
Tackling Provably Hard Representative Selection via Graph Neural Networks

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 *Representative Selection* (RS) is the problem of finding a small subset of exemplars
2 from a dataset that is representative of the dataset. In this paper, we study RS for
3 unlabeled datasets and focus on finding representatives that optimize the accuracy
4 of a model trained on the selected representatives. Theoretically, we establish a
5 new hardness result for RS by proving that a particular, highly practical variant of it
6 (RS for Learning) is hard to approximate in polynomial time within any reasonable
7 factor, which implies a significant potential gap between the optimum solution of
8 widely-used surrogate functions and the actual accuracy of the model. We then
9 study a setting where additional information in the form of a (homophilous) graph
10 structure is available, or can be constructed, between the data points. We show that
11 with an appropriate modeling approach, the presence of such a structure can turn a
12 hard RS (for learning) problem into one that can be effectively solved. To this end,
13 we develop *RS-GNN*, a representation learning-based **RS** model based on Graph
14 Neural Networks. Empirically, we demonstrate the effectiveness of RS-GNN on
15 problems with predefined graph structures as well as problems with graphs induced
16 from node feature similarities, by showing that RS-GNN achieves significant
17 improvements over established baselines on a suite of eight benchmarks.

18 1 Introduction

19 In the age of massive data, having access to tools that can select exemplar data points representative
20 of an entire dataset is of crucial importance. *Representative selection* (RS) [35], finding a small subset
21 of exemplars from an unlabeled dataset that transmits maximal information for a certain objective,
22 has numerous applications in summarization, active learning, data compression, model training cost
23 reduction, and many other domains (see, e.g., [8, 63, 78, 106, 114, 24, 66, 34]).

24 We first study the computational complexity of a specific but widely-applicable formulation of the
25 RS problem, where we attempt to find a fixed-size subset of representative exemplars from a dataset
26 that can be used to train a model with the best possible *accuracy* on the entire dataset. We show it
27 is impossible to provide a polynomial-time RS algorithm with an approximation factor better than
28 $\omega(n^{-1/\text{poly log log } n})$, unless the Exponential Time Hypothesis (ETH) fails. ETH is a widely-believed
29 assumption in the domain of parameterized complexity which states that the 3-SAT problem cannot be
30 solved in subexponential time in the worst case. Note that $\omega(n^{-1/\text{poly log log } n})$ is almost polynomial,
31 ruling out the existence of any constant approximation or even poly-logarithmic approximation.

32 Our subconstant hardness result is of particular importance because several previous works find
33 representatives by optimizing *surrogate functions*—these can be approximated well in theory—
34 instead of the actual model accuracy. For instance, they consider a submodular surrogate function
35 which can be approximated within a factor $1 - 1/e$ in polynomial time [47, 105, 83, 23, 31]. Our
36 hardness result implies that, in the worst case, there is a significant gap between the optimum solution

37 of such surrogate functions and the actual accuracy of the model, rendering the surrogate functions
38 poor estimators for the quality of the model. To the best of our knowledge, this is the first subconstant
39 hardness result for the RS problem. This motivates us to deviate from directly defining proxy
40 functions, and make use of learning-based approaches that discover the hidden structure of the data
41 to guide the selection. This is in line with the recent attempts to solve computationally hard problems
42 with neural networks, e.g., [107, 27].

43 We therefore study a setting where besides having access to data point features, we also have access to
44 additional information about the data points that can help guide the selection process. Specifically, in
45 this paper we assume the extra information is in the form of a (homophilous) graph structure between
46 the data points. We show empirically that with an appropriate modeling approach, the presence of
47 such a graph structure can turn an originally hard RS problem into one that can be effectively solved.
48 To this end, we develop RS-GNN: a learning-based model for **R**epresentatives **S**election via **G**raph
49 **N**eural Networks. We first demonstrate the effectiveness of RS-GNN for selecting representative
50 nodes from datasets where a natural graph can be accessed, i.e., where edges may be specified by
51 some natural property of the data (e.g., paper citations). Then, we demonstrate that even when
52 a natural graph is not available, creating a similarity graph of the input data points and applying
53 RS-GNN can still select high-quality representatives. We conduct experiments on eight datasets
54 with different sizes and properties, and in both settings where we have and do not have access to
55 a graph structure. Our results show that our model provides significant improvements over three
56 kinds of baselines: 1) well-established baselines that optimize predefined surrogate functions, 2)
57 learning-based methods utilizing graph clustering/pooling and 3) baselines based on active learning.

58 Our main contributions are: 1) Providing a hardness result establishing that, under a standard
59 computational-complexity assumption, RS is hard to approximate in polynomial time within any
60 reasonable factor (this is the *first* subconstant hardness result for RS, to the best of our knowledge),
61 2) Demonstrating empirically that the existence of additional information in the form of a graph
62 structure can make hard RS problems effectively solvable, and 3) Developing RS-GNN for effective
63 RS when one has access to such a graph structure and showing its merit for datasets with natural
64 and/or similarity graphs.

65 2 Related Work

66 We group the existing work that relates to our paper as follows (see Appendix B for more).

67 **Active learning:** In active learning [95, 26] we have an unlabeled set of data points that we can
68 request to label. Since labeling is an expensive task, we usually have a limited budget, say, we can
69 label up to k data points, which are then used to predict the labels of all data points. The goal is to
70 select the set of data points to label in such a way as to maximize the accuracy of the final model. The
71 data points can be iteratively selected in mini-batches (select a mini-batch, label the data points in the
72 batch, update the model, and repeat), or in one-shot [54, 48, 19, 25, 5]. The latter is typically used
73 when model training is time-consuming. RS can be used in the context of one-shot active learning,
74 or for selecting the first mini-batch in the context of mini-batch active learning. In these contexts, a
75 common approach to active learning is to use unsupervised surrogate functions such as KMediod
76 [93] and MaxCover [53] to select a set of data points that maximally cover the dataset with respect to
77 some objective. Moreover, active learning models have been developed for attributed graphs both
78 for mini-batched labeling [17, 41] and for one-shot [109, 118]. We compare against many surrogate
79 functions as well as one-shot graph active learning models in our experiments.

80 **Hardness of Clustering:** The hardness problem studied in this work is distantly related to clustering,
81 which has come in many flavors and shapes: flat vs. hierarchical, partitioning vs. overlapping, graph-
82 based vs. embedding-based vs. time-series-based, supervised vs unsupervised, etc. The interested
83 reader may refer to references (e.g., [50, 62, 39, 111, 58, 110, 32]) for further information. We do
84 emphasize here, though, that the most similar clustering objectives to what we study here are the
85 center-based clustering problems such as k -center and k -means. In these settings, the hardness results
86 and known algorithmic guarantees (i.e., lower and upper bounds) are not far from each other. For
87 example, while non-metric k -center cannot be approximated to within any constant, the metric special
88 case (generalizing the ubiquitous Euclidean setting) admits a 2-approximation [45] and cannot be
89 approximated to within better than a factor 2 [102]. On the other hand, the k -means objective admits
90 a constant-factor approximation [4, 61] and the best hardness results are 1.0013 [7, 71]. In contrast

91 to all these results, we present a superconstant hardness for the RS problem, stressing the big gap
92 between any optimizable surrogate function and the true objective.

93 **Graph clustering (community detection):** Early approaches only considered the graph structure
94 and disregarded the node features. These approaches typically learn an embedding for each node
95 (e.g., the spectral features, random walk embeddings, or auto-encoder based embeddings) and then
96 feed these node embeddings into a clustering algorithm such as kmeans (see, e.g., [18, 84, 46, 112]).
97 Recently, approaches based on GNNs, which take both graph structure and node features into account,
98 have gained more popularity and success [119, 91, 76, 104, 12, 60, 100]. RS and clustering are two
99 highly related tasks: many models developed for RS can (in theory) be used for clustering and vice
100 versa. However, due to the distinct properties of the two tasks, a model that works well for one task
101 may not necessarily work well for the other. While our main focus is on RS, we also experiment with
102 graph clustering and compare against several existing approaches.

103 3 Notation & Problem Definition

104 We use bold lowercase letters to denote vectors and bold uppercase letters to denote matrices. Let
105 x_i represent the i^{th} element of x and M_i represent the i^{th} row of M . For a function $f : A \mapsto B$
106 and a subset $A' \subseteq A$, we use $f|_{A'}$ to denote the restriction of the domain of f to A' . For a dataset,
107 we use $\mathcal{V} = \{v_1, \dots, v_m\}$ to represent the set of data points (of size m) and \mathbf{X} to represent the data
108 matrix, such that \mathbf{X}_i corresponds to the features of v_i . When data points have class labels, we use c
109 to represent the number of classes. First, we define a general framework for Representative Selection
110 (RS) as follows:

111 **Definition 3.1 (RS).** *Given a set of data points \mathcal{V} , their features \mathbf{X} , a number $0 < k \leq |\mathcal{V}|$, and a*
112 *utility function $u : 2^{\mathcal{V}} \mapsto \mathbb{R}$, the representative selection problem is to select a subset $\mathcal{S} \subseteq \mathcal{V}$ of k*
113 *representatives that maximize the utility $u(\mathcal{S})$.*¹

114 The applicability and tractability of an RS problem depends on the utility function, u . Intuitively, u
115 should capture the usefulness of the subset \mathcal{S} as a representative of \mathcal{V} ; more precisely, if there is a
116 particular application of the full dataset \mathcal{V} , u quantifies the degree to which \mathcal{S} can be used instead.
117 This can vary with the particular application considered; in this paper we are mostly concerned with
118 a particular utility model that associates representativeness with *learnability*. In the *Representative*
119 *Selection for Learning (RSL)* problem, we get to see the labels of the selected representatives, based
120 on which we aim to train a classifier with highest predictive performance on the entire dataset.

121 **Definition 3.2 (RSL).** *Let \mathcal{V} be a set of data points with observed features \mathbf{X} and with labels Y*
122 *generated from an oracle function $\phi^* : \mathcal{V} \mapsto \{1, \dots, c\}$, and Φ be a class of predictor functions (not*
123 *necessarily containing ϕ^*). Given \mathcal{V} , \mathbf{X} , Φ , y , and a number $0 < k \leq |\mathcal{V}|$, the goal of RSL is to*
124 *select a subset \mathcal{S} of k representatives from \mathcal{V} such that training a classifier $\phi \in \Phi$ based on \mathbf{X} and*
125 *$Y|_{\mathcal{S}}$ maximizes the normalized accuracy of ϕ on the entire set \mathcal{V} .*

126 The model class Φ defines a suitable inductive hypothesis for the learning problem so that RS is
127 well-defined. The predictor ϕ is chosen to maximize the normalized accuracy, defined as $\overline{\text{Acc}} =$
128 $c(\text{accuracy} - 1/c)$, on the entire dataset; in other words it is a *transductive learning problem* [38].
129 RSL is a natural problem that has multiple real-world applications in dataset selection, active learning,
130 efficient ML and other areas (see section 7 for examples). Further, it can be expected to correspond to
131 a general notion of the representativeness of a subset, even outside a learning use case. In the rest of
132 this paper we will focus on this version of the RS problem, and defer generalizations to other utility
133 models to future work.

134 4 Theoretical Findings: Hardness Results for RSL

135 We describe a theoretical finding that explains why the common practice of hand designing and
136 optimizing surrogate functions may not be a good approach for RSL. In Definition 3.2, let $u(\mathcal{S}) =$
137 $\overline{\text{Acc}}(\phi(\mathcal{S}))$ represent the normalized accuracy of the classifier ϕ when trained on a subset \mathcal{S} of data
138 points. Let $\mathcal{S}^* = \arg \max_{\mathcal{S}} u(\mathcal{S})$ be the optimal set of representatives. Ideally, one would optimize
139 $u(\mathcal{S})$ and find \mathcal{S}^* . This may, however, be impossible without a-priori having access to the labels for

¹Note that \mathcal{V} is implicit in the definition of u , hence in the definition of RS.

140 all data points. Alternatively, most existing works on RSL (and RS in general) focus on defining an
 141 intuitive surrogate function Ω and find a solution \mathcal{S}^Ω by optimizing $\Omega(\mathcal{S})$ in the hope that $u(\mathcal{S}^\Omega)$ is a
 142 good approximator for $u(\mathcal{S}^*)$. In this section, we establish an inapproximability result demonstrating
 143 that $u(\mathcal{S}^\Omega)$ may not be a good estimator of $u(\mathcal{S}^*)$ for any polynomial-time-computable surrogate
 144 function Ω . We do this by showing that there are naturally defined learning tasks for which there
 145 exists a significant gap between $u(\mathcal{S}^\Omega)$ and $u(\mathcal{S}^*)$.

146 We start by studying the computational hardness of RSL on an end-to-end binary classification
 147 task, from which the aforementioned claims follow. We say an RS algorithm \mathcal{A} (e.g., optimizing
 148 a surrogate function) approximates the optimal solution with an approximation factor α if we can
 149 establish that $u(\mathcal{S}^{\mathcal{A}})$ is within a multiplicative factor α of the optimal solution $u(\mathcal{S}^*)$ for any learning
 150 problem, where $\mathcal{S}^{\mathcal{A}}$ is the output of the RS algorithm.

151 Under the Exponential-Time Hypothesis (ETH) assumption², we show that there is no polynomial-
 152 time RS algorithm with an approximation factor better than $\omega(n^{-1/\text{poly log log } n})$. In other words,
 153 we show that there is an instance of RSL for which the gap between $u(\mathcal{S}^{\mathcal{A}})$ and $u(\mathcal{S}^*)$ for any
 154 polynomial-time RSL algorithm \mathcal{A} is at least $\omega(n^{-1/\text{poly log log } n})$, unless ETH fails. A similar
 155 approximation gap exists between the best polynomial-time and best exponential-time algorithm. Note
 156 that $\omega(n^{-1/\text{poly log log } n})$ is almost polynomial, ruling out the existence of any constant approximation
 157 or even poly-logarithmic approximation. Also note that the best solution may not give 100% accuracy,
 158 nor does it necessarily match the accuracy obtained by using all labels (since we only use k).

