

Finding the Best Not the Most: Regularized Loss Minimization Subgraph Selection for Graph Classification

Shirui Pan^{a,*}, Jia Wu^a, Xingquan Zhu^b, Guodong Long^a, Chengqi Zhang^a

^aCentre for Quantum Computation & Intelligent Systems, FEIT, University of Technology Sydney, Australia

^bDepartment of Computer and Electrical Engineering & Computer Science, Florida Atlantic University, USA

Abstract

Classification on structure data, such as graphs, has drawn wide interest in recent years. Due to the lack of explicit features to represent graphs for training classification models, extensive studies have been focused on extracting the most discriminative subgraphs features from the training graph dataset to transfer graphs into vector data. However, such filter-based methods suffer from two major disadvantages: (1) the subgraph feature selection is separated from the model learning process, so the selected most discriminative subgraphs may not best fit the subsequent learning model, resulting in deteriorated classification results; (2) all these methods rely on users to specify the number of subgraph features K , and suboptimally specified K values often result in significantly reduced classification accuracy.

In this paper, we propose a new graph classification paradigm which overcomes the above disadvantages by formulating subgraph feature selection as learning a K -dimensional feature space from an *implicit* and *large* subgraph space, with the optimal K value being automatically determined. To achieve the goal, we propose a regularized loss minimization-driven (RLMD) feature selection method for graph classification. RLMD integrates subgraph selection and model learning into a unified framework to find discriminative subgraphs with guaranteed minimum loss *w.r.t.* the objective func-

*Corresponding author. Address: Centre for Quantum Computation & Intelligent Systems, FEIT, University of Technology, Sydney, Ultimo, NSW 2007, Australia. Tel: +61 450768511. Fax: +61 2 9514 4535.

Email addresses: shirui.pan@uts.edu.au (Shirui Pan), jia.wu@student.uts.edu.au (Jia Wu), xzhu3@fau.edu (Xingquan Zhu), guodong.long@uts.edu.au (Guodong Long), chengqi.zhang@uts.edu.au (Chengqi Zhang)

tion. To automatically determine the optimal number of subgraphs K from the exponentially large subgraph space, an effective *elastic net* and a subgradient method are proposed to derive the stopping criterion, so that K can be automatically obtained once RLMD converges. The proposed RLMD method enjoys gratifying property including proved convergence and applicability to various loss functions. Experimental results on real-life graph datasets demonstrate significant performance gain.

Keywords: Feature Selection, Classification, Graph Classification, Sparse Learning

1. Introduction

Recent years have witnessed an increasing number of applications involving objects with structural relationships, including chemical compounds in Bioinformatics [1], brain networks [2], image structures [3], and academic citation networks [4]. For these applications, graph is a natural and powerful tool for modeling and capturing dependency relationships between objects.

Unlike conventional data, where each instance is represented in a feature-value vector format, graphs exhibit node-edge structural relationships and have no natural vector representation¹. As a result, a common practice is to transfer graphs into vectors [5, 6, 7, 8, 9] in structure space or in Euclidean space, so that traditional machine learning algorithms such as Support Vector Machines (SVM) and Decision Tree can be applied. In the structure space (also referred to as quotient space) [7, 8], the distance relations and nature of the original data are preserved, and some geometrical and analytical concepts such as derivatives of functions on structures can be determined, so that it can be applied to solve problems in structural pattern recognition. In the Euclidean space, the structural relations may be lost, but it provides simpler and more powerful analytical techniques for data analysis. Therefore, numerous approaches [10, 9, 11, 12, 13, 14, 15, 16, 17, 18] have been proposed to represent graphs in Euclidean space. The key idea of transferring graphs into vectors in Euclidean space is to extract a set of subgraphs as features and use the presence/absence of the features

¹In this paper, we only consider graphs with labels but no other feature values on nodes and edges.

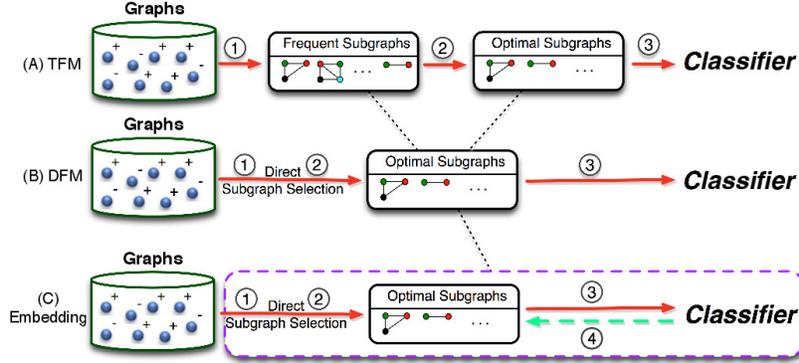


Figure 1: Subgraph-based methods for graph classification from the feature selection perspective. TFM methods (A) sequentially perform frequent subgraph mining ①, optimal feature selection ②, and classifier learning process ③. DFM methods (B) integrate the feature selection ② into the frequent subgraph mining ① process. Our embedding method RLMD (C) unifies all steps (①②③) into a whole framework, and iterates until convergence ④.

21 to represent each graph. From a feature selection perspective [19], these subgraph-
 22 based algorithms follow a filter approach for graph classification, *i.e.*, the subgraph
 23 feature selection and the subsequent model training are separated into two steps. In
 24 summary, existing filter-based graph classification methods roughly fall into the fol-
 25 lowing two categories:

26 **Two-step Filter Methods (TFMs):** This type of method first mines a set of frequent
 27 subgraphs as features and then applies a feature selection procedure to the discovered
 28 subgraphs, and uses the selected subgraph features to learn a classifier (*e.g.*, an SVM or
 29 Naive Bayes), as shown in Fig. 1 (A). An early study [9] has shown that using frequent
 30 subgraphs as features can achieve reasonable good classification results. However, be-
 31 cause TFMs separate subgraph feature discovery and feature evaluation into two steps,
 32 they may suffer from severe disadvantage in that the number of discovered subgraphs
 33 will grow exponentially when the minimum support value for subgraph mining is low.
 34 As a result, it will make the feature selection step heavily time-consuming. On the other
 35 hand, for relatively high minimum support values, many good subgraphs are pruned out
 36 because they do not meet the frequency requirement, so cannot be found to represent
 37 graphs.

38 **Direct Filter Methods (DFMs):** To improve the subgraph feature selection efficiency,
 39 numerous approaches [11, 12, 15, 16, 17] have been proposed to combine subgraph
 40 mining and feature selection into one step, representing a *direct discriminative feature*
 41 *selection* [18] scheme. So the feature selection is integrated into a subgraph mining
 42 process (Fig. 1 (B)), with pruning rules derived from the anti-monotone property of
 43 the significance (p -value) of each graph being used to reduce the search space. While
 44 DFMs substantially overcome the subgraph feature selection bottleneck, they also have
 45 a number of major disadvantages: (1) The subgraph selection is separated from the
 46 model learning process, so the selected subgraphs features may not best fit the under-
 47 lying learning model, and (2) All these methods require users to specify the number of
 48 subgraph features K , whereas the optimal number of subgraphs K required for training
 49 a good classifier for graph classification is unknown and difficult to determine. Al-
 50 though subgraphs are selected using optimized measures, due to the redundancy inside
 51 the feature set, the accuracy of the classifiers, when varying the number of selected
 52 subgraph features K , is highly variable, as shown in Fig. 2. This is a common problem
 53 for all existing filter-based graph classification methods.

54 The above observations motivate the proposed research which **aims** to *integrate*
 55 subgraph mining, feature selection, and model training into one single framework (Fig.
 56 1 (C)) with the optimal number of subgraphs K being automatically determined for
 57 graph classification. To achieve this goal, we formulate subgraph feature selection as
 58 the problem of learning a K -dimensional feature space from a **huge subgraph space**
 59 in order to result in minimum regularized loss on the training data as follows:

$$\min_{\mathbf{w}} \frac{1}{n} \sum_{i=1}^n \mathcal{L}(y_i, f(\mathbf{x}_i)) + \gamma R(\mathbf{w}) \quad (1)$$

60 where $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ are the vector representations of the training graphs, \mathcal{L} is a loss
 61 function measuring the difference between the prediction $f(\mathbf{x}_i)$ and the true label y_i ,
 62 and $R(\mathbf{w})$ is a regularization term on parameters \mathbf{w} to avoid over-fitting.

63 Indeed, the optimization in Eq.(1) has been widely studied [20, 21, 22] in machine
 64 learning community, but mainly for data with vector format. Several significant chal-
 65 lenges remain for graph data:

66 1. **Implicit Subgraph Features:** For graph classification, no subgraph features are

67 readily available (*i.e.*, x_i is unknown) for training the model in Eq.(1). Instead,
68 the feature space used to represent graphs is implicit and needs to be discovered
69 by subgraph mining procedure as needed.

70 **2. K -dimensional Features from Huge Subgraph Space:** The number of sub-
71 graph candidates representing graphs is exponentially large. Finding an optimal
72 number of K subgraphs for different graph datasets (in order to result in best
73 classifiers), is crucial but has not been addressed by existing research.

74 In this paper, we propose a *unified* regularized loss minimization-driven (RLMD) graph
75 classification framework. Our theme is to progressively select the most discriminative
76 subgraph features from the training data in order to achieve minimum regularized loss
77 for a well defined objective function. To integrate subgraph selection into the model
78 learning process (**Challenge 1**), we formulate an objective function and design a sub-
79 gradient method to induce a measurement to assess the utility of each subgraph, so that
80 the best subgraph features can be identified and incrementally included to optimize
81 the objective function for maximum performance gain. To determine the optimal K
82 value for each dataset (**Challenge 2**), we use an *elastic net* [21] and derive a stopping
83 condition, so that the K value can be automatically obtained when the algorithm con-
84 verges. By using the automatically determined optimal K value, as shown in Fig. 2,
85 RLMD finds 180 best subgraphs and achieves the best performance, which is 6% more
86 accurate than the second best method.

87 The main contributions of this paper are summarized as follows:

- 88 • We propose an *embedded* and *theoretically convergent* graph classification algo-
89 rithm, which can automatically determine the optimal number of subgraphs K
90 for graph classification. This is a *unified* approach in the sense that (1) it can em-
91 ploy any differentiable loss function (including least squares, exponential, and
92 logistic loss functions) for graph classification; and (2) it integrates subgraph
93 mining, feature selection, and model learning into one single framework.
- 94 • We generalize the *column generation technique* of gBoost [23] for graph classi-
95 fication, and demonstrate that gBoost [23] is a special case of our loss minimiza-
96 tion algorithm.

