
Bachelor Research Project: Graph Machine Learning

Jannick Stuby
University of Stuttgart
st160888@stud.uni-stuttgart.de

Christian Staib
University of Stuttgart
staib.christian@hotmail.de

Lukas Zeil
University of Stuttgart
st161467@stud.uni-stuttgart.de

Abstract

We present our research on graph machine learning that we have made as part of the bachelor research project. We started our project examining DeepChem and MoleculeNet but later decided to focus on graph machine learning in general. During our project we had two main goals: to learn about graph machine learning in general and implement some graph machine learning algorithms ourselves. Therefore our report is structured in two parts: first we give an overview of graph machine learning and then we present our own implementations of some graph machine learning algorithms.

1 Introduction

1.1 Problem Domain

Neural Networks are thriving in the modern society with software ranging from facial recognition to language translation augmenting humans day-to-day. In this aspect, Machine Learning (ML) for chemical sciences gains more and more value and therefore is the focus of an increasing amount of research. In consequence, an increasing amount of data regarding molecules are gathered and labeled, multiplying the amount of information that can be used to train ML models to help chemical discovery and molecule analysis to aid chemical discovery for drug development and molecule investigation.

To aid in the improvement of ML models, benchmarks have been developed. A benchmark in a ML setting, is a site where models are compared against an accumulation of datasets, where the datasets are usually split into fixed training- and testing sets. The benchmark around which this project will be largely revolving is MoleculeNet. They tried a standardized approach for splitting their data, by creating more ways in which these datasets can be split, that can be applied to every dataset they have.

As the MoleculeNet paper was published in 2018, we try to find and use methods already used in other context but not in the paper. We will set up a dataset-to-metric model for each benchmark using new tools in four areas: splitters, featurizers, deep learning, and transfer learning. These results will be compared against the best-known results from MoleculeNet. While we may not be able to enhance the best-known results, we will formulate our accumulated knowledge for future exploration.

1.2 Definitions

A graph is a mathematical structure that consists of a set of vertices V and a set of edges E . Similarly a attributed graph is a mathematical structure that consists of a set of vertices and a set of edges, where each vertex and edge has an attribute, e.g. $A(V \cup E) \rightarrow \Sigma^*$. While this basic definition does

Table 1: Node feature vector

Index	Feature	Value Range
0	atom number (H excluded, C=5)	1 - 119
1	chirality	unspecified, Tetrahedral CW, Tetrahedral CCW, other and misc
2	degree	0 - 10, misc
3	formal charge	-5 - 5, misc
4	number of H-atoms	0 - 8, misc
5	number of radicals	0 - 4, misc
6	hybridization	SP, SP2, SP3, SP3D, SP3D2, misc
7	is aromatic	true, false
8	is in ring	true, false

Table 2: Edge feature vector

Index	Feature	Value Range
0	bond type	single, double, triple, aromatic, misc
1	bond stereo	none, z, e, cis, trans, any
2	is conjugated	true, false

not specify any conditions on the attributes, it is possible to limit the attributes to a certain set of values. Some graph machine learning methods require the attributes to be of certain types or to provide certain functionality.

1.3 MoleculeNet

MoleculeNet [11] was introduced as a comprehensive benchmark for comparing ML models in molecular sciences.

2 Preliminary

As we started our project examining DeepChem, we first looked into handcrafted graph to vector methods.

Later we looked into graph machine learning.

2.1 Open Graph Benchmark

Open Graph Benchmark(OGB) [3] is a framework for benchmarking machine learning on graphs, providing large-scale, real world datasets as well as an easy to use software framework for loading datasets and evaluating model performance. For comparing model performances, OGB provides leaderboards for different datasets and challenges, split into the sections Node Property Prediction, Link Property Prediction and Graph Property Prediction, as well as a large-scale challenge [2] leaderboard, which aids comparison of performance on very large graphs.

Our main focus will be Graph Property Prediction, mainly on the datasets ogb-molhiv and ogb-molpcba. Both these datasets are derived from MoleculeNet and respectively provide a large amount of molecules for molecular property prediction. Molecules are represented by graphs, where atoms are the nodes and corresponding chemical bonds are represented by the edges. Node features are represented as a 9-dimensional vector, edge feature vectors are 3-dimensional, the respective value ranges can be seen in Table 1 and Table 2. For ogb-molhiv the model is evaluated using ROC-AUC, for ogb-molpcba Average Precision is used.

