the inference that hyperedges incident are related diminishes. These high degree vertices can be interpreted as highly popular person participating in a huge number of teams, but should the embeddings of these teams be related to each other is difficult to infer. In text data these high degree nodes correspond to words occurring in a large number of phrases, but this doesn't mean the sentences are related, rather its just that this word is quiet common.

In case of vertex embeddings, for each hyperedge $g_i \in G$, we need to initialize $|g_i|!$ elements of $\mathcal{A}_{\text{hyp}}^{|g_i|}$, with a worst case enumeration cost of $\sum_{g_i \in G} |g_i|!$. Given this cost can be prohibitive as the hyperedge cardinalities increase, we propose an augmented initialization scheme for hypergraph tensors where we shall restrict ourselves to hyperedges of cardinality $\leq \beta$. For hyperedges g_i with $|g_i| = k \leq \beta$ we initialize $\mathcal{A}_{\text{hyp}}^k = (a_{p_1, p_2, \dots, p_k}) \in \mathbb{R}^{[k, n]}$ same as Equation 2. For hyperedge with cardinality $|g_i| > \beta$, we initialize the elements of $\mathcal{A}_{\text{hyp}}^\beta \in \mathbb{R}^{[\beta, n]}$ as follows:

$$a_{p_1^l, p_2^l, \dots, p_\beta^l} = \frac{R(g_i)}{\gamma} \quad (17)$$

where $(p_1^l, p_2^l, \dots, p_\beta^l)$ is the l -th permutation among $l = \{1, 2, \dots, \gamma\}$, where $\gamma = \binom{|g_i|}{\beta}$. Enumeration cost now decreases to: $\Delta_{\text{hyp}} = \sum_{g_i \in G_a} |g_i|! + \sum_{g_i \in G_b} \binom{|g_i|}{\beta} \beta!$, where $G_a = \{g_i | g_i \in G, |g_i| \leq \beta\}$ and $G_b = (G - G_a)$. The choice of cut-off (β) is therefore, dictated by the computation resources available. Notice that although we

perform this cardinality thresholding we still get embedding for all the vertices. However, what we sacrifice is the higher order information that a given hyperedge connects a set of incident vertices jointly. Now, as the k increases, the number of k cardinality hyperedges decrease as well as the $\mathcal{A}_{\text{hyp}}^k$ tensor grows exponentially making it increasingly sparse and less informative (in information theoretic sense). Therefore, restricting to β cardinality hyperedges is a trade-off between enumeration cost and hyperedge level joint information loss. We observe that enumeration involved in both hyperedge and vertex embeddings, Equations(15- 17), are either *vertex or hyperedge centric computations* and can be performed in a scalable manner using the generic hyperedge-centric distributed computation libraries like MESH [29, 30] and HyperX [31].

Now we consider the tensor decomposition complexity and its scalability. For hyperedge embeddings, the per iteration complexity of the algorithm is $O(\sum_{\delta=1}^\alpha \text{nnz}(\mathcal{A}_{\text{dual}}^\delta))$. Notice, $(\sum_{\delta=1}^\alpha \text{nnz}(\mathcal{A}_{\text{dual}}^\delta)) \leq \Delta_{\text{dual}}$. Similarly, in case of vertex embeddings the per iteration complexity of the algorithm is $O(\sum_{k=c_{\min}}^\beta \text{nnz}(\mathcal{A}_{\text{hyp}}^k))$. Notice, $(\sum_{k=c_{\min}}^\beta \text{nnz}(\mathcal{A}_{\text{hyp}}^k)) \leq \Delta_{\text{hyp}}$. Hypergraph Tensor Decomposition basically involves learning vertex or hyperedge embeddings, i.e. parameters which are vertex or hyperedge centric. We can therefore, convert the tensor decompositions into equivalent hyperedge or vertex centric message passing algorithm and use the generic hyperedge or vertex centric distributed computation libraries like MESH or HyperX, as mentioned previously. Although, we highlight the possible directions for scalability, but we propose them as a future work.

