
AdaGrid: Adaptive Grid Search for Link Prediction Training Objective

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 One of the most important factors which contribute to the success of a machine
2 learning model is a good training objective. Training objective crucially influences
3 a model’s performance and generalization capabilities. The automated process
4 of designing a good training objective involves optimizing a machine learning
5 process, therefore can be viewed as a meta-learning problem. In this paper, we
6 specifically focus on graph neural network training objectives for link prediction,
7 which has not been explored in the existing literature. Here, the training objective
8 includes, among others, training mode, negative sampling strategy, and various
9 hyperparameters, such as edge message ratio. Commonly, these hyperparame-
10 ters are fine-tuned by complete grid search, which is very time-consuming and
11 model-dependent. To mitigate these limitations, we propose Adaptive Grid Search
12 (AdaGrid), which dynamically adjusts the edge message ratio during training. It is
13 model agnostic and highly scalable with a fully customizable computational budget.
14 Through extensive experiments, we show that AdaGrid can boost the performance
15 of the models up to 1.9%, while can be nine times more efficient than a complete
16 search. Overall, AdaGrid represents an effective automated algorithm for designing
17 machine learning training objectives.

18 1 Introduction

19 Link prediction is one of the most important tasks on graph-structured data. For a given pair of
20 entities, the goal of link prediction is predicting whether they are going to interact. Applications
21 of link prediction are found in various fields, such as social networks, recommender systems, and
22 biology. There have been many strategies to cope with the link prediction task [19], where the state-
23 of-the-art approaches use Graph Neural Networks (GNNs) [12, 10, 36]. GNNs capture dependencies
24 via message passing between the nodes of graphs and are suitable for various graph-related tasks.
25 To train GNNs, an appropriate training objective has to be selected, which includes choice of the
26 objective function, evaluation metric, and training strategy.

27 The automated process of designing a good training objective involves optimizing a machine learning
28 process, therefore can be viewed as a meta-learning problem. In this paper, we specifically focus
29 on graph neural network training objectives for link prediction, which has not been explored in
30 the existing literature. Here, link prediction training objectives encompass training mode, negative
31 sampling strategy, and hyperparameters such as edge message ratio; they crucially influence the
32 performance of a model and its generalization capabilities on the link prediction task, so they have
33 to be chosen carefully. Notably, setting an appropriate edge message ratio is notoriously hard since
34 there are many relevant factors that determine the optimal edge message ratio: data split ratio, model
35 configuration, the average clustering coefficient of the network, and the average shortest path length
36 of the network.

37 The standard approach to fine-tune link prediction objectives is a complete search over some hyperpa-
38 rameter space [37, 29]. Complete search trains more models for full epochs, each with a different edge
39 message ratio from a set of predefined values, and selects the one with the maximal final validation
40 AUC. Even though it exhaustively searches hyperparameter space, which is very time-consuming, it
41 still obtains non-optimal performance. To make matters worse, even when the configuration of the
42 model slightly changes, a complete search has to be repeated.

43 In this paper, we propose Adaptive Grid Search (AdaGrid) which adjusts edge message ratio during
44 training to each specific model and dataset, therefore alleviates the above-mentioned limitations.
45 AdaGrid is model agnostic and highly scalable with a fully customizable computational budget.
46 Through extensive experiments, we show that AdaGrid can boost the performance of the models
47 up to 1.9%, while can be nine times more efficient than a complete search. We also propose a new
48 negative sampling strategy, which is based on the community structure of the network. Unlike the
49 standard uniform negative sampling, it samples more difficult negative instances and prevents abuse
50 of community structure of networks.

51 Our key contributions include: **(1)** We define GNN training modes for link prediction. **(2)** We explore
52 the GNN training objective for link prediction. **(3)** We propose Adaptive Grid Search (AdaGrid),
53 which resolves limitations of complete search. **(4)** We introduce community ratio-based negative
54 sampling, which samples more difficult negative instances and prevents the abuse of the community
55 structure of networks.

56 2 Related Work

57 Central to our research are graph neural network models, to which we tailor training objectives for
58 the link prediction tasks. Thus, we provide a brief overview of general graph neural networks and
59 their application to link prediction tasks.

60 **General graph neural networks** In recent years, numerous graph neural network (GNN) models
61 have been proposed. Some of them are inspired by convolutional neural networks [3, 6, 7, 10, 12,
62 25, 15, 21, 28], recurrent neural networks [16, 31], recursive neural networks [2, 8], and loopy
63 belief propagation [5]. The majority of these approaches can be united under the Graph Network
64 framework proposed by Battaglia et al. [1]. In the Graph Network framework, a GNN is considered
65 as a message-passing algorithm, where representations of nodes, edges, and the global attribute are
66 iteratively updated from neighboring nodes, edges, and the global attribute using a differentiable,
67 order-invariant aggregation function. A recent extensive overview of recent GNN models, general
68 design pipeline for constructing such models and their applications is provided by Zhou et al. [36].

69 **Link prediction with graph neural networks** GNNs have been applied to a wide variety of tasks,
70 including, but not limited to, node classification [10, 12], link prediction [22], graph classification
71 [5, 7, 35], and chemoinformatics [20, 18, 9, 11]. In the context of link prediction, a few approaches
72 have been introduced recently [14]. Schlichtkrul et al. [22] proposed a relational graph convolutional
73 neural network model (R-GCN) which incorporates GCN [12] as the building block to model
74 relational data, knowledge graph in particular. It encodes nodes such that a node’s latent representation
75 depends on all neighboring nodes as well as the node itself. Kipf et al. [13] introduced a variational
76 graph autoencoder (VGAE) framework that learns latent representation on graph-structured data. The
77 encoder part, based on GCNs, embeds nodes in latent space, then, the decoder part tries to reconstruct
78 the adjacency matrix of the graph by the inner product of latent representations. Zhang et al. [32]
79 proposed a novel framework called Weisfeiler-Lehman Neural Machine (WLNLM), which is based
80 on the Weisfeiler-Lehman algorithm that labels nodes of the graph according to the topology of the
81 underlying graph. For each prospective node pair, WLNLM extracts a subgraph in their neighborhood
82 and encodes it as an adjacency matrix. Based on these adjacency matrices, WLNLM then trains a
83 classifier. Furthermore, Zhang et al. [33] introduced SEAL, a new heuristic learning paradigm, which
84 captures first, second, and higher-order structural information in the form of local subgraphs. It
85 unifies local subgraph, embedding, and attribute information using the graph convolutional network.

86 Although numerous approaches for link prediction with GNN models have been proposed, their
87 training objective is usually vaguely specified. This inspired us to define graph neural network
88 training objectives for link prediction and research their properties.

89 **3 Preliminaries**

90 We represent graph as $G = (V, E)$, where $V = \{v_1, \dots, v_n\}$ is the node set and $E \subseteq V \times V$ is the
 91 edge set. Nodes can be paired with features $X = \{x_v : v \in V\}$. Given a set of labeled node pairs
 92 $D = \{(u_1, v_1, y_1), (u_2, v_2, y_2), \dots\}$, where y_i denotes presence/absence of the edge between the
 93 node pair $(u_i, v_i) \in V \times V$, the goal of link prediction is to learn a mapping $f : V \times V \rightarrow \{0, 1\}$
 94 that for a given pair of nodes predicts whether they are likely to be connected by an edge.

