
Condensing Graphs via One-Step Gradient Matching

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 As training deep learning models on large dataset takes a lot of time and resources,
2 it is desired to construct a small synthetic dataset with which we can train deep
3 learning models sufficiently. There are recent works that have explored solutions
4 on condensing image datasets through complex bi-level optimization. For instance,
5 dataset condensation (DC) matches network gradients w.r.t. large-real data and
6 small-synthetic data, where the network weights are optimized for multiple steps at
7 each outer iteration. However, existing approaches have their inherent limitations:
8 (1) they are not directly applicable to graphs where the data is discrete; and (2) the
9 condensation process is computationally expensive due to the involved nested opti-
10 mization. To bridge the gap, we investigate efficient dataset condensation tailored
11 for graph datasets where we model the discrete graph structure as a probabilistic
12 model. We further propose a one-step gradient matching scheme, which performs
13 gradient matching for only one single step without training the network weights.
14 Our theoretical analysis shows this strategy can generate synthetic graphs that lead
15 to lower classification loss on real graphs. Extensive experiments on various graph
16 datasets demonstrate the effectiveness and efficiency of the proposed method. In
17 particular, we are able to reduce the dataset size by 90% while approximating up
18 to 98% of the original performance and our method is significantly faster than
19 multi-step gradient matching (e.g. 15× in CIFAR10 for synthesizing 500 graphs).

20 1 Introduction

21 Graph-structured data plays a key role in various real-world applications. For example, by exploiting
22 graph structural information, we can predict the chemical property of a given molecular graph [1],
23 detect fraud activities in a financial transaction graph [2], or recommend new friends to users
24 in a social network [3]. Due to its prevalence, graph neural networks (GNNs) [4, 5, 6, 7] have
25 been developed to effectively extract meaningful patterns from graph data and thus tremendously
26 facilitate computational tasks on graphs. Despite their effectiveness, GNNs are notoriously data-
27 hungry like traditional deep neural networks: they usually require massive datasets to learn powerful
28 representations. Thus, training GNNs is often computationally expensive. Such cost even becomes
29 prohibitive when we need to repeatedly train GNNs, e.g., in neural architecture search [8] and
30 continual learning [9].

31 One potential solution to alleviate the aforementioned issue is *dataset condensation* or *dataset*
32 *distillation*. It targets at constructing a small-synthetic training set that can provide sufficient
33 information to train neural networks [10, 11, 12, 13, 14, 15, 16]. In particular, one of the representative
34 methods, DC [11], formulates the condensation goal as matching the gradients of the network
35 parameters between small-synthetic and large-real training data. It has been demonstrated that such a
36 solution can greatly reduce the training set size of image datasets without significantly sacrificing
37 model performance. For example, using 100 images generated by DC can achieve 97.4% test accuracy
38 on MNIST compared with 99.6% on the original dataset (60, 000 images). These condensed samples

39 can significantly save space for storing datasets and speed up retraining neural networks in many
40 critical applications, e.g., continual learning and neural architecture search. In spite of the recent
41 advances in dataset distillation/condensation for images, limited attention has been paid on domains
42 involving graph structures.

43 To bridge this gap, we investigate the problem of condensing graphs such that GNNs trained on
44 condensed graphs can achieve comparable performance to those trained on the original dataset.
45 However, directly applying existing solutions for dataset condensation [10, 11, 12, 13] to graph
46 domain faces some challenges. First, existing solutions have been designed for images where the data
47 is continuous and they cannot output binary values to form the discrete graph structure. Thus, we
48 need to develop a strategy that can handle the discrete nature of graphs. Second, they usually involve
49 a complex bi-level problem that is computationally expensive to optimize: they require multiple
50 iterations (inner iterations) of updating neural network parameters before updating the synthetic data
51 for multiple iterations (outer iterations). It can be catastrophically inefficient for learning pairwise
52 relations for nodes, of which the complexity is quadratic to the number of nodes.

53 To address the aforementioned challenges, we propose an efficient condensation method for graphs,
54 where we follow DC [11] to match the gradients of GNNs between synthetic graphs and real graphs.
55 In order to produce discrete values, we model the graph structure as a probabilistic graph model and
56 optimize the discrete structures in a differentiable manner. Based on this formulation, we further
57 propose a *one-step gradient matching* strategy which only performs gradient matching for one single
58 step. Consequently, the advantages of the proposed strategy are twofold. First, it significantly speeds
59 up the condensation process while providing reasonable guidance for synthesizing condensed graphs.
60 Second, it removes the burden of tuning hyper-parameters such as the number of outer/inner iterations
61 of the bi-level optimization as required by DC. Furthermore, we demonstrate the effectiveness of the
62 proposed one-step gradient matching strategy both theoretically and empirically. Our contributions
63 can be summarized as follows:

- 64 1. We study a novel problem of learning discrete synthetic graphs for condensing graph datasets,
65 where the discrete structure is captured via a graph probabilistic model that can be learned in a
66 differentiable manner.
- 67 2. We propose a one-step gradient matching scheme that significantly accelerates the vanilla gradient
68 matching process. Theoretical analysis is provided to understand the rationality of the proposed
69 one-step gradient matching. We show that learning with one-step matching produces synthetic
70 graphs that lead to a smaller classification loss on real graphs.
- 71 3. Extensive experiments have demonstrated the effectiveness and efficiency of the proposed method.
72 Particularly, we are able to reduce the dataset size by 90% while approximating up to 98% of the
73 original performance and our method is significantly faster than multi-step gradient matching (e.g.
74 15 \times in CIFAR10 for synthesizing 500 graphs).

75 2 The Proposed Framework

76 Before detailing the framework, we first introduce the main notations used in this paper. We
77 majorly focus on the graph classification task where the goal is to predict the labels of given graphs.
78 Specifically, we denote a graph dataset as $\mathcal{T} = \{G_1, \dots, G_N\}$ with ground-truth label set \mathcal{Y} . Each
79 graph in \mathcal{T} is associated with a discrete adjacency matrix and a node feature matrix. Let $\mathbf{A}_{(i)}$, $\mathbf{X}_{(i)}$
80 represent the adjacency matrix and the feature matrix of i -th real graph, respectively. Similarly, we
81 use $\mathcal{S} = \{G'_1, \dots, G'_{N'}\}$ and \mathcal{Y}' to indicate the synthetic graphs and their labels, respectively. Note
82 that the number of synthetic graphs N' is essentially much smaller than that of real graphs N . We
83 use d and n to denote the number of feature dimensions the number of nodes in each synthetic graph,
84 respectively¹. Let C denote the number of classes and ℓ denote the cross entropy loss. The goal of
85 our work is to learn a set of synthetic graphs \mathcal{S} such that a GNN trained on \mathcal{S} can achieve comparable
86 performance to the one trained on the much larger dataset \mathcal{T} .

87 2.1 Gradient Matching as the Objective

88 Since we aim at learning synthetic graphs that are highly informative, one solution is to allow GNNs
89 trained on synthetic graphs to imitate the training trajectory on the original large dataset. Dataset
90 condensation [11, 12] introduces a gradient matching scheme to achieve this goal. Concretely, it tries
91 to reduce the difference of model gradients w.r.t. large-real data and small-synthetic data for model

¹We set n to the average number of nodes in original dataset.

92 parameters at every training epoch. Hence, the model parameters trained on synthetic data will be
 93 close to these trained on real data at every training epoch. Let θ_t denote the network parameters at
 94 the t -th epoch and f_{θ_t} indicate the neural network parameterized by θ_t . The condensation objective
 95 is expressed as:

$$\min_S \sum_{t=0}^{T-1} D(\nabla_{\theta} \ell(f_{\theta_t}(\mathcal{S}), \mathcal{Y}'), \nabla_{\theta} \ell(f_{\theta_t}(\mathcal{T}), \mathcal{Y})), \quad \text{s.t. } \theta_{t+1} = \text{opt}_{\theta}(\theta_t, \mathcal{S}), \quad (1)$$

96 where $D(\cdot, \cdot)$ is a distance function, T is the number of steps of the whole training trajectory and
 97 $\text{opt}_{\theta}(\cdot)$ is the optimization operator for updating parameter θ . Note that Eq. (1) is a bi-level problem
 98 where we need to learn the synthetic graphs \mathcal{S} at the outer optimization and update model parameters
 99 θ_t at the inner optimization. To learn synthetic graphs that generalize to a distribution of model
 100 parameters P_{θ_0} , we sample $\theta_0 \sim P_{\theta_0}$ and rewrite Eq. (1) as:

$$\min_S E_{\theta_0 \sim P_{\theta_0}} \left[\sum_{t=0}^{T-1} D(\nabla_{\theta} \ell(f_{\theta_t}(\mathcal{S}), \mathcal{Y}'), \nabla_{\theta} \ell(f_{\theta_t}(\mathcal{T}), \mathcal{Y})) \right], \quad \text{s.t. } \theta_{t+1} = \text{opt}_{\theta}(\theta_t, \mathcal{S}). \quad (2)$$

101 **Discussion.** The aforementioned strategy has demonstrated promising performance on condensing
 102 image datasets [11, 12]. However, it is not clear how to model the discrete graph structure. Moreover,
 103 the inherent bi-level optimization inevitably hinders its scalability. To tackle these shortcomings, we
 104 propose *DosCond* that models the structure as a probabilistic graph model and is optimized through
 105 one-step gradient matching. In the following subsections, we introduce the details of *DosCond*.