Scalability without partitioning for t2v As discussed previously, as we are performing hyperedge cardinality as well as vertex degree thresholding and for that reason we developed augmented enumeration schemes in Equation (15- 17). The partitioning involved in this augmentation ensures that all the hyperedges and all the vertices are linked by some vertex of degree $\leq \alpha$ and some hyperedge of cardinality $\leq \beta$. If this augmentation was not performed, and we simply use the enumeration of Equations (2- 4), it was possible for example, that some hyperedge which has no vertex whose degree is $\leq \alpha$, and therefore, would be left out and no embedding would be learned for this hyperedge (because of cold start in tensor decomposition). Rather than performing the augmentation we can use semi-supervised learning to learn the embeddings for *critical hyperedges* (hyperedges with all vertices with degree $> \alpha$) using the embeddings of the *non-critical hyperedges* (hyperedges with at least one vertex with degree $\leq \alpha$). As all the hyperedges (critical or not) are linked using the hypergraph topology, we can perform semi-supervised learning using topology based regularization while performing tensor decomposition. This can be done using a graph regularization term $((_k \mathbf{U}^{(i)})^T \mathbf{L} (_k \mathbf{U}^{(i)}))$ in our tensor decomposition objective function ???. Here, \mathbf{L} is the graph laplacian of our choice. For example we can use \mathbf{L}_{hyp} or $\mathbf{L}_{\text{graph}}$ when we perform tensor decomposition for vertex embeddings and use \mathbf{L}_{inv} or \mathbf{L}_{dual} for hyperedge embeddings. Appendix A.2 mentions the details about these laplacians.

Scalability for h2v-auto In case of auto-encoder method the main scalability challenge lies in generating the noisy hyperedges. For a given hyperedge the intermediate sub-lattice from which p samples are drawn is of size $O(n^h)$, where h is max distance from

the hyperedge considered. Once this sampling is performed we have total mp hyperedges as input for the auto-encoder, which is linear in the number of hyperedge as p will be a constant. The main challenge therefore, is to generate the noisy samples where we have to tackle exponential size sub-lattice per hyperedge. Again, as this sampling is done for each hyperedge separately (hyperedge-centric), we can therefore, use distributed computing for hypergraphs using MESH and HyperX, as mentioned before.

Issues with aggregating vertex embedding Furthermore, we also observe that even simple element-wise summation or averaging of node embeddings for the nodes (in a given hyperedge) also perform comparably when compared to hyperedge embedding alone. From this we can infer that depending upon the dataset, if we have less hyperedges and more nodes, than we would rather prefer to simply learn the hyperedge embedding directly rather than learning node embeddings and then performing aggregating operation over them. Aggregation might turn out be costly specially if average hyperedge size is large and the choice of aggregation function is an issue. Therefore, learning hyperedge embeddings directly seems to be escape the problem of choosing the aggregation function all together.

Scalability issues for baselines Another thing we notice, is that in case of **LangNet** dataset, that node2vec based **h2v-inv** and **h2v-dual** methods are simply unable to run (see N/A in Table 2). It seems that in **LangNet** (and possibly in text data in general) the hypergraph dual, A_{dual} , which contain phrase to phrase edges turns out be containing significantly more edges than the number of node to node edges in A_{hyp} . Therefore, performing context generation (using node2vec) over A_{inv} and A_{dual} graphs turns out be very costly and we were unable to get hyperedge embeddings for **h2v-inv** and **h2v-dual** methods. Note, we however do get vertex embedding (as mentioned with a * mark below Table 3). But as you can see that even the time taken for vertex embedding is similar or more than the total time taken by spectral methods for learning both hyperedge as well as vertex embedding, together. This indicates the robustness of the spectral methods, specially the **e2v-hyp** for different kinds of hypergraph data and edge densities (sparsity).

5 RELATED WORKS

Hypergraphs were studied rigorously by [11, 12] as a generalization of graphs and *directed* hypergraphs have been introduced by [13]. Hypergraph were argued for the first time as a model to naturally capture higher-order relationships between sets of objects across variety of domains by [23]. Hypergraphs have been used to model complex networks in different fields including biology [33], databases [24] and data mining [28, 62]. Within machine learning, algorithms guided by the structure of hypergraph were introduced by [61] and have found applications in a variety of domains [25, 37, 45, 46, 56]. Simplicial complex [40] based view of hypergraph using *hasse lattice* [48] within machine learning has recently been proposed by [46].