159 **Definition 4.1** (Fit-or-Not (FoN) Learning Problem). *We have m data points and n binary features.*
 160 *Each feature is associated with one of two types: red or blue. The types are generated independently*
 161 *and uniformly at random, they are consistent across data points, and are hidden from the algorithm.*
 162 *Each data point has value 1 for two features and 0 for the rest. The label of a data point is 1 if the*
 163 *type of its features of value 1 are the same, and the label is 0 otherwise. The goal is to maximize the*
 164 *normalized accuracy for all the data points given labels only on selected data points.*

165 Figure 1 visually presents an instance of the FoN problem. Each
 166 row represents a data point and each column represents a feature.
 167 There are $m = 10$ data points, each having $n = 5$ binary
 168 features. Each data point has a value of 1 exactly for two of
 169 the features, and a value of 0 for the rest. The types of the first
 170 and the third features are *blue* and the types of the other three
 171 features are *red*. These types are hidden from the algorithm.
 172 For the first data point, the two values of 1 are for the first and
 173 second features. Since these features have different types, the
 174 label for this data point is 0. For the second data point, however,
 175 the two values of 1 are for the first and the third features that
 176 have the same type, therefore the label for this data point is 1.
 177 The labels of the other data points are determined similarly.

178 The goal of the algorithm is to select a subset with $k < m$ data
 179 points in such a way that a model trained on the labels of those k
 180 data points makes accurate predictions for all the m data points
 181 (i.e. it generalizes well to the other $(m - k)$ data points).

182 FoN can be naturally framed as RSL by defining the components
 183 from Definition 3.2 as follows: let \mathcal{V} be the set of m data points,
 184 $\mathbf{X} \in \mathbb{R}^{m \times n}$ be the features matrix, Φ be the class of models that correspond to the data generation
 185 process in Definition 4.1 (i.e. constructed from particular partitions of features into red/blue etc.),
 186 k be the budget for selecting representatives, and let the labels ϕ^* be as defined above; the latent
 187 information consists of the types of the n features. The simplicity of FoN shows that RS is hard in a
 188 very broad form. The next theorem is the main result of this section.

189 **Theorem 4.2.** *There is no polynomial-time RSL algorithm for FoN with an approximation factor*
 190 *better than $\omega(n^{-1/\text{poly log log } n})$, unless the exponential-time hypothesis fails.*

		n=5 binary features					Labels
m=10 data points		1	1	0	0	0	0
		1	0	1	0	0	1
		0	0	0	1	1	1
		0	0	1	1	0	0
		0	1	0	1	0	1
		1	0	0	0	1	0
		1	0	0	1	0	0
		0	1	0	0	1	1
		0	0	1	0	1	0
		0	1	1	0	0	0

Figure 1: An instance of the FoN problem with $m = 10$ data points and $n = 5$ binary features.

²A widely-believed assumption (e.g., see [2, 15, 16, 22, 21, 28, 56, 59, 79, 80]) in the domain of parameterized complexity which states that 3-SAT cannot be solved in subexponential time in the worst case.

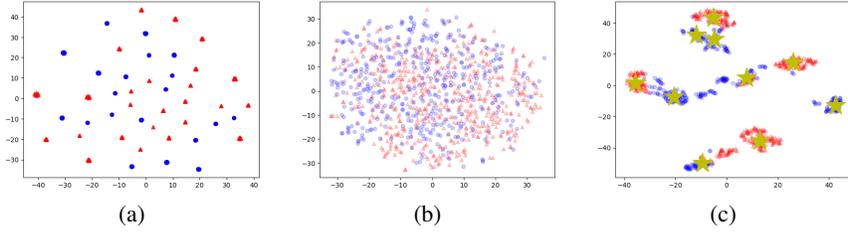


Figure 2: A t-SNE visualization of an instance of the FoN problem (the colors show the classes) when using (a) the original features, (b) adding a graph context and applying graph convolution, and (c) adding a graph context and applying a variant of RS-GNN to embed the nodes and select representatives (the yellow stars represent the selected representatives).

191 To the best of our knowledge, this is the first subconstant hardness result for any RS problem.
 192 This is of particular importance because several previous works have chosen to optimize *surrogate*
 193 *functions*—these can be approximated well in theory—instead of the actual model accuracy. For
 194 instance, previous work considers a submodular surrogate function which can be approximated within
 195 a factor $1 - 1/e$ in polynomial time [47, 105, 83, 23, 31]. Our hardness result implies the following.

196 **Corollary 4.3.** *In the worst case, there is a significant gap of $\omega(n^{-1/\text{poly log log } n})$ between $u(\mathcal{S}^*)$*
 197 *and the solution of any polynomial-time approximable surrogate function that estimates $u(\mathcal{S}^*)$.*

198 The corollary follows because such surrogate functions can be optimized or approximated in polyno-
 199 mial time, but Theorem 4.2 shows that even for simple learning problems, approximating the accuracy
 200 is not possible in polynomial time (assuming the ETH); therefore, there are certain instances of the
 201 problem where optimizing for the surrogate functions does not optimize for the learning accuracy
 202 within the given approximation factor.

203 It is worth noting that many ML tasks are (known to be) hard to optimize, hence surrogate loss
 204 functions are commonly used in the context of deep learning: For example, finding a linear classifier
 205 with minimum 0-1 loss is NP-complete [81], and convex relaxations of this loss are typically used as
 206 surrogates in practice. What we show in this section is that fairly simple and natural instances of the
 207 representation selection problem are not only NP-complete but also hard to approximate. Thus the
 208 optimal solution of any surrogate function will be far from the actual optimum. This gap is prominent
 209 especially when we disentangle the task of finding the set to label from the task of training a model.
 210 We will see in the next sections that combining the two remedies the problem and produces very good
 211 results, which is in line with the applicability of surrogate loss functions within deep learning models.

212 5 RSL in Presence of a Graph Structure

213 In the previous section, we established the hardness of RSL by finding a natural problem called FoN
 214 for which RSL is hard to approximate in polynomial time within any reasonable factor. The hardness
 215 of RSL motivates seeking additional information about the problem that may help better guide the
 216 selection process. We next study whether the presence of additional information can help tackle the
 217 hard RSL problem effectively. In particular, we study the case where the additional information is
 218 provided in terms of graph structure of the data points with some degree of homophily; that is, nodes
 219 belonging to the same class are more likely to be connected to each other than nodes belonging to
 220 different classes. We present an algorithm for doing RSL in this setting and provide an empirical
 221 study and leave a theoretical analysis of this setting as well as settings with other types of additional
 222 information available as future work.

223 Let us start by providing a visual understanding of the FoN problem. We construct a version of the
 224 FoN problem with $m = 1000$ data points $n = 10$ features, where we assume the first 5 features
 225 have type red and the next 5 features have type blue. For each data point, we select two of its
 226 features uniformly at random and set their values to a number from $[0.9, 1.0]$ (other values are 0).
 227 In Figure 2(a), we present a visualization of the data points using t-SNE [101], where the colors
 228 represent the classes. From the visualization, we observe that there are 45 dense blocks of nodes
 229 (each having the same label) corresponding to $\binom{10}{2}$ different ways of selecting the position for the

230 non-zero elements. The blocks are scattered in such a way that it is difficult for an RSL algorithm to
 231 select a small subset such that a model trained on the labels of the nodes in that subset generalizes
 232 well to the other data points.

233 We now consider a variant of the FoN problem where an additional homophilous graph structure is
 234 available among the data points.

235 **Definition 5.1 (GFoN).** *GFoN is a variant FoN where for each pair of data points we add an edge*
 236 *between them with probability p if they belong to the same class and p' if they belong to different*
 237 *classes. Letting $p > p'$ ensures the graph has some degree of homophily.*

238 A popular modeling choice to use when graph structure is available are graph convolution operations,
 239 where the features for each node are replaced by a weighted average of the features of their neighbors,
 240 with weights proportional to the degrees. If we apply a graph convolution operation on GFoN with
 241 $p = 0.05$ and $p' = 0.01$, we get the updated node features in Figure 2(b). One can observe that
 242 the nodes from the two classes have a high overlap and so any RSL algorithm may still fail. To
 243 understand why this happens, consider four data points whose non-zero values are in positions 1 and
 244 2, 9 and 10, 1 and 9, and 2 and 10 respectively. The first two nodes will have a label of 1 and the
 245 other two nodes will have a label of 0. However, if we disregard the node degrees, aggregating the
 246 first two nodes will give the same representation as aggregating the second two nodes.

247 According to Figure 2(b), in presence of a graph structure, a naive application of graph convolution
 248 may not necessarily work best, and motivates the development of better modeling techniques. We
 249 show in our experiments that this also holds for many existing graph clustering/pooling approaches.

250 To better leverage the graph, in the next section we develop an approach that works by learning a
 251 mapping of the data points to a latent space where 1- we can group the nodes and select representatives
 252 that cover groups of nodes, and 2- we can better distinguish nodes belonging to each class. Figure 2(c)
 253 shows a t-SNE visualization of the data points in the latent space and the selected representatives.
 254 One can observe that with appropriate modeling, the presence of a homophilous graph structure turns
 255 the provably hard FoN problem into an RSL problem that can be effectively solved.

256 We next describe our approach for taking advantage of the graph context in improving RSL. In our
 257 experiments, we show that even when such a graph is not available, a similarity graph of the input
 258 node features may still be quite effective. Our theoretical and empirical results motivate obtaining or
 259 constructing graph contexts for RSL problems when possible.

260 6 RS-GNN: A Representation Learning-based Model for RSL

261 In this section we develop **RS-GNN**, a representation learning-based approach for **RS** via **GNNs**.
 262 We start with defining notation and provide necessary background. We denote an attributed graph
 263 as $\mathcal{G} = \{\mathcal{V}, \mathbf{A}, \mathbf{X}\}$ where $\mathcal{V} = \{v_1, \dots, v_m\}$ represents the set of nodes, $\mathbf{A} \in \mathbb{R}^{m \times m}$ represents the
 264 adjacency matrix, and $\mathbf{X} \in \mathbb{R}^{m \times n}$ represents the matrix of node features (m nodes and n features).

265 **Graph Neural Networks (GNNs)** encode graph-structured data in continuous space [20]. Graph
 266 convolutional networks (GCNs) [67] are a powerful variant of GNNs. Let $\mathcal{G} = \{\mathcal{V}, \mathbf{A}, \mathbf{X}\}$ be
 267 an attributed graph, $\hat{\mathbf{A}} = \mathbf{A} + \mathbf{I}$ be the adjacency matrix of \mathcal{G} with self-loops included, and \mathbf{D}
 268 be the degree matrix of $\hat{\mathbf{A}}$, where $D_{ii} = \sum_j \hat{A}_{ij}$ and $D_{ij} = 0$ for $i \neq j$. The l^{th} layer of
 269 an L-layer GCN model with parameters $\Theta = \{\mathbf{W}^{(1)}, \dots, \mathbf{W}^{(L)}\}$ can be defined as: $\mathbf{H}^{(l)} =$
 270 $\sigma(\mathbf{D}^{-\frac{1}{2}} \hat{\mathbf{A}} \mathbf{D}^{-\frac{1}{2}} \mathbf{H}^{(l-1)} \mathbf{W}^{(l)})$, where $\mathbf{H}^{(l)}$ represents node embeddings in the l^{th} layer ($\mathbf{H}^{(0)} = \mathbf{X}$),
 271 $\mathbf{W}^{(l)}$ is a weight matrix, and σ is an activation function. In the rest of the paper, we use $\text{GCN}(\mathcal{G}; \Theta)$
 272 to show the application of a GCN function with parameters Θ on a graph \mathcal{G} .

273 **Deep Graph Infomax (DGI)** [103] is an approach for unsupervised representation-learning in
 274 attributed graphs. Given an attributed graph \mathcal{G} , in each iteration DGI creates a corrupted graph \mathcal{G}'
 275 from \mathcal{G} . Then, it computes node embeddings \mathbf{H} and \mathbf{H}' for the two graphs by applying a GNN
 276 model on them, and a summary vector \mathbf{s} based on the node embeddings \mathbf{H} . Finally, a discriminator
 277 is simultaneously trained to separate the node embeddings of the original graph (i.e., \mathbf{H}) from those
 278 of the corrupted graph (i.e., \mathbf{H}') based on the summary vector \mathbf{s} . We describe the details of each step
 279 later when we define our final model.

280 **6.1 Representative Selection via GNNs**

281 An attributed graph $\mathcal{G} = \{\mathcal{V}, \mathbf{A}, \mathbf{X}\}$ has two modalities, the node features \mathbf{X} and the graph structure
 282 \mathbf{A} . To be able to exploit both modalities in selecting good representatives, we employ a function
 283 EMB that combines the two modalities into a single embedding matrix of the graph. Concretely,
 284 $\text{EMB}(\mathcal{G}) = \mathbf{H}$, where $\mathbf{H} \in \mathbb{R}^{m \times d}$ and d represents the embedding dimension. We also employ a
 285 differentiable function SEL that receives \mathbf{H} as input and selects k nodes as representatives. That is,
 286 $\text{SEL}(\mathbf{H}) = \mathcal{S}$. The two functions are optimized in a multi-task setting with the loss function $\mathcal{L} =$
 287 $\mathcal{L}_{\text{EMB}} + \lambda \mathcal{L}_{\text{SEL}}$, where \mathcal{L}_{EMB} encourages learning informative node embeddings and \mathcal{L}_{SEL} encourages
 288 selecting good representatives. One can create different RSL models with different choices of EMB,
 289 SEL, \mathcal{L}_{EMB} and \mathcal{L}_{SEL} . GNNs have prove effective in learning node embeddings, so we use GNNs
 290 as our embedding function EMB. For \mathcal{L}_{EMB} , we use the DGI objective which has shown to provide
 291 high-quality embeddings in unsupervised settings. For SEL, we consider a representative embedding
 292 matrix $\mathbf{R} \in \mathbb{R}^{k \times d}$ with learnable parameters where \mathbf{R} is initialized randomly and \mathbf{R}_j represents the
 293 embedding for the j^{th} representative. We let: $\mathcal{L}_{\text{SEL}} = \sum_i \min_j (\text{Dist}(\mathbf{H}_i, \mathbf{R}_j))$, where $\text{Dist}(\mathbf{H}_i, \mathbf{R}_j)$
 294 is the distance between the i^{th} node’s embedding \mathbf{H}_i and the j^{th} representative’s embedding \mathbf{R}_j . We
 295 use Euclidean distance as the distance function. We select the representative corresponding to each
 296 \mathbf{R}_j by finding the closest node embedding from \mathbf{H} to \mathbf{R}_j , i.e., $\text{argmin}_i (\text{Dist}(\mathbf{H}_i, \mathbf{R}_j))$.