We additionally used ogb-molfreesolv as a smaller dataset which is also derived from MoleculeNet. FreeSolv contains SMILES strings of molecules from the Free Solvation Database and consists of "experimental and calculated hydration free energy of small molecules in water" [11]. The model is evaluated using Root Mean Square Error.

We decided to use OGB as a benchmarking framework instead of the MoleculeNet benchmark results, because OGB submissions are much more versatile and active, containing results from many research teams with the latest submission for ogb-molhiv dating to May 2022, contrary to the MoleculeNet results which were provided by the authors of the paper only and are dated to January 2018.

2.2 Graph Embedding

Graph embedding is a technique in graph machine learning that maps nodes and edges of a graph into a low-dimensional vector space while preserving the graph structure and semantic information. This technique has become increasingly popular due to its usefulness in various applications such as social network analysis, recommendation systems, and molecular sciences.

The goal for this technique is to produce a vector which can be used as input for various existing machine-learning models and techniques, elevating their usefulness to graphs. To achieve a good representation for those graphs, local and global structural information has to be encoded into such vectors.

There is a multitude of different ideas and methods to create graph embeddings, we will focus mainly on random walk embeddings and deep learning approaches, to be exact most methods are combinations of those ideas [6, 1, 12].

2.2.1 doc2vec

Doc2vec is a widely used method for obtaining distributed representations of documents as continuous vectors in a low-dimensional space, also known as paragraph embeddings. Doc2vec, introduced by Le and Mikolov [4], extends the popular Word2vec algorithm [5] to the domain of entire documents and together with Word2vec builds the base for many ideas regarding graph embedding.

Doc2vec learns a vector representation for each document in a corpus by utilizing a neural network architecture called Paragraph Vector. The architecture is similar to that of Word2vec but includes an additional vector representation for each document. There are two ways to obtain document vectors: Distributed Memory (DM) and Distributed Bag of Words (DBOW).

In the Distributed Memory model, the document vector is trained to predict the next word in a randomly sampled sequence of words, with the document vector and a window of surrounding words as input. The document vector is then updated alongside the word vectors during backpropagation. This model captures the overall topic and meaning of the document, as it has access to all words in the document.

In contrast, the Distributed Bag of Words model disregards the order of words and only utilizes the occurrence information of the words in the document. The input to the model is a bag of words from a randomly sampled window, and the objective is to predict the words in the window. The document vector is trained alongside the word vectors, and it aims to represent the document’s overall semantic content.

In both models, the document vector is learned by optimizing a loss function, typically the negative log-likelihood of predicting the target words. The resulting document vectors are continuous, dense, and low-dimensional, making them suitable for various downstream tasks, such as document classification, clustering, and information retrieval.

Doc2vec has been shown to outperform other methods for document embeddings in several benchmark datasets and downstream tasks. The embeddings capture the semantics and meaning of the documents and can be used to infer similarities between documents, making them useful for document clustering and information retrieval.

2.3 node2vec

Node2Vec, presented by Grover and Leskovec in 2016, is a framework for learning low-dimensional representations of nodes in large-scale networks [1]. The goal of node2vec is to learn embeddings that capture the structural information of nodes, such as their connectivity patterns, in a way that is useful for downstream tasks, such as node classification, clustering, and link prediction.

Node2vec builds on the skip-gram model used in natural language processing, which learns word embeddings by predicting the context of a given word. Similarly, node2vec learns node embeddings by training a neural network to predict the context of a given node, where the context of a node is defined as the set of nodes that are likely to appear in random walks starting from that node.

To generate random walks, node2vec uses a biased random walk strategy that balances the trade-off between exploring the network and exploiting the local neighborhood of a node. The random walk starts at a given node choosing for each step to either revisit the previous node or move to the next one based on two parameters: the in-out degree bias parameter p and the return parameter q , where p controls the likelihood of staying in the same community, while q controls the likelihood of visiting a node that is far away from the current one. By tuning these parameters, node2vec can generate different types of random walks that capture different types of structural information in the network.

Once the random walks are generated, node2vec uses the skip-gram model to learn node embeddings by maximizing the likelihood of predicting the context nodes given the current node. The skip-gram model is trained using negative sampling, which samples negative nodes that are not in the context set and adjusts the embedding vectors to minimize the difference between the predicted and actual context sets.

Node2vec has several advantages over previous methods for network embedding. While being scalable to large networks and being able to handle heterogeneous networks with multiple types of nodes and edges, it can also capture both the local and global structure of the network by tuning the p and q parameters. Additionally, it can handle noisy or incomplete networks by incorporating information from the context nodes. Finally, it can be applied to a wide range of downstream tasks, such as link prediction, node classification, and community detection.