106 2.2 Learning Discrete Graph Structure

107 For graph classification, each graph in the dataset is composed of an adjacency matrix and a feature
 108 matrix. For simplicity, we use $\mathbf{X}' \in R^{N' \times n \times d}$ to denote the node features in all synthetic graphs
 109 \mathcal{S} and $\mathbf{A}' \in \{0, 1\}^{N' \times n \times n}$ to indicate the graph structure information in \mathcal{S} . Note that f_{θ_t} can be
 110 instantiated as any graph neural network and it takes both graph structure and node features as input.
 111 Then we rewrite the objective in Eq. (2) as follows:

$$\min_{\mathbf{A}', \mathbf{X}'} E_{\theta_0 \sim P_{\theta_0}} \left[\sum_{t=0}^{T-1} D(\nabla_{\theta} \ell(f_{\theta_t}(\mathbf{A}', \mathbf{X}'), \mathcal{Y}'), \nabla_{\theta} \ell(f_{\theta_t}(\mathcal{T}), \mathcal{Y})) \right], \quad \text{s.t. } \theta_{t+1} = \text{opt}_{\theta}(\theta_t, \mathcal{S}), \quad (3)$$

112 where we aim to learn both graph structure \mathbf{A}' and node features \mathbf{X}' . However, Eq. (3) is challenging
 113 to optimize as it requires a function that outputs binary values. To address this issue, we propose
 114 to model the graph structure as a probabilistic graph model with Bernoulli distribution. Note that
 115 in the following, we reshape \mathbf{A}' from $N' \times n \times n$ to $N' \times n^2$ for the purpose of demonstration
 116 only. Specifically, for each entry $\mathbf{A}'_{ij} \in \{0, 1\}$ in the adjacency matrix \mathbf{A}' , it follows a Bernoulli
 117 distribution: $P_{\Omega_{ij}}(\mathbf{A}'_{ij}) = \mathbf{A}'_{ij} \sigma(\Omega_{ij}) + (1 - \mathbf{A}'_{ij}) \sigma(-\Omega_{ij})$, where $\sigma(\cdot)$ is the sigmoid function;
 118 $\Omega_{ij} \in R$ is the success probability of the Bernoulli distribution and also the parameter to be
 119 learned. Since \mathbf{A}'_{ij} is independent of all other entries, the distribution of \mathbf{A}' can be modeled
 120 as: $P_{\Omega}(\mathbf{A}') = \prod_{i=1}^{N'} \prod_{j=1}^{n^2} P_{\Omega_{ij}}(\mathbf{A}'_{ij})$. Then, the objective in Eq. (2) needs to be modified to

$$\min_{\mathbf{A}', \mathbf{X}'} E_{\theta_0 \sim P_{\theta_0}} \left[E_{\mathbf{A}' \sim P_{\Omega}} [\ell(\mathbf{A}'(\Omega), \mathbf{X}', \theta_0)] \right]. \quad (4)$$

121 With the new parameterization, we obtain a function that outputs discrete values but it is not differentiable
 122 due to the involved sampling process. Thus, we employ the reparameterization method [17],
 123 binary concrete distribution, to refactor the discrete random variable into a differentiable function
 124 of its parameters and a random variable with fixed distribution. Specifically, we first sample
 125 $\alpha \sim \text{Uniform}(0, 1)$, and edge weight $\mathbf{A}'_{ij} \in [0, 1]$ is calculated by:

$$\mathbf{A}'_{ij} = \sigma((\log \alpha - \log(1 - \alpha) + \Omega_{ij}) / \tau), \quad (5)$$

126 where $\tau \in (0, \infty)$ is the temperature parameter that controls the continuous relaxation. As $\tau \rightarrow 0$,
 127 the random variable \mathbf{A}'_{ij} smoothly approaches the Bernoulli distribution. In other words, we have
 128 $\lim_{\tau \rightarrow 0} P(\mathbf{A}'_{ij} = 1) = \sigma(\Omega_{ij})$. While small τ is necessary for obtaining discrete samples, large

129 τ is useful in getting large gradients as suggested by [17]. In practice, we employ an annealing
 130 schedule [18] to gradually decrease the value of τ in training. With the reparameterization trick, the
 131 objective function becomes differentiable w.r.t. $\Omega_{i,j}$ with well-defined gradients. Then we rewrite our
 132 objective as:

$$\min_{\Omega, \mathbf{X}'_{\theta_0} \sim P_{\theta_0}} \left[E_{\alpha \sim \text{Uniform}(0,1)} [\ell(\mathbf{A}'(\Omega), \mathbf{X}', \theta_0)] \right] = \quad (6)$$

$$E_{\theta_0} \left[E_{\alpha} \left[\sum_{t=0}^{T-1} D(\nabla_{\theta} \ell(f_{\theta_t}(\mathbf{A}'(\Omega), \mathbf{X}'), \mathcal{Y}'), \nabla_{\theta} \ell(f_{\theta_t}(\mathcal{T}), \mathcal{Y})) \right] \right], \quad \text{s.t. } \theta_{t+1} = \text{opt}_{\theta}(\theta_t, \mathcal{S}).$$

133 2.3 One-Step Gradient Matching

134 The vanilla gradient matching scheme in Eq. (2) presents a bi-level optimization problem. To solve
 135 this problem, we need to update the synthetic graphs \mathcal{S} at the outer loop and then optimize the
 136 network parameters θ_t at the inner loop. The nested loops heavily impede the scalability of the
 137 condensation method, which motivates us to design a new strategy for efficient condensation. In this
 138 work, we propose a *one-step gradient matching* scheme where we only match the network gradients
 139 for the model initializations θ_0 while discarding the training trajectory of θ_t . Essentially, this strategy
 140 approximates the overall gradient matching loss for θ_t with the initial matching loss at the first epoch,
 141 which we term as *one-step matching loss*. The intuition is: the one-step matching loss informs us
 142 about the direction to update the synthetic data, in which, we have empirically observed a strong
 143 decrease in the cross-entropy loss (on real samples) obtained from the model trained on synthetic
 144 data. Hence, we can drop the summation symbol $\sum_{t=0}^{T-1}$ in Eq. (6) and simplify Eq. (6) as follows:

$$\min_{\Omega, \mathbf{X}'_{\theta_0}} E_{\alpha} \left[E_{\alpha} [D(\nabla_{\theta} \ell(f_{\theta_0}(\mathbf{A}'(\Omega), \mathbf{X}'), \mathcal{Y}'), \nabla_{\theta} \ell(f_{\theta_0}(\mathcal{T}), \mathcal{Y}))] \right], \quad (7)$$

145 where we sample $\theta_0 \sim P_{\theta_0}$ and $\alpha \sim \text{Uniform}(0, 1)$. Compared with Eq. (6), one-step gradient
 146 matching avoids the expensive nested-loop optimization and directly updates the synthetic graph
 147 \mathcal{S} . It greatly simplifies the condensation process. In practice, as shown in Section 3.3, we find this
 148 strategy yields comparable performance to its bi-level counterpart while enabling much more efficient
 149 condensation. Next, we provide theoretical analysis to understand the rationality of the proposed
 150 one-step gradient matching scheme.

151 **Theoretical Understanding.** We denote the cross entropy loss on the real graphs as $\ell_{\mathcal{T}}(\theta) =$
 152 $\sum_i \ell_i(\mathbf{A}_{(i)}, \mathbf{X}_{(i)}, \theta)$ and that on synthetic graphs as $\ell_{\mathcal{S}}(\theta) = \ell_{\mathcal{S}}(\mathbf{A}'_{(i)}, \mathbf{X}'_{(i)}, \theta)$. Let θ^* denote the
 153 optimal parameter and θ_t be the parameter trained on \mathcal{S} at the t -th epoch by optimizing $\ell_{\mathcal{S}}(\theta)$. For
 154 notation simplicity, we assume that \mathbf{A} and \mathbf{A}' are already normalized. The matrix norm $\|\cdot\|$ is the
 155 Frobenius norm. We focus on the GNN of Simple Graph Convolutions (SGC) [19] to study our
 156 problem since SGC has a simpler architecture but shares a similar filtering pattern as GCN.

157 **Theorem 1** *When we use a K -layer SGC as the GNN used in condensation, i.e., $f_{\theta}(\mathbf{A}_{(i)}, \mathbf{X}_{(i)}) =$
 158 $\text{Pool}(\mathbf{A}_{(i)}^K \mathbf{X}_{(i)} \mathbf{W}_1) \mathbf{W}_2$ with $\theta = [\mathbf{W}_1; \mathbf{W}_2]$ and assume that all network parameters satisfy $\|\theta\|^2 \leq$
 159 $M^2 (M > 0)$, we have*

$$\min_t \ell_{\mathcal{T}}(\theta_t) - \ell_{\mathcal{T}}(\theta^*) \leq \sum_{t=0}^{T-1} \frac{\sqrt{2}M}{T} \|\nabla_{\theta} \ell_{\mathcal{T}}(\theta_t) - \nabla_{\theta} \ell_{\mathcal{S}}(\theta_t)\| \frac{3M}{2\sqrt{T}} \frac{C-1}{CN'} \sqrt{\sum_i \gamma_i \|\mathbf{1}^{\top} \mathbf{A}'_{(i)}^K \mathbf{X}'_{(i)}\|^2} \quad (8)$$

160 where $\gamma_i = 1$ if we use sum pooling in f_{θ} ; $\gamma_i = \frac{1}{n_i}$ if we use mean pooling, with n_i as the number of
 161 nodes in the i -th synthetic graph.

162 We provide the proof of Theorem 1 in Appendix C.1. Theorem 1 suggests that the smallest gap
 163 between the resulted loss (by training on synthetic graphs) and the optimal loss has an upper bound.
 164 This upper bound depends on two terms: (1) the difference of gradients w.r.t. real data and synthetic
 165 data and (2) the norm of input matrices. Thus, the theorem justifies that reducing the gradient
 166 difference w.r.t. real and synthetic graphs can help learn desirable synthetic data that preserves
 167 sufficient information to train GNNs well. Based on Theorem 1, we have the following proposition.

168 **Proposition 1** Assume the largest gradient gap happens at 0-th epoch, i.e., $\|\nabla_{\theta} \ell_{\mathcal{T}}(\theta_0) -$
 169 $\nabla_{\theta} \ell_S(\theta_0)\| = \max_t \|\nabla_{\theta} \ell_{\mathcal{T}}(\theta_t) - \nabla_{\theta} \ell_S(\theta_t)\|$ with $t = 0, 1, \dots, T-1$, we have

$$\min_t \ell_{\mathcal{T}}(\theta_t) - \ell_{\mathcal{T}}(\theta^*) \leq \sqrt{2}M \|\nabla_{\theta} \ell_{\mathcal{T}}(\theta_0) - \nabla_{\theta} \ell_S(\theta_0)\| + \frac{3M}{2\sqrt{T}} \frac{C-1}{CN'} \sqrt{\sum_i \gamma_i \|\mathbf{1}^{\top} \mathbf{A}'_{(i)} \mathbf{X}'_{(i)}\|^2}. \quad (9)$$

170 We omit the proof for the proposition since it is straightforward. The above proposition suggests
 171 that the smallest gap between the $\ell_{\mathcal{T}}(\theta_t)$ and $\ell_{\mathcal{T}}(\theta^*)$ is bounded by the one-step matching loss and
 172 the norm $\|\mathbf{1}^{\top} \mathbf{A}'_{(i)} \mathbf{X}'_{(i)}\|^2$. As we will show in Section B.1, when using mean pooling, the second
 173 term tend to have a smaller scale than the first one and can be neglected; the second term matters
 174 more when we use sum pooling. Hence, we solely optimize the one-step gradient matching loss for
 175 GNNs with mean pooling and additionally include the second term (the norm of input matrices) as a
 176 regularization for GNNs with sum pooling. As such, when we consider the optimal loss $\ell_{\mathcal{T}}(\theta^*)$ as a
 177 constant, reducing the one-step matching loss indeed learns synthetic graphs that lead to a smaller
 178 loss on real graphs. This demonstrates the rationality of one-step gradient matching theoretically.