Representation learning (RL) *Node Representations in Graph:* Traditionally unsupervised node embedding learning has been done using latent models like matrix factorization [3] or by community detection [53] based techniques for networks [8, 21, 44, 54]. In each case there is a vector of features learned for a node, each of whose

entries reflects node's association with some latent dimension or a network community. More recently, there has been a revived interest in graph embedding in form of *context* oriented techniques [27, 41, 52]. These techniques are inspired by recent unsupervised RL methods in NLP [36, 39] where word embeddings are learned that are similar to words in a given neighborhood or *context*. These techniques differ in the manner they generate this context as well as in the objective which they optimize. Also there are supervised algorithms learn embeddings which are optimal for the specific task at hand. This results in high accuracy but incurs significant computational cost for training. Recently, several supervised learning algorithms have been proposed for network analysis [55, 60] and for text networks in a semi-supervised setting [51]. Finally, we refer readers to a very recent and comprehensive survey on graph embedding methods by [15].

Node Representations in Hypergraph: Learning embeddings for nodes within a hypergraph while incorporating the hypergraph topology using proxy graphs is introduced by [61]. Using graph proxy destroys the hyperedge-level joint information and thus, incur loss of information. Also, [2] squarely criticize that such representations can be learned by constructing graphs, which are proxies for the hypergraph structure.

Set Representations: RL for sets using neural networks has been proposed recently [59], where a memory network is used to compose features sequentially but in an order invariant manner. In their very recent paper, [43] have tried to answer this set ordering issue by the use of *random set theory*. However, they do not consider embedding but focus on learning set-level probabilities. More importantly, both of these works, do not consider the hypergraph structure of overlapping sets which is the main focus of this paper.

Tensors For comprehensive view of tensors, tensor decomposition as well as applications we refer to the survey by [35]. The connection between k -way tensor and k -uniform hypergraph eigen values was established by [42]. The use of k -way symmetric tensor and their non-negative decomposition for uniform hypergraph partitioning was first introduced by [47]. But they are again restricted to uniform hypergraphs.

6 CONCLUSION

In this paper we have proposed two methods to generate higher-order representations for both hyperedges (representing sets of nodes) and hypergraph nodes (that also take into account the hypergraph structure). First, is an auto-encoder based deep-learning method and second, is a tensor-based algebraic method. Both learning models are unique in the manner they leverage the existing structure present in network data as context. While introducing a new idea of a dual tensors corresponding to the hypergraph dual, we develop a novel approach of using factors from joint decomposition of k -way tensors corresponding to k -uniform sub-hypergraphs, as generic node & hyperedge representations. We show that that both methods perform comparably with several other baselines in terms of accuracy. We also observe that the proposed tensor based methods are more efficient and also have a natural hypergraph theoretic interpretation; unlike deep learning based black-box approach. We therefore, argue that we have proposed more general methods which are principally suited for hypergraphs (and therefore also for graphs) while maintaining accuracy and efficiency.

7 ACKNOWLEDGEMENTS

This work has been supported in part by the NSF Award IIS-1422802.