2.3.1 graph2vec

Graph2vec is a method proposed by Narayanan et al. [6] for learning distributed representations of graphs. This approach builds on top of the previously described word2vec and Doc2vec methods, but instead of learning representations for words or paragraphs, it learns representations for graphs.

The Graph2vec method is based on two main ideas: Firstly, a set of subgraphs is generated for each graph in the dataset, which can be seen as 'context graphs'. These subgraphs are used to train a skip-gram neural network to learn embeddings for all graphs.

Context graphs are generated by performing a random walk on the graph. The next node to be selected during each step of the walk is chosen according to its distance to the currently selected one. The resulting sequence of nodes is then used to generate a set of subgraphs which are added to a pool of context graphs.

After that, the skip-gram model is used to learn embeddings for the graphs. Each graph is represented as a bag-of-context, i.e., a set of subgraphs that appear in the context of the graph. The skip-gram model then learns to predict the context graphs of a given graph, given the graph's embedding.

The skip-gram model used in Graph2vec is similar to the one used in word2vec, but instead of predicting the neighboring words of a given word, it predicts the context graphs of a given graph. The model takes as input a graph embedding and a context graph, and predicts whether the context graph appears in the bag-of-context-graphs of the input graph.

Graph2vec uses a negative sampling technique that is similar to the one used in word2vec to train the model. The objective is to maximize the probability of observing the context graphs of a given graph, while minimizing the probability of observing randomly sampled negative context graphs.

The resulting embeddings capture the structural properties of the graphs, allowing them to be used for various downstream tasks. In our case, we will mainly use Graph2vec as a benchmark and reference point.

2.4 Transformers

TODO: add Transformer [8] section: add explanation, formulas, visualization

In 2017 Vaswani et al. proposed a new network architecture for NLP problems, called the Transformer [8]. The Transformer architecture is build of many layers, each containing a self-attention

module and a position-wise feed-forward network. Self-attention is then calculated as described in [12]:

$$Q = HW_Q, \quad K = HW_K, \quad V = HW_V, \quad (1)$$

$$A = \frac{QK^\top}{\sqrt{d_K}}, \quad \text{Attn}(H) = \text{softmax}(A)V, \quad (2)$$

Here Q , K , and V are calculated by multiplying the input H of the attention module with the respective matrices $W_Q \in \mathbb{R}^{d \times d_K}$, $W_K \in \mathbb{R}^{d \times d_K}$ and $W_V \in \mathbb{R}^{d \times d_V}$ which then represent Queries, Keys and Values of the attention module. Those matrices are learnable and are produced during the training of the model. d denotes the hidden dimensions, which in case of single-head self attention is the same as d_K and d_V . The hidden dimension is the dimension of the key vectors, in case of [8] this was 64. This leads to the equation for the attention matrix above and its adaption for Graphormer in Equation 3. The result of this equation is the softmax score, which is then multiplied with the value vectors. The sum of those weighted vectors is the output of the self-attention module. By doing this, other tokens of the input which are seen as relevant, will continue upwards while non relevant tokens are drowned with a small softmax score.

2.4.1 Graphormer

As an adaption of Transformers [8] for graph representation learning, Graphormer [12] was presented by Ying et al., providing a performant Transformer model, excelling in various graph representation learning tasks, particularly for complex structures.

Graphormer achieves such good results by introducing three types of encodings as a bias to the self-attention mechanism used in Transformer architecture. While Transformer models work by computing weighted sum of input features based on attention scores between different tokens during self-attention, Graphormer incorporates structural information of different nodes in pairwise relation to each other into the self-attention mechanism, thus increasing the information encapsulated in the self-attention module, which is held by graph structures and features present.

To capture information about the importance of a node in a given network, Graphormer uses a measure called node centrality by utilizing degree centrality of a given node. This is done by introducing learnable embedding vectors, one for the degree of ingoing edges and one for the outgoing ones. Those vectors are scalars which are added to the features of each node, so node importance and semantic correlation are respected during self attention. Those scalars are indexed from two vectors according to the respective degrees as seen in Figure 1.

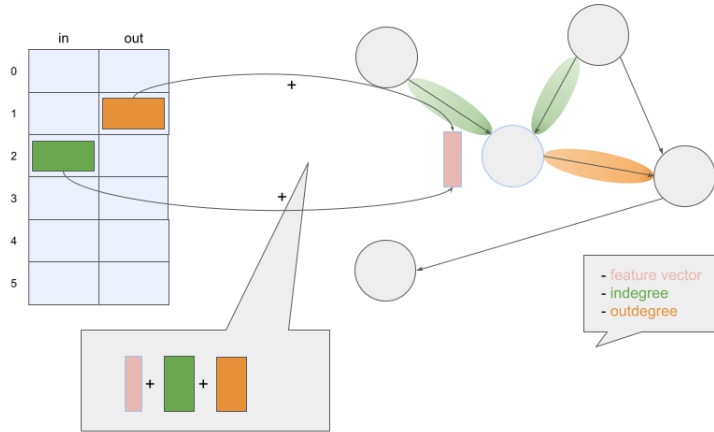


Figure 1: Centrality Encoding: learnable embedding vectors are added to feature vectors of nodes

The other two encodings are added as bias to the Attention Module used in Transformers as described in Figure 2 and the following equation.

$$A_{ij} = \frac{(h_i W_Q)(h_j W_K)^T}{\sqrt{d}} + b_{\Phi(v_i, v_j)} + c_{ij} \quad (3)$$

The bias $b_{\Phi(v_i, v_j)}$ is produced by the so called spatial encoding, which captures positional information for each node with respect to the network. This is achieved by using the distance of the shortest path between two connected nodes, described by the function Φ , combined with a learnable vector which is indexed according to Φ . The resulting scalar is added as bias in the self-attention module.

By doing that, the model can attend to a more relevant subset of the network according to the task. An example given in the paper is the possibility to increase attention to nodes surrounding a given node instead of nodes which lay further away.

Finally Graphormer introduces a method to include edge information in the model by adding an the additional bias term c_{ij} to the attention module. This bias is created by averaging the dot-products of edge features (x_{e_n}) and a learnable scalar ($(w_n^E)^T$), which is indexed according to the distance of the edge to the start node for each edge along the shortest path between each two connected nodes.

$$c_{ij} = \frac{1}{N} \sum_{n=1}^N x_{e_n} (w_n^E)^T \quad (4)$$

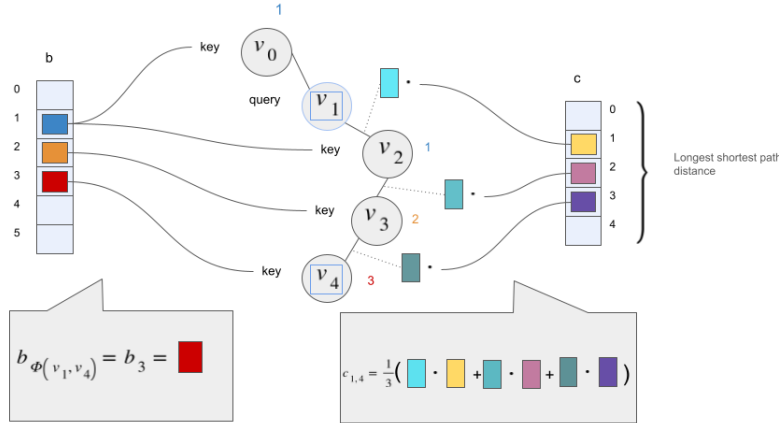


Figure 2: Spatial Encoding (left) and Edge Encoding (right) added as bias to Attention Module

2.4.2 GraphGPS

L. Rampásek et al. introduced GraphGPS in their paper 'Recipe for a General, Powerful, Scalable Graph Transformer.' [7] The objective of this model recipe is to provide a foundation for incorporating structural and positional information into node and edge features, while maintaining the flexibility of the model that implements these features.

GraphGPS serves as a preprocessor for the actual model learning on the dataset. The authors tested GraphGPS with several models, such as Transformer, Performer, BigBird, Gine, and more. They designed it to allow models that permit only Node Features, known as the Global Attention layer, as well as models that permit both Node and Edge Features, known as the MPNN layer. Both layers can also be used in combination to integrate the distinct advantages of two models.

The authors categorized structural and positional information into two types since the distance between nodes, although indicating the graph's structure, is insufficient to capture structural similarities. They incorporated this information by utilizing structural and positional encodings (SE and PE), which are

divided into Local, Global, and Relative categories. Local and Global encodings are implemented as node features, while Relative encodings are implemented as edge features.

