
BN-Pool: a Bayesian Nonparametric Approach to Graph Pooling

Daniele Castellana¹, Filippo Maria Bianchi^{2,3}

¹ Dept. of Statistics, Computer Science and Applications, Università degli Studi di Firenze

² Dept. of Mathematics and Statistics, UiT the Arctic University of Norway

³ NORCE, Norwegian Research Centre AS

daniele.castellana@unifi.it, filippo.m.bianchi@uit.no

Abstract

We introduce BN-Pool, the first clustering-based pooling method for Graph Neural Networks (GNNs) that adaptively determines the number of supernodes in a coarsened graph. By leveraging a Bayesian non-parametric framework, BN-Pool employs a generative model capable of partitioning graph nodes into an unbounded number of clusters. During training, we learn the node-to-cluster assignments by combining the supervised loss of the downstream task with an unsupervised auxiliary term, which encourages the reconstruction of the original graph topology while penalizing unnecessary proliferation of clusters. This adaptive strategy allows BN-Pool to automatically discover an optimal coarsening level, offering enhanced flexibility and removing the need to specify sensitive pooling ratios. We show that BN-Pool achieves superior performance across diverse benchmarks.

graph properties through a hierarchy of coarsened graphs. Those, are crucial for building deep GNNs for tasks such as graph classification [1], node classification [2, 3], graph matching [4], and spatio-temporal forecasting [5, 6].

While pooling in GNNs serves a similar purpose as in CNNs, its implementation is more challenging due to the irregular and non-Euclidean structure of the graphs. Popular (and often better performing [7]) graph pooling methods generate coarsened graphs by aggregating nodes into clusters. However, these approaches typically require a predefined number of clusters, i.e., nodes in the coarsened graph, which is difficult to set in advance. Moreover, enforcing the same fixed number of clusters across all graphs, regardless of their original size, results in all coarsened graphs having the same number of nodes. This rigidity hinders the ability of the model to adapt dynamically to the graph structure, resulting in redundancies and reducing its effectiveness in datasets with significant variability in the size of the graphs.

1 Introduction

Graph Neural Networks (GNNs) have emerged as powerful tools to solve various tasks involving graph-structured data, such as node classification, graph classification, and link prediction. Despite their success, one of the persistent challenges in GNNs is efficiently handling large-scale graphs while preserving their structural and feature information.

Pooling is a widely used technique in deep learning architectures such as Convolutional Neural Networks (CNNs) to progressively distill global properties from the data by summarizing spatially contiguous information. Similarly, GNNs use pooling layers to summarize the information on the graph. Of particular interest to us are the pooling layers that gradually extract global

To overcome these limitations, we introduce Bayesian Non-parametric Pooling (BN-Pool): a novel pooling technique for GNN based on a Bayesian Non-Parametric (BNP) approach. Our method defines a generative process for the adjacency matrix of the input graph where the probability of having a link between two nodes depends on their cluster membership, ensuring that clusters reflect the graph topology. Thanks to the BNP approach, the number of clusters is not fixed in advance but is adapted to the input graph. Within our Bayesian framework, the clustering function is the posterior of the cluster membership given the input graph. In this work, we approximate the posterior by employing a GNN; on the one hand, this permits to capture complex relations that usually appear between the hidden and the observable variables; on the other hand, we can condition the posterior on the node (and potentially edge) features, and on the downstream task.

Preprint.

The GNN parameters are trained by optimizing two complementary objectives: one defined by the losses of the downstream task (e.g., graph classification), the other defined by an unsupervised auxiliary loss that derives from the probabilistic nature of the model.

In the following sections, we will detail the theoretical foundations of our approach, demonstrate its efficacy through empirical evaluation, and compare it against state-of-the-art methods to highlight its unique advantages.

2 Background

2.1 BAYESIAN NON-PARAMETRIC

The BNP framework [8] aims to build non-parametric models by applying Bayesian techniques. The term *non-parametric* indicates the ability of a model to adapt its size (i.e., the number of parameters) directly to data. In contrast, in the *parametric* approach the model size is fixed in advance by setting some hyper-parameters.

The BNP literature relevant to our work relates to the families of Dirichlet Process (DP) [9]. In its most essential definition, a DP is a stochastic process whose samples are categorical distributions of infinite size. Thus, in the same way as the Dirichlet distribution is the conjugate prior for the categorical distribution, the DP is the conjugate prior for infinite discrete distributions.

A classical usage of DP is in the definition of mixture models which allow an infinite number of components, where the DP is used as the prior distribution over the mixture weights. The key of DP is its clusterisation property: even if there is an infinite number of components available, the DP tends to use the components that have been already used. We refer the reader to [Appendix A](#) for an introduction to the DP.

2.2 GRAPH NEURAL NETWORKS

Let $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ be a graph with node features $\mathbf{X}^0 \in \mathbb{R}^{N \times F}$, where $|\mathcal{V}| = N$. Each row $\mathbf{x}_i^0 \in \mathbb{R}^F$ of the matrix \mathbf{X}^0 represents the initial node feature of the node i , $\forall i = 1, \dots, N$. Through the MP layers, a GNN implements a local computational mechanism to process graphs [10]. Specifically, each feature vector \mathbf{x}_v is updated by combining the features of the neighboring nodes. After l iterations, \mathbf{x}_v^l embeds both the structural information and the content of the nodes in the l -hop neighborhood of v . With enough iterations, the feature vectors can be used to classify the nodes or the entire graph. More rigorously, the output of the l -th layer of a MP-GNN is:

$$\mathbf{x}_v^l = \text{COMB}^{(l)} \left(\mathbf{x}_v^{l-1}, \text{AGGR}^{(l)}(\{\mathbf{x}_u^{l-1}, u \in \mathcal{N}[v]\}) \right) \quad (1)$$

where $\text{AGGR}^{(l)}$ is a function that aggregates the node features from the neighborhood $\mathcal{N}[v]$ at the $(l-1)$ -th

iteration, and $\text{COMB}^{(l)}$ is a function that combines the own features with those of the neighbors.

