

Graph Representation Learning

Notes¹

Winter 2021

¹ Authors: Eeshaan Jain
Notes of the book Graph Representation Learning by William L. Hamilton

Keyphrases: Graph Representation Learning. Graph Statistics. Graph Kernel Methods. Graph Embeddings. Graph Neural Networks. Graph Generative Models.

This set of notes begin by introducing graphs, and learning tasks on it. Then we discuss the traditional approaches to infer graphs at the node-level and graph-level and further see inter-node relations. This extends to Spectral Graph Theory and use of Graph Laplacians to perform optimal clustering. An overview of shallow embeddings is done using matrix factorization and random walks on simple graphs, whose idea is extended to multi-relational graphs.[Introduction will be updated as file is updated.]

Contents

1	Introduction to Graphs and Learning on Graphs	2
1.1	Basic definitions	2
1.2	Multi-relational Graphs	2
1.3	Information representation	3
1.4	Node Classification	3
1.5	Link Prediction	4
1.6	Community Detection	4
1.7	Other popular graphical tasks	4
2	Traditional Approaches to Graph Learning	5
2.1	Node-level statistics	5
2.2	Graph Kernel Methods	6
2.3	Measures of Neighborhood Overlap	7
2.4	Graph Laplacians	8
2.5	Clustering using Laplacians	9
2.6	Spectral Clustering	9
2.7	Learning Representations	10
3	Learning Node Embeddings	11
3.1	Encoder-Decoder Model	11
3.2	Drawbacks of Shallow Embeddings	13
3.3	Reconstruction of Multi-Relational Data	13
3.4	Representation Abilities of Multi-Relational Decoders	15

1 Introduction to Graphs and Learning on Graphs

1.1 Basic definitions

Definition 1.1 (Graph)

A graph $G = (V, E)$ is defined by a set of nodes V and a set of edges between these nodes. We denote an edge going from $u \in V$ to $v \in V$ as $(u, v) \in E$.

Note that definition 1.1 is general, and most of the times we are only concerned with *simple graphs*, which have the following properties:

- ◇ For any two nodes in the set V , there exists at most one edge between them.
- ◇ There are no self-loops, i.e for any $u \in V$, $(u, u) \notin E$.
- ◇ All edges of the graph are *undirected*, i.e for $u, v \in V$, $(u, v) \in E \Leftrightarrow (v, u) \in E$.

Definition 1.2 (Adjacency Matrix)

For a simple graph $G = (V, E)$, the adjacency matrix $\mathbf{A} \in \mathbb{R}^{|V| \times |V|}$ is a square matrix with $A_{uv} = 1$ if there is an edge from vertex u to v , and $A_{uv} = 0$ otherwise.

Since we are dealing with simple graphs, some immediate consequences are that

- ◇ \mathbf{A} is symmetric, i.e $\mathbf{A}^T = \mathbf{A}$.
- ◇ $A_{uu} = 0 \forall u \in V$, i.e all diagonal elements are zero.

Notice that since we are defining the graph in terms of the matrix \mathbf{A} , we need to order the nodes prior to the representation (See Fig 1). An extension to definition 1.2 occurs when we allow the use of weighted edges, and in that case for $u, v \in V$, $A_{uv} \in [0, 1]$.

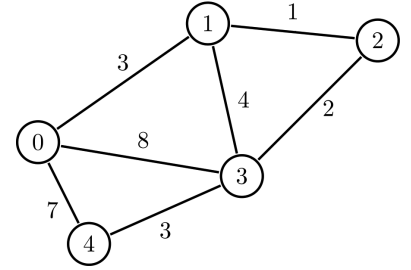


Figure 1: Undirected graph with ordered edges and nodes

$$\mathbf{A} = \begin{bmatrix} 0 & 1 & 0 & 1 & 1 \\ 1 & 0 & 1 & 1 & 0 \\ 0 & 1 & 0 & 1 & 0 \\ 1 & 1 & 1 & 0 & 1 \\ 1 & 0 & 0 & 1 & 0 \end{bmatrix}$$

The adjacency matrix for Fig 1.

1.2 Multi-relational Graphs

We have noted three types of edges till now: undirected, directed and weighted. But for graphs, it is possible to have different types of edges. In such cases, we extend the notation to include relation type $\tau \in \mathcal{R}$ and write that $(u, \tau, v) \in E$ along with defining an adjacency matrix \mathbf{A}_τ for each type. The above class of graphs are called multi-relational, and are overall represented by an adjacency tensor $\mathcal{A} \in \mathbb{R}^{|V| \times |\mathcal{R}| \times |V|}$. Two major subsets of such graphs are:

1. **Heterogeneous graphs:** We can partition the set of nodes into k disjoint sets $V = \bigcup_{i=1}^k V_i$, and thus we have types imbued on

nodes. Usually in such graphs, we pose constraints on edges, for example they connect nodes of different types only (*multipartite graphs*).

2. **Multiplex graphs:** We can decompose such graphs into a set of layers (see Fig 2), and each node is in every layer with each layer possessing a unique relation. Edges will be present intra-layer and inter-layer.

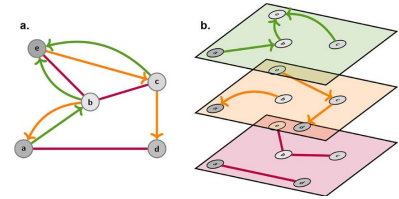


Figure 2: Layered multiplex graph. Source: [Gaudelet et al., 2021]

1.3 Information representation

Usually we have features representing graph-level features - may it be on nodes, edges or the entire graph itself. The most common one is node-level, and we associate an m -dimensional feature vector with each node, thus overall having a matrix $\mathbf{X} \in \mathbb{R}^{|V| \times m}$.

Now, let's go over some common graphical learning tasks.