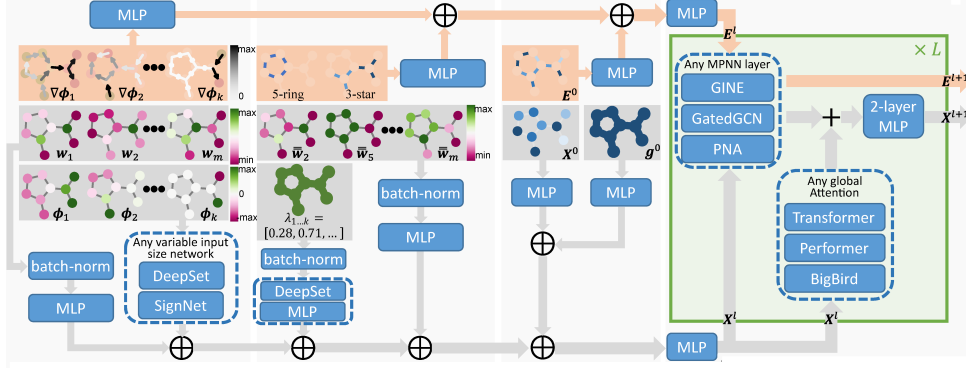


Figure 3: General overlook of GraphGPS taken from [7]

2.4.3 Heterogeneous Interpolation on Graphs

A winning entry for the dataset ogb-molpcba, called HIG-GraphClassification [9, 10] was submitted to the OGB leaderboard in 2021 by Wang et al., providing an implementation and a brief technical report describing the research. By using their method combined with Graphormer, they were able to achieve better average precision than all previous submissions.

The report introduces heterogeneous interpolation, which is done by dropping the feature vectors of several randomly selected nodes and replacing them by the interpolated feature mix of all neighboring nodes. By using a mixing ratio, the influence of each neighbors features can be adapted. To account for possible information loss, KL-Divergence constraint loss is added. By doing that, the distributions of two identical graphs remain similar after interpolating some feature vectors. The general idea was visualized by Wang et al. in Figure 4.

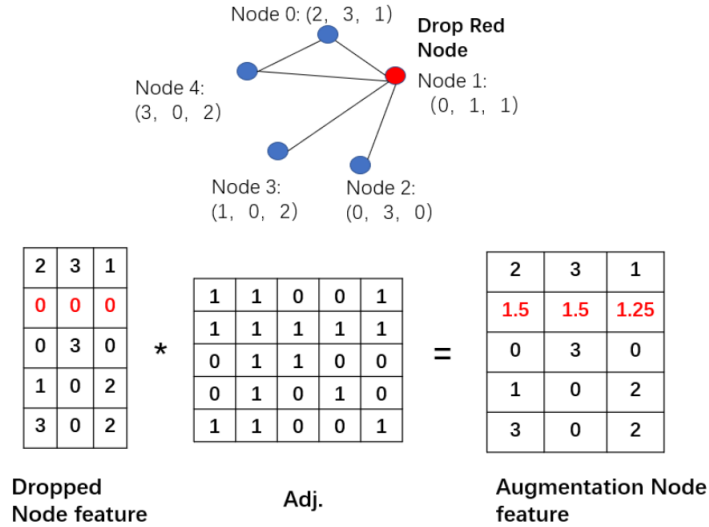


Figure 4: Visualization of HIG-GraphClassification by Wang et al., source: [10]

We combine the idea of interpolating feature vectors of neighboring nodes with GraphGPS to test performance with different embedding methods.

3 Our approach

In order to further strengthen our understanding, we decided to try to implement a graph machine learning method ourselves. As we knew we wouldn't be able to invent a game-changing method, we decided to try the simplest method we could imagine: A walk based embedding for whole graphs.

3.1 Walk based graph embeddings

One of the primary challenges in graph machine learning, as opposed to text-based machine learning, is the multidimensionality of graphs. Each node can have varying numbers of neighbors and different features and feature types. Walk based embeddings reduce this multidimensionality.

For our walk based embedding we focused on the embedding of whole graphs, as a similar method to embed nodes already exists (see [1]). The downstream task after the embedding of the graph is either regression or classification. For this we use sklearn.

3.1.1 walks to vector

A walk on graph is a unambiguously sequence of vertices and edges, where each vertex is connected to the next vertex by an edge. Given a graph g , we can write a walk as list of nodes (denoted by their id), e.g. $(0, 2, 3)$ is a walk that start from the node with the id 0, goes to 2 and end at 3. For a non-trivial graph there are many different possible walks, therefore we can create a set of walks w that all walk on the same graph. At this moment the walk doesn't contain any information relevant for our task, for this we need to substitute the node ids with the corresponding features. We substitute each id with the feature value prefixed with index of the feature separated by an otherwise unused character. For this we require that the features of a node are ordered. Therefore our schema for a feature in the walk is `FeatureIndex_FeatureValue`. We use this prefix in order to differentiate between same values of different feature types.