The most simple GNN architectures are “flat” and consist of a stack of Message Passing (MP) layers followed by a final readout [11]. For graph-level tasks, such as graph classification and regression, the readout includes a global pooling layer that combines all the node features at once by taking e.g., their sum or average. Such an aggressive pooling operation often fails to effectively extract the global graph properties necessary for the downstream task. On the other hand, GNN architectures that alternate MP with graph pooling layers can gradually distill information into “hierarchical” graph representations.

2.3 GRAPH POOLING

Graph pooling methods can be broadly described through Select-Reduce-Connect (SRC), which provides a general framework to describe different graph pooling operators [12]. According to SRC, a pooling operator, denoted as $\text{POOL} : (\mathbf{A}, \mathbf{X}) \rightarrow (\mathbf{A}_p, \mathbf{X}_p)$, is decomposed into three sub-operators:

- **Select (SEL)**: maps the original nodes of the graph to a reduced set of nodes, called supernodes. Often, the mapping can be represented by a selection matrix $\mathbf{S} \in \mathbb{R}^{N \times K}$, where N and K are the number of nodes and supernodes, respectively.
- **Reduce (RED)**: generates the features $\mathbf{X}_p \in \mathbb{R}^{K \times F}$ of the supernodes based on the selection matrix and the original node features.
- **Connect (CON)**: constructs the new adjacency matrix $\mathbf{A}_p \in \mathbb{R}_{\geq 0}^{K \times K}$ based on the selection matrix and the original topology.

Different pooling methods are obtained by a specific implementation of these operators and can be broadly categorized into three main families: score-based, one-every- K , and soft-clustering methods.

Score-Based methods compute a score for each node using a trainable function in their SEL operator. Nodes with the highest scores become the supernodes of the pooled graph. Representatives such as Top- k Pooling (Top- k) [2, 13], ASAPool [14], SAGPool [15], PanPool [3], TAPool [16], CGIPool [17], and IPool [18] primarily differ in how they compute the scores or in the auxiliary tasks they optimize to improve the quality of the pooled graph. These methods are computationally efficient and can dynamically adapt the size of the pooled graph, e.g., $K_i = \kappa N_i$, where κ is the pooling ratio and N_i and K_i are the sizes of the i -th graph before and after pooling. Score-based methods tend to retain neighbouring nodes that have similar features. As a result, entire parts of the graph are under-represented after pooling, reducing

the performance in tasks where all the graph structure should be preserved.

One-Every- K methods pool the graph by uniformly subsampling nodes, extending the concept of one-every- K to irregular graph structures. They are typically efficient and perform pooling by optimizing graph-theoretical objectives, such as spectral clustering [19], maxcut [20], and maximal independent sets [21]. However, these methods lack flexibility because their SEL operator neither accounts for node or edge features nor can be influenced by the downstream task’s loss. Despite they can adapt the size of the pooled graph K_i to the original graph size N_i , they cannot specify the pooling ratio κ explicitly, which is determined *a-priori* by the graph-theoretical objective.

Soft-Clustering methods have a SEL operator that computes a soft-clustering matrix \mathbf{S} , which assigns each node to multiple supernodes with different degrees of membership. Representatives such as Diffpool [22], MinCut Pool (MinCut) [23], and Structpool [24], leverage flexible trainable functions guided by auxiliary losses to compute the soft assignments from the node features. The auxiliary losses ensure that the partition is consistent with the graph topology and that the clusters are well-formed, e.g., the assignments are sharp and the clusters balanced. While soft-clustering methods generally achieve high performance due to their flexibility and ability to retain information from the entire graph, they face a primary limitation. They require to predefine the size K of *every* pooled graphs, which is fixed for each graph i regardless of its size N_i . A typical choice is to set $K = \kappa N$, where \bar{N} is the average size of all the graphs in the dataset. Clearly, this might not work well in datasets where the graphs’ size varies too much, especially if there are graphs where $N_i < \kappa \bar{N}$. In those cases, the pooling operator *expands* rather than coarsening the graph, which goes against the principle of graph pooling. Finally, while the possibility of specifying the pooling ratio κ offers a greater flexibility to soft-clustering and score-based methods, it might be a difficult hyperparameter to tune.

3 Method

We propose a novel soft-clustering pooling operator whose SEL function addresses the main drawbacks of existing soft-clustering methods by learning, for each graph i , a pooled graph with a variable number of supernodes K_i . We refer to our proposal as BN-Pool since it is grounded in the Bayesian non-parametric theory. In the following, we present the method by considering only a single graph to ease the notation.

BN-Pool assumes that the adjacency matrix \mathbf{A} of the input graph is generated by a process similar to the Stochastic Block Model (SBM): each node u is associated with a vector π_u whose entries indicate the probability

that u belongs to a given cluster. The edges are generated according to a block matrix \mathbf{K} whose entry K_{ij} represents the unnormalised log-probability of occurrence of a directed edge from a node in cluster i to a node in cluster j . Differently from the SBM, we relax the requirement of specifying the number of clusters in advance and leverage the DP to define a prior over an infinite number of clusters. It is worth mentioning that, even if there is an infinite number of clusters, only a "small" number of them is used due to the clusterisation property of the DP (see Appendix A for details).

