

KS-GNNExplainer: Global Model Interpretation Through Instance Explanations On Histopathology images

Sina Abdous, Reza Abdollahzadeh, Mohammad Hossein Rohban

Sharif University of Technology
`{sina.abdous,re.abd,rohban}@sharif.edu`

Abstract. Instance-level graph neural network explainers have proven beneficial for explaining such networks on histopathology images. However, there has been few methods that provide model explanations, which are common patterns among samples within the same class. We envision that graph-based histopathological image analysis can benefit significantly from such explanations. On the other hand, current model-level explainers are based on graph generation methods that are not applicable in this domain because of no corresponding image for their generated graphs in real world. Therefore, such explanations are communicable to the experts. To follow this vision, we developed KS-GNNExplainer, the first instance-level graph neural network explainer that leverages current instance-level approaches in an effective manner to provide more informative and reliable explainable outputs, which are crucial for applied AI in the health domain. Our experiments on various datasets, and based on both quantitative and qualitative measures, demonstrate that the proposed explainer is capable of being a global pattern extractor, which is a fundamental limitation of current instance-level approaches in this domain.

Keywords: Graph Neural Networks · Explainability · Diagnosis · Histopathology.

1 Introduction

Recent advances in machine learning, and particularly deep learning, have transformed histopathological image understanding. However, such advantages come at the cost of a reduced transparency in decision-making processes [8]. Given the importance of reasoning in any clinical decision, pathologists often seek explanations for the deep learning decisions. Therefore, inspired by the explainability techniques on natural images [16], several explainers have been implemented in the digital pathology, such as feature attribution based methods [12] but, inter-entity interactions are ignored using these methods, which is destructive in the pathologist’s diagnosis process.

To address the aforementioned concerns, a histological image is represented as an entity graph, with nodes and edges indicating biological entities and inter-entity interactions, respectively to make explainability possible in the entity

space, such as cells. The graph is then used to form a Graph Neural Network (GNN) [11], which aims to solve the graph classification problem. Subsequently, graph explainers [4] are employed, which could highlight the responsible entities for the final diagnosis, providing pathologists with intuitive explanations.

Recently several graph explainers have been proposed, which can be classified into two categories: instance-level and model-level explanations [19]. Instance-level explanations focus on explaining the model prediction for a given graph instance, which means their inability to consider a batch of samples and extract common patterns like tumor subregions among similar samples. Meanwhile, current model-level explanation methods aim at understanding the general behavior of the model using graph generation methods, which is not applicable in the context of explaining GNNs for histopathology images, as there is no corresponding image for the generated graphs in real world.

In this paper, first, a framework to compare explainers for the histopathology images is presented. Then, based on our findings from the proposed framework, we present a novel model-level explainer, which addresses the limitations of methods in both categories. The core idea of the proposed benchmark framework is to gradually remove the nodes that are declared as the most important ones in making the decision by an explainer, and assess whether the model output changes in response to such removals. By doing this, we get a cumulative distribution on the number of instances that their predicted labels change. Repeating the same procedure for the least important nodes, we get a second cumulative distribution. We then propose to use the Kolmogorov-Smirnov (KS) test on the two empirical cumulative distributions to quantify the effectiveness of the explainer. One would expect a shift between these two distributions in case that the most and least important nodes have been correctly identified. We show that our validation framework is sound in the sense that the explainers with high proposed KS score could be leveraged to enhance the diagnosis accuracy much better than other explainers.

Then, we present KS-GNNExplainer which is built upon GNNExplainer, a competent instance-level explainer of GNNs [10], and extend it to be a model-based explainer by incorporating (1) pairwise embedding similarities among explanation graphs with the same label; and (2) the mentioned KS-score in its objective function. Here, we seek subgraphs, one in each instance of a batch of data, with similar structures (global model explainers) that optimize the KS score, once removed from the graphs. We show the effectiveness of our approach on a variety of datasets.