95 **Graph neural networks** In this paper we approach the link prediction problem by message passing
 96 GNNs [27]. The goal of GNNs is to learn meaningful node embeddings h_v , which are based on
 97 iterative aggregation of local neighborhoods and should be informative for the task in question. The
 98 k -th iteration/layer of message passing can be written as [30]:

$$m_u^{(k)} = \text{MSG}^{(k)}(h_u^{(k-1)}), \tag{1}$$

$$h_v^{(k)} = \text{AGG}^{(k)}(\{m_u^{(k)}, u \in \mathcal{N}(v)\}, h_v^{(k-1)}), \tag{2}$$

99 where $h_v^{(k)}$ is the node embedding after k iterations, $h_v^0 = x_v$, $m_v^{(k)}$ is message embedding, and $\mathcal{N}(v)$
 100 is the local neighborhood of node v . Definitions of $\text{MSG}^{(k)}(\cdot)$ and $\text{AGG}^{(k)}(\cdot)$ depend on version of
 101 the GNN. The final node embeddings $h_v = h_v^{(K)}$ can be then used for various node, edge, and graph
 102 level tasks. In the case of link prediction, the embeddings can be optimized so that the following
 103 equation resembles the probability of an edge between nodes u and v :

$$P((u, v) \in E) = \sigma(h_u \cdot h_v), \tag{3}$$

104 where σ represents sigmoid function and \cdot denotes dot product.

105 **Standard learning-based link prediction experimental setting** Learning-based link prediction
 106 task is frequently formulated as binary classification, where potential edges are classified as true or
 107 false. The standard learning-based link prediction setting [33] first splits the graph’s edge set E into a
 108 training set, validation set, and test set. The ratio between the cardinalities of these sets is called split
 109 ratio (i.e. training/validation/test split ratio or training/validation split ratio). For link prediction using
 110 GNN models, each set of edges have to be further divided into message-passing edges and objective
 111 edges. Message-passing edges are used for the propagation of information between nodes, while
 112 loss is calculated based on objective edges. During the evaluation of the model’s performance, all
 113 validation and test edges are used only for the estimation of loss. For propagation of information,
 114 validation utilizes training message-passing edges, while testing exploits all training and validation
 115 edges. So far, objective edges comprise only of edges, which are present in the graph – they are
 116 positive instances or positive edges for the link prediction task. Therefore, the node pairs that are
 117 not connected by an edge are required for a fair binary classification task. These so-called negative
 118 instances or negative edges are usually obtained by uniformly randomly sampling nonexistent edges
 119 (unconnected node pairs). The standard learning-based link prediction setting samples negative
 120 instances until each set of objective edges is balanced – there is the same number of positive and
 121 negative instances.

122 **Training objective for link prediction** When training a machine learning model, one has to specify
 123 model evaluation function, objective function for optimization, as well as strategy how to use the
 124 available data. All of these components are captured under the term training objective. Models
 125 are usually evaluated by different metrics, while optimization is performed based on loss functions.
 126 On the other hand, strategies define how data is split, which instances are presented to the model
 127 in which epoch, how to sample negative instances, and so on. In the case of link prediction using
 128 GNNs, the latter includes also the choice of appropriate training mode, edge message ratio, and
 129 negative sampling. While evaluation and objective functions are well studied, not much attention has
 130 been given to training objective strategies. Therefore, we provide the *first* categorization of different
 131 training modes, which are also supported by DeepSNAP [24] – a popular network library for GNNs.

132 **Training modes** For link prediction using GNN models, edges have to be split into message-passing
 133 edges and objective edges. This split can be performed in different ways, each corresponding to a
 134 different training mode. We recognize two main training modes for link prediction: `all` mode and
 135 `disjoint` mode. The `all` mode exploits the whole edge set for message passing and calculation

136 of loss, so each edge is used for propagation of information as well as for updating weights of the
137 model. Such training mode is employed by graph autoencoders [13]. However, prior researchers have
138 identified shortcomings of the autoencoder approaches [34]: The setup potentially leads to overfitting
139 of the GNN model, since the model is never asked to predict edges that are not part of input, and thus
140 tend to only memorize the message-passing edges. The `disjoint` mode resolves this by dividing
141 edges in a set for message passing, which is disjoint from the set of objective edges. In this case, each
142 edge belongs to only one of the roles. Consequentially, `disjoint` mode introduces an additional
143 hyperparameter – *edge message ratio*, which control the proportion of message-passing edges. For
144 instance, if there are 70 message-passing edges and 30 objective edges, the edge message ratio equals
145 to 0.7. Normally, the split into objective and message-passing edges is performed prior to the training.
146 Nevertheless, there exists an extension to the `disjoint` mode – the `resample disjoint` mode,
147 which resamples both sets every fixed number of epochs to allow training to take full advantage of all
148 edges for both purposes.

149 4 Proposed Method: Adaptive Grid Search (AdaGrid)

150 The key idea behind our approach is that we dynamically adapt training objectives during training
151 to each configuration of the model and each dataset. We first compare training modes, highlight
152 limitations of training objective, and expose disadvantages of complete search, which is the standard
153 approach for selection of GNN training objective for link prediction (Section 4.1). We then propose
154 Adaptive Grid Search (AdaGrid), which resolves drawbacks of complete search (Section 4.2). At the
155 end, we introduce a novel negative sampling method, which takes into consideration the community
156 structure of the network and creates more challenging negative instances (Section 4.3).

157 4.1 Limitations of Training Objective

158 We first compare training modes to establish which one yields the best results. As is shown in
159 Appendix C.1, the `resample disjoint` mode consistently outperforms the `all` mode and the
160 `vanilla disjoint` mode according to final validation AUC. However, additional experiments disclose
161 that the `resample disjoint` mode is superior only when combined with an appropriate edge
162 message ratio. Selection of bad edge message ratio results in performance considerably worse than
163 the `all` mode. So the `resample disjoint` mode improves the model but in exchange for some
164 extra work with an additional hyperparameter tuning. This suggests that the edge message ratio has
165 to be fine-tuned carefully. The most straightforward and standard approach is to do a complete search
166 on some predefined values, but this is very time-consuming. Based on the above observations, it is
167 beneficial to estimate a good edge message ratio quickly. However, it turns out that the edge message
168 ratio is notoriously hard to set. Appendix C.2 illustrates that there are no patterns that depend on: data
169 split ratio, model configuration, the average clustering coefficient of the network, and the average
170 shortest path length of the network.

171 **Complete search** Since a good training objective is specific to each configuration of the model and
172 dataset, the standard approach to finding a good edge message ratio is a complete search over some
173 predefined set of values Q with cardinality $L = |Q|$ [29]. Complete search is very time-consuming
174 and computationally demanding because the model has to be trained multiple times with different
175 edge message ratios. A possible speedup is to train each version of the model only for part of all
176 epochs and then take the one which performs best at that time. This relaxation assumes that the
177 best model is already evident before all models are trained for full epochs. While this seems as an
178 adequate solution, there are still certain limitations. Even by slightly changing the configuration of
179 the model, the optimal edge message ratio changes and complete search has to be repeated, which is
180 extremely inconvenient. Another, possible false, assumption of this approach is that the optimal edge
181 message ratio does not change during training. To resolve these drawbacks of the standard approach,
182 we propose Adaptive Grid Search (AdaGrid).

183 4.2 Adaptive Grid Search (AdaGrid)

184 **Description** The key feature of AdaGrid is its ability to adapt edge message ratio during training to
185 each configuration of the model and each dataset. AdaGrid changes the edge message ratio every
186 *adapt_epochs*. Then, it trains L copies of the model with different edge message ratios in parallel,
187 where each copy is trained for *try_epochs*. A set of predefined edge message ratios Q which are
188 taken into consideration can be chosen according to preference. After *try_epochs* of training, the

189 edge message ratio is selected based on one of two possible criteria: *validation criterion* and *gap*
 190 *criterion*. The validation criterion selects the edge message ratio which corresponds to the model with
 191 the highest final validation AUC. Since sometimes validation AUC is a bit unstable during training
 192 also smoothing can be performed – instead of the final validation AUC, rather average of a few last
 193 validation AUCs is taken. On the other hand, the gap criterion chooses the edge message ratio of the
 194 model with minimal absolute difference between the final training and validation AUC. Analogously,
 195 also gap criterion can utilize smoothing. One of the main upsides of AdaGrid is its flexibility since
 196 it can be adjusted to each application separately. The *adapt_epochs* and the *try_epochs* regulate
 197 training time, while selection criterion and set of considered edge message ratios Q can be tailored to
 198 each specific task.