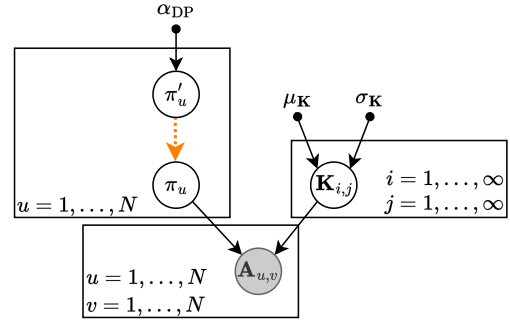


Figure 1: Graphical representation of BN-Pool in plate notation. The orange dashed arrow represents a deterministic computation.

By exploiting the stick-breaking construction of DPs, we define the generative process of BN-Pool as:

$$K_{ij} \sim p(K_{ij}) = \begin{cases} \mathcal{N}(\mu_{\mathbf{K}}, \sigma_{\mathbf{K}}) & \text{if } i = j \\ \mathcal{N}(-\mu_{\mathbf{K}}, \sigma_{\mathbf{K}}) & \text{if } i \neq j \end{cases}, \quad (2)$$

$$\pi'_{ui} \sim p(\pi'_{ui}) = \text{Beta}(1, \alpha_{\text{DP}}), \quad (3)$$

$$\pi_{ui} = \pi'_{ui} \prod_{j=1}^{i-1} (1 - \pi'_{uj}), \quad p_{uv} = \sigma(\pi_u^\top \mathbf{K} \pi_v), \quad (4)$$

$$\mathbf{A}_{uv} \sim p(\mathbf{A}_{uv}) = \text{Bernoulli}(p_{uv}), \quad (5)$$

where $u, v \in \mathcal{V}$ are nodes in the input graph, $i, j \in \mathbb{N}$ are cluster indexes, and $\sigma(\cdot)$ is the sigmoid function; the hyper-parameters $\alpha_{\text{DP}} \in \mathbb{R}^+$, $\mu_{\mathbf{K}} \in \mathbb{R}^+$, $\sigma_{\mathbf{K}} \in \mathbb{R}^+$ define the shape of the prior distributions. The prior distribution on the matrix \mathbf{K} defined in equation 2 encodes our assumption that most of the edges link nodes of the same group. The generative process is schematized in Figure 1.

The BNP setting makes the posterior computation intractable and approximations are required to perform training. We rely on a truncated variational approximation of the posterior [25]: even if there is an infinite number of clusters, we truncate the posterior by considering a finite value C representing the maximum number of clusters. It is worth highlighting that this does not imply that the model has a fixed number of clusters but, rather, that the model will choose a suitable number of non-empty (i.e., active) clusters $K_i < C$ for the i -th graph.

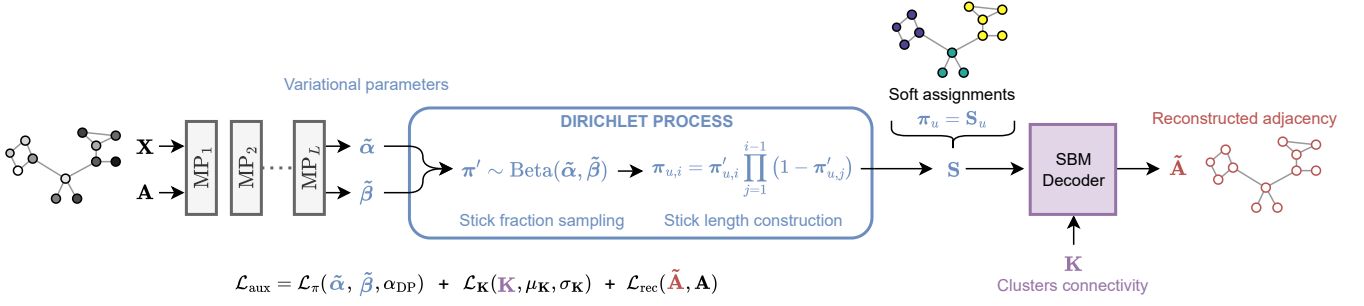


Figure 2: The SEL operation of and the components of the auxiliary loss.

We follow the classical mean-field approximation¹ and define two variational distributions: one to model the posterior of the stick fractions π' , and one to model the posterior of the model parameter \mathbf{K} . Note that the cluster assignment vector π is fully determined by the stick-breaking construction given the stick fractions π' . The posterior approximation can be detailed as follows:

$$q(\pi'_{ui}) = \text{Beta}(\tilde{\alpha}_{ui}, \tilde{\beta}_{vi}), \quad (6)$$

$$q(\mathbf{K}_{ij}) = \mathcal{N}(\tilde{\mu}_{ij}, \epsilon), \quad (7)$$

where $\tilde{\alpha}_{ui}, \tilde{\beta}_{vi} \in \mathbb{R}^+$ for all $u \in \mathcal{V}, i \in \{1, \dots, C\}$, and $\tilde{\mu}_{ij}$ for all $i, j \in \{1, \dots, C\}$ are the variational parameters. The value of ϵ is fixed a priori and it is not optimised during the training.

While $\tilde{\mu}_{ij}$ are free parameters that we optimize directly, we employ a GNN with parameters Θ to estimate $\tilde{\alpha}$ and $\tilde{\beta}$:

$$\tilde{\alpha}, \tilde{\beta} = \text{GNN}_{\Theta}(\mathbf{X}, \mathbf{A}). \quad (8)$$

On the one hand, the GNN allows for representing complex relations between hidden and observable variables that usually appear in the posterior distribution. On the other hand, we can condition the posterior on the graph topology, on the node (and potentially edge) features, and on the downstream task at hand that drives the GNN optimization.

We summarize the proposed architecture in Figure 2: the GNN employed to estimate the posterior acts as the *encoder* in the classic Variational Auto-Encoder (VAE) approach, while the SBM is the *decoder* which reconstructs the adjacency matrix of the input graph. The soft assignments in \mathbf{S} are the latent representation π_u for each node u , which follow a DP by allowing an infinite number of clusters.