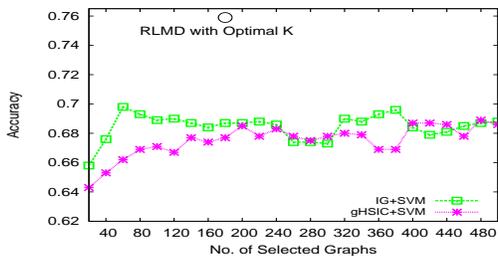


Figure 2: Classification accuracy for filter subgraph-based methods *w.r.t.* different numbers of subgraphs on the NCI-1 chemical compound dataset. IG is a TFM method which uses information gain to select subgraphs, whereas gHSIC [12] is a DFM method. All methods use SVM as a base classifier. The optimal number of subgraph features K is crucial, but difficult to decide for filter methods. In comparison, the proposed method (RLMD) automatically finds 180 best subgraphs and achieves the highest accuracy, which is 6% more accurate than the second best method.

- 97 • We propose the use of *elastic net* (which integrates two sparsity-inducing regu-
- 98 larization norms, ℓ_1 -norm and ℓ_2 -norm) to produce a sparse and robust solution
- 99 for discriminative subgraph selection.
- 100 • We derive a branch-and-bound rule according to the subgradient of our objective
- 101 function to prune search space for optimal subgraph mining.
- 102 • Experimental results show that our algorithm RLMD outperforms two-step filter
- 103 methods (TFMs), direct filter methods (DFMs), and gBoost algorithm with
- 104 significant performance gain.

105 2. Related Work

106 Our work is closely related to graph-based learning and graph classification.

107 2.1. Graph-based learning

108 Learning from data with dependency structures has been commonly addressed by

109 existing research. Instead of considering samples as *I.I.D* observations, graph-based

110 learning takes the relationships/correlations between samples to ensure effective learn-

111 ing. For example, graph-based approaches have been popularly used to propagate la-

112 bels in semi-supervised learning [24, 25, 26], where training samples are connected

113 through one or multiple graphs. A recent method [27] considers preserving global and

114 local structures inside the training data for feature selection. For large scale networks,
115 predicting linkage relationships between nodes (*i.e.* link prediction) can be used for
116 friendship recommendation in social networks [28], or suggesting potential interac-
117 tions between proteins in bioinformatics research. A recent work [29] proposed to use
118 latent feature kernels to support link prediction on sparse graphs. All above methods
119 consider a large scale network with thousands (or millions) of integer-connected nodes
120 in the network. In contrast, we consider small graph classification problem, in which
121 each graph has a label indicating the property of the graph, and the graph normally
122 contains tens or several hundreds of nodes. The purpose is to predict the label of the
123 graph by using node and structure information inside the graphs, for purposes such
124 as chemical compound activity prediction [1] and gender classification using magnetic
125 resonance connectome (*i.e.* brain-graph) [2].

126 2.2. Graph Classification

127 Existing methods for graph classification [18, 10, 9, 11, 12, 13, 14, 15, 16, 17,
128 23, 30, 31] can be roughly categorized into two groups: similarity-based methods and
129 vector representation-based methods.

130 2.2.1. Similarity-based methods

131 These approaches aim to directly learn global similarities between graphs by using
132 graph kernels [9, 32, 33, 34] or graph embedding [35]. Global similarities are then fed
133 to similarity-based classifiers, such as KNN or SVM, for learning. One clear drawback
134 of global similarity-based approaches is that the similarity is calculated based on global
135 graph structures, such as random walks or embedding space. Therefore, it is not clear
136 which substructures are more important for classifying graphs into different classes.

137 2.2.2. Vector representation-based methods

138 Another branch of methods transfer graphs into vector representations in structure
139 space or in Euclidean space. In structure space [7, 8], geometrical and analytical con-
140 cepts such as the angle between structures and the derivatives of functions on structures
141 can be obtained, so that the structural pattern recognition problems can be formulated
142 as optimization problems with certain cost functions. In Euclidean space, the goal is

143 to transfer graphs into vector representations in Euclidean space so existing analytical
144 techniques can be applied for data analysis. Methods in this category are mainly filter-
145 based approaches, including two-step filter methods (**TFMs**) or direct filter methods
146 (**DFMs**).

147 TFMs are straightforward approaches for graph classification which simply de-
148 compose frequent subgraph generation and selection as two separated steps. An early
149 work [9] has shown that learning an SVM classifier based on the discovered frequent
150 subgraphs can achieve reasonably good accuracy for graph classification. On the other
151 hand, research [16, 15] also indicates that TFM methods may result in a bottleneck
152 for the subsequent feature selection module. Specifically, the number of frequent sub-
153 graphs will grow exponentially if the minimum support threshold is low, which imposes
154 a great challenge for the subsequent feature selection task. This challenge has moti-
155 vated many direct filter methods (**DFMs**), which seek to integrate subgraph discovery
156 and feature selection into one step.

157 For DFMs (a review on this category can be found in [18]), a key issue is to
158 define a proper measurement to assess the utility of each subgraph. Yan *et al* [17] pro-
159 posed a LEAP algorithm to exploit the correlation between structural similarity and
160 significance similarity, so that a branch-and-bound rule can be derived to prune out
161 unpromising searching space efficiently. Ranu and Singh [16] proposed a scalable
162 GraphSig algorithm, which is able to mine significant subgraphs with low frequencies.
163 Thoma *et al.* [15] propose a CORK algorithm to find subgraph features. Recently, re-
164 searchers have extended DFM to other graph applications, and have proposed effective
165 algorithms such as gSemi [11] for the semi-supervised setting, gCGVFL [36] for multi-
166 view learning, gHSIC [12] for multi-label classification, and our recent multi-graph
167 classification for classifying graph bags, each containing multiple graphs [37, 38].

168 Although filter methods for graph classification have been extensively studied, they
169 all suffer from two major disadvantages: (1) the feature selection is not linked to the
170 model learning process. As a result, the selected subgraph features may not best fit
171 the underlying learning algorithms; and (2) the optimal number of subgraphs K for
172 graph classification is difficult to decide and often varies from dataset to dataset, and
173 inappropriately specified K value often results in significantly reduced classification

174 accuracy. This is the common drawback for filter-based methods [19].

175 *Embedded Methods.* Our algorithm belongs to the embedded approach which inte-
 176 grates the subgraph selection into the model training process. In this subcategory, Saigo
 177 *et al* [23] proposed a gBoost (its variants for imbalanced graph classification [39, 40]
 178 and cost-sensitive learning [41] are proposed recently) algorithm which formulates the
 179 graph classification as a linear program. As will be elaborated in Sec 4.6, our algorithm
 180 is more general in the sense that it can adopt any differentiable loss function and use
 181 more robust regularization to produce better performance. In fact, gBoost [23] can be
 182 considered as a special case of our loss minimization problem.

183 3. Problem Definition

184 **Definition 1. Connected Graph:** A graph is denoted by $G = (\mathcal{V}, E, L = \{L_1, L_2\}, \mathcal{A} =$
 185 $\{\mathcal{A}_1, \mathcal{A}_2\})$, where \mathcal{V} is the vertex set, $E \subseteq \mathcal{V} \times \mathcal{V}$ is the edge set, $\mathcal{A} = \{\mathcal{A}_1, \mathcal{A}_2\}$ with \mathcal{A}_1
 186 and \mathcal{A}_2 being the set of labels for vertices and edges, respectively; and $L = \{L_1, L_2\}$,
 187 $L_1 : \mathcal{V} \rightarrow \mathcal{A}_1$, $L_2 : E \rightarrow \mathcal{A}_2$ are labeling functions that assigns labels to a node or an
 188 edge, respectively. A connected graph is a graph such that there is a path between any
 189 pair of vertices.

190 In this paper, we focus on connected graphs and assume that each graph G has a class
 191 label y , $y \in \mathcal{Y} = \{-1, +1\}$. We only focus on binary-class classification tasks, but our
 192 methods can be easily extended to multi-class tasks.

193 **Definition 2. Subgraph:** Given two graphs $G = (\mathcal{V}, E, L = \{L_1, L_2\}, \mathcal{A} = \{\mathcal{A}_1, \mathcal{A}_2\})$
 194 and $g_k = (\mathcal{V}', E', L' = \{L'_1, L'_2\}, \mathcal{A}' = \{\mathcal{A}'_1, \mathcal{A}'_2\})$, g_k is a subgraph of G (i.e., $g_k \subseteq G$) if
 195 there is an injective function $\hat{f}: \mathcal{V}' \rightarrow \mathcal{V}$, such that $\forall (a, b) \in E'$, we have $(\hat{f}(a), \hat{f}(b)) \in$
 196 E , $L'_1(a) = L_1(\hat{f}(a))$, $L'_1(b) = L_1(\hat{f}(b))$, $L'_2(a, b) = L_2(\hat{f}(a), \hat{f}(b))$. If g_k is a subgraph of
 197 G ($g_k \subseteq G$), G is a supergraph of g_k ($G \supseteq g_k$).

198 **Subgraph-based Graph Classification:** Given a set of labeled graphs $\mathcal{T} = \{(G_1, y_1),$
 199 $\dots, (G_n, y_n)\}$, subgraph-based graph classification **aims** to select an optimal set of dis-
 200 criminative subgraphs \mathcal{F}_1 from \mathcal{T} , and learn a classification model from the reduced
 201 subgraph space \mathcal{F}_1 to predict previously unseen test graphs with a maximum accuracy.
 202 Set \mathcal{F}_1 is optimal if the classifier learned from \mathcal{F}_1 has the highest classification accuracy,

203 compared to classifiers trained from any subset of \mathcal{T} . A major feature of our method is
 204 that it can automatically determine the best set of subgraph features to represent each
 205 graph datasets without requiring users to specify the number of subgraph features. This
 206 essentially advances the existing subgraph feature-based graph classification methods
 207 from finding the most discriminative subgraph features to finding the best subgraph set
 208 for maximum accuracy gain.

209 4. Regularized Loss Minimization for Graph Classification

210 To support graph classification, state-of-the-art algorithms [23, 10] use a set of
 211 subgraphs discovered from the training graphs as features, where each subgraph g_k can
 212 be used to represent a graph G_i as follows:

$$\tilde{h}_{g_k}(G_i) = 2I(g_k \subseteq G_i) - 1; \quad (2)$$

213 Here $I(a) = 1$ if a holds, and 0 otherwise. This rule simply maps a graph G_i into +1 if
 214 $g_k \subseteq G_i$, or -1 otherwise.

215 Let $\mathcal{F} = \{g_1, \dots, g_m\}$ be the full set of subgraphs for the training graphs. We can use
 216 \mathcal{F} as features to represent each graph G_i into a vector space as $\mathbf{x}_i = \{\tilde{h}_{g_1}(G_i), \dots, \tilde{h}_{g_m}(G_i)\}$,
 217 with $\mathbf{x}_i^k = \tilde{h}_{g_k}(G_i)$. In the following subsection, G_i and \mathbf{x}_i are used interchangeably as
 218 they both refer to the same graph. Given the full subgraph features \mathcal{F} , the prediction
 219 function for the graph \mathbf{x}_i is a linear classifier:

$$f(\mathbf{x}_i) = \mathbf{x}_i \cdot \mathbf{w} + b = \sum_{g_k \in \mathcal{F}} w_k \tilde{h}_{g_k}(G_i) + b \quad (3)$$

220 where $\mathbf{w} = [w_1, \dots, w_m]'$ is the weight vector for all features \mathcal{F} , and b is the bias of
 221 the model. The predicted class of \mathbf{x}_i is +1 if $f(\mathbf{x}_i) > 0$ or -1 otherwise. Note that
 222 in practice, subgraph space \mathcal{F} is **implicit** and exponentially **large**, *i.e.*, the number of
 223 subgraphs grows exponentially with respect to the number of nodes.