REFERENCES

- [1] PubMed Data. <https://www.ncbi.nlm.nih.gov/pubmed/>. (????).
- [2] Sameer Agarwal, Kristin Branson, and Serge Belongie. 2006. Higher order learning with graphs. In *Proceedings of the 23rd international conference on Machine learning*. ACM, 17–24.
- [3] Amr Ahmed, Nino Shervashidze, Shrawan Narayanamurthy, Vanja Josifovski, and Alexander J Smola. 2013. Distributed large-scale natural graph factorization. In *Proceedings of the 22nd international conference on World Wide Web*. ACM, 37–48.
- [4] Iftexhar Ahmed, Channing Brown, Andrew Pilny, Dora Cai, Yannick Atouba Ada, and Marshall Scott Poole. 2011. Identification of groups in online environments: The twist and turns of grouping groups. In *Privacy, Security, Risk and Trust (PASSAT) and 2011 IEEE Third International Conference on Social Computing (SocialCom)*, 2011 IEEE Third International Conference on. IEEE, 629–632.
- [5] Brett W Bader and Tamara G Kolda. 2007. Efficient MATLAB computations with sparse and factored tensors. *SIAM Journal on Scientific Computing* 30, 1 (2007), 205–231.
- [6] Grey Ballard, Tamara Kolda, and Todd Plantenga. 2011. Efficiently computing tensor eigenvalues on a GPU. In *Parallel and Distributed Processing Workshops and Phd Forum (IPDPSW)*, 2011 IEEE International Symposium on. IEEE, 1340–1348.
- [7] Anthony F Bartholomay. 1960. Molecular set theory: A mathematical representation for chemical reaction mechanisms. *Bulletin of Mathematical Biology* 22, 3 (1960), 285–307.
- [8] Mikhail Belkin and Partha Niyogi. 2001. Laplacian eigenmaps and spectral techniques for embedding and clustering. In *NIPS*, Vol. 14. 585–591.
- [9] Yoshua Bengio and others. 2009. Learning deep architectures for AI. *Foundations and trends® in Machine Learning* 2, 1 (2009), 1–127.
- [10] Yoshua Bengio and Samy Bengio. 2000. Modeling high-dimensional discrete data with multi-layer neural networks. In *Advances in Neural Information Processing Systems*. 400–406.
- [11] C. Berge. 1976. *Graphs and hypergraphs*. Vol. 6. Elsevier.
- [12] Claude Berge. 1984. *Hypergraphs: combinatorics of finite sets*. Vol. 45. Elsevier.
- [13] Alain Bretto. 2013. Hypergraph theory. *An introduction. Mathematical Engineering. Cham: Springer* (2013).
- [14] Samuel R Bulò and Marcello Pelillo. 2009. A game-theoretic approach to hypergraph clustering. In *Advances in neural information processing systems*. 1571–1579.
- [15] Hongyun Cai, Vincent W Zheng, and Kevin Chen-Chuan Chang. 2017. A Comprehensive Survey of Graph Embedding: Problems, Techniques and Applications. *arXiv preprint arXiv:1709.07604* (2017).
- [16] Eric C Chi and Tamara G Kolda. 2012. On tensors, sparsity, and nonnegative factorizations. *SIAM J. Matrix Anal. Appl.* 33, 4 (2012), 1272–1299.
- [17] Evangelia Christakopoulou and George Karypis. 2014. HOSLIM: higher-order sparse linear method for top-n recommender systems. In *Pacific-Asia Conference on Knowledge Discovery and Data Mining*. Springer, 38–49.
- [18] Fan RK Chung, Ronald L Graham, Peter Frankl, and James B Shearer. 1986. Some intersection theorems for ordered sets and graphs. *Journal of Combinatorial Theory, Series A* 43, 1 (1986), 23–37.
- [19] Pierre Comon, Gene Golub, Lek-Heng Lim, and Bernard Mourrain. 2008. Symmetric tensors and symmetric tensor rank. *SIAM J. Matrix Anal. Appl.* 30, 3 (2008), 1254–1279.
- [20] Joshua Cooper and Aaron Dutle. 2012. Spectra of uniform hypergraphs. *Linear Algebra Appl.* 436, 9 (2012), 3268–3292.
- [21] Trevor F Cox and Michael AA Cox. 2000. *Multidimensional scaling*. CRC press.
- [22] Mukund Deshpande and George Karypis. 2004. Item-based top-n recommendation algorithms. *ACM Transactions on Information Systems (TOIS)* 22, 1 (2004), 143–177.
- [23] Ernesto Estrada and Juan A Rodriguez-Velazquez. 2005. Complex networks as hypergraphs. *arXiv preprint physics/0505137* (2005).
- [24] R. Fagin. 1983. Degrees of acyclicity for hypergraphs and relational database schemes. *Journal of the ACM (JACM)* 30, 3 (1983), 514–550.
- [25] Shenghua Gao, Ivor Wai-Hung Tsang, and Liang-Tien Chia. 2013. Laplacian sparse coding, hypergraph laplacian sparse coding, and applications. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 35, 1 (2013), 92–104.
- [26] Debarghya Ghoshdastidar and Ambedkar Dukkipati. 2015. A provable generalized tensor spectral method for uniform hypergraph partitioning. In *International Conference on Machine Learning*. 400–409.
- [27] Aditya Grover and Jure Leskovec. 2016. node2vec: Scalable Feature Learning for Networks. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, San Francisco, CA, USA, August 13-17, 2016*. 855–864. DOI : <http://dx.doi.org/10.1145/2939672.2939754>