¹It assumes the variational distribution factorises over the latent variables: $p(\pi', \mathbf{K} | \mathbf{X}) \approx q(\pi', \mathbf{K}) \approx q(\pi')q(\mathbf{K})$.

3.1 TRAINING PROCEDURE

The parameters $\{\Theta, \tilde{\mu}\}$ are learned by maximising the Evidence Lower-BOund (ELBO):

$$\begin{aligned}
 \log p(\mathbf{A}) \geq & \underbrace{\sum_u \sum_v \mathbb{E}_{q(\pi')q(\mathbf{K})} [\log p(\mathbf{A}_{uv} | \pi, \mathbf{K})]}_{-\mathcal{L}_{\text{rec}}} \\
 & - \underbrace{\sum_u \sum_i D_{\text{KL}}(q(\pi'_{ui}) | p(\pi'_{ui}))}_{-\mathcal{L}_{\pi}} \\
 & - \underbrace{\sum_i \sum_j D_{\text{KL}}(q(\mathbf{K}_{ij}) | p(\mathbf{K}_{ij}))}_{-\mathcal{L}_{\mathbf{K}}}.
 \end{aligned} \quad (9)$$

The first term in equation 9 is the *reconstruction loss* that measures how good is the model at reconstructing the adjacency matrix. The last two terms instead measure the distances between the prior and the variational distributions, and they act as a regularisation. While the reconstruction loss \mathcal{L}_{rec} has a straightforward interpretation, we can think of \mathcal{L}_{π} as the total cost to pay to have a certain number of clusters active. Hence, \mathcal{L}_{π} reflects the clusterisation property of the DP in reusing non-empty clusters. On the other hand, $\mathcal{L}_{\mathbf{K}}$ penalizes the discrepancy from the connectivity across clusters described by the SBM prior.

In practice, instead of maximising the ELBO in Eq. 9, we train the model by minimising the following loss:

$$\mathcal{L}_{\text{aux}} = \frac{1}{N} \mathcal{L}_{\text{rec}} + \eta \frac{1}{N} \mathcal{L}_{\pi} + \frac{1}{N} \mathcal{L}_{\mathbf{K}}, \quad (10)$$

where N is the number of nodes in the input graph and it is used to rescale the losses, while η is a hyperparameter which balances the contrasting effect of \mathcal{L}_{rec} and \mathcal{L}_{π} . The interplay between all the loss terms is crucial for an effective adaptive nonparametric method. The normalization and scaling parameters avoid a dominance of the KL divergence and have already been applied on VAEs [26, 27]. We refer to the loss in equation 10 as *auxiliary* since during the pooling it will be combined with the supervised loss of the downstream task.

The training is performed by employing the Stochastic Gradient Variational Bayes (SGVB) framework [28], where the expectation in the reconstruction loss is approximated with a Monte Carlo estimate of the binary cross-entropy between the true edges and the probabilities predicted by the model:

$$\mathcal{L}_{\text{rec}} \approx \sum_{t=1}^T \sum_u \sum_v -\mathbf{A}_{uv} \log p_{uv}^t - (1 - \mathbf{A}_{uv}) \log(1 - p_{uv}^t), \quad (11)$$

where T is the number of samples used for the Monte Carlo approximation and $p_{uv}^t = \sigma(\sum_i \sum_j \pi_{ui}^t \tilde{\mu}_{ij} \pi_{vj}^t)$ being the values π_u^t and π_v^t the t -th samples of the soft assignments for the node u and v . The sampling step to approximate \mathcal{L}_{rec} is not differentiable and prevents the gradient to be back-propagated to the GNN parameters Θ . A common approach to solve this issue is the reparametrisation trick [28], which, however, cannot be applied to the Beta distribution [29]. In BN-Pool, we back-propagate the information by approximating the pathwise gradient of the sampled values w.r.t. the distribution parameters² [30].

To reduce the stochasticity of the approximation, we assume that the variational distribution $q(\mathbf{K})$ has a low variance (i.e., $\varepsilon \rightarrow 0$ in Eq. 7) and directly use the variational parameter $\tilde{\mu}$ rather than sampling the cluster connectivity from its variational distribution. Finally, we initialise the GNN parameters Θ by using the default initialisation of the backend, while the variational parameter $\tilde{\mu}$ of the cluster connectivity matrix is initialised by setting the element on-diagonal (off-diagonal) equals to $\eta_{\mathbf{K}}$ ($-\eta_{\mathbf{K}}$), where $\eta_{\mathbf{K}}$ is an hyperparameter.

3.2 PRIOR HYPERPARAMETERS INTERPRETATION

To fully define BN-Pool model, we have to specify three hyperparameters: α_{DP} , $\mu_{\mathbf{K}}$ and $\sigma_{\mathbf{K}}$. The probabilistic nature of our method allows for a direct interpretation that facilitates their tuning.

The value of $\alpha_{\text{DP}} \in \mathbb{R}^+$ defines the shape of the prior over the cluster assignments; in particular, it specifies the concentration of the DP. To understand the effect of α_{DP} , we recall that the loss \mathcal{L}_{π} is the cost to pay to have a certain number of clusters active. The value α_{DP} is inversely proportional to the price to activate a new cluster: low values force the model to use a few clusters (only one in the extreme case). Conversely, higher values do not penalise the model when it uses more clusters to reduce the reconstruction loss. Note that in practice we truncate the posterior to at most C clusters, meaning that too high values of α_{DP} create degenerate solutions where the last cluster is always used.

²This approximation is already implemented in the PyTorch library. See Appendix B for more details about our implementation.

