

Efficient Subsampling for GNN Downstream Tasks

Author Name1

ABC@SAMPLE.COM and **Author Name2**

XYZ@SAMPLE.COM

Address

Editors: Hung-yi Lee and Tongliang Liu

Abstract

While Graph Neural Networks (GNNs) have shown significant promise for data integration using graph structures, methods to support subsampling graph data are lagging. To address this gap, in this paper, we propose a novel importance-based data subsampling framework. This framework strategically identifies inputs from a primary graph dataset based on their impact on the model’s learning of downstream tasks, such as graph or node classification. Our measure of impact is the predictive uncertainty of each data point. To ensure the subsample is well-representative of the original sample, we cluster them based on their learned graph representation. Finally, subsampling is performed from these identified clusters. The process favours selecting data points with greater prediction uncertainty, while preserving the diversity of the overall sample. We evaluate our approach using a multi-source, real-world dataset on child and youth mental health, comprising emergency department (ED) admissions and mental health questionnaire data. Our experimental results demonstrate that training a GNN with samples identified by the proposed framework yields a statistically significant improvement (on average, 10.13% improvement across metrics from the baseline approach) in predictive performance compared to training on a randomly selected subset of patients. The code is available at https://anonymous.4open.science/r/graph_subsampling-5343.

Keywords: Graph Subsampling; Multi-Dataset Data Integration; Uncertainty Quantification.

1. Introduction

Graph Neural Networks (GNNs) have emerged as a powerful class of deep learning models designed to handle data with complex relational structures. While GNNs have demonstrated strong potential for integrating data through graph-based representations, techniques for effectively subsampling graph data remain underdeveloped. An important application area of graph data subsampling is in medical AI, where data is often multi-modal and collected from diverse data stores. The following real-world case demonstrates the challenges of graph data subsampling.

The authors of this article, in collaboration with a clinical team, are designing a clinical decision support system for mental health using GNN where the relational structure arises from the electronic health records (EHRs) of the patients and their responses to mental health questionnaire, where the data collection scale is on the order of thousands. To improve the prediction performance of the network, the study also needs to collect audio data from a subset of patients on the order of a hundred, as the process of audio data collection is resource-intensive. As a result, finding patients who are well-representative of

the original patient data and at the same time will have the most influence on the network’s performance is a major challenge this research project is facing.

Several researchers have studied the problem of graph data subsampling with the goal of improving network performance (Xu et al., 2023; Jin et al., 2022; Gupta et al., 2024; Georgiev et al., 2023; Jain et al., 2025). These studies primarily focus on compressing graphs or sampling subgraphs, thereby reducing both dataset size and the size of individual graphs. However, graph compression is generally not applicable to diverse graph integration, where data originates from various data stores, such as EHRs and mental health questionnaires, each with its own distinct graph structure. Unlike graph compression, which aims to reduce the size of a graph, graph integration often seeks to construct larger, unified graphs. A promising approach with potential for graph data subsampling is coreset selection, which is primarily used in active learning. Coreset selection offers an approach to select informative samples from unlabeled data for human labelling, which can subsequently improve network performance (Yoo and Kweon, 2019; Katharopoulos and Fleuret, 2018).

In this paper, we propose an efficient graph subsampling approach tailored for downstream prediction tasks, such as graph classification. We are inspired by the coreset selection approach in identifying samples that are most influential for training a predictive GNN. We use the network’s prediction uncertainty as a metric for identifying samples that are highly informative for network training (Katharopoulos and Fleuret, 2018; Jeong et al., 2023; Chen et al., 2024). However, to improve efficiency, we employ self-distillation in the quantification of uncertainty. Self-distillation allows for the training of a multi-classifier GNN, thereby significantly reducing the computational costs of both training and inference for uncertainty estimation (Daneshvar and Samavi, 2025). Selecting data points solely based on the highest uncertainty can lead to the undesirable outcome of mostly selecting outliers (Katharopoulos and Fleuret, 2018; Jeong et al., 2023; Settles and Craven, 2008; Chen et al., 2024). To counter this problem and promote diversity among selected samples, we utilize K-means clustering based on graph representations. Within each cluster, data points are further grouped into percentiles according to their uncertainty values. Our diversity-enabled grouping strategy ensures the inclusion of samples across the full range of uncertainty, preventing bias towards extreme outliers. Ultimately, samples are selected with a weighted emphasis on those exhibiting higher prediction uncertainty.

The key contributions of this paper are as follows. First, we propose a diversity-enabled subsampling method for GNNs specifically designed for graph classification—an underexplored direction, particularly in the context of heterogeneous graphs. Second, our method leverages self-distillation to train a multi-classifier GNN, followed by uncertainty quantification, substantially reducing the computational cost typically associated with uncertainty-based subsampling. Third, we empirically demonstrate that our framework effectively selects representative samples from a primary dataset to augment with a secondary one, thereby improving training performance in scenarios involving diverse data integration under limited-data conditions.

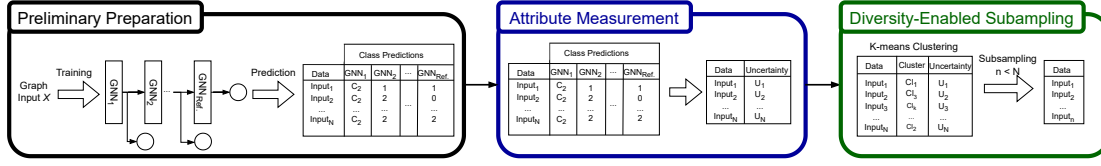


Figure 1: Overview of the proposed subsampling approach.