199 **Computational budget of AdaGrid** By setting the *adapt_epochs* and the *try_epochs* appropri-
 200 ately, AdaGrid’s computational budget can be fully customizable. If model is trained for *full_epochs*
 201 and L edge message ratios are considered, overall number of training epochs of AdaGrid is:

$$\text{number of epochs} = \text{full_epochs} \cdot \left(1 + \frac{(L - 1) \cdot \text{try_epochs}}{\text{adapt_epochs}} \right), \quad (4)$$

202 while complete search requires:

$$\text{number of epochs} = \text{full_epochs} \cdot L. \quad (5)$$

203 For instance, if *try_epochs* = *adapt_epochs* both required the same number of epochs, while if
 204 *try_epochs* = 5, *adapt_epochs* = 50, and $L = 9$, AdaGrid needs five times fewer training epochs.

205 4.3 Community Ratio-based Negative Sampling

206 Uniform negative sampling seems to be an appropriate approach for obtaining negative instances,
 207 however, after careful analysis, it turns out that these instances are rather simple negative examples
 208 for link prediction. A lot of networks inherently display some kind of community structure – the
 209 network can be partitioned into disjoint communities so that connections within communities are
 210 denser than the connections with the rest of the network. Let us define edges that have both endpoints
 211 in the same community as the *within community edges*, and edges that have endpoints in the different
 212 communities as the *between communities edges*.

213 **Limitations of uniform negative sampling** In Appendix B.1, we have theoretically shown that,
 214 under certain assumptions, uniform negative sampling yields mostly between community edges. The
 215 same findings are confirmed empirically in Appendix B.2. This can be exploited by a naive model
 216 which does not even consider nodes’ features to perform unfairly well. By always predicting there is
 217 an edge between u and v for node pairs from the same community, and always predicting there is not
 218 one for node pairs from different communities, classification accuracy far beyond baseline 0.5 can be
 219 obtained.

220 **Community ratio-based negative sampling** To improve uniform negative sampling, we propose
 221 a negative sampling which also considers the community structure of the network. First, let us
 222 define *community ratio* as the proportion of node pairs with nodes from the same community. Our
 223 community ratio-based negative sampling obtains negative edges in such a way that sets of negative
 224 and positive instances can not be distinguished based on community ratio. It first performs community
 225 detection on the graph with all training edges, then it measures the community ratio on validation
 226 edges. Afterward, negative instances for all three sets are sampled in compliance with the gauged
 227 community ratio. Our approach is beneficial because community structure can not be abused anymore.
 228 At the same time, it generates more challenging negative instances, which can better differentiate the
 229 performance of models.

230 5 Experiments

231 We test the performance of AdaGrid on various model configurations, datasets, data split ratios,
 232 negative sampling strategies, and hyperparameter settings of AdaGrid. Firstly, we describe datasets,
 233 the model configuration used in experiments, baselines against which AdaGrid is compared, and
 234 experimental set-up. In Section 5.1, we interpret the results of AdaGrid on the standard uniform
 235 negative sampling setting, while Section 5.2 explains the benefits of both, community ratio-based
 236 negative sampling and AdaGrid. In the end, in Section 5.3 we unfold why AdaGrid performs better
 237 than the baselines.

Table 1: AUC in percent for AdaGrid, complete search, and random search with uniform negative sampling evaluation. Model has $K = 2$, $o = 64$, and is trained in the `resample disjoint` mode.

Methods	Datasets						
	Cora			CiteSeer			
	20/40/40	50/25/25	80/10/10	20/40/40	50/25/25	80/10/10	
Complete search	94.85	96.65	97.17	95.84	97.78	98.46	
Random search	94.71	96.40	97.02	95.88	97.67	98.45	
AdaGrid							
<i>adapt_epochs</i>	<i>try_epochs</i>						
100	1	95.03	96.89	97.54	95.96	97.82	98.69
100	5	95.03	96.87	97.63	95.91	97.80	98.64
100	100	94.93	96.98	97.59	95.97	97.79	98.67
50	1	95.01	96.98	97.68	95.94	97.81	98.72
50	5	95.05	97.04	97.75	95.87	97.81	98.68
50	50	95.03	97.06	97.76	95.97	97.83	98.71
10	1	95.01	97.07	97.90	95.79	97.89	98.70
10	5	95.14	97.10	97.83	95.81	97.83	98.71
10	10	95.10	97.09	97.86	95.94	97.88	98.71
Gain		0.29	0.45	0.73	0.09	0.11	0.26

238 **Datasets** All experiments in the paper are conducted on Cora and CiteSeer datasets [23]. Additionally, limitations of training objective experiments use also the PubMed dataset [23]. All three datasets
 239 are well-known citation networks (additional details are available in the Appendix A).
 240

241 **Model configuration** Experiments were conducted using the following model configuration. Our
 242 model consists of K GCN layers applied sequentially, where before each GCN layer there is a dropout
 243 of 0.2, and after each layer but the last there is a ReLU activation. As input it accepts nodes' features
 244 $h_v^{(0)} = x_v \in \mathbb{R}^d$ and it outputs final hidden representations $h_v = h_v^{(K)} \in \mathbb{R}^o$. All intermediate
 245 hidden representations have the same dimensionality as the final representation: $h_v^{(i)} \in \mathbb{R}^o$ for
 246 $i = 1, 2, \dots, K - 1$. The model predicts the probability of an edge between nodes u and v according
 247 to Equation (3). In experiments, d depends solely on the dimensionality of the graph's features, while
 248 o is a hyperparameter. The model is always trained for 500 epochs using binary cross-entropy as
 249 loss function, however, the quality of the model is rather measured by AUC metric. Its parameters
 250 are optimized by stochastic gradient descent (SGD) with the learning rate of 0.1, the momentum of
 251 0.9, and the weight decay of $5 \cdot 10^{-4}$. During training, a cosine annealing [17] schedule is used for
 252 alternation of the learning rate.

253 **Baselines** To contextualize the empirical results of our approach, we compare AdaGrid against two
 254 baselines: complete search and random search. Complete search trains more models for full epochs,
 255 each with a different edge message ratio from a set of predefined values, and selects the one with the
 256 maximal final validation AUC. It exhaustively searches hyperparameter space, which makes it very
 257 time-consuming. It also does not change the edge message ratio during training. On the other hand,
 258 random search modifies edge message ratio to a random value from $[0.1, 0.9]$ interval after every
 259 training epoch. This makes it very fast, however, different edge message ratios are not inspected. We
 260 show that both baselines perform inferior to AdaGrid: complete search has static edge message ratio
 261 and is slow, while random search does not explore edge message ratio space.

262 **Experimental set-up** When comparing AdaGrid with the standard complete search and random
 263 search, we are interested principally in absolute performance and the trade-off between training time
 264 and performance. Experiments are systematically conducted over various settings, which include
 265 different model configurations, datasets, data split ratios, and negative samplings. We use Cora and
 266 CiteSeer datasets. To get more representative results we test the model with two configurations:
 267 $K = 2$, $o = 64$ and $K = 3$, $o = 128$. The model is always trained in the `resample disjoint`
 268 mode, however, data is every time divided according to other split ratios: 20/40/40, 50/25/25, and
 269 80/10/10. Models are trained and evaluated using the standard uniform sampling as well as commu-
 270 nity ratio-based negative sampling, proposed in Section 4.3. In a community ratio-based negative
 271 sampling setting, community detection is performed using the Clauset-Newman-Moore greedy modu-

Table 2: AUC in percent for AdaGrid, complete search, and random search with community ratio-based negative sampling evaluation. Model has $K = 2$, $o = 64$, and is trained in the `resample disjoint` mode.