The other two hyperparameters $\mu_{\mathbf{K}} \in \mathbb{R}^+$ and $\sigma_{\mathbf{K}} \in \mathbb{R}^+$ specify the prior over the cluster connectivity matrix \mathbf{K} which affects the reconstruction loss. Again, the most intuitive way to understand the effect of \mathbf{K} is in terms of costs: if the value \mathbf{K}_{ij} is positive (negative), the price of creating an edge between a node in cluster i and a node in cluster j is low (high). Thus, to encode our prior belief that most of the edges appear between nodes in the same cluster, we impose that the elements on the diagonal are positives with value $\mu_{\mathbf{K}}$ (i.e., intra-cluster edges are cheap), while the off-diagonal elements are negatives with value $-\mu_{\mathbf{K}}$ (i.e., inter-cluster edges are costly). The hyperparameter $\sigma_{\mathbf{K}}$ controls the strength of the prior: the lower the more the posterior matches the prior rather than the data.

The values of $\mu_{\mathbf{K}}$ and $\sigma_{\mathbf{K}}$ also affect the number of active clusters. For example, the degenerate solution that assigns all the nodes to the first cluster satisfies the clusterisation property of the DP. However, by referring at Eq. 11, this means paying $-\log(1 - \sigma(\tilde{\mu}_{11})) = -\log \sigma(-\tilde{\mu}_{11})$ every time $\mathbf{A}_{uv} = 0$. If the posterior matches our prior (i.e., $\tilde{\mu}_{11} \approx \mu_{\mathbf{K}}$), this results in a great cost since $\mu_{\mathbf{K}} \gg 0$ implies $-\log \sigma(-\mu_{\mathbf{K}}) \gg 0$; thus, the model will likely prefer to reduce \mathcal{L}_{rec} at the price of having more clusters, i.e., a larger \mathcal{L}_{π} .

Finally, we note that while the other hyperparameters (truncation level C , number of samples T and initialisation of the variational parameters Θ and $\eta_{\mathbf{K}}$) influence the training procedure, they do not affect the model definition.

3.3 SEL, RED, CON

We conclude by casting BN-Pool into the SRC framework. For each graph i the SEL operator generates a cluster assignment matrix $\mathbf{S}_i \in \mathbb{R}^{N \times C}$ where the first K_i columns contain non-zero values. The entry $s_{uj} = \pi_{uj}$ represents the membership of node u to cluster j . The RED and CON functions are implemented as in other soft-clustering methods. The RED function is $\mathbf{X}_p = \mathbf{S}^T \mathbf{X}$, where \mathbf{X} is the feature matrix of the original graph, and \mathbf{X}_p are the features of the pooled graph. The CON function is implemented as $\tilde{\mathbf{A}}_p = \mathbf{S}^T \mathbf{A} \mathbf{S}$. Following [23], we set the diagonal elements of $\tilde{\mathbf{A}}_p$ to zero to prevent that self-loops dominate the propagation in the MP layers after pooling and we symmetrically normalize it by the nodes' degree: $\mathbf{A}_p = \tilde{\mathbf{D}}_p^{-1/2} \tilde{\mathbf{A}}_p \tilde{\mathbf{D}}_p^{-1/2}$.

4 Related work

BN-Pool belongs to the family of Soft-Clustering pooling methods and the closest approach is Diffpool [22], which employs an auxiliary loss $\|\mathbf{A} - \mathbf{S} \mathbf{S}^T\|_F$ to align the assignments to the graph topology. In this work, we go

beyond the formulation of such a simple loss and define a whole generative process for the adjacency matrix.

Similar to our work is the Dirichlet Graph Variational Auto-Encoder (DGVAE) [31], which defines a VAE with a Dirichlet prior over the latent variables to cluster graph nodes. We extend DGVAE in two ways. First, we define a more flexible generative process for the adjacency matrix thanks to the block matrix \mathbf{K} . Second, we allow an infinite number of clusters by specifying a DP prior over the latent variables. Moreover, we do not rely on the Laplace approximation of the Dirichlet distribution, whose behaviour is similar to a Gaussian prior [32].

The Stick-Breaking Variational Auto-Encoder (SB-VAE) [33] shares our idea of specifying a non-parametric prior over the hidden variables by using a DP prior that leverages the stick-breaking construction, but does not focus on graphs. We also employ a different approximation of the posterior, which is based on pathwise gradients rather than the Kumaraswamy distribution.

Another work which shares some similarities with our method is [34], which introduces a sparse VAE for overlapping SBM. They also allow an infinite number of clusters, but use a different nonparametric prior: the Indian Buffet Process (IBP) [35]. The IBP is suitable to model multiple cluster membership, i.e., a node can belong to more than one cluster, which is not desirable in the context of pooling. Moreover, Mehta et al. define another dense latent variable with a Gaussian prior for each node to gain more flexibility during the generation process of the adjacency matrix. Instead, in BN-Pool all the information useful for the generation is encoded in the soft cluster assignments \mathbf{S} .

5 Experiments

The purpose of our experiments is twofold. Being BN-Pool the first BNP pooling method, we first analyse its ability in detecting communities on a single graph. Then, we test the effectiveness of BN-Pool in GNNs for graph classification, showing that it can achieve competitive performance w.r.t. other pooling methods. In all experiments, we use very simple GNN models to better appreciate the differences between the pooling methods. While GNNs with larger capacity can achieve SOTA performance, in a more complex model it is harder to disentangle the actual contribution of the pooling method. Details about the architectures, the training procedures, and the datasets are in Appendices C and D. The code is available online³.

³ <https://github.com/NGMLGroup/Bayesian-Nonparametric-Graph-Pooling>

5.1 COMMUNITY DETECTION

This task consists in learning a partition of the graph nodes in an unsupervised fashion, only based on the node features and the graph topology. The architecture used for clustering consists of a stack of MP layers that generate the feature vectors \mathbf{X}' . Those are processed by the SEL operator that produces the cluster assignments \mathbf{S} . Since clustering is an unsupervised task, the GNN is trained by minimizing only the auxiliary losses. Even if our primary focus is on graph pooling, this experiment allows us to evaluate the consistency between the node labels \mathbf{y} and the cluster assignments learned by minimizing only the auxiliary losses.