1.4 Node Classification

The goal of node classification is to predict a label y_u associated with all nodes $u \in V$ when we have been only given the true labels of a subset of the nodes $V_{train} \subset V$. A few examples in literature of the task are:

1. Interactome protein function classification: [Hamilton et al., 2017a]
2. Document classification based on citation graphs: [Kipf and Welling, 2016]

Though it seems similar to supervised classification, there is a difficulty that the nodes in a graph are not *i.i.d.* and thus we need to model dependencies between data points. This is exactly what node classification takes a benefit of through

1. **Homophily** [Mcpherson et al., 2001]: Tendency of nodes to share attributes with their neighbors
2. **Structural Equivalence** [Donnat et al., 2017]: Idea that nodes with similar local neighborhood structures will have similar labels
3. **Heterophily**: Presumption that nodes will be preferentially connected to differently-labeled nodes

Note that, the only part of the network we don't access is the labels of the test nodes. We still have access to the structural information, and thus our task is **semi-supervised**.

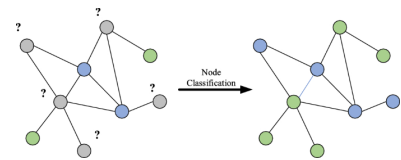


Figure 3: Example of node classification. Source: [Xiao et al., 2022]

1.5 Link Prediction

This task is related to the inference of edges between nodes in a graph. Given a set of nodes V and an incomplete set of edges $E_{train} \subset E$, we want to infer the set $E \setminus E_{train}$. A few examples in literature of the task are:

1. Content recommendation on social media: [Ying et al., 2018]
2. Modeling side-effects of drugs: [Zitnik et al., 2018]

Although for simple graphs, we can have simple heuristics to give decent results [Lü and Zhou, 2011], the task on multi-relational graphs is much more complex.

Much like node classification, the task boundary is smudged, and it falls under both supervised and unsupervised.

1.6 Community Detection

This task is the graphical analog of unsupervised learning. We want to infer latent community structures given only the input graph $G = (V, E)$. A few examples in literature of the task are:

1. Uncovering functional modules in genetic interaction networks: [Agrawal et al., 2017]
2. Fraud detection in financial transaction networks: [Pandit et al., 2007]

1.7 Other popular graphical tasks

- ◇ **Graph Regression / Graph Classification:** In these tasks, we want to learn over graph data, and our dataset consists of multiple different graphs and we predict graph-wise. We can treat each graph as i.i.d and thus, the only challenge is to define useful features over graphs.
- ◇ **Graph Clustering** In this task, we learn an unsupervised measure of similarity between pairs of graphs.

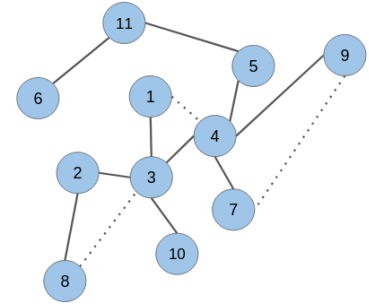


Figure 4: Example of node classification. Source: [AnalyticsVidhya](#)

2 Traditional Approaches to Graph Learning

2.1 Node-level statistics

Prior to advent of deep learning on graphs, a large part of machine learning on graphs was directed towards extraction of statistics/features based on heuristic functions and domain knowledge from graphs, and passing this information to a machine learning classifier. Now, we look over common node-level statistics.

Definition 2.1 (Node degree)

The degree d_u of a node $u \in V$ is the number of edges incident to a node, and in terms of \mathbf{A} for a simple graph

$$d_u = \sum_{v \in V} A_{uv} \quad (1)$$

If our graph is directed or weighted, we can have the notion of in-degree and out-degree by summing over rows or columns in Equation (1).

Definition 2.2 (Eigenvector Centrality)

The eigenvector centrality e_u of a node $u \in V$ is defined by the recurrence relation

$$e_u = \frac{1}{\lambda} \sum_{v \in V} A_{uv} e_v \quad \forall u \in V \quad (2)$$

We can rewrite the above equation to notice that $\mathbf{A}\mathbf{e} = \lambda\mathbf{e}$ where \mathbf{e} is the vector of node centralities.

Perron-Frobenius Theorem: A real square matrix with positive entries has a unique largest real eigenvalue, and the corresponding eigenvector can be chosen to have strictly positive components.

Note that, we want the centrality values to be positive, and by the **Perron-Frobenius Theorem**. In all, \mathbf{e} is given by the eigenvector of the largest eigenvalue of \mathbf{A} . Now, since λ is the leading eigenvalue, e_u ranks the likelihood that a node is visited on a random-walk of infinite length on the graph as we can use power iteration method ($\mathbf{e}^{(t)} = \mathbf{A}\mathbf{e}^{(t-1)}$) to find \mathbf{e} .

Note that the matrix $\mathbf{B}^{(n)} = \mathbf{A}^n$ has the elements such that $\mathbf{B}_{uv}^{(n)}$ denotes the number of n -length paths from $u \rightarrow v$.

By centrality, we can also measure how often a node lies on the shortest path between two nodes (*betweenness centrality*), or average shortest path length between a node and all other nodes (*closeness centrality*).

Definition 2.3 (Clustering Coefficient)

The clustering coefficient [Watts et al., 2006] is given as

$$c_u = \frac{|(v_1, v_2) \in E : v_1, v_2 \in \mathcal{N}(u)|}{\binom{d_u}{2}} \quad (3)$$

It essentially measures the proportion of closed triangles in a node's local neighborhood.

We can view clustering coefficient as the ratio between actual number

of triangles and total possible number of triangles within a node's *ego graph*: subgraph having the node, neighbors and all edges between nodes in its neighborhood.

2.2 Graph Kernel Methods

These methods are used for generating graph-level features for tasks such as graph classification.

Definition 2.4 (Bag of Nodes)

Bag of nodes refers to aggregates of node-level statistics such as histograms based on clustering coefficients, and the use of such aggregates as graph-level features.