- [28] E.H. Han, G. Karypis, V. Kumar, and B. Mobasher. 1997. *Clustering based on association rule hypergraphs*. University of Minnesota, Department of Computer Science.
- [29] Benjamin Heintz and Abhishek Chandra. 2015. Enabling Scalable Social Group Analytics via Hypergraph Analysis Systems. In *HotCloud*.
- [30] Benjamin Heintz, Shivangi Singh, Rankyung Hong, Guarav Khandelwal, Corey Tesdahl, and Abhishek Chandra. 2017. MESH: A Flexible Distributed Hypergraph Processing System. (2017).
- [31] Jin Huang, Rui Zhang, and Jeffrey Xu Yu. 2015. Scalable hypergraph learning and processing. In *Data Mining (ICDM), 2015 IEEE International Conference on. IEEE*, 775–780.
- [32] TaeHyun Hwang, Ze Tian, Rui Kuangy, and Jean-Pierre Kocher. 2008. Learning on weighted hypergraphs to integrate protein interactions and gene expressions for cancer outcome prediction. In *Data Mining, 2008. ICDM'08. Eighth IEEE International Conference on. IEEE*, 293–302.
- [33] S. Klamt, U.U. Haus, and F. Theis. 2009. Hypergraphs and cellular networks. *PLoS computational biology* 5, 5 (2009), e1000385.
- [34] Tamara G Kolda. 2015. Numerical optimization for symmetric tensor decomposition. *Mathematical Programming* 151, 1 (2015), 225–248.
- [35] Tamara G Kolda and Brett W Bader. 2009. Tensor decompositions and applications. *SIAM review* 51, 3 (2009), 455–500.
- [36] Quoc V Le and Tomas Mikolov. 2014. Distributed Representations of Sentences and Documents. In *ICML*, Vol. 14. 1188–1196.
- [37] Lei Li and Tao Li. 2013. News recommendation via hypergraph learning: encapsulation of user behavior and news content. In *Proceedings of the sixth ACM international conference on Web search and data mining*. ACM, 305–314.
- [38] Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781* (2013).
- [39] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013. Distributed Representations of Words and Phrases and Their Compositionality. In *Proceedings of the 26th International Conference on Neural Information Processing Systems (NIPS'13)*. 3111–3119.
- [40] James R Munkres. 1984. *Elements of algebraic topology*. Vol. 2. Addison-Wesley Menlo Park.
- [41] Bryan Perozzi, Rami Al-Rfou, and Steven Skiena. 2014. Deepwalk: Online learning of social representations. In *Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining*. ACM, 701–710.
- [42] Liqun Qi. 2005. Eigenvalues of a real supersymmetric tensor. *Journal of Symbolic Computation* 40, 6 (2005), 1302–1324.
- [43] Seyed Hamid Rezaatofighi, Anton Milan, Ehsan Abbasnejad, Anthony Dick, Ian Reid, and others. 2016. DeepSetNet: Predicting Sets with Deep Neural Networks. *arXiv preprint arXiv:1611.08998* (2016).
- [44] Sam T Roweis and Lawrence K Saul. 2000. Nonlinear dimensionality reduction by locally linear embedding. *Science* 290, 5500 (2000), 2323–2326.
- [45] Ankit Sharma, Rui Kuang, Jaideep Srivastava, Xiaodong Feng, and Kartik Singhal. 2015. Predicting small group accretion in social networks: A topology based incremental approach. In *2015 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining (ASONAM)*. IEEE, 408–415.
- [46] Ankit Sharma, Terrence J Moore, Ananthram Swami, and Jaideep Srivastava. 2017. Weighted Simplicial Complex: A Novel Approach for Predicting Small Group Evolution. In *Pacific-Asia Conference on Knowledge Discovery and Data Mining*. Springer, 511–523.
- [47] Amnon Shashua, Ron Zass, and Tamir Hazan. 2006. Multi-way clustering using super-symmetric non-negative tensor factorization. *Computer Vision—ECCV 2006* (2006), 595–608.
- [48] Steven Skiena. 1990. Hasse Diagrams. *Implementing Discrete Mathematics: Combinatorics and Graph Theory With Mathematica* (1990), 163.
- [49] Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher Manning, Andrew Ng, and Christopher Potts. 2013. Parsing With Compositional Vector Grammars. In *EMNLP*.
- [50] Karthik Subbian, Charu Aggarwal, and Jaideep Srivastava. 2013. Content-centric flow mining for influence analysis in social streams. In *Proceedings of the 22nd ACM international conference on Conference on information & knowledge management*. ACM, 841–846.
- [51] Jian Tang, Meng Qu, and Qiaozhu Mei. 2015. 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*. ACM, 1165–1174.
- [52] Jian Tang, Meng Qu, Mingzhe Wang, Ming Zhang, Jun Yan, and Qiaozhu Mei. 2015. Line: Large-scale information network embedding. In *Proceedings of the 24th International Conference on World Wide Web*. ACM, 1067–1077.
- [53] Lei Tang and Huan Liu. 2011. Leveraging social media networks for classification. *Data Mining and Knowledge Discovery* 23, 3 (2011), 447–478.
- [54] Joshua B Tenenbaum, Vin De Silva, and John C Langford. 2000. A global geometric framework for nonlinear dimensionality reduction. *science* 290, 5500