Clustering performance is commonly measured with Normalized Mutual Information (NMI), Completeness, and Homogeneity scores, which only work with hard cluster assignments. While the latter can be obtained by taking the argmax of a soft assignment, the discretisation process can discard useful information. Consider for example a case where two nodes u and v have assignment vectors $\mathbf{s}_u = [.0, .5, .5, .0]$ and $\mathbf{s}_v = [0, .5, 0, .5]$. Taking the argmax would map both nodes in the 2nd cluster, even if the two assignment vectors are clearly distinguishable. This problem is exacerbated when we do not fix the number of clusters K equal to the true number of classes; in this case, there is no direct correspondence between the clusters and the classes and nothing prevents different classes to be represented by partially overlapping assignment vectors with multiple non-zero entries.

Therefore, to measure the agreement between \mathbf{S} and \mathbf{y} we first consider the cosine similarity between the cluster assignments and the one-hot representation of the node labels:

$$\text{COS} = \frac{\sum_{i,j} [\mathbf{S}\mathbf{S}^\top \odot \mathbf{Y}\mathbf{Y}^\top]_{i,j}}{\sqrt{\sum_{i,j} [\mathbf{S}\mathbf{S}^\top]_{i,j} + \sum_{i,j} [\mathbf{Y}\mathbf{Y}^\top]_{i,j}}} \quad (12)$$

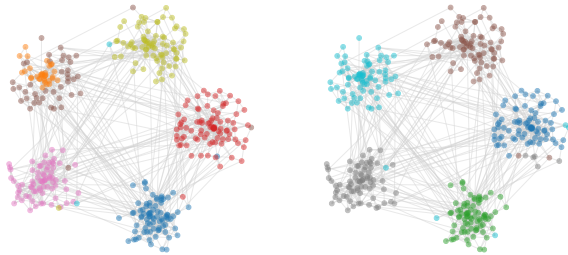
where $\mathbf{Y} = \text{one-hot}(\mathbf{y})$. As a second measure, we consider the accuracy (ACC) obtained by training a simple logistic regression classifier to predict \mathbf{y} from \mathbf{S} .

We compare the performance of BN-Pool with the assignments obtained by four other soft-clustering pooling methods, DiffPool [22], MinCut [23], Deep Modularity Network (DMoN) [36], and Just-balance Graph Neural Network (JBGNN) [37], which are optimized by minimizing their own auxiliary losses. Importantly, we note that the other methods leverage supervised information by setting the number of clusters K equal to the number of node classes, while BN-Pool is completely unsupervised.

As datasets, we consider *Community*, a synthetic dataset generated from a SBM, and four real-world citation networks. Table 1 reports the results and show

Table 1: Mean and standard deviations of ACC and COS for vertex clustering.

Method	Community		Cora		Citeseer		Pubmed		DBLP	
	ACC	COS	ACC	COS	ACC	COS	ACC	COS	ACC	COS
DiffPool	81.9 \pm 1.3	62.9 \pm 0.6	50.4 \pm 1.1	43.3 \pm 0.0	37.9 \pm 1.4	42.4 \pm 0.0	52.4 \pm 0.7	59.8 \pm 0.0	49.5 \pm 4.9	57.4 \pm 0.0
MinCut Pool	97.1 \pm 0.3	94.3 \pm 0.5	57.0 \pm 2.1	40.1 \pm 1.8	54.3 \pm 5.0	36.9 \pm 3.8	61.3 \pm 0.2	46.6 \pm 0.3	69.2 \pm 3.4	52.5 \pm 3.9
DMoN	96.2 \pm 0.9	92.5 \pm 1.6	57.9 \pm 3.8	40.1 \pm 2.3	50.7 \pm 2.4	34.6 \pm 1.6	59.6 \pm 1.4	45.5 \pm 0.7	63.7 \pm 3.2	45.4 \pm 1.3
JBGNN	83.9 \pm 8.7	83.0 \pm 8.9	55.4 \pm 2.4	39.0 \pm 2.8	48.1 \pm 5.0	36.1 \pm 3.3	55.8 \pm 3.8	44.6 \pm 2.0	68.6 \pm 1.8	53.0 \pm 4.4
BN-Pool	98.5 \pm 0.5	83.0 \pm 1.4	66.8 \pm 1.0	47.7 \pm 1.3	47.9 \pm 1.7	37.8 \pm 0.3	81.3 \pm 0.5	62.5 \pm 0.7	75.2 \pm 0.7	58.5 \pm 0.7



(a) BN-Pool

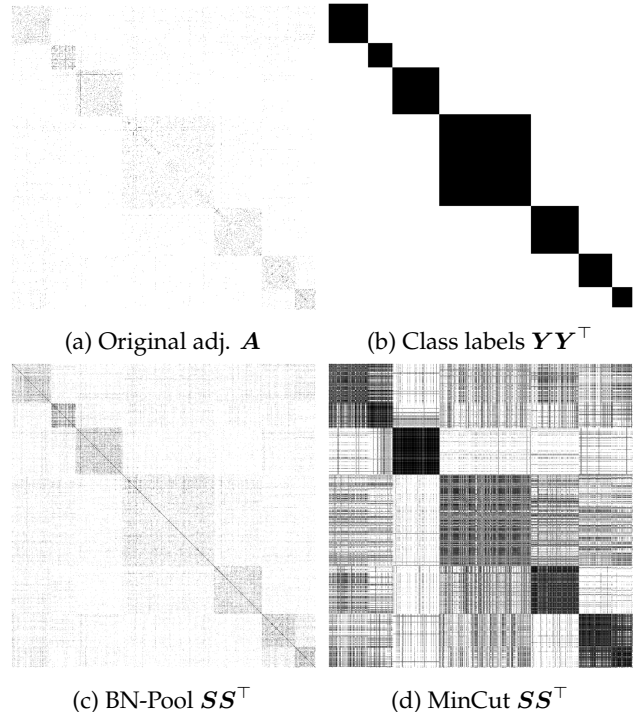
(b) MinCut

Figure 3: Clusters found on a graph with five communities.

that, despite not knowing the real number of classes, BN-Pool achieves good clustering performance.