179 **Remark 1.** Note that the spectral analysis from [19] demonstrated that both GCN and SGC share
 180 similar graph filtering behaviors. Thus practically, we extend the one-step gradient matching loss
 181 from K -layer SGC to K -layer GCN and observe that it works well under the non-linear scenario.

182 **Remark 2.** While we focus on the graph classification task, it is straightforward to extend our
 183 framework to node classification. We obtain similar conclusions for node classification as shown in
 184 Theorem 2 in Appendix C.2 and achieve impressive empirical performance in Appendix B.2.

185 2.4 Final Objective and Training Algorithm

186 In this subsection, we describe the final objective function and the detailed training algorithm. We note
 187 that the objective in Eq. (6) involves two nested expectations, we adopt Monte Carlo to approximately
 188 optimize the objective function. Together with one-step gradient matching, we have

$$\min_{\Omega, \mathbf{X}'_{\theta_0}} E_{\alpha \sim \text{Uniform}(0,1)} E [[\ell(\mathbf{A}'(\Omega), \mathbf{X}', \theta_0)]] \approx \sum_{k_1=1}^{K_1} \sum_{k_2=1}^{K_2} D(\nabla_{\theta} \ell(f_{\theta_0}(\mathbf{A}'(\Omega), \mathbf{X}'), \mathcal{Y}'), \nabla_{\theta} \ell(f_{\theta_0}(T), \mathcal{Y}))$$

189 where K_1 is the number of sampled model initializations and K_2 is the number of sampled graphs.
 190 We find that $K_2 = 1$ is able to yield good performance in our experiments.

191 **Regularization.** In addition to the one-step gradient matching loss, we note that the proposed
 192 *DosCond* can be easily integrated with various priors as regularization terms. In this work, we focus
 193 on exerting sparsity regularization on the adjacency matrix, since a denser adjacency matrix will
 194 lead to higher cost for training graph neural networks. Specifically, we penalize the difference of the
 195 sparsity between $\sigma(\Omega)$ and a given sparsity ϵ :

$$\ell_{\text{reg}} = \max\left(\frac{1}{|\Omega|} \sum_{i,j} \sigma(\Omega_{ij}) - \epsilon, 0\right). \quad (10)$$

196 We initialize $\sigma(\Omega)$ and \mathbf{X}' as randomly sampled training graphs and set ϵ to the average sparsity of
 197 initialized $\sigma(\Omega)$ so as to maintain a low sparsity. On top of that, as we discussed earlier in Section 2.3,
 198 we include the following regularization for GNNs with sum pooling:

$$\ell_{\text{reg2}} = \frac{3}{2\sqrt{2T}} \cdot \frac{C-1}{CN'} \sqrt{\sum_i \|\mathbf{1}^{\top} \mathbf{A}'_{(i)} \mathbf{X}'_{(i)}\|^2} \quad (11)$$

199 **Training Algorithm.** We provide the details of our proposed framework in Algorithm 1 in Ap-
 200 pendix A.3. Specifically, we sample K_1 model initializations θ_0 to perform one-step gradient
 201 matching. Following the convention in DC [11], we match gradients and update synthetic graphs
 202 for each class separately in order to make matching easier. For class c , we first retrieve the synthetic
 203 graphs of that class, denoted as $(\mathbf{A}'_c, \mathbf{X}'_c, \mathcal{Y}'_c) \sim \mathcal{S}$, and sample a batch of real graphs $(\mathbf{A}_c, \mathbf{X}_c, \mathcal{Y}_c)$.
 204 We then forward them to the graph neural network and calculate the one-step gradient matching loss
 205 together with the regularization term. Afterwards, Ω and \mathbf{X}' are updated via gradient descent. It is
 206 worth noting that the training process for each class can be run in parallel since the graph updates for
 207 one class is independent of another class.

208 **Comparison with DC.** Recall that the gradient matching scheme in DC involves a complex bi-level
209 optimization. If we denote the number of inner-iterations as τ_i and that of outer-iterations as τ_o ,
210 its computational complexity can be $\tau_i \times \tau_o$ of our method. Thus DC is significantly slower than
211 *DosCond*. In addition to speeding up condensation, *DosCond* removes the burden of tuning some
212 hyper-parameters, i.e., the number of iterations for outer/inner optimization and learning rate for
213 updating f_θ , which can save us enormous training time when learning larger synthetic sets.

214 **Comparison with Coreset Methods.** Coreset methods [20, 21] select representative data samples
215 based on some heuristics calculated on the pre-trained embedding. Thus, it requires training the
216 model first. Given the cheap cost on calculating and ranking heuristics, the major computational
217 bottleneck for coreset method is on pre-training the neural network for a certain number of iterations.
218 Likewise, our proposed *DosCond* has comparable complexity because it also needs to forward and
219 backward the neural network for multiple iterations. Thus, their efficiency difference majorly depends
220 on how many epochs we run for learning synthetic graphs in *DosCond* and for pre-training the model
221 embedding in coreset methods. In practice, we find that *DosCond* even requires less training cost
222 than the coreset methods as shown in Section 3.2.

223 3 Experiment

224 3.1 Experimental settings

225 **Datasets.** To evaluate the performance of our method, we use multiple molecular datasets from
226 Open Graph Benchmark (OGB) [22] and TU Datasets (DD, MUTAG and NCI1) [23] for graph-level
227 property classification, and one superpixel dataset CIFAR10 [24]. We also introduce a real-world
228 e-commerce dataset. In particular, we randomly sample 1,109 sub-graphs from a large, anonymized
229 internal knowledge graph. Each sub-graph is created from the ego network of a random selected
230 product on the e-commerce website. We form a binary classification problem aiming at predicting
231 the product category of the central product node in each sub-graph. We use the public splits for OGB
232 datasets and CIFAR10. For TU Datasets and the e-commerce dataset, we randomly split the graphs
233 into 80%/10%/10% for training/validation/test. Detailed dataset statistics are shown in Appendix A.2.

234 **Baselines.** We compare our proposed methods with four baselines that produce discrete structures:
235 three coreset methods (*Random*, *Herding* [20] and *K-Center* [25, 21]), and a *dataset condensation*
236 method DCG [11]: (a) *Random*: it randomly picks graphs from the training dataset. (b) *Herding*: it
237 selects samples that are closest to the cluster center. *Herding* is often used in replay-based methods
238 for continual learning [26, 27]. (c) *K-Center*: it selects the center samples to minimize the largest
239 distance between a sample and its nearest center. (d) *DCG*: As vanilla DC [11] cannot generate
240 discrete structure, we randomly select graphs from training and apply DC to learn the features
241 for them, which we term as DCG. We use the implementations provided by DC [11] for *Herding*,
242 *K-Center* and *DCG*. Note that coreset methods only select existing samples from training while *DCG*
243 learns the node features.

244 **Evaluation Protocol.** To evaluate the effectiveness of the proposed method, we test the classification
245 performance of GNNs trained with condensed graphs on the aforementioned graph datasets. Con-
246 cretely, it involves three stages: (1) learning synthetic graphs, (2) training a GCN on the synthetic
247 graphs and (3) test the performance of GCN. We first generate the condensed graphs following the
248 procedure in Algorithm 1. Then we train a GCN classifier with the condensed graphs. Finally we
249 evaluate its classification performance on the real graphs from test set. For baseline methods, we first
250 get the selected/condensed graphs and then follow the same procedure. We repeat the generation
251 process of condensed graphs 5 times with different random seeds and train GCN on these graphs with
252 10 different random seeds. We report the mean and standard deviation of these results.

253 3.2 Performance with Condensed Graphs

254 **Classification Performance Comparison.** To validate the effectiveness of the proposed framework,
255 we measure the classification performance of GCN trained on condensed graphs. Specifically, we vary
256 the number of learned synthetic graphs per class in the range of $\{1, 10, 50\}$ ($\{1, 10, 20\}$ for MUTAG
257 and E-commerce) and train a GCN on these graphs. Then we evaluate the classification performance
258 of the trained GCN on the original test graphs. Following the convention in OGB [22], we report the
259 ROC-AUC metric for ogbg-molbace, ogbg-molbbbp and ogbg-molhiv; for other datasets we report
260 the classification accuracy (%). The results are summarized in Table 1. Note that the *Ratio* column
261 presents the ratio of synthetic graphs to original graphs and we name it as *condensation ratio*; the

Table 1: The classification performance comparison. We report the ROC-AUC for the first three datasets and accuracies (%) for others. *Whole Dataset* indicates the performance with original dataset.

	Graphs/Cls.	Ratio	Random	Herding	K-Center	DCG	<i>DosCond</i>	Whole Dataset
ogbg-molbace (ROC-AUC)	1	0.2%	0.580±0.067	0.548±0.034	0.548±0.034	0.623±0.046	0.657±0.034	0.714±0.005
	10	1.7%	0.598±0.073	0.639±0.039	0.591±0.056	0.655±0.033	0.674±0.035	
	50	8.3%	0.632±0.047	0.683±0.022	0.589±0.025	0.652±0.013	0.688±0.012	
ogbg-molbbbp (ROC-AUC)	1	0.1%	0.519±0.016	0.546±0.019	0.546±0.019	0.559±0.044	0.581±0.005	0.646±0.004
	10	1.2%	0.586±0.040	0.605±0.019	0.530±0.039	0.568±0.032	0.605±0.008	
	50	6.1%	0.606±0.020	0.617±0.003	0.576±0.019	0.579±0.032	0.620±0.007	
ogbg-molhiv (ROC-AUC)	1	0.01%	0.719±0.009	0.721±0.002	0.721±0.002	0.718±0.013	0.726±0.003	0.757±0.007
	10	0.06%	0.720±0.011	0.725±0.006	0.713±0.009	0.728±0.002	0.728±0.005	
	50	0.3%	0.721±0.014	0.725±0.003	0.725±0.006	0.726±0.010	0.731±0.004	
DD (Accuracy)	1	0.2%	57.69±4.92	61.97±1.32	61.97±1.32	58.81±2.90	70.42±2.21	78.92±0.64
	10	2.1%	64.69±2.55	69.79±2.30	63.46±2.38	61.84±1.44	73.53±1.13	
	50	10.6%	67.29±1.53	73.95±1.70	67.41±0.92	61.27±1.01	77.04±1.86	
MUTAG (Accuracy)	1	1.3%	67.47±9.74	70.84±7.71	70.84±7.71	75.00±8.16	82.21±1.61	88.63±1.44
	10	13.3%	77.89±7.55	80.42±1.89	81.00±2.51	82.66±0.68	82.76±2.31	
	20	26.7%	78.21±5.13	80.00±1.10	82.97±4.91	82.89±1.03	83.26±2.34	
NC11 (Accuracy)	1	0.1%	51.27±1.22	53.98±0.67	53.98±0.67	51.14±1.08	56.58±0.48	71.70±0.20
	10	0.6%	54.33±3.14	57.11±0.56	53.21±1.44	51.86±0.81	58.02±1.05	
	50	3.0%	58.51±1.73	58.94±0.83	56.58±3.08	52.17±1.90	60.07±1.58	
CIFAR10 (Accuracy)	1	0.06%	15.61±0.52	22.38±0.49	22.37±0.50	21.60±0.42	24.70±0.70	50.75±0.14
	10	0.2%	23.07±0.76	28.81±0.35	20.93±0.62	29.27±0.77	30.70±0.23	
	50	1.1%	30.56±0.81	33.94±0.37	24.17±0.51	34.47±0.52	35.34±0.14	
E-commerce (Accuracy)	1	0.2%	51.31±2.89	52.18±0.25	52.36±0.38	57.14±1.72	60.82±1.23	69.25±0.50
	10	0.9%	54.99±2.74	56.83±0.87	56.49±0.36	61.03±1.32	64.73±1.34	
	20	3.6%	57.80±3.58	62.56±0.71	62.76±0.45	64.92±1.35	67.71±1.22	

