

A Other related work

Imbalanced Learning in Vision Domain. Most real-world data is naturally imbalanced, presenting a significant challenge in training fair models that are not biased towards majority classes. To address this problem, various approaches have been commonly utilized. Ensemble learning [7, 18, 49, 38, 20, 3] combines the outputs of multiple weak classifiers. Data re-sampling methods [5, 10, 29, 27, 14, 36] aim to adjust the label distribution in the training set by synthesizing or duplicating samples from the minority class. Another approach tackles the imbalance issue by modifying the loss function, assigning larger weights to minority classes or adjusting the margins between different classes [50, 32, 4, 31, 40, 26, 37]. Post-hoc correction methods compensate for the imbalanced classes during the inference step, after completing the model training [14, 33, 21, 12]. Although these techniques have been extensively applied to i.i.d. data, extending them to graph-structured data poses non-trivial challenges.

Graph Contrastive Learning. Contrastive methods, which have proven effective for unsupervised learning in vision, have also been adapted for graph data. One notable approach is DGI [35], which presents a framework for unsupervised node-level representation learning that maximizes global mutual information. Other approaches, such as GRACE [51], GCA [52], and GraphCL [45], utilize augmented graphs to optimize the similarity between positive node pairs and minimize negative pairs. CCA-SSG [46] introduces an efficient loss function based on canonical correlation analysis, eliminating the need for negative samples. Incorporating community information, gCool [17] enhances node representations and downstream task performance. GGD [48] simplifies the mutual information loss function by directly discriminating between two sets of node samples, resulting in faster computation and lower memory usage. These contrastive methods exhibit potential for improving unsupervised learning on graph data. However, it is important to note that our model operates within the context of semi-supervised learning, which significantly differs from the mechanisms employed by these models.

464 B Proofs

465 B.1 Proofs of Theorem 1

466 **Theorem 1** Under the condition that $\sum_{i=1}^c n_i$ is a constant, the variance
 467 $\sum_{i=1}^c \mathbb{E}_x \left[\frac{1}{n_i} h^T(x) \Lambda^i h(x) \right]$ reach its minimum when all n_i equal.

468 *Proof* The expression $\sum_{i=1}^c \mathbb{E}_x \left[\frac{1}{n_i} h^T(x) \Lambda^i h(x) \right]$ can be equivalently expressed as
 469 $\sum_{i=1}^c \frac{1}{n_i} \mathbb{E}_x [h^T(x) \Lambda^i h(x)]$. As previously assumed, the $\mathbb{E}_x [h^T(x) \Lambda^i h(x)]$ is the same for dif-
 470 ferent i , which implies that our goal is to demonstrate that $\sum_{i=1}^c \frac{1}{n_i}$ is minimized when all n_i are
 471 equal.

472 Let m be $\sum_{i=1}^c n_i$. We wish to find the extremum of the sum of their reciprocals, which is given by

$$S = \frac{1}{n_1} + \frac{1}{n_2} + \cdots + \frac{1}{n_c}. \quad (11)$$

473 Using the inequality of arithmetic and harmonic means, we have

$$\frac{c}{\frac{1}{n_1} + \frac{1}{n_2} + \cdots + \frac{1}{n_c}} \leq \frac{n_1 + n_2 + \cdots + n_c}{c}, \quad (12)$$

474 with equality if and only if $n_1 = n_2 = \cdots = n_c$. Rearranging, we get

$$\frac{c^2}{n_1 + n_2 + \cdots + n_c} \leq S, \quad (13)$$

475 with equality if and only if $n_1 = n_2 = \cdots = n_c$. Since $m = n_1 + n_2 + \cdots + n_c$, we have

$$\frac{c^2}{m} \leq S, \quad (14)$$

476 with equality if and only if $n_1 = n_2 = \cdots = n_c = \frac{m}{c}$. Therefore, when all the c numbers are equal,
 477 the sum of their reciprocals is minimized and given by $S = \frac{c^2}{m}$.

478 We can also use the method of Lagrange multipliers. Let $f(n_1, n_2, \dots, n_c) = \frac{1}{n_1} + \frac{1}{n_2} + \cdots + \frac{1}{n_c}$
 479 be the function that we want to extremize. Then, the Lagrangian is:

$$\mathcal{L}(n_1, n_2, \dots, n_c, \lambda) = f(n_1, n_2, \dots, n_c) + \lambda(n_1 + n_2 + \cdots + n_c - m). \quad (15)$$

480 Taking the partial derivatives of \mathcal{L} with respect to a_i and λ , we get:

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial n_i} &= -\frac{1}{n_i^2} + \lambda \\ \frac{\partial \mathcal{L}}{\partial \lambda} &= n_1 + n_2 + \cdots + n_c - m. \end{aligned} \quad (16)$$

481 Setting these partial derivatives to zero, we get:

$$n_1 = n_2 = \cdots = n_c = \frac{m}{c}, \lambda = \frac{c^2}{m^2}. \quad (17)$$

482 Thus, when the c numbers are equal, their reciprocal sum is minimized and is equal to $\frac{c^2}{m}$. Moreover,
 483 since S is a continuously differentiable function of n_1, n_2, \dots, n_c , this extremum is a minimum.

484 Therefore, we have shown that the reciprocal sum of c numbers with a sum of m is minimized and
 485 equal to $\frac{c^2}{m}$ when the c numbers are equal.

486 □

487 **B.2 Proofs of Lemma 1 and More Details for Estimating the Expectation with Unlabeled**
 488 **Nodes (Section 3.2)**

489 **B.2.1 Proofs of Lemma 1**

490 **Lemma 1** Under the above assumption for $h^i \sim N(\mu^i, \Lambda^i)$, $C^i \sim N\left(\mu^i, \frac{1}{n_i}\Lambda^i\right)$, minimizing the
 491 $\sum_{i=1}^c \mathbb{E}_x [\text{Var}(x)]$ is equivalent to minimizing Equation 18:

$$\frac{1}{N} \sum_{x \in G} \sum_{i=1}^c \left(h(x)^T \frac{1}{\sqrt{2n_i}} (h_1^i - h_2^i) \right)^2 = \frac{1}{N} \sum_{k=1}^{n_j} \sum_{j=1}^c \sum_{i=1}^c \left((\mu_k^j + \epsilon_k^j)^T \frac{1}{\sqrt{2n_i}} (\epsilon_1^i - \epsilon_2^i) \right)^2. \quad (18)$$

492 *Proof* Let ϵ^i denote $h^i - \mu^i$, which follows the distribution $\epsilon^i \sim N(0, \Lambda^i)$. Similarly, we denote
 493 $e^i = C^i - \mu^i$ and it follows that $e^i \sim N\left(0, \frac{1}{n_i}\Lambda^i\right)$.

