

# Rethinking the Effectiveness of Graph Classification Datasets in Benchmarks for Assessing GNNs

Zhengdao Li<sup>1,2</sup>, Yong Cao<sup>3</sup>, Kefan Shuai<sup>1</sup>, Yiming Miao<sup>1\*</sup> and Kai Hwang<sup>1</sup>

<sup>1</sup>The Chinese University of Hong Kong, Shenzhen, China

<sup>2</sup>Guangzhou University, Guangzhou, China

<sup>3</sup>Huazhong University of Science and Technology, Wuhan, China

{zhengdaoli, kefanshuai}@link.cuhk.edu.cn

yongcao2018@gmail.com

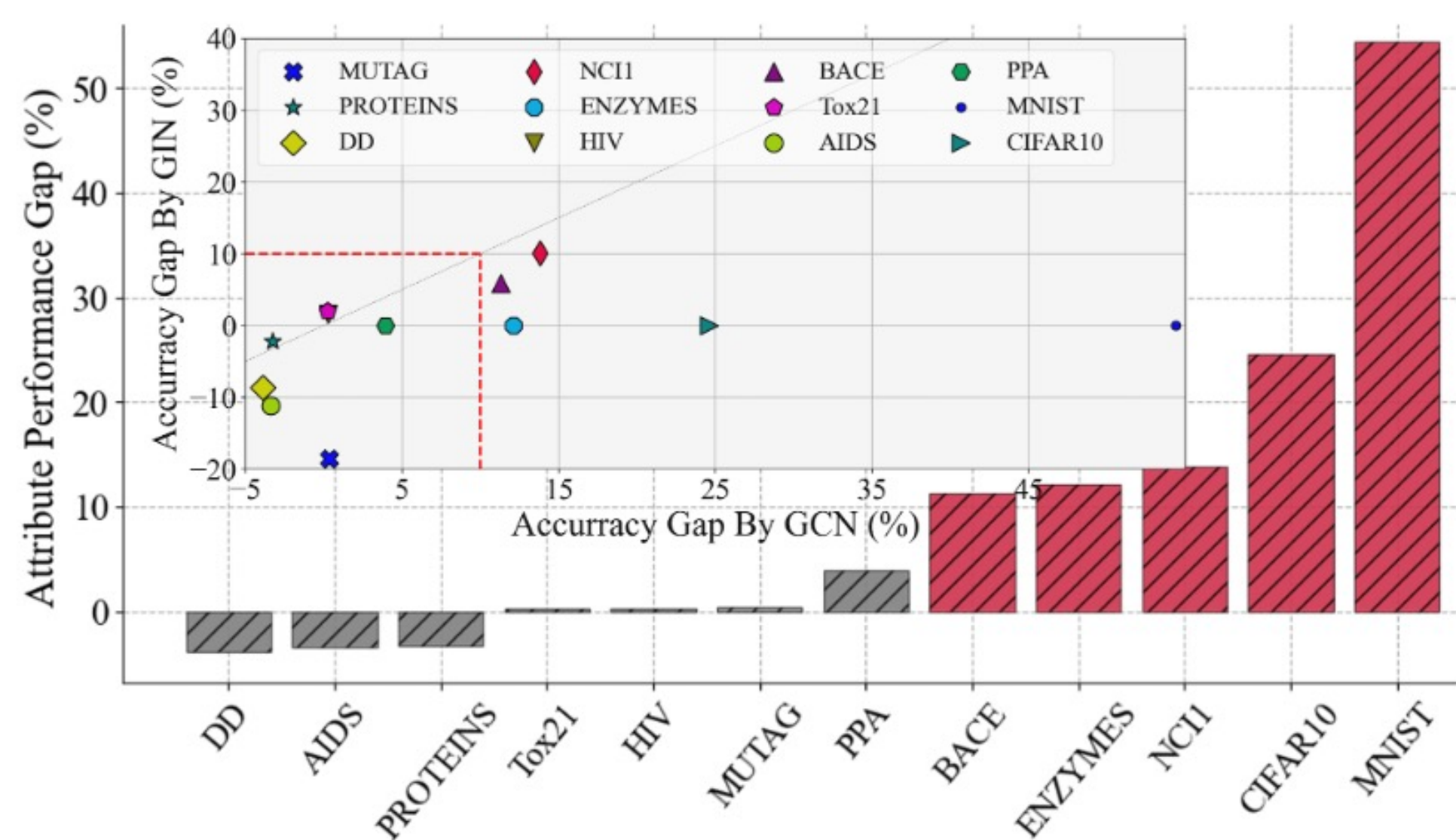
{miaoyiming, hwangkai}@cuhk.edu.cn

## Research Questions:

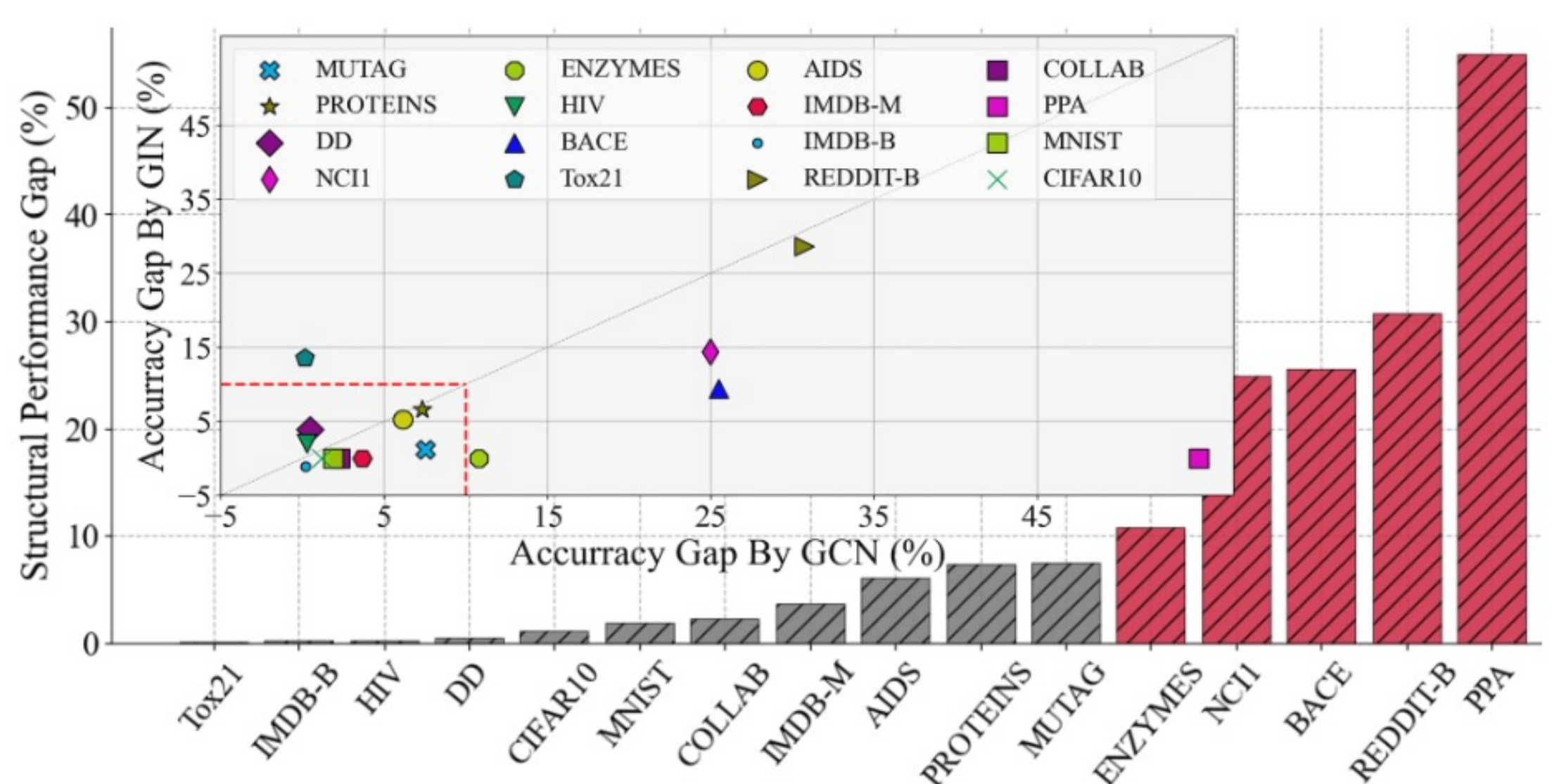
**RQ1:** Can commonly used graph classification datasets serve the benchmarking purpose which is to effectively distinguish advancements of GNNs compared with other methods?