Methods	Datasets						
	Cora			CiteSeer			
	20/40/40	50/25/25	80/10/10	20/40/40	50/25/25	80/10/10	
Complete search	84.01	82.98	84.62	83.87	82.22	83.96	
Random search	84.00	82.05	83.82	83.87	82.61	83.65	
AdaGrid							
<i>adapt_epochs</i>	<i>try_epochs</i>						
100	1	84.35	83.05	85.27	84.01	83.20	84.58
100	5	84.52	82.88	85.43	83.92	82.86	84.34
100	100	84.49	83.61	85.57	84.00	83.12	84.75
50	1	84.51	83.48	86.02	83.88	83.11	84.79
50	5	84.59	83.56	85.98	83.90	82.94	84.84
50	50	84.59	83.74	85.74	84.00	82.99	84.81
10	1	84.49	83.69	86.27	83.65	83.11	84.90
10	5	84.65	83.82	86.54	83.75	83.37	84.89
10	10	84.59	84.07	86.56	83.94	83.04	84.83
Gain		0.64	1.09	1.94	0.14	0.76	0.94

272 larity maximization algorithm [4]. AdaGrid and complete search both consider only the following set
 273 of edge message ratios: $Q = \{0.1, 0.2, \dots, 0.9\}$. To examine power of AdaGrid, it is assessed with
 274 various configurations of $(adapt_epochs, try_epochs) \in \{10, 50, 100\} \times \{1, 5, adapt_epochs\}$ and
 275 both criteria from Section 4.2. Both criteria utilize smoothing. The selection criterion is considered
 276 a hyperparameter of AdaGrid, so we present results for the better of the two. Each experiment is
 277 repeated three times to mitigate the effect of randomness and the average performance is reported.

278 5.1 AdaGrid and Uniform Negative Sampling

279 We first evaluate AdaGrid on a uniform negative sampling setting, because this is the standard
 280 approach for link prediction. Table 1 contains results of AdaGrid and baselines on $K = 2$ and $o = 64$
 281 model configuration, however, $K = 3$ and $o = 128$ displays similar performance. The table shows
 282 that AdaGrid consistently performs better than the best baseline approach, no matter its configuration.
 283 Improvement of more than 0.7% is especially evident for the 80/10/10 split ratio, which is the most
 284 similar to the usual experimental data splits. Another crucial aspect of AdaGrid is its adjustability
 285 in terms of computational budget. Even when AdaGrid is trained for considerably fewer training epochs
 286 than complete search, it constantly performs better than the best baseline approach. If AdaGrid is
 287 trained with $adapt_epochs = 10$ and $try_epochs = 1$, it requires almost five times fewer epochs
 288 than a complete search and twice as many as a random search, while it can surpass both for more than
 289 0.7%. AdaGrid even outperforms both of them, when $adapt_epochs = 100$ and $try_epochs = 1$.
 290 In this case, AdaGrid is about nine times faster than a complete search and is computationally
 291 almost equivalent to the training of a single model. When decreasing the number of try_epochs , the
 292 performance also slowly decreases, however, a similar performance can be obtained in only a fraction
 293 of computational time. This can be particularly useful for huge models.

294 5.2 AdaGrid and Community Ratio-based Negative Sampling

295 We also evaluate AdaGrid on the proposed community ratio-based negative sampling setting to
 296 display its advantages, as well as advantages of AdaGrid. Table 2 shows that AdaGrid nearly always
 297 performs better than the best baseline approach, regardless of its configuration. By comparing
 298 results of negative sampling strategies, it is evident that community ratio-based negative sampling
 299 creates a more challenging evaluation setting since AUC scores are considerably lower. It also better
 300 differentiates the performance of models, because gains of AdaGrid are always bigger than the ones
 301 of uniform negative sampling. Especially outstanding is the gain of more than 1.9%, which again
 302 confirms the benefits of AdaGrid.

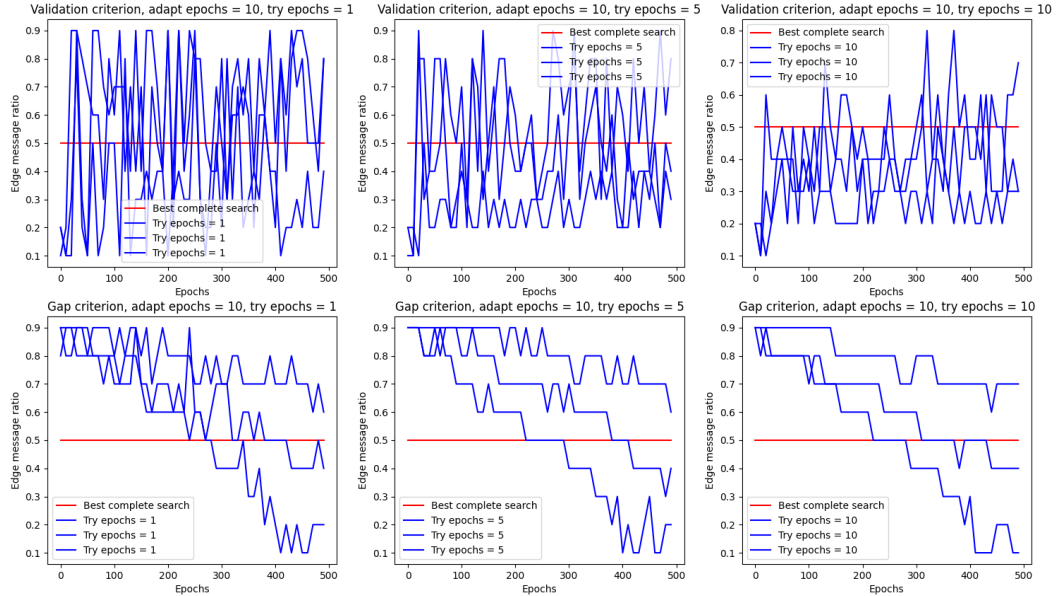


Figure 1: Edge message ratio during training, regulated by validation criterion and gap criterion. Model has $K = 2$ and $o = 64$. It is trained on CiteSeer with 80/10/10 split ratio in the resample disjoint mode.

303 5.3 AdaGrid and Edge Message Ratio

304 We have investigated the properties of AdaGrid, which could explain its improvement of performance.
 305 The success of AdaGrid probably stems from its ability to change edge message ratio during training
 306 in an informed way. According to Figure 1 AdaGrid modifies edge message ratio almost every
 307 *adapt_epochs* in conjunction with both criteria. Additionally, Tables 1 and 2 show that AdaGrid
 308 performs better with a lower number of *adapt_epochs* – it is beneficial to be capable of changing
 309 edge message ratio more frequently. Therefore, complete search results in non-optimal performance,
 310 because it assumes the edge message ratio is static. On the other hand, random search does not
 311 explore edge message ratio space, so it incorrectly alters it.

312 6 Conclusion

313 The aim of the paper is to explore graph neural network training objectives for link prediction. We first
 314 propose community-based negative sampling, which is fairer and harder than the standard uniform
 315 negative sampling. It makes distributions of negative and positive instances indistinguishable based
 316 only on community ratio, as well as disables some naive models which rely solely on community
 317 detection algorithms. We also show that it is very hard to find a suitable training objective, since it is
 318 specific to each dataset and model configuration. To diminish the inconvenience of edge message
 319 ratio fine-tuning, we propose AdaGrid, which adapts edge message ratio “on-the-fly” during training
 320 and overcomes limitations of complete search. It is model agnostic and has a fully customizable com-
 321 putational budget. More importantly, AdaGrid can also reduce training time and boost performance
 322 at once. It can improve the performance of the models up to 1.9%, while can be nine times more
 323 efficient than a complete search.

324 **Future work** Since AdaGrid performs well on link prediction, it would be interesting to apply
 325 AdaGrid to other graph learning tasks. Due to its generality, it would be suitable even for other deep
 326 learning tasks. Instead of edge message ratio, AdaGrid can optimize any parameter which can be
 327 dynamically changed during training. To even further speed up AdaGrid, instead of considering
 328 all edge message ratios only those which are adjacent to the current one can be considered during
 329 adaptation phases. Another intriguing idea is to eliminate a predefined set of values from which edge
 330 message ratios are selected. By locally interpolating validation AUCs near the current edge message
 331 ratio, a new edge message ratio can be chosen so that it maximizes validation AUC.

References

- 332
- 333 [1] Peter W. Battaglia, Jessica B. Hamrick, Victor Bapst, Alvaro Sanchez-Gonzalez, Vinicius Zam-
334 baldi, Mateusz Malinowski, Andrea Tacchetti, David Raposo, Adam Santoro, Ryan Faulkner,
335 Caglar Gulcehre, Francis Song, Andrew Ballard, Justin Gilmer, George Dahl, Ashish Vaswani,
336 Kelsey Allen, Charles Nash, Victoria Langston, Chris Dyer, Nicolas Heess, Daan Wierstra,
337 Pushmeet Kohli, Matt Botvinick, Oriol Vinyals, Yujia Li, and Razvan Pascanu. Relational
338 inductive biases, deep learning, and graph networks, 2018.
- 339 [2] M. Bianchini, M. Gori, and F. Scarselli. Processing directed acyclic graphs with recursive
340 neural networks. *IEEE Transactions on Neural Networks*, 12(6):1464–1470, 2001. doi:
341 10.1109/72.963781.
- 342 [3] Joan Bruna, Wojciech Zaremba, Arthur Szlam, and Yann Lecun. Spectral networks and locally
343 connected networks on graphs. In *International Conference on Learning Representations*
344 (*ICLR2014*), *CBLIS, April 2014*, 2014.
- 345 [4] Aaron Clauset, M. E. J. Newman, and Cristopher Moore. Finding community structure in very
346 large networks. *Phys. Rev. E*, 70:066111, Dec 2004. doi: 10.1103/PhysRevE.70.066111. URL
347 <https://link.aps.org/doi/10.1103/PhysRevE.70.066111>.
- 348 [5] Hanjun Dai, Bo Dai, and Le Song. Discriminative embeddings of latent variable models for
349 structured data. In Maria Florina Balcan and Kilian Q. Weinberger, editors, *Proceedings of The*
350 *33rd International Conference on Machine Learning*, volume 48 of *Proceedings of Machine*
351 *Learning Research*, pages 2702–2711, New York, New York, USA, 20–22 Jun 2016. PMLR.
352 URL <https://proceedings.mlr.press/v48/daib16.html>.
- 353 [6] M. Defferrard, X. Bresson, and P. Vandergheynst. Convolutional neural networks on graphs
354 with fast localized spectral filtering. In *NIPS*, 2016.
- 355 [7] David K Duvenaud, Dougal Maclaurin, Jorge Iparraguirre, Rafael Bombarell, Timothy
356 Hirzel, Alan Aspuru-Guzik, and Ryan P Adams. Convolutional networks on graphs for
357 learning molecular fingerprints. In C. Cortes, N. Lawrence, D. Lee, M. Sugiyama, and
358 R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 28. Cur-
359 ran Associates, Inc., 2015. URL [https://proceedings.neurips.cc/paper/2015/file/](https://proceedings.neurips.cc/paper/2015/file/f9be311e65d81a9ad8150a60844bb94c-Paper.pdf)
360 [f9be311e65d81a9ad8150a60844bb94c-Paper.pdf](https://proceedings.neurips.cc/paper/2015/file/f9be311e65d81a9ad8150a60844bb94c-Paper.pdf).
- 361 [8] Matthias Fey, J. E. Lenssen, F. Weichert, and H. Müller. Splinecnn: Fast geometric deep
362 learning with continuous b-spline kernels. *2018 IEEE/CVF Conference on Computer Vision*
363 *and Pattern Recognition*, pages 869–877, 2018.
- 364 [9] A. Fout, Jonathon Byrd, B. Shariat, and A. Ben-Hur. Protein interface prediction using graph
365 convolutional networks. In *NIPS*, 2017.
- 366 [10] William L Hamilton, Rex Ying, and Jure Leskovec. Inductive representation learning on large
367 graphs. In *Proceedings of the 31st International Conference on Neural Information Processing*
368 *Systems*, pages 1025–1035, 2017.
- 369 [11] Wengong Jin, Connor W. Coley, R. Barzilay, and T. Jaakkola. Predicting organic reaction
370 outcomes with weisfeiler-lehman network. In *NIPS*, 2017.
- 371 [12] Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional
372 networks. *arXiv preprint arXiv:1609.02907*, 2016.
- 373 [13] Thomas N. Kipf and Max Welling. Variational graph auto-encoders, 2016.
- 374 [14] Ajay Kumar, Shashank Sheshar Singh, Kuldeep Singh, and Bhaskar Biswas. Link prediction
375 techniques, applications, and performance: A survey. *Physica A: Statistical Mechanics and its*
376 *Applications*, 553:124289, 2020.
- 377 [15] Tao Lei, Wengong Jin, Regina Barzilay, and Tommi Jaakkola. Deriving neural architectures
378 from sequence and graph kernels. In Doina Precup and Yee Whye Teh, editors, *Proceedings of*
379 *the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine*
380 *Learning Research*, pages 2024–2033. PMLR, 06–11 Aug 2017. URL [https://proceedings.](https://proceedings.mlr.press/v70/lei17a.html)
381 [mlr.press/v70/lei17a.html](https://proceedings.mlr.press/v70/lei17a.html).

- 382 [16] Yujia Li, Daniel Tarlow, Marc Brockschmidt, and Richard Zemel. Gated graph sequence neural
383 networks. *arXiv preprint arXiv:1511.05493*, 2015.
- 384 [17] Ilya Loshchilov and Frank Hutter. SGDR: stochastic gradient descent with warm restarts. In *5th
385 International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26,
386 2017, Conference Track Proceedings*. OpenReview.net, 2017. URL [https://openreview.
387 net/forum?id=Skq89Scxx](https://openreview.net/forum?id=Skq89Scxx).
- 388 [18] A. Lusci, G. Pollastri, and P. Baldi. Deep architectures and deep learning in chemoinformatics:
389 The prediction of aqueous solubility for drug-like molecules. *Journal of chemical information
390 and modeling*, 53 7:1563–75, 2013.
- 391 [19] Linyuan Lü and Tao Zhou. Link prediction in complex networks: A survey. *Physica A: Statistical
392 Mechanics and its Applications*, 390(6):1150–1170, 2011. ISSN 0378-4371. doi: <https://doi.org/10.1016/j.physa.2010.11.027>. URL [https://www.sciencedirect.com/science/
393 //doi.org/10.1016/j.physa.2010.11.027](https://www.sciencedirect.com/science/article/pii/S037843711000991X). URL [https://www.sciencedirect.com/science/
394 article/pii/S037843711000991X](https://www.sciencedirect.com/science/article/pii/S037843711000991X).
- 395 [20] C. Merkwirth and Thomas Lengauer. Automatic generation of complementary descriptors with
396 molecular graph networks. *Journal of chemical information and modeling*, 45 5:1159–68, 2005.
- 397 [21] Mathias Niepert, Mohamed Ahmed, and Konstantin Kutzkov. Learning convolutional neural
398 networks for graphs. In Maria Florina Balcan and Kilian Q. Weinberger, editors, *Proceedings of
399 The 33rd International Conference on Machine Learning*, volume 48 of *Proceedings of Machine
400 Learning Research*, pages 2014–2023, New York, New York, USA, 20–22 Jun 2016. PMLR.
401 URL <https://proceedings.mlr.press/v48/niepert16.html>.
- 402 [22] M. Schlichtkrull, Thomas Kipf, Peter Bloem, Rianne van den Berg, Ivan Titov, and M. Welling.
403 Modeling relational data with graph convolutional networks. *ArXiv*, abs/1703.06103, 2018.
- 404 [23] P. Sen, Galileo Namata, M. Bilgic, L. Getoor, B. Gallagher, and Tina Eliassi-Rad. Collective
405 classification in network data. *AI Mag.*, 29:93–106, 2008.
- 406 [24] Stanford. DeepSNAP, 2021. URL <https://github.com/snap-stanford/deepsnap>.
- 407 [25] Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio, and Yoshua
408 Bengio. Graph attention networks. *arXiv preprint arXiv:1710.10903*, 2017.
- 409 [26] Duncan J Watts and Steven H Strogatz. Collective dynamics of ‘small-world’ networks. *nature*,
410 393(6684):440–442, 1998.
- 411 [27] Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural
412 networks? In *International Conference on Learning Representations*, 2019.
- 413 [28] Rex Ying, Jiaxuan You, Christopher Morris, Xiang Ren, William L Hamilton, and Jure Leskovec.
414 Hierarchical graph representation learning with differentiable pooling. *NeurIPS*, 2018.
- 415 [29] Jiaxuan You, Zhitao Ying, and Jure Leskovec. Design space for graph neural networks. *Advances
416 in Neural Information Processing Systems*, 33, 2020.
- 417 [30] Jiaxuan You, Jonathan Gomes-Selman, Rex Ying, and Jure Leskovec. Identity-aware graph
418 neural networks. *AAAI Conference on Artificial Intelligence*, 2021.
- 419 [31] Victoria Zayats and Mari Ostendorf. Conversation Modeling on Reddit Using a Graph-
420 Structured LSTM. *Transactions of the Association for Computational Linguistics*, 6:121–132,
421 02 2018. ISSN 2307-387X. doi: 10.1162/tacl_a_00009. URL [https://doi.org/10.1162/
422 tacl_a_00009](https://doi.org/10.1162/tacl_a_00009).
- 423 [32] Muhan Zhang and Yixin Chen. Weisfeiler-lehman neural machine for link prediction. In
424 *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery
425 and Data Mining, KDD ’17*, page 575–583, New York, NY, USA, 2017. Association for
426 Computing Machinery. ISBN 9781450348874. doi: 10.1145/3097983.3097996. URL <https://doi.org/10.1145/3097983.3097996>.
427 <https://doi.org/10.1145/3097983.3097996>.
- 428 [33] Muhan Zhang and Yixin Chen. Link prediction based on graph neural networks. In *Advances
429 in Neural Information Processing Systems*, pages 5165–5175, 2018.

- 430 [34] Muhan Zhang and Yixin Chen. Link prediction based on graph neural networks. *Advances in*
431 *Neural Information Processing Systems*, 31:5165–5175, 2018.
- 432 [35] Muhan Zhang, Zhicheng Cui, Marion Neumann, and Yixin Chen. An end-to-end deep learning
433 architecture for graph classification. In *AAAI*, 2018.
- 434 [36] Jie Zhou, Ganqu Cui, Shengding Hu, Zhengyan Zhang, Cheng Yang, Zhiyuan Liu, Lifeng
435 Wang, Changcheng Li, and Maosong Sun. Graph neural networks: A review of methods
436 and applications. *AI Open*, 1:57–81, 2020. ISSN 2666-6510. doi: [https://doi.org/10.1016/](https://doi.org/10.1016/j.aiopen.2021.01.001)
437 [j.aiopen.2021.01.001](https://doi.org/10.1016/j.aiopen.2021.01.001). URL [https://www.sciencedirect.com/science/article/pii/](https://www.sciencedirect.com/science/article/pii/S2666651021000012)
438 [S2666651021000012](https://www.sciencedirect.com/science/article/pii/S2666651021000012).
- 439 [37] Kaixiong Zhou, Qingquan Song, Xiao Huang, and Xia Hu. Auto-gnn: Neural architecture
440 search of graph neural networks. *arXiv preprint arXiv:1909.03184*, 2019.

441 **A Dataset Details**

442 **Cora** Cora dataset [23] consists of 2708 scientific publications from Cora and 5429 citation links.
 443 Each scientific publication belongs to one of the seven classes and is described by a binary-valued
 444 word vector, which indicates the absence/presence of the corresponding words from the dictionary.
 445 The dictionary contains 1433 unique words.

446 **CiteSeer** CiteSeer dataset [23] consists of 3312 scientific publications from CiteSeer and 4732
 447 citation links. Each scientific publication belongs to one of the six classes and is described by a
 448 binary-valued word vector, which indicates the absence/presence of the corresponding words from
 449 the dictionary. The dictionary contains 3703 unique words.

450 **PubMed** PubMed Diabetes dataset [23] consists of 19717 scientific publications from PubMed
 451 related to diabetes and 44338 citation links. Each scientific publication belongs to one of the three
 452 classes and is described by a TF-IDF weighted word vector from a dictionary with a size of 500.

453 **B Limitations of Uniform Negative Sampling**

454 Although uniform negative sampling is the standard approach for the evaluation of link prediction
 455 tasks, it has a serious drawback. We first theoretically show that under certain assumptions uniform
 456 negative sampling creates simple link prediction evaluation (Appendix B.1). Then we confirm
 457 theoretical results also empirically on a real network (Appendix B.2).

458 **B.1 Theoretical Analysis of Uniform Negative Sampling**

459 Imagine that a network with N nodes is partitioned into $R \ll N$ communities of equal size, where
 460 N is a multiple of R . Then for a random node pair (u, v) , it is much more likely that nodes are
 461 from different communities than conversely. Let A denote the event that u and v are from the same
 462 community, and B that they are from different communities:

$$P(B) = 1 - P(A) = 1 - \frac{R \cdot (\frac{N}{R} \cdot (\frac{N}{R} - 1)/2)}{N \cdot (N - 1)/2} \quad (6)$$

$$= 1 - \frac{(\frac{N}{R} - 1)}{(N - 1)} \quad (7)$$

$$\approx 1 - \frac{1}{R}, \quad (8)$$

463 if N and $\frac{N}{R}$ are large enough. Since it makes sense to have multiple communities, it is safe to assume
 464 $P(B) \geq \frac{1}{2}$. The latter estimate is not tight for real networks, because they usually have R at least a
 465 bit larger than 2. For example, Equation (8) yields probability $P(B) = 0.9$ for $R = 10$. Therefore, if
 466 node pairs are sampled uniformly, the majority of negative instance node pairs correspond to nodes
 467 from different communities.

468 The above observation can be exploited so that a baseline model based only on community structure
 469 can already get high classification accuracy. For the following analysis let's assume that positive
 470 edges are equally likely to be within community and between communities edges. This assumption
 471 provides a lower bound on classification accuracy because positive edges should contain way more
 472 within community edges. So by always predicting there is an edge between u and v for node pairs
 473 from the same community, and always predicting there is not one for node pairs from different
 474 communities, the following accuracy score (ACC) is obtained:

$$ACC = P(\text{correct} | NEG) \cdot P(NEG) + P(\text{correct} | POS) \cdot P(POS) \quad (9)$$

$$\approx (1 - \frac{1}{R}) \cdot \frac{1}{2} + \frac{1}{2} \cdot \frac{1}{2}, \quad (10)$$

475 where POS and NEG denote events that an edge belongs to positive instances or negative instances.
 476 Accuracy above is considerably larger than expected 0.5 (e.g. ≈ 0.7 for $R = 10$), which is not
 477 desired.

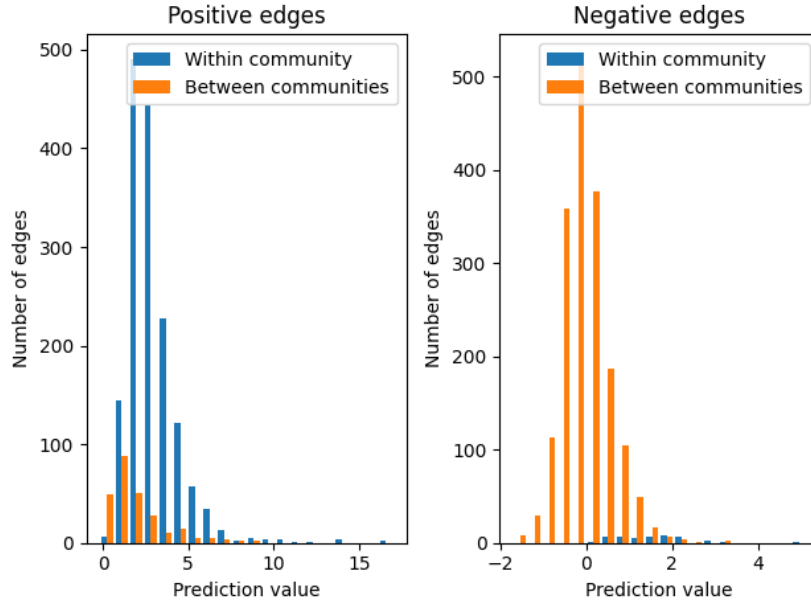


Figure 2: Prediction values (dot product of nodes’ hidden representations) of negative and positive edges before sigmoid activation on CiteSeer dataset. To find edges within and between communities, community detection is performed using the Clauset-Newman-Moore greedy modularity maximization algorithm [4]. Model has $K = 3$, $o = 128$, and $edge_message_ratio = 0.6$. It is trained with 80/20 split ratio in the `resample disjoint` mode.

478 B.2 Empirical Analysis of Uniform Negative Sampling

479 The assumption that all communities are of equal size is usually not true for real networks, however,
 480 on real networks similar findings are discovered empirically. From Figure 2, it is evident that on the
 481 CiteSeer dataset the vast majority of negative instances corresponds to nodes from different commu-
 482 nities and that positive instances are mostly within community edges. It can also be observed that the
 483 model has difficulties with the prediction of negative edges within communities because it assigns
 484 them considerably above-average prediction values. These edges are probably also misclassified by
 485 the model.

486 C Limitations of Training Objective

487 In this section, we explore the GNN training objective for link prediction. We compare different
 488 training modes (Appendix C.1) and highlight limitations of the `resample disjoint` mode, which
 489 is the best-performing training mode (Appendix C.2).

490 C.1 Comparison of Training Modes

491 In this section we compare the following training modes: `all mode`, `disjoint mode`, and
 492 `resample disjoint mode`. These training modes are presented in detail in Section 3. Figure
 493 3 shows that `resample disjoint mode` is superior according to final validation AUC, however,
 494 additional experiments reveal that only if combined with an appropriate edge message ratio. The
 495 `resample disjoint mode` has also the smallest difference between training and validation AUC,
 496 since `resample disjoint mode` significantly reduces overfitting. These findings are consistent
 497 overall datasets from Section 5, so we recommend using the `resample disjoint mode`.

498 C.2 Limitations of the Resample Disjoint Mode

499 Another important parameter that can drastically boost performance is the edge message ratio. In
 500 Figure 4 can be observed that the difference between the best and the worst edge message ratio is
 501 more than 1.4%. This indicates that the edge message ratio should be fine-tuned carefully. The most

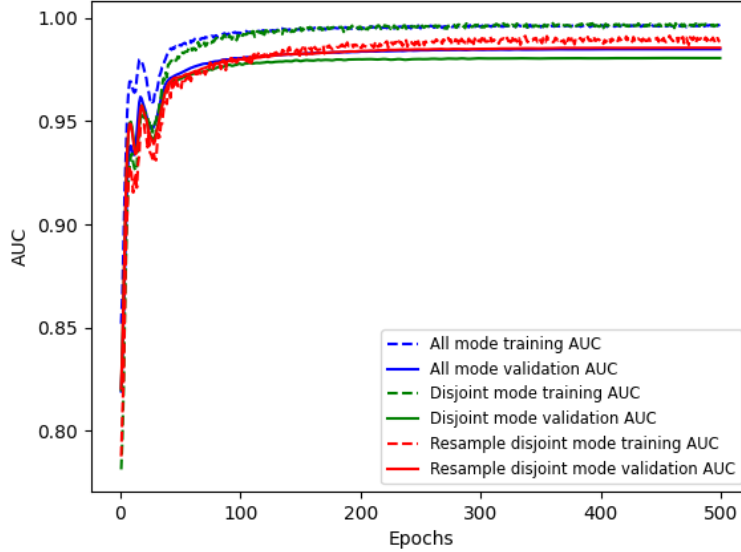


Figure 3: AUC scores on training and validation sets for all mode, disjoint mode, and resample disjoint mode. Model has $K = 2$, $o = 128$, and $edge_message_ratio = 0.8$ for the both disjoint modes. It is trained on CiteSeer dataset with 80/20 split ratio.

502 straightforward approach is to do a complete search on some predefined values, but this is very time
 503 consuming. Another strategy is to randomly change the edge message ratio after each epoch. It
 504 surprisingly achieves almost as good results as a complete search (Figure 4), even though, it has
 505 much lower computational complexity. Based on the above observations, it is beneficial to estimate a
 506 good edge message ratio quickly. However, it turns out that the edge message ratio is notoriously
 507 hard to set. In the following paragraphs is illustrated that there are no patterns that depend on: data
 508 split ratio, model configuration, the average clustering coefficient of the network, and the average
 509 shortest path length of the network.

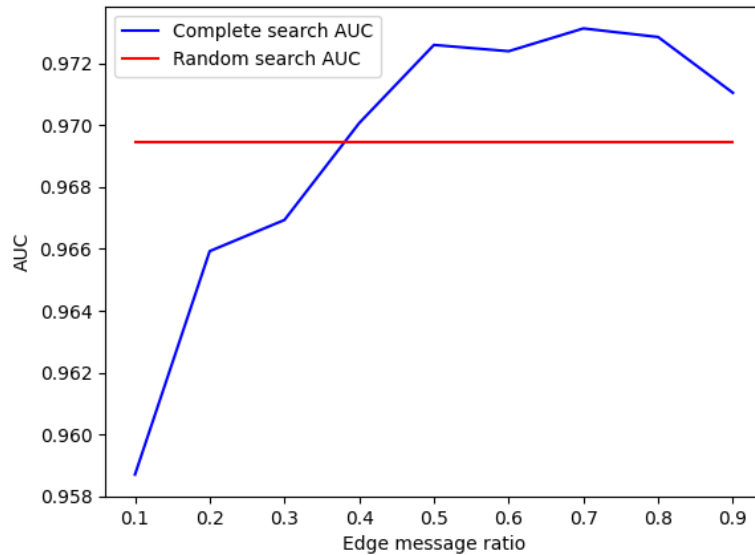


Figure 4: AUC scores for different edge message ratios and selection strategies on Cora dataset. Model has $K = 3$ and $o = 128$. It is trained with 80/10/10 split ratio in the resample disjoint mode.

510 **Data split ratio and model configuration** The data split ratio controls how much information is
511 available for training, so the model should prefer a different edge message ratio when given a new data
512 split ratio. Although the optimal edge message ratio indeed changes, patterns are inconsistent across
513 various model configurations and networks. In Figures 5 and 6, the optimal edge message ratios are
514 plotted with the same configuration of the model but on different networks. The trend of the optimal
515 edge message ratios fundamentally differs, although both datasets correspond to citation networks.
516 Even on the same network, trends are often not coherent across distinct model configurations. Figures
517 6 and 7 display discrepancy between the optimal edge message ratio trends on PubMed network
518 when the models are configured differently.

519 **Average clustering coefficient and average shortest path length of the network** The average
520 clustering coefficient is a local measure, while the average shortest path length is a global one. These
521 two characteristics should impact the optimal edge message ratio. For example, a higher average
522 clustering coefficient implies there are a lot of edge triangles, so by removing some of the edges, it is
523 expected that the distance between a pair of nodes does not drastically increase. This might indicate
524 that networks with a higher average clustering coefficient would not need that high edge message
525 ratio. By similar consideration, a higher average shortest path length should intuitively require a
526 higher edge message ratio. However, Figures 8 and 9 depict that the average clustering coefficient
527 and average shortest path length are uncorrelated with the optimal edge message ratio.