In summary, our key contributions are as follows:

1. Introducing a simple method based on two sample Kolmogorov-Smirnov (KS) test for benchmarking graph neural networks explainers on histopathology images.
2. Introducing a model-level explainer (KS-GNNExplainer), which distinguishes itself from all previous works by introducing a novel objective function that suits histopathology images.

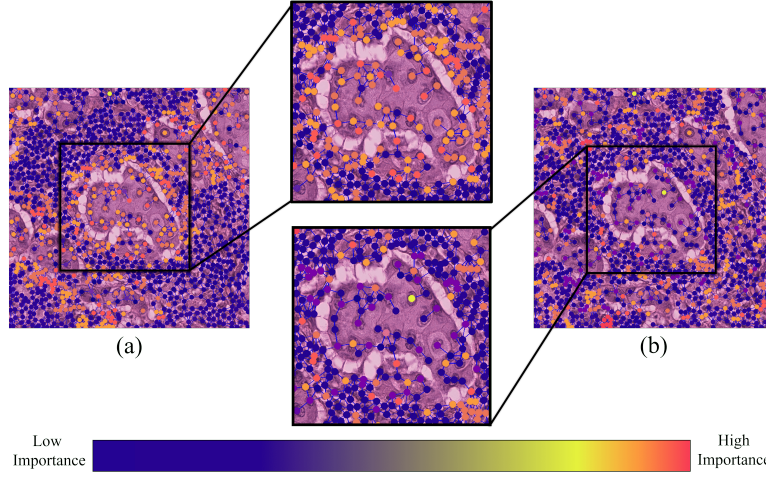


Fig. 1. Explanation graph produced for a ductal carcinoma in situ RoI by (a) KS-GNNExplainer and (b) GNNExplainer as current state-of-the-art method [10]. Nuclei involved in lymph-vascular invasion (confirmed by a pathologist) are detected in a more integrated appearance using our proposed method.

2 Background

2.1 Notations

Cell-graphs (CG) created from the pathology images are defined as an undirected graph $G_{CG} = (V, E, H)$ composed of vertices V , and edges E . An embedding $h \in \mathbb{R}^d$ describes each vertex, which is collectively expressed in its matrix form as $H \in \mathbb{R}^{|V| \times d}$. A symmetric adjacency matrix $A \in \mathbb{R}^{|V| \times |V|}$ describes the graph topology, where $A_{u,v} = 1$ if an edge exists between the vertices u and v .

We use the Hover-Net [7] for nuclei segmentation. The cell graphs are then made based on the segmented nuclei as their nodes. Each node is then connected to its k -NNs. The node features are image patches of a fixed size around the nuclei. Finally, HACT-Net [13], is employed to build fixed-size graph embeddings from the CGs, which are fed into a Multi-Layer Perceptron (MLP) to predict the cancer stage labels.

2.2 Metrics

Here, we used Fidelity+^{acc} score [19], and call it Fidelity onward, for studying faithfulness of different methods by masking the nodes that are detected as the most important ones in the explanation and expect that the predictions of model would change under the mentioned masking.

$$\text{Fidelity} = \frac{1}{N} \sum_{i=1}^N (\mathbb{1}(\hat{y}_i = y_i) - \mathbb{1}(\hat{y}_i^{1-m_i} = y_i)), \quad (1)$$

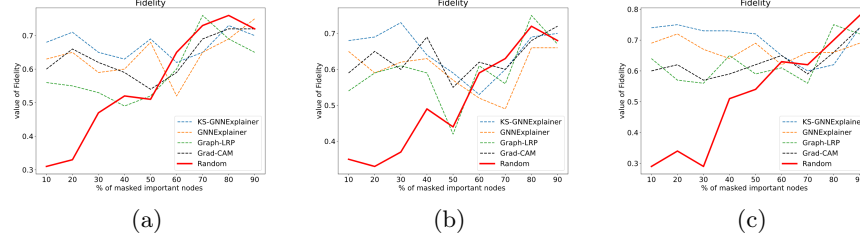


Fig. 2. Performance of studied explainers in addition to random explainer using Fidelity on (a) BRACS, (b) BACH, and (c) CRC. The effectiveness of KS-GNNE explainer on masking up to nearly 50% of important nodes can be demonstrate.