2. Methodology

2.1. Problem Setup

Our proposed approach is inspired by *coreset selection*, a technique primarily used in active learning to efficiently select representative samples from an unlabeled dataset for human experts to label (Yoo and Kweon, 2019). By adopting similar principles, we aim to propose a subsampling methodology for graphs that selects samples highly relevant to a given downstream task, such as graph classification, thereby making the sample more effective for GNN training.

There are three approaches of *coreset selection*: (1) uncertainty-based, (2) diversity-based, and (3) approaches based on expected changes in the network (Yoo and Kweon, 2019). The uncertainty-based approaches utilize the quantified uncertainty of the data points to select samples with higher uncertainty. Diversity-based approaches select diverse samples to represent the underlying distribution of the data. Finally, the approaches based on the expected changes in the network select samples with greater impact on network parameters or outputs (Yoo and Kweon, 2019). Key challenges in *coreset selection* include achieving computational efficiency, a common disadvantage of uncertainty-based approaches, and simultaneously ensuring sample diversity.

Our proposed approach addresses both key challenges in identifying the most influential samples for network training. Specifically: (1) the method achieves computational efficiency in uncertainty quantification by utilizing only one network instead of multiple networks (deep ensemble) or multiple passes of the input through the same network (MC Dropout). Furthermore, (2) the subsampling strategy promotes diversity among samples by grouping the data based on their graph structure and sampling from all groups. As shown in Figure 1, the approach consists of three pipelines: (1) preliminary preparation (Section 2.2), (2) attribute measurement, in which we have chosen prediction uncertainty as the attribute (Section 2.3), and (3) diversity-enabled subsampling (Section 2.4). We use GNNs for graph classification throughout this work; however, the approach is generalizable to any classification task involving GNNs.

2.2. Preliminary Preparation

The goal of the preliminary preparation phase is to train a network that can be effectively utilized to quantify the prediction uncertainty for all data points. We utilize the *self-distillation* approach to train multiple GNN networks simultaneously. Self-distillation is a special case of knowledge distillation, where the same network serves as both the teacher and the student simultaneously, facilitating knowledge transfer within the network (Zhang

et al., 2022; Gou et al., 2021). Self-distillation is especially helpful when used with over-parameterized GNNs as it reduces the computational complexity of training multiple networks separately (Chen et al., 2021). We follow the same strategy as Zhang et al. (2022) by adding a classifier after each layer of the network and utilizing the deepest classifier as the teacher.

During training, each classifier learns to mimic the deepest classifier, i.e., the teacher. As a result, the training process involves reducing the total distillation loss, which is a combination of distillation loss and feature penalty (Zhang et al., 2022; Daneshvar and Samavi, 2025). The total distillation loss for all n training graphs can be defined as

$$L = \frac{1}{n} \sum_{i=1}^n L_{\text{dist}}^i + L_{\text{pen}}^i, \quad (1)$$

where L_{dist} and L_{pen} are the distillation loss and feature penalty respectively.

The distillation loss is the mean of layer-wise weighted sums of two components, cross-entropy and KL-divergence, and can be computed as

$$L_{\text{dis}}^i = \frac{1}{m} \sum_{l=1}^m \left((1 - \alpha_l) L_{\text{CE}_l}^i + \alpha_l L_{\text{KL}_l}^i \right), \quad (2)$$

where $\alpha_l \in [0, 1]$ is the imitation parameter for each classifier. The L_{CE} and L_{KL} are the cross-entropy of the classifier output at layer l with the true label, and the KL-divergence between the student’s soft labels at layer l and the teacher’s soft labels, respectively. The feature penalty can be computed as

$$L_{\text{pen}}^i = \frac{1}{m} \sum_{l=1}^m \lambda_l L_2(h_l^i, h_t^i), \quad (3)$$

where $\lambda_l > 0$ is the trade-off parameter for each classifier and L_2 is the squared ℓ_2 -norm loss. h_l^i and h_t^i are features extracted in layer l and features extracted by the teacher network, respectively. Both the imitation parameter and the trade-off parameter are set to zero for the teacher network.

To make predictions for all data points, we utilize K-fold cross-validation. As a result, each data point will be included in the test dataset exactly once. The predictions of the students and the teacher network are recorded for each data point in the test dataset, which will later be used to distinguish between hard and easy examples.

2.3. Attribute Measurement

During the attribute measurement phase, we utilize the previously trained network to quantify the prediction uncertainty of each data point in the dataset. A data point would be harder for the network to classify if the prediction of shallower classifiers, i.e., classifiers after shallower layers, don’t match that of the deepest classifier (Zhang et al., 2022; Daneshvar and Samavi, 2025). This is because deeper classifiers utilize deeper feature extractors, providing them with more informative information. This is especially the case with GNNs, as deeper graph layers, i.e., deeper feature extractors, aggregate information from more distant nodes (Hamilton, 2020). To rank the samples based on how hard they are for the network

to classify, we utilize a metric proposed by Daneshvar and Samavi (2025) that captures disagreement between predictions of each classifier. The disagreement between the classifiers' predictions is referred to as the network's prediction uncertainty.

