Graph Mixer Networks

Ahmet Sarıgün

Middle East Technical University ahmet.sarigun@metu.edu.tr

Abstract

In recent years, the attention mechanism has demonstrated superior performance in various tasks, leading to the emergence of GAT and Graph Transformer models that utilize this mechanism to extract relational information from graph-structured data. However, the high computational cost associated with the Transformer block, as seen in Vision Transformers, has motivated the development of alternative architectures such as MLP-Mixers, which have been shown to improve performance in image tasks while reducing the computational cost. Despite the effectiveness of Transformers in graph-based tasks, their computational efficiency remains a concern. The logic behind MLP-Mixers, which addresses this issue in image tasks, has the potential to be applied to graph-structured data as well. In this paper, we propose the Graph Mixer Network (GMN), also referred to as Graph Nasreddin Nets (GNasNets), a framework that incorporates the principles of MLP-Mixers for graph-structured data. Using a PNA model with multiple aggregators as the foundation, our proposed GMN has demonstrated improved performance compared to Graph Transformers. The source code is available publicly at https: //github.com/asarigun/GraphMixerNetworks.

1 Introduction

Graph Neural Networks (GNNs) are a powerful tool for working with graph-structured data, which is data that is made up of entities and their relationships. GNNs have been used to solve a wide range of problems, such as node classification, link prediction, graph generation, and many others. They have attracted great interest in recent years due to their performance and the ability to extract complex information.

Graph Convolutional Networks [1] (GCN) is a type of graph neural network (GNN) that uses graph convolutional layers to process data represented as graphs. GCNs can be used for various tasks such as node classification [2], graph classification [3], and link prediction. In each graph convolutional layer, the node features are updated by aggregating the features of their neighboring nodes. This is done through a convolution operation, where a linear combination of the neighboring node features is applied to each node, followed by a non-linear activation function. Graph Isomorphism Network (GIN) [4] is another type of GNN. GIN are able to distinguish non-isomorphic graphs. GIN consists of multiple layers of neural networks, where each layer aggregates the features of the neighboring nodes using a sum pooling operation, followed by a multi-layer perceptron (MLP). GIN can be used for various tasks such as node classification, graph classification, and link prediction. Graph Attention Networks (GAT) [5] is a GNN that uses attention mechanisms to assign different importance to different neighboring nodes when aggregating their features. Each node in GAT has a self-attention mechanism that allows it to weigh its own features and the features of its neighboring nodes in a learnable way. GAT can be used for various tasks such as node classification, graph classification, and link prediction. Message Passing Neural Networks (MPNN) [6] is a class of GNNs that generalize the idea of message passing between nodes in a graph. In MPNNs, messages are passed between nodes in the graph, and the node updates its state based on the messages received from its neighbors. MPNNs can be used for various tasks such as node classification, graph classification, and link prediction.

GCN, GIN, GAT, and MPNN are types of GNNs, each with their own characteristics and capabilities. GCN uses graph convolutional layers, GIN uses a sum pooling operation, GAT uses attention mechanisms to assign importance to different neighboring nodes, and MPNNs generalize the idea

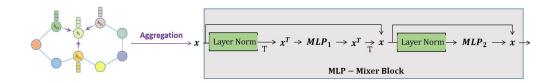


Figure 1: Mixer Layer of Graph Mixer Network where T denotes transpose operator.

of message passing between nodes in a graph. They can all be used for various tasks such as node classification, graph classification, and link prediction.

Transformer networks have been particularly successful in natural language processing tasks such as machine translation, text summarization, and language understanding.

The key innovation of transformer networks is the use of self-attention mechanisms. In a transformer network, each element in a sequence (such as a word in a sentence) is processed by attending to all the other elements in the sequence. This allows the network to weigh the importance of different elements in the sequence when making a prediction.

Vision transformers [7] are a type of transformer network designed for computer vision tasks, adapting the transformer architecture from natural language processing to handle image data. The input image is divided into non-overlapping patches and processed individually using self-attention mechanisms. The key advantage of vision transformers is their ability to handle images of arbitrary size, unlike traditional convolutional neural networks [8] which require a fixed image size.

Although transformers have shown superior performance on the image task, they have a computational quadratic complexity. MLP-Mixers [9, 10] have been used for the first time in the image task, and it has been shown to perform well without the quadratic complexity like in the attention mechanism and using only MLP.

In this work, we propose Graph Mixer Networks (or Graph Nasreddin Networks) which uses MLP-Mixer on Graphs as a novel Graph Neural Network. They are shown to have less computational complexity than transformers and comparable performance to baseline models.

The source code is available publicly at https://github.com/asarigun/GraphMixerNetworks.

2 Graph Mixer Networks

The proposed Graph Mixer Networks (or Graph Nasreddin Networks) method leverages the increased expressivity from the multi-aggregators models such as PNA [11] and MLP-Mixer [9, 10]. GMN architecture is given in Figure 1.

2.1 Motivation

Graph Transformers use the attention mechanism. Within this mechanism, in the self-attention mechanism, the dot product between Query and Key makes the computational complexity by $O(n^2)$. For this reason, MLP-Mixers, which are applied in image tasks, are seen to be more computationally efficient and perform better than Vision Transformers [9, 10]. That is why MLP-Mixer has considered solving computational efficiency in graphs in this study.

2.2 Multi-aggregators

Multiple aggregations are generalizations of sum aggregation within single aggregators and have been shown theoretically and empirically to be better at discriminating graphs. Therefore, in this study, we use multiple aggregators such as min, max, and mean aggregators used in the PNA model simultaneously.

Degree Scalers. In GMN, degree scalers amplify or decrease signals based on a node's degree, allowing for more flexibility in the model. The formula for this is S (scaling factor) = α (amplification

factor) x d (node degree) / delta (average degree in the training set). In our research, we tested amplification factors of -1, 0, and 1, which respectively result in attenuation, no change, and signal amplification based on the node's degree.

$$S(d,\alpha) = \left(\frac{\log(d+1)}{\delta}\right)^{\alpha}, d > 0, -1 < \alpha < 1 \tag{1}$$

Combining Aggregators. Like in the PNA, we incorporate various aggregators and degree scalers using the equation \otimes (Tensor product) and \oplus (general aggregation function) in the GMN framework to enhance the network's flexibility.

$$\oplus = \begin{pmatrix} I \\ S(D, \alpha = 1) \\ S(D, \alpha = -1) \end{pmatrix} \otimes \begin{pmatrix} Max \\ Min \\ Mean \end{pmatrix}$$
(2)

2.3 MLP-Mixer

In this work, we used the Mixer block recommended in MLP-Mixer, which is used for the image task, as the Mixer layer in GMN. Input,x first goes through the layer norm operation, then transposes this matrix by the T operator. The transposed x passes through the Linear layer as x^T and is again transposed and transformed into x. It is then merged with the input x with the residual connection. Then it passes through Layer Norm, linear layer, and residual connection again. MLP-Mixer block can be summarized as the following equation 3:

$$Mix = MLP_2(LayerNorm((MLP_1((LayerNorm(x))^T))^T + x)) + x$$
 (3)

2.4 Graph Mixer Network

As in PNA, after h_i^k and h_j^k , which are the features of the neighboring node, are concatenated and pass through the linear layer containing learnable parameters, they pass through multiple aggregation and scalars. We take the embedding features obtained as a result of this process here as x. This x enters the Mixer block, just like the MLP-Mixers performed in the image task. Figure 1 shows the components of the Mixer Blocks, and Equation 3 shows the Mix operator. as a result, the update function h_i^{k+1} is obtained as in Equation 4.

$$h_i^{k+1} = Mix(h_i^k, \oplus (h_i^k, h_j^k, e_{ij}^k)) \tag{4} \label{eq:4}$$

3 Experiments

We evaluated the performance of GMN models on ZINC [3] dataset. The performance results of GMN were compared with the Message Passing Neural Networks (MPNN) [6] and Graph Transformer.

3.1 Dataset

Our method was trained using the ZINC dataset, which is a dataset for predicting the solubility of chemical compounds through graph regression. The compounds in the dataset are represented as graphs, with atoms as nodes and bonds between atoms as edges. The ZINC dataset includes 12,000 molecules, with atom numbers ranging from 9 to 37. The performance of the method was evaluated using the mean absolute error (MAE) metric.

3.2 Results

We trained models using the multiple aggregator(s) such as mean, max, min. Finally, we compared our results with the well-known baseline methods in the literature. The results are given in Table 1. Our results have shown improved performance over attention based mechanism methods.

Table 1: Benchmarking GMN on ZINC dataset.

Models	ZINC(MAE)
GCN [1]	0.367
GAT [5]	0.384
MPNN [6]	0.288
Graph Transformer [12]	0.226
GMN (ours)	0.212
GraphGPS [13]	0.070

4 Discussion and Conclusion

The experiment shows that it is possible to create a powerful transformer-style graph regressor without using attention layers. Additionally, the MLP-Mixer model has a significant advantage over the Graph Transformer in terms of complexity [9, 10], as it is linear in relation to the sequence length instead of quadratic. This is achieved through the use of an intermediate projection dimension within the feed-forward layer applied to aggregated learnable embeddings.

The performance of the GMN is lower than GraphGPS because the positional encoding [14] and the number of parameters are quite low. In addition to its lower performance, the main disadvantage of the MLP-Mixer model is that it can only operate on sequences of a fixed length (as a result of the feed-forward layer applied to aggregated learnable embeddings). While this is not a problem in the image domain, it can be a limitation for graph neural networks because graphs do not have a fixed data structure.

This study shows that MLP-Mixers are effective for graph regression. Future research should focus on understanding the specific roles of other parts of the MLP-Mixer, such as interpretability or initialization scheme. Additionally, the report hopes to inspire further investigation into the underlying reasons for the performance of current models.

Acknowledgements

The author express his gratitude to Ahmet S. Rifaioğlu, Gökhan Özsarı and Mehmet Volkan Atalay not just only for writing and encouragin him to write this paper but also giving or the valuable insights and discussions.

References

- [1] Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. *arXiv preprint arXiv:1609.02907*, 2016. 1, 4
- [2] Prithviraj Sen, Galileo Namata, Mustafa Bilgic, Lise Getoor, Brian Galligher, and Tina Eliassi-Rad. Collective classification in network data. *AI magazine*, 29(3):93–93, 2008. 1
- [3] John J Irwin, Teague Sterling, Michael M Mysinger, Erin S Bolstad, and Ryan G Coleman. Zinc: a free tool to discover chemistry for biology. *Journal of chemical information and modeling*, 52 (7):1757–1768, 2012. 1, 3
- [4] Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural networks? *arXiv preprint arXiv:1810.00826*, 2018. 1
- [5] Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio, and Yoshua Bengio. Graph attention networks. *arXiv preprint arXiv:1710.10903*, 2017. 1, 4
- [6] Justin Gilmer, Samuel S Schoenholz, Patrick F Riley, Oriol Vinyals, and George E Dahl. Neural message passing for quantum chemistry. In *International conference on machine learning*, pages 1263–1272. PMLR, 2017. 1, 3, 4
- [7] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al.

- An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020. 2
- [8] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 2
- [9] Ilya O Tolstikhin, Neil Houlsby, Alexander Kolesnikov, Lucas Beyer, Xiaohua Zhai, Thomas Unterthiner, Jessica Yung, Andreas Steiner, Daniel Keysers, Jakob Uszkoreit, et al. Mlp-mixer: An all-mlp architecture for vision. *Advances in Neural Information Processing Systems*, 34: 24261–24272, 2021. 2, 4
- [10] Luke Melas-Kyriazi. Do you even need attention? a stack of feed-forward layers does surprisingly well on imagenet. *arXiv preprint arXiv:2105.02723*, 2021. 2, 4
- [11] Gabriele Corso, Luca Cavalleri, Dominique Beaini, Pietro Liò, and Petar Veličković. Principal neighbourhood aggregation for graph nets. *arXiv preprint arXiv:2004.05718*, 2020. 2
- [12] Vijay Prakash Dwivedi and Xavier Bresson. A generalization of transformer networks to graphs. *arXiv preprint arXiv:2012.09699*, 2020. 4
- [13] Ladislav Rampášek, Mikhail Galkin, Vijay Prakash Dwivedi, Anh Tuan Luu, Guy Wolf, and Dominique Beaini. Recipe for a general, powerful, scalable graph transformer. *arXiv preprint arXiv:2205.12454*, 2022. 4
- [14] Vijay Prakash Dwivedi, Anh Tuan Luu, Thomas Laurent, Yoshua Bengio, and Xavier Bresson. Graph neural networks with learnable structural and positional representations. arXiv preprint arXiv:2110.07875, 2021. 4