224 4.1. Regularized Loss Minimization Formulation

225 In this paper, we propose to learn a K -dimensional feature space from the *implicit*
 226 and *large* subgraph space \mathcal{F} to achieve the lowest regularized empirical risks for the

227 graph dataset, with K being automatically determined. Eq.(1) can be reformulated as
 228 the following objective function:

$$\min_{\mathbf{w}, b} \mathcal{J}(\mathbf{w}, b) = \underbrace{\min_{\mathbf{w}, b} \frac{1}{n} \sum_{i=1}^n \mathcal{L}(y_i, \mathbf{x}_i \cdot \mathbf{w} + b)}_C + \underbrace{\gamma_1 \|\mathbf{w}\|_1 + \gamma_2 \|\mathbf{w}\|_2^2}_R \quad (4)$$

229 The first term C measures the loss on the training graphs, where $\mathcal{L}(y_i, f(\mathbf{x}_i))$ can be any
 230 loss function measuring the misclassification penalty of a graph G_i . The second part R
 231 consists of regularization terms to enforce sparse and robust solutions. Parameters γ_1
 232 and γ_2 are used to trade-off these parts ($\gamma_1 \geq 0, \gamma_2 \geq 0$). For the regularization, our
 233 objective is to obtain a sparse and stable solution on \mathbf{w} , *i.e.*, low dimensional subgraph
 234 features for final graph classification. Here, we combine both ℓ_1 and ℓ_2 norm, which is
 235 known as *elastic net* in machine learning [21]. The motivation of our regularization is
 236 as follows:

237 The ℓ_1 -norm regularizer ($\sum_k |w_k|$) can produce solutions with many coefficients be-
 238 ing 0, which is known as lasso [20] and has been widely applied for variable selections.

239 Although ℓ_1 regularization can produce a sparse solution, it suffers from two major
 240 disadvantages: (1) the number of selected variables is limited by the number of obser-
 241 vations; and (2) the lasso penalized model can only select one variable from a group
 242 of correlated variables and does not care which one is selected [21]. In contrast, ℓ_2
 243 regularization, which is widely used in SVM formulation ($(\|\mathbf{w}\|_2^2 = \sum_{k=1}^m |w_k|^2)$), can
 244 produce more stable and robust classification results. However, ℓ_2 formulation cannot
 245 produce a sparse solution. By combining ℓ_1 and ℓ_2 norm, known as elastic net [21], we
 246 can overcome these issues and enjoy the sparse and stable properties.

247 4.2. Sparse Subgraph Learning: Challenges and Solution Overview

248 **Challenges:** For explicit vector data with moderate feature size, the problem defined in
 249 Eq.(4) can be effectively solved in traditional supervised learning. However, for graph
 250 data the challenges are evident: (1) the feature set \mathcal{F} is unavailable (**implicit**) unless
 251 we enumerate all subgraphs from the training graphs, which is NP-complete; and (2)
 252 the whole subgraph set is exponentially **large**, and only a small subset of subgraphs are
 253 useful for classifiers to achieve maximum graph classification accuracy.

254 **Solution Overview:** To solve the aforementioned challenges, we propose a regularized
 255 loss minimization-driven (RLMD) subgraph selection method for graph classification.
 256 Driven by our formulation in Eq.(4), our principle is to iteratively mine the best sub-
 257 graph feature to reduce the empirical loss on the training graphs. To this end, we resort
 258 to the subgradient method in the functional space to define the utility of each sub-
 259 graph, and embed the feature selection/ranking into the subgraph mining/enumeration
 260 process. To handle the exponentially large subgraph space, we derive an effective
 261 branch-and-bound pruning scheme to reduce the search space. After a new subgraph
 262 is selected, we include and re-solve the new restricted objective function of Eq.(4) by
 263 using currently selected subgraphs. To find optimal K value, we derive a stopping cri-
 264 terion for our feature selection procedure based on the subgradient in the functional
 265 space, so that K can be automatically obtained once the algorithm converges.

266 **Logistic Loss Function: A Running Example.** Our method is based on the gradi-
 267 ent/subgradient on the functional space of the objective function Eq.(4). In this paper,
 268 we use the logistic loss function as an example to illustrate how subgraph selection is
 269 performed by subgradient methods, and the logistic loss function is given as follows:

$$\mathcal{L}(y_i, f(x_i)) = \log \left(1 + \exp \{ -y_i f(x_i) \} \right) \quad (5)$$

270 Note that our algorithm is a general method in the sense that any other differentiable
 271 loss function, such as least square loss $\mathcal{L}(y, f_i) = \frac{1}{2}(y - f_i)^2$ or exponential loss $\mathcal{L}(y, f_i) =$
 272 $\exp\{-y f_i\}$, can be directly used in our algorithm. As discussed latter in Section 4.6,
 273 our method is also applicable to convex locally Lipschitz but non-differentiable loss
 274 functions such as the hinge loss used by the (margin) perceptron and linear SVM.

275 The partial derivative of the loss term C in Eq.(4) on the subgraph feature g_k is
 276 defined as $\frac{\partial C}{\partial w_k}$. For logistic loss function,

$$\begin{aligned} \frac{\partial C}{\partial w_k} &= \frac{1}{n} \sum_{i=1}^n \frac{\partial \mathcal{L}(y_i, f(x_i))}{\partial f(x_i)} \frac{\partial f(x_i)}{\partial w_k} \\ &= -\frac{1}{n} \sum_{i=1}^n \frac{y_i \mathbf{x}_i^k}{1 + e^{y_i f(x_i)}} = \sum_{i=1}^n y_i \alpha_i \mathbf{x}_i^k \end{aligned} \quad (6)$$

277 Here, $\alpha_i = -\frac{1}{n(1 + e^{y_i f(x_i)})}$ can be regarded as a weight associated with graph G_i for the
 278 subgraph mining process.

279 *4.3. RLMD Subgraph Selection for Graph Classification*

280 Because we aim to learn a sparse solution of subgraph features (K -dimensional
 281 feature space) from graph data, some subgraphs/features g_k with zero weights, *i.e.*,
 282 $w_k = 0$ will not be used for learning the classification model. Thus it makes sense
 283 to partition the subgraph features \mathcal{F} into two disjoint subsets \mathcal{F}_1 and \mathcal{F}_2 . \mathcal{F}_1 stores
 284 active features which are used to learn the classification model and this set is frequently
 285 updated as desired, and \mathcal{F}_2 includes unselected graphs with 0 weights (*i.e.*, for $g_k \in$
 286 $\mathcal{F}_2, w_k = 0$). Then we can iteratively select the best feature from \mathcal{F}_2 to \mathcal{F}_1 , and solve
 287 the following restricted subproblem:

$$\begin{aligned} \min_{\mathbf{w}, b} \mathcal{J}_t(\mathbf{w}, b) &= \min_{\mathbf{w}, b} \frac{1}{n} \sum_{i=1}^n \underbrace{\mathcal{L}(y_i, \mathbf{x}_i(t) \cdot \mathbf{w} + b)}_{C'} + \underbrace{\gamma_1 \|\mathbf{w}\|_1 + \gamma_2 \|\mathbf{w}\|_2^2}_{R'} \\ &= \mathcal{J}_t(\mathbf{w}(t), b(t)) \end{aligned} \quad (7)$$

288 where $\mathbf{x}_i(t)$ is the feature representations for graph G_i based on the active set \mathcal{F}_1 in the
 289 t -th iteration, and $(\mathbf{w}(t), b(t))$ is the solution of Eq.(7). Note that $\mathcal{J}_t(\mathbf{w}, b)$ is used here to
 290 denote the restricted subproblem Eq.(7) while $\mathcal{J}(\mathbf{w}, b)$ is referred to original problem
 291 in Eq.(4).

292 The optimal number of subgraphs K can be automatically determined by setting
 293 $K = |\mathcal{F}_1|$ once the algorithm converges. Note that given a solution $(\mathbf{w}(t), b(t))$, the
 294 loss term C' in Eq.(7) equals to C in Eq.(4) because the prediction of of graph mainly
 295 depends on the active set \mathcal{F}_1 , *i.e.*, $C' = C$. In the following, we will derive the stopping
 296 condition of our algorithm, and prove its convergence.

297 **Stopping Condition for Optimal K value:** Our objective function Eq.(4) is con-
 298 vex and non-smooth, *i.e.*, it may be non-differentiable at a point \mathbf{w} . When it is non-
 299 differentiable at \mathbf{w} , we can compute its generalized gradient (*i.e.*, subgradient) instead.
 300 According to the optimization conditions, when reaching the optimum, we will have

$$0 \in \frac{\partial C}{\partial w_k} + \gamma_1 o_k + 2\gamma_2 w_k; \quad (8)$$

301 where o_k is the subgradient with respect to w_k

$$o_k \in \begin{cases} \text{sign}(w_k) & : w_k \neq 0 \\ [-1, 1] & : w_k = 0 \end{cases} \quad (9)$$

302 where $\text{sign}(a) = 1$ if $a > 0$ otherwise -1 .

303 According to Eq.(8) and Eq.(9), we can now state the optimal condition for our
304 sparse subgraph learning problem.

305 **Proposition 1. Optimal Solution:** Let $\hat{\mathbf{w}} = [\hat{w}_1, \dots, \hat{w}_m]$. Suppose that $(\hat{\mathbf{w}}, \hat{\mathbf{b}})$ is the
306 optimal solution of our objective function Eq.(4), then Eq.(10) and Eq.(11) hold.

$$\frac{\partial \mathcal{C}}{\partial \hat{w}_k} + \gamma_1 \text{sign}(\hat{w}_k) + 2\gamma_2 \hat{w}_k = 0 \quad \text{if } \hat{w}_k \neq 0 \quad (10)$$

$$\left| \frac{\partial \mathcal{C}}{\partial \hat{w}_k} \right| \leq \gamma_1 \quad \text{if } \hat{w}_k = 0 \quad (11)$$

307 Eq.(11) holds because for $\hat{w}_k = 0$, the third term of Eq.(8) disappears. Combining
308 Eq.(8) and Eq.(9) will result in Eq.(11).

309 To reduce the objective value \mathcal{J}_t in Eq.(7), we propose to select a subgraph in \mathcal{F}_2
310 whose weight violates Eq.(11), and update the selected active set \mathcal{F}_1 with the newly
311 selected feature and re-optimize the restricted subproblem Eq.(7) with current features.
312 This process will repeat until no candidate violates Eq.(11). In other words, Eq.(11)
313 is a stopping condition and determines the number of subgraphs being selected for
314 RLMD's subgraph selection process.