To quantify the network's prediction uncertainty, the metric utilizes a weighted disagreement metric based on a normalized weighted Jensen–Shannon divergence (JSD) (Daneshvar and Samavi, 2025). The uncertainty metric can be computed as

$$UC = \sum_{l=1}^m W(l) \times JSD(P_l || P_{teacher}), \quad (4)$$

where $W(l)$ is a bounded weight function ($1 \leq W(l) \leq 2$) that assigns a weight to layer l . We utilize the same nonlinear weight function proposed by Daneshvar and Samavi (2025).

The goal of the weight function is to assign higher weights to classifiers based on their depth relative to the deepest classifier. The weight of a classifier at layer l is computed by

$$W(l) = (-(\exp(D(l) - L)) + 2)^{\mathbb{1}_{\{y_l \neq y_{teacher}\}}}, \quad (5)$$

where L and $D(i)$ are the total number of network layers and classifier i 's distance from the deepest layer, respectively. $\mathbb{1}$ is the indicator function, which returns one if the y_l differs from $y_{teacher}$; otherwise, it returns zero. y_l and $y_{teacher}$ are the predicted class by the l th classifier and the predicted class by the teacher classifier, respectively.

JSD can be computed as

$$JSD(P_l || P_{outcome}) = \frac{1}{2} \left(KL(P_l || M) + KL(P_{outcome} || M) \right), \quad (6)$$

where $M = \frac{1}{2}(P_l + P_{outcome})$ is a mixture distribution of P_l and $P_{outcome}$. JSD is bounded ($0 \leq JSD \leq \log_b(2)$) as the mixture distribution helps with averaging and smoothing out the values. The uncertainty quantification metric needs to be normalized. The upper bound of UC for a network with m layers, including the teacher network, can be computed as

$$UC_{max} = \sum_{l=1}^{m-1} W(l) \times \log_e(2), \quad (7)$$

where \log_e is the natural logarithm and $W(l)$ computes layer l 's. Finally, the normalized uncertainty metric can be computed as

$$UC_{norm} = \frac{UC}{UC_{max}}. \quad (8)$$

Algorithm 1 outlines the overall steps in quantifying the uncertainty of each data point. It begins by training a network for each split (lines 2-6) and then quantifying the uncertainty of the test set for each split (lines 7-11).

2.4. Diversity-Enabled Subsampling

The diversity-enabled subsampling phase ensures that the selected samples are not only representative of the dataset but also more influential in training the final predictive network.

Algorithm 1: Computing uncertainty of all data points

Input: Graph data points X , a GNN network M , total number of layers l

Output: Prediction uncertainty of the data points U

```

1  $U \leftarrow \emptyset$ 
2 // Step 1: Split data using K-Fold Cross Validation.
3 foreach ( $trainSet, testSet$ )  $\in k\_fold(X)$  do
4   // Step 2: Train the network on  $trainSet$ .
5    $M \leftarrow train(M, trainSet)$ 
6   // Step 3: Compute uncertainty of  $testSet$ .
7   foreach  $testData \in testSet$  do
8      $u \leftarrow 0$ 
9      $u \leftarrow uncertainty(M, testData, l)$  // Compute uncertainty.
10     $U \leftarrow U \cup \{u\}$  // Add  $u$  to the set  $U$ 
11  end
12 end
13 return  $U$ 

```

If the subsampling strategy selects only from the data points with the highest uncertainty, we risk selecting only from the outliers (Jeong et al., 2023; Settles and Craven, 2008). Therefore, to promote diversity among samples, we will first group the graphs and then sample graphs proportional to the size of each group. We can utilize unsupervised learning to find hidden patterns in the data without the need for labelled data. One fundamental unsupervised learning approach is clustering. Clustering aims to group data samples based on similarity without requiring labelled data by finding hidden patterns that might exist in the data (Xu and Wunsch, 2005; Huang, 1998). The clustering algorithm will provide disjoint clusters, each containing data that is similar (Na et al., 2010). One of the widely used algorithms in clustering is the K-means clustering algorithm (MacQueen, 1967), which is known for its simplicity.

The clustering approach in K-means clustering begins by selecting k cluster centers, usually chosen randomly. Then the algorithm assigns each data point to a cluster based on its distance to the chosen centers. Euclidean distance is a widely used metric to compute the distance of each data point to the cluster centers (Xu and Wunsch, 2005; Na et al., 2010). The value of k is arbitrary and fixed at the beginning of the algorithm. After assigning each data point to a cluster, the algorithm recalculates the cluster centers by minimizing an objective function. The objective function reduces the distance of each item in the cluster to the cluster average (cluster center), and is computed as (Na et al., 2010):

$$E = \sum_{i=1}^k \sum_{x \in C_i} |x - x_i|^2, \quad (9)$$

where C_i , x , and x_i are cluster i , data point x in cluster C_i , and the average of cluster C_i respectively. The process is repeated until there is no change in the cluster centers.

In our approach, we utilize the graphs' information for clustering. Graphs are grouped based on their graph representation. To this end, the graph representations created by the teacher network are utilized. The intuition behind grouping graphs based on their

representations is to ensure the sample is representative of the underlying data distribution while taking the graph structure into account.