Clearly bag of nodes misses global graph properties, and one way to improve it is to use iterative neighborhood aggregation.

Weisfieler–Lehman Algorithm

1. Assign an initial label $\ell^{(0)}(v)$ to each node. A usual practice is to just assign the degree.
2. Iteratively assign a new label by hashing the multiset of current labels within the neighborhood

$$\ell^{(i)}(v) = \text{HASH}(\{\{\ell^{(i-1)}(u) \mid u \in \mathcal{N}(v)\}\})$$

3. After running K iterations of step-2, our label summarizes a K -hop neighborhood. Now use statistics of labels as features for the graphs.

The WL Kernel is computed by measuring the difference between the label sets of two graphs.

Definition 2.5 (Graphlets)

Graphlets are small non-isomorphic induced subgraphs representing connected patterns in a network, and their frequency can be used to assess network structure.

Graphlet kernels enumerate all graph structures of a possible size and count their occurrence in a graph. This is combinatorially difficult and an alternative is to use path-based methods which examine different kinds of paths in the graph. For example, the *random-walk kernel* [Kashima et al., 2003] involves running random walks all over the graph and then counting the occurrence of different degree sequences, while the *shortest-path kernel* [Borgwardt and Kriegel, 2005] uses shortest paths between nodes instead of random-walks.

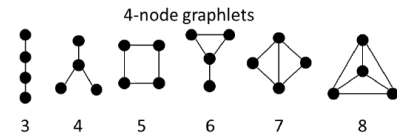


Figure 5: Graphlets with 4 nodes

2.3 Measures of Neighborhood Overlap

Now we try to quantify the relation between nodes, which cannot be done using node-level or graph-level kernels. We denote $\mathbf{S} \in \mathbb{R}^{|V| \times |V|}$ as the similarity matrix summarizing the statistics. We commonly take the probability of an edge occurring proportional to the statistical measure, i.e $\mathbb{P}(A_{uv} = 1) \propto \mathbf{S}[u, v]$.

1. Local overlap measures

- (a) **Common Neighbors (CN)**: It is the simplest neighborhood overlap measure and counts the number of neighbors two nodes share.

$$\mathbf{S}_{CN}[u, v] = |\mathcal{N}(u) \cap \mathcal{N}(v)| \quad (4)$$

We can extend the idea to local overlap measures defined over CN, denoted by the next three quantifiers.

- (b) **Sorensen index**: This normalizes the counts of common neighbours by the sum of node-degrees.

$$\mathbf{S}_{Sorensen}[u, v] = \frac{2|\mathcal{N}(u) \cap \mathcal{N}(v)|}{d_u + d_v} \quad (5)$$

- (c) **Salton index**: This normalizes using the product of degrees.

$$\mathbf{S}_{Salton}[u, v] = \frac{2|\mathcal{N}(u) \cap \mathcal{N}(v)|}{\sqrt{d_u d_v}} \quad (6)$$

- (d) **Jaccard index**: This is an IOU measure.

$$\mathbf{S}_{Jaccard}[u, v] = \frac{|\mathcal{N}(u) \cap \mathcal{N}(v)|}{|\mathcal{N}(u) \cup \mathcal{N}(v)|} \quad (7)$$

- (e) **Resource Allocation (RA)**: It places importance on the common neighbors by summing the inverse degrees, i.e

$$\mathbf{S}_{RA}[u, v] = \sum_{w \in \mathcal{N}(u) \cap \mathcal{N}(v)} \frac{1}{d_w} \quad (8)$$

- (f) **Adamic-Adar (AA)**: It is similar to RA, just utilizes inverse log of the degrees, i.e

$$\mathbf{S}_{AA}[u, v] = \sum_{w \in \mathcal{N}(u) \cap \mathcal{N}(v)} \frac{1}{\log d_w} \quad (9)$$

Both these measures infer that common neighbors with low degrees are more informative.

2. Global overlap measures

- (a) **Katz Index:** We take a discounted count of number of paths of all lengths between a pair of nodes, i.e

$$\mathbf{S}_{Katz}[u, v] = \sum_{i=1}^{\infty} \beta^i A_{uv}^i = \left((\mathbf{I} - \beta \mathbf{A})^{-1} - \mathbf{I} \right)[u, v] \quad (10)$$

Here $\beta \in \mathbb{R}^+$ is our discount, and a small $\beta < 1$ down-weights the importance of long-paths.

- (b) **Leicht, Holme, and Newman Similarity:** The issue with Katz index is that it is biased towards high-degree nodes. To avoid that, we normalize the number of paths by its expectation, i.e $\frac{A_{uv}^i}{\mathbb{E}[A_{uv}^i]}$. The expectations can be computed by drawing a random graph with the same set of degrees as the graph under test, and obtain

$$\mathbb{E}[A_{uv}] = \frac{d_u d_v}{2|E|} \quad \mathbb{E}[A_{uv}^2] = \frac{d_u d_v}{(2|E|)^2} \sum_{w \in V} d_w (d_w - 1) \quad (11)$$

Calculation for $i > 2$ needs to be approximated, and to do so we utilize the fact mentioned earlier (2.1) to notice that the growth is by a factor λ for large i , and thus $\mathbb{E}[A_{uv}^i] = \frac{d_u d_v \lambda^{i-1}}{2|E|}$. Finally, the index is given as

$$\mathbf{S}_{LHN}[u, v] = I_{uv} + \frac{2|E|}{d_u d_v} \sum_{i=0}^{\infty} \beta \lambda^{i-1} A_{uv}^i \quad (12)$$

- (c) **Random Walk Methods:** Take \mathbf{D} to be a diagonal matrix with entries as the node degrees. Looking at the Personalized PageRank algorithm [Rajaraman and Ullman, 2012], we define $\mathbf{P} = \mathbf{A}\mathbf{D}^{-1}$ and define vector \mathbf{q}_u as the probability vector of visiting all nodes starting at node u . In a recursive fashion we incorporate the probability of restarting the random walk with probability c and write $\mathbf{q}_u = c\mathbf{P}\mathbf{q}_u + (1 - c)\mathbf{e}_u$ where \mathbf{e}_u is a one-hot vector at node u . It's solution is

$$\mathbf{q}_u = (1 - c)(1 - c\mathbf{P})^{-1}\mathbf{e}_u \quad (13)$$

The similarity measure is

$$\mathbf{S}_{RW}[u, v] = \mathbf{q}_u[v] + \mathbf{q}_v[u] \quad (14)$$

2.4 Graph Laplacians

Laplacians are transformations of the adjacency matrix which have useful properties.