- (2000), 2319–2323.
- [55] Fei Tian, Bin Gao, Qing Cui, Enhong Chen, and Tie-Yan Liu. 2014. Learning Deep Representations for Graph Clustering. In *AAAI*. 1293–1299.
- [56] Ze Tian, TaeHyun Hwang, and Rui Kuang. 2009. A hypergraph-based learning algorithm for classifying gene expression and arrayCGH data with prior knowledge. *Bioinformatics* 25, 21 (2009), 2831–2838.
- [57] Pascal Vincent, Hugo Larochelle, Yoshua Bengio, and Pierre-Antoine Manzagol. 2008. Extracting and composing robust features with denoising autoencoders. In *Proceedings of the 25th international conference on Machine learning*. ACM, 1096–1103.
- [58] Pascal Vincent, Hugo Larochelle, Isabelle Lajoie, Yoshua Bengio, and Pierre-Antoine Manzagol. 2010. Stacked denoising autoencoders: Learning useful representations in a deep network with a local denoising criterion. *Journal of Machine Learning Research* 11, Dec (2010), 3371–3408.
- [59] Oriol Vinyals, Samy Bengio, and Manjunath Kudlur. 2016. Order Matters: Sequence to sequence for sets. In *ICLR*.
- [60] Li Xiaoyi, Li Hui Du Nan, and others. 2013. A deep learning approach to link prediction in dynamic networks. In *Proceedings of the 2013 SIAM International Conference on Data Mining, Philadelphia, PA, USA: SIAM*.
- [61] Dengyong Zhou, Jiayuan Huang, and Bernhard Schölkopf. 2006. Learning with hypergraphs: Clustering, classification, and embedding. In *Advances in neural information processing systems*. 1601–1608.
- [62] D. Zhou, J. Huang, and B. Schölkopf. 2007. Learning with hypergraphs: Clustering, classification, and embedding. *Advances in Neural Information Processing Systems* 19 (2007), 1601.