The sum of the squared distances of samples to their closest cluster center, also known as inertia, is used to optimize the number of clusters (k). However, having too many clusters will result in poor subsampling, as there is a risk of outliers being assigned to the same cluster, and fewer data points in each cluster, resulting in a sample that is not representative of the dataset. Therefore, k should satisfy two conditions: (1) k should be a small number, i.e., limited number of cluster, and (2) the drop in the sum of the squared distances of samples to their closest cluster center computed with k clusters, should be significant compared to having only one cluster.

To reduce the risk of only selecting outliers in each cluster, the framework groups graphs in each cluster into four percentiles based on the uncertainty values computed by Equation (8), sorted in descending order, 0-25%, 25-50%, 50-75%, and 75-100%. The approach samples the same number of graphs from each cluster with an emphasis on the harder examples, i.e., 50-75% and 75-100% percentiles. To achieve this, in each cluster, random sampling is performed as follows: 45% of samples are drawn from the 75-100% percentile range, 30% from 50-75%, 15% from 25-50%, and 10% from 0-25%. Algorithm 2 outlines the subsampling strategy. It starts by grouping the graphs into k clusters using individual graphs (line 2). Then, for each group, samples will be chosen from each percentile with a focus on samples with higher uncertainty, i.e., upper percentiles, proportionate to the size of the cluster (lines 5-10).

Algorithm 2: Subsampling data points with higher uncertainty

Input: Graph representations provided by the teacher network G_{rep} , uncertainty of samples U , number of clusters k , and number of samples needed m

Output: Subset of the data points X' with higher uncertainty

```

1 // Step 1: Grouping graphs.
2 groups = k_means( $G_{rep}$ ,  $k$ )
3 // Step 2: Select a subsample of size  $m$ .
4  $X' \leftarrow \emptyset$ 
5 foreach group  $\in$  groups do
6    $n = \text{round}((\text{size}(\text{group})/\text{size}(X)) * m)$  // Samples needed from the group.
7   foreach  $(p_1, p_2, r) \in \{(0, 25, 0.1), (25, 50, 0.15), (50, 75, 0.3), (75, 100, 0.45)\}$  do
8      $X' \leftarrow X' \cup \text{subsample}(\text{percentile}(U, \text{group}, p_1, p_2), \text{round}(r * n))$ 
9   end
10 end
11 return  $X'$ 

```

3. Experimental Evaluations

The proposed framework identifies and selects samples that provide superior utility for network training compared to random sampling. We assess the impact of our importance-based subsampling method on network performance, comparing it directly with random sampling. To do so, we evaluate the performance of an identical network trained under

three conditions: (1) on the entire dataset using K-fold cross-validation, (2) on randomly sampled data, and (3) on data sampled using the proposed framework. We anticipate observing improved network performance when training with the limited data selected by our subsampling framework, compared to random sampling, and achieving comparable performance to training the network on all the training data.

3.1. Experimental Setup

Dataset: We have utilized two linked datasets containing data from medical health records and mental health questionnaires from real-world patients. The prediction task involves determining whether a patient will be admitted to the emergency department (ED) within 180 days of completing the mental health questionnaire. The datasets comprise ED visits of 1,086 unique patients, along with their responses to the mental health outpatient questionnaire. There are 281 patients who have visited the ED within 180 days of their initial outpatient visit. To address dataset imbalance, we balanced the "Not Admitted" and "Admitted" groups by repeated undersampling the former to 281 patients. We utilized the available information to create a graph for each patient based on their medical history. We then utilized the mental health questionnaire responses to integrate additional data into each patient's graph. Appendix A provides more details on the dataset and the generated graph for each patient.

Network: We utilized two networks. The first network is a multi-classifier GNN, used by the subsampling framework to sample patients from the main dataset, i.e., the medical records dataset, trained using the self-distillation approach. The network consists of three GraphConv (Morris et al., 2021) layers, followed by a ReLU activation function and a batch normalization layer. A final readout layer is applied before each classifier, which consists of a global mean pooling layer. The purpose of utilizing a readout layer is to combine node representations into a single, final graph representation for use in graph classification. The network is trained on the same 180-day ED admission prediction task. It is worth noting that the first network is only utilized for uncertainty quantification; therefore, the network's performance is not a concern.

The second network utilizes the augmented patient graph, generated by integrating information from both datasets, to predict whether a patient would be admitted to the ED within 180 days of their initial mental health assessment. The network consists of a single layer of GraphConv (Morris et al., 2021), followed by similar ReLU and batch normalization layers. A similar readout layer is applied to the node representations to create a graph representation, which is then fed to the classifier.

Training and Hardware Specifications: All networks were implemented using Python version 3.9 and PyTorch version 1.13.1. The networks have been trained on a GPU (NVIDIA GeForce RTX 3050) with CUDA version 12.5. For optimization, we used Adam Optimizer.