494 We know

$$\begin{aligned} \text{Var}(x) &= \sum_{i=1}^c \mathbb{E}_D \left[\left(h^T(x) C^i - \mathbb{E}_D [h^T(x) C^i] \right)^2 \right] \\ &= \sum_{i=1}^c \mathbb{E}_{e^i} \left[\left(h^T(x) (\mu^i + e^i) - h^T(x) \mu^i \right)^2 \right] \\ &= \sum_{i=1}^c \mathbb{E}_{e^i} \left[\left(h^T(x) e^i \right)^2 \right] = \sum_{i=1}^c \frac{1}{n_i} h^T(x) \Lambda^i h(x), \end{aligned} \quad (19)$$

495 so we get

$$\sum_{i=1}^c \mathbb{E}_x [\text{Var}(x)] = \sum_{i=1}^c \mathbb{E}_x \left[\frac{1}{n_i} h^T(x) \Lambda^i h(x) \right]. \quad (20)$$

496 Motivated by stochastic gradient descent, we opt to sample an e^i on each occasion and compute, as
 497 opposed to directly calculating the expectation $\mathbb{E}_{D \subset G}$.

498 At present, we remain uncertain about how to sample the noise term e^i associated with the class
 499 center C_i . We put forth a proposition to estimate the class center noise e utilizing the feature noise ϵ .
 500 Under the supposition that when $v \in C_k$, the j_{th} element of the feature $f(v)$ adheres to a Gaussian
 501 distribution:

$$\epsilon^i \sim N(0, \Lambda^i), \quad e^i \sim N\left(0, \frac{1}{n_i}\Lambda^i\right). \quad (21)$$

502 Consequently, multiplying the noise term ϵ^i in Equation 21 by $\frac{1}{\sqrt{n_i}}$ yields a random variable that
 503 exhibits an identical distribution to e^i :

$$\frac{1}{\sqrt{n_i}} \epsilon^i \sim N\left(0, \frac{1}{n_i}\Lambda^i\right) \quad (22)$$

504 Moreover, how might we compute ϵ_i ? In practical scenarios, we merely possess the feature h .
 505 However, we are able to calculate the disparity between two features, h_1^i and h_2^i , originating from the
 506 same class i :

$$\frac{1}{\sqrt{2n_i}} (h_1^i - h_2^i) = \frac{1}{\sqrt{2n_i}} (\epsilon_1 - \epsilon_2) \sim N\left(0, \frac{1}{n_i}\Lambda^i\right). \quad (23)$$

507 Incorporating this into Equation 20, the loss is expressed as:

$$\frac{1}{N} \sum_{x \in G} \sum_{i=1}^c \left(h(x)^T \frac{1}{\sqrt{2n_i}} (h_1^i - h_2^i) \right)^2. \quad (24)$$

508 \square

509 B.2.2 More Details for Estimating the Expectation with Unlabeled Nodes (Section 3.2)

510 Lemma 1 suggests that $\sum_{i=1}^c \mathbb{E}_x [\text{Var}(x)]$ can be estimated by sampling labeled node pairs (h_1^i, h_2^i)
 511 from the same class i . However, due to the scarcity of data in the minority class, this can often
 512 be challenging. In Section 3.2, we address this issue by utilizing random data augmentation and
 513 defining nodes from different views as members of the same class, enabling us to sample node pairs
 514 for Equation 3 and perform embedding subtraction.

515 However, as we mentioned, the pseudo node pairs lack information regarding their class membership,
 516 making it impossible to assign appropriate $\frac{1}{\sqrt{2n_i}}$ values for different i in Equation 3. These coefficients
 517 play a crucial role in compensating for the minority class, making it imperative to reintroduce
 518 information about the class number when constructing variance regularization. Therefore, in the
 519 second step, we replace h in Equation 3 with the class center C^i , as demonstrated in Equation 4.
 520 As $C^i = \mu^i + e^i$ and the variance of e^i is proportional to $\frac{1}{n_i}$, using Equation 4 to estimate Lemma
 521 1 can provide similar compensatory effects for the minority class. We present the proof for the
 522 aforementioned propositions below.

523 *Proof* Firstly, the $\sum_{i=1}^c \mathbb{E}_x [\text{Var}(x)]$ can be decomposed as follows:

$$\begin{aligned}
 \sum_{i=1}^c \mathbb{E}_x [\text{Var}(x)] &= \frac{1}{N} \sum_{x \in G} \sum_{i=1}^c \left(h(x)^T \frac{1}{\sqrt{2n_i}} (h_1^i - h_2^i) \right)^2 \\
 &= \frac{1}{N} \sum_{k=1}^{n_j} \sum_{j=1}^c \sum_{i=1}^c \left((\mu_k^j + \epsilon_k^j)^T \frac{1}{\sqrt{2n_i}} (h_1^i - h_2^i) \right)^2 \\
 &= \frac{1}{N} \sum_{k=1}^{n_j} \sum_{j=1}^c \sum_{i=1}^c \left((\mu_k^j)^T \frac{1}{\sqrt{2n_i}} (\epsilon_1^i - \epsilon_2^i) + (\epsilon_k^j)^T \frac{1}{\sqrt{2n_i}} (\epsilon_1^i - \epsilon_2^i) \right)^2 \\
 &= \frac{1}{2N} \sum_{k=1}^{n_j} \sum_{j=1}^c \sum_{i=1}^c \left((\mu_k^j)^T \frac{1}{\sqrt{n_i}} (\epsilon_1^i - \epsilon_2^i) + (\epsilon_k^j)^T \frac{1}{\sqrt{n_i}} (\epsilon_1^i - \epsilon_2^i) \right)^2 \\
 &= \frac{1}{2N} \sum_{k=1}^{n_j} \sum_{j=1}^c \sum_{i=1}^c \left(\frac{1}{n_i} [(\mu_k^j)^T (\epsilon_1^i - \epsilon_2^i)]^2 \right) + \frac{1}{2N} \sum_{k=1}^{n_j} \sum_{j=1}^c \sum_{i=1}^c \left(\frac{1}{n_i} [(\epsilon_k^j)^T (\epsilon_1^i - \epsilon_2^i)]^2 \right) \\
 &\quad + \frac{1}{N} \sum_{k=1}^{n_j} \sum_{j=1}^c \sum_{i=1}^c \left(\frac{1}{n_i} [(\mu_k^j)^T (\epsilon_1^i - \epsilon_2^i) (\epsilon_k^j)^T (\epsilon_1^i - \epsilon_2^i)] \right)
 \end{aligned} \tag{25}$$