297 With the above loss function, the model
 298 can trivially reduce \mathcal{L}_{SEL} by making the
 299 values in \mathbf{H} arbitrarily small. That is be-
 300 cause multiplying \mathbf{H} by a small constant
 301 may not change \mathcal{L}_{EMB} substantially, but it
 302 can make the distances between the nodes
 303 arbitrarily small, resulting in a low value
 304 for \mathcal{L}_{SEL} even for a random representative
 305 embedding matrix \mathbf{R} . We next describe a
 306 normalization scheme, *CenterNorm*, that
 307 is applied to \mathbf{H} before \mathbf{H} is used, which
 308 helps avoid this problem.

309 **CenterNorm:** As we show in Ap-
 310 pendix C, using the DGI loss function re-
 311 sults in corrupted node embeddings that
 312 form a dense cluster in some part of the
 313 embedding space, and node embeddings
 314 (from the actual graph) that arrange them-
 315 selves in subclusters around (and outside)
 316 this dense cluster of negative examples.
 317 Based on the this observation, we propose
 318 an ℓ_2 normalization of the node embeddings in \mathbf{H} with respect to the center of the node embeddings:

$$\boldsymbol{\mu} = \frac{1}{n} \sum_i \mathbf{H}_i, \boldsymbol{\zeta} = \|\mathbf{H} - \boldsymbol{\mu}\|, \tilde{\mathbf{H}} = (\mathbf{H} - \boldsymbol{\mu})/\boldsymbol{\zeta}, \quad (1)$$

319 where $\boldsymbol{\mu}$ is the center of the embeddings \mathbf{H} , $\boldsymbol{\zeta}$ is the ℓ_2 norms of the nodes with respect to the center,
 320 and $\tilde{\mathbf{H}}$ represents the normalized embeddings. With CenterNorm, the model can no longer decrease
 321 \mathcal{L}_{SEL} simply by making the values \mathbf{H} smaller. Note that since the embedding clusters in \mathbf{H} are at
 322 a large angle from each other (cf. supplementary material), ℓ_2 normalization has a low chance of
 323 collapsing two clusters. Furthermore, ℓ_2 normalization helps bring the nodes within one cluster closer
 324 to each other, which helps in identifying clusters and selecting representatives.

325 **The Final RS-GNN Model:** The full RS-GNN model is described in Algorithm 1. The input is an
 326 attributed graph $\mathcal{G} = (\mathcal{V}, \mathbf{A}, \mathbf{X})$ and a number k corresponding to the number of representative nodes
 327 that must be selected from the graph. The model initializes \mathbf{R} (the representative embeddings), $\boldsymbol{\Theta}$
 328 (the GCN parameters), and \mathbf{U} (the parameters for the DGI discriminator). Lines 3 to 7 compute node
 329 embeddings \mathbf{H} using a GCN³ model and compute a DGI loss as \mathcal{L}_{EMB} . Here, $\text{bilinear}(\mathbf{H}, \mathbf{s}; \mathbf{U}) =$

³While we use GCN for direct comparability to existing work, note that one can use any other GNN model for RS-GNN. For example, for the FoN problem in Figure 2(c), we used a variant of the GraphSage [49] model as it better matched the problem.

Algorithm 1 The training procedure of RS-GNN.

Input: $\mathcal{G} = (\mathcal{V}, \mathbf{A}, \mathbf{X})$, k

- 1: Initialize \mathbf{R} , $\boldsymbol{\Theta}$, and \mathbf{U}
- 2: **for** epoch=1 **to** #epochs **do**
- 3: $\mathcal{G}' = (\mathcal{V}, \mathbf{A}, \text{shuffle}(\mathbf{X}))$
- 4: $\mathbf{H}' = \text{GCN}(\mathcal{G}; \boldsymbol{\Theta})$, $\mathbf{H}' = \text{GCN}(\mathcal{G}'; \boldsymbol{\Theta})$
- 5: $\mathbf{s} = \text{sigmoid}(\frac{1}{n} \sum_i \mathbf{H}_i)$
- 6: $\mathbf{p} = \text{bilinear}(\mathbf{H}, \mathbf{s}; \mathbf{U})$, $\mathbf{p}' = \text{bilinear}(\mathbf{H}', \mathbf{s}; \mathbf{U})$
- 7: $\mathcal{L}_{\text{EMB}} = -\sum_i (\log(\mathbf{p}_i) + \log(1 - \mathbf{p}'_i))$
- 8: $\boldsymbol{\mu} = \frac{1}{n} \sum_i \mathbf{H}_i$, $\boldsymbol{\zeta} = \|\mathbf{H} - \boldsymbol{\mu}\|$
- 9: $\tilde{\mathbf{H}} = \text{CenterNorm}(\mathbf{H}) = (\mathbf{H} - \boldsymbol{\mu})/\boldsymbol{\zeta}$
- 10: $\mathcal{L}_{\text{SEL}} = \sum_i \min_j \text{Dist}(\tilde{\mathbf{H}}_i, \mathbf{R}_j)$
- 11: $\mathcal{L} = \mathcal{L}_{\text{EMB}} + \lambda \mathcal{L}_{\text{SEL}}$
- 12: Compute gradients for \mathcal{L} , upd. params.
- 13: Let $\hat{\mathbf{R}}$ and $\hat{\mathbf{H}}$ be the representative and normalized node embeddings with minimum \mathcal{L} during training.
- 14: **for** j=1 **to** k **do**
- 15: The j^{th} representative = $\text{argmin}_i \text{Dist}(\hat{\mathbf{H}}_i, \hat{\mathbf{R}}_j)$

330 $\text{sigmoid}(\mathbf{H}^T \mathbf{U} \mathbf{s})$ where \mathbf{H}^T indicates the transpose of \mathbf{H} . Lines 8 and 9 apply CenterNorm. Lines
331 10 to 12 compute \mathcal{L}_{SEL} and then \mathcal{L} and update the parameters accordingly.

332 We keep track of the epoch with minimum joint loss \mathcal{L} . Let $\hat{\mathbf{R}}$ and $\hat{\mathbf{H}}$ be the corresponding
333 representative and normalized node embeddings in the best epoch. We select the j^{th} representative to
334 be the node whose normalized embedding is closest to the j^{th} representative embedding. Lines 13 to
335 15 select the representatives based on $\hat{\mathbf{R}}$ and $\hat{\mathbf{H}}$.

336 7 Empirical Results

337 We describe our baselines, datasets, and metrics, and defer implementation details to supplementary
338 material⁴.

339 **Baselines:** We compare against a representative set of baselines from different categories, as
340 described below. The details for each baseline can be found in Appendix E. *Random:* Selects k
341 nodes uniformly at random. *Popular:* Selects the k nodes with the maximum degree from the
342 graph. *Surrogate functions on node features:* We test a representative set of surrogate functions:
343 KMedoid, KMeans, Farthest First Search (FFS) [114] and (Greedy) MaxCover [53]. For KMeans, we
344 select the closest node to each cluster center as a representative. For FFS and MaxCover, we select
345 representatives sequentially. In FFS, the next representative is the node farthest away (by Euclidean
346 distance) from the closest representative in the current set. In MaxCover, the next representative
347 is the node that increases the coverage of the non-selected nodes the most. For MaxCover, we
348 experiment with RBF kernel and cosine similarities; we use MC-RBF and MC-Cos to refer to the
349 two versions respectively. Note that the sequential nature of FFS and MC makes them less amenable
350 to parallelization. *Surrogate functions on node embeddings:* We use similar functions as above but
351 apply them on DGI node embeddings as opposed to on the initial node features. Note that when we
352 run these baselines using DGI embeddings as context, their selections are informed by both node
353 features and the graph structure. *Graph clustering/pooling/active learning:* We compare against a
354 number of graph clustering/pooling/active learning approaches from different categories. Specifically,
355 we compare against MinCut [11] which is a well-established pooling approach, FeatProp [109] which
356 is successful graph active learning approaches, SDCN [12] which is a well-established auto-encoder
357 based graph clustering model, EGAE [36] and GCC [36] which are recent joint representation
358 learning and clustering approaches, and DMoN [100] which is a state-of-the-art graph clustering
359 approach based on modularity maximization.

360 **Datasets:** We use eight established benchmarks in the GNN literature: three citation networks namely
361 Cora, CiteSeer, and Pubmed [94, 55], a citation network named OGBN-Arxiv [55] which is orders of
362 magnitude larger than the previous three, two datasets from Amazon products (Photos and PC) [97],
363 and two datasets from Microsoft Academic (CS and physics) [97]. Supplementary material offers a
364 more detailed description of datasets and their statistics. Our datasets have a wide range in terms of
365 the number of nodes (from 2K to 170K), edges (from 4.5K to 1.1M), features (from 100 to 8.5K) and
366 classes (from 3 to 40).

367 **Measures:** We measure the quality of the selected representatives $\mathcal{S} \subseteq \mathcal{V}$ using the following
368 transductive semi-supervised node-classification problem. We train a GCN model on the dataset
369 where the parameters of the GCN are learned only based on the labels of the nodes in \mathcal{S} . Note that
370 this GCN is completely independent of the internal GCN model used in RS-GNN. We randomly split
371 the remaining nodes into validation and test sets. The validation set is used for early stopping. The
372 classification accuracy on the test set is used as the metric for measuring the quality of the selected
373 representatives. Considering a validation set for early stopping reduces the chances of overfitting for
374 the classifier and makes the reported test accuracy mainly a function of the quality of the selected
375 representatives.

376 **RSL in the Presence of a Graph Structure:** The results are presented in Table 1. For the results
377 in this table, we set k for each dataset to be $2c$, where c represents the number of classes. We
378 found this to be a small enough number for a meaningful comparison of the quality of the selected
379 representatives⁵, and high enough for the classification GCN model to learn appropriate functions of

⁴The code is available at: https://github.com/google-research/google-research/rs_gnn

⁵Note that if $k \approx n$, all models may perform equally well.

Table 1: For each dataset, each algorithm selects $2c$ representatives. Then, we train a GCN model on the labels of the selected representatives. The reported metric is the test accuracy of the GCN models. Bold numbers indicate statistically significant winner(s) following a t-test (p-value=0.05). The color-codes and symbols represent \clubsuit surrogate functions on features, \spadesuit surrogate functions on embeddings, and \heartsuit attributed graph clustering/pooling or active learning approaches.

Selector	Cora	CiteSeer	Pubmed	Photos	PC	CS	Physics	Arxiv	Avg.
Random	49.1±6.9	33.1±8.3	52.0±8.1	70.2±6.4	65.4±6.1	72.9±5.0	73.6±6.8	49.0±1.4	58.2
Popular	59.2±1.3	35.5±0.8	63.3±0.3	34.9±0.5	49.9±2.1	73.1±1.1	50.5±0.0	31.3±0.8	49.7
\clubsuit KMedoid	53.7±5.4	40.3±2.8	53.2±1.0	65.8±1.2	66.9±1.8	51.8±1.0	66.6±1.6	43.9±1.5	55.3
\clubsuit KMeans	32.5±5.8	35.4±1.2	50.6±0.5	72.5±0.9	70.9±1.0	66.5±0.5	73.0±1.4	48.4±0.5	56.2
\clubsuit FFS	48.5±8.0	39.3±6.8	43.1±4.9	80.0±5.0	71.5±3.4	54.6±2.2	76.6±2.8	45.2±0.8	57.3
\clubsuit MC-RBF	45.5±2.7	25.8±3.7	53.0±0.2	78.4±1.2	65.5±0.9	66.4±1.2	58.1±0.3	51.4±0.6	55.5
\clubsuit MC-Cos	49.7±9.5	50.2±3.1	66.6±0.5	77.2±1.0	74.0±3.7	87.3±1.4	81.3±3.4	47.6±1.6	70.0
\spadesuit KMedoid	48.4±4.4	34.1±1.9	60.9±5.3	81.5±2.6	69.8±3.4	82.5±2.9	81.2±7.0	OOM	—
\spadesuit KMeans	62.6±9.3	42.7±6.3	60.5±6.3	83.6±2.9	74.8±2.7	86.9±1.8	90.6±2.4	51.2±1.0	69.1
\spadesuit FFS	62.6±4.5	50.4±5.7	46.7±7.2	73.4±5.7	63.5±6.4	84.8±5.3	83.4±5.0	48.6±2.2	64.2
\spadesuit MC-RBF	66.3±2.6	35.3±4.5	54.9±5.3	37.4±3.7	50.7±2.5	65.2±1.2	59.8±5.1	41.2±1.6	51.4
\spadesuit MC-Cos	67.3±5.2	49.0±4.1	67.3±1.1	84.4±1.0	74.0±3.7	87.3±1.4	81.3±3.4	47.6±1.6	70.0
\heartsuit MinCUT	51.9±7.5	37.3±8.0	59.5±6.1	14.4±7.5	18.3±8.7	85.5±1.4	86.0±3.3	32.4±5.2	48.2
\heartsuit FeatProp	56.6±1.7	37.8±1.6	65.2±0.6	78.2±1.8	68.5±1.1	74.7±0.3	81.4±0.6	47.8±1.0	63.8
\heartsuit SDCN	41.6±9.5	33.8±9.3	47.8±8.5	61.0±10.8	54.2±8.7	66.4±6.7	77.5±9.2	38.4±4.4	52.6
\heartsuit EGAE	64.4±3.8	45.0±5.8	57.7±5.1	83.5±3.0	75.4±3.2	79.9±3.2	80.4±4.0	50.6±1.0	67.1
\heartsuit GCC	68.7±2.2	49.3±6.0	63.9±5.1	84.2±1.4	72.1±2.1	85.8±1.5	89.6±0.9	52.1±1.4	70.7
\heartsuit DMoN	58.0±7.1	40.5±7.4	55.3±7.5	78.6±9.2	70.3±3.3	84.3±1.4	85.9±3.8	52.5±1.9	65.7
RS-GNN	72.4±3.7	54.7±3.9	65.8±3.0	86.3±1.4	74.3±1.7	89.3±0.8	90.0±2.6	52.6±1.2	73.2

380 the data. Since c is different for each dataset, making k a function of c also provides the opportunity
381 to compare performance not only in terms of variation in the datasets, but also in terms of variation in
382 the number of selected representatives.

383 RS-GNN performs well across all datasets and consistently outperforms (or matches) the baselines
384 (with the exception of Pubmed). It has a low variance across different runs making it a reliable model.
385 Among the surrogate functions, MC-Cos and KMeans perform best. We found FFS to be sensitive
386 to outliers. We also found it difficult to select a set of hyperparameters for MC-RBF that work well
387 across datasets. Among graph clustering/pooling and active learning approaches, we found GCC to
388 perform best and be the only model that (overall) outperforms surrogate functions when applied on
389 DGI embeddings. Many of the other graph clustering/pooling or active learning approaches even fall
390 short of the MC-Cos model when applied on the node features alone, thus showing the importance
391 of developing appropriate models for taking advantage of the graph structure and confirming our
392 finding in Figure 2. MinCut produced degenerate solutions for Photos and PC datasets in many runs
393 (assigning all nodes to one cluster), hence performing poorly on them.

394 More results, analysis, and visualizations is provided in Appendix A.

395 8 Conclusion

396 In this paper, we studied the representative selection (RSL) problem theoretically and empirically.
397 We proved new hardness results showing it is impossible to provide a polynomial-time algorithm for
398 RSL with an approximation within any reasonable factor, unless the exponential time hypothesis fails.
399 The hardness result explains the significant gap between the accuracy of models trained on optimal
400 representatives and the widely-used surrogate functions for RS problems, and, in turn, justifies new
401 techniques to solve this problem. In light of this result we proposed RS-GNN to optimize the RSL
402 task via graph neural networks, and showed its effectiveness on a suite of different datasets and tasks.

403 References

- 404 [1] Zeinab Abbassi, Vahab S. Mirrokni, and Mayur Thakur. Diversity maximization under matroid
405 constraints. In Inderjit S. Dhillon, Yehuda Koren, Rayid Ghani, Ted E. Senator, Paul Bradley,
406 Rajesh Parekh, Jingrui He, Robert L. Grossman, and Ramasamy Uthurusamy, editors, *The*

- 407 *19th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining,*
408 *KDD 2013, Chicago, IL, USA, August 11-14, 2013*, pages 32–40. ACM, 2013.
- 409 [2] Divesh Aggarwal and Noah Stephens-Davidowitz. (gap/s) eth hardness of svp. In *Proceedings*
410 *of the 50th Annual ACM SIGACT Symposium on Theory of Computing*, pages 228–238, 2018.
- 411 [3] Amir Hosein Khas Ahmadi. *Memory-based graph networks*. PhD thesis, University of Toronto
412 (Canada), 2020.
- 413 [4] Sara Ahmadian, Ashkan Norouzi-Fard, Ola Svensson, and Justin Ward. Better guarantees for
414 k-means and euclidean k-median by primal-dual algorithms. *SIAM J. Comput.*, 49(4), 2020.
- 415 [5] Kareem Amin, Corinna Cortes, Giulia DeSalvo, and Afshin Rostamizadeh. Understanding
416 the effects of batching in online active learning. In *International Conference on Artificial*
417 *Intelligence and Statistics*, pages 3482–3492. PMLR, 2020.
- 418 [6] Martin Aumüller, Sariel Har-Peled, Sepideh Mahabadi, Rasmus Pagh, and Francesco Silvestri.
419 Fair near neighbor search via sampling. *SIGMOD Rec.*, 50(1):42–49, 2021.
- 420 [7] Pranjal Awasthi, Moses Charikar, Ravishankar Krishnaswamy, and Ali Kemal Sinop. The hard-
421 ness of approximation of euclidean k-means. In *International Symposium on Computational*
422 *Geometry*, 2015.
- 423 [8] Ramakrishna Bairi, Rishabh Iyer, Ganesh Ramakrishnan, and Jeff Bilmes. Summarization
424 of multi-document topic hierarchies using submodular mixtures. In *Proceedings of the 53rd*
425 *Annual Meeting of the Association for Computational Linguistics and the 7th International*
426 *Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 553–563,
427 2015.
- 428 [9] MohammadHossein Bateni, Hossein Esfandiari, and Vahab S. Mirrokni. Optimal distributed
429 submodular optimization via sketching. In Yike Guo and Faisal Farooq, editors, *Proceedings*
430 *of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining,*
431 *KDD 2018, London, UK, August 19-23, 2018*, pages 1138–1147. ACM, 2018.
- 432 [10] Aditya Bhaskara, Mehrdad Ghadiri, Vahab S. Mirrokni, and Ola Svensson. Linear relaxations
433 for finding diverse elements in metric spaces. In Daniel D. Lee, Masashi Sugiyama, Ulrike
434 von Luxburg, Isabelle Guyon, and Roman Garnett, editors, *Advances in Neural Information*
435 *Processing Systems 29: Annual Conference on Neural Information Processing Systems 2016,*
436 *December 5-10, 2016, Barcelona, Spain*, pages 4098–4106, 2016.
- 437 [11] Filippo Maria Bianchi, Daniele Grattarola, and Cesare Alippi. Spectral clustering with graph
438 neural networks for graph pooling. In *ICML*, pages 874–883. PMLR, 2020.
- 439 [12] Deyu Bo, Xiao Wang, Chuan Shi, Meiqi Zhu, Emiao Lu, and Peng Cui. Structural deep
440 clustering network. In *Proceedings of The Web Conference 2020*, pages 1400–1410, 2020.
- 441 [13] Verónica Bolón-Canedo, Noelia Sánchez-Marroño, and Amparo Alonso-Betanzos. Recent
442 advances and emerging challenges of feature selection in the context of big data. *Knowl. Based*
443 *Syst.*, 86:33–45, 2015.
- 444 [14] James Bradbury, Roy Frostig, Peter Hawkins, Matthew James Johnson, Chris Leary, Dougal
445 Maclaurin, George Necula, Adam Paszke, Jake VanderPlas, Skye Wanderman-Milne, and Qiao
446 Zhang. JAX: composable transformations of Python+NumPy programs, 2018.
- 447 [15] Mark Braverman, Young Kun Ko, Aviad Rubinfeld, and Omri Weinstein. Eth hardness for
448 densest-k-subgraph with perfect completeness. In *Proceedings of the Twenty-Eighth Annual*
449 *ACM-SIAM Symposium on Discrete Algorithms*, pages 1326–1341. SIAM, 2017.
- 450 [16] Mark Braverman, Young Kun Ko, and Omri Weinstein. Approximating the best nash equi-
451 librium in no (log n)-time breaks the exponential time hypothesis. In *Proceedings of the*
452 *twenty-sixth annual ACM-SIAM symposium on Discrete algorithms*, pages 970–982. SIAM,
453 2014.

- 454 [17] Hongyun Cai, Vincent W Zheng, and Kevin Chen-Chuan Chang. Active learning for graph
455 embedding. *arXiv preprint arXiv:1705.05085*, 2017.
- 456 [18] Shaosheng Cao, Wei Lu, and Qiongkai Xu. Grarep: Learning graph representations with
457 global structural information. In *CIKM*, pages 891–900, 2015.
- 458 [19] Shayok Chakraborty, Vineeth Balasubramanian, and Sethuraman Panchanathan. Adaptive
459 batch mode active learning. *IEEE transactions on neural networks and learning systems*,
460 26(8):1747–1760, 2014.
- 461 [20] Ines Chami, Sami Abu-El-Haija, Bryan Perozzi, Christopher Ré, and Kevin Murphy. Machine
462 learning on graphs: A model and comprehensive taxonomy. *arXiv preprint arXiv:2005.03675*,
463 2020.
- 464 [21] Lin Chen, Hossein Esfandiari, Gang Fu, Vahab S Mirrokni, and Qian Yu. Feature cross search
465 via submodular optimization. *arXiv preprint arXiv:2107.02139*, 2021.
- 466 [22] Yijia Chen, Kord Eickmeyer, and Jörg Flum. The exponential time hypothesis and the
467 parameterized clique problem. In *International Symposium on Parameterized and Exact*
468 *Computation*, pages 13–24. Springer, 2012.
- 469 [23] Yuxin Chen and Andreas Krause. Near-optimal batch mode active learning and adaptive
470 submodular optimization. In *ICML*, pages 160–168. PMLR, 2013.
- 471 [24] Richard Church and Charles ReVelle. The maximal covering location problem. In *Papers of*
472 *the regional science association*, volume 32, pages 101–118. Springer-Verlag, 1974.
- 473 [25] Gui Citovsky, Giulia DeSalvo, Claudio Gentile, Lazaros Karydas, Anand Rajagopalan, Afshin
474 Rostamizadeh, and Sanjiv Kumar. Batch active learning at scale. *Advances in Neural*
475 *Information Processing Systems*, 34, 2021.
- 476 [26] David A Cohn, Zoubin Ghahramani, and Michael I Jordan. Active learning with statistical
477 models. *Journal of artificial intelligence research*, 4:129–145, 1996.
- 478 [27] Cody Coleman, Christopher Yeh, Stephen Musmann, Baharan Mirzasoleiman, Peter Bailis,
479 Percy Liang, Jure Leskovec, and Matei Zaharia. Selection via proxy: Efficient data selection
480 for deep learning. *arXiv preprint arXiv:1906.11829*, 2019.
- 481 [28] Marek Cygan, Fedor V Fomin, Łukasz Kowalik, Daniel Lokshantov, Dániel Marx, Marcin
482 Pilipczuk, Michał Pilipczuk, and Saket Saurabh. Lower bounds based on the exponential-time
483 hypothesis. In *Parameterized Algorithms*, pages 467–521. Springer, 2015.
- 484 [29] Chris H. Q. Ding and Hanchuan Peng. Minimum redundancy feature selection from microarray
485 gene expression data. In *2nd IEEE Computer Society Bioinformatics Conference, CSB 2003,*
486 *Stanford, CA, USA, August 11-14, 2003*, pages 523–529. IEEE Computer Society, 2003.
- 487 [30] S Durga, Rishabh Iyer, Ganesh Ramakrishnan, and Abir De. Training data subset selection for
488 regression with controlled generalization error. In *ICML*, pages 9202–9212. PMLR, 2021.
- 489 [31] Hossein Esfandiari, Amin Karbasi, and Vahab Mirrokni. Adaptivity in adaptive submodularity.
490 In *Conference on Learning Theory*, pages 1823–1846. PMLR, 2021.
- 491 [32] Absalom E. Ezugwu, Abiodun M. Ikotun, Olaide O. Oyelade, Laith Abualigah, Jeffery O.
492 Agushaka, Christopher I. Eke, and Andronicus A. Akinyelu. A comprehensive survey of
493 clustering algorithms: State-of-the-art machine learning applications, taxonomy, challenges,
494 and future research prospects. *Engineering Applications of Artificial Intelligence*, 110:104743,
495 2022.
- 496 [33] Bahare Fatemi, Layla El Asri, and Seyed Mehran Kazemi. Slaps: Self-supervision improves
497 structure learning for graph neural networks. In *Neural Information Processing Systems*
498 *(NeurIPS)*, 2021.
- 499 [34] Dan Feldman. *Coresets and Their Applications*. Citeseer, 2010.

- 500 [35] Dan Feldman. Core-sets: Updated survey. In *Sampling techniques for supervised or unsuper-*
501 *vised tasks*, pages 23–44. Springer, 2020.
- 502 [36] Chakib Fettal, Lazhar Labiod, and Mohamed Nadif. Efficient graph convolution for joint node
503 representation learning and clustering. In *Proceedings of the Fifteenth ACM International*
504 *Conference on Web Search and Data Mining*, pages 289–297, 2022.
- 505 [37] Fabrizio Frasca, Emanuele Rossi, Davide Eynard, Ben Chamberlain, Michael Bronstein,
506 and Federico Monti. Sign: Scalable inception graph neural networks. *arXiv preprint*
507 *arXiv:2004.11198*, 2020.
- 508 [38] A. Gammerman, V. Vovk, and V. Vapnik. Learning by transduction. In *Proceedings of the*
509 *Fourteenth Conference on Uncertainty in Artificial Intelligence, UAI’98*, page 148–155, San
510 Francisco, CA, USA, 1998. Morgan Kaufmann Publishers Inc.
- 511 [39] Guojun Gan, Chaoqun Ma, and Jianhong Wu. *Data Clustering: Theory, Algorithms, and*
512 *Applications*. Society for Industrial and Applied Mathematics, 2007.
- 513 [40] Hongyang Gao and Shuiwang Ji. Graph u-nets. In *international conference on machine*
514 *learning*, pages 2083–2092. PMLR, 2019.
- 515 [41] Li Gao, Hong Yang, Chuan Zhou, Jia Wu, Shirui Pan, and Yue Hu. Active discriminative
516 network representation learning. In *IJCAI*, 2018.
- 517 [42] Xing Gao, Wenrui Dai, Chenglin Li, Hongkai Xiong, and Pascal Frossard. ipool–information-
518 based pooling in hierarchical graph neural networks. *IEEE Transactions on Neural Networks*
519 *and Learning Systems*, 2021.
- 520 [43] Yunhao Ge, Yunkui Pang, Linwei Li, and Laurent Itti. Graph autoencoder for graph com-
521 pression and representation learning. In *Neural Compression: From Information Theory to*
522 *Applications–Workshop@ ICLR 2021*, 2021.
- 523 [44] Jonathan Godwin*, Thomas Keck*, Peter Battaglia, Victor Bapst, Thomas Kipf, Yujia Li,
524 Kimberly Stachenfeld, Petar Veličković, and Alvaro Sanchez-Gonzalez. Jraph: A library for
525 graph neural networks in jax., 2020.
- 526 [45] Teofilo F. Gonzalez. Clustering to minimize the maximum intercluster distance. *Theoretical*
527 *Computer Science*, 38:293–306, 1985.
- 528 [46] Aditya Grover and Jure Leskovec. node2vec: Scalable feature learning for networks. In
529 *Proceedings of the 22nd ACM SIGKDD international conference on Knowledge discovery and*
530 *data mining*, pages 855–864, 2016.
- 531 [47] Andrew Guillory. *Active Learning and Submodular Functions*. University of Washington,
532 2012.
- 533 [48] Yuhong Guo and Dale Schuurmans. Discriminative batch mode active learning. In *NIPS*,
534 pages 593–600. Citeseer, 2007.
- 535 [49] Will Hamilton, Zhitao Ying, and Jure Leskovec. Inductive representation learning on large
536 graphs. *NeurIPS*, 30, 2017.
- 537 [50] T. Hastie, R. Tibshirani, and J.H. Friedman. *The Elements of Statistical Learning: Data*
538 *Mining, Inference, and Prediction*. Springer series in statistics. Springer, 2009.
- 539 [51] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers:
540 Surpassing human-level performance on imagenet classification. In *Proceedings of the IEEE*
541 *international conference on computer vision*, pages 1026–1034, 2015.
- 542 [52] Jonathan Heek, Anselm Levskaya, Avital Oliver, Marvin Ritter, Bertrand Rondepierre, Andreas
543 Steiner, and Marc van Zee. Flax: A neural network library and ecosystem for JAX, 2020.
- 544 [53] Dorit S Hochbaum and Anu Pathria. Analysis of the greedy approach in problems of maximum
545 k-coverage. *Naval Research Logistics (NRL)*, 45(6):615–627, 1998.

- 546 [54] Steven CH Hoi, Rong Jin, Jianke Zhu, and Michael R Lyu. Batch mode active learning
547 and its application to medical image classification. In *Proceedings of the 23rd international*
548 *conference on Machine learning*, pages 417–424, 2006.
- 549 [55] Weihua Hu, Matthias Fey, Marinka Zitnik, Yuxiao Dong, Hongyu Ren, Bowen Liu, Michele
550 Catasta, and Jure Leskovec. Open graph benchmark: Datasets for machine learning on graphs.
551 *arXiv preprint arXiv:2005.00687*, 2020.
- 552 [56] Joey Huchette, Haihao Lu, Hossein Esfandiari, and Vahab Mirrokni. Contextual reserve price
553 optimization in auctions via mixed-integer programming. *arXiv preprint arXiv:2002.08841*,
554 2020.
- 555 [57] Piotr Indyk, Sepideh Mahabadi, Mohammad Mahdian, and Vahab S. Mirrokni. Composable
556 core-sets for diversity and coverage maximization. In Richard Hull and Martin Grohe, edi-
557 tors, *Proceedings of the 33rd ACM SIGMOD-SIGACT-SIGART Symposium on Principles of*
558 *Database Systems, PODS’14, Snowbird, UT, USA, June 22-27, 2014*, pages 100–108. ACM,
559 2014.
- 560 [58] A.K. Jain and R.C. Dubes. *Algorithms for Clustering Data*. Prentice Hall advanced reference
561 series. Prentice Hall, 1988.
- 562 [59] Peter Jonsson, Victor Lagerkvist, Gustav Nordh, and Bruno Zanuttini. Complexity of sat
563 problems, clone theory and the exponential time hypothesis. In *Proceedings of the twenty-*
564 *fourth annual ACM-SIAM symposium on Discrete algorithms*, pages 1264–1277. SIAM, 2013.
- 565 [60] Barakeel Fansu Kamhoua, Lin Zhang, Kaili Ma, James Cheng, Bo Li, and Bo Han. Grace:
566 A general graph convolution framework for attributed graph clustering. *ACM Journal of the*
567 *ACM (JACM)*, 2022.
- 568 [61] Tapas Kanungo, David M. Mount, Nathan S. Netanyahu, Christine D. Piatko, Ruth Silverman,
569 and Angela Y. Wu. A local search approximation algorithm for k-means clustering. In
570 *Proceedings of the Eighteenth Annual Symposium on Computational Geometry, SCG ’02*,
571 pages 10–18, New York, NY, USA, 2002. Association for Computing Machinery.
- 572 [62] L. Kaufman and P.J. Rousseeuw. *Finding Groups in Data: An Introduction to Cluster Analysis*.
573 Wiley Series in Probability and Statistics. Wiley, 2009.
- 574 [63] Vishal Kaushal, Rishabh Iyer, Suraj Kothawade, Rohan Mahadev, Khoshrav Doctor, and
575 Ganesh Ramakrishnan. Learning from less data: A unified data subset selection and active
576 learning framework for computer vision. In *2019 IEEE Winter Conference on Applications of*
577 *Computer Vision (WACV)*, pages 1289–1299. IEEE, 2019.
- 578 [64] Krishnateja Killamsetty, S Durga, Ganesh Ramakrishnan, Abir De, and Rishabh Iyer. Grad-
579 match: Gradient matching based data subset selection for efficient deep model training. In
580 *ICML*, pages 5464–5474. PMLR, 2021.
- 581 [65] Krishnateja Killamsetty, Durga Sivasubramanian, Ganesh Ramakrishnan, and Rishabh Iyer.
582 Glistar: Generalization based data subset selection for efficient and robust learning. In *AAAI*,
583 volume 35, pages 8110–8118, 2021.
- 584 [66] Krishnateja Killamsetty, Xujiang Zhao, Feng Chen, and Rishabh Iyer. Retrieve: Coreset
585 selection for efficient and robust semi-supervised learning. *NeurIPS*, 34, 2021.
- 586 [67] Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional
587 networks. *arXiv preprint arXiv:1609.02907*, 2016.
- 588 [68] Günter Klambauer, Thomas Unterthiner, Andreas Mayr, and Sepp Hochreiter. Self-normalizing
589 neural networks. In *NeurIPS*, pages 972–981, 2017.
- 590 [69] Ziyi Kou, Yang Zhang, Lanyu Shang, and Dong Wang. Faircrowd: Fair human face dataset
591 sampling via batch-level crowdsourcing bias inference. In *2021 IEEE/ACM 29th International*
592 *Symposium on Quality of Service (IWQOS)*, pages 1–10, 2021.

- 593 [70] Dongha Lee, Su Kim, Seonghyeon Lee, Chanyoung Park, and Hwanjo Yu. Learnable structural
594 semantic readout for graph classification. In *ICDM*, pages 1180–1185. IEEE, 2021.
- 595 [71] Euiwoong Lee, Melanie Schmidt, and John Wright. Improved and simplified inapproximability
596 for k-means. *Information Processing Letters*, 120:40–43, 2017.
- 597 [72] Jae-Gil Lee, Yuji Roh, Hwanjun Song, and Steven Euijong Whang. Machine learning ro-
598 bustness, fairness, and their convergence. In Feida Zhu, Beng Chin Ooi, and Chunyan Miao,
599 editors, *KDD '21: The 27th ACM SIGKDD Conference on Knowledge Discovery and Data
600 Mining, Virtual Event, Singapore, August 14-18, 2021*, pages 4046–4047. ACM, 2021.
- 601 [73] J.G. Lee and C.G. Chung. An optimal representative set selection method. *Information and
602 Software Technology*, 42(1):17–25, 2000.
- 603 [74] Junhyun Lee, Inyeop Lee, and Jaewoo Kang. Self-attention graph pooling. In *International
604 conference on machine learning*, pages 3734–3743. PMLR, 2019.
- 605 [75] Chuang Liu, Yibing Zhan, Chang Li, Bo Du, Jia Wu, Wenbin Hu, Tongliang Liu, and Dacheng
606 Tao. Graph pooling for graph neural networks: Progress, challenges, and opportunities. *arXiv
607 preprint arXiv:2204.07321*, 2022.
- 608 [76] Liang Liu, Zhao Kang, Jiajia Ruan, and Xixu He. Multilayer graph contrastive clustering
609 network. *Information Sciences*, 613:256–267, 2022.
- 610 [77] Ning Liu, Songlei Jian, Dongsheng Li, Yiming Zhang, Zhiquan Lai, and Hongzuo Xu. Hierar-
611 chical adaptive pooling by capturing high-order dependency for graph representation learning.
612 *IEEE Transactions on Knowledge and Data Engineering*, 2021.
- 613 [78] Yuzong Liu, Rishabh Iyer, Katrin Kirchhoff, and Jeff Bilmes. Svitchboard ii and fisver i:
614 High-quality limited-complexity corpora of conversational english speech. In *Sixteenth Annual
615 Conference of the International Speech Communication Association*, 2015.
- 616 [79] Daniel Lokshtanov, Dániel Marx, Saket Saurabh, et al. Lower bounds based on the exponential
617 time hypothesis. *Bulletin of the EATCS*, (105):41–72, 2011.
- 618 [80] Pasin Manurangsi. Almost-polynomial ratio eth-hardness of approximating densest k-subgraph.
619 In *Proceedings of the 49th Annual ACM SIGACT Symposium on Theory of Computing*, pages
620 954–961, 2017.
- 621 [81] Patrice Marcotte and Gilles Savard. Novel approaches to the discrimination problem. *ZOR —
622 Methods and Models of Operations Research*, 36, 1992.
- 623 [82] Leland McInnes, John Healy, and James Melville. Umap: Uniform manifold approximation
624 and projection for dimension reduction. *arXiv preprint arXiv:1802.03426*, 2018.
- 625 [83] Baharan Mirzasoleiman, Jeff Bilmes, and Jure Leskovec. Coresets for data-efficient training
626 of machine learning models. In *ICML*, pages 6950–6960. PMLR, 2020.
- 627 [84] Giannis Nikolentzos, Polykarpos Meladianos, and Michalis Vazirgiannis. Matching node
628 embeddings for graph similarity. In *Thirty-first AAAI conference on artificial intelligence*,
629 2017.
- 630 [85] Jana Novovicová, Petr Somol, Michal Haindl, and Pavel Pudil. Conditional mutual information
631 based feature selection for classification task. In Luis Rueda, Domingo Mery, and Josef
632 Kittler, editors, *Progress in Pattern Recognition, Image Analysis and Applications, 12th
633 Iberoamericann Congress on Pattern Recognition, CIARP 2007, Valparaiso, Chile, November
634 13-16, 2007, Proceedings*, volume 4756 of *Lecture Notes in Computer Science*, pages 417–426.
635 Springer, 2007.
- 636 [86] Feng Pan, Wei Wang, Anthony K. H. Tung, and Jiong Yang. Finding representative set from
637 massive data. In *Proceedings of the 5th IEEE International Conference on Data Mining
638 (ICDM 2005), 27-30 November 2005, Houston, Texas, USA*, pages 338–345. IEEE Computer
639 Society, 2005.

- 640 [87] Yunsheng Pang, Yunxiang Zhao, and Dongsheng Li. Graph pooling via coarsened graph
641 infomax. In *Proceedings of the 44th International ACM SIGIR Conference on Research and*
642 *Development in Information Retrieval*, pages 2177–2181, 2021.
- 643 [88] Mansheej Paul, Surya Ganguli, and Gintare Karolina Dziugaite. Deep learning on a data diet:
644 Finding important examples early in training. *NeurIPS*, 34:20596–20607, 2021.
- 645 [89] Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion,
646 Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, et al.
647 Scikit-learn: Machine learning in python. *JMLR*, 12:2825–2830, 2011.
- 648 [90] Tiago P Peixoto. Efficient monte carlo and greedy heuristic for the inference of stochastic
649 block models. *Physical Review E*, 89(1):012804, 2014.
- 650 [91] Zhihao Peng, Hui Liu, Yuheng Jia, and Junhui Hou. Attention-driven graph clustering network.
651 In *Proceedings of the 29th ACM International Conference on Multimedia*, pages 935–943,
652 2021.
- 653 [92] Yuji Roh, Kangwook Lee, Steven Euijong Whang, and Changho Suh. Fairbatch: Batch
654 selection for model fairness. In *ICLR*, 2021.
- 655 [93] Erich Schubert and Peter J Rousseeuw. Faster k-medoids clustering: improving the pam, clara,
656 and clarans algorithms. In *Similarity Search and Applications: 12th International Conference,*
657 *SISAP 2019, Newark, NJ, USA, October 2–4, 2019, Proceedings 12*, pages 171–187. Springer,
658 2019.
- 659 [94] Prithviraj Sen, Galileo Namata, Mustafa Bilgic, Lise Getoor, Brian Galligher, and Tina Eliassi-
660 Rad. Collective classification in network data. *AI magazine*, 29(3):93–93, 2008.
- 661 [95] Burr Settles. Active learning literature survey. 2009.
- 662 [96] Oleksandr Shchur and Stephan Günnemann. Overlapping community detection with graph
663 neural networks. *arXiv preprint arXiv:1909.12201*, 2019.
- 664 [97] Oleksandr Shchur, Maximilian Mumme, Aleksandar Bojchevski, and Stephan Günnemann.
665 Pitfalls of graph neural network evaluation. *arXiv preprint arXiv:1811.05868*, 2018.
- 666 [98] Shubhanshu Shekhar, Greg Fields, Mohammad Ghavamzadeh, and Tara Javidi. Adaptive
667 sampling for minimax fair classification. In A. Beygelzimer, Y. Dauphin, P. Liang, and
668 J. Wortman Vaughan, editors, *NeurIPS*, 2021.
- 669 [99] Tom AB Snijders and Krzysztof Nowicki. Estimation and prediction for stochastic blockmodels
670 for graphs with latent block structure. *Journal of classification*, 1997.
- 671 [100] Anton Tsitsulin, John Palowitch, Bryan Perozzi, and Emmanuel Müller. Graph clustering with
672 graph neural networks. *JMLR*, 2023.
- 673 [101] Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *JMLR*, 9(11),
674 2008.
- 675 [102] V.V. Vazirani. *Approximation Algorithms*. Springer Berlin Heidelberg, 2013.
- 676 [103] Petar Velickovic, William Fedus, William L Hamilton, Pietro Liò, Yoshua Bengio, and R Devon
677 Hjelm. Deep graph infomax. *ICLR*, 2(3):4, 2019.
- 678 [104] Chun Wang, Shirui Pan, Ruiqi Hu, Guodong Long, Jing Jiang, and Chengqi Zhang. Attributed
679 graph clustering: A deep attentional embedding approach. *arXiv preprint arXiv:1906.06532*,
680 2019.
- 681 [105] Kai Wei, Rishabh Iyer, and Jeff Bilmes. Submodularity in data subset selection and active
682 learning. In *ICML*, pages 1954–1963. PMLR, 2015.
- 683 [106] Kai Wei, Yuzong Liu, Katrin Kirchhoff, Chris Bartels, and Jeff Bilmes. Submodular subset
684 selection for large-scale speech training data. In *2014 IEEE International Conference on*
685 *Acoustics, Speech and Signal Processing (ICASSP)*, pages 3311–3315. IEEE, 2014.