3.2. Results and Discussion

To select a suitable value for k , we employed the *Elbow Method*. During each iteration, we used 33% of the data to fit a K-means clustering model and computed the inertia. Figure 2 illustrates how different values of k ($1 \leq k \leq 20$) affect the inertia. As can be observed, $k =$

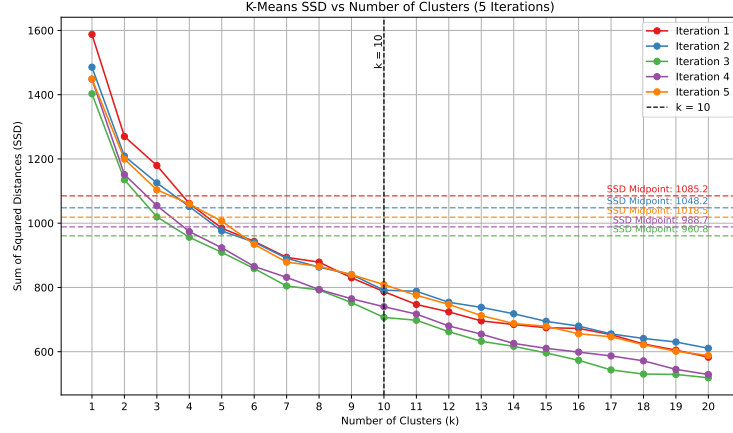


Figure 2: Candidate k versus sum of squared distances of samples to their cluster center, repeated for 5 iterations.

10 represents an 'elbow' point where the rate of decrease in inertia significantly diminishes. This observation was further validated by experimentally evaluating the effect of different k values on the overall subsampling framework. As shown in Table 1, increasing the number of clusters from 5 to 10 results in improved performance; however, performance decreases as the number of clusters approaches 20. The decline in performance occurs because a larger number of clusters can reduce the number of samples in each cluster, thereby increasing the risk of selecting more outliers and diminishing overall sample diversity. Thus, $k = 10$ yields the best performance among other candidates, as it (1) is small enough for computational efficiency yet large enough to provide sufficient diversity, and (2) substantially reduces the clustering error (by more than half), satisfying both conditions for selecting k as discussed in Section 2.4.

Table 2 shows the performance of the network. The same number of samples was selected for both the random sampling and the proposed framework (roughly 278 samples). Additionally, we trained the same GNN network without subsampling to compare the performance of the network trained with and without subsampling. The GNN network takes the augmented patient graph as input and predicts whether a patient would be at risk of admission to the ED within 180 days of their initial mental health assessment. We have measured accuracy, ROC AUC, precision, recall, and F1 score to compare the performance

Table 1: Comparison of different values of k on final network performance.

Number of Clusters	Accuracy	ROC AUC	F1 Score
5	0.5 ± 0.03	0.55 ± 0.06	0.51 ± 0.05
10	0.61 ± 0.04	0.63 ± 0.04	0.64 ± 0.03
15	0.52 ± 0.08	0.57 ± 0.06	0.53 ± 0.07
20	0.5 ± 0.04	0.52 ± 0.11	0.54 ± 0.07

of the approaches. As shown in the table, subsampling using the proposed framework has improved the results across all metrics. Note that training the network with a smaller sample size, when samples are selected using the proposed framework, yields a performance comparable to (and even slightly improved over) training the network using the entire training dataset.

The improvement in Recall and F1 Score (marked with an asterisk in Table 2) is statistically significant when comparing the proposed framework with the random sampling strategy. This improvement has been verified under both Bonferroni and Benjamini-Hochberg corrections, indicating that the network’s ability to identify positive cases and maintain a strong precision-recall balance correctly is a statistically reliable improvement. Our proposed subsampling framework demonstrates a consistent, albeit slight, improvement across most evaluation metrics when compared to training with the complete dataset. Of these improvements, only the gain in Recall (marked with a double dagger in Table 2) was found to be statistically significant.

Finding: The proposed subsampling framework allows for identifying samples that are more informative in training the network. This approach enables researchers to purposefully select samples from the dataset, selecting those most helpful for training the networks, rather than relying on random selection. This is particularly important in utilizing linked datasets, as it facilitates subsampling for data integration.

4. Related Work

We can categorize methods for graph subsampling into two main categories: *coreset selection* and *graph dataset compression*. In this section, we provide a review of methods in each category and conclude by outlining how our approach differs from these methods.

Coreset selection is an approach primarily used in active learning to efficiently select representative samples from an unlabeled dataset for human experts to label (Yoo and Kweon, 2019). Coreset selection is extensively studied for subsampling purposes, as it provides methods to identify and sample data points that significantly improve a network’s learning ability, stemming from the understanding that not all samples contribute equally to model performance (Katharopoulos and Fleuret, 2018; Jeong et al., 2023; Xie et al., 2023; Joshi and Mirzasoleiman, 2023; Yoo and Kweon, 2019; Citovsky et al., 2021; Ash et al., 2020; Caramalau et al., 2021). Principles of coreset selection have been widely applied to various deep learning domains, including large language models (Xie et al., 2023) and computer vision (Jeong et al., 2023). However, applying coreset selection to graphs remains an underexplored topic, particularly for graph subsampling. For instance, Ding et al. (2024) proposed a method for neighbourhood subgraph selection around a node, which differs from our focus on selecting graphs for dataset-level subsampling.

Table 2: Results of performance improvement using the proposed framework.