524 The above Equation can be decomposed into three parts, namely \mathcal{T}_1 , \mathcal{T}_2 , and \mathcal{T}_3 . Notably, each of
 525 these parts is associated with weight $\frac{1}{n_i}$. This observation supports Theorem 1 and Lemma 1, which
 526 suggest that the variance of a model is highly dependent on the distribution of the dataset's sample
 527 size, and that the extent of sample imbalance can significantly increase the model's variance.

528 Note that $\epsilon^i \sim N(0, \Lambda^i)$, $e^i \sim N(0, \frac{1}{n_i} \Lambda^i)$, and so $\frac{1}{\sqrt{n_i}} \epsilon^i \sim N(0, \frac{1}{n_i} \Lambda^i)$. So, if we use the class
 529 center C^i to replace h in Equation 3, then we can get the following expression,

$$\begin{aligned}
& \frac{1}{N} \sum_{x \in G} \sum_{i=1}^c \left((C^i)^T h(x) - (C^{i'})^T h'(x) \right)^2 \\
&= \frac{1}{N} \sum_{k=1}^{n_j} \sum_{j=1}^c \sum_{i=1}^c \left((\mu^i)^T (\epsilon_k^j - \epsilon_k^{j'}) + (e^i)^T \epsilon_k^j - (e^{i'})^T \epsilon_k^{j'} \right)^2 \\
&= \frac{1}{N} \sum_{k=1}^{n_j} \sum_{j=1}^c \sum_{i=1}^c \left((\mu^i)^T (\epsilon_k^j - \epsilon_k^{j'}) + \left(\frac{1}{\sqrt{n_i}} \epsilon^i \right)^T \epsilon_k^j - \left(\frac{1}{\sqrt{n_i}} \epsilon^i \right)^T \epsilon_k^{j'} \right)^2 \\
&= \frac{1}{N} \sum_{k=1}^{n_i} \sum_{i=1}^c \sum_{j=1}^c \left((\mu^j)^T (\epsilon_k^i - \epsilon_k^{i'}) + \left(\frac{1}{\sqrt{n_j}} \epsilon^j \right)^T \epsilon_k^i - \left(\frac{1}{\sqrt{n_j}} \epsilon^j \right)^T \epsilon_k^{i'} \right)^2 \\
&= \frac{1}{N} \sum_{k=1}^{n_i} \sum_{i=1}^c \sum_{j=1}^c \left((\mu^j)^T (\epsilon_k^i - \epsilon_k^{i'}) + \left(\frac{1}{\sqrt{n_j}} \epsilon^j \right)^T (\epsilon_k^i - \epsilon_k^{i'}) \right)^2 \\
&= \frac{1}{N} \sum_{k=1}^{n_i} \sum_{i=1}^c \sum_{j=1}^c \left(\underbrace{\left[(\mu^j)^T (\epsilon_k^i - \epsilon_k^{i'}) \right]^2}_{\mathcal{L}_1} + \underbrace{\frac{1}{n_j} \left[(\epsilon^j)^T (\epsilon_k^i - \epsilon_k^{i'}) \right]^2}_{\mathcal{L}_2} \right. \\
&\quad \left. + \underbrace{\frac{2}{\sqrt{n_j}} \left[(\mu^j)^T (\epsilon_k^i - \epsilon_k^{i'}) \right] (\epsilon^j)^T (\epsilon_k^i - \epsilon_k^{i'})}_{\mathcal{L}_3} \right)
\end{aligned} \tag{26}$$

530 Like the previous Equation 25, Equation 26 can be decomposed into three parts: \mathcal{L}_1 , \mathcal{L}_2 , and \mathcal{L}_3 . If
 531 we substitute the class center C^i for h , we can observe that even though \mathcal{L}_1 is insensitive to class
 532 imbalance, variables \mathcal{L}_2 and \mathcal{L}_3 , which respectively incorporate weights $\frac{1}{n_j}$ and $\frac{1}{\sqrt{n_j}}$, can still provide
 533 the following support: when optimizing the variance of the model, more attention can be given to
 534 the variance introduced by the minority classes, which also provides an innovative perspective for
 535 understanding class imbalance on graphs. \square

C More Results

C.1 More results Under Traditional Semi-supervised Settings

In Table 3, we present the results of RVGNN and other algorithms on Cora-Semi ($\rho=10$). The experimental results demonstrate that TAM has a significant gain effect on all models (BalancedSoftmax, Renode, GraphENS), while GraphENS achieved the best performance among all current baselines on various models. By comparing with GraphSMOTE, we can conclude that addressing the node classification imbalance issue should focus on the topological characteristics of the graph. Finally, and most importantly, RVGNN achieved state-of-the-art results under all basic models, which also verifies the superiority of our model in traditional imbalanced segmentation.

Table 3: Experimental results of our method RVGNN and other baselines on *Cora-Semi*. We report averaged balanced accuracy (bAcc.,%) and F1-score (%) with the standard errors over 5 repetitions on three representative GNN architectures. Highlighted are the top **first** and **second**. Δ is the margin by which our method leads state-of-the-art method.