Examining Eq. 1, y_i is original prediction of the i -th input graph, and N is the total number of input graphs. Furthermore, $\hat{y}_i^{1-m_i}$ denotes prediction of i -th graph after masking its m_i most important nodes. Higher values for Fidelity indicate more discriminative features being identified. As the generated importance map for nuclei distribution using studied explainers are continuous, we plot the Fidelity score based on different thresholds of importance in Fig. 2.

3 KS-Bench framework

As the first step, we propose a new framework to benchmark explainers on the histopathology images. Unlike previous work [9], our benchmark evaluates explanation qualities, which is then used to improve the classification. For this purpose, we have considered graph-based explainers such as GNNE explainers [18], GraphLRP [14], Grad-CAM [15], and Grad-CAM++ [6].

The main idea in KS-Bench is the gradual removal of the most and the least important vertices that are reported by each explainer, and then observing changes of the output label of the model. Note that unlike the image inputs, node and edge removal often do not result in artifacts and distribution shifts in the input graphs, and hence does not prohibit the use of the original GNN model in assessing the prediction of masked inputs. This is also due to the fact that the high variety of graph structures and sizes are used in the training phase. More specifically, we remove the most, and the least important nodes reported by each explainer in two separate phases, gradually (from 0% to 100% by 5% increments). Then, we obtain two cumulative distributions for each explainer, corresponding to each phase. At any percentage point j , these distributions represent the proportion of the inputs that has experienced any change in the model output label after masking lower than j percent of their nodes so far. We envision KS statistics compared to existing metrics like Fidelity, gives a higher discrimination power to evaluate explainers as it checks the effect of masking both most and least important components simultaneously while these two type of components are assessed separately in the case of Fidelity (Fidelity+ and Fidelity-),

Table 1. The KS statistic, p-value, and class accuracies of HACT-Net when incorporating the explainer importance score in the top 30% nodes in the BRACS dataset. Note that the first row is baseline accuracy of HACT-Net on BRACS dataset.

Explainer	KS value	p-value	Normal	Benign	UDH	ADH	FEA	DCIS	Invasive
None	-	-	60.2%	53.3%	52.5%	51.4%	65.1%	56.2%	67.7%
GNNExplainer	0.748	2.04e−10	68.2 %	58.5%	56.1%	59.6%	63.7%	68.9%	72.4%
GraphLRP	0.552	2.17e−12	57.9%	52.0%	54.1%	54.2%	56.1%	60.2%	64.9%
Grad-CAM	0.615	2.21e−7	63.5%	54.1%	54.7%	57.9%	58.2%	64.7%	65.4%
Grad-CAM++	0.658	1.93e−8	66.9%	57.2%	55.3%	61.0%	60.8%	65.3%	68.1%

which makes it hard to interpret if only either of the metrics improves. We used the two-sample Kolmogorov–Smirnov (KS) test to compare these two cumulative distributions, which will assign a KS value to each explainer. Using KS test, the ideal explainer is the one that the prediction of the trained model changes for a significant number of sample images after removing first few most important nodes. Also, the label is expected to remain almost unchanged for most of images by removing the least important nodes.

To validate our benchmarking scheme, we incorporate the node importance score given by each explainer as a feature in the HACT-Net to potentially improve the accurate classification. We then correlate the performance of the modified HACT-Net with the proposed KS score. More specifically, in the modified HACT-Net, we extend the node features to include a flag of whether that node is among the 30% of most important nodes. The results are reported in Table 1. An important observation here is that GNNExplainer achieves the highest KS value, and also its importance scores improve the HACT-Net by a significant amount in most classes. We see a similar trend in the Grad-CAM++, which is the runner up according to the KS values. This observation shows that our KS benchmark gives an indication on the usefulness of an explainer in improving the classification .