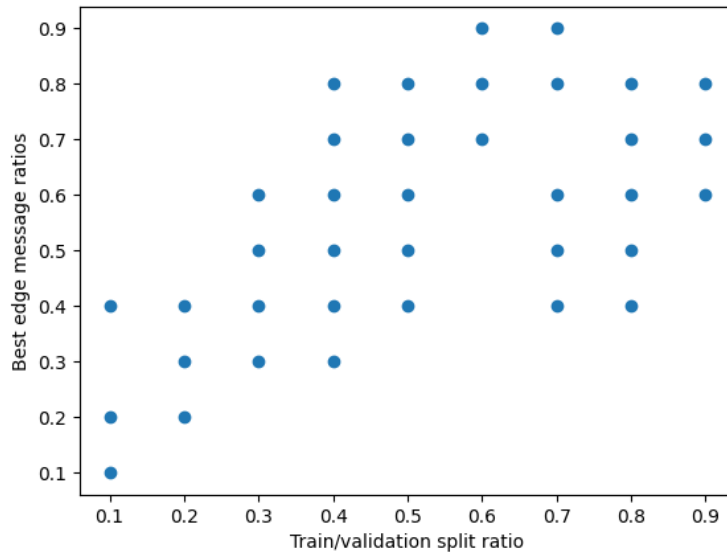


Figure 5: The optimal edge message ratios within 0.1% margin for different training/validation splits. Model has $K = 2$ and $o = 64$. It is trained on CiteSeer dataset in the `resample disjoint` mode.

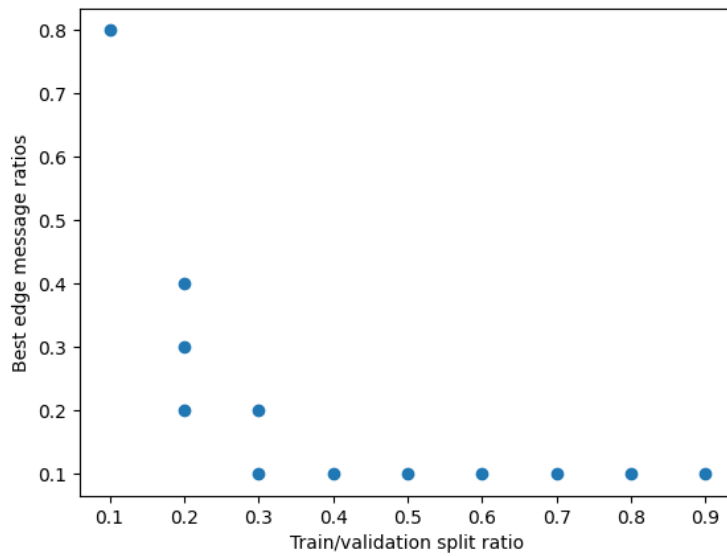


Figure 6: The optimal edge message ratios within 0.1% margin for different training/validation splits. Model has $K = 2$ and $o = 64$. It is trained on PubMed dataset in the `resample disjoint` mode.

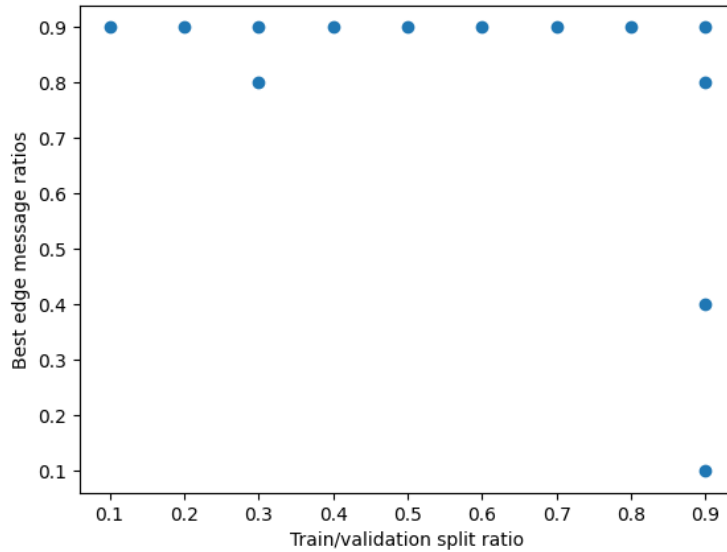


Figure 7: The optimal edge message ratios within 0.1% margin for different training/validation splits. Model has $K = 3$ and $o = 128$. It is trained on PubMed dataset in the `resample disjoint` mode.

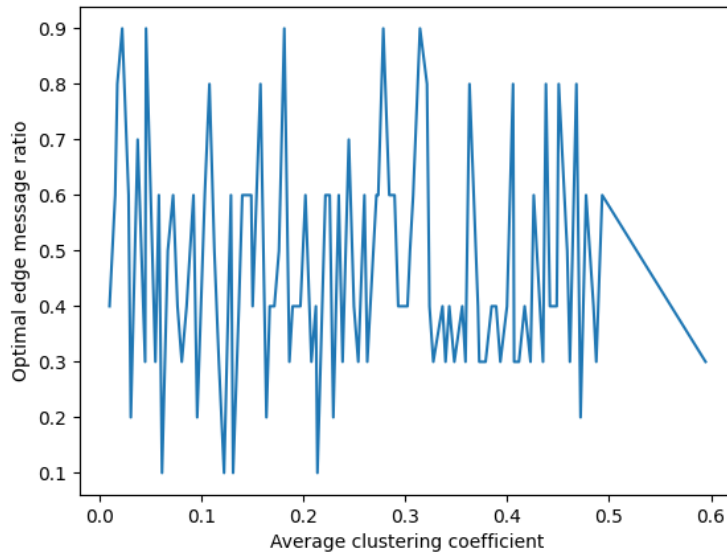


Figure 8: The optimal edge message ratio for different average clustering coefficients on artificial Watts-Strogatz small-world graphs [26] with 256 nodes and average degree 6. Model has $K = 2$ and $o = 32$. It is trained with 80/20 split ratio in the `resample disjoint` mode.

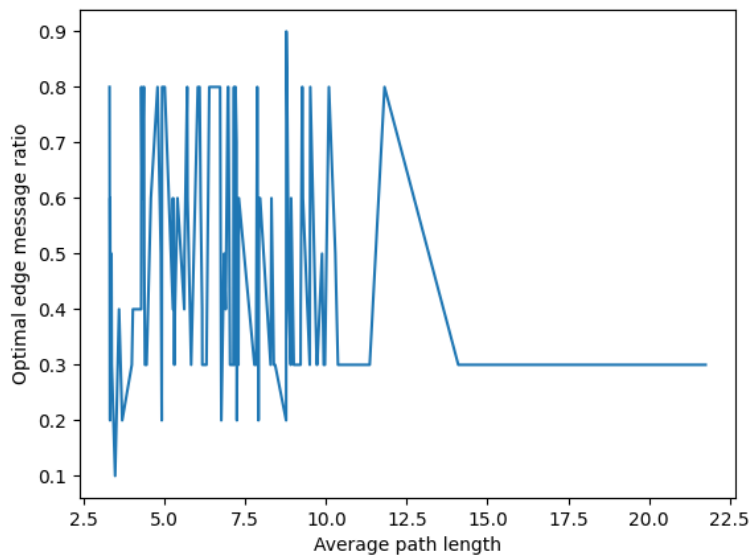


Figure 9: The optimal edge message ratio for different average shortest path lengths on artificial Watts-Strogatz small-world graphs [26] with 256 nodes and average degree 6. Model has $K = 2$ and $o = 32$. It is trained with 80/20 split ratio in the `resample disjoint` mode.