

1 We thank the reviewers for their time, valuable feedback, and recommendations for improvements. Overall, the
 2 reviewers found our methodology interesting, novel and technically sound, and our contributions to be very timely.
 3 However, a key point of clarification was raised regarding the selection of hyperparameters and the effects of subgraph
 4 constraints on generated explanations. We address this key point of clarification in detail below. The four reviewers also
 5 raised important clarification points on the motivation for the use of mutual information (R5), synthetic datasets (R2,
 6 R3, R5) and quantitative experiments (R3, R5, R6), and we provide clarification on these issues below as well. These
 7 clarification issues arise—in large part—because certain details were omitted from the paper due to space constraints.
 8 However, an extra page and lengthened appendix will allow us to address these clarification points in the revised version.

9 **Parameters & effects of constraints on explanations.** R2, R3 and R5 rightly point out the need for more
 10 investigation into hyperparameters and regularization constraints. In new
 11 experiments, we observe that varying regularization strength of constraints
 12 can affect explanations (Fig. S1). However, this gives GNNEXPLAINER
 13 flexibility to encode domain-specific priors into constraints and, crucially,
 14 allows the explainer to balance verbosity and completeness of the explana-
 15 tion. In particular, without regularization, the explainer will include many
 16 edges in the explanation even though those edges only negligibly increase
 17 confidence of a GNN’s prediction. Conversely, imposing a very small size
 18 constraint K_M will produce meaningless explanations, e.g., single edges.
 19 In practice, we observe that prior knowledge about a task (e.g., select K_M
 20 to be the size of a chemical functional group) or commonly used heuris-
 21 tics (e.g., select K_M that results in largest increase of GNN’s confidence
 22 score) determine an appropriate value for K_M . Furthermore, the range of
 23 acceptable explanation sizes is also important in determining the threshold
 24 for cutting off low-importance edges (e.g. how much time a user has to
 25 examine an explanation). In all experiments, we ensure that the setup is
 26 fair for all methods and use the same value for explanation size.

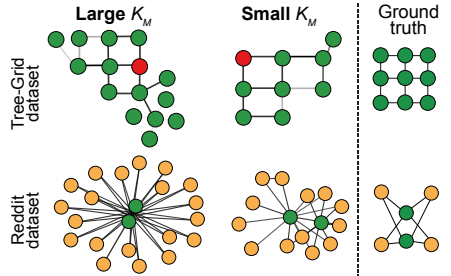


Figure S1: Effects of constraint K_M on explanations. On Reddit, we find that modifying K_M doesn’t have strong effects, one can easily see the explanatory structure. Effects are stronger when these structures are of fixed size (e.g., 3x3 grid) rather than general shapes (e.g., star or cycle of any size). Further, on Tree-Grid dataset, GNNEXPLAINER is robust to selection of hyperparameter K_M ; AUCs for $K_M = 6, 8, 10, 12, 14, 16$ are 0.71, 0.83, 0.88, 0.89, 0.86, 0.81, respectively.

27 **Mutual information.** R5 raises a key point about motivation behind using mutual information. This is indeed
 28 fundamental, but there is intuition for using mutual information: The goal is to identify a small subgraph together
 29 with a small subset of node features that maximize confidence of a GNN’s prediction. Mutual information finds
 30 such a subgraph as a connected subgraph of GNN’s computation graph. Here, the restriction to be “connected” is
 31 on computation graph rather than the original graph, which is intuitive. Even GNNs that capture information from
 32 far-away nodes through message passing (e.g., Deep Graph Infomax, Jumping Knowledge Networks) have computation
 33 graphs that are different (i.e., multi-hop neighborhoods) but are always connected even though corresponding nodes in
 34 the original graph are not necessarily connected. Further, unlike counterfactual reasoning, mutual information-based
 35 objective allows users to understand what graph structure is critical for a particular prediction in a way that gives a
 36 concise, semantically meaningful explanation. For instance, in Question-Answer graphs of Reddit threads (Fig. S1)
 37 we see that explanations have “2-3 high degree nodes that simultaneously connect to many low degree nodes,” which
 38 translates to the explanation of “2-3 experts who all answer many different questions in a QA thread on Reddit.”

39 **Synthetic datasets.** R3 and R5 raise an important clarification point regarding the description of synthetic datasets.
 40 This is important as these datasets allow us to quantify the quality of explanations without necessitating manual curation,
 41 and thus they represent an advance over the prior art on explaining predictions, which often only provide hand-picked
 42 real-world examples. We acknowledge that the data generation procedure was not adequately explained in the draft and
 43 we will include more information on datasets/groundtruth in the main text with details in the Appendix. Further, we
 44 will open-source the code for data generation together with GNNEXPLAINER’s code. Briefly, we first generate a base
 45 graph (e.g., Barabási-Albert graph). For a random set of nodes, we then attach a particular structure (e.g., a house- or
 46 grid-shaped motif) to each of the nodes. These nodes will have labels that are different others, which gives us a dataset
 47 in which absence/presence of the structure indicates a label. The groundtruth thus corresponds to these structures, and
 48 we calculate precision/recall of edges in an explanation relative to edges in the groundtruth.

49 **New experiments.** In response to the constructive feedback we conducted further experiments whose results confirm
 50 our findings and increase our confidence that GNNEXPLAINER is a general approach for explaining GNNs. Among
 51 others, for R5, we used GNNEXPLAINER to explain GNN models that use simple attention mechanisms and we
 52 observed similarly good performance as when explaining GCNs. Second, in addition to explaining node classification
 53 and graph classification, we also tested GNNEXPLAINER on link prediction (for R3). Experiments on a 2D-Grid dataset,
 54 where each grid graph has 20% of random edges removed, show that explanation for a predicted edge are typically
 55 edges forming a 4-cycle with the predicted edge, consistent with groundtruth. Third, R6 alludes to a hypothetical
 56 situation in which explanation is predictive of a wrong label with high probability. This is very interesting and related to
 57 adversarial attacks, however, we note that this situation was never realized in any of our experiments. That is because, at
 58 GNNEXPLAINER’s initialization, the probability of GNN’s predicted label is already the highest and thus it is unlikely
 59 that explainer would optimize the subgraph for false prediction. We will include all new analyses in Appendix.

60 Finally, reviewers raise minor points, such as a typo “ K_M nodes” by R5 (it should be K_M edges) and a suggestion
 61 by R6 to investigate how explanations look if their size is constrained by the number of nodes instead of edges (our
 62 experiments show such constraints result in explanations being fully induced subgraphs). We will add formal discussions
 63 to the Appendix, including explicit discussion on future directions and adversarial attacks.