4 KS-GNNExplainer

Based on our findings from proposed benchmark in previous section, we reformulate the objective function of GNNExplainer as our baseline method to interpret the whole model while maintaining its simplicity as an instance-based explainer. GNNExplainer as a perturbation-based method, tries to find an explanation subgraph which maximizes mutual information between the true label and the subgraph using an iterative mask generation algorithm. So for the task of generating explanation for a single graph, we considered a batch of graphs with the same label, and leveraged the objective function from a simple optimization problem which maximizes mutual information between the explanation subgraphs and label to a general optimization problem as illustrated in Fig. 3.

$$\max_{G_1, \dots, G_k} \sum_{i=1}^k MI(y_i, G_i) + \lambda_1 \sum_{i,j=1}^k P(G_i, G_j) + \lambda_2 \sum_{i=1}^k ks_i - \lambda_3 Var, \quad (2)$$

where $Var = variance(ks_1, \dots, ks_k)$ and KS value for every graph in each iteration is calculated using two-sample KS test and is as follows:

$$ks_i = \sup_x |F_{1,n}(x) - F_{2,m}(x)|, \quad (3)$$

where n_1 and n_2 are observations, and $F_{1,n}$ and $F_{2,m}$ are the empirical cumulative distribution functions of the first and the second distributions, respectively.

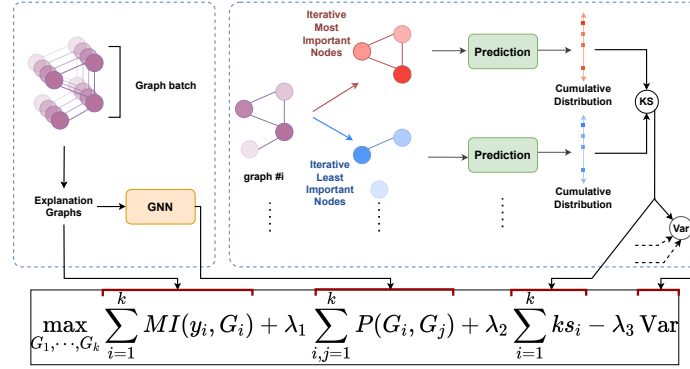


Fig. 3. The objective function of KS-GNNExplainer. It maximizes three components (mutual information, pairwise similarity, and summation of KS value) among all samples, in addition to minimizing the variance of obtained KS values.

Examining Eq. (2), it maximizes summation of four terms. First term is sum of mutual information for all samples in the batch, which is the main idea of GNNExplainer at the single instance level. Continuous relaxation suggested in [18] has applied here to overcome the discrete optimization difficulty. Inspired by objective function of [17], the second term is sum of pairwise similarity for each two subgraphs, based on their cosine similarity of their embeddings from the HACT-Net model, denoted by $P(.,.)$. This term ensures that explanation subgraphs of inputs from the same class are similar, which is step towards a model-level explanation. The third term is sum of KS value for each subgraph in each step. Our objective function also tries to minimize variance between all KS values as we would expect a uniformly large KS values across samples of batch.

λ_1 , λ_2 , and λ_3 are coefficients that determine the relative importance of different terms, and optimal values for them are empirically set using the validation set of the BRACS.

5 Experiments

5.1 Datasets

We focus on prediction of the Breast cancer subtypes. So we evaluate our framework on public datasets, BRACS [5] and BACH [2], as well as a nuclei annotated dataset, BreCaHAD [1]. To further examine our method, we used Colorectal Cancer Grading (CRC) dataset [3] as well.