- 686 [107] Bryan Wilder, Bistra Dilkina, and Milind Tambe. Melding the data-decisions pipeline:
687 Decision-focused learning for combinatorial optimization. In *Proceedings of the AAAI Confer-*
688 *ence on Artificial Intelligence*, volume 33, pages 1658–1665, 2019.
- 689 [108] Felix Wu, Amauri Souza, Tianyi Zhang, Christopher Fifty, Tao Yu, and Kilian Weinberger.
690 Simplifying graph convolutional networks. In *ICML*, pages 6861–6871. PMLR, 2019.
- 691 [109] Yuexin Wu, Yichong Xu, Aarti Singh, Yiming Yang, and Artur Dubrawski. Active learning for
692 graph neural networks via node feature propagation. *arXiv preprint arXiv:1910.07567*, 2019.
- 693 [110] Dongkuan Xu and Yingjie Tian. A comprehensive survey of clustering algorithms. *Annals of*
694 *Data Science*, 2, 2015.
- 695 [111] Rui Xu and Donald Wunsch. Survey of clustering algorithms. *IEEE Transactions on Neural*
696 *Networks*, 16(3):645–678, 2005.
- 697 [112] Fanghua Ye, Chuan Chen, and Zibin Zheng. Deep autoencoder-like nonnegative matrix factor-
698 ization for community detection. In *Proceedings of the 27th ACM international conference on*
699 *information and knowledge management*, pages 1393–1402, 2018.
- 700 [113] Zhitao Ying, Jiaxuan You, Christopher Morris, Xiang Ren, Will Hamilton, and Jure Leskovec.
701 Hierarchical graph representation learning with differentiable pooling. *NeurIPS*, 31, 2018.
- 702 [114] Chong You, Chi Li, Daniel Robinson, and Rene Vidal. Self-representation based unsupervised
703 exemplar selection in a union of subspaces. *IEEE Transactions on Pattern Analysis and*
704 *Machine Intelligence*, 2020.
- 705 [115] Hao Yuan and Shuiwang Ji. Structpool: Structured graph pooling via conditional random
706 fields. In *ICLR*, 2020.
- 707 [116] Sepehr Abbasi Zadeh, Mehrdad Ghadiri, Vahab S. Mirrokni, and Morteza Zadimoghaddam.
708 Scalable feature selection via distributed diversity maximization. In Satinder P. Singh and
709 Shaul Markovitch, editors, *AAAI*, pages 2876–2883, 2017.
- 710 [117] Liang Zhang, Xudong Wang, Hongsheng Li, Guangming Zhu, Peiyi Shen, Ping Li, Xiaoyuan
711 Lu, Syed Afaq Ali Shah, and Mohammed Bennamoun. Structure-feature based graph self-
712 adaptive pooling. In *Proceedings of The Web Conference*, pages 3098–3104, 2020.
- 713 [118] Wentao Zhang, Zhi Yang, Yexin Wang, Yu Shen, Yang Li, Liang Wang, and Bin Cui. Grain:
714 Improving data efficiency of graph neural networks via diversified influence maximization.
715 *arXiv preprint arXiv:2108.00219*, 2021.
- 716 [119] Xiaotong Zhang, Han Liu, Qimai Li, and Xiao-Ming Wu. Attributed graph clustering via
717 adaptive graph convolution. *arXiv preprint arXiv:1906.01210*, 2019.

718 A More results

719 A.1 Selecting More Representatives

720 We already presented results for the case where we select $2c$ representatives from each dataset, where c
721 is the number of classes. Here, we present more results for the case where we select $5c$ representatives
722 from each class. We limit our baselines to the ones that were either more competitive/informative
723 or took less time to run. The results are in Table 2. We can observe that the performance for all
724 models improves when more representatives are selected. We can also observe that the gap between
725 the models shrinks when more representatives are selected (even the random baseline now shows
726 a competitive performance). Nevertheless, RS-GNN shows a similar trend as when we selected $2c$
727 representatives and outperforms the baselines when averaged across datasets.

Table 2: For each dataset, each algorithm selects $5c$ representatives. Then, we train a GCN model on the labels of the selected representatives. The reported metric is the test accuracy of the GCN models. Surrogate function selectors mark in blue use DGI embeddings.

Selector	Cora	CiteSeer	Pubmed	Photos	PC	CS	Physics	Arxiv	Avg.
Random	66.3 \pm 4.6	47.3 \pm 6.6	62.0 \pm 8.3	84.0 \pm 3.5	76.8 \pm 3.4	83.9 \pm 2.3	84.4 \pm 5.8	54.6 \pm 1.2	69.9
Popular	64.1 \pm 0.6	43.2 \pm 1.4	63.4 \pm 0.9	43.9 \pm 1.1	54.1 \pm 1.4	78.6 \pm 0.5	50.5 \pm 0.0	42.6 \pm 1.4	55.0
KMedoid	62.2 \pm 1.3	42.1 \pm 3.1	60.8 \pm 0.3	78.8 \pm 0.6	74.4 \pm 1.2	63.9 \pm 1.7	64.3 \pm 0.6	50.6 \pm 0.5	62.1
KMeans	66.7 \pm 1.1	53.8 \pm 1.3	71.6 \pm 0.5	84.8 \pm 1.0	81.4 \pm 0.7	76.7 \pm 1.0	80.6 \pm 0.3	56.8 \pm 0.3	71.6
MC-Cos	71.3 \pm 2.8	59.0 \pm 2.8	64.8 \pm 0.3	87.5 \pm 0.4	78.3 \pm 0.7	88.4 \pm 0.2	92.4 \pm 0.4	56.2 \pm 0.5	74.7
KMeans	74.9 \pm 3.1	54.3 \pm 4.8	69.3 \pm 3.2	89.4 \pm 1.3	82.0 \pm 1.4	89.2 \pm 0.8	92.1 \pm 0.7	55.1 \pm 1.1	75.8
MC-Cos	75.8 \pm 1.7	59.8 \pm 2.5	70.6 \pm 2.4	90.0 \pm 1.2	82.3 \pm 1.1	89.4 \pm 0.8	86.1 \pm 1.9	55.8 \pm 0.7	76.2
MinCUT	66.8 \pm 4.3	48.7 \pm 6.6	69.0 \pm 4.2	16.7 \pm 8.1	16.5 \pm 10.2	88.6 \pm 1.2	92.3 \pm 1.0	55.3 \pm 7.5	56.7
DMoN	70.1 \pm 3.8	54.9 \pm 4.6	70.2 \pm 4.0	86.7 \pm 2.4	74.8 \pm 5.8	89.7 \pm 0.7	90.9 \pm 1.8	57.1 \pm 1.0	74.3
RS-GNN	77.3 \pm 1.9	62.7 \pm 2.3	68.7 \pm 2.4	90.6 \pm 0.5	83.0 \pm 1.6	90.1 \pm 0.6	92.4 \pm 0.8	56.6 \pm 1.2	77.7

728 A.2 Results: RSL in the Absence of a Graph Structure

729 In several applications, a natural graph may not be available. For semi-supervised node classification,
730 it has recently been shown that even without a natural graph, one can still leverage GNNs by learning
731 both a graph structure and GNN parameters simultaneously [33]. We extend the aforementioned
732 results to the RSL problem. In particular, we assume we have access *only* to the node features of our
733 datasets, and not to their graph structures. The baselines select representatives only based on the
734 node features. For RS-GNN, we first create a kNN similarity graph between the nodes and then run
735 the model on the node features and the created graph. The kNN graph is computed once and then
736 fixed during training; we leave further experiments with learning the graph structure as future work.
737 Notice that all models have access to the same context.

738 Once the representatives are selected, we train GCN classifiers on the labels of the selected nodes.
739 For the baselines, we run our classification GCN in two settings: (1) the only edges in the graph are
740 self-loops, (2) the edges in the graph are those from a kNN similarity graph. We call the former a
741 multi-layer perceptron (MLP) and the latter a (kNN-)GCN. We report results for the top baselines
742 that operate on node features *only*.

743 Since operating without graph structure reduces the signal in the datasets, we allow each model to
744 select $5c$ representatives for the experiments in this section. The results are presented in Table 5.
745 The results show that selecting representatives using RS-GNN performs consistently well on all
746 datasets and provides a boost compared to our baselines on many of the datasets. This confirms that
747 even a similarity graph can still be helpful in improving RSL and that RS-GNN is an effective RSL
748 algorithm for datasets where a graph structure is not available.

749 A.3 Label Coverage

750 Once a set of representatives are selected from a dataset, we define a specific class label to be *covered*
751 if at least one of the representatives belongs to that class. We define *label coverage* as the percentage
752 of class labels that are covered by the selected representatives. In Table 3, we report the label coverage
753 results for the baselines and our model when selecting $k = 2c$ representatives. RS-GNN performs
754 well in terms of selecting points that cover all labels, with a coverage of 100.0% on four of the
755 datasets. We observe an interesting phenomenon on the Arxiv dataset where (unlike other datasets)
756 many baselines outperform RS-GNN in terms of label coverage (especially GCC), but RS-GNN
757 shows a higher accuracy compared to the baselines. We believe this is due to the high label imbalance
758 in this dataset (the largest class has 27321 examples and the smallest has 29 examples). Future work
759 can extend our work for optimizing macro accuracy across classes.

760 A.4 Visualization

761 In Figure 3(a), we use UMAP [82] to visualize the learned embeddings and the selected representatives
762 of RS-GNN for the Cora dataset. The colors in the plot represent the class to which the node or the

Table 3: Label coverage of different RS algorithms when $k = 2c$. The surrogate function selectors marked in blue use DGI embeddings.

Selector	Cora	CiteSeer	Pubmed	Photos	PC	CS	Physics	Arxiv
Random	86.4±10.8	84.2±12.7	85.0±17.0	80.6±11.1	64.0±9.4	72.7±9.4	75.0±14.3	55.0±4.9
Popular	85.7±0.0	50.0±0.0	66.7±0.0	37.5±0.0	30.0±0.0	66.7±0.0	20.0±0.0	15.0±0.0
KMedoid	100.0±0.0	100.0±0.0	66.7±0.0	87.5±0.0	80.0±0.0	73.3±0.0	60.0±0.0	50.0±0.0
KMeans	100.0±0.0	83.3±0.0	100.0±0.0	75.0±0.0	70.0±0.0	46.7±0.0	60.0±0.0	65.0±0.0
FFS	84.3±11.3	88.3±9.5	71.7±12.2	93.1±8.6	74.5±9.4	36.3±3.4	67.0±9.8	63.0±2.3
MC-RBF	100.0±0.0	100.0±0.0	66.7±0.0	75.0±0.0	60.0±0.0	66.7±0.0	60.0±0.0	52.5±0.0
MC-Cos	97.1±5.9	98.3±5.1	100.0±0.0	75.0±0.0	80.0±0.0	66.7±0.0	100.0±0.0	55.2±0.8
KMedoid	79.3±11.8	58.3±10.1	93.3±13.7	100.0±0.0	80.0±9.2	83.7±7.0	87.0±14.9	OOM
KMeans	100.0±0.0	90.0±8.4	100.0±0.0	96.9±5.5	82.5±6.4	99.3±2.0	100.0±0.0	54.2±5.0
FFS	96.4±6.3	86.7±10.3	66.7±18.7	77.5±7.7	65.0±6.1	97.0±3.4	95.0±8.9	52.0±5.1
MC-RBF	71.4±0.0	81.7±5.1	66.7±0.0	75.0±0.0	60.0±0.0	52.0±4.6	85.0±8.9	28.1±4.0
MC-Cos	100.0±0.0	100.0±0.0	100.0±0.0	87.5±0.0	82.5±5.5	93.7±1.5	85.0±11.0	37.5±2.4
MinCUT	78.6±8.7	83.3±15.3	95.0±12.2	12.5±0.0	10.0±0.0	85.7±5.0	99.0±4.5	42.2±4.6
FeatProp	85.7±0.0	83.3±0.0	100.0±0.0	100.0±0.0	70.0±0.0	60.0±0.0	80.0±0.0	45.0±0.0
SDCN	85.7±9.5	86.7±10.5	93.3±14.0	82.5±8.7	67.5±10.4	72.7±9.1	87.5±10.4	56.7±2.9
EGAE	98.6±4.5	93.3±8.6	86.7±17.2	88.8±4.0	85.0±8.5	76.7±11.4	86.0±9.7	59.2±5.8
GCC	100.0±0.0	95.0±8.0	100.0±0.0	87.5±0.0	85.0±7.1	92.7±3.8	100.0±0.0	72.5±3.4
DMoN	90.0±9.4	90.0±10.0	91.7±14.8	87.5±5.7	76.0±9.9	90.3±4.6	97.0±7.3	60.2±6.3
RS-GNN	100.0±0.0	90.0±8.4	100.0±0.0	98.8±3.9	82.5±9.7	100.0±0.0	100.0±0.0	54.8±3.6

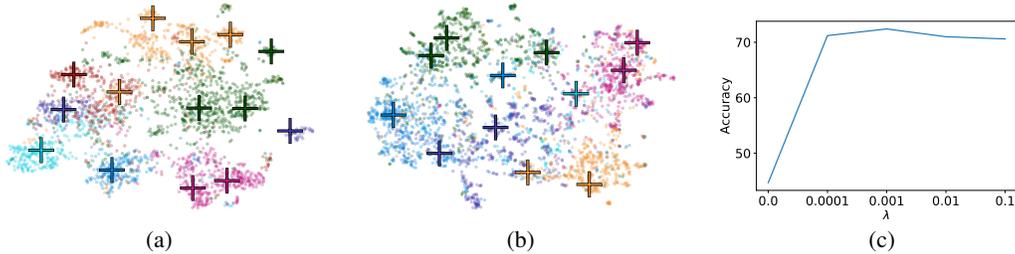


Figure 3: (a) A UMAP visualization of the node embeddings and the selected representatives for Cora (colors represent the class to which the nodes/representatives belong), (b) for CiteSeer, (c) RS-GNN results on Cora for different values of λ .

763 representative belongs. According to the visualization, the nodes from each class have formed one
 764 or more dense clusters and our model has selected a representative from (almost all of) these dense
 765 regions of points. More visualizations can be found in supplementary material.

766 A.5 Time Complexity and Memory Scalability

767 Let m represent the number of nodes, d represent the average degree of nodes, n represent the
 768 number of features (and, for simplicity, the embedding dimension), and k represent the number
 769 of representatives. The time complexity of each epoch in Algorithm 1 is governed by $O(mnd)$
 770 for computing graph convolutions, $O(mn^2)$ for computing node projections and bilinear functions,
 771 and $O(mnk)$ for assigning nodes to representatives. Overall, this gives a time complexity of
 772 $O(mn(d + n + k))$ for each epoch. For large datasets where a large number of representatives is
 773 needed (i.e. $k > n$ and $k > d$), the time complexity becomes $O(mnk)$. The memory complexity is
 774 $O(mk)$ for constructing and storing the matrix of distances from each node to its representatives.

775 When m and k are both large, Algorithm 1 may exhaust the accelerator memory due to its $O(mk)$
 776 complexity. To reduce the memory usage and allow for applying RS-GNN to such settings, we
 777 modify Algorithm 1 to Algorithm 2 (in the Appendix). The main modifications include: 1- as is
 778 common in the GNN scalability literature (see, e.g., [37]), we replace the GCN modules with a

Table 4: Ablation study results for RS-GNN.

Ablation	Cora	CiteSeer	PubMed	Photos	PC	CS	Physics	Avg.
NoNorm	61.6 \pm 7.7	41.5 \pm 3.5	59.6 \pm 4.7	82.5 \pm 2.1	71.3 \pm 2.8	86.6 \pm 2.4	90.8 \pm 1.5	70.6
ConstNorm	59.3 \pm 6.9	44.3 \pm 4.2	61.7 \pm 3.9	84.2 \pm 1.4	75.0 \pm 3.1	86.7 \pm 1.3	90.5 \pm 2.1	71.7
Full Model	72.4 \pm 3.0	54.7 \pm 3.9	65.8 \pm 3.0	86.3 \pm 1.4	74.3 \pm 1.7	89.3 \pm 0.8	90.0 \pm 2.6	76.1

Table 5: Classification accuracies when a graph structure is not provided as input. Selecting 5c representatives. Bold numbers indicate statistically significant winner(s) following a t-test (p-value=0.05).

Selector	Model	Cora	Citeseer	Pubmed	Photos	PC	CS	Physics	Avg.
Random	MLP	41.9 \pm 2.9	37.4 \pm 4.0	51.9 \pm 5.0	57.5 \pm 4.0	55.2 \pm 4.4	76.1 \pm 2.7	76.7 \pm 5.2	56.7
Random	GCN	57.7 \pm 4.1	57.7 \pm 4.5	59.6 \pm 4.6	79.2 \pm 3.4	71.8 \pm 4.4	84.8 \pm 2.2	86.0 \pm 2.9	71.0
KMedoid	MLP	39.3 \pm 0.9	33.4 \pm 0.7	46.1 \pm 1.0	40.6 \pm 0.9	47.5 \pm 0.8	62.8 \pm 1.6	67.0 \pm 0.8	48.1
KMedoid	GCN	52.5 \pm 1.4	57.4 \pm 0.8	55.0 \pm 0.7	72.3 \pm 0.6	63.4 \pm 1.8	66.2 \pm 2.9	70.8 \pm 0.3	62.5
KMeans	MLP	42.4 \pm 0.9	40.1 \pm 1.1	58.5 \pm 1.3	70.0 \pm 1.0	63.7 \pm 0.8	68.5 \pm 0.7	77.4 \pm 0.3	60.1
KMeans	GCN	56.5 \pm 0.8	55.9 \pm 0.8	72.7 \pm 0.3	80.9 \pm 0.6	75.8 \pm 0.6	80.0 \pm 0.4	83.7 \pm 0.8	72.2
FFS	MLP	39.2 \pm 3.7	44.0 \pm 3.1	43.1 \pm 3.0	60.2 \pm 3.8	55.8 \pm 1.9	56.3 \pm 0.9	77.7 \pm 2.5	53.8
FFS	GCN	56.9 \pm 2.6	62.3 \pm 3.7	48.8 \pm 3.9	80.7 \pm 2.0	74.8 \pm 2.5	59.8 \pm 2.5	84.0 \pm 2.1	66.8
MC-Cos	MLP	46.5 \pm 2.2	47.5 \pm 2.3	52.1 \pm 0.7	54.0 \pm 2.4	58.2 \pm 1.2	83.0 \pm 0.5	86.8 \pm 0.7	61.2
MC-Cos	GCN	62.0 \pm 1.7	63.0 \pm 2.5	59.3 \pm 1.3	79.5 \pm 0.5	75.4 \pm 1.1	87.6 \pm 0.3	92.8 \pm 0.3	74.2
RS-GNN	GCN	64.6 \pm 2.4	64.3 \pm 2.1	65.1 \pm 3.0	82.2 \pm 2.0	75.5 \pm 2.1	88.3 \pm 1.6	89.9 \pm 1.9	75.7

779 SGC module [108] and pre-compute the graph convolutions F for the original graph, 2- we create
780 corrupted graphs $nCorrupt$ times and pre-compute the graph convolutions F' for the corrupted graphs,
781 3- at the beginning of each epoch, we randomly select a subset F'' of F' to make the size match
782 that of F , 4- we batch the data and compute the loss and gradients for each batch separately and
783 then aggregate the gradients and update the parameters. Assuming the batch size is b , the memory
784 complexity for Algorithm 2 reduces to $O(bk)$ as each batch can be computed separately. Moreover,
785 Algorithm 2 is also amenable to parallelization on multiple accelerators by distributing the batches
786 across the accelerators.

787 Note that computing the loss separately for each batch corresponds to approximating s and μ with
788 the batch data as opposed to the entire data, but the approximation is expected to be close to the
789 true value if the batch sizes are large enough. To this end, the batch size can be set to the largest
790 number that does not exhaust the memory. We tested the Algorithm 2 version of RS-GNN on the
791 Arxiv dataset when setting $nCorrupt$ to 10 and $b = m/8$ (distributing over 8 accelerators). In terms
792 of performance, we obtained an accuracy of 52.9 ± 1.6 (compared to 52.6 ± 1.2 for the original
793 algorithm) showing that the approximations do not result in a performance degradation.

794 A.6 Ablation Study: CenterNorm and the Value of λ

795 We conduct an ablation study to verify the role of CenterNorm. To ablate CenterNorm, we run our
796 model under two settings: (1) we do not normalize (we refer to this as “NoNorm”), and (2) we
797 divide all the values in the embedding matrix by a constant number corresponding to the mean of the
798 ℓ_2 -norms of the embeddings (we refer to this as “ConstNorm”). According to the results in Table 4,
799 our model benefits from CenterNorm and CenterNorm is more effective than ConstNorm because the
800 per-node ℓ_2 normalization of CenterNorm helps bring the nodes within one cluster closer to each
801 other which helps in identifying clusters and selecting representatives.

802 To verify the sensitivity of RS-GNN to the value of λ used in equation ??, we ran RS-GNN on Cora
803 with different values for λ . The results are presented in Figure 3(b). When $\lambda = 0$, the representatives
804 will not receive gradients and hence the model ends up selecting random representatives. Therefore,
805 the accuracy is quite low. For non-zero values of λ , we observe that the model is not highly sensitive
806 to the value of λ and achieves good results for values in a large range. The model reaches its highest
807 performance around $\lambda = 0.001$, and then the performance starts to slightly decrease for larger values
808 of λ .

Table 6: Normalized mutual information (NMI) scores between the ground-truth labels of nodes and their cluster assignments. We take the results of the baselines from [100]. Bold numbers indicate the winners. We use — to indicate that the results were not reported (either because the model did not converge, or because the model did not scale to the dataset, or because the model was not tested on the dataset).

Selector	Context	Cora	CiteSeer	Pubmed	Photos	PC	CS	Physics	Avg.
KMeans	X	18.5	24.5	19.4	28.8	21.1	35.7	30.6	25.5
SBM	A	36.2	15.3	16.4	59.3	48.4	58.0	45.4	39.9
MinCut	X, A	35.8	25.9	25.4	—	—	64.6	48.3	—
AGC	X, A	34.1	25.5	18.2	59.0	51.3	43.3	—	—
DAEGC	X, A	8.3	4.3	4.4	47.6	42.5	36.3	—	—
SDCN	X, A	27.9	31.4	19.5	41.7	24.9	59.3	50.4	36.4
NOCD	X, A	46.3	20.0	25.5	62.3	44.8	70.5	51.9	45.9
DMoN	X, A	48.8	33.7	29.8	63.3	49.3	69.1	51.9	49.4
RS-GNN	X, A	55.4\pm0.8	41.3 \pm 1.0	26.1 \pm 3.6	58.3 \pm 1.6	50.1 \pm 0.9	75.8\pm1.2	56.9\pm3.3	52.0

809 **Attributed Graph Clustering:** RS-GNN can also be used for attributed graph clustering: we cluster
810 the nodes in each dataset into c groups (recall that c is the number of classes) by assigning each
811 node to its closest representative. While our main motivation is RSL, for completeness we show the
812 performance of our model for attributed graph clustering as well. We note that different works on
813 attributed graph clustering use different settings that are not directly comparable (e.g., some works use
814 the same hyperparameters for all datasets, whereas some other works optimize the hyperparameters
815 for each dataset and report the best test performance). Since our setting is similar to that of [100]
816 and the results for many models have been provided in that work, we compare RS-GNN against the
817 model proposed in [100] and their baselines. This includes KMeans, SBM (this works by estimating
818 [90] a constrained Stochastic Block Model [99] with given number of k clusters), MinCut, AGC
819 [119], DAEGC [104], SDCN, NOCD [96], and DMoN [100].

820 Table 6 shows a comparison of RS-GNN with the baselines in terms of the normalized mutual
821 information (NMI) score (an established score for measuring and comparing clustering algorithms)
822 between the cluster assignments and the node labels. From the results, we can observe that RS-GNN
823 also shows a good performance for attributed graph clustering.

824 B More Related Work

825 Other works that are related to our work can be grouped as follows.

826 **Feature selection:** RS and feature selection are transposed views of a similar problem when it
827 comes to compressing or summarizing datasets. Both have been studied extensively via *filter* and
828 *wrapper* methods: an evaluation based on final task performance or on some proxy metric such
829 as correlation, redundancy, coverage (or more general submodular functions), or the distance of
830 selected entities [73, 86, 13, 9] as well as their mutual information or correlation with the prediction
831 labels [29, 85]. While these methods tackle the diversity of the sample set [1, 57, 116, 10], there has
832 also been extensive attention on taking fairness constraints into account as well [69, 92, 72, 98, 6].

833 **Supervised data subset selection:** Given a large dataset of labeled training examples, a class of RS
834 models aim at selecting a small representative set from the dataset to reduce training time without
835 substantially sacrificing model accuracy (see, e.g., [65, 64, 105, 63, 30, 88, 83]). For example, [65]
836 aim at selecting a set of training data points such that a model trained on these examples generalizes
837 well to the validation set and [83] aim at selecting a set of training data points whose gradients
838 approximate the gradient of the full dataset. These models have been also applied to mini-batch active
839 learning where a small set of labeled data points are assumed to be initially provided and then the
840 next batches are selected based on pseudo-labels predicted by the model trained on the data available
841 so far. While these models assume the data labels are available when selecting representatives, in this
842 paper we assume no labels are provided as input.

843 **Graph pooling:** A technique commonly used in graph representation learning (especially for learning
844 a representation for the entire graph) is graph pooling [75], where the nodes of the graph are iteratively

845 coarsened into “super-nodes”. The parameters for the pooling operation are trained with the rest of the
846 model parameters either to minimize a supervised loss (e.g., graph classification) or an unsupervised
847 loss (e.g., graph reconstruction) [40, 43]. Graph pooling techniques can be classified into two
848 categories: 1- clustering pooling (these are in the same vein as the graph clustering algorithms
849 discussed earlier) and 2- node drop pooling. Clustering pooling approaches [113, 70, 115, 3, 11, 77]
850 employ a differentiable graph clustering algorithm and consider each cluster to be a super-node;
851 these approaches can be re-purposed for RS by selecting the node closest in the latent space to each
852 super-node as a representative. Node drop clustering approaches [40, 74, 87, 117, 42] operate by
853 dropping unimportant nodes and retaining the important nodes as the super-nodes. Graph pooling
854 approaches have been mostly developed for smaller-sized graphs such as molecules that exhibit
855 specific properties. For example, many of these approaches employ a GNN that assigns importance
856 scores to each node, and then select the top-k most similar nodes. Such an approach assigns similar
857 importance scores to highly similar nodes and results in sampling only from some parts of the graph.
858 While this might be a reasonable approach for molecule classification tasks (as it makes the model
859 focus on a few important sub-structures), it may not select a subset of the nodes that cover the entire
860 graph (which is the desired property in our work). Nevertheless, in our experiments, we compare
861 against several graph pooling approaches from both categories, both for RS and clustering tasks.

862 C CenterNorm Motivation

863 With the joint loss function used in the main text, the model can trivially reduce \mathcal{L}_{SEL} by making the
864 values in the embedding matrix \mathbf{H} arbitrarily small. That is because multiplying a small constant to
865 \mathbf{H} may not change \mathcal{L}_{EMB} substantially, but it can make the distances between the nodes arbitrarily
866 small, resulting in a low value for \mathcal{L}_{SEL} even for a random representative embedding matrix \mathbf{R} .

867 One way to avoid the aforementioned problem is by normalizing the embedding matrix \mathbf{H} before
868 using it for selection. However, one should be careful about the choice of the normalization to avoid
869 losing useful information. Before explaining how we normalize the embeddings, we describe a
870 property of DGI embeddings that motivates our normalization.

871 Figure 4(b) shows a UMAP plot of the DGI embeddings \mathbf{H} for the nodes in the original graph of Cora
872 and the embeddings \mathbf{H}' for the nodes in a corrupted Cora graph (red). The \mathbf{H}' embeddings form a
873 large cluster that is mostly in between the \mathbf{H} embeddings and the \mathbf{H} embeddings form small size
874 clusters that are placed around the large cluster of \mathbf{H}' . To understand why this happens, notice that
875 when we shuffle the node features for creating corrupted graphs for DGI training, the node features
876 of the neighbors of each node are a random subsample of the node features in the graph. Therefore,
877 the GNN aggregation function applied on the projected node embeddings makes the embeddings go
878 toward the mean of projected embeddings. This makes the corrupted node embeddings \mathbf{H}' form a
879 large cluster in the middle and the embeddings \mathbf{H} be placed outside and around this cluster.

880 Besides visual inspection, we also cluster the node embeddings \mathbf{H} for Cora into 7 clusters using
881 KMeans (7 is the number of classes in Cora; this provides good clusters with a normalized mutual
882 information of 55.95 with the node labels). Then we compute the distance between each cluster
883 center and the mean of these centers and obtain the following seven distances: 1.4, 1.2, 1.3, 1.2, 0.6,
884 1.4, 1.5. All cluster centers are at a good distance from the mean and, with the exception of one
885 cluster, they are at a similar distance from the mean. We then subtract the mean and compute the
886 angle between the cluster centers. We observe that the minimum angle between two cluster centers is
887 60.1 degrees and the average angle is 76.6 degrees.

888 The above analysis motivates the CenterNorm normalization outlined in the main text. Furthermore,
889 the analysis shows that DGI is a good candidate for RSL as it groups data points into small-sized
890 dense clusters in the latent space, thus an RSL algorithm can select representatives from each of the
891 dense clusters.

892 D Proof of the Theorem

893 **Theorem D.1.** *There is no polynomial-time representative selection algorithm for FoN with an*
894 *approximation factor better than $\omega(n^{-1/\text{poly} \log \log n})$, unless the ETH fails.*

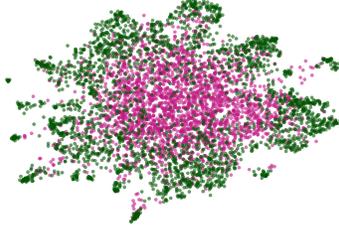


Figure 4: A UMAP plot of the DGI node embeddings for the nodes in the original graph of Cora (green) and the nodes in a corrupted Cora graph (red).

895 We start with two lemmas before proceeding to prove this result. We use (i, j) to represent a data
896 point in the FoN problem with value 1 on features i and j .

897 **Lemma D.2.** *Let \mathcal{S} be the set of data points selected by an RSL algorithm. Let (i, j) be a data point
898 such that for all t we have $(i, t) \notin \mathcal{S}$, or for all t we have $(t, j) \notin \mathcal{S}$. Then the label of (i, j) is
899 independent of the labels of \mathcal{S} .*

900 *Proof.* Let \mathcal{S} be the set of data points selected by an RSL algorithm. Let (i, j) be a data point such
901 that for all t we have $(i, t) \notin \mathcal{S}$. This means that the labels of data points in \mathcal{S} are independent
902 of the type of feature i . Recall that the type of feature i is chosen independently and uniformly at
903 random. Hence, conditioned on the labels in \mathcal{S} , the label of (i, j) is 0 with probability $1/2$ and 1 with
904 probability $1/2$. Similar argument holds when for all t we have $(t, j) \notin \mathcal{S}$. \square

905 **Lemma D.3.** *Let $(i_0, i_1), (i_1, i_2), (i_2, i_3), \dots, (i_{l-1}, i_l)$ be a sequence of data points such that for
906 all $t \in \{1, \dots, l\}$ we have $(i_{t-1}, i_t) \in \mathcal{S}$. Given the labels of the data points in \mathcal{S} we can infer the
907 label of (i_0, i_l) .*

908 *Proof.* Consider two data points (i, j) and (j, t) . If the labels of both data points are 1, then the
909 features i, j and t have the same type. Hence, the label of (i, t) is 1 too. If the labels of both of them
910 are 0, then the type of features i and j are different, and the type of features j and t are different.
911 Hence, the type of features i and t are the same, which means the label of (i, t) is 1. A similar
912 argument shows that if either (i, j) or (j, t) has label 1 and the other has label 0, then the label of
913 (i, t) is 0. Thus, knowing the labels of (i, j) and (j, t) determines the label of (i, t) . Applying this
914 inductively proves the lemma. \square

915 *Proof of Theorem 4.2.* The proof goes by reducing the densest k -subgraph problem to FoN. In the
916 densest k -subgraph problem, we have an unweighted graph \mathbb{G} , and the goal is to find a subgraph of
917 \mathbb{G} with k vertices and the maximum number of edges. We say an algorithm is an α -approximation
918 algorithm for the densest k -subgraph problem if it returns a subgraph with k vertices where the
919 number of edges is at least α times that of the densest k -subgraph. It is known that there is no
920 $\omega(n^{-1/\text{poly log log } n})$ -approximation polynomial-time algorithm for the densest k -subgraph problem
921 unless ETH fails [80].

922 Next, we show how to transform an input of the densest k -subgraph problem to an input of FoN,
923 and then show how to transform an approximate solution for FoN to an approximate solution
924 for the densest k -subgraph problem while only increasing the approximation factor by a constant.
925 Therefore an $\omega(n^{-1/\text{poly log log } n})$ -approximation polynomial-time algorithm for the FoN implies an
926 $\omega(n^{-1/\text{poly log log } n})$ -approximation polynomial-time algorithm for the densest k -subgraph problem,
927 which does not exist unless ETH fails.

928 Let $\mathbb{G} = (V, E)$ be an input to the densest k -subgraph problem.⁶ For each vertex in \mathbb{G} we define
929 a feature and for each edge in \mathbb{G} we construct a data point. For each data point corresponding to
930 an edge (u, v) , the value of the features corresponding to vertices u and v are 1 and the value of all
931 other features are 0. As defined in the FoN problem the type (red or blue) of each feature is chosen
932 independently and uniformly at random.

⁶Note that graph \mathbb{G} is not an attributed graph, rather a simple graph which is an input to the densest k -subgraph problem.

933 Let $\mathbb{H} = (V_{\mathbb{H}}, E_{\mathbb{H}})$ be a densest k -subgraph of \mathbb{G} and let \mathbb{F} be a maximal spanning forest of \mathbb{H} .
 934 Note that since there is no cycle in \mathbb{F} , the number of edges in \mathbb{F} is at most $k - 1$. Moreover, since
 935 \mathbb{F} is a maximal forest of \mathbb{H} , for each edge e in \mathbb{H} , there is a path between the endpoints of e in \mathbb{F}
 936 (otherwise we could add e to \mathbb{F}). Hence, if we query the data points corresponding to the edges of \mathbb{F} ,
 937 by Lemma D.3, we can determine the label of all the edges in \mathbb{H} , which is an $\frac{|E_{\mathbb{H}}|}{|E|}$ fraction of all data
 938 points. This gives us a solution with $\overline{\text{Acc}} \geq \Omega\left(\frac{|E_{\mathbb{H}}|}{|E|}\right)$.

939 Let \mathcal{S} be the set of data points selected by an α -approximation RSL algorithm, and $V_{\mathcal{S}}$ be the set
 940 of vertices adjacent to the edges corresponding to the data points in \mathcal{S} . By Lemma D.2, if the edge
 941 corresponding to a data point has one (or two) endpoints in $V \setminus V_{\mathcal{S}}$, then the label of that data point
 942 is independent of the labels of \mathcal{S} . Hence, the number of data points whose label is not independent
 943 of the labels in \mathcal{S} is at most the number of edges induced by $V_{\mathcal{S}}$. We denote this edge set by $E_{\mathcal{S}}$.
 944 Recall that \mathcal{S} is an α -approximate solution, i.e., $|E_{\mathcal{S}}| = \Omega(\alpha|E_{\mathbb{H}}|)$. On the other hand, $|\mathcal{S}| \leq k$ and
 945 hence $|V_{\mathcal{S}}| \leq 2k$. One can decompose the induced subgraph of $V_{\mathcal{S}}$ into $\binom{4}{2} = 6$ subgraphs each with
 946 k vertices, and pick the one with the maximum number of edges. This gives an $\Omega(\alpha)$ -approximate
 947 solution to the densest k -subgraph problem. \square

948 E Implementation Details

949 **Baselines:** For KMeans, we select the
 950 closest node to each cluster center as
 951 a representative. For FFS and Max-
 952 Cover, we select representatives se-
 953 quentially. In FFS, the next repre-
 954 sentative is the node farthest away (by
 955 Euclidean distance) from the closest
 956 representative in the current set. In
 957 MaxCover, the next representative is
 958 the node that increases the coverage of
 959 the non-selected nodes the most. Note
 960 that the sequential nature of FFS and
 961 MC makes them less amenable to par-
 962 allelization. Also note that when we
 963 run surrogate function baselines us-
 964 ing DGI embeddings as context, their
 965 selections are informed by both node
 966 features and the graph structure. For
 967 DMoN and MinCut models, we com-
 968 pute cluster centers by averaging the
 969 node embeddings with respect to the
 970 (hard) cluster assignments, and then
 971 select the closest point to each cluster
 972 center as a representative.

973 We implemented our model and the
 974 baselines in Jax/Flax [14, 52] and
 975 used the Jraph library [44] for our
 976 GNN operations. Our experiments were done on a JellyFish TPU for all datasets except for the Arxiv
 977 dataset where we used a DragonFish TPU as the experiments with the Arxiv dataset require more
 978 memory. For our DGI model, we used a single-layer GCN model with SeLU activations [68]. For
 979 the experiments that had access to the original graph structure, we set the DGI hidden dimension to
 980 512 for all datasets except for the Arxiv dataset where we set it to 256 to reduce memory usage. For
 981 the experiments with no access to the original graph structure, we set the DGI hidden dimension to
 982 128 as there exists less signal in this case. We trained the DGI models for 2000 epochs both for our
 983 model and the baselines. For KMeans and KMedoid, we used the implementation in scikit-learn [89]
 984 and scikit-learn-extra⁷ respectively. To reduce the quadratic time complexity of MaxCover, we apply
 985 MaxCover on a k -nearest neighbors similarity graph in the input features/embeddings as opposed

Algorithm 2 Memory-Efficient RS-GNN.

Input: $\mathcal{G} = (\mathcal{V}, \mathbf{A}, \mathbf{X})$, k

- 1: Initialize \mathbf{R} , Θ , and \mathbf{U}
- 2: $\mathbf{F} = \text{GC}(\mathcal{G})$, $\mathbf{F}' = \square$
- 3: **for** $i=1$ **to** $n\text{Corrupt}$ **do**
- 4: $\mathcal{G}' = (\mathcal{V}, \mathbf{A}, \text{shuffle}(\mathbf{X}))$
- 5: $\mathbf{F}' = \text{concat}(\mathbf{F}', \text{GC}(\mathcal{G}'))$
- 6: **for** $\text{epoch}=1$ **to** $\#\text{epochs}$ **do**
- 7: $\mathbf{F}'' = \text{subsample}(\mathbf{F}', \text{len}(\mathbf{F}'))$
- 8: $\nabla = \mathbf{0}$
- 9: **for** $\mathbf{F}^{(b)}, \mathbf{F}^{(b')} \in \text{batch}(\mathbf{F}, \mathbf{F}'')$ **do**
- 10: $\mathbf{H} = \mathbf{W}\mathbf{F}^{(b)}$, $\mathbf{H}' = \mathbf{W}\mathbf{F}^{(b')}$
- 11: Compute \mathcal{L}_{EMB} based on \mathbf{H} and \mathbf{H}'
- 12: $\mu = \frac{1}{n} \sum_i \mathbf{H}_i$, $\zeta = \|\mathbf{H} - \mu\|$
- 13: $\tilde{\mathbf{H}} = \text{CenterNorm}(\mathbf{H}) = (\mathbf{H} - \mu)/\zeta$
- 14: $\mathcal{L}_{\text{SEL}} = \sum_i \min_j \text{Dist}(\tilde{\mathbf{H}}_i, \mathbf{R}_j)$
- 15: $\mathcal{L} = \mathcal{L}_{\text{EMB}} + \lambda \mathcal{L}_{\text{SEL}}$
- 16: Compute gradients for \mathcal{L} and add to ∇
- 17: Update parameters based on ∇
- 18: Let $\hat{\mathbf{R}}$ and $\hat{\mathbf{H}}$ be the representative and normalized node embeddings with minimum \mathcal{L} during training.
- 19: **for** $j=1$ **to** k **do**
- 20: The j^{th} representative = $\text{argmin}_i \text{Dist}(\hat{\mathbf{H}}_i, \hat{\mathbf{R}}_j)$

⁷<https://github.com/scikit-learn-contrib/scikit-learn-extra>

Table 7: Dataset statistics.

Dataset	Nodes	Edges	Features	Classes
Cora	2,708	5,278	1433	7
Citeseer	3,312	4,536	3703	6
Pubmed	19,717	44,324	500	3
Amazon Photo	7,650	119,081	745	8
Amazon PC	13,752	245,861	767	10
Coauthor CS	18,333	81,894	6,805	15
Coauthor PHY	34,493	247,962	8,415	5
OGBN-Arxiv	169,343	1,157,799	128	40

986 to the full graph. We used different hyperparameters for the RBF kernel (in the case of MaxCover
987 with RBF similarities) and kNN and reported the values that resulted in the best overall accuracy
988 across models. For MinCUT and DMoN, we used the implementation from the DMoN paper. For
989 our model, we set λ in the main loss function to 0.001 for all datasets. Also, for the experiments
990 where a graph structure is not provided as input, to create a kNN graph we connect each node to its
991 closest 15 nodes for all the datasets. For the one-shot graph active learning models, unfortunately we
992 did not find the code to be able to test the models in our setting. Therefore, we For the graph active
993 learning baseline, unfortunately we did not find source codes to be able to test them in our setting.
994 Therefore, we re-implemented FeatProp [109] and included the results of our implementation in the
995 experiments. For SDCN, EGAE, and GCC, we used the public codes released by the authors to select
996 representatives.

997 For the classification GCN model, we used a two-layer GCN model with PReLU activations [51] and
998 with a hidden dimension of 32. We added a dropout layer after the first layer with a drop rate of 0.5.
999 The weight decay was set to $5e^{-4}$. The GCN is trained based on the nodes in the selected set \mathcal{S} of
1000 representatives. We randomly split the remaining nodes in $(\mathcal{V} - \mathcal{S})$ into validation and test sets by
1001 selecting 500 nodes for validation and the rest for testing.

1002 We ran all the experiments 20 times (except for Arxiv where we ran it 10 times) with different random
1003 seeds and reported the mean and standard deviation of the runs. Our code will be released upon the
1004 acceptance of the paper.

1005 F Datasets

1006 We used eight established benchmarks in the GNN literature. A summary of our dataset statistics
1007 are provided in Table 7. The first three datasets are Cora, Citeseer, and Pubmed [94, 55]. These
1008 datasets are citation networks in which nodes represent papers, edges represent citations, features
1009 are bag-of-word abstracts, and the labels represent paper topics. The next two datasets are Amazon
1010 Photo and Amazon PC [97]. These two datasets correspond to photo and computers subgraphs of the
1011 Amazon copurchase graph. In these graphs, the nodes represent goods with an edge between two
1012 nodes representing that they have been frequently purchased together. Node features are bag-of-word
1013 reviews and class labels are product categories. The next two datasets are Coauthor CS and Coauthor
1014 Physics [97]. These are co-authorship networks for the computer science and physics fields based on
1015 the Microsoft Academic Graph respectively. The nodes in these two datasets represent authors, edges
1016 represent co-authorship, node features are a collection of paper keywords from author’s papers, and
1017 he class labels are the most common fields of study. Our last dataset is OGBN-Arxiv [55] which is
1018 also a citation dataset similar to Cora, Citeseer, and Pubmed, but orders of magnitude bigger than the
1019 three. The features in this dataset are average word embeddings of the paper abstracts.

1020 G Limitations

1021 We identify the following limitations with our current work:

- 1022 • Both our model and baselines optimize for micro-average classification accuracy; optimizing for
1023 macro-average classification accuracy may require extra terms in the loss function or architectural
1024 modifications.

- 1025 • While optimizing for micro-average accuracy is common in various domains, it raises the risk of
1026 being unfair to smaller sub-populations by not selecting any representatives from them. One must
1027 be cautious when using our model or any other model that optimizes for micro-average accuracy
1028 in applications when such a fairness is important.