Definition 2.6 (Unnormalized Laplacian)

¹ The unnormalized Laplacian is given by $\mathbf{L} = \mathbf{D} - \mathbf{A}$ and for simple

graphs has the following properties:

◇ It is symmetric and positive semi-definite.

◇ For $\mathbf{x} \in \mathbb{R}^{|V|}$,

$$\mathbf{x}^T \mathbf{L} \mathbf{x} = \sum_{(u,v) \in E} (\mathbf{x}[u] - \mathbf{x}[v])^2 \quad (15)$$

◇ The eigenvalues of \mathbf{L} are non-negative.

◇ The geometric multiplicity of \mathbf{L} denotes the number of connected components in the graphs.

We extend the definition to normalized Laplacians in two ways:

$$\mathbf{L}_{sym} = \mathbf{D}^{-1/2} \mathbf{L} \mathbf{D}^{-1/2} \quad \mathbf{L}_{RW} = \mathbf{D}^{-1} \mathbf{L} \quad (16)$$

2.5 Clustering using Laplacians

We can extend property 4 of Laplacians to give optimal clustering within a fully-connected graph.

Definition 2.7 (Graph Cut)

Let $U \subset V$ and let \bar{U} be its complement. Given K non-overlapping partitions of a graph $\{U_i\}_{i=1}^K$, we define the cut value as

$$cut(U_1, \dots, U_K) = \frac{1}{2} \sum_{k=1}^K |(u, v) \in E : u \in U_k, v \in \bar{U}_k| \quad (17)$$

It essentially counts the number of edges crossing the boundary between partitions.

We simply can't minimize this as it yields single-node clusters.

Hence, we normalize it in two ways:

$$RatioCut(U_1, \dots, U_K) = \frac{1}{2} \sum_{k=1}^K \frac{|(u, v) \in E : u \in U_k, v \in \bar{U}_k|}{|U_k|} \quad (18)$$

$$NCut(U_1, \dots, U_K) = \frac{1}{2} \sum_{k=1}^K \frac{|(u, v) \in E : u \in U_k, v \in \bar{U}_k|}{\sum_{w \in U_k} d_w} \quad (19)$$

The former penalizes small cluster sizes, while the latter forces similar number of edges incident on nodes in all clusters. Minimizing them can be defined as optimal clustering. The optimal solution for 2-clustering is derived here 2.5. The above derivation can be done for $NCut$ but the result is has v_2 for \mathbf{L}_{RW} .

2.6 Spectral Clustering

We can use the principle of spectral clustering by representing the nodes as a spectrum of \mathbf{L} , and do an approximate optimal clustering

2-clustering over RatioCut: We solve

$$\min_{U \in V} RatioCut(U, \bar{U})$$

Define $\mathbf{a} \in \mathbb{R}^{|V|}$ and $\xi = \sqrt{\frac{|U|}{|\bar{U}|}}$ such that $a_u = \xi$ if $u \in U$ else $a_u = -\frac{1}{\xi}$. With some calculation, we can show that

$$\mathbf{a}^T \mathbf{L} \mathbf{a} = |V| RatioCut(U, \bar{U})$$

We can also show that

$$1. \sum_{u \in V} a_u = 0 \Leftrightarrow \mathbf{a} \perp \mathbf{1}$$

$$2. \|\mathbf{a}\|^2 = |V|$$

Solving $\min_{U \in V} \mathbf{a}^T \mathbf{L} \mathbf{a}$ and the given constraints is NP-hard as our set is discrete. We remove discreteness and solve

$$\min_{U \in \mathbb{R}^{|V|}} \mathbf{a}^T \mathbf{L} \mathbf{a}$$

using the **Rayleigh-Ritz Theorem** to get that $\mathbf{a} = v_2$, i.e the second smallest eigenvector of \mathbf{L} (note that the smallest eigenvector is just $\mathbf{1}$). Finally, we discretize the nodes based on the sign of a_u , i.e $u \in U$ if $a_u \geq 0$ else $u \in \bar{U}$.

(Note: Smallest eigenvector denotes the eigenvector for the smallest eigenvalue).

as follows:

General Spectral Clustering Algorithm

1. Find K smallest eigenvectors of \mathbf{L} except the smallest.
2. Let $\mathbf{U} \in \mathbb{R}^{|V| \times (K-1)}$ and have columns as the eigenvectors from above.
3. Let $\mathbf{z}_u = \mathbf{U}[u]$ be the row for node u .
4. Run K-means clustering on the embeddings $\mathbf{z}_u \forall u \in V$.

2.7 Learning Representations

Clearly, although the traditional approaches are useful many times, they are hand-engineered and thus are inflexible, i.e cannot adapt to changes. Also, designing them is time-consuming and expensive. Thus, we would benefit if we could learn structural representations over graphs.

3 Learning Node Embeddings

Following [Hoff et al., 2002], we project nodes into a latent space and the geometric relations in that space correspond to those in the original graph.

3.1 Encoder-Decoder Model

The task is essentially to *encode* each node into a low-dimensional vector, and finally *decode* the embedding and use it to reconstruct information [Hamilton et al., 2017b].