Method	Accuracy	ROC AUC	Precision	Recall	F1 Score
Without Subsampling	0.59 ± 0.07	0.63 ± 0.06	0.59 ± 0.07	0.61 ± 0.1	0.59 ± 0.07
Random Sampling	0.57 ± 0.08	0.59 ± 0.07	0.57 ± 0.08	0.57 ± 0.1	0.57 ± 0.08
Proposed Framework	0.61 ± 0.04	0.63 ± 0.04	0.6 ± 0.05	$0.68 \pm 0.03^{*\dagger}$	$0.64 \pm 0.03^*$

There are limited studies on subsampling and coreset selection for graph datasets. The primary focus of existing graph subsampling studies, such as the KiDD (Xu et al., 2023) and DosCond (Jin et al., 2022) approaches, is **graph dataset compression**, which involves compressing the GNN training dataset by creating a synthetic, smaller representation of the original. Both KiDD and DosCond notably utilize gradient matching for graph compression. Additionally, MIRAGE (Gupta et al., 2024) further leverages computation trees of training graphs for dataset compression. Beyond the mentioned studies, researchers have employed tree mover’s distance (TMD) either to select specific training sets for tasks related to neural algorithmic reasoning (Georgiev et al., 2023) or to aim to reduce both the number of graphs and the size of individual graphs for GNN training (Jain et al., 2025). However, this graph compression paradigm, which prioritizes creating smaller datasets with smaller graphs, is not ideal for scenarios involving the integration of data from multiple datasets, especially when the datasets contain a limited number of data points.

Our goal is to sample graphs that are most influential in training a GNN for a specific downstream task, such as graph classification. Our framework will enable researchers to leverage insights from our subsampling approach to strategically collect more diverse and potentially multi-source data, facilitating its integration into a unified graph structure that can improve the network’s performance.

5. Conclusion

In this paper, we present an efficient framework for data subsampling based on the prediction uncertainty of the available data points. The purpose of the approach is to help sample data points from a primary dataset for additional data integration from other available, and often limited, datasets. The additional data will help augment the primary dataset and potentially improve network performance. Using the primary dataset, the framework utilizes self-distillation to quantify the prediction uncertainty of the data points. Using K-means clustering, the data points are grouped into k clusters based on their graph representation. In each cluster, based on the data point’s prediction uncertainty, the data is grouped into 0-25%, 25-50%, 50-75%, and 75-100% percentiles. Finally, data is sampled from each cluster with an emphasis on higher percentiles. The evaluations have shown that when sampling data using the proposed framework, the network trained on the augmented dataset has a statistically significant improvement in performance compared to a network trained on randomly sampled augmented data. Additionally, the network exhibits comparable performance to one trained on the whole augmented training dataset.

As a future direction, we plan to incorporate explainability into the uncertainty quantification approach to help users better understand the data selection process. Additionally, we plan to develop an approach to find the right number of samples. Finally, we plan to evaluate the approach on additional datasets.

References

Jordan T. Ash, Chicheng Zhang, Akshay Krishnamurthy, John Langford, and Alekh Agarwal. Deep batch active learning by diverse, uncertain gradient lower bounds. In *Internation-*

- 355 *tional Conference on Learning Representations*, 2020. URL [https://openreview.net/](https://openreview.net/forum?id=ryghZJBKPS)
356 [forum?id=ryghZJBKPS](https://openreview.net/forum?id=ryghZJBKPS).
- 357 Razvan Caramalau, Binod Bhattarai, and Tae-Kyun Kim. Sequential graph convolutional
358 network for active learning. In *Proceedings of the IEEE/CVF Conference on Computer*
359 *Vision and Pattern Recognition (CVPR)*, pages 9583–9592, June 2021.
- 360 Jiayi Chen, Benteng Ma, Hengfei Cui, and Yong Xia. Think twice before selection: Fed-
361 erated evidential active learning for medical image analysis with domain shifts. In *Pro-*
362 *ceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*
363 *(CVPR)*, pages 11439–11449, June 2024.
- 364 Yuzhao Chen, Yatao Bian, Xi Xiao, Yu Rong, Tingyang Xu, and Junzhou Huang. On
365 self-distilling graph neural network, 2021. URL <https://arxiv.org/abs/2011.02255>.
- 366 Gui Citovsky, Giulia DeSalvo, Claudio Gentile, Lazaros Karydas, Anand Rajagopalan, Af-
367 shin Rostamizadeh, and Sanjiv Kumar. Batch active learning at scale. In M. Ranzato,
368 A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan, editors, *Advances*
369 *in Neural Information Processing Systems*, volume 34, pages 11933–11944. Curran Asso-
370 ciates, Inc., 2021. URL [https://proceedings.neurips.cc/paper_files/paper/2021/](https://proceedings.neurips.cc/paper_files/paper/2021/file/64254db8396e404d9223914a0bd355d2-Paper.pdf)
371 [file/64254db8396e404d9223914a0bd355d2-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2021/file/64254db8396e404d9223914a0bd355d2-Paper.pdf).
- 372 Hiran Daneshvar and Reza Samavi. Gnn’s uncertainty quantification using self-distillation,
373 2025. URL <https://arxiv.org/abs/2506.20046>.
- 374 Mucong Ding, Yinhan He, Jundong Li, and Furong Huang. Spectral greedy coresets for
375 graph neural networks, 2024. URL <https://arxiv.org/abs/2405.17404>.
- 376 Dobrik Georgiev Georgiev, Pietro Lio, Jakub Bachurski, Junhua Chen, Tunan Shi, and
377 Lorenzo Giusti. Beyond erdos-renyi: Generalization in algorithmic reasoning on graphs.
378 In *The Second Learning on Graphs Conference*, 2023. URL [https://openreview.net/](https://openreview.net/forum?id=TTxQAkg9QG)
379 [forum?id=TTxQAkg9QG](https://openreview.net/forum?id=TTxQAkg9QG).
- 380 Jianping Gou, Baosheng Yu, Stephen J. Maybank, and Dacheng Tao. Knowledge distil-
381 lation: A survey. *International Journal of Computer Vision*, 129(6):1789–1819, 6 2021.
382 ISSN 1573-1405. doi: 10.1007/s11263-021-01453-z.
- 383 Mridul Gupta, Sahil Manchanda, Hariprasad Kodamana, and Sayan Ranu. Mirage: Model-
384 agnostic graph distillation for graph classification, 2024. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2310.09486)
385 [2310.09486](https://arxiv.org/abs/2310.09486).
- 386 William L Hamilton. Graph representation learning. *Synthesis Lectures on Artificial Intel-*
387 *ligence and Machine Learning*, 14(3):1–159, 2020.
- 388 Zhexue Huang. Extensions to the k-means algorithm for clustering large data sets with
389 categorical values. *Data Mining and Knowledge Discovery*, 2(3):283–304, 09 1998.
- 390 Mika Sarkin Jain, Stefanie Jegelka, Ishani Karmarkar, Luana Ruiz, and Ellen Vitercik.
391 Subsampling graphs with gnn performance guarantees, 2025. URL [https://arxiv.org/](https://arxiv.org/abs/2502.16703)
392 [abs/2502.16703](https://arxiv.org/abs/2502.16703).

- 393 Yuna Jeong, Myunggwon Hwang, and Wonkyung Sung. Training data selection based on
 394 dataset distillation for rapid deployment in machine-learning workflows. *Multimedia Tools*
 395 *and Applications*, 82(7):9855–9870, 2023.
- 396 Wei Jin, Xianfeng Tang, Haoming Jiang, Zheng Li, Danqing Zhang, Jiliang Tang, and Bing
 397 Yin. Condensing graphs via one-step gradient matching. In *Proceedings of the 28th ACM*
 398 *SIGKDD Conference on Knowledge Discovery and Data Mining*, KDD ’22, page 720–730,
 399 New York, NY, USA, 2022. Association for Computing Machinery. ISBN 9781450393850.
 400 doi: 10.1145/3534678.3539429. URL <https://doi.org/10.1145/3534678.3539429>.
- 401 Siddharth Joshi and Baharan Mirzasoleiman. Data-efficient contrastive self-supervised
 402 learning: Most beneficial examples for supervised learning contribute the least. In An-
 403 dreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and
 404 Jonathan Scarlett, editors, *Proceedings of the 40th International Conference on Machine*
 405 *Learning*, volume 202 of *Proceedings of Machine Learning Research*, pages 15356–15370.
 406 PMLR, 07 2023. URL <https://proceedings.mlr.press/v202/joshi23b.html>.
- 407 Angelos Katharopoulos and Francois Fleuret. Not all samples are created equal: Deep
 408 learning with importance sampling. In Jennifer Dy and Andreas Krause, editors,
 409 *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of
 410 *Proceedings of Machine Learning Research*, pages 2525–2534. PMLR, 07 2018. URL
 411 <https://proceedings.mlr.press/v80/katharopoulos18a.html>.
- 412 James MacQueen. Some methods for classification and analysis of multivariate observa-
 413 tions. In *Proceedings of the Fifth Berkeley Symposium on Mathematical Statistics and*
 414 *Probability, Volume 1: Statistics*, volume 5, pages 281–298. University of California press,
 415 1967.
- 416 Christopher Morris, Martin Ritzert, Matthias Fey, William L Hamilton, Jan Eric Lenssen,
 417 Gaurav Rattan, and Martin Grohe. Weisfeiler and Leman Go Neural: Higher-order Graph
 418 Neural Networks, 2021.
- 419 Shi Na, Liu Xumin, and Guan Yong. Research on k-means clustering algorithm: An
 420 improved k-means clustering algorithm. In *2010 Third International Symposium on*
 421 *Intelligent Information Technology and Security Informatics*, pages 63–67, 2010. doi:
 422 10.1109/IITSI.2010.74.
- 423 Burr Settles and Mark Craven. An analysis of active learning strategies for sequence labeling
 424 tasks. In *proceedings of the 2008 conference on empirical methods in natural language*
 425 *processing*, pages 1070–1079, 2008.
- 426 Sang Michael Xie, Shibani Santurkar, Tengyu Ma, and Percy S Liang. Data se-
 427 lection for language models via importance resampling. In A. Oh, T. Naumann,
 428 A. Globerson, K. Saenko, M. Hardt, and S. Levine, editors, *Advances in Neural*
 429 *Information Processing Systems*, volume 36, pages 34201–34227. Curran Associates,
 430 Inc., 2023. URL [https://proceedings.neurips.cc/paper_files/paper/2023/file/](https://proceedings.neurips.cc/paper_files/paper/2023/file/6b9aa8f418bde2840d5f4ab7a02f663b-Paper-Conference.pdf)
 431 [6b9aa8f418bde2840d5f4ab7a02f663b-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2023/file/6b9aa8f418bde2840d5f4ab7a02f663b-Paper-Conference.pdf).

- 432 Rui Xu and D. Wunsch. Survey of clustering algorithms. *IEEE Transactions on Neural*
 433 *Networks*, 16(3):645–678, 2005. doi: 10.1109/TNN.2005.845141.
- 434 Zhe Xu, Yuzhong Chen, Menghai Pan, Huiyuan Chen, Mahashweta Das, Hao Yang, and
 435 Hanghang Tong. Kernel ridge regression-based graph dataset distillation. In *Proceedings*
 436 *of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, KDD
 437 '23, page 2850–2861, New York, NY, USA, 2023. Association for Computing Machinery.
 438 ISBN 9798400701030. doi: 10.1145/3580305.3599398. URL [https://doi.org/10.1145/](https://doi.org/10.1145/3580305.3599398)
 439 [3580305.3599398](https://doi.org/10.1145/3580305.3599398).
- 440 Donggeun Yoo and In So Kweon. Learning loss for active learning. In *Proceedings of*
 441 *the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June
 442 2019.
- 443 Linfeng Zhang, Chenglong Bao, and Kaisheng Ma. Self-distillation: Towards efficient and
 444 compact neural networks. *IEEE Transactions on Pattern Analysis and Machine Intelli-*
 445 *gence*, 44(8):4388–4403, 2022. doi: 10.1109/TPAMI.2021.3067100.

446 Appendix A. Clinical Dataset

447 We utilized two linked datasets, a medical records dataset and a mental health outpatient
 448 questionnaire dataset. Table 3 shows available information in the medical record’s dataset.
 449 The dataset includes several categorical features, which are **triage level**, describing the type
 450 and severity of a patient’s initial symptoms, **visit disposition code**, which indicates how a
 451 patient was discharged from ambulatory care after registration, **service utilization**, which
 452 refers to the specific health professional services accessed during the patient’s visit, and
 453 **most responsible diagnosis code** that identifies the primary, most clinically significant
 454 problem determined by the healthcare provider during service utilization.

455 The patient graph consists of three feature categories: visit (including age, triage, and
 456 disposition), service utilization (based on service code), and diagnosis (using diagnosis code).
 457 Every visit includes at least one diagnosis, and each diagnosis is associated with at least
 458 one service. Connections extracted for each visit are: *visit-visit* (chronologically ordered),
 459 *visit-diagnosis*, and *diagnosis-service*. Figure 3 shows a sample EHR graph.

460 Table 4 shows the subset of the questions used in the mental health questionnaires
 461 dataset. The questions were selected through an ablation study designed to identify the
 462 most influential questions for the downstream task. The question categories have been used
 463 to create nodes and connections between them. We utilized the timeline of events in the
 464 questionnaire to augment the EHR patient graph with the questionnaire responses as shown
 465 in Figure 4.

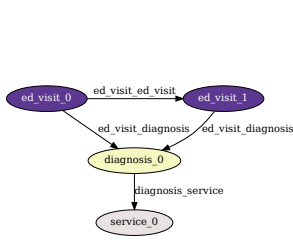


Figure 3: A sample patient graph using EHR data.

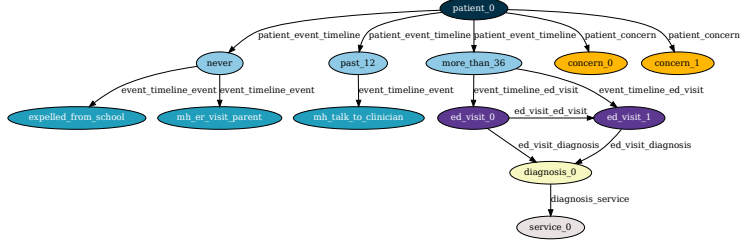


Figure 4: A sample augmented patient graph using EHR and questionnaire data.

Table 3: Medical dataset information.

Attribute	Type	Insight	Example
Patient age at the time of visit	Numeric	Min: 4 Max: 17 Mean: 12.5 Std: 3.9	12
ED admission date	Date	Min: 2006-10-16 Max: 2021-07-31	-
Triage level	Categorical	6 Categories	Emergent
Visit disposition code	Categorical	16 Categories	Intra-facility transfer to the ED
Utilized service code	Categorical	35 Categories	Orthopedic Surgery
Diagnosis code	Categorical	707 Categories	Allergic Purpura

Table 4: The subset of selected questionnaires, along with their category.

Question Category	Question
Patient Specific	Are you taking any medication or pills prescribed for mental health concerns?
	What sex were you assigned at birth, on your original birth certificate?
	How do you describe yourself (Male, Female, Transgender, Do not identify as female, male or transgender, Don't know)?
	Are you currently living in a friend's house or apartment?
	Do you think of yourself as (Straight, Gay or Lesbian, Bisexual, Transgender, transsexual or gender non-conforming, Don't know)?
Events	Were you born in Canada?
	During the past 12 months, did you visit an emergency room about concerns regarding your mental health?
	During the past 12 months, have you been suspended or expelled from school?
Concerns	During the past 12 months, did you see or talk to a doctor, psychologist, psychiatrist, or counsellor about concerns regarding your mental health?
	I am concerned about my physical appearance I am concerned about physically hurting myself