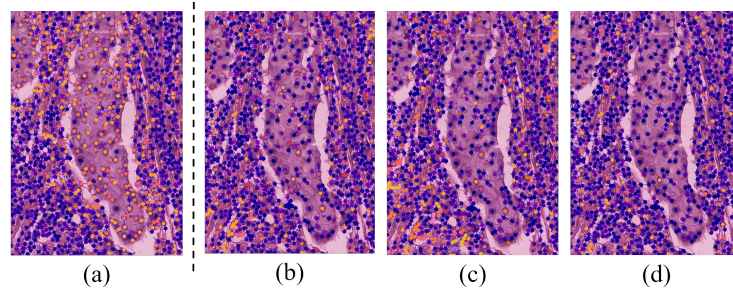


Fig. 4. Qualitative results on ductal carcinoma in situ RoI sample of BRACS. Darker colors correspond to lower importance. Each column correspond to an explainer: (a) KS-GNNExplainer, (b) GNNExplainer, (c) Grad-CAM, and (d) GraphLRP.

5.2 Qualitative analysis

We conduct a set of experiments to determine how each explainer performs in recognizing lymph-vascular invasion as a critical sign in cancer grading as well as discriminating tumor from non-tumor nuclei. Fig. 4 shows the nuclei importance map of studied explainers. It shows effectiveness of our method over others in detecting lymph-vascular invasion as a critical sign of DCI subtype of breast cancer. We envision the superior performance in extracting common patterns among instances is due to our model-level objective function through applying both pairwise similarity and KS test.

5.3 Quantitative results

We also use the mentioned metrics to further analyze previous explainers quantitatively. First, we plot changes in the Fidelity for all explainers, as well as a random one, for different thresholds on each dataset in Fig. 2. As can be seen, superior performance of the KS-GNNExplainer over other methods is evident.

In addition, the macro-averaged F1 score results are shown in Table 2. Here, we considered a binary classification of tumor vs. non-tumor nuclei through the importance score. We assign nuclei with the importance score of 0.5 and greater as tumor, and else as non-tumor, and compare against the ground truth nuclei label. Here, our KS-GNNExplainer method shows superior macro averaged F1 score compared to the other explainers.

Table 2. The macro F1 score for the nuclei classification on BreCaHAD.

EXPLAINER	PER-CLASS F1 SCORE		MACRO-AVERAGED F1 SCORE
	TUMOR	NON-TUMOR	
KS-GNNEXPLAINER	0.725	0.669	0.697
GNNEXPLAINER	0.692	0.682	0.687
GRAPHLRP	0.426	0.553	0.489
GRAD-CAM	0.595	0.503	0.549

Ablation study: We conduct a thorough study to show how each term contributes to our objective function. We draw the following observations from Table 3: (1) Considering our objective function, by removing the KS test part (third and fourth terms), a marginal decrease in Fidelity score can be seen and is a proof of our effective method (Exceptions like the one for BRACS, are a fine case study for future works); (2) The Fidelity score dropped slightly by removing the second term, shows the idea of making the embeddings of subgraphs belonging to the same subclass is also beneficial. Our findings indicate that unifying all terms in our method, significantly leverages the GNNExplainer by through extracting common patterns from samples.

Table 3. Ablation study on different components of the proposed objective function.

Dataset	$term_{MI}$	$term_P$	$term_{KS}$	$sumation$	$term_{KS}$	$variance$	Fidelity score
BACH	✓	✓				✓	0.68
	✓	✗			✗		0.32 (↓ 0.36)
	✓	✓	✗		✗		0.49 (↓ 0.19)
	✓	✓	✓		✗		0.52 (↓ 0.16)
	✓	✗	✓		✓		0.56 (↓ 0.12)
BRACS	✓	✓				✓	0.63
	✓	✗			✗		0.28 (↓ 0.35)
	✓	✓	✗		✗		0.40 (↓ 0.23)
	✓	✓	✓		✗		0.67 (↑ 0.04)
	✓	✗	✓		✓		0.45 (↓ 0.18)

6 Conclusion

We proposed KS-GNNExplainer, the first model-level graph neural network explainer inspired by instance-level methods. The key contribution of the KS-GNNExplainer arises from the insight that current instance-level explainers, lack a comprehensive approach to look for similar information in different samples simultaneously without examining the model directly. We have proved its efficiency through our experiments on different datasets.