315 **Utility of Subgraphs:** Eq.(11) can be used naturally to induce a criterion for quantify-
316 ing the utility value of a subgraph. The larger $|\frac{\partial \mathcal{C}}{\partial w_k}|$ is, the more informative it will be
317 for reducing the objective function. Accordingly, we formally define the informative
318 score as follows:

319 **Definition 3. Informative Score:** For a subgraph pattern g_k , its informative score for
320 graph classification is defined as follows:

$$\Theta(g_k) = \left| \frac{\partial \mathcal{C}}{\partial w_k} \right| = \left| \frac{\partial \mathcal{C}'}{\partial w_k} \right| = \left| \sum_{i=1}^n y_i \alpha_i \mathbf{x}_i^k \right| \quad (12)$$

321 where $\alpha_i = -\frac{1}{n(1+e^{y_i f(\mathbf{x}_i)})}$.

322 Note that the informative score directly depends on the weight of each graph α_i , which
323 is calculated based on the active set \mathcal{F}_1 . Intuitively, the best subgraph of \mathcal{F}_2 is the one
324 with the maximum informative score, because it is more likely to violate the stopping
325 condition Eq.(11).

Algorithm 1 Regularized Loss Minimization-Driven Subgraph Selection (RLMD) for Graph Classification

Require:

$\{(G_1, y_1), \dots, (G_n, y_n)\}$: Training Graphs;
 S_{max} : Maximum number of iterations;

Ensure:

w, b : Parameters for classifier model

```

1:  $\alpha_i = 1/n$ ;  $\mathcal{F}_1 \leftarrow \emptyset$ ;  $t \leftarrow 0$ ;
2: while  $t < S_{max}$  do
3:   Mine an optimal subgraph features  $g^*$  with maximum informative score defined by Eq.12; //Algorithm 2;
4:   if  $\Theta(g^*) \leq \gamma_1 + \varepsilon$  then
5:     break;
6:   end if
7:    $\mathcal{F}_1 \leftarrow \mathcal{F}_1 \cup g^*$ ;
8:   Solve Eq.(7) based on  $\mathcal{F}_1$  to get the new solution  $(w(t), b(t))$ ;
9:   Update the graph weights on each training graph
                                     
$$\alpha_i = -\frac{1}{n(1+e^{\beta_i f(x_i(t))})}$$

10:   $t \leftarrow t + 1$ ;
11: end while
12:  $K = |\mathcal{F}_1|$ ;
13: return  $w, b$ ;
```

326 **RLMD Algorithm:** Algorithm 1 illustrates the detailed steps of RLMD for graph
327 classification. Initially, the weights for all training graphs are equally set as $1/n$, and
328 the active set \mathcal{F}_1 is initialized to be empty.

329 In the next step, the algorithm mines an optimal subgraph g^* from \mathcal{F}_2 which has
330 the highest informative scores defined by Eq.(12). This step involves a subgraph min-
331 ing procedure, which will be addressed in the next subsection. On steps 4-5, if current
332 optimal subgraph no longer violates the optimal condition Eq.(11), the algorithm termi-
333 nates. Here, we have relaxed the convergence condition to ε tolerance; this is because
334 in the last few iterations, the maximum score will only change subtly (we set $\varepsilon = 0.005$
335 in our experiments).

336 On step 7, we add the newly selected subgraph g^* to the existing subgraph set
337 \mathcal{F}_1 , and re-solve the following restricted subproblem Eq.(7). To solve this restricted

338 objective function, we use the MALSAR toolbox ² in our experiments. It is worth
 339 noting that because this step only involves a very small number of features, it is very
 340 efficient in practice.

341 Subsequently, the algorithm updates the weight α_i for each graph G_i . This will help
 342 compute the derivative of $\frac{\partial C}{\partial w_k}$ for subgraph mining in next round. After the algorithm
 343 terminates, the optimal number of subgraphs K can be easily obtained as $K = |\mathcal{F}_1| = t$
 344 on step 12.

345 Note that our algorithm 1 generalizes the column generation technique in gBoost [23]
 346 by iteratively selecting the most violated subgraph in each iteration until convergence.
 347 Our algorithm 1 relies on ε and γ_1 , which serve as a stopping condition to determine K .
 348 In practice, ε is a subtle value insensitive to the algorithm performance. Meanwhile,
 349 γ_1 is much easier to set than asking users to specify K values because γ_1 is chosen in a
 350 much smaller range, as we will demonstrate in Section 5.2.4.

351 4.4. Theoretical Study

352 **Theorem 1. (Convergence Property:)** *Algorithm 1 guarantees that the restricted ob-*
 353 *jective function Eq.(7) will monotonically decrease.*

PROOF. *Suppose in the t -th iteration, the optimal objective value based on current t*
features (i.e., $|\mathcal{F}_1| = t$) is obtained with $(\mathbf{w}(t), b(t))$, i.e.,

$$\mathcal{J}_t(\mathbf{w}(t), b(t)) = \underbrace{\frac{1}{n} \sum_{i=1}^n \mathcal{L}(y_i, \mathbf{x}_i(t) \cdot \mathbf{w}(t) + b(t))}_{C=C'} + \underbrace{\gamma_1 \|\mathbf{w}(t)\|_1 + \gamma_2 \|\mathbf{w}(t)\|_2^2}_{\mathcal{R}=\mathcal{R}} = (C+\mathcal{R})_{(\mathbf{w}(t), b(t))}$$

then in the $t + 1$ -th iterations, the optimal objective value for Eq.(7) is

$$\min \mathcal{J}_{t+1}(\mathbf{w}, b) = \min(C + \mathcal{R}) \leq (C + \mathcal{R})_{(\mathbf{w}(t), 0), b(t)} = \mathcal{J}_t(\mathbf{w}(t), b(t))$$

354 *Here $[\mathbf{w}(t), 0]$ means that the weights for subgraphs selected in the t -th iteration remain*
 355 *unchanged while the weight for newly selected subgraph in the $t+1$ -th iteration is 0.*

²<http://www.MALSAR.org>

356 Thus the objective value of the restricted problem Eq.(7) based on the currently
 357 selected features \mathcal{F}_1 always monotonously decreases in two successive iterations. Be-
 358 cause the objective function value is non-negative (bounded), we can ensure that it will
 359 finally converge as iteration continues. The proof is complete.

360 Suppose the algorithm converges in the K -th iteration with a solution $(\mathbf{w}(K), b(K))$,
 361 and the objective value for Eq.(7) is $\mathcal{J}_K(\mathbf{w}(K), b(K))$. By adding $m - K$ zeros for
 362 subgraphs in \mathcal{F}_2 to $\mathbf{w}(K)$, i.e., $\hat{\mathbf{w}}(K) = [\mathbf{w}(K), 0 \cdots]$, we obtain a solution $(\hat{\mathbf{w}}(K), b(K))$
 363 for Eq.(4).

364 **Corollary 1. Optimal Solution Guarantees:** *If Algorithm 1 converges with solution*
 365 *$(\mathbf{w}(K), b(K))$ for Eq.(7), then $(\hat{\mathbf{w}}(K), b(K))$ is an optimal solution for Eq.(4).*

PROOF. According to our Proposition 1, $(\hat{\mathbf{w}}(K), b(K))$ is an optimal solution of Eq.(4),
 because $\forall w_k = 0 (g_k \in \mathcal{F}_2)$, we have $\Theta(g_k) < \gamma_1$ based on our stopping condition. Thus
 we will have

$$\mathcal{J}_K(\mathbf{w}(K), b(K)) = \mathcal{J}(\hat{\mathbf{w}}(K), b(K)) = \min \mathcal{J}(\mathbf{w}, b)$$

366 where $\mathcal{J}_K(\mathbf{w}(K), b(K))$ and $\mathcal{J}(\hat{\mathbf{w}}(K), b(K))$ refer to the objective values of the restricted
 367 subproblem Eq.(7) and Eq.(4), respectively.

368 We have proved that objective value for Eq.(7) is monotonously decreasing (Theo-
 369 rem 1) and its recovered solution $(\hat{\mathbf{w}}(K), b(K))$ is an optimal solution to Eq.(4) (Corol-
 370 lary 1).

371 4.5. Optimal Subgraph Mining

372 In order to mine optimal subgraph g^* on step 3 of Algorithm 1, we need to perform
 373 the subgraph enumeration procedure. In RLMD, we employ the frequent subgraph
 374 mining-based algorithm gSpan [42]. The key idea of gSpan is that each subgraph has
 375 a unique DFS Code, which is defined by a lexicographic order of the discovery time
 376 during the search process. By employing a depth first search strategy on the DFS Code
 377 tree (where each node is a subgraph), gSpan can enumerate all frequent subgraphs
 378 efficiently.

379 During the subgraph mining process, the search space is exponentially large, which
 380 requires an effective pruning scheme to reduce the search space. In this subsection, we
 381 will derive the upper-bound of the informative score for each subgraph, which helps
 382 prune the search space and speed up the subgraph mining.

383 **Theorem 2. (Upper-bound Score:)** Let g and g' be two subgraph patterns, and $g \subseteq g'$,
 384 for the subgraph g , we define

$$\begin{aligned} A_1(g) &= 2 \sum_{\{i|y_i=+1, g \in G_i\}} \alpha_i \\ A_2(g) &= 2 \sum_{\{i|y_i=-1, g \in G_i\}} \alpha_i \\ A_3 &= \sum_{i=1}^n \alpha_i y_i \end{aligned}$$

385

$$\hat{\Theta}(g) = \begin{cases} \max\{|A_1(g) - A_3|, |A_2(g)|\} & : A_3 \geq 0 \\ \max\{|A_2(g) + A_3|, |A_1(g)|\} & : A_3 < 0 \end{cases}$$

386 then $\Theta(g') \leq \hat{\Theta}(g)$, where $\Theta(g')$ is defined in Eq.(12).

387 **PROOF.** We start with the definition of $\Theta(g')$:

$$\begin{aligned} \Theta(g') &= \left| \sum_{i=1}^n y_i \alpha_i \mathbf{x}'_i \right| \\ &= \left| \sum_{i=1}^n y_i \alpha_i \cdot [2I(g' \subseteq G_i) - 1] \right| \\ &= \left| 2 \sum_{g' \subseteq G_i} y_i \alpha_i - \sum_{i=1}^n \alpha_i y_i \right| \\ &= |A_1(g') - A_2(g') - A_3| \\ &\leq \begin{cases} \max\{|A_1(g') - A_3|, |A_2(g')|\} & : A_3 \geq 0 \\ \max\{|A_2(g') + A_3|, |A_1(g')|\} & : A_3 < 0 \end{cases} \\ &\leq \begin{cases} \max\{|A_1(g) - A_3|, |A_2(g)|\} & : A_3 \geq 0 \\ \max\{|A_2(g) + A_3|, |A_1(g)|\} & : A_3 < 0 \end{cases} \\ &= \hat{\Theta}(g) \end{aligned}$$

388 The first inequality holds because for $\alpha_i < 0$, $A_1(g') \leq 0$ and $A_2(g') \leq 0$, so the upper-
 389 bound depends on A_3 . If $A_3 \geq 0$, $A_1(g')$ and A_3 will have different signs, then the
 390 upper-bound is the maximum between $\{|A_1(g') - A_3|, |A_2(g')|\}$. The case is similar for
 391 $A_3 < 0$. The second inequality holds because $|A_1(g')| \leq |A_1(g)|$ and $|A_2(g')| \leq |A_2(g)|$
 392 for $g \subseteq g'$.