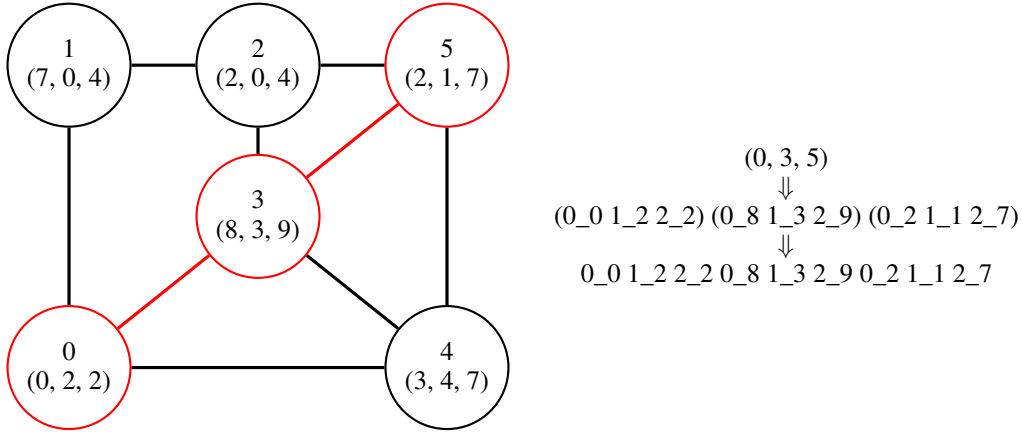


Figure 5: Example of graph on which a walk is done that is converted into a sentence

Lets talk about a concrete example. In the left side of figure 5 is an undirected graph with six vertices given. Each vertex has three numerical features. These features are given as list of integers below the node id. The red marked nodes form a walk from node 0 to node 5. On the right side of the figure, we listed the steps needed for each walk to generate the document needed for the downstream training: First we generate a walk, then we substitute the node id with the prefixed features and at last we concatenate the features into a sentence.

After this overview, we can formalize this method. Given a set of graphs $G = \{g_i \mid i \in 0 \dots n\}$, where n is the number of graphs, we can generate a set of sets of walks $W = \{w_i \mid i \in 0 \dots n\}$. Now we create a set W' where for each walk we substitute the node id with the features of the node. If we view each walk as sentence we get a set of documents (multiple sentences). These documents are feed into a text embedding method.

Algorithm 1: basic idea of our walk based embedding

```
1 def get_vectors(graphs)
2   documents = [ ]
3   forall graph in graphs do
4     walks = generate_walks(graph)
5     document = walks_to_document(walks)
6     documents.append(document)
7   end
8   model = Text_Embedding_Model()
9   model.fit(documents)
10  return model.get_document_vectors()
```

Pseudocode for our first implementation of our graph embedding method is given in algorithm 1. We identified two ways to influence the embedding in a meaningful way:

1. change the way walks are generated
2. change the text embedding model

3.1.2 change the way walks are generated

To accurately depict the neighborhood of a vertex, it is essential to represent the neighborhood within the generated document. This can be achieved when every combination of adjacent vertices is represented in the document multiple times.

We have considered two approaches for generating walks: All Pair Shortest Paths (APSP) and Random Walks. Each method possesses distinct advantages and drawbacks. APSP tends to have more hotspots of frequently visited graphs, which may introduce bias in the document; however, it generates similar walks for structurally analogous graphs. On the other hand, Random Walks can be set to a specific length and do not exhibit the same issue with frequently visited graphs. Nevertheless, they do not guarantee the generation of similar walks for structurally similar graphs.

One reason to consider using graph walks for analysis is that, given a sufficient number of walks over a graph, it becomes possible to reconstruct the original graph to some degree. These walks therefore provide insight into the graph's structure, capturing information about the relationships between vertices and the paths connecting them. By aggregating data from numerous walks, a comprehensive representation of the graph can be formed, enabling its reconstruction for further evaluation.

Reconstructing a graph from walks With an adequate quantity of walks over a graph, one can partially reconstruct the graph. However, there are limitations: For example, distinguishing between 'line' and 'circle' subgraphs is not always possible without any specific preparation of the graph. To illustrate this, refer to Figure 6. It is evident that each potential walk over one two graphs can be accommodated within the other graph as well.

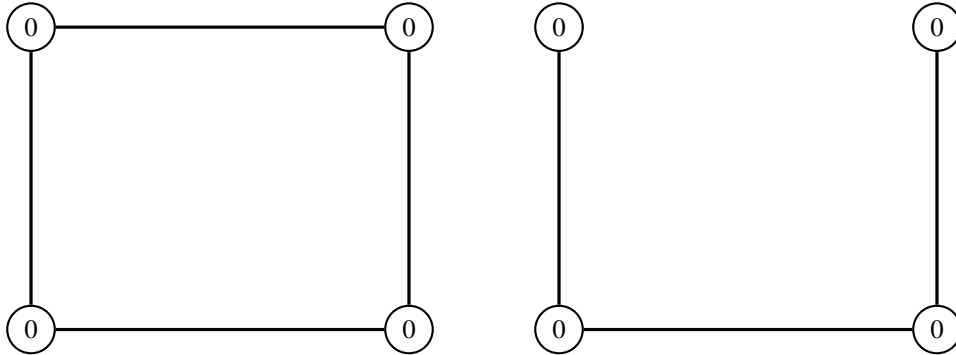


Figure 6: Example of two graphs whose walks can not be differentiated

Similar, in the case of regular graphs (i.e., graphs wherein all vertices exhibit the same degree), nodes with identical features tend to appear indistinguishable during walks. Consequently, due to these

inherent constraints, it was not anticipated that walk-based whole graph embeddings would yield exceptional results.

change the text embedding model Throughout the entirety of our project, we employed Doc2Vec[4] for transforming our walks into vector representations. Doc2Vec facilitates the training of a custom model using our specific vocabulary, allowing for the embedding of documents into vectors within this vocabulary. Despite our efforts to identify an alternative document-to-vector model that was user-friendly, we were unable to locate a suitable candidate.

3.1.3 results

Given the substantial amount of data generated by the random walks, which in turn resulted in a significant volume of training data for Doc2Vec, our focus was directed towards the smaller datasets from the Open Graph Benchmark (OGB). Specifically, we concentrated on ogb-molfreesolv, a regression dataset comprising 642 graphs. As previously noted, we were aware of the limitations associated with our methodology. Consequently, we established a baseline by comparing the vectors generated through our approach with randomly generated ones.

ogbg-molfreesolv	all	pair	shortest	path:6.528969941013406	[6.609125555671113,
6.57278283862703,		6.456680877236287,		6.539184198294062,	6.593562123437148,
6.417587074655182,		6.570877169599866,		6.478911030725958,	6.496761022461019,
6.5542275194263935]					

3.1.4 change the text embedding model

3.2 GraphGPS with HIG

GraphGPS is by design a highly malleable model, which allows for easy additions to the code. Thus we decided to implement HIG-GraphClassification in GraphGPS, to combine two already existing papers. For this combination we specifically targeted the Positional Encodings 7

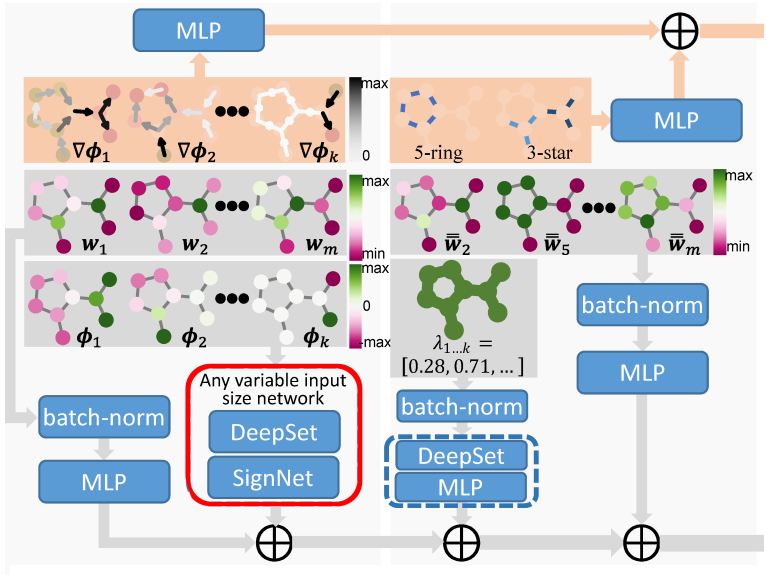


Figure 7: Position of HIG in GraphGPS structure

3.2.1 Implementation

The implementation of the heterogeneous interpolation, which involves dropping the feature vector of several randomly selected nodes and replacing them with an interpolated feature mix, can be found in algorithm 2. While our implementation does not include the mixing ratio, we have added a loss parameter and a min_size parameter. The loss parameter enables us to specify the probability with

which the graph replaces a node with an interpolated feature mix. On the other hand, the `min_size` parameter sets a minimum size for the graph. This is because a smaller graph tends to lose more information if a node is replaced, which could be disadvantageous. Additionally, we have included the `nodes_interpolated` parameter, which allows us to specify the number of nodes that get dropped instead of randomly choosing how many get dropped.

Algorithm 2: HiG-Code in GraphGPS

```

1 if 'HiG' in pe_types then
2   interpolation_chance = cfg.posenc_HIG.loss
3   is_interpolating = np.random.choice(2, 1, p=[1-interpolation_chance,
4     interpolation_chance])[0]
5   minimum_node_size = cfg.posenc_HIG.minimum_node_size
6   nodes_interpolated = cfg.posenc_HIG.nodes_interpolated
7   minimum_node_size = max(minimum_node_size, nodes_interpolated)
8   if data.num_nodes > minimum_node_size and is_interpolating == 1 then
9     rand_ints = np.random.choice(data.num_nodes, nodes_interpolated)
10    sum = {}
11    for rand_int in rand_ints do
12      data.x[rand_int] = 0
13      sum[rand_int] = 0
14    end
15    for i in data.edge_index[0] do
16      if i in rand_ints then
17        relevant_node = data.edge_index[1][i]
18        data.x[i] = torch.add(data.x[i], data.x[relevant_node])
19        sum[i] += 1
20      end
21    end
22    for rand_int in rand_ints do
23      data.x[rand_int] = data.x[rand_int] / sum[rand_int]
24    end
25  end

```

3.2.2 Results

Table 3: Results for different HIG Parameters

Model Focus	Loss	Nodes Interpolated	min.graph size	Test-AUC
No HIG				0.7710
Standard HIG	1.0	1	5	0.75274
Loss	0.5	1	5	0.77894
	0.1	1	5	0.77033
	0.01	1	5	0.7625
Nodes Interpolated	1.0	2	5	0.73967
	1.0	3	5	0.73882
	1.0	5	10	0.71498
min. graph size	1.0	1	10	0.76795
	1.0	1	20	0.77084
	1.0	1	50	0.7566
Combination	0.5	1	20	0.77164
	0.75	1	20	0.78321
HIG + RWSE	0.75	1	20	0.77924
HIG + LapPE	0.75	1	20	0.78412

References

- [1] Aditya Grover and Jure Leskovec. *node2vec: Scalable Feature Learning for Networks*. 2016. DOI: 10.48550/arXiv.1607.00653.
- [2] Weihua Hu et al. *OGB-LSC: A Large-Scale Challenge for Machine Learning on Graphs*. 2021. DOI: 10.48550/arXiv.2103.09430.
- [3] Weihua Hu et al. *Open Graph Benchmark: Datasets for Machine Learning on Graphs*. 2021. DOI: 10.48550/arXiv.2005.00687.
- [4] Quoc V. Le and Tomas Mikolov. *Distributed Representations of Sentences and Documents*. 2014. DOI: 10.48550/arXiv.1405.4053.
- [5] Tomas Mikolov et al. *Distributed Representations of Words and Phrases and their Compositionality*. 2013. DOI: 10.48550/arXiv.1310.4546.
- [6] Annamalai Narayanan et al. *graph2vec: Learning Distributed Representations of Graphs*. 2017. DOI: 10.48550/arXiv.1707.05005.
- [7] Ladislav Rampásek et al. *Recipe for a General, Powerful, Scalable Graph Transformer*. 2023. DOI: 10.48550/arXiv.2205.12454.
- [8] Ashish Vaswani et al. *Attention Is All You Need*. 2017. DOI: 10.48550/arXiv.1706.03762.
- [9] Yan Wang et al. *HIG-GraphClassification*. <https://github.com/TencentYoutuResearch/HIG-GraphClassification.git>. 2021.
- [10] Yan Wang et al. *Technical Report for OGB Graph Property Prediction*. Tech. rep. Shanghai, China: Tencent Youtu Lab, Dec. 2021. URL: https://github.com/TencentYoutuResearch/HIG-GraphClassification/blob/main/report/Report_HIG_OGB.pdf.
- [11] Zhenqin Wu et al. *MoleculeNet: A Benchmark for Molecular Machine Learning*. 2018. DOI: 10.48550/arXiv.1703.00564.
- [12] Chengxuan Ying et al. *Do Transformers Really Perform Bad for Graph Representation?* 2021. DOI: 10.48550/arXiv.2106.05234.