References

1. Aksac, A., Demetrick, D.J., Ozyer, T., Alhajj, R.: Brecahad: a dataset for breast cancer histopathological annotation and diagnosis. *BMC research notes* **12**(1), 1–3 (2019) [7](#)
2. Aresta, G., Araújo, T., Kwok, S., Chennamsetty, S.S., Safwan, M., Alex, V., Marami, B., Prastawa, M., Chan, M., Donovan, M., Fernandez, G., Zeineh, J., Kohl, M., Walz, C., Ludwig, F., Braunewell, S., Baust, M., Vu, Q.D., To, M.N.N., Kim, E., Kwak, J.T., Galal, S., Sanchez-Freire, V., Brancati, N., Frucci, M., Riccio, D., Wang, Y., Sun, L., Ma, K., Fang, J., Kone, I., Boulmane, L., Campilho, A., Eloy, C., Polónia, A., Aguiar, P.: Bach: Grand challenge on breast cancer histology images. *Medical Image Analysis* **56**, 122–139 (2019). <https://doi.org/10.1016/j.media.2019.05.010>, <https://www.sciencedirect.com/science/article/pii/S1361841518307941> [7](#)
3. Awan, R., Sirinukunwattana, K., Epstein, D., Jefferyes, S., Qidwai, U., Aftab, Z., Mujeeb, I., Snead, D., Rajpoot, N.: Glandular morphometrics for objective grading of colorectal adenocarcinoma histology images. *Scientific Reports* **7**(1), 16852 (Dec 2017). <https://doi.org/10.1038/s41598-017-16516-w>, <https://doi.org/10.1038/s41598-017-16516-w> [7](#)
4. Baldassarre, F., Azizpour, H.: Explainability techniques for graph convolutional networks. *International Conference on Machine Learning Workshops* (2019) [2](#)
5. Brancati, N., Anniciello, A.M., Pati, P., Riccio, D., Scognamiglio, G., Jaume, G., De Pietro, G., Di Bonito, M., Foncubierta, A., Botti, G., Gabrani, M., Feroce, F., Frucci, M.: Bracs: A dataset for breast carcinoma subtyping in h&e histology images (2021). <https://doi.org/10.48550/ARXIV.2111.04740>, <https://arxiv.org/abs/2111.04740> [7](#)
6. Chattopadhyay, A., Sarkar, A., Howlader, P., Balasubramanian, V.N.: Gradcam++: Generalized gradient-based visual explanations for deep convolutional networks. In: *2018 IEEE Winter Conference on Applications of Computer Vision (WACV)*. pp. 839–847 (2018). <https://doi.org/10.1109/WACV.2018.00097> [4](#)
7. Graham, S., Vu, Q.D., Raza, S.E.A., Azam, A., Tsang, Y.W., Kwak, J.T., Rajpoot, N.: Hover-net: Simultaneous segmentation and classification of nuclei in multi-tissue histology images (2018). <https://doi.org/10.48550/ARXIV.1812.06499>, <https://arxiv.org/abs/1812.06499> [3](#)
8. Holzinger, A., Malle, B., Kieseberg, P., Roth, P.M., Müller, H., Reihs, R., Zatloukal, K.: Towards the augmented pathologist: Challenges of explainable-ai in digital pathology. *ArXiv abs/1712.06657* (2017) [1](#)
9. Jaume, G., Pati, P., Anklin, V., Foncubierta, A., Gabrani, M.: Histocartography: A toolkit for graph analytics in digital pathology. In: Atzori, M., Burlutskiy, N., Ciompi, F., Li, Z., Minhas, F., Müller, H., Peng, T., Rajpoot, N., Torben-Nielsen, B., van der Laak, J., Veta, M., Yuan, Y., Zlobec, I. (eds.) *Proceedings of the MICCAI Workshop on Computational Pathology*. *Proceedings of Machine Learning Research*, vol. 156, pp. 117–128. PMLR (27 Sep 2021), <https://proceedings.mlr.press/v156/jaume21a.html> [4](#)
10. Jaume, G., Pati, P., Bozorgtabar, B., Foncubierta-Rodríguez, A., Feroce, F., Anniciello, A.M., Rau, T., Thiran, J., Gabrani, M., Goksel, O.: Quantifying explainers of graph neural networks in computational pathology. *CoRR abs/2011.12646* (2020), <https://arxiv.org/abs/2011.12646> [2](#), [3](#)
11. Kipf, T.N., Welling, M.: Semi-Supervised Classification with Graph Convolutional Networks. In: *Proceedings of the 5th International Conference on Learning Representations*. *ICLR '17* (2017), <https://openreview.net/forum?id=SJU4ayYgl> [2](#)

12. Korbar, B., Olofson, A.M., Miraflor, A.P., Nicka, C.M., Suriawinata, M.A., Torresani, L., Suriawinata, A.A., Hassanpour, S.: Looking under the hood: Deep neural network visualization to interpret whole-slide image analysis outcomes for colorectal polyps. In: 2017 IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW). pp. 821–827 (2017). <https://doi.org/10.1109/CVPRW.2017.114> 1
13. Pati, P., Jaume, G., Fernandes, L.A., Foncubierto-Rodríguez, A., Feroce, F., Annicciello, A.M., Scognamiglio, G., Brancati, N., Riccio, D., Di Bonito, M., De Pietro, G., Botti, G., Goksel, O., Thiran, J.P., Frucci, M., Gabrani, M.: Hact-net: A hierarchical cell-to-tissue graph neural network for histopathological image classification. In: Sudre, C.H., Fehri, H., Arbel, T., Baumgartner, C.F., Dalca, A., Tanno, R., Van Leemput, K., Wells, W.M., Sotiras, A., Papiez, B., Ferrante, E., Parisot, S. (eds.) *Uncertainty for Safe Utilization of Machine Learning in Medical Imaging, and Graphs in Biomedical Image Analysis*. pp. 208–219. Springer International Publishing, Cham (2020) 3
14. Schwarzenberg, R., Hübner, M., Harbecke, D., Alt, C., Hennig, L.: Layerwise relevance visualization in convolutional text graph classifiers. In: Ustalov, D., Somasundaran, S., Jansen, P., Glavas, G., Riedl, M., Surdeanu, M., Vazirgiannis, M. (eds.) *Proceedings of the Thirteenth Workshop on Graph-Based Methods for Natural Language Processing, TextGraphs@EMNLP 2019*, Hong Kong, November 4, 2019. pp. 58–62. Association for Computational Linguistics (2019). <https://doi.org/10.18653/v1/D19-5308>, <https://doi.org/10.18653/v1/D19-5308> 4
15. Selvaraju, R.R., Cogswell, M., Das, A., Vedantam, R., Parikh, D., Batra, D.: Grad-cam: Visual explanations from deep networks via gradient-based localization. In: 2017 IEEE International Conference on Computer Vision (ICCV). pp. 618–626 (2017). <https://doi.org/10.1109/ICCV.2017.74> 4
16. Simonyan, K., Vedaldi, A., Zisserman, A.: Deep inside convolutional networks: Visualising image classification models and saliency maps. *CoRR abs/1312.6034* (2013), <http://dblp.uni-trier.de/db/journals/corr/corr1312.html#SimonyanVZ13> 1
17. Wang, X., Shen, H.W.: GNNInterpreter: A Probabilistic Generative Model-Level Explanation for Graph Neural Networks. In: *Proceedings of the 5th International Conference on Learning Representations. ICLR '23* (2023), <https://openreview.net/forum?id=rqq6Dh8t4d> 6
18. Ying, R., Bourgeois, D., You, J., Zitnik, M., Leskovec, J.: GNNExplainer: Generating Explanations for Graph Neural Networks. Curran Associates Inc., Red Hook, NY, USA (2019) 4, 6
19. Yuan, H., Yu, H., Gui, S., Ji, S.: Explainability in graph neural networks: A taxonomic survey. *CoRR abs/2012.15445* (2020), <https://arxiv.org/abs/2012.15445> 2, 3