**RQ2:** How to measure the effectiveness of existing graph classification datasets?

## Observations:



(a) Attributed accuracy gap



(b) Structural accuracy gap

Figure 1: The performance gaps on 16 graph classification datasets are categorized into two types: *Ineffective* (gray) and *Effective* (red) benchmarks. These are sorted in ascending order based on the size of the performance gap. An empirical threshold of 10% is used for categorization, as observed in the inner box of each figure. This box represents the distribution of the accuracy gap for GCN and GIN.

## Definition1 ---- Performance gap

Given a dataset  $D$ , a baseline method  $\mathcal{M}_{\text{type}}^{\text{Baseline}}$ , and a graph-based method  $\mathcal{M}_{\text{type}}^{\text{Graph}}$ , the performance gap  $\delta_{\text{type}}(D, R, \mathcal{M}_{\text{type}}^{\text{Graph}}, \mathcal{M}_{\text{type}}^{\text{Baseline}})$  (simply denoted by  $\delta_{\text{type}}$ ) between baseline and graph-based method is defined as:

$$\delta_{\text{type}} \triangleq R(D, \mathcal{M}_{\text{type}}^{\text{Graph}}) - R(D, \mathcal{M}_{\text{type}}^{\text{Baseline}}), \text{ type} \in \{\text{S}, \text{A}\},$$

## Limitations of the performance gap:

$$R(D_1, \mathcal{M}^{\text{GNN}}) = 100\% \\ = 10\%$$

$$R(D_1, \mathcal{M}^{\text{Baseline}}) = 90\%$$

$$R(D_2, \mathcal{M}^{\text{GNN}}) = 60\% \\ = 10\%$$

$$R(D_2, \mathcal{M}^{\text{Baseline}}) = 50\%$$



## Proposed metric:

### Definition 2 ---- Dataset Effectiveness

Given a graph classification dataset  $D$  which has  $|Y|$  classes, and the performance gap  $\delta_{\text{type}}(D)$  between two methods  $\mathcal{M}^1$  and  $\mathcal{M}^2$ , the  $\mathcal{E}$  to quantify the discriminating degree of  $\mathcal{M}^1$  and  $\mathcal{M}^2$  is defined as follows:

$$\mathcal{E}(D) = \sum_{\text{type} \in \{S, A\}} \frac{|\delta_{\text{type}}(D)|}{R^*(|Y| - 1)} \cdot \frac{1 - R^*}{1 - |Y|^{-1}},$$

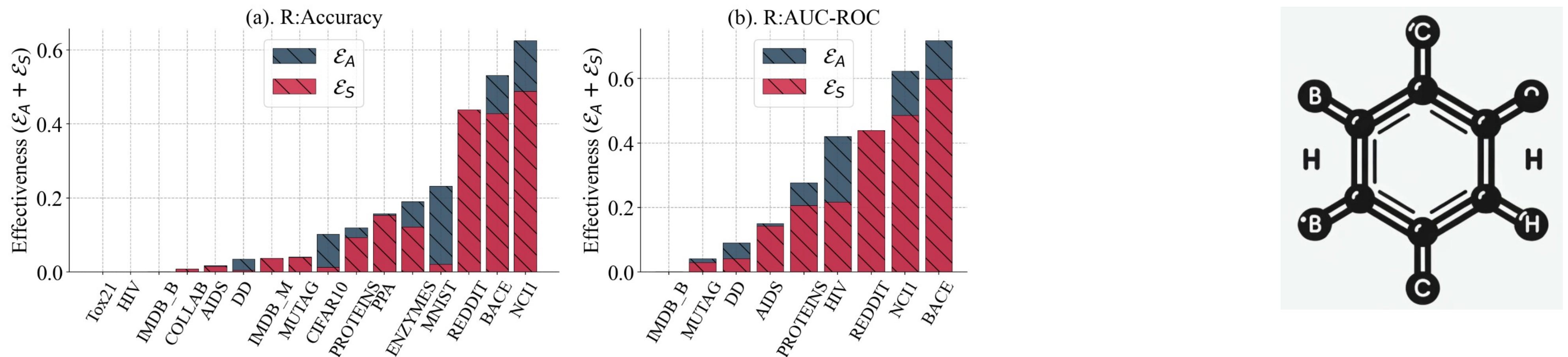


Figure 3: Effectiveness using Accuracy metric and AUC-ROC metric in terms of structural type and attributed type.