Definition 3.1 (Encoder)

The encoder is a function $ENC: V \rightarrow \mathbb{R}^d$ which maps nodes $v \in V$ to vector embeddings $\mathbf{z}_v \in \mathbb{R}^d$. Commonly the approach is a shallow embedding one, and thus the encoder function is an embedding lookup on the node ID, i.e $ENC(v) = \mathbf{Z}[v]$ where $\mathbf{Z} \in \mathbb{R}^{|V| \times d}$ is the embedding matrix.

Definition 3.2 (Decoder)

The decoder reconstructs certain graph statistics from node embeddings generated by the encoder. It is common to use pairwise decoders, i.e $DEC: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}^+$ and they are similarity measures.

The goal is to minimize reconstruction loss, i.e

$$DEC(ENC(u), ENC(v)) = DEC(\mathbf{z}_u, \mathbf{z}_v) \approx \mathbf{S}[u, v] \quad (20)$$

To achieve the goal above, we minimize an empirical reconstruction loss function

$$\mathcal{L} = \sum_{(u,v) \in \mathcal{D}} \ell(DEC(\mathbf{z}_u, \mathbf{z}_v), \mathbf{S}[u, v]) \quad (21)$$

where \mathcal{D} is a subset of node pairs, and ℓ is a loss function such as cross-entropy or mean-squared error loss. We look at some common shallow embedding algorithms:

1. Factorization based Encoder-Decoder

The loss function for such approaches can be minimized using factorization algorithms such as SVD.

- (a) **Laplacian Eigenmaps** [Belkin and Niyogi, 2001]: It is an old but influential approach built upon spectral clustering. The decoder is

$$DEC(\mathbf{z}_u, \mathbf{z}_v) = \|\mathbf{z}_u - \mathbf{z}_v\|_2^2 \quad (22)$$

and the loss function is

$$\mathcal{L} = \sum_{(u,v) \in \mathcal{D}} DEC(\mathbf{z}_u, \mathbf{z}_v) \mathbf{S}[u, v] \quad (23)$$

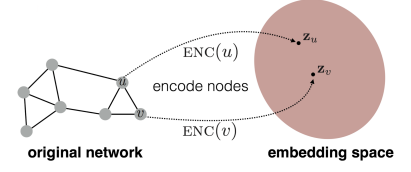


Figure 6: Node Representation Learning Source: [Stanford](#)

Mean Squared Error loss:

$$MSE(y, \tilde{y}) = \frac{1}{n} \sum_{i=1}^n (y_i - \tilde{y}_i)^2$$

Cross Entropy loss:

$$CE(y, \tilde{y}) = -\frac{1}{n} \sum_{i=1}^n y_i \log(\tilde{y}_i)$$

If we construct \mathbf{S} to satisfy properties of \mathbf{L} , and \mathbf{z}_u are d -dimensional, then the optimal solution for minimizing the loss is given by the d smallest eigenvectors of \mathbf{L} excluding the smallest one, i.e $\mathbf{1}$

(b) **Inner-product methods:** The decoder in these is

$$\text{DEC}(\mathbf{z}_u, \mathbf{z}_v) = \mathbf{z}_u^T \mathbf{z}_v \quad (24)$$

and the loss function is

$$\mathcal{L} = \sum_{(u,v) \in \mathcal{D}} \|\text{DEC}(\mathbf{z}_u, \mathbf{z}_v) - \mathbf{S}[u, v]\|_2^2 \quad (25)$$

- i. *Graph Factorization* [Ahmed et al., 2013]: It uses $\mathbf{S} \triangleq \mathbf{A}$.
- ii. *GraRep* [Cao et al., 2015]: It defines \mathbf{S} based on powers of \mathbf{A} .
- iii. *HOPE* [Ou et al., 2016]: It defines \mathbf{S} based on neighborhood overlap measures.

2. Random Walk Embeddings

The notion in such embeddings are that two nodes have similar embeddings if they tend to co-occur on short random walks over the graph. We learn embeddings to get

$$\text{DEC}(\mathbf{z}_u, \mathbf{z}_v) \triangleq \frac{e^{\mathbf{z}_u^T \mathbf{z}_v}}{\sum_{w \in V} e^{\mathbf{z}_u^T \mathbf{z}_w}} \approx p_{G,T}(v|u) \quad (26)$$

where $p_{G,T}(v|u)$ is the probability of reaching node v starting at u on a T -length random walk ($T \in \llbracket 2..10 \rrbracket$). Clearly the decoder is learning an asymmetric stochastic measure. To train the embeddings, the loss function employed is

$$\mathcal{L} = \sum_{(u,v) \in \mathcal{D}} -\log(\text{DEC}(\mathbf{z}_u, \mathbf{z}_v)) \quad (27)$$

But the issue here is that evaluating the loss function has time complexity $\mathcal{O}(|V||\mathcal{D}|) \sim \mathcal{O}(|V|^2)$.

- (a) **DeepWalk** [Perozzi et al., 2014]: It employs uniformly random walks to define $p_{G,T}(v|u)$ and employs hierarchical softmax [Morin and Bengio, 2005] to approximate Equation (26)
- (b) **node2vec** [Grover and Leskovec, 2016]: It introduces hyperparameters to allow random walk probabilities to interpolate between walks more akin to BFS or DFS. It uses noise contrastive estimation along with negative sampling to rewrite \mathcal{L} as

$$\mathcal{L} = \sum_{(u,v) \in \mathcal{D}} -\log(\sigma(\mathbf{z}_u^T \mathbf{z}_v)) - \gamma \mathbb{E}_{v_n \sim P_n(V)} [\log(-\sigma(\mathbf{z}_u^T \mathbf{z}_{v_n}))] \quad (28)$$

with $\gamma > 0$ and $P_n(V)$ being the uniform distribution. Monte Carlo Sampling is used to calculate the expectation.

- (c) **LINE** [Tang et al., 2015]: It proposes a first objective

$$\text{DEC}(\mathbf{z}_u, \mathbf{z}_v) = \sigma(\mathbf{z}_u^T \mathbf{z}_v) \quad (29)$$

The notation $\llbracket x..y \rrbracket$ denotes the integer interval spanning from x to y , i.e $\{x, \dots, y\}$

It has been shown in [Qiu et al., 2017] that if we define

$$\mathbf{S}_{DW} = \log \left(\frac{\sum_{w \in V} d_w}{T} \left(\sum_{t=1}^T \mathbf{P}^t \right) \mathbf{D}^{-1} \right) - \log(b)$$

then, the embeddings \mathbf{Z} by DeepWalk satisfy $\mathbf{Z}^T \mathbf{Z} \approx \mathbf{S}_{DW}$. What is interesting to note that is we can eigendecompose the interior of above equation, and note that DeepWalk is closely related to \mathbf{L}_{sym} and hence spectral clustering, but differ by the fact that influence of eigenvalues is controlled by DeepWalk through T .

and $\mathbf{S}[u, v] = A_{uv}$. Along with that, a second objective is the decoder same as DeepWalk, but trained using KL-Divergence to encode 2-hop information (i.e \mathbf{A}^2), and it explicitly reconstructs first and second-order neighborhood information.

KL-Divergence

$$D_{KL}(p \parallel q) = \sum_{i=1}^n p(x_i) \log \left(\frac{p(x_i)}{q(x_i)} \right)$$

3.2 Drawbacks of Shallow Embeddings

- ◇ Lack of parameter sharing between nodes is a hinder to statistic and computational efficiency. We want embeddings to be independent of $|V|$.
- ◇ They do not leverage node features in ENC.
- ◇ They are transductive, and thus it is difficult to generate embeddings for newer nodes and cannot be used in inductive applications which involve generalization.

3.3 Reconstruction of Multi-Relational Data

Multi-relational graphs add a relation type, and to address that, we make the decoder multi-relational, i.e $\text{DEC} : \mathbb{R}^d \times \mathcal{R} \times \mathbb{R}^d \rightarrow \mathbb{R}^+$.

RESCAL [Nickel et al., 2011] defines the multi-relational decoder as $\text{DEC}(u, \tau, v) = \mathbf{z}_u^T \mathbf{R}_\tau \mathbf{z}_v$, where \mathbf{R}_τ is a learnable matrix for relation $\tau \in \mathcal{R}$, and we can learn \mathbf{Z} and \mathbf{R}_τ using a simple reconstruction loss

$$\mathcal{L} = \sum_{u \in V} \sum_{v \in V} \sum_{\tau \in \mathcal{R}} \|\mathbf{z}_u^T \mathbf{R}_\tau \mathbf{z}_v - \mathcal{A}_{u\tau v}\|^2 \quad (30)$$

As seen above, it is a common practice to define the similarity measure based on the adjacency tensor, and thus we are trying to reconstruct immediate neighbors from the embeddings.

1. **Loss Functions:** The loss defined in Equation (30) requires $\mathcal{O}(|V|^2|\mathcal{R}|)$ time complexity, and is extremely expensive. Since it is realistically seen that multi-relational graphs are sparse, we would want a loss function of the order $\mathcal{O}(|E|)$. Also, since our target is to decode the adjacency tensor, which has binary values, we can't use the simple MSE loss.

- (a) **Cross-Entropy with Negative Sampling:** Defining σ as the sigmoid function, $P_{n,u}(V)$ as the negative sampling distribution, and $\gamma > 0$, we write

$$\mathcal{L} = \sum_{(u, \tau, v) \in E} \left(-\log(\sigma(\text{DEC}(\mathbf{z}_u, \tau, \mathbf{z}_v))) \right. \\ \left. - \gamma \mathbb{E}_{v_n \sim P_{n,u}(V)} [\log(\sigma(-\text{DEC}(\mathbf{z}_u, \tau, \mathbf{z}_{v_n}))) \right] \quad (31)$$

Knowledge Graphs: They are multi-relational graphs $G = (V, E)$ with edges defined as tuples $e = (u, \tau, v)$ indicating the presence of a relation $\tau \in \mathcal{R}$ between two nodes. The task of knowledge graph completion is generally to predict the missing edges in the graph.

If we calculate the expectation using Monte Carlo Sampling, we can write by sampling over a small subset $\mathcal{P}_{n,u} \subset P_{n,u}(V)$

$$\mathcal{L} = - \sum_{(u,\tau,v) \in E} \log(\sigma(\text{DEC}(\mathbf{z}_u, \tau, \mathbf{z}_v))) - \sum_{v_n \in \mathcal{P}_{n,u}} \log(\sigma(-\text{DEC}(\mathbf{z}_u, \tau, \mathbf{z}_{v_n}))) \quad (32)$$

- (b) **Max-Margin Loss:** Also called the *hinge loss*, we set a margin Δ and use *contrastive estimation* (but with direct comparison of output of the decoders) to write the loss function as

$$\mathcal{L} = \sum_{(u,\tau,v) \in E} \sum_{v_n \in \mathcal{P}_{n,u}} \max(0, -\text{DEC}(\mathbf{z}_u, \tau, \mathbf{z}_v) + \text{DEC}(\mathbf{z}_u, \tau, \mathbf{z}_{v_n}) + \Delta) \quad (33)$$

2. Decoders:

- (a) **RESCAL:** In the RESCAL decoder, a trainable matrix $\mathbf{R}_\tau \in \mathbb{R}^{d \times d}$ is associated with each relation

$$\text{DEC}(\mathbf{z}_u, \tau, \mathbf{z}_v) = \mathbf{z}_u^T \mathbf{R}_\tau \mathbf{z}_v \quad (34)$$

This gives rise to $\mathcal{O}(d^2)$ number of parameters for each relation, and that has high computational cost. Recent decoders aim to get $\mathcal{O}(d)$ number of parameters.

- (b) **Translational Decoders:** Such decoders represent relations as translations in embedding space.

- i. **TransE** [Bordes et al., 2013]: Using a d -dimensional embedding, the decoder is

$$\text{DEC}(\mathbf{z}_u, \tau, \mathbf{z}_v) = -\|\mathbf{z}_u + \mathbf{r}_\tau - \mathbf{z}_v\| \quad (35)$$

It is a simple decoder, and the likelihood of an edge is proportional to the distance translated according to the relation.

- ii. **TransX:** These are a class of models built upon the TransE model, and have the general form

$$\text{DEC}(\mathbf{z}_u, \tau, \mathbf{z}_v) = -\|g_{1,\tau}(\mathbf{z}_u) + \mathbf{r}_\tau - g_{2,\tau}(\mathbf{z}_v)\| \quad (36)$$

where $g_{i,\tau}$ are trainable transformations. [Wang et al., 2014] proposed the TransH model which projects the entity embeddings onto a learnable relation-specific hyperplane defined by the normal vector \mathbf{w}_r and then translating it.

$$\text{DEC}(\mathbf{z}_u, \tau, \mathbf{z}_v) = -\|(\mathbf{z}_u - \mathbf{w}_r^T \mathbf{z}_u \mathbf{w}_r) + \mathbf{r}_\tau - (\mathbf{z}_v - \mathbf{w}_r^T \mathbf{z}_v \mathbf{w}_r)\| \quad (37)$$

- (c) **DistMult** [Yang et al., 2014]: This generalizes the dot-product decoder from simple graphs, thus

$$\text{DEC}(\mathbf{z}_u, \tau, \mathbf{z}_v) = \langle \mathbf{z}_u, \mathbf{r}_\tau, \mathbf{z}_v \rangle = \sum_{i=1}^d \mathbf{z}_u[i] \times \mathbf{r}_\tau[i] \times \mathbf{z}_v[i] \quad (38)$$

This has a restriction, that it can only be used for symmetric relations, but majority of the relation are assymmetric.

- (d) **ComplEx** [Trouillon et al., 2016]: This extends the above idea and uses complex-valued embeddings, to bring in asymmetry, i.e $\mathbf{z}_u, \mathbf{r}_\tau, \mathbf{z}_v \in \mathbb{C}^d$.

$$\text{DEC}(\mathbf{z}_u, \tau, \mathbf{z}_v) = \Re \langle \mathbf{z}_u, \mathbf{r}_\tau, \bar{\mathbf{z}}_v \rangle = \Re \left(\sum_{i=1}^d \mathbf{z}_u[i] \times \mathbf{r}_\tau[i] \times \bar{\mathbf{z}}_v[i] \right) \quad (39)$$

- (e) **RotatE** [Sun et al., 2019]: This is again a complex-valued decoder, with the constraint that $|\mathbf{r}_\tau[i]| = 1 \implies \mathbf{r}_\tau[i] = e^{i\theta_{r,i}}$. This brings in rotations in the complex plane, and

$$\text{DEC}(\mathbf{z}_u, \tau, \mathbf{z}_v) = -\|\mathbf{z}_u \circ \mathbf{r}_\tau - \mathbf{z}_v\| \quad (40)$$

3.4 Representation Abilities of Multi-Relational Decoders

We characterize the representation abilities of the multi-relational decoders in three ways:

- ◇ **Symmetry**: Can different decoders model both symmetric and anti-symmetric relations? DistMult can only model symmetric, while TransE can only model anti-symmetric, while all others are seen to model both.
- ◇ **Inversion**: Most decoders can model inversion, though TransX and DistMult fail in this respect.
- ◇ **Compositionality**: It is seen that in RESCAL, defining $\mathbf{R}_{\tau_3} = \mathbf{R}_{\tau_1} \mathbf{R}_{\tau_2}$ and in TransE defining $\mathbf{r}_{\tau_3} = \mathbf{r}_{\tau_1} + \mathbf{r}_{\tau_2}$ allows modeling of compositionality. RotatE also models it, while the rest three decoders we saw fail to do so.

Symmetric relations mean that

$$(u, \tau, v) \in E \leftrightarrow (v, \tau, u) \in E$$

Anti-symmetric relations mean that

$$(u, \tau, v) \in E \rightarrow (v, \tau, u) \notin E$$

Inversion means that

$$(u, \tau_1, v) \in E \leftrightarrow (v, \tau_2, u) \in E$$

Compositionality means that

$$(u, \tau_1, y) \in E \wedge (y, \tau_2, v) \in E \rightarrow (u, \tau_3, v) \in E$$

References

- [Agrawal et al., 2017] Agrawal, M., Zitnik, M., and Leskovec, J. (2017). Large-scale analysis of disease pathways in the human interactome.
- [Ahmed et al., 2013] Ahmed, A., Shervashidze, N., Narayanamurthy, S., Josifovski, V., and Smola, A. (2013). Distributed large-scale natural graph factorization. pages 37–48.
- [Belkin and Niyogi, 2001] Belkin, M. and Niyogi, P. (2001). Laplacian eigenmaps and spectral techniques for embedding and clustering. In *NIPS*.
- [Bordes et al., 2013] Bordes, A., Usunier, N., Garcia-Duran, A., Weston, J., and Yakhnenko, O. (2013). Translating embeddings for modeling multi-relational data. In Burges, C. J. C., Bottou, L., Welling, M., Ghahramani, Z., and Weinberger, K. Q., editors, *Advances in Neural Information Processing Systems*, volume 26. Curran Associates, Inc.
- [Borgwardt and Kriegel, 2005] Borgwardt, K. and Kriegel, H. (2005). Shortest-path kernels on graphs. In *Fifth IEEE International Conference on Data Mining (ICDM'05)*, pages 8 pp.–.
- [Cao et al., 2015] Cao, S., Lu, W., and Xu, Q. (2015). 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. Association for Computing Machinery.
- [Donnat et al., 2017] Donnat, C., Zitnik, M., Hallac, D., and Leskovec, J. (2017). Spectral graph wavelets for structural role similarity in networks.
- [Gaudeflet et al., 2021] Gaudeflet, T., Day, B., Jamasb, A., Soman, J., Regep, C., Liu, G., Hayter, J., Vickers, R., Roberts, C., Tang, J., Roblin, D., Blundell, T., Bronstein, M., and Taylor-King, J. (2021). Utilizing graph machine learning within drug discovery and development. *Briefings in Bioinformatics*, 22.
- [Grover and Leskovec, 2016] Grover, A. and Leskovec, J. (2016). node2vec: Scalable feature learning for networks. *CoRR*, abs/1607.00653.
- [Hamilton et al., 2017a] Hamilton, W. L., Ying, R., and Leskovec, J. (2017a). Inductive representation learning on large graphs. *CoRR*, abs/1706.02216.
- [Hamilton et al., 2017b] Hamilton, W. L., Ying, R., and Leskovec, J. (2017b). Representation learning on graphs: Methods and applications. *CoRR*, abs/1709.05584.
- [Hoff et al., 2002] Hoff, P., Raftery, A., and Handcock, M. (2002). Latent space approaches to social network analysis. *Journal of the American Statistical Association*, 97:1090–1098.
- [Kashima et al., 2003] Kashima, H., Tsuda, K., and Inokuchi, A. (2003). Marginalized kernels between labeled graphs. volume 1, pages 321–328.
- [Kipf and Welling, 2016] Kipf, T. N. and Welling, M. (2016). Semi-supervised classification with graph convolutional networks. *CoRR*, abs/1609.02907.
- [Lü and Zhou, 2011] Lü, L. and Zhou, T. (2011). Link prediction in complex networks: A survey. *Physica A: Statistical Mechanics and its Applications*, 390(6):1150–1170.
- [Mcpherson et al., 2001] Mcpherson, M., Smith-Lovin, L., and Cook, J. (2001). Birds of a feather: Homophily in social networks. *Annual Review of Sociology*, 27:415–.
- [Morin and Bengio, 2005] Morin, F. and Bengio, Y. (2005). Hierarchical probabilistic neural network language model. In Cowell, R. G. and Ghahramani, Z., editors, *Proceedings of the Tenth International Workshop on Artificial Intelligence and Statistics*, volume R5 of *Proceedings of Machine Learning Research*, pages 246–252. PMLR. Reissued by PMLR on 30 March 2021.
- [Nickel et al., 2011] Nickel, M., Tresp, V., and Kriegel, H.-P. (2011). A three-way model for collective learning on multi-relational data. In *Proceedings of the 28th International Conference on Machine Learning, ICML'11*, page 809–816, Madison, WI, USA. Omnipress.

- [Ou et al., 2016] Ou, M., Cui, P., Pei, J., Zhang, Z., and Zhu, W. (2016). Asymmetric transitivity preserving graph embedding. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '16*, page 1105–1114, New York, NY, USA. Association for Computing Machinery.
- [Pandit et al., 2007] Pandit, S., Chau, D. H., Wang, S., and Faloutsos, C. (2007). Net-probe: A fast and scalable system for fraud detection in online auction networks.
- [Perozzi et al., 2014] Perozzi, B., Al-Rfou, R., and Skiena, S. (2014). Deepwalk: Online learning of social representations. *CoRR*, abs/1403.6652.
- [Qiu et al., 2017] Qiu, J., Dong, Y., Ma, H., Li, J., Wang, K., and Tang, J. (2017). Network embedding as matrix factorization: Unifying deepwalk, line, pte, and node2vec. *CoRR*, abs/1710.02971.
- [Rajaraman and Ullman, 2012] Rajaraman, A. and Ullman, J. (2012). *Mining of Massive Datasets*. Cambridge University Press.
- [Sun et al., 2019] Sun, Z., Deng, Z., Nie, J., and Tang, J. (2019). Rotate: Knowledge graph embedding by relational rotation in complex space. *CoRR*, abs/1902.10197.
- [Tang et al., 2015] Tang, J., Qu, M., Wang, M., Zhang, M., Yan, J., and Mei, Q. (2015). LINE: large-scale information network embedding. *CoRR*, abs/1503.03578.
- [Trouillon et al., 2016] Trouillon, T., Welbl, J., Riedel, S., Gaussier, E., and Bouchard, G. (2016). Complex embeddings for simple link prediction. In Balcan, M. F. and Weinberger, K. Q., editors, *Proceedings of The 33rd International Conference on Machine Learning*, volume 48 of *Proceedings of Machine Learning Research*, pages 2071–2080, New York, New York, USA. PMLR.
- [Wang et al., 2014] Wang, Z., Zhang, J., Feng, J., and Chen, Z. (2014). Knowledge graph embedding by translating on hyperplanes. *Proceedings of the AAAI Conference on Artificial Intelligence*, 28(1).
- [Watts et al., 2006] Watts, D., J. D., Strogatz, and H, S. (2006). *Collective dynamics of 'small world' networks*, pages 301–303.
- [Xiao et al., 2022] Xiao, S., Wang, S., Dai, Y., and Guo, W. (2022). Graph neural networks in node classification: survey and evaluation. *Machine Vision and Applications*, 33.
- [Yang et al., 2014] Yang, B., Yih, W.-t., He, X., Gao, J., and Deng, I. (2014). Embedding entities and relations for learning and inference in knowledge bases.
- [Ying et al., 2018] Ying, R., He, R., Chen, K., Eksombatchai, P., Hamilton, W. L., and Leskovec, J. (2018). Graph convolutional neural networks for web-scale recommender systems. *CoRR*, abs/1806.01973.
- [Zitnik et al., 2018] Zitnik, M., Agrawal, M., and Leskovec, J. (2018). Modeling polypharmacy side effects with graph convolutional networks.