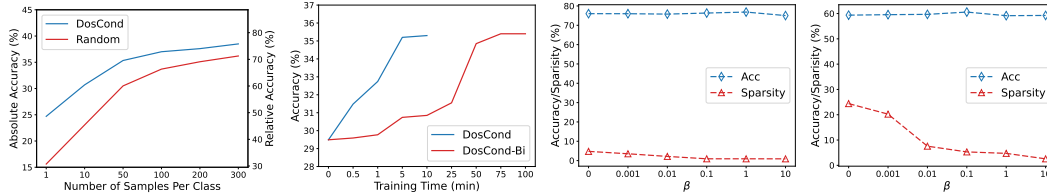
262 *Whole Dataset* column shows the GCN performance achieved by training on the original dataset.
 263 From the table, we make four observations:

- 264 (a) The proposed *DosCond* consistently achieves better performance than the baseline methods under
 265 different condensation ratios and different datasets. Notably, when generating only 2 graphs
 266 on ogbg-molbace dataset (0.2%), we achieve an ROC-AUC of 0.657 while the performance on
 267 full training set is 0.714, which means we approximate 92% of the original performance with
 268 only 0.2% data. Likewise, we are able to approximate 96.5% of the original performance on
 269 ogbg-molhiv with 0.3% data. By contrast, baselines underperform our method by a large margin.
 270 Similar observations can be made on other datasets, which demonstrates the effectiveness of
 271 learned synthetic graphs in preserving the information of the original dataset.
- 272 (b) Increasing the number of synthetic graphs can improve the classification performance. For
 273 example, we can approximate the original performance by 89%/93%/98% with 0.2%/2.1%/10.6%
 274 data on DD. More synthetic samples indicate more learnable parameters that can preserve the
 275 information residing in the original dataset and present more diverse patterns that can help train
 276 GNNs better. This observation is in line with our experimental results in Section 3.3.
- 277 (c) The performance on CIFAR10 is less promising due to the limit number of synthetic graphs.
 278 We posit that the dataset has more complex topology and feature information and thus requires
 279 more parameters to preserve sufficient information. However, we note that our method still
 280 outperforms the baseline methods especially when producing only 1 sample per class, which
 281 suggests that our method is much more data-efficient. Moreover, we are able to promote the
 282 performance on CIFAR10 by learning a larger synthetic set as shown in Section 3.3.
- 283 (d) Learning both synthetic graph structure and node features is necessary for preserving the infor-
 284 mation in original graph datasets. By checking the performance DCG, which only learns node
 285 features based on randomly selected graph structure, we see that DCG underperforms *DosCond*
 286 by a large margin in most cases. This indicates that learning node features solely is sub-optimal.

287 **Efficiency Comparison.** Since one of our goals is to enable scalable dataset condensation, we now
 288 evaluate the efficiency of *DosCond*. We compare *DosCond* with the coreset method Herding, as it
 289 is less time-consuming than DCG and generally achieves better performance than other baselines.
 290 We adopt the same setting as in Table 1: 1000 iterations for *DosCond*, i.e., $K_1 = 1000$, and 500
 291 epochs (100 epochs for ogbg-molhiv) for pre-training the graph convolutional network as required by
 292 Herding. We also note that pre-training the neural network need to go over the whole dataset at every

Table 2: Comparison of running time (minutes).

	CIFAR10		ogbg-molhiv		DD	
G./Cls.	Herding	<i>DosCond</i>	Herding	<i>DosCond</i>	Herding	<i>DosCond</i>
1	44.5m	4.7m	4.3m	0.66m	1.6m	1.5m
10	44.5m	4.9m	4.3m	0.67m	1.6m	1.5m
50	44.5m	5.7m	4.3m	0.68m	1.6m	2.0m



(a) Larger synthetic set. (b) One-Step v.s. bi-Level (c) Varying β on DD (d) Varying β on NCI1
 Figure 1: Algorithm analysis and parameter analysis w.r.t. the sparsity regularization.

epoch while *DosCond* only processes a batch of graphs. In Table 2, we report the running time on an NVIDIA V100 GPU for CIFAR10, ogbg-molhiv and DD. We make three observations:

- (a) *DosCond* can be faster than Herding. In fact, *DosCond* requires less training time in all the cases except in DD with 50 graphs per class. Herding needs to fully training the model on the whole dataset to obtain good-quality embedding, which can be quite time-consuming. On the contrary, *DosCond* only requires matching gradients for K_1 initializations and does not need to fully train the model on the large real dataset.
- (b) The running time of *DosCond* increases with the increase of the number of synthetic graphs N' . It is because *DosCond* processes the condensed graphs at each iteration, of which the time complexity is $O(N'L(n^2d + nd^2))$ for an L -layer GCN. Thus, the additional complexity depends on N' . By contrast, the increase of N' has little impact on Herding since the process of selecting samples based on pre-defined heuristic is very fast.
- (c) The average nodes in synthetic graph n also impacts the training cost of *DosCond*. For instance, the training cost on ogbg-molhiv ($n=26$) is much lower than that on DD ($n=285$), and the gap of cost between the two methods on ogbg-molhiv and DD is very different. As mentioned earlier, it is because the complexity of the forward process in GCN is $O(N'L(n^2d + nd^2))$ for N' condensed graphs with node size of n .

3.3 Further Investigation

Increasing the Number of Synthetic Graphs. We study whether the classification performance can be further boosted when using larger synthetic size. Concretely, we vary the size of the learned graphs from 1 to 300 and report the results of absolute and relative accuracy w.r.t. whole dataset training accuracy for CIFAR10 in Figure 1a. It is clear to see that both Random and *DosCond* achieve better performance when we increase the number of samples used for training. Moreover, our method outperforms the random baseline under different condensed dataset sizes. Note that the performance gap between the two methods diminishes with the increase of the number of samples. This is because the random baseline will finally approach the whole dataset training if we continue to enlarge the size of the condensed set, in which the performance can be viewed as the upper bound of *DosCond*.

Ablation Study. We perform ablation study on the proposed one-step gradient matching and regularization terms. We create an ablation of our method, namely *DosCond-Bi*, which adopts the vanilla gradient matching scheme that involves a bi-level optimization. Without loss of generality, we compare the training time and classification accuracy of *DosCond* and *DosCond-Bi* in the setting of learning 50 graphs/class synthetic graphs on CIFAR10 dataset. The results are summarized in Figure 1b and we can see that *DosCond* needs approximately 5 minutes to reach the performance of *DosCond-Bi* trained for 75 minutes, which indicates that *DosCond* only requires 6.7% training cost. It further demonstrates the efficiency of the proposed one-step gradient matching strategy.

Next we study the effect of sparsity regularization on *DosCond*. Specifically, we vary the sparsity coefficient β in the range of $\{0, 0.001, 0.01, 0.1, 1, 10\}$ and report the classification accuracy and graph sparsity on DD and NCI datasets in Figure 1c and 1d. As shown in the figure, when β gets larger, we exert a stronger regularization on the learned graphs and the graphs become more sparse. Furthermore, the increased sparsity does not affect the classification performance. This is a desired

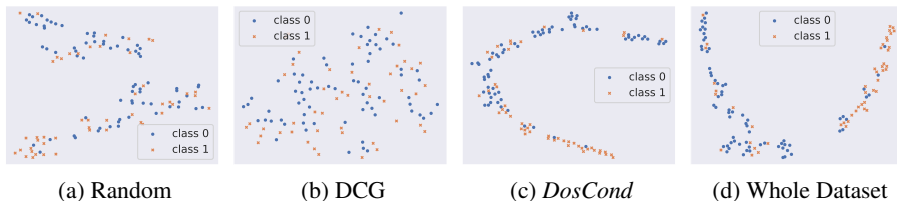


Figure 2: T-SNE visualizations of embedding learned with condensed graphs on DD.

333 property since sparse graphs can save much space for storage and reduce training cost for GNNs.
 334 We also remove the regularization of Eq. (11) for ogbg-molhiv, we obtain the performance of 0.724/
 335 0.727/0.731 for 1/10/50 graphs per class, slightly worse than the one with this regularization.

336 **Visualization.** We further investigate whether GCN can learn discriminative representations from
 337 the synthetic graphs learned by *DosCond*. Specifically, we use t-SNE [28] to visualize the learned
 338 graph representation from GCN trained on different condensed graphs. We train a GCN on graphs
 339 produced by different methods and use it to extract the latent representation for real graphs from
 340 test set. Without loss of generality, we provide the t-SNE plots on DD dataset with 50 graphs per
 341 class in Figure 2. It is observed that the graph representations learned with randomly selected graphs
 342 are mixed for different classes. Similarly, DCG graphs also resulted in poorly trained GCN that
 343 outputs indistinguishable graph representations. By contrast, the representations are well separated
 344 for different classes when learned with *DosCond* graphs (Figure 2c) and they are as discriminative as
 345 those learned on the whole dataset (Figure 2d). This demonstrates that the graphs learned by *DosCond*
 346 preserve sufficient information of the original dataset so as to recover the original performance.

347 4 Related Work

348 **Graph Neural Networks.** As the generalization of deep neural network to graph data, graph neural
 349 networks (GNNs) [4, 29, 5, 7, 19, 30, 31, 32, 33] have revolutionized the field of graph representation
 350 learning through effectively exploiting graph structural information. GNNs have achieved remarkable
 351 performances in basic graph-related tasks such as graph classification [34, 35], link prediction [3]
 352 and node classification [4]. Recent years have also witnessed their great success achieved in many
 353 real-world applications such as recommender systems [3], computer vision [36], drug discovery [37]
 354 and etc. GNNs take both adjacency matrix and node feature matrix as input and output node-level
 355 representations or graph-level representations. Essentially, they follow a message-passing scheme [38]
 356 where each node first aggregates the information from its neighborhood and then transforms the
 357 aggregated information to update its representation.

358 **Dataset Distillation & Dataset Condensation.** It is widely received that training neural networks on
 359 large datasets can be prohibitively costly. To alleviate this issue, dataset distillation (DD) [10] aims to
 360 distill knowledge of a large training dataset into a small number of synthetic samples. DD formulates
 361 the distillation process as a learning-to-learning problem and solves it through bi-level optimization.
 362 To improve the efficiency of DD, dataset condensation (DC) [11, 12] is proposed to learn the small
 363 synthetic dataset by matching the gradients of the network parameters w.r.t. large-real and small-
 364 synthetic training data. It has been demonstrated that these condensed samples can facilitate critical
 365 applications such as continual learning [11, 12, 39, 40, 41], neural architecture search [13, 14, 42]
 366 and privacy-preserving scenarios [43] Recently, following the gradient matching scheme in DC,
 367 *GCond* proposes a condensation method to condense a large-scale graph to a small graph for node
 368 classification. Different from *GCond*, we aim to solve the challenge of learning discrete structure and
 369 we majorly target at graph classification. Our method avoids the costly bi-level optimization and is
 370 much more efficient than the previous work. A detailed comparison is included in Appendix B.2.

371 5 Conclusion

372 Training graph neural networks on a large-scale graph dataset consumes high computational cost. One
 373 solution to alleviate this issue is to condense the large graph dataset into a small synthetic dataset. In
 374 this work, we propose a novel framework *DosCond* that adopts a one-step gradient matching strategy
 375 to efficiently condenses real graphs into a small number of informative graphs with discrete structures.
 376 We further justify the proposed method from both theoretical and empirical perspectives. Notably,
 377 our experiments show that we are able to reduce the dataset size by 90% while approximating up
 378 to 98% of the original performance. In the future, we plan to investigate interpretable condensation
 379 methods and diverse applications of the condensed graphs.

References

- 380
- 381 [1] Zhitao Ying, Jiaxuan You, Christopher Morris, Xiang Ren, William L. Hamilton, and Jure
382 Leskovec. Hierarchical graph representation learning with differentiable pooling. In *NeurIPS*,
383 2018.
- 384 [2] Daixin Wang, Jianbin Lin, Peng Cui, Quanhui Jia, Zhen Wang, Yanming Fang, Quan Yu, Jun
385 Zhou, Shuang Yang, and Yuan Qi. A semi-supervised graph attentive network for financial
386 fraud detection. In *ICDM*, 2019.
- 387 [3] Wenqi Fan, Yao Ma, Qing Li, Yuan He, Yihong Eric Zhao, Jiliang Tang, and Dawei Yin. Graph
388 neural networks for social recommendation. In *WWW*, 2019.
- 389 [4] Thomas N. Kipf and Max Welling. Semi-supervised classification with graph convolutional
390 networks. In *ICLR*, 2017.
- 391 [5] Petar Velickovic, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua
392 Bengio. Graph attention networks. In *ICLR*, 2018.
- 393 [6] Peter W Battaglia, Jessica B Hamrick, Victor Bapst, Alvaro Sanchez-Gonzalez, Vinicius Zam-
394 baldi, Mateusz Malinowski, Andrea Tacchetti, David Raposo, Adam Santoro, Ryan Faulkner,
395 et al. Relational inductive biases, deep learning, and graph networks. *ArXiv preprint*, 2018.
- 396 [7] Zonghan Wu, Shirui Pan, Fengwen Chen, Guodong Long, Chengqi Zhang, and Philip S Yu. A
397 comprehensive survey on graph neural networks. *ArXiv preprint*, 2019.
- 398 [8] Hanxiao Liu, Karen Simonyan, and Yiming Yang. DARTS: differentiable architecture search.
399 In *ICLR*, 2019.
- 400 [9] Zhizhong Li and Derek Hoiem. Learning without forgetting. *IEEE transactions on pattern*
401 *analysis and machine intelligence*, 40(12):2935–2947, 2017.
- 402 [10] Tongzhou Wang, Jun-Yan Zhu, Antonio Torralba, and Alexei A Efros. Dataset distillation.
403 *ArXiv preprint*, 2018.
- 404 [11] Bo Zhao, Konda Reddy Mopuri, and Hakan Bilen. Dataset condensation with gradient matching.
405 In *ICLR*, 2021.
- 406 [12] Bo Zhao and Hakan Bilen. Dataset condensation with differentiable siamese augmentation. In
407 *ICML*, Proceedings of Machine Learning Research, 2021.
- 408 [13] Timothy Nguyen, Zhoung Chen, and Jaehoon Lee. Dataset meta-learning from kernel
409 ridge-regression. In *ICLR*, 2021.
- 410 [14] Timothy Nguyen, Roman Novak, Lechao Xiao, and Jaehoon Lee. Dataset distillation with
411 infinitely wide convolutional networks. *NeurIPS*, 34, 2021.
- 412 [15] George Cazenavette, Tongzhou Wang, Antonio Torralba, Alexei A Efros, and Jun-Yan Zhu.
413 Dataset distillation by matching training trajectories. In *CVPR*, 2022.
- 414 [16] Kai Wang, Bo Zhao, Xiangyu Peng, Zheng Zhu, Shuo Yang, Shuo Wang, Guan Huang, Hakan
415 Bilen, Xinchao Wang, and Yang You. Cafe: Learning to condense dataset by aligning features.
416 In *CVPR*, 2022.
- 417 [17] Chris J Maddison, Andriy Mnih, and Yee Whye Teh. The concrete distribution: A continuous
418 relaxation of discrete random variables. *arXiv preprint arXiv:1611.00712*, 2016.
- 419 [18] Abubakar Abid, Muhammad Fatih Balin, and James Zou. Concrete autoencoders for differen-
420 tiable feature selection and reconstruction. *arXiv preprint arXiv:1901.09346*, 2019.
- 421 [19] Felix Wu, Amauri H. Souza Jr., Tianyi Zhang, Christopher Fifty, Tao Yu, and Kilian Q. Wein-
422 berger. Simplifying graph convolutional networks. In *ICML*, 2019.
- 423 [20] Max Welling. Herding dynamical weights to learn. In *ICML*, 2009.

- 424 [21] Ozan Sener and Silvio Savarese. Active learning for convolutional neural networks: A core-set
425 approach. In *ICLR*, 2018.
- 426 [22] Weihua Hu, Matthias Fey, Marinka Zitnik, Yuxiao Dong, Hongyu Ren, Bowen Liu, Michele
427 Catasta, and Jure Leskovec. Open graph benchmark: Datasets for machine learning on graphs.
428 In *NeurIPS*, 2020.
- 429 [23] Christopher Morris, Nils M Kriege, Franka Bause, Kristian Kersting, Petra Mutzel, and Marion
430 Neumann. Tudataset: A collection of benchmark datasets for learning with graphs. *arXiv*
431 *preprint arXiv:2007.08663*, 2020.
- 432 [24] Vijay Prakash Dwivedi, Chaitanya K Joshi, Thomas Laurent, Yoshua Bengio, and Xavier
433 Bresson. Benchmarking graph neural networks. *arXiv preprint arXiv:2003.00982*, 2020.
- 434 [25] Reza Zanjirani Farahani and Masoud Hekmatfar. *Facility location: concepts, models, algorithms*
435 *and case studies*. 2009.
- 436 [26] Sylvestre-Alvise Rebuffi, Alexander Kolesnikov, Georg Sperl, and Christoph H. Lampert. icarl:
437 Incremental classifier and representation learning. In *CVPR*, 2017.
- 438 [27] Francisco M Castro, Manuel J Marín-Jiménez, Nicolás Guil, Cordelia Schmid, and Karteek
439 Alahari. End-to-end incremental learning. In *ECCV*, 2018.
- 440 [28] Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of machine*
441 *learning research*, (11), 2008.
- 442 [29] Johannes Klicpera, Aleksandar Bojchevski, and Stephan Günnemann. Predict then propagate:
443 Graph neural networks meet personalized pagerank. In *ICLR 2019*, 2019.
- 444 [30] Xianfeng Tang, Yandong Li, Yiwei Sun, Huaxiu Yao, Prasenjit Mitra, and Suhang Wang.
445 Transferring robustness for graph neural network against poisoning attacks. In *WSDM*, 2020.
- 446 [31] Wei Jin, Yao Ma, Xiaorui Liu, Xianfeng Tang, Suhang Wang, and Jiliang Tang. Graph structure
447 learning for robust graph neural networks. In *KDD*, 2020.
- 448 [32] Meng Liu, Youzhi Luo, Kanji Uchino, Koji Maruhashi, and Shuiwang Ji. Generating 3d
449 molecules for target protein binding. In *ICML*, 2022.
- 450 [33] Yu Wang, Yuying Zhao, Yushun Dong, Huiyuan Chen, Jundong Li, and Tyler Derr. Improving
451 fairness in graph neural networks via mitigating sensitive attribute leakage. In *KDD*, 2022.
- 452 [34] Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural
453 networks? In *ICLR*, 2019.
- 454 [35] Zhichun Guo, Chuxu Zhang, Wenhao Yu, John Herr, Olaf Wiest, Meng Jiang, and Nitesh V
455 Chawla. Few-shot graph learning for molecular property prediction. In *Proceedings of the Web*
456 *Conference 2021*, pages 2559–2567, 2021.
- 457 [36] Guohao Li, Matthias Müller, Ali K. Thabet, and Bernard Ghanem. Deepgcns: Can gcns go as
458 deep as cnns? In *ICCV*, 2019.
- 459 [37] David Duvenaud, Dougal Maclaurin, Jorge Aguilera-Iparraguirre, Rafael Gómez-Bombarelli,
460 Timothy Hirzel, Alán Aspuru-Guzik, and Ryan P. Adams. Convolutional networks on graphs
461 for learning molecular fingerprints. In *NeurIPS*, 2015.
- 462 [38] Justin Gilmer, Samuel S. Schoenholz, Patrick F. Riley, Oriol Vinyals, and George E. Dahl.
463 Neural message passing for quantum chemistry. In *ICML*, 2017.
- 464 [39] Jang-Hyun Kim, Jinuk Kim, Seong Joon Oh, Sangdoon Yun, Hwanjun Song, Joonhyun Jeong,
465 Jung-Woo Ha, and Hyun Oh Song. Dataset condensation via efficient synthetic-data parameteri-
466 zation. *arXiv:2205.14959*, 2022.
- 467 [40] Saehyung Lee, Sanghyuk Chun, Sangwon Jung, Sangdoon Yun, and Sungroh Yoon. Dataset
468 condensation with contrastive signals. In *ICML*, 2022.

- 469 [41] Bo Zhao and Hakan Bilen. Dataset condensation with distribution matching. *arXiv preprint*
470 *arXiv:2110.04181*, 2021.
- 471 [42] Shuo Yang, Zeke Xie, Hanyu Peng, Min Xu, Mingming Sun, and Ping Li. Dataset pruning:
472 Reducing training data by examining generalization influence. *arXiv preprint arXiv:2205.09329*,
473 2022.
- 474 [43] Tian Dong, Bo Zhao, and Lingjuan Lyu. Privacy for free: How does dataset condensation help
475 privacy? In *ICML*, 2022.
- 476 [44] Wei Jin, Lingxiao Zhao, Shichang Zhang, Yozen Liu, Jiliang Tang, and Neil Shah. Graph
477 condensation for graph neural networks. In *ICLR 2022*, 2022.
- 478 [45] Hanqing Zeng, Hongkuan Zhou, Ajitesh Srivastava, Rajgopal Kannan, and Viktor K. Prasanna.
479 Graphsaint: Graph sampling based inductive learning method. In *ICLR*, 2020.
- 480 [46] Krishnateja Killamsetty, Durga S, Ganesh Ramakrishnan, Abir De, and Rishabh Iyer. Grad-
481 match: Gradient matching based data subset selection for efficient deep model training. In
482 *ICML*. PMLR, 2021.
- 483 [47] Jeremy Watt, Reza Borhani, and Aggelos K Katsaggelos. *Machine learning refined: Founda-*
484 *tions, algorithms, and applications*. Cambridge University Press, 2020.
- 485 [48] Rahul Yedida, Snehanshu Saha, and Tejas Prashanth. Lipschitzlr: Using theoretically computed
486 adaptive learning rates for fast convergence. *Applied Intelligence*, 51(3):1460–1478, 2021.

487 A Experimental Setup

488 A.1 Parameter Settings.

489 When learning the synthetic graphs, we adopt 3-layer GCN with 128 hidden units as the model for
490 gradient matching. The learning rates for structure and feature parameters are set to 1.0 (0.01 for
491 ogbg-molbace and CIFAR10) and 0.01, respectively. We set K_1 to 1000 and β to 0.1. Additionally,
492 we use mean pooling to obtain graph representation for all datasets except ogbg-molhiv. We use sum
493 pooling for ogbg-molhiv as it achieves better classification performance on the real dataset. During
494 the test stage, we use GCN with the same architecture and we train the model for 500 epochs (100
495 epochs for ogbg-molhiv) with an initial learning rate of 0.001.

496 A.2 Dataset Statistics

497 Dataset statistics for node classification and graph classification are shown in Table 3 and 4, respec-
498 tively.

Table 3: Graph classification dataset statistics.

Dataset	Type	#Classes	#Graphs	Avg. Nodes	Avg. Edges
CIFAR10	Superpixel	10	60,000	117.6	941.07
ogbg-molhiv	Molecule	2	41,127	25.5	54.9
ogbg-molbace	Molecule	2	1,513	34.1	36.9
ogbg-molbbbp	Molecule	2	2,039	24.1	26.0
MUTAG	Molecule	2	188	17.93	19.79
NCI1	Molecule	2	4,110	29.87	32.30
DD	Molecule	2	1,178	284.32	715.66
E-commerce	Transaction	2	1,109	33.7	56.3

Table 4: Node classification dataset statistics.

Dataset	#Nodes	#Edges	#Classes	#Features
Cora	2,708	5,429	7	1,433
Citeseer	3,327	4,732	6	3,703
Pubmed	19,717	44,338	3	500
Arxiv	169,343	1,166,243	40	128
Flickr	89,250	899,756	7	500

Algorithm 1: *DosCond* for Condensing Graphs

-
- 1: **Input:** Training data $\mathcal{T} = (\mathbf{A}, \mathbf{X}, \mathcal{Y})$
 - 2: **Required:** Pre-defined condensed labels \mathcal{Y}' , graph neural network f_θ , temperature τ , desired sparsity ϵ , regularization coefficient β , learning rates η_1, η_2 , number of epochs K_1 .
 - 3: Initialize Ω, \mathbf{X}'
 - 4: **for** $k = 0, \dots, K_1 - 1$ **do**
 - 5: Sample $\theta_0 \sim P_{\theta_0}$
 - 6: Sample $\alpha \sim \text{Uniform}(0, 1)$
 - 7: Compute $\mathbf{A}' = \sigma((\log \alpha - \log(1 - \alpha) + \Omega) / \tau)$
 - 8: **for** $c = 0, \dots, C - 1$ **do**
 - 9: Sample $(\mathbf{A}_c, \mathbf{X}_c, \mathcal{Y}_c) \sim \mathcal{T}$ and $(\mathbf{A}'_c, \mathbf{X}'_c, \mathcal{Y}'_c) \sim \mathcal{S}$
 - 10: Compute $\ell_T = \ell(f_{\theta_0}(\mathbf{A}_c, \mathbf{X}_c), \mathcal{Y}_c)$
 - 11: Compute $\ell_S = \ell(f_{\theta_0}(\mathbf{A}'_c, \mathbf{X}'_c), \mathcal{Y}'_c)$
 - 12: Compute $\ell_{\text{reg}} = \max(\sum_{i,j} \sigma(\Omega_{ij}) - \epsilon, 0)$
 - 13: Update $\Omega \leftarrow \Omega - \eta_1 \nabla_{\Omega} (D(\nabla_{\theta_0} \ell_T, \nabla_{\theta_0} \ell_S) + \beta \ell_{\text{reg}})$
 - 14: Update $\mathbf{X}' \leftarrow \mathbf{X}' - \eta_2 \nabla_{\mathbf{X}'} (D(\nabla_{\theta_0} \ell_T, \nabla_{\theta_0} \ell_S) + \beta \ell_{\text{reg}})$
 - 15: **end for**
 - 16: **end for**
 - 17: **Return:** $(\Omega, \mathbf{X}', \mathcal{Y}')$
-

499 **A.3 Algorithm**

500 We provide the details of our proposed framework in Algorithm 1. Specifically, we sample K_1 model
501 initializations θ_0 to perform one-step gradient matching. Following the convention in DC [11], we
502 match gradients and update synthetic graphs for each class separately in order to make matching
503 easier. For class c , we first retrieve the synthetic graphs of that class, denoted as $(\mathbf{A}'_c, \mathbf{X}'_c, \mathcal{Y}'_c) \sim \mathcal{S}$,
504 and sample a batch of real graphs $(\mathbf{A}_c, \mathbf{X}_c, \mathcal{Y}_c)$. We then forward them to the graph neural network
505 and calculate the one-step gradient matching loss together with the regularization term. Afterwards,
506 Ω and \mathbf{X}' are updated via gradient descent. It is worth noting that the training process for each class
507 can be run in parallel since the graph updates for one class is independent of another class.

508 **B More Experiments**509 **B.1 Scale of the two terms in Eq. (9).**

510 As mentioned earlier in Section 2.3, the scale of the first term is essentially larger than the second
511 term in Eq. (9). We now perform empirical study to verify this statement. Since both terms contain
512 the factor M , we simply drop it and focus on studying $\ell_1 = \sqrt{2} \|\nabla_{\theta} \ell_{\mathcal{T}}(\theta_0) - \nabla_{\theta} \ell_{\mathcal{S}}(\theta_0)\|$ and
513 $\ell_2 = \frac{3}{2\sqrt{T}} \cdot \frac{C-1}{CN'} \sqrt{\sum_i \gamma_i \|\mathbf{1}^\top \mathbf{A}'_{(i)} \mathbf{X}'_{(i)}\|^2}$. Specifically, we set T to 500 and N' to 50, and plot the
514 changes of these two terms during the training process of *DosCond*. The results on DD (with mean
515 pooling) and ogbg-molhiv (with sum pooling) are shown in Figure 3. We can observe that the scale
516 of ℓ_1 is much larger than ℓ_2 at the first few epochs when using mean pooling as shown in Figure 3a.
517 By contrast, ℓ_2 is not negligible when using sum pooling as shown in Figure 3b and it is desired to
518 include it as a regularization term in this case. These observations provide support for our discussion
519 of theoretical analysis in Section 2.3.

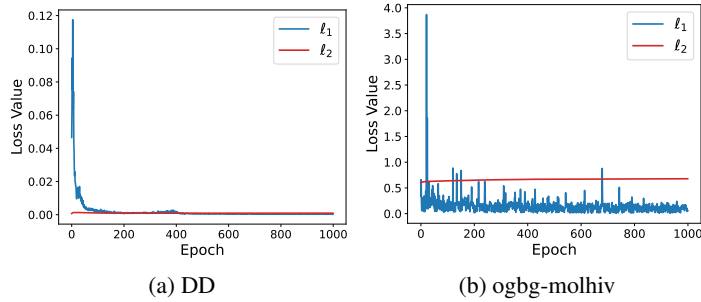


Figure 3: Scale of the two terms in Eq. (11).

Table 5: Node classification accuracy (%) comparison. The numbers in parentheses indicate the running time for 100 epochs and r indicates the ratio of number of nodes in the condensed graph to that in the original graph.

	Cora, $r=2.6\%$	Citeseer, $r=1.8\%$	Pubmed, $r=0.3\%$	Arxiv, $r=0.25\%$	Flickr, $r=0.1\%$
<i>GCond</i>	80.1 (75.9s)	70.6 (71.8s)	77.9 (51.7s)	59.2 (494.3s)	46.5 (51.9s)
<i>DosCond</i>	80.0 (3.5s)	71.0 (2.8s)	76.0 (1.3s)	59.0 (32.9s)	46.1 (14.3s)
Whole Dataset	81.5	71.7	79.3	71.4	47.2

520 B.2 Node Classification

521 Next, we investigate whether the proposed method works well in node classification so as to support
 522 our analysis in Theorem 2 in Appendix C.2. Specifically, following *GCond* [44], a condensation
 523 method for node classification, we use 5 node classification datasets: Cora, Citeseer, Pubmed [4],
 524 ogbn-arxiv [22] and Flickr [45]. The dataset statistics are shown in 4. We follow the settings in
 525 *GCond* to generate one condensed graph for each dataset, train a GCN on the condensed graph,
 526 and evaluate its classification performance on the original test nodes. To adopt *DosCond* into node
 527 classification, we replace the bi-level gradient matching scheme in *GCond* with our proposed one-step
 528 gradient matching. The results of classification accuracy and running time per epoch are summarized
 529 in Table 5. From the table, we make the following observations:

- 530 (a) The proposed *DosCond* achieves similar performance as *GCond* and the performance is also
 531 comparable to the original dataset. For example, we are able to approximate the original training
 532 performance by 99% with only 2.6% data on Cora. It demonstrates the effectiveness of *DosCond*
 533 in the node classification case and justifies Theorem 2 from an empirical perspective.
 534 (b) The training cost of *DosCond* is essentially lower than *GCond* as *DosCond* avoids the expensive
 535 bi-level optimization. By examining their running time, we can see that *DosCond* is up to 40
 536 times faster than *GCond*.

537 We further note that *GCond* produces weighted graphs which require storing the edge weights in
 538 float formats, while *DosCond* outputs discrete graph structure which can be stored as binary values.
 539 Hence, the graphs learned by *DosCond* are more memory-efficient.

540 C Proofs

541 C.1 Proof of Theorem 1

542 Let $\mathbf{A}_{(i)}$, $\mathbf{X}_{(i)}$ denote the adjacency matrix and the feature matrix of i -th real graph, respectively. We
 543 denote the cross entropy loss on the real samples as $\ell_{\mathcal{T}}(\theta) = \sum_i \ell_i(\mathbf{A}_{(i)}, \mathbf{X}_{(i)}, \theta)$ and denote that on
 544 synthetic samples as $\ell_S(\theta) = \ell_S(\mathbf{A}'_{(i)}, \mathbf{X}'_{(i)}, \theta)$. Let θ^* denote the optimal parameter and let θ_t be the
 545 parameter trained on condensed data at t -th epoch by optimizing $\ell_S(\theta)$. For simplicity of notations,
 546 we assume \mathbf{A} and \mathbf{A}' are already normalized. Part of the proof is inspired from the work [46].

547 **Theorem 1** When we use a linearized K -layer SGC as the GNN used in condensation, i.e.,
 548 $f_\theta(\mathbf{A}_{(i)}, \mathbf{X}_{(i)}) = \text{Pool}(\mathbf{A}_{(i)}^K \mathbf{X}_{(i)} \mathbf{W}_1) \mathbf{W}_2$ with $\theta = [\mathbf{W}_1; \mathbf{W}_2]$ and assume that all network pa-
 549 rameters satisfy $\|\theta\|^2 \leq M^2 (M > 0)$, we have

$$\begin{aligned} \min_{t=0,1,\dots,T-1} \ell_{\mathcal{T}}(\theta_t) - \ell_{\mathcal{T}}(\theta^*) &\leq \sum_{t=0}^{T-1} \frac{\sqrt{2}M}{T} \|\nabla_{\theta} \ell_{\mathcal{T}}(\theta_t) - \nabla_{\theta} \ell_S(\theta_t)\| \\ &+ \frac{3M}{2\sqrt{T}} \cdot \frac{C-1}{CN'} \sqrt{\sum_i \gamma_i \|\mathbf{1}^\top \mathbf{A}_{(i)}^K \mathbf{X}'_{(i)}\|^2} \end{aligned} \quad (12)$$

550 where $\gamma_i = 1$ if we use sum pooling in f_θ ; $\gamma_i = \frac{1}{n_i}$ if we use mean pooling, with n_i being the number
 551 of nodes in i -th synthetic graph.

552 We start by proving that $\ell_{\mathcal{T}}(\theta)$ is convex and $\ell_S(\theta)$ is lipschitz continuous when we use
 553 $f_\theta(\mathbf{A}_{(i)}, \mathbf{X}_{(i)}) = \text{Pool}(\mathbf{A}_{(i)}^K \mathbf{X}_{(i)} \mathbf{W}_1) \mathbf{W}_2$ as the mapping function. Before proving these two proper-
 554 ties, we first rewrite $f_\theta(\mathbf{A}_{(i)}, \mathbf{X}_{(i)})$ as:

$$f_\theta(\mathbf{A}_{(i)}, \mathbf{X}_{(i)}) = \begin{cases} \mathbf{1}^\top \mathbf{A}_{(i)}^K \mathbf{X}_{(i)} \mathbf{W}_1 \mathbf{W}_2 & \text{if use sum pooling,} \\ \frac{1}{n_i} \mathbf{1}^\top \mathbf{A}_{(i)}^K \mathbf{X}_{(i)} \mathbf{W}_1 \mathbf{W}_2 & \text{if use mean pooling,} \end{cases} \quad (13)$$

555 where n is the number of nodes in $\mathbf{A}_{(i)}$ and $\mathbf{1}$ is an $n_i \times 1$ matrix filled with constant one. From
 556 the above equation we can see that f_θ with different pooling methods only differ in a multiplication
 557 factor $\frac{1}{n_i}$. Thus, in the following we focus on f_θ with sum pooling to derive the major proof.

558 I. For f_θ with sum pooling:

559 Substitute \mathbf{W} for $\mathbf{W}_1 \mathbf{W}_2$ and we have $f_\theta(\mathbf{A}_{(i)}, \mathbf{X}_{(i)}) = \mathbf{1}^\top \mathbf{A}_{(i)}^K \mathbf{X}_{(i)} \mathbf{W}$ for the case with sum
 560 pooling. Next we show that $\ell_{\mathcal{T}}(\theta)$ is convex and $\ell_S(\theta)$ is lipschitz continuous when we use
 561 $f_\theta(\mathbf{A}_{(i)}, \mathbf{X}_{(i)}) = \mathbf{1}^\top \mathbf{A}_{(i)}^K \mathbf{X}_{(i)} \mathbf{W}$ with $\theta = \mathbf{W}$.

562 (a) Convexity of $\ell_{\mathcal{T}}(\theta)$. From chapter 4 of the book [47], we know that softmax classification
 563 $f(\mathbf{W}) = \mathbf{X}\mathbf{W}$ with cross entropy loss is convex w.r.t. the parameters \mathbf{W} . In our case, the mapping
 564 function $f_\theta(\mathbf{A}_{(i)}, \mathbf{X}_{(i)}) = \mathbf{1}^\top \mathbf{A}_{(i)}^K \mathbf{X}_{(i)} \mathbf{W}$ applies an affine function on $\mathbf{X}\mathbf{W}$. Given that applying
 565 affine function does not change the convexity, we know that $\ell_{\mathcal{T}}(\theta)$ is convex.

566 (b) Lipschitz continuity of $\ell_S(\theta)$. In [48], it shows that the lipschitz constant of softmax regression
 567 with cross entropy loss is $\frac{C-1}{Cm} \|\mathbf{X}\|$, where \mathbf{X} is the input feature matrix, C is the number of classes
 568 and m is the number of samples. Since $\ell_S(\theta)$ is cross entropy loss and f_θ is linear, we know that the
 569 f_θ is lipschitz continuous and it satisfies:

$$\nabla_{\theta} \ell_S(\theta) \leq \frac{C-1}{CN'} \sqrt{\sum_i \|\mathbf{1}^\top \mathbf{A}_{(i)}^K \mathbf{X}'_{(i)}\|^2} \quad (14)$$

570 With (a) and (b), we are able to proceed our proof. First, from the convexity of $\ell_{\mathcal{T}}(\theta)$ we have

$$\ell_{\mathcal{T}}(\theta_t) - \ell_{\mathcal{T}}(\theta^*) \leq \nabla_{\theta} \ell_{\mathcal{T}}(\theta_t)^T (\theta_t - \theta^*) \quad (15)$$

571 We can rewrite $\nabla_{\theta} \ell_{\mathcal{T}}(\theta_t)^T (\theta_t - \theta^*)$ as follows:

$$\begin{aligned} \nabla_{\theta} \ell_{\mathcal{T}}(\theta_t)^T (\theta_t - \theta^*) &= (\nabla_{\theta} \ell_{\mathcal{T}}(\theta_t)^T - \nabla_{\theta} \ell_S(\theta_t)^T + \nabla_{\theta} \ell_S(\theta_t)^T) (\theta_t - \theta^*) \\ &= (\nabla_{\theta} \ell_{\mathcal{T}}(\theta_t)^T - \nabla_{\theta} \ell_S(\theta_t)^T) (\theta_t - \theta^*) + \nabla_{\theta} \ell_S(\theta_t)^T (\theta_t - \theta^*) \end{aligned} \quad (16)$$

572 Given that we use gradient descent to update network parameters, we have $\nabla_{\theta} \ell_S(\theta_t) = \frac{1}{\eta} (\theta_t - \theta_{t+1})$
 573 where η is the learning rate. Then we have,

$$\begin{aligned} \nabla_{\theta} \ell_S(\theta_t)^T (\theta_t - \theta^*) &= \frac{1}{\eta} (\theta_t - \theta_{t+1})^T (\theta_t - \theta^*) \\ &= \frac{1}{2\eta} \left(\|\theta_t - \theta_{t+1}\|^2 + \|\theta_t - \theta^*\|^2 - \|\theta_{t+1} - \theta^*\|^2 \right) \\ &= \frac{1}{2\eta} \left(\|\eta \nabla_{\theta} \ell_S(\theta_t)\|^2 + \|\theta_t - \theta^*\|^2 - \|\theta_{t+1} - \theta^*\|^2 \right) \end{aligned} \quad (17)$$

574 Combining Eq. (15) and Eq. (17) we have,

$$\begin{aligned} \ell_{\mathcal{T}}(\theta_t) - \ell_{\mathcal{T}}(\theta^*) &\leq (\nabla_{\theta} \ell_{\mathcal{T}}(\theta_t)^T - \nabla_{\theta} \ell_{\mathcal{S}}(\theta_t)^T)(\theta_t - \theta^*) \\ &\quad + \frac{1}{2\eta} \left(\|\eta \nabla_{\theta} \ell_{\mathcal{S}}(\theta_t)\|^2 + \|\theta_t - \theta^*\|^2 - \|\theta_{t+1} - \theta^*\|^2 \right) \end{aligned} \quad (18)$$