Algorithm 2 Optimal Subgraph Mining

Require:

$\{(G_1, y_1), \dots, (G_n, y_n)\}$: Training graphs;
 α_i : Weight for each graph example;
 \mathcal{F}_1 : Already selected subgraph set;

Ensure:

g^* : The optimal subgraph;
1: $\eta = 0$;
2: **while** Recursively visit the DFS Code Tree in gSpan **do**
3: $g_p \leftarrow$ current visited subgraph in DFS Code Tree;
4: **if** g_p has been examined **then**
5: **continue**;
6: **end if**
7: Compute scores $\Theta(g_p)$ for subgraph g_p according Eq.(12);
8: **if** $g_p \notin \mathcal{F}_1$ & $\Theta(g_p) > \eta$ **then**
9: $\eta = \Theta(g_p)$;
10: $g^* \leftarrow g_p$;
11: **end if**
12: **if** $\hat{\Theta}(g_p) > \eta$ **then**
13: Depth-first search the subtree rooted from node g_p ;
14: **end if**
15: **end while**
16: **return** g^* ;

393 Theorem 2 states that for any super graph of a subgraph g , its informative score is
394 upper-bounded by $\hat{\Theta}(g)$. This rule can prune unpromising candidates effectively.

395 **Optimal Subgraph Exploration Algorithm:** Our optimal subgraph mining algorithm
396 is listed in Algorithm 2. The minimum value η in the optimal set is initialized on step 1.
397 Duplicated subgraph features are pruned on steps 4-5, and the informative score $\Theta(g_p)$
398 for g_p is calculated on step 7. If g_p is not selected before ($g_p \notin \mathcal{F}_1$) and $\Theta(g_p)$ is larger
399 than η , we replace the optimal subgraph g^* with the current g_p and update the optimal
400 score η (steps 8-11).

401 A branch-and-bound pruning rule, according to Theorem 2, is subsequently used to
402 prune the search space on step 12. Lastly, the optimal subgraph g^* is returned on step
403 16.

404 The above pruning process is a key feature of our algorithm, because we do not
 405 require a support threshold value for subgraph mining (whereas all filter subgraph min-
 406 ing methods will require users to predefine a threshold value in order to discover sub-
 407 graphs).

408 4.6. Relation to gBoost

409 Our RLMD subgraph selection algorithm advances the existing *column generation*
 410 style techniques employed in gBoost [23] for graph classification. The learning objec-
 411 tive function for gBoost is

$$\begin{aligned}
 & \max_{\rho, \mathbf{w}, \xi} \quad \rho - \frac{1}{vn} \sum_{i=1}^n \xi_i \\
 & \text{s.t.} \quad y_i \sum_{k=1}^m \tilde{h}_{g_k}(G_i) w_k + \xi_i \geq \rho; \\
 & \quad \quad \sum_{k=1}^m w_k = 1; \\
 & \quad \quad w_k \geq 0, \xi_i \geq 0;
 \end{aligned} \tag{13}$$

412 From [43], we know that this formula is equivalent to the following linear program-
 413 ming:

$$\begin{aligned}
 & \min_{\mathbf{w}, \xi} \quad \sum_{k=1}^m w_k + C \sum_{i=1}^n \xi_i \\
 & \text{s.t.} \quad y_i \sum_{k=1}^m \tilde{h}_{g_k}(G_i) w_k + \xi_i \geq 1; \\
 & \quad \quad w_k \geq 0, \xi_i \geq 0;
 \end{aligned} \tag{14}$$

414 Eq.(14) is actually a ℓ_1 svm formulation, and can also be formulated as the regularized
 415 loss minimization formulation problem:

$$\min \|\mathbf{w}\|_1 + C \sum_{i=1}^n \mathcal{L}_h(y_i, f(\mathbf{x}_i)) \tag{15}$$

416 Here, $\mathcal{L}_h(y_i, f(\mathbf{x}_i)) = \max(1 - y_i f(\mathbf{x}_i), 0)$, which is known as hinge loss in machine
 417 learning.

418 Compared to our objective function in Eq.(4), we can find that gBoost (Eq.15) is
 419 a special case of Eq.(4), with the ℓ_2 regularization term being 0. Although the hinge
 420 loss function is non-differentiable, our subgradient method still applies, as long as $\frac{\partial C}{\partial w_k}$
 421 in Eq.(6) is properly defined. This observation shows the following advantages of our
 422 algorithm: (1) gBoost employs a hinge loss function which is similar to SVM and

423 requires the problem to be formulated as a linear programming. Our algorithm gen-
424 eralizes and advances gBoost by removing the linear programming constraint and can
425 employ any differentiable loss function, in addition to the logistic loss function con-
426 sidered in our paper. This generalization has great attractiveness in many applications,
427 especially when the probability estimation for classification is required (the logistic
428 function can provide some probabilistic information compared to the hinge loss func-
429 tion); (2) while gBoost employs ℓ_1 norm regularization to obtain a sparse solution, our
430 algorithm considers an additional norm ℓ_2 . This combined norm (known as *elastic net*)
431 enables a sparse and more stable solution.

432 5. Experiment

433 5.1. Experimental Settings

434 **Benchmark Data:** We validate the performance of the proposed algorithm on two
435 types of graph classification datasets.

436 *Anti-cancer activity prediction (NCI):* The NCI graph collection³ is a benchmark
437 for predicting the biological activity of small molecules for different types of cancers.
438 Each NCI dataset belongs to a bioassay task for anticancer activity prediction, such
439 as Breast cancer or Leukemia cancer. Each molecule is represented as a graph, with
440 atoms representing nodes and bonds denoting edges. A molecule is positive if it is
441 active against a certain type of cancer, or negative otherwise. Table 1 summarizes nine
442 NCI graph classification tasks used in our experiments, where columns 2-4 denote the
443 number of positive molecules, the total number of graphs, and the type of cancer of
444 each dataset. In our experiments, we randomly select 1000 graphs from each dataset
445 with balanced class distributions for graph classification.

446 *Predictive Toxicology Challenge Dataset (PTC):* The PTC challenge includes a
447 number of carcinogenicity classifications for the toxicology prediction of chemical
448 compounds⁴. The dataset we selected contains 417 compounds with four types of
449 test animals: MM (male mouse), FM (female mouse), MR (male rat), and FR (female

³<http://pubchem.ncbi.nlm.nih.gov>

⁴<http://www.predictive-toxicology.org/ptc/>

Table 1: Datasets Used in Experiments

ID	#Pos	#Total	Learning tasks
1	1793	37349	Non-Small Cell Lung
33	1467	37022	Melanoma
41	1350	25336	Prostate
47	1735	37298	Central Nerv Sys
81	2081	37549	Colon
83	1959	25550	Breast
109	1773	37518	Ovarian
123	2715	36903	Leukemia
145	1641	37043	Renal

450 rat). Each compound has labels selected from {CE, SE, P, E, EE, IS, NE, N}. Similar
 451 to [44], we set {CE, SE, P} as positive labels, and {NE,N} as negative labels.

452 **Baseline Methods:** In our experiments, we consider three types of baseline methods,
 453 namely the two-step filter methods (TFMs), direct filter methods (DFMs), and embed-
 454 ded methods, as follows:

- 455 • **IG+SVM** is a TFM method that simply mines a set of frequent subgraphs, and
 456 then performs feature selection by using Information Gain. A SVM classifier is
 457 trained by using selected subgraph features for graph classification.
- 458 • **TOP+SVM** is similar to IG+SVM except that it selects the top K subgraphs
 459 based on their frequency rather than their information gain values.
- 460 • **gSemi+SVM [11]** is a DFM method, which integrates the feature selection into
 461 the subgraph mining process. The measurement for feature selection mainly con-
 462 sideres the *must-link* and *cannot link* constraints between graph samples within
 463 the same or between different classes.
- 464 • **gHSIC+SVM [12]** is another DFM method which exploits the correlations be-
 465 tween features and labels.
- 466 • **gBoost [23]** is a state-of-the-art embedded method which formulates the feature
 467 selection as a linear problem and selects subgraph features which best fit the

468 objective function.

- 469 • **RLMD** is our proposed method which employ a logistic loss function together
470 with an *elastic net* for regularization, and automatically determines optimal number
471 of subgraphs K .

472 We conduct 10-fold cross-validation on all graph datasets and report the average
473 results and *standard errors* of 10 folds in the final result. The parameters for γ_1 are
474 selected from $\{0.005, 0.01, 0.03, 0.05\}$, and γ_2 is selected from $\{0.01, 0.03, 0.05\}$. We
475 will further analyze the impact of γ_1 and γ_2 in wider ranges in Section 5.2.4. For the
476 filter methods (IG+SVM, TOP+SVM, gSemi+SVM, and gHSIC+SVM), the minimum
477 support for frequent subgraph mining is set to 10% on NCI graph datasets and 1%
478 on PTC classification tasks, and an SVM classifier is trained with C parameter from
479 the range $\{0.1, 1, 10, 100, 1000, 10000\}$. For the gBoost algorithm, the parameter ν is
480 selected from $\{0.1, 0.2, 0.3, 0.4\}$. Following [23], we select the best average results
481 of 10-fold cross-validation for each baseline algorithm by varying these parameters,
482 which represents the best performance each baseline can achieve.

483 For fairness of comparison, we increase the number of features to be selected for
484 the filter methods (IG+SVM, TOP+SVM, gSemi+SVM, and gHSIC+SVM), and in-
485 crease the iterations for the embedded methods (gBoost and RLMD), then collect and
486 compare the performance of all algorithms under the same number of features. We set
487 $S_{max} = 200$, which defines the maximum number of features used to learn the classifier
488 models. Note that for RLMD, the algorithm may stop before reaching the maximum
489 iterations/subgraphs we set, *i.e.*, the optimal K is obtained. When RLMD stops, the
490 optimal number of subgraph features has been discovered and RLMD will not add
491 additional subgraphs to the feature set. We also compare RLMD under the optimal
492 subgraph value to other baselines with the same number of K features (the purpose is
493 to show that the optimal subgraph features discovered by RLMD are indeed optimal
494 for graph classification).

495 5.2. Experimental Results

496 5.2.1. Results on NCI Graph Dataset

497 For the NCI graph datasets, we vary the number of selected subgraph features from
498 20 to 200 for filter methods, and the number of iterations for gBoost and RLMD from
499 1 to 200. The accuracies and AUC values are shown in Fig. 3.

500 **Comparison with Filter methods:** The results in Fig. 3 show that with the in-
501 crease in the number of features/iterations, the filter methods (TOP+SVM, IG+SVM,
502 gSemi+SVM, and gHSIC+SVM) are inferior to RLMD. This is because filter methods
503 separate the feature selection module from the model learning process. The subgraph
504 features selected from filter methods may not fit the underlying learning model very
505 well (we use SVM in our experiments). This is actually an observed common draw-
506 back of filter methods [19]. The performance among these filter-based methods varies
507 from one graph dataset to another, and none of them significantly outperforms others.
508 For instance, gSemi+SVM outperforms TOP+SVM, IG+SVM, and gHSIC+SVM on
509 NCI-1 (Fig. 3.A) when the number of selected graphs is considerably large (≥ 160),
510 but is worse than gHSIC on NCI-109 (Fig. 3.G). This may be attributed to the inherent
511 differences underneath the graph datasets.

512 **How many subgraphs to select:** Another drawback of filter methods, shown in our
513 experiments, is that the performance of filter methods varies significantly *w.r.t.* dif-
514 ferent numbers of selected features (K). Indeed, all these filter methods only select
515 subgraph features with maximum discriminative score regardless of the redundancy
516 among the features. Adding redundant features may decrease the performance of an
517 algorithm. Further analysis of subgraph features is presented with a case study in the
518 next subsection.

519 In contrast, for embedded methods, the above drawbacks can be handled effec-
520 tively. Our algorithm RLMD unifies the feature selection and model learning into a
521 whole framework, so that the feature selection process is driven by the well-defined
522 objective function, and the selected features can further enhance the learning models.
523 At the same time, RLMD is guaranteed to be convergent given an appropriate γ_1 value,
524 which means that we do not need to specify the total number of selected graphs K . For

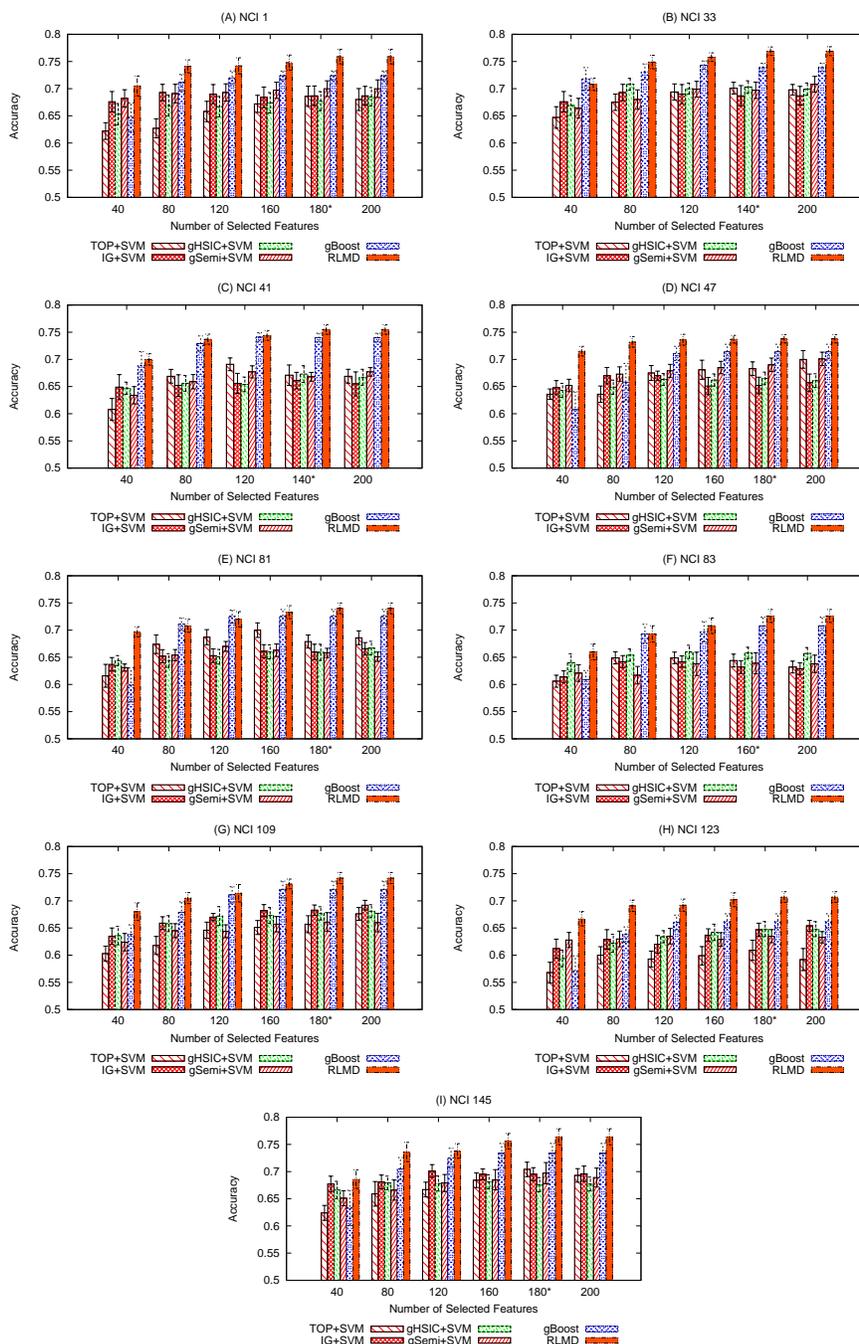


Figure 3: The classification accuracy and standard error on NCI dataset *w.r.t.* the number of selected graphs. The optimal number of subgraphs selected by RLMD (once it converges for all 10 folds experiments) is marked by a star * at x-axis.

instance, in Fig. (3).A, RLMD reaches convergence with 180 features.

RLMD vs. gBoost: It is evident that RLMD outperforms gBoost for most datasets. This is mainly because gBoost only uses ℓ_1 -norm regularization to produce a sparse solution. As pointed out in [21], the lasso (ℓ_1 -norm) has several drawbacks: (1) when the number of features (m , which is exponentially huge) is much bigger than the number of observations (n), the ℓ_1 norm selects at most n variables before it saturates; and (2) when the pairwise correlations in a group of variables are very high, lasso tends to select only one variable from the group and does not discriminate which one it selects. By contrast, RLMD uses an elastic net (combination of ℓ_1 and ℓ_2 norm), which encourages a grouping effect, where strongly correlated features will be included/excluded. As a result, RLMD results in a similar sparsity of representation to gBoost, but often outperforms gBoost.

Table 2: Averaged accuracies and standard errors on NCI Graphs with Optimal K .

ID	RLMD	TOP +SVM	IG +SVM	gHSIC +SVM	gSemi +SVM	gBoost
1	0.759 \pm 0.014	0.686 \pm 0.019	0.687 \pm 0.018	0.677 \pm 0.018	0.700 \pm 0.015	0.724 \pm 0.009
33	0.769 \pm 0.008	0.701 \pm 0.011	0.687 \pm 0.019	0.703 \pm 0.012	0.697 \pm 0.015	0.739 \pm 0.008
41	0.755 \pm 0.009	0.677 \pm 0.017	0.649 \pm 0.015	0.672 \pm 0.015	0.670 \pm 0.009	0.740 \pm 0.008
47	0.738 \pm 0.009	0.683 \pm 0.013	0.652 \pm 0.015	0.665 \pm 0.012	0.690 \pm 0.012	0.714 \pm 0.014
81	0.740 \pm 0.010	0.679 \pm 0.012	0.660 \pm 0.014	0.659 \pm 0.015	0.658 \pm 0.009	0.725 \pm 0.013
83	0.726 \pm 0.012	0.644 \pm 0.012	0.632 \pm 0.012	0.658 \pm 0.011	0.639 \pm 0.020	0.708 \pm 0.016
109	0.742 \pm 0.010	0.657 \pm 0.016	0.683 \pm 0.009	0.677 \pm 0.012	0.661 \pm 0.017	0.721 \pm 0.015
123	0.707 \pm 0.010	0.609 \pm 0.019	0.647 \pm 0.012	0.648 \pm 0.014	0.635 \pm 0.012	0.663 \pm 0.013
145	0.764 \pm 0.015	0.704 \pm 0.013	0.695 \pm 0.012	0.676 \pm 0.013	0.697 \pm 0.020	0.734 \pm 0.018

Overall Performance with Optimal K : In Table 2, we summarize the performance of our algorithm under optimal K value with other methods, where filter methods use the same number of subgraphs (K) for graph classification, and gBoost runs until convergence. The result in Table 2 clearly demonstrates that RLMD outperforms two-step filter methods (TOP+SVM and IG+SVM), direct filter methods (gHSIC+SVM and gSemi+SVM), and gBoost algorithm in NCI datasets.

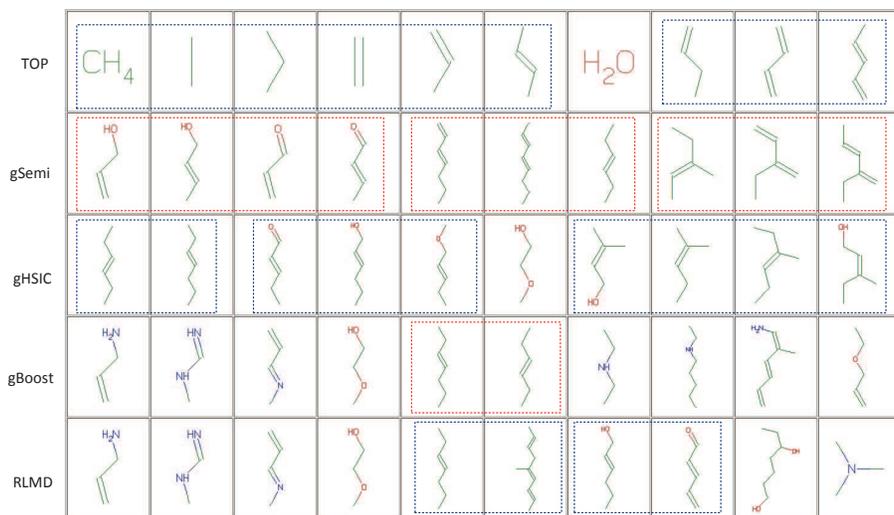


Figure 4: Case Study: comparison of discriminative subgraph features discovered by different algorithms. Subgraph features with high similarities are grouped and highlighted in the dashed rectangles. Subgraphs mined by filter methods are similar to each other and share high redundancy.

5.2.2. Case Study: Subgraph Feature Comparison

In this subsection, we use NCI-1 dataset as a case study to investigate subgraphs discovered by different algorithms. In our experiments, the top-10 subgraph features are discovered and illustrated in Fig. 4.

It is evident that the features for all filter methods (TOP, gSemi, and gHSIC) share high correlations. For the gSemi algorithm, for instance, the top-10 subgraphs form 3 groups. In each group, the subgraph features are very similar to each other. This is because the subgraph mining algorithm follows the depth-first-search (DFS) scheme, and subgraphs from the same sub-tree are very close to each other in terms of their geometrical structure. Because these methods consider each subgraph independently, the selected subgraphs may have high redundancy, which imposes a great challenge in determining the optimal K subgraphs for graph classification and also causes fluctuating results when the K values are varied.

In contrast, the subgraph correlations for gBoost and RLMD are much smaller.

557 The subgraphs discovered by gBoost and RLMD are highly overlapping (the first 5
558 subgraphs are identical). As pointed out by [21], ℓ_1 regularization tends to select only
559 one subgraph from a group of features and is not selective about which one is included,
560 thus the redundancy among the features in gBoost is minimal. By using elastic norm,
561 RLMD retains several group effects (some discriminative features may be included and
562 excluded simultaneously), and usually achieves better results. This result is consistent
563 with observations reported in [21] for vector data.

564 5.2.3. Results on PTC Tasks

565 We also conducted extensive experiments on the PTC datasets. The accuracies and
566 AUC values (*i.e.*, the area under ROC curves) are reported in Tables 3 and 4, where the
567 results are obtained after RLMD converges, and $K = 200$ for all filter methods.

Table 3: Accuracies and standard errors on PTC graphs with Optimal K for 10-fold cross-validation.

ID	RLMD	TOP +SVM	IG +SVM	gHSIC +SVM	gSemi +SVM	gBoost
MR	0.655 \pm 0.013	0.606 \pm 0.024	0.607 \pm 0.025	0.596 \pm 0.019	0.601 \pm 0.020	0.608 \pm 0.027
MM	0.664 \pm 0.019	0.060 \pm 0.025	0.599 \pm 0.023	0.613 \pm 0.021	0.603 \pm 0.019	0.622 \pm 0.018
FR	0.704 \pm 0.018	0.607 \pm 0.029	0.584 \pm 0.029	0.635 \pm 0.026	0.619 \pm 0.023	0.678 \pm 0.008
FM	0.615 \pm 0.018	0.592 \pm 0.031	0.594 \pm 0.026	0.581 \pm 0.021	0.575 \pm 0.018	0.603 \pm 0.017

Table 4: AUC values and standard errors on PTC graphs with Optimal K for 10-fold cross-validation.

ID	RLMD	TOP +SVM	IG +SVM	gHSIC +SVM	gSemi +SVM	gBoost
MR	0.681 \pm 0.021	0.597 \pm 0.025	0.596 \pm 0.025	0.560 \pm 0.021	0.580 \pm 0.021	0.649 \pm 0.033
MM	0.680 \pm 0.020	0.593 \pm 0.025	0.590 \pm 0.025	0.583 \pm 0.024	0.600 \pm 0.022	0.600 \pm 0.034
FR	0.673 \pm 0.022	0.600 \pm 0.029	0.575 \pm 0.029	0.618 \pm 0.021	0.623 \pm 0.021	0.640 \pm 0.026
FM	0.614 \pm 0.013	0.585 \pm 0.032	0.587 \pm 0.024	0.582 \pm 0.021	0.580 \pm 0.017	0.583 \pm 0.017

568 The results in Tables 3 and 4 show that RLMD achieves considerable performance
569 gains over all filter methods (TFM and DFM) and gBoost algorithm for all PTC datasets.

570 Note that for PTC classifications, AUC values are more important because they are all
 571 imbalanced classification tasks.

572 5.2.4. Parameter Analysis

573 In this subsection, we study the impact of parameters γ_1 and γ_2 on algorithm per-
 574 formance. Both γ_1 and γ_2 values are selected from $\{0, 0.01, 0.03, 0.05, 0.07, 0.09,$
 575 $0.11, 0.15, 0.2\}$, and the results under 10-fold cross-validation on NCI-1 and NCI-33
 576 are shown in Fig. 5.

577 **Impact of γ_1 values:** The experimental results in Fig. 5 show that γ_1 plays a more
 578 important role for the final classification model. With the increase of γ_1 from 0 to 0.2,
 the classification performance drops rapidly in terms of accuracy.

Table 5: Impact of different γ_1 values on NCI-1 dataset with $\gamma_2 = 0.03, S_{max} = 200$.

γ_1	0	0.01	0.03	0.05	0.15
#Selected Subgraphs	200	180	130	65	0
Accuracy	0.775	0.759	0.711	0.679	0.5
AUC	0.831	0.811	0.795	0.713	0.5

579 To better understand the impact of γ_1 , we also summarize the number of subgraphs
 580 selected with different γ_1 values in Table 5. The results show that increasing γ_1 values
 581 will result in fewer subgraphs being selected for the final classifier model, because a
 582 larger ℓ_1 norm regularizes more elements as 0. For $\gamma_1 = 0$, there is no sparse solu-
 583 tion. In other words, every subgraph should be used for graph classification. In this
 584 case, RLMD will only terminate when all subgraphs are incorporated for learning the
 585 model, or the maximum number of iterations S_{max} is reached. As the subgraph space
 586 is exponentially large, it is impractical to use all subgraph features to learn the model.
 587 The algorithm relies on S_{max} to terminate (200 is set in our experiment). The result
 588 shows that $\gamma_1=0$ even achieves better classification result, which is attributed to the fact
 589 that although ℓ_1 regularization introduces a sparse solution, it may be biased in some
 590 applications [45], so the accuracy may drop. For other cases with γ_1 being considerably
 591 large ($\gamma_1 = 0.5$), the regularization term dominates the objective function Eq.(4) with
 592 no subgraph being used for classification, which results in poor classification accuracy.
 593

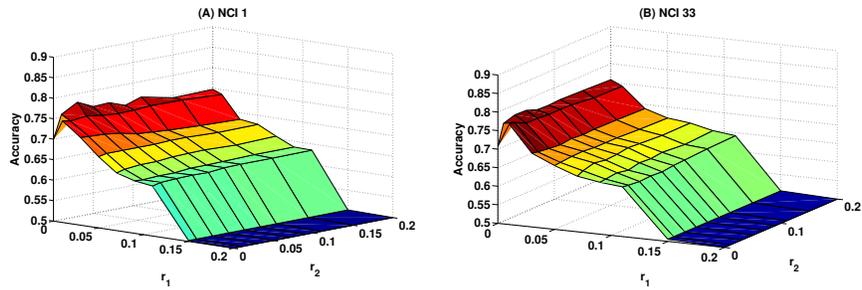


Figure 5: The accuracies with different γ_1 and γ_2 values.

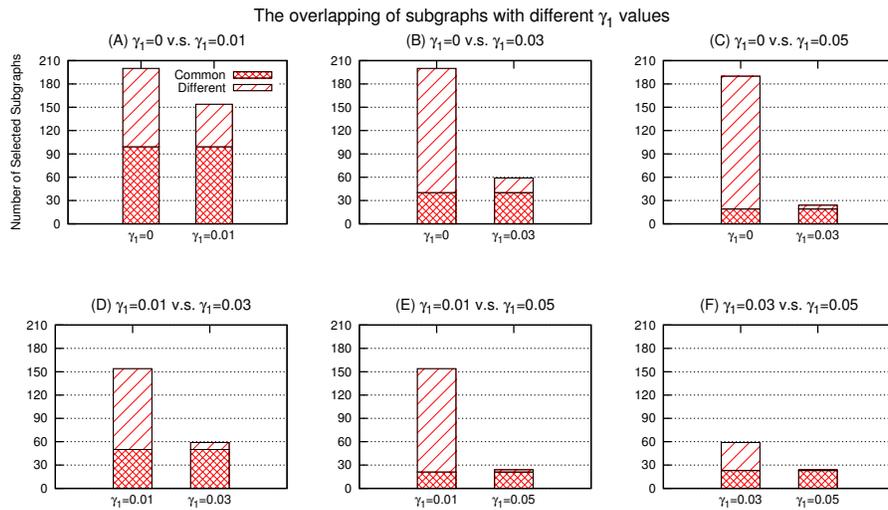


Figure 6: The overlapping of subgraphs (common subgraphs vs. different subgraphs) for different γ_1 values on NCI-1 dataset with a 70%-30% splitting on the NCI-1 dataset, *i.e.*, 70% graphs are randomly selected as training graphs, and 30% are used as test graphs. $\gamma_2 = 0.03$ and $S_{max} = 200$.

Table 6: Impact of different γ_2 values on NCI-1 dataset with $\gamma_1 = 0.01$.

γ_2	0	0.03	0.07	0.15	0.2
Accuracy	0.748	0.759	0.74	0.741	0.72
AUC	0.810	0.811	0.806	0.811	0.78

Table 7: Impacts of different ϵ values on NCI-1 dataset.

ϵ	0.01	0.005	0.001	0.0001
Accuracy	0.748	0.759	0.761	0.760
AUC	0.803	0.811	0.814	0.812

594 Note that the convergence property of our algorithm is dependent on γ_1 . In our
 595 experiments, we notice that γ_1 is very easy to set (in a small range [0.01,0.03]) for
 596 obtaining satisfactory results. This is much easier than requiring users to specify the
 597 number of subgraph features K needed for each graph dataset, because users may not
 598 have any prior knowledge about the selection of K values for different datasets, and
 599 different K values often result in significant changes in the algorithm performance.

600 **Interplay between γ_1 and subgraph selection:** We further compare the common sub-
 601 graphs selected by different γ_1 values, and report the results in Fig. 6. The results show
 602 that the subgraphs selected by using a smaller γ_1 values contain many subgraphs which
 603 are selected by using a larger γ_1 value. This observation is further evident in Fig.6.(E)
 604 and (F). The reason is that a smaller γ_1 value will result in more subgraph features to be
 605 selected, which increases the possibility of covering a small subgraph set selected by
 606 using a larger γ_1 value. In other words, a slightly smaller γ_1 value will result in more
 607 subgraph feature candidates to be explored and be beneficial for the classification task.

608 **Impact of γ_2 values:** We also vary γ_2 from 0 to 0.2, and report the results in Table
 609 6. The results show that a small regularization value $\gamma_2 = 0.03$ outperforms the case
 610 of $\gamma_2 = 0$, where the ℓ_2 regularization effect disappears (only ℓ_1 is used). This result
 611 is consistent with observations from a previous study [21]. This may be because ℓ_1
 612 ignores the correlated subgraphs in a group of features. When γ_2 keeps increasing,
 613 the classification performance drops because the larger ℓ_2 regularization dominates the
 614 objective function and the loss minimization term has less effect.

615 **Impact of ϵ values:** We vary the ϵ values from 0.01 to 0.0001 to study the final classifi-
 616 cation performance of our algorithm, and report the final classification results in Table
 617 7. The results show that as long as ϵ is subtle (from 0.001 to 0.0001), our classifica-
 618 tion can achieve similar classification results, *i.e.*, an ϵ -tolerance accuracy result to the

619 optimal solution.

620 **6. Conclusion**

621 In this paper, we proposed a regularized loss minimization-driven (RLMD) graph
622 classification method. We argued that existing filter-based subgraph selection methods
623 simply focus on finding most discriminative subgraph features, and suffer severe disad-
624 vantages in determining the optimal number of subgraphs for graph classification and
625 separating feature selection from the model learning phase. As a result, they might be
626 able to find most discriminative subgraph features, but cannot form high accuracy clas-
627 sifiers because they cannot determine how many discriminative features are needed to
628 train classifiers with the best performance gain. By integrating subgraph mining, dis-
629 criminative subgraph selection, and model learning into one unified framework, RMLD
630 is able to automatically determine the optimal number of discriminative subgraphs for
631 best graph classification results. Our algorithm generalizes the state-of-the-art gBoost
632 algorithm in the sense that it can employ any differentiable loss function and achieve
633 better classification accuracy by using an elastic net regularization. Experimental re-
634 sults on real-world graph datasets show a clear performance gain over existing two-step
635 filter methods (TFMs), direct filter methods (DFMs), and embedding methods.

636 **Acknowledgment**

637 We thank anonymous reviewers for their constructive comments. This work is par-
638 tially supported by Australia ARC Discovery Project (DP140102206), and National
639 Scholarship for Building High Level Universities, China Scholarship Council (No.
640 2011630030).

641 **References**

- 642 [1] M. Deshpande, M. Kuramochi, N. Wale, G. Karypis, Frequent substructure-based ap-
643 proaches for classifying chemical compounds, IEEE Transactions on Knowledge and Data
644 Engineering (2005) 1036–1050.

- 645 [2] J. Vogelstein, W. Roncal, R. Vogelstein, C. Priebe, Graph classification using signal-
646 subgraphs: Applications in statistical connectomics, *IEEE Transactions on Pattern Analy-
647 sis and Machine Intelligence* 35 (7) (2013) 1539–1551.
- 648 [3] A. Morales-Gonzlez, N. Acosta-Mendoza, A. Gago-Alonso, E. B. Garca-Reyes, J. E.
649 Medina-Pagola, A new proposal for graph-based image classification using frequent ap-
650 proximate subgraphs, *Pattern Recognition* 47 (1) (2014) 169 – 177.
- 651 [4] S. Pan, X. Zhu, C. Zhang, P. S. Yu, Graph stream classification using labeled and unlabeled
652 graphs, in: 2013 IEEE 29th International Conference on Data Engineering (ICDE), IEEE,
653 2013, pp. 398–409.
- 654 [5] F. Wyszotki, W. Kolbe, J. Selbig, Concept learning by structured examples-an algebraic
655 approach., in: *International Joint Conferences on Artificial Intelligence*, 1981, pp. 153–
656 158.
- 657 [6] P. Geibel, F. Wyszotki, Learning relational concepts with decision trees, in: *International
658 Conference on Machine Learning*, Vol. 96, 1996, pp. 166–174.
- 659 [7] B. J. Jain, K. Obermayer, Structure spaces, *Journal of Machine Learning Research* 10
660 (2009) 2667–2714.
- 661 [8] B. J. Jain, K. Obermayer, Learning in Riemannian orbifolds, *arXiv preprint
662 arXiv:1204.4294* (2012).
- 663 [9] H. Kashima, K. Tsuda, A. Inokuchi, *Kernels for graphs*, MIT Press, Cambridge (Mas-
664 sachusetts), 2004, Ch. In: Schlkopf B, Tsuda K, Vert JP, editors. *Kernel methods in com-
665 putational biology*, pp. 101–113.
- 666 [10] H. Fei, J. Huan, Boosting with structure information in the functional space: an application
667 to graph classification, in: *ACM SIGKDD Conference on Knowledge Discovery and Data
668 Mining (KDD)*, Washington DC, USA, 2010.
- 669 [11] X. Kong, P. Yu, Semi-supervised feature selection for graph classification, in: *Proceedings
670 of the 16th ACM SIGKDD international conference on Knowledge discovery and data
671 mining*, 2010, pp. 793–802.
- 672 [12] X. Kong, P. S. Yu, Multi-label feature selection for graph classification, in: *IEEE Interna-
673 tional Conference on Data Mining*, IEEE, 2010, pp. 274–283.

- 674 [13] Y. Zhu, J. Yu, H. Cheng, L. Qin, Graph classification: A diversified discriminative feature
675 selection approach, in: Conference on Information and Knowledge Management (CIKM),
676 ACM, 2012, pp. 205–214.
- 677 [14] N. Jin, C. Young, W. Wang, Graph classification based on pattern co-occurrence, in: Con-
678 ference on Information and Knowledge Management (CIKM), 2009, pp. 573–582.
- 679 [15] M. Thoma, H. Cheng, A. Gretton, J. Han, H. Kriegel, A. Smola, L. Song, P. Yu, X. Yan,
680 K. Borgwardt, Near-optimal supervised feature selection among frequent subgraphs, in:
681 SIAM International Conference on Data Mining, USA, 2009.
- 682 [16] S. Ranu, A. K. Singh, Graphsig: A scalable approach to mining significant subgraphs in
683 large graph databases, in: International Conference on Data Engineering (ICDE), IEEE,
684 2009, pp. 844–855.
- 685 [17] X. Yan, H. Cheng, J. Han, P. S. Yu, Mining significant graph patterns by leap search, in:
686 ACM SIGMOD Conference, ACM, 2008, pp. 433–444.
- 687 [18] H. Cheng, X. Yan, J. Han, Discriminative frequent pattern-based graph classification, *Link*
688 *Mining: Models, Algorithms, and Applications* (2010) 237–262.
- 689 [19] I. Guyon, A. Elisseeff, An introduction to variable and feature selection, *Journal of Ma-*
690 *chine Learning Research* 3 (2003) 1157–1182.
- 691 [20] R. Tibshirani, Regression shrinkage and selection via the lasso, *Journal of the Royal Sta-*
692 *tistical Society. Series B (Methodological)* (1996) 267–288.
- 693 [21] H. Zou, T. Hastie, Regularization and variable selection via the elastic net, *Journal of the*
694 *Royal Statistical Society: Series B (Statistical Methodology)* 67 (2) (2005) 301–320.
- 695 [22] T. Zhou, D. Tao, X. Wu, Manifold elastic net: A unified framework for sparse dimension
696 reduction, *Data Mining and Knowledge Discovery* 22 (3) (2011) 340–371.
- 697 [23] H. Saigo, S. Nowozin, T. Kadowaki, T. Kudo, K. Tsuda, gBoost: A mathematical pro-
698 gramming approach to graph classification and regression, *Machine Learning* 75 (2009)
699 69–89.
- 700 [24] M. Culp, G. Michailidis, Graph-based semisupervised learning, *IEEE Transactions on Pat-*
701 *tern Analysis and Machine Intelligence* 30 (1) (2008) 174–179.

- 702 [25] A. Mantrach, N. van Zeebroeck, P. Francq, M. Shimbo, H. Bersini, M. Saerens, Semi-
703 supervised classification and betweenness computation on large, sparse, directed graphs,
704 *Pattern Recognition* 44 (6) (2011) 1212 – 1224.
- 705 [26] M. Karasuyama, H. Mamitsuka, Multiple graph label propagation by sparse integration,
706 *IEEE Transactions on Neural Networks and Learning Systems* 24 (12) (2013) 1999–2012.
- 707 [27] X. Liu, L. Wang, J. Zhang, J. Yin, H. Liu, Global and local structure preservation for
708 feature selection, *IEEE Transactions on Neural Networks and Learning Systems* 25 (6)
709 (2014) 1083–1095.
- 710 [28] D. Liben-Nowell, J. Kleinberg, The link-prediction problem for social networks, *Journal of*
711 *the American Society for Information Science and Technology* 58 (7) (2007) 1019–1031.
- 712 [29] C. Nguyen, H. Mamitsuka, Latent feature kernels for link prediction on sparse graphs,
713 *IEEE Transactions on Neural Networks and Learning Systems* 23 (11) (2012) 1793–1804.
- 714 [30] N. Jin, C. Young, W. Wang, GAIA: graph classification using evolutionary computation,
715 in: *ACM SIGMOD Conference*, ACM, 2010, pp. 879–890.
- 716 [31] J. Wu, X. Zhu, C. Zhang, P. Yu, Bag constrained structure pattern mining for multi-graph
717 classification, *IEEE Transactions on Knowledge and Data Engineering* 26 (10) (2014) 2382
718 – 2396.
- 719 [32] L. Bai, L. Rossi, A. Torsello, E. R. Hancock, A quantum Jensen-Shannon graph kernel for
720 unattributed graphs, *Pattern Recognition* 48 (2) (2015) 344 – 355.
- 721 [33] N. Shervashidze, P. Schweitzer, E. J. Van Leeuwen, K. Mehlhorn, K. M. Borgwardt,
722 Weisfeiler-Lehman graph kernels, *Journal of Machine Learning Research* 12 (2011) 2539–
723 2561.
- 724 [34] M. Neumann, N. Patricia, R. Garnett, K. Kersting, Efficient graph kernels by randomiza-
725 tion, in: *Machine Learning and Knowledge Discovery in Databases*, Springer, 2012, pp.
726 378–393.
- 727 [35] K. Riesen, H. Bunke, Graph classification by means of Lipschitz embedding, *IEEE Trans-*
728 *actions on SMC, Part B: Cybernetics* 39 (2009) 1472–1483.
- 729 [36] J. Wu, Z. Hong, S. Pan, X. Zhu, C. Zhang, Multi-graph-view learning for graph classifica-
730 tion, in: *IEEE International Conference on Data Mining*, 2014, pp. 590–599.

- 731 [37] J. Wu, S. Pan, X. Zhu, Z. Cai, Boosting for multi-graph classification, *IEEE Transactions*
732 *on Cybernetics* 45 (3) (2015) 430–443.
- 733 [38] J. Wu, Z. Hong, S. Pan, X. Zhu, C. Zhang, Z. Cai, Multi-graph learning with positive
734 and unlabeled bags, in: *SIAM International Conference on Data Mining*, SIAM, 2014, pp.
735 217–225.
- 736 [39] S. Pan, X. Zhu, Graph classification with imbalanced class distributions and noise, in:
737 *International Joint Conferences on Artificial Intelligence*, 2013, pp. 1586–1592.
- 738 [40] S. Pan, J. Wu, X. Zhu, C. Zhang, Graph ensemble boosting for imbalanced noisy graph
739 stream classification, *IEEE Transactions on Cybernetics* 45 (5) (2015) 940–954.
- 740 [41] S. Pan, J. Wu, X. Zhu, Cogboost: Boosting for fast cost-sensitive graph classification, *IEEE*
741 *Transactions on Knowledge and Data Engineering (In Press)* (2015) 1. [doi:10.1109/
742 TKDE.2015.2391115](https://doi.org/10.1109/TKDE.2015.2391115).
- 743 [42] X. Yan, J. Han, gSpan: Graph-based substructure pattern mining, in: *IEEE International*
744 *Conference on Data Mining*, 2002, pp. 721–724.
- 745 [43] A. Demiriz, K. Bennett, J. Shawe-Taylor, Linear programming boosting via column gener-
746 ation, *Machine Learning* (2002) 225–254.
- 747 [44] T. Kudo, E. Maeda, Y. Matsumoto, An application of boosting to graph classification, in:
748 *Neural Information Processing Systems (NIPS)*, 2004, pp. 729–736.
- 749 [45] M. Tan, I. W. Tsang, L. Wang, Towards ultrahigh dimensional feature selection for big
750 data, *Journal of Machine Learning Research* 15 (2014) 1371–1429.