A APPENDIX

A.1 Hyperedge2Vec Using Sentence Embeddings

Recently, [36] proposed two representation learning methods for sentences: **DM** model and **DBOW** model. Both the methods take sentences as input and return embeddings for words as well as the sentence. To apply these sentence embedding models directly to hyperedges we generate a *proxy* sentence for each hyperedge $g_i \in G$ as a sequence made by concatenating all the permutations of the nodes (as words) in the hyperedge and further repeating this sequence as many time this hyperedge occurred. For example, for a three node hyperedge $g_i = \{1, 4, 7\}$ which has occurred two times ($R(g_i) = 2$) we make the following proxy sentence: $\{1, 4, 7, 1, 7, 4, 7, 1, 4, 7, 4, 1, 4, 1, 7, 4, 7, 1, 1, 4, 7, 1, 7, 4, 7, 1, 4, 7, 4, 1, 4, 1, 7, 4, 7, 1\}$ of length 12 (6 permutations times 2 occurrence). We refer to the node and hyperedge embedding pairs resulting from **DM**, **DBOW** models as **h2v-dm** and **h2v-dbow**, respectively.

A.2 Hyperedge2vec Using Spectral Embeddings and node2vec

These set of methods extract embeddings as the eigenvectors associated with Laplacian matrices corresponding the following four adjacency matrices. (1) Adjacency matrix associated with a hypergraph [61], which is defined as: $A_{\text{hyp}} = (\mathbf{H}^T \mathbf{W}_e \mathbf{H} - \mathbf{D}_v)$ and laplacian: $L_{\text{hyp}} = (\mathbf{I} - \mathbf{D}_v^{-1/2} \mathbf{A}_{\text{hyp}} \mathbf{D}_v^{-1/2})$ where \mathbf{W}_e ($\mathbf{W}_e(i, i) = R(g_i)$) is a diagonal matrix containing the weights of each hyperedge and \mathbf{D}_v is a diagonal matrix containing the degree of each vertex (2) $\mathbf{A}_{\text{graph}} = \mathbf{H}^T \mathbf{W}_e \mathbf{H}$ is the weighted graph associated with the hypergraph (\mathbf{A}_{hyp}) and its laplacian: $L_{\text{graph}} = (\mathbf{I} - \mathbf{D}_v^{-1/2} \mathbf{A}_{\text{graph}} \mathbf{D}_v^{-1/2})$ (3) we consider the following adjacency matrix $\mathbf{A}_{\text{inv}} = \mathbf{H} \mathbf{H}^T$ associated with what we refer to as the *inverted hypergraph*. This inverted hypergraph is a graph (unlike the dual which is a hypergraph) and there is an edge between two nodes if the hyperedges corresponding to the nodes in the original hypergraph have nodes in common. Weight of this edge is the number of common nodes. Its laplacian is:

$L_{\text{hyp}} = (\mathbf{I} - \mathbf{D}_v^{-1/2} \mathbf{A}_{\text{hyp}} \mathbf{D}_v^{-1/2})$ (4) The adjacency matrix associated with the hypergraph dual is: $\mathbf{A}_{\text{dual}} = (\mathbf{H}_{\text{dual}}^T \mathbf{W}_v \mathbf{H}_{\text{dual}} - \mathbf{D}_e) = (\mathbf{H} \mathbf{W}_v \mathbf{H}^T - \mathbf{D}_e)$ where \mathbf{W}_v is a diagonal matrix containing the weights of each node and \mathbf{D}_e is a diagonal matrix containing the degree of each hyperedge. We assume no weights on the nodes and take $\mathbf{W}_v = \mathbf{I}$. Its laplacian is: $L_{\text{dual}} = (\mathbf{I} - \mathbf{D}_e^{-1/2} \mathbf{A}_{\text{dual}} \mathbf{D}_e^{-1/2})$. We get vertex embeddings using of (1) and hyperedge embedding using (3) via eigen-decomposition, and we refer to them together as **e2v**. Similarly, we get another pair of vertex embeddings using (2) and hyperedge embedding using (4) which we refer as **e2v-hyp**.

[27] propose a representation method for nodes in a graph called *node2vec*, which uses the skip-gram model [39] to obtain node embeddings for an input graph. We get vertex embeddings using of (1) and hyperedge embedding using (3) via *node2vec*, and we refer to them together as **h2v-hyp**. Similarly, we get another pair of vertex embeddings using (2) and hyperedge embedding using (4), via *node2vec* which we refer as **h2v-dual**.