575 We sum up the two sides of the above inequality for different values of $t \in [0, T-1]$:

$$\begin{aligned} \sum_{t=0}^{T-1} \ell_{\mathcal{T}}(\theta_t) - \ell_{\mathcal{T}}(\theta^*) &\leq \sum_{t=0}^{T-1} (\nabla_{\theta} \ell_{\mathcal{T}}(\theta_t)^T - \nabla_{\theta} \ell_{\mathcal{S}}(\theta_t)^T)(\theta_t - \theta^*) \\ &\quad + \frac{1}{2\eta} \sum_{t=0}^{T-1} \|\eta \nabla_{\theta} \ell_{\mathcal{S}}(\theta_t)\|^2 + \frac{1}{2\eta} \|\theta_0 - \theta^*\|^2 - \frac{1}{2\eta} \|\theta_T - \theta^*\|^2 \end{aligned} \quad (19)$$

576 Since $\frac{1}{2\eta} \|\theta_T - \theta^*\|^2 \geq 0$, we have

$$\begin{aligned} \sum_{t=0}^{T-1} \ell_{\mathcal{T}}(\theta_t) - \ell_{\mathcal{T}}(\theta^*) &\leq \sum_{t=0}^{T-1} (\nabla_{\theta} \ell_{\mathcal{T}}(\theta_t)^T - \nabla_{\theta} \ell_{\mathcal{S}}(\theta_t)^T)(\theta_t - \theta^*) \\ &\quad + \frac{1}{2\eta} \sum_{t=0}^{T-1} \|\eta \nabla_{\theta} \ell_{\mathcal{S}}(\theta_t)\|^2 + \frac{1}{2\eta} \|\theta_0 - \theta^*\|^2 \end{aligned} \quad (20)$$

577 As we assume that $\|\theta\|^2 \leq M^2$, we have $\|\theta - \theta^*\|^2 \leq 2\|\theta\|^2 = 2M^2$. Then Eq. (20) can be rewritten
578 as,

$$\begin{aligned} \sum_{t=0}^{T-1} \ell_{\mathcal{T}}(\theta_t) - \ell_{\mathcal{T}}(\theta^*) &\leq \sum_{t=0}^{T-1} \sqrt{2}M \|\nabla_{\theta} \ell_{\mathcal{T}}(\theta_t) - \nabla_{\theta} \ell_{\mathcal{S}}(\theta_t)\| \\ &\quad + \frac{1}{2\eta} \sum_{t=0}^{T-1} \|\eta \nabla_{\theta} \ell_{\mathcal{S}}(\theta_t)\|^2 + \frac{M^2}{\eta} \end{aligned} \quad (21)$$

579 Recall that $\ell_{\mathcal{S}}(\theta)$ is lipschitz continuous as shown in Eq. (14), and combine

580 $\min_{t=0,1,\dots,T-1} (\ell_{\mathcal{T}}(\theta_t) - \ell_{\mathcal{T}}(\theta^*)) \leq \frac{\sum_{t=0}^{T-1} \ell_{\mathcal{T}}(\theta_t) - \ell_{\mathcal{T}}(\theta^*)}{T}$:

$$\begin{aligned} \min_{t=0,1,\dots,T-1} \ell_{\mathcal{T}}(\theta_t) - \ell_{\mathcal{T}}(\theta^*) &\leq \sum_{t=0}^{T-1} \frac{\sqrt{2}M}{T} \|\nabla_{\theta} \ell_{\mathcal{T}}(\theta_t) - \nabla_{\theta} \ell_{\mathcal{S}}(\theta_t)\| \\ &\quad + \frac{\eta(C-1)^2}{2C^2N'^2} \sum_i \|\mathbf{1}^\top \mathbf{A}'_{(i)} \mathbf{X}'_{(i)}\|^2 + \frac{M^2}{T\eta} \end{aligned} \quad (22)$$

581 Then we choose $\eta = \frac{M}{\sqrt{T} \sqrt{\sum_i \|\mathbf{1}^\top \mathbf{A}'_{(i)} \mathbf{X}'_{(i)}\|^2}}$ and we can get:

$$\begin{aligned} \min_{t=0,1,\dots,T-1} \ell_{\mathcal{T}}(\theta_t) - \ell_{\mathcal{T}}(\theta^*) &\leq \sum_{t=0}^{T-1} \frac{\sqrt{2}M}{T} \|\nabla_{\theta} \ell_{\mathcal{T}}(\theta_t) - \nabla_{\theta} \ell_{\mathcal{S}}(\theta_t)\| \\ &\quad + \frac{3M}{2\sqrt{T}} \cdot \frac{C-1}{CN'} \sqrt{\sum_i \|\mathbf{1}^\top \mathbf{A}'_{(i)} \mathbf{X}'_{(i)}\|^2} \end{aligned} \quad (23)$$

582 II. For f_{θ} with mean pooling:

583 Following similar derivation as in the case of sum pooling, we have

$$\begin{aligned} \min_{t=0,1,\dots,T-1} \ell_{\mathcal{T}}(\theta_t) - \ell_{\mathcal{T}}(\theta^*) &\leq \sum_{t=0}^{T-1} \frac{\sqrt{2}M}{T} \|\nabla_{\theta} \ell_{\mathcal{T}}(\theta_t) - \nabla_{\theta} \ell_{\mathcal{S}}(\theta_t)\| \\ &\quad + \frac{3M}{2\sqrt{T}} \cdot \frac{C-1}{CN'} \sqrt{\sum_i \frac{1}{n_i} \|\mathbf{1}^\top \mathbf{A}'_{(i)} \mathbf{X}'_{(i)}\|^2} \end{aligned} \quad (24)$$

584 where n_i is the number of nodes in i -th synthetic graph.

585 **C.2 Theorem for Node Classification Case**

586 We adopt similar notations for representing the data in node classification but note that there is only
 587 one graph for node classification task. Let $\mathbf{A} \in \{0, 1\}^{N \times N}$, $\mathbf{A}' \in \{0, 1\}^{N' \times N'}$ denote the adjacency
 588 matrix for real graph and synthetic graph, respectively. Let $\mathbf{X} \in R^{N \times d}$, $\mathbf{X}' \in R^{N' \times d}$ denote the
 589 feature matrix for real graph and synthetic graph, respectively. We denote the cross entropy loss on
 590 the real samples as $\ell_{\mathcal{T}}(\theta)$ and denote that on synthetic samples as $\ell_S(\theta)$.

591 **Theorem 2** When we use a K -layer SGC as the model used in condensation, i.e., $f_{\theta}(\mathbf{A}, \mathbf{X}, \theta) =$
 592 $\mathbf{A}^K \mathbf{X} \mathbf{W}$ with $\theta = \mathbf{W}$ and assume that all network parameters satisfy $\|\theta\|^2 \leq M^2 (M > 0)$, we
 593 have

$$\begin{aligned} \min_{t=0,1,\dots,T-1} \ell_{\mathcal{T}}(\theta_t) - \ell_{\mathcal{T}}(\theta^*) &\leq \sum_{t=0}^{T-1} \frac{\sqrt{2}M}{T} \|\nabla_{\theta} \ell_{\mathcal{T}}(\theta_t) - \nabla_{\theta} \ell_S(\theta_t)\| \\ &+ \frac{3M}{2\sqrt{T}} \cdot \frac{C-1}{CN'} \|\mathbf{A}'^K \mathbf{X}'\| \end{aligned} \quad (25)$$

594 We start by proving that $\ell_{\mathcal{T}}(\theta)$ is convex and $\ell_S(\theta)$ is lipschitz continuous when $f_{\theta}(\mathbf{A}, \mathbf{X}, \theta) =$
 595 $\mathbf{A}^K \mathbf{X} \mathbf{W}$.

596 (a) Convexity of $\ell_{\mathcal{T}}(\theta)$: Similar to the graph classification case, the Hessian matrix of $\ell_{\mathcal{T}}(\theta)$ in node
 597 classification is positive semidefinite and thus $\ell_{\mathcal{T}}(\theta)$ is convex.

598 (b) Lipschitz continuity of $\ell_S(\theta)$: As shown in [48], the lipschitz constant of softmax regression
 599 with cross entropy loss is $\frac{C-1}{Cm} \|\mathbf{X}\|$ with C being the number of classes and m being the number
 600 of samples. Thus, we know that the lipschitz constant of $\ell_S(\theta)$ is $\frac{C-1}{CN'} \|\mathbf{A}'^K \mathbf{X}'\|$, which indicates
 601 $\nabla_{\theta} \ell_S(\theta) \leq \frac{C-1}{CN'} \|\mathbf{A}'^K \mathbf{X}'\|$.

602 From the convexity of $\ell_{\mathcal{T}}(\theta)$, we still have the following inequality (see Eq. (21)). Then
 603 recall that $\ell_S(\theta)$ is lipschitz continuous and $\nabla_{\theta} \ell_S(\theta) \leq \frac{C-1}{CN'} \|\mathbf{A}'^K \mathbf{X}'\|$, and combine
 604 $\min_t (\ell_{\mathcal{T}}(\theta_t) - \ell_{\mathcal{T}}(\theta^*)) \leq \frac{\sum_{t=0}^{T-1} \ell_{\mathcal{T}}(\theta_t) - \ell_{\mathcal{T}}(\theta^*)}{T}$:

$$\begin{aligned} \min_{t=0,1,\dots,T-1} \ell_{\mathcal{T}}(\theta_t) - \ell_{\mathcal{T}}(\theta^*) &\leq \sum_{t=0}^{T-1} \frac{\sqrt{2}M}{T} \|\nabla_{\theta} \ell_{\mathcal{T}}(\theta_t) - \nabla_{\theta} \ell_S(\theta_t)\| \\ &+ \frac{\eta(C-1)^2}{2C^2 N'^2} \|\mathbf{A}'^K \mathbf{X}'\|^2 + \frac{M^2}{T\eta} \end{aligned} \quad (26)$$

605 Then we choose $\eta = \frac{M}{\sqrt{T} \|\mathbf{A}'^K \mathbf{X}'\|}$ and we can get:

$$\begin{aligned} \min_{t=0,1,\dots,T-1} \ell_{\mathcal{T}}(\theta_t) - \ell_{\mathcal{T}}(\theta^*) &\leq \sum_{t=0}^{T-1} \frac{\sqrt{2}M}{T} \|\nabla_{\theta} \ell_{\mathcal{T}}(\theta_t) - \nabla_{\theta} \ell_S(\theta_t)\| \\ &+ \frac{3M}{2\sqrt{T}} \cdot \frac{C-1}{CN'} \|\mathbf{A}'^K \mathbf{X}'\| \end{aligned} \quad (27)$$

606