Dataset(<i>Cora-Semi</i>)	GCN		GAT		SAGE	
Imbalance Ratio ($\rho = 10$)	bAcc.	F1	bAcc.	F1	bAcc.	F1
Vanilla	62.82 \pm 1.43	61.67 \pm 1.59	62.33 \pm 1.56	61.82 \pm 1.84	61.82 \pm 0.97	60.97 \pm 1.07
Re-Weight	65.36 \pm 1.15	64.97 \pm 1.39	66.87 \pm 0.97	66.62 \pm 1.13	63.94 \pm 1.07	63.82 \pm 1.30
PC Softmax	68.04 \pm 0.82	67.84 \pm 0.81	66.69 \pm 0.79	66.04 \pm 1.10	65.79 \pm 0.70	66.04 \pm 0.92
GraphSMOTE	66.39 \pm 0.56	65.49 \pm 0.93	66.71 \pm 0.32	65.01 \pm 1.21	61.65 \pm 0.34	60.97 \pm 0.98
BalancedSoftmax	69.98 \pm 0.58	68.68 \pm 0.55	67.89 \pm 0.36	67.96 \pm 0.41	67.43 \pm 0.61	67.66 \pm 0.69
+ TAM	69.94 \pm 0.45	69.54 \pm 0.47	69.16 \pm 0.27	69.39 \pm 0.37	69.03 \pm 0.92	69.03 \pm 0.97
Renode	67.03 \pm 1.41	67.16 \pm 1.67	67.33 \pm 0.79	68.08 \pm 1.16	66.84 \pm 1.78	67.08 \pm 1.75
+ TAM	68.26 \pm 1.84	68.11 \pm 1.97	67.50 \pm 0.67	68.06 \pm 0.96	67.28 \pm 1.11	67.15 \pm 1.11
GraphENS	70.89 \pm 0.71	70.90 \pm 0.81	70.45 \pm 1.25	69.87 \pm 1.32	68.74 \pm 0.46	68.34 \pm 0.33
+ TAM	71.69 \pm 0.36	72.14 \pm 0.51	70.15 \pm 0.18	70.00 \pm 0.40	70.45 \pm 0.74	70.40 \pm 0.75
RVGNN	72.92 \pm 2.27	72.60 \pm 2.26	74.56 \pm 0.96	74.61 \pm 0.96	73.32 \pm 3.02	68.91 \pm 3.13
Δ	+ 1.23	+ 0.46	+ 4.11	+ 4.61	+ 2.87	- 1.49

Table 4: Experimental results of our method RVGNN and other baselines on *Computers-Random*. We report averaged balanced accuracy (bAcc.,%) and F1-score (%) with the standard errors over 5 repetitions on three representative GNN architectures. Highlighted are the top **first** and **second**. Δ is the margin by which our method leads state-of-the-art method.

Dataset(<i>Computers-Random</i>)	GCN		GAT		SAGE	
Imbalance Ratio ($\rho = 25.50$)	bAcc.	F1	bAcc.	F1	bAcc.	F1
Vanilla	78.43 \pm 0.41	77.14 \pm 0.39	71.35 \pm 1.18	69.60 \pm 1.11	65.30 \pm 1.07	64.77 \pm 1.19
Re-Weight	80.49 \pm 0.44	75.07 \pm 0.58	71.95 \pm 0.80	70.67 \pm 0.51	66.50 \pm 1.47	66.10 \pm 1.46
PC Softmax	81.34 \pm 0.55	75.17 \pm 0.57	70.56 \pm 1.46	67.26 \pm 1.48	69.73 \pm 0.53	67.03 \pm 0.6
GraphSMOTE	80.50 \pm 1.11	73.79 \pm 0.14	71.98 \pm 0.21	67.98 \pm 0.31	72.69 \pm 0.82	68.73 \pm 1.01
BalancedSoftmax	81.39 \pm 0.25	74.54 \pm 0.64	72.09 \pm 0.31	68.38 \pm 0.69	73.80 \pm 1.06	69.74 \pm 0.60
+ TAM	81.64 \pm 0.48	75.59 \pm 0.83	74.00 \pm 0.77	70.72 \pm 0.50	73.77 \pm 1.26	71.03 \pm 0.69
Renode	81.64 \pm 0.34	76.87 \pm 0.32	72.80 \pm 0.94	71.40 \pm 0.97	70.94 \pm 1.50	70.04 \pm 1.16
+ TAM	80.50 \pm 1.11	75.79 \pm 0.14	71.98 \pm 0.21	70.98 \pm 0.31	72.69 \pm 0.82	70.73 \pm 1.01
GraphENS	82.66 \pm 0.61	76.55 \pm 0.17	75.25 \pm 0.85	71.49 \pm 0.54	77.64 \pm 0.52	72.65 \pm 0.53
+ TAM	82.83 \pm 0.68	76.76 \pm 0.39	75.81 \pm 0.72	72.62 \pm 0.57	78.98 \pm 0.60	73.59 \pm 0.55
RVGNN	85.00 \pm 0.07	82.35 \pm 0.08	81.94 \pm 0.54	80.94 \pm 0.25	80.61 \pm 0.11	77.49 \pm 0.09
Δ	+ 2.17	+ 5.48	+ 6.13	+ 8.32	+ 1.63	+ 3.90

C.2 More results On Naturally Imbalanced Datasets

One of our major highlights is testing our algorithms on naturally imbalanced datasets, which is more representative of real-world scenarios. This means that in a semi-supervised setting, our training set and unlabeled data are both imbalanced, which poses a significant challenge for oversampling methods. In Table 4, we present the results of RVGNN and all baselines on *Computers-Random*,

550 demonstrating that our model still achieves state-of-the-art performance. In contrast, oversampling
551 methods such as GraphSMOTE perform poorly on this dataset due to overfitting limitations. We
552 believe that our model’s ability to regularize variance from a bottom-up perspective contributes to
553 its superior performance on naturally imbalanced graphs, while also demonstrating its powerful
554 generalization capabilities.

555 D More Analysis

556 D.1 More Ablation Analysis for the Loss Function

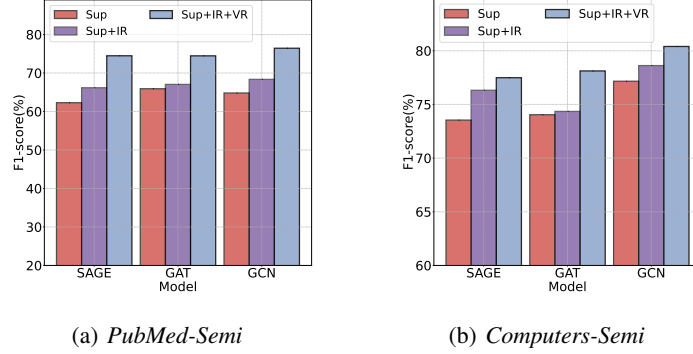


Figure 4: More Ablation Analysis for the Loss Function.

557 **Details of Experimental Setup.** We conduct more experiments on imbalance ratio 10 ($\rho = 10$) for
 558 the datasets PubMed and Amazon-Computers. For each GNN network, we add one loss regularization
 559 at a time, i.e. for GCN, we gradually add \mathcal{L}_{IR} , \mathcal{L}_{VR} from the initial \mathcal{L}_{sup} . The architecture we
 560 employed consisted of a 2-layers graph neural network (GNN) with 128 hidden dimensions, using
 561 GCN [16], GAT [34], and GraphSAGE [9]. The models were trained for 2000 epochs.

562 **Analysis.** Figure 4 reports more experiments for each component in the loss function. As we have
 563 seen, each component of the loss function can bring an improvement in the training effect. It is worth
 564 noting that in most cases \mathcal{L}_{VR} has a larger effect boost than \mathcal{L}_{IR} . We believe that \mathcal{L}_{VR} plays a more
 565 important role in the loss function.

566 D.2 More Experiments for Variance and Imbalance Ratio Correlation in Theorem 1

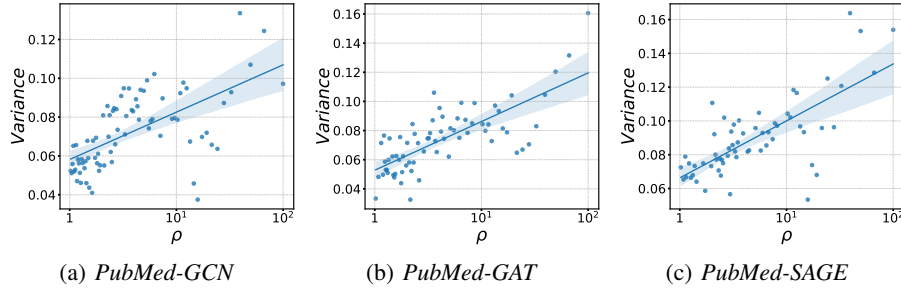


Figure 5: More experiments for variance and imbalance ratio correlation in Theorem 1.

567 **Details of Experimental Setup.** In this experiment, the classifier is not MLP, which means that
 568 the probability of a node v being classified into class i depends on the distance of this node from the
 569 center of class i . We classify half of the classes as majority classes and the other half as minority
 570 classes. Initially, each class in majority classes and minority classes has 200 labeled nodes. To
 571 generate different imbalance scenarios (ρ), we reduce the number of labeled nodes in the minority
 572 class by 1 and add 1 to the number of labeled nodes in the majority class each time, so as to keep the
 573 number of samples in the training set consistent and eliminate the effect of the sample size variance
 574 in the training set. To calculate the variance of the model under a certain imbalance ratio (ρ), we
 575 repeat 20 times to randomly select different but the same number of training sets, and train 2000
 576 epochs for each fixed training set. The architecture we employed consisted of a 2-layers graph neural
 577 network (GNN) with 128 hidden dimensions, using GCN [16], GAT [34], and GraphSAGE [9].

Analysis. Figure 5 reveals an intriguing pattern, as the majority of data points closely align along a regression curve exhibiting a positive slope. This observation provides substantial evidence that establishes a strong and direct linear relationship between the imbalance ratio (ρ) and the associated variance. Consequently, our findings provide compelling support for the hypothesis postulated in Theorem 1.

D.3 More Experiments for Hyperparameter Sensitivity Analysis of RVGNN

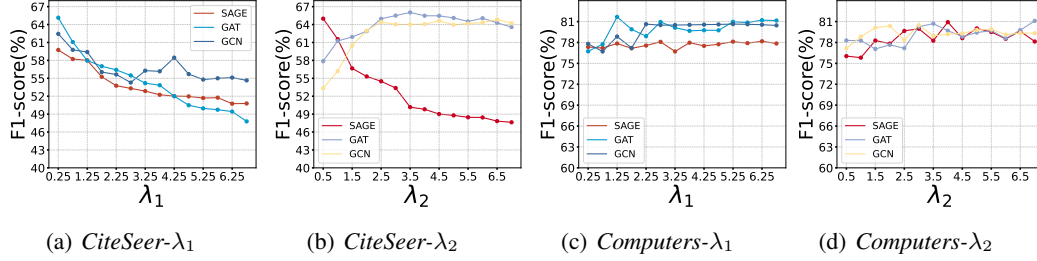


Figure 6: More Experiments of Hyperparameter Sensitivity Analysis for RVGNN.

Details of Experimental Setup. We implemented experiments on the two datasets, CiteSeer and Amazon-Computers. In this experiment, we evaluate the F1 score when we fix one λ while the other changes. The architecture we employed consisted of a 2-layer graph neural network (GNN) with 128 hidden dimensions, using GCN [16], GAT [34], and GraphSAGE [9]. The models were trained for 2000 epochs.

Analysis. In Figure 6, we present the sensitivity analysis of RVGNN to two weight hyperparameters (λ_1, λ_2) in the loss function on the CiteSeer and Amazon-Computers. It is evident that RVGNN exhibits varying degrees of sensitivity to hyperparameters on different datasets. Specifically, the model demonstrates lower sensitivity to a and b on CiteSeer, while the opposite is observed on Amazon-Computers. We postulate that the number of nodes (i.e., dataset size) may be a crucial factor in this discrepancy. Nevertheless, the optimal range of hyperparameter selection appears to be relatively narrow, indicating that significant efforts are not be required for model fine-tuning.

D.4 More Results of loss curve and F1 score

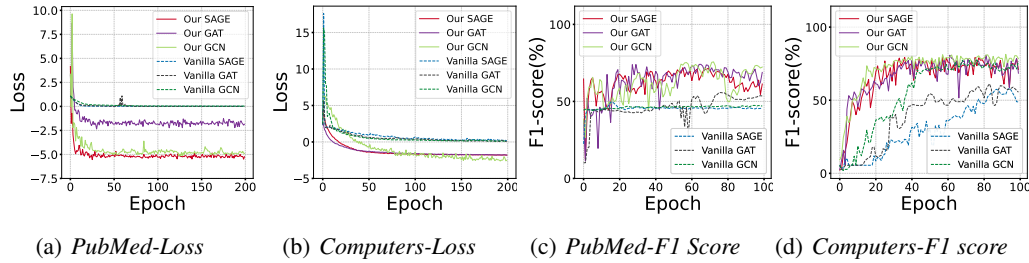


Figure 7: More analysis of loss curve and F1 score.

Details of Experimental Setup. We implemented experiments on three datasets PubMed, CiteSeer and Amazon-Computers. We compare our method RVGNN and other vanilla models (only using cross-entropy and not using graph augmentation). The architecture we employed for RVGNN and vanilla model consisted of a 2-layer graph neural network (GNN) with 128 hidden dimensions, using GCN [16], GAT [34], and GraphSAGE [9]. The models were trained for 2000 epochs.

602 **Analysis.** We plot the loss curve and F1 score with epoch on the PubMed and Computers datasets.
603 By analyzing the plotted loss curves, we observe that our model demonstrates notable advantages,
604 including faster convergence of the loss curves and improved performance across multiple datasets.
605 For the F1 score, our model consistently outperforms alternative approaches, showcasing its superior
606 capability in effectively capturing the complex relationships within the data and making accurate
607 predictions.

E Elaboration of the experimental setup

In this section, we present our approach for constructing imbalanced datasets, describe our evaluation protocol, and provide comprehensive details on our algorithm as well as the baseline methods utilized in our study. To achieve this, we leverage sophisticated techniques and utilize advanced metrics to ensure the reliability and relevance of our results.

E.1 Construction for Imbalanced datasets

The detailed descriptions of the datasets are shown in Table 5. The details of label distribution in the training set of the five imbalanced benchmark datasets are in Table 7, and the label distribution of the full graph is provided in Table 7.

Table 5: Summary of the datasets used in this work.

Dataset	Nodes	Edges	Features	Classes
<i>Cora</i>	2,708	5,429	1,433	7
<i>Citeseer</i>	3,327	4,732	3,703	6
<i>Pubmed</i>	19,717	44,338	500	3
<i>Amazon-Computers</i>	13,752	491,722	767	10
<i>Coauthor-CS</i>	18,333	163,788	6,805	15

Imbalanced datasets construction for Traditional Semi-supervised Settings. For each citation dataset, we adopt the "public" split and apply a random undersampling technique to make the class distribution imbalanced until the target imbalance ratio ρ is achieved. Specifically, we convert the minority class nodes to unlabeled nodes in a random manner. Regarding the co-purchased networks Amazon-Computers, we conduct replicated experiments by randomly selecting nodes as the training set in each trial. We create a random validation set that contains 30 nodes in each class, and the remaining nodes are used as the testing set.

Imbalanced datasets construction for Naturally Imbalanced Datasets. For *Computers-Random* and *CS-Random*, we constructed a training set with equal proportions based on the label distribution of the complete graph (Amazon-Computers). The label distributions for the training sets of *Computers-Random* and *CS-Random* are presented in Table 7. To achieve this, we employed a stratified sampling approach that ensured an unbiased representation of all labels in the training set.

Table 6: Label distributions in the training sets

Dataset	C_0	C_1	C_2	C_3	C_4	C_5	C_6	C_7	C_8	C_9	C_{10}	C_{11}	C_{12}	C_{13}	C_{14}
<i>Cora-Semi</i> ($\rho = 10$)	23.26%	23.26%	23.26%	23.26%	2.32%	2.32%	2.32%	-	-	-	-	-	-	-	-
<i>CiteSeer-Semi</i> ($\rho = 10$)	30.30%	30.30%	30.30%	3.03%	3.03%	3.03%	-	-	-	-	-	-	-	-	-
<i>PubMed-Semi</i> ($\rho = 10$)	47.62%	47.62%	4.76%	-	-	-	-	-	-	-	-	-	-	-	-
<i>Computers-Semi</i> ($\rho = 10$)	18.18%	18.18%	18.18%	18.18%	18.18%	1.82%	1.82%	1.82%	1.82%	1.82%	-	-	-	-	-
<i>Computers-Random</i> ($\rho = 25.50$)	3.01%	15.79%	10.53%	3.76%	38.35%	2.26%	3.01%	6.02%	15.79%	1.50%	-	-	-	-	-
<i>CS-Random</i> ($\rho = 41.00$)	3.98%	2.27%	11.36%	2.27%	7.39%	11.93%	1.70%	5.11%	3.98%	0.57%	7.95%	11.36%	2.27%	23.30%	4.55%

Table 7: Label distributions on the whole graphs

Dataset	C_0	C_1	C_2	C_3	C_4	C_5	C_6	C_7	C_8	C_9	C_{10}	C_{11}	C_{12}	C_{13}	C_{14}
<i>Cora</i> ($\rho \approx 4.54$)	351	217	418	818	426	298	180	-	-	-	-	-	-	-	-
<i>CiteSeer</i> ($\rho \approx 2.66$)	264	590	668	701	696	508	-	-	-	-	-	-	-	-	-
<i>PubMed</i> ($\rho \approx 1.91$)	4103	7739	7835	-	-	-	-	-	-	-	-	-	-	-	-
<i>Amazon-Computers</i> ($\rho \approx 17.73$)	436	2142	1414	542	5158	308	487	818	2156	291	-	-	-	-	-
<i>Coauthor-CS</i> ($\rho \approx 35.05$)	708	462	2050	429	1394	2193	371	924	775	118	1444	2033	420	4136	876

E.2 Architecture of GNNs

We conducted evaluations using three classic GNN architectures, namely GCN [16], GAT [34], and GraphSAGE [9]. The GNN models were constructed with different numbers of layers, namely $L = 1, 2, 3$. To enhance the learning process, each GNN layer was accompanied by a BatchNorm layer with a momentum of 0.99, followed by a PRelu activation function [11]. For the GAT architecture, we employed multi-head attention with 8 heads. We search for the best model on the validation set. The available choices for the hidden unit sizes were 64, 128, and 256.

E.3 Evaluation Protocol

We utilized the Adam optimizer [15] with an initial learning rate of 0.01 or 0.005. To manage the learning rate, we employed a scheduler based on the approach outlined in [30], which reduced the learning rate by half when there was no decrease in validation loss for 100 consecutive epochs. Weight decay with a rate of 0.0005 was applied to all learnable parameters in the model. In the initial training iteration, we trained the model for 200 epochs using the original training set for Cora, CiteSeer, PubMed, or Amazon-Computers. However, for Flickr, the training was extended to 2000 epochs in the first iteration. Subsequently, in the remaining iterations, we trained the models for 2000 epochs using the aforementioned optimizer and scheduler. The best models were selected based on validation accuracy, and we employed early stopping with a patience of 300 epochs.

E.4 Technical Details of RVGNN

For all datasets *Cora-Semi*, *CiteSeer-Semi*, *PubMed-Semi*, *Computers-Semi*, *Computers-Random* and *CS-Random*, the learning rate η is chosen from $\{0.0001, 0.0005, 0.001, 0.005, 0.01, 0.1\}$. The temperature hyperparameter τ is chosen from $\{0.05, 0.08, 0.13, 0.16, 0.21, 0.23, 0.26\}$. The threshold v is chosen from $\{0.6, 0.63, 0.66, 0.7, 0.8, 0.83, 0.9, 0.93, 0.96, 0.99\}$. The factor λ_1 of \mathcal{L}_{VR} is chosen from $\{0.25, 0.35, 0.5, 0.85, 1, 1.5, 2, 2.15, 2.65, 3\}$. The factor λ_2 of \mathcal{L}_{IR} is chosen from $\{0.35, 0.5, 1, 1.25, 1.5, 2.85, 3\}$. The factors of mask node properties and edges in \tilde{G} are chosen from $\{0.4, 0.45, 0.5, 0.6, 0.65, 0.7\}$ and $\{0.4, 0.45, 0.5, 0.55, 0.6, 0.65, 0.7\}$. The factors of mask node properties and edges in \tilde{G}' are chosen from $\{0.1, 0.15, 0.2, 0.3, 0.4, 0.45\}$ and $\{0.1, 0.15, 0.2, 0.3, 0.35, 0.4, 0.45\}$.

E.5 Technical Details of Baselines

For the GraphSMOTE method [47], we employed the branched algorithms whose edge predictions are discrete-valued, which have achieved superior performance over other variants in most experiments. Regarding the ReNode method [6], we search hyperparameters within the lower bound of cosine annealing, $w_{\min} \in 0.25, 0.5, 0.75$, and the upper bound range, $w_{\max} \in 1.25, 1.5, 1.75$, as suggested in [6]. The PageRank teleport probability was fixed at $\alpha = 0.15$, which is the default setting in the released codes. As for TAM [30], we performed a hyperparameter search for the coefficient of the ACM term, $\alpha \in 1.25, 1.5, 1.75$, the coefficient of the ADM term, $\beta \in 0.125, 0.25, 0.5$, and the minimum temperature of the class-wise temperature, $\phi \in 0.8, 1.2$, following the approach described in [30]. The sensitivity to the imbalance ratio of the class-wise temperature, δ , was fixed at 0.4 for all the main experiments. Additionally, consistent with [30], we implemented a warmup phase for 5 iterations, as we utilized model predictions for unlabeled nodes.

E.6 Configuration

All the algorithms and models are implemented in Python and PyTorch Geometric. Experiments are conducted on a server with an NVIDIA 3090 GPU (24 GB memory) and an Intel(R) Xeon(R) Silver 4210R CPU @ 2.40GHz.

Algorithm 1 *RVGNN*

Input: Imbalanced Graph $\mathbf{G} = (\mathbf{V}, \mathbf{E}, \mathbf{V}_L)$, feature matrix \mathbf{X} , adjacency matrix \mathbf{A} , unlabeled set $\mathbf{V}_U = \mathbf{V} - \mathbf{V}_L$, label matrix \mathbf{Y} , feature matrix of labeled nodes and unlabeled nodes \mathbf{X}_L and \mathbf{X}_U , the number of classes k , the threshold v for determining whether a node has confident prediction, temperature hyperparameter τ , GNN model f_θ , the classifier h_θ , learning rate η , The total number T of epochs for model training. The $\text{sim}(\cdot, \cdot)$ computes the cosine similarity between two vectors, CE represents the cross-entropy loss function, and $\mathbb{I}(\cdot)$ is an indicator function.

- 1: **for** $t = 0, 1, \dots, T$ **do**
- 2: Generate two differently augmented views $\tilde{\mathbf{G}} = (\tilde{\mathbf{A}}, \tilde{\mathbf{X}})$ and $\tilde{\mathbf{G}}' = (\tilde{\mathbf{A}}', \tilde{\mathbf{X}}')$ from original graph \mathbf{G} .
- 3: $\tilde{\mathbf{O}}_L \leftarrow f_\theta(\tilde{\mathbf{A}}, \tilde{\mathbf{X}}_L)$, $\tilde{\mathbf{O}}_U \leftarrow f_\theta(\tilde{\mathbf{A}}, \tilde{\mathbf{X}}_U)$
- 4: $\tilde{\mathbf{O}}'_L \leftarrow f_\theta(\tilde{\mathbf{A}}', \tilde{\mathbf{X}}'_L)$, $\tilde{\mathbf{O}}'_U \leftarrow f_\theta(\tilde{\mathbf{A}}', \tilde{\mathbf{X}}'_U)$
- 5: **% Component 1: Supervised Loss**
- 6: $\mathcal{L}_{sup} \leftarrow \frac{1}{2}CE(h_\theta(\tilde{\mathbf{O}}_L), \mathbf{Y}) + \frac{1}{2}CE(h_\theta(\tilde{\mathbf{O}}'_U), \mathbf{Y})$
- 7: **% Component 2: Intra-Class Aggregation Regularization**
- 8: $\mathcal{L}_{IR} \leftarrow -\frac{1}{N_U} \sum_{i=1}^{N_U} \text{sim}(\tilde{\mathbf{O}}_i, \tilde{\mathbf{O}}'_i) - \frac{1}{N_{all}} (\sum_{i=1}^{N_L} \sum_{\substack{j=1 \\ (i,j) \in same}}^{N_L} \text{sim}(\tilde{\mathbf{O}}_i, \tilde{\mathbf{O}}'_j) + \sum_{i=1}^{N_L} \sum_{\substack{j=1, j \neq i \\ (i,j) \in same}}^{N_L} \text{sim}(\tilde{\mathbf{O}}_i, \tilde{\mathbf{O}}_j))$
- 9: **% Component 3: Variance-constrained Optimization with Adaptive Regularization**
- 10: **for** $i = 1, 2, \dots, k$ **do**
- 11: Compute the class centers $\tilde{\mathbf{C}}_i$ and $\tilde{\mathbf{C}}'_i$ for class i
- 12: **end for**
- 13: ▷ Obtain the Label Probability for Each Node
- 14: **for** $i = 1, 2, \dots, |\mathbf{V}|$ **do**
- 15: $\tilde{\mathbf{p}}_i \leftarrow \text{Softmax}(\tilde{\mathbf{O}}_i \cdot [\tilde{\mathbf{C}}_1, \dots, \tilde{\mathbf{C}}_k]^T)$
- 16: $\tilde{\mathbf{p}}'_i \leftarrow \text{Softmax}(\tilde{\mathbf{O}}'_i \cdot [\tilde{\mathbf{C}}'_1, \dots, \tilde{\mathbf{C}}'_k]^T)$
- 17: **end for**
- 18: ▷ Eliminate Nodes with Low Confidence.
- 19: $\mathbf{V}_{\text{conf}} \leftarrow \{i \mid \mathbb{I}(\max(\tilde{\mathbf{p}}'_i) > v) = 1, \forall i \in V_U\}$
- 20: ▷ Replace Labeled Node Prediction
- 21: **for** $i = 1, 2, \dots, |\mathbf{V}_L|$ **do**
- 22: $\tilde{\mathbf{p}}'_i \leftarrow \mathbf{Y}_i$
- 23: **end for**
- 24: $\mathcal{L}_{VR} \leftarrow \frac{1}{|\mathbf{V}_{\text{conf}}|} \sum_{i \in \mathbf{V}_{\text{conf}}} CE(\tilde{\mathbf{p}}'_i, \tilde{\mathbf{p}}_i) + \frac{1}{|\mathbf{V}_L|} \sum_{v_i \in \mathbf{V}_L} CE(\mathbf{Y}_i, \tilde{\mathbf{p}}_i)$
- 25: $\mathcal{L}_{Training} \leftarrow \mathcal{L}_{sup} + \lambda_1 \mathcal{L}_{VR} + \lambda_2 \mathcal{L}_{IR}$
- 26: $\theta \leftarrow \theta - \eta (\nabla \mathcal{L}_{Training})$
- 27: **end for**

Table 8: Notation table.

Indices	
n	The number of nodes, $n = \mathbf{V} $.
f	The node feature dimension.
k	The number of different classes.
D	The dimension of the embedding space, or the dimension of the last layer of GNNs.
Parameters	
\mathbf{G}	An undirected and unweighted graph.
\mathbf{V}	The node set of \mathbf{G} .
\mathbf{E}	The edge set of \mathbf{G} .
\mathbf{X}	The feature matrix of \mathbf{G} , $\mathbf{X} \in \mathbb{R}^{n \times f}$.
\mathbf{V}_L	The set of labeled nodes of \mathbf{G} .
\mathbf{A}	The adjacency matrix of \mathbf{G} , $\mathbf{A} \in \{0, 1\}^{n \times n}$.
$\mathbf{N}(v)$	The set of 1-hop neighbors for node v .
\mathbf{V}_U	The set of unlabeled nodes, $\mathbf{V}_U := \mathbf{V} \setminus \mathbf{V}_L$.
\mathbf{C}_i	The labeled set for class i .
ρ	Imbalance ratio of a dataset, $\rho = \max_i \mathbf{C}_i / \min_i \mathbf{C}_i $.
$\tilde{\mathbf{G}}$	A view from the original graph \mathbf{G} , $\tilde{\mathbf{G}} = (\tilde{\mathbf{A}}, \tilde{\mathbf{X}})$.
$\tilde{\mathbf{G}}'$	Another view from the original graph \mathbf{G} , $\tilde{\mathbf{G}}' = (\tilde{\mathbf{A}}', \tilde{\mathbf{X}}')$.
\mathbf{h}	Output of GNN network, $\mathbf{h} = f_\theta(\mathbf{A}, \mathbf{X})$.
\mathcal{S}	The probability distribution reference matrix, $\mathcal{S} := [C_1; C_2; \dots; C_k]$.
\mathbf{C}	The collection of class centers, $\mathbf{C} := (C_1; C_2; \dots; C_k)$.
p_i	The probability distribution for node i .
y_i	The one-hot label vector for node i .
\mathcal{L}_{sup}	Supervised loss.
\mathcal{L}_{IR}	Intra-class aggregation regularization.
\mathcal{L}_{VR}	Variance-constrained optimization with adaptive regularization.
$\mathcal{L}_{\text{final}}$	The sum of all losses.
V_{conf}	The set of nodes with confident predictions.
λ_1	The weight for \mathcal{L}_{VR} .
λ_2	The weight for \mathcal{L}_{IR} .
v	The threshold for determining whether a node has confident prediction.
τ	The temperature hyperparameter.
α	The size threshold of nodes being added in each class per round.
η	Learning rate of GNN model.
Functions	
$\mathbb{I}(\cdot)$	The indicator function.
$\text{sim}(\cdot)$	The function that computes the cosine similarity between two vectors.
$CE(\cdot)$	The cross-entropy function.
f_θ	GNN model.