Figure 3a shows a typical situation where BN-Pool splits a community in two. This happens if there are a few edges within the community and increasing K yield more compact clusters. This cannot occur in other soft-clustering methods such as MinCut. The latter always find the same pre-defined number of clusters ($K = 5$ in this case, see Figure 3b) but create clusters that are more spurious.

Figure 4 shows the original adjacency matrix of the Cora dataset, a visualization of the class labels (YY^T), and the adjacency matrix reconstruction SS^T , where S is the assignment matrix obtained by BN-Pool and MinCut, respectively. While the SS^T produced by BN-Pool follows more closely the actual sparsity pattern of the adjacency matrix, in MinCut SS^T has a block structure. This difference is explained by the different optimisation objectives: while BN-Pool aims to reconstruct the whole adjacency matrix, MinCut recover the communities by cutting the smallest number of edges. In addition, MinCut uses a regularization to encourage clusters to have the same size. This makes it difficult to isolate the smallest clusters (bottom-right and top-left part of the matrix) that, instead, are distinguishable in BN-Pool. Given that in Cora the average edge density between nodes of the same class is only 0.001, a natural way for BN-Pool to lower \mathcal{L}_{rec} is to activate new clusters and generate assignments with multiple non-zero,



(a) Original adj. A

(b) Class labels YY^T

(c) BN-Pool SS^T

(d) MinCut SS^T

Figure 4: Adjacency matrix of Cora, class labels visualization, and adjacency matrix reconstruction by BN-Pool and MinCut.

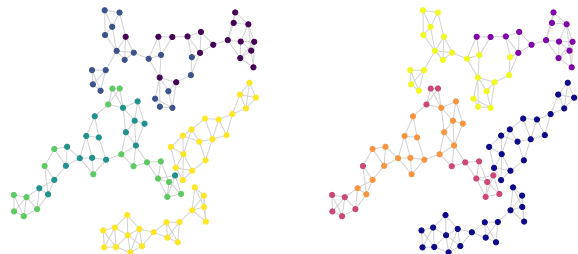
yet low, membership values. See Appendix E for a discussion.

5.2 GRAPH CLASSIFICATION

In graph classification a class label y_i is assigned to the i -th graph $\{A_i, X_i\}$. Differently from the community detection task, here the GNN is optimized by jointly minimizing the cross-entropy loss between true and predicted class labels and the auxiliary loss \mathcal{L}_{aux} . For this task, we compare also against Score-Based and One-Every- K pooling operators, such as Top- k , Edge-Contraction Pooling (ECPool), k Maximal Independent Sets Pooling (k -MIS), and Graclus, which have no auxiliary losses. As datasets, we consider TUData [38] including Colors3 [13], GCB-H [39], and ogbg-molhiv [40]. We report the results in Table 2.

Table 2: Mean and standard deviations of the graph classification accuracy (ROC-AUC for molhiv).

Pooler	GCB-H	Collab	Colors3	IMDB	Mutag.	NCI1	RedditB	DD	MUTAG	Enzymes	Proteins	molhiv
Graclus	75 \pm 3	72 \pm 3	68 \pm 1	77 \pm 6	80 \pm 2	77 \pm 2	90 \pm 3	73 \pm 4	82 \pm 12	33 \pm 7	73 \pm 4	74 \pm 3
ECPool	75 \pm 4	72 \pm 3	69 \pm 2	75 \pm 7	80 \pm 2	77 \pm 3	91 \pm 2	73 \pm 5	84 \pm 12	35 \pm 8	74 \pm 5	74 \pm 1
k -MIS	75 \pm 4	71 \pm 2	84 \pm 1	74 \pm 7	79 \pm 2	75 \pm 3	90 \pm 2	75 \pm 3	83 \pm 10	33 \pm 8	73 \pm 5	74 \pm 2
Top- k	56 \pm 5	72 \pm 2	78 \pm 1	74 \pm 5	75 \pm 3	73 \pm 2	77 \pm 2	72 \pm 5	82 \pm 10	29 \pm 7	74 \pm 5	76 \pm 1
DiffPool	51 \pm 8	70 \pm 2	65 \pm 1	72 \pm 6	78 \pm 2	75 \pm 2	90 \pm 2	75 \pm 4	81 \pm 11	36 \pm 7	75 \pm 3	70 \pm 4
MinCut	75 \pm 5	70 \pm 2	69 \pm 1	73 \pm 6	78 \pm 3	73 \pm 3	87 \pm 2	78 \pm 5	81 \pm 12	34 \pm 9	77 \pm 5	76 \pm 1
DMoN	74 \pm 3	68 \pm 2	69 \pm 2	73 \pm 6	80 \pm 2	73 \pm 3	88 \pm 2	78 \pm 5	82 \pm 11	37 \pm 7	76 \pm 4	77 \pm 1
JBGNN	75 \pm 4	72 \pm 2	68 \pm 2	75 \pm 6	80 \pm 2	78 \pm 3	90 \pm 1	79 \pm 4	87 \pm 14	39 \pm 6	75 \pm 5	73 \pm 2
BN-Pool	74 \pm 3	74 \pm 2	93 \pm 1	75 \pm 8	81 \pm 1	78 \pm 3	90 \pm 2	76 \pm 5	91 \pm 8	52 \pm 8	76 \pm 5	77 \pm 1



(a) Features X

(b) Assignments S

Figure 5: Original node features and assignments S of BN-Pool

In general, BN-Pool performs on par with the best performing pooling operator among those in the Soft-Clustering family. This indicates that BN-Pool can effectively 1) find a meaningful number of clusters, and 2) aggregate nodes without sacrificing useful information. Notable exceptions are the results obtained on the datasets Colors-3 and Enzymes, where BN-Pool outperforms any other pooling method by a significant margin.