## Investigation of low effectiveness

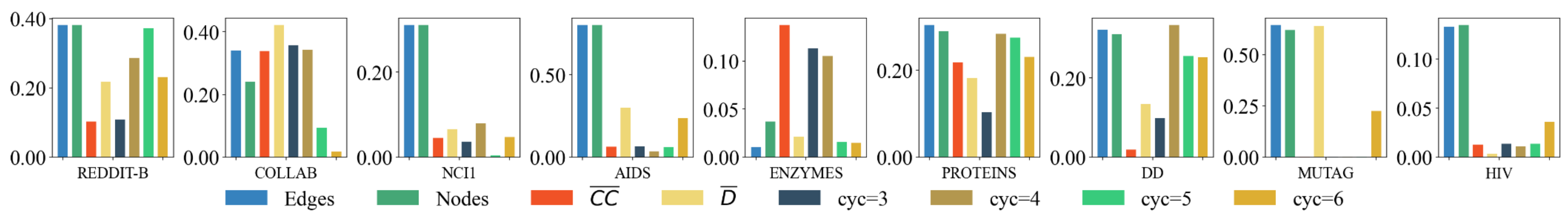


Figure 4: Correlations between graph property sequences and class labels on 9 real-world datasets.

## Synthetic dataset generation

### Theorem 1:

Given a set of property variables  $\{\mathcal{P}_i\}_{i=1}^K$ , each  $\mathcal{P}_i$  follows a Gaussian distribution  $\mathcal{N}(\mu_k, \sigma_k)$  or Uniform distribution  $\mathcal{U}(a_k, b_k)$ , and given corresponding Pearson correlation coefficients  $\{r_i\}_{i=1}^K$  with label variable  $\mathcal{Y}$ , with the constraint  $\sum_{i=1}^K r_i^2 \leq 1$ , then we have:

$$\mathcal{Y} = \sigma_{\mathcal{Y}} \left( \sum_{i=1}^K n_i r_i + n_0 \sqrt{1 - \sum_{i=1}^K r_i^2} \right), \quad (1)$$

where  $\sigma_{\mathcal{Y}}$  is any desired standard deviation, and each  $n_i$  is mutually independent and follows the same distribution as the corresponding  $\mathcal{P}_i$  with the same mean value  $\mu_i$  but with standard deviation equals to 1. (The proof is based on Cholesky decomposition of a given covariance matrix.)

## Prediction of Effectiveness

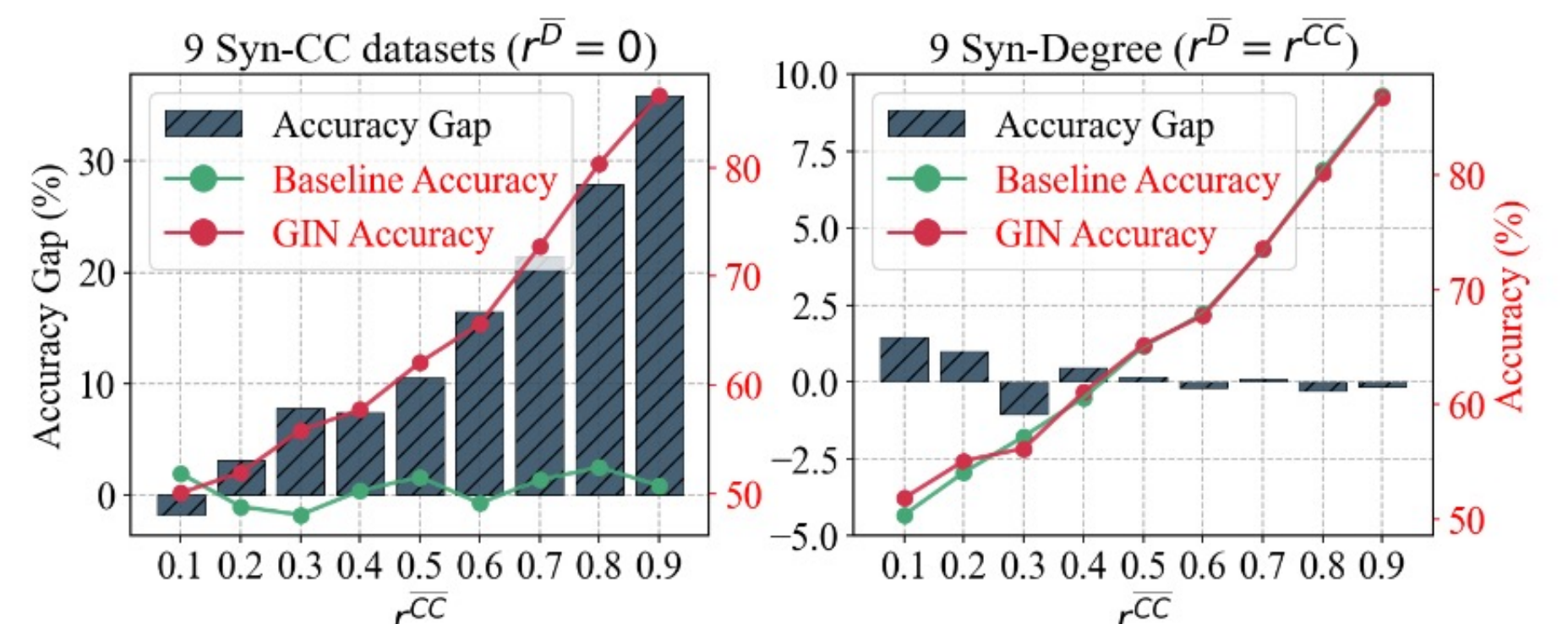


Figure 6: Controllable performance gaps by two types of synthetic datasets.

### Algorithm 1: Controllable dataset construction

- 1 **Input:**  $\{r_k\}_{k=1}^K$ , number of labels  $C$ ,  $\{\mathcal{P}_k\}_{k=1}^K \sim \mathcal{N}(\mu_k, \sigma_k)$  or  $\mathcal{U}(a_k, b_k)$ ;
- 2 **Output:** Dataset  $\mathbb{D}$  with size  $N$ ;
- 3 **for**  $k = 1$  **to**  $K$  **do**
- 4     Sample  $n_k \sim \mathcal{N}(0, \sigma_k)$  or  $\mathcal{U}(-\sqrt{3}, \sqrt{3})$ ;
- 5      $\mathbb{P}_k \leftarrow \mu_k + \sigma_k n_k$  or  $\frac{a_k + b_k}{2} + \sqrt{\frac{b_k - a_k}{12}} n_k$ ;
- 6 **end**
- 7 Calculate  $\mathcal{Y}$  by the Eq. 2 ;
- 8  $\mathbb{Y} \leftarrow \text{ROUND}(\text{NORM}(\mathcal{Y}) * C)$ ;
- 9  $\mathbb{D} \leftarrow \{(g_i, y_i)\}_{i=1}^N$ ,
- 10 where each graph  $g_i$  has properties  $\{\mathbb{P}_k[i]\}_{k=1}^K$ , and corresponding label  $y_i = \mathbb{Y}[i]$ ;

Regressor	Real-world datasets		Synthetic-CC datasets	
	Pearson	P-Value	Pearson	P-Value
Ridge	$0.80 \pm 0.09$	$\leq 1 \times 10^{-6}$	$0.87 \pm 0.03$	$\leq 1 \times 10^{-6}$
SVR	$0.80 \pm 0.09$	$\leq 1 \times 10^{-6}$	$0.89 \pm 0.04$	$\leq 1 \times 10^{-6}$
RF	$0.89 \pm 0.03$	$\leq 1 \times 10^{-6}$	$0.87 \pm 0.06$	$\leq 1 \times 10^{-6}$

Table 3: Summary of regression results