Fig. 5 shows the actual node features from a graph from the GCB-H dataset and the node-to-supernodes assignments according to the S found by BN-Pool. Interestingly, there is a very precise matching. GCB-H is a very homophilic dataset, where the nodes can assume only 1 of 5 possible features and nodes with the same features are strongly connected, making it perfectly reasonable to assign node with the same features to the same supernode when learning the pooled graph. Additional examples of this kind of result for the other pooling methods, are reported in [Appendix E](#).

The other Soft-Clustering pooling methods pool each graph in the same predefined number of supernodes K . Instead, BN-Pool does not require to specify K and finds a different K_i for each graph, resulting in a non-trivial distribution pooled graphs' sizes. Fig. 6 shows the distributions of non-empty clusters found

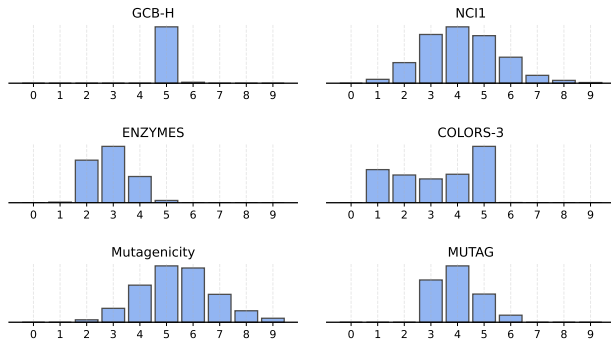


Figure 6: Distribution of non-empty clusters.

by BN-Pool on different datasets, which also allows us to gain further insights about the optimal number of pooled nodes in each dataset.

6 Conclusions

We introduced BN-Pool, a novel graph pooling method that automatically discovers the number of supernodes for each input graph. BN-Pool defines a SBM-like generative process for the input adjacency matrix. By specifying a DP prior over the cluster memberships, our model can handle (theoretically) an infinite number of clusters. Due to the probabilistic nature of BN-Pool, training is performed through the variational inference framework. We employ a GNN to approximate the posterior of the node cluster membership, which allows conditioning the posterior on the node (and potentially edge) features, and on the downstream task at hand. To the extent of our knowledge, this is the first attempt to employ BNP techniques to define a graph pooling method.

Experiments showed that BN-Pool can effectively find a meaningful number of clusters, both to solve unsupervised node clustering and supervised graph classification tasks. Notably, on two graph classification datasets, it outperforms any other pooling method by a significant margin. While we focus on the homophilic

setting, we believe that BN-Pool can be successfully applied also on heterophilic setting by 1) changing the underlying GNN that computes the posterior, and 2) specifying a different prior for the matrix \mathbf{K} which reflect the heterophily in the data.

Acknowledgements

This work was partially supported by the Norwegian Research Council projects 345017: *RELAY: Relational Deep Learning for Energy Analytics*. We gratefully thank Nvidia Corporation for the donation of some of the GPUs used in this project. We also wish to thank Davide Bacciu for the initial discussions and for helping to establish the collaboration that made this work possible.

References

- [1] Amir Hosein Khasahmadi, Kaveh Hassani, Parsa Moradi, Leo Lee, and Quaid Morris. Memory-based graph networks. In *International Conference on Learning Representations*, 2020.
- [2] Hongyang Gao and Shuiwang Ji. Graph u-nets. In *international conference on machine learning*, pages 2083–2092. PMLR, 2019.
- [3] Zheng Ma, Junyu Xuan, Yu Guang Wang, Ming Li, and Pietro Liò. Path integral based convolution and pooling for graph neural networks. *Advances in Neural Information Processing Systems*, 33:16421–16433, 2020.
- [4] Ning Liu, Songlei Jian, Dongsheng Li, Yiming Zhang, Zhiquan Lai, and Hongzuo Xu. Hierarchical adaptive pooling by capturing high-order dependency for graph representation learning. *IEEE Transactions on Knowledge and Data Engineering*, 35(4):3952–3965, 2021.
- [5] Andrea Cini, Danilo Mandic, and Cesare Alippi. Graph-based Time Series Clustering for End-to-End Hierarchical Forecasting. *International Conference on Machine Learning*, 2024.
- [6] Ivan Marisca, Cesare Alippi, and Filippo Maria Bianchi. Graph-based forecasting with missing data through spatiotemporal downsampling. In *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*, pages 34846–34865. PMLR, 2024.
- [7] Jie Zhou, Ganqu Cui, Shengding Hu, Zhengyan Zhang, Cheng Yang, Zhiyuan Liu, Lifeng Wang, Changcheng Li, and Maosong Sun. Graph neural networks: A review of methods and applications. *AI Open*, 1:57–81, 2020. ISSN 2666-6510.
- [8] Peter Orbanz and Yee Whye Teh. Bayesian non-parametric models. *Encyclopedia of machine learning*, 1, 2010.
- [9] Samuel J Gershman and David M Blei. A tutorial on bayesian nonparametric models. *Journal of Mathematical Psychology*, 56(1):1–12, 2012.
- [10] Justin Gilmer, Samuel S Schoenholz, Patrick F Riley, Oriol Vinyals, and George E Dahl. Neural message passing for quantum chemistry. In *International conference on machine learning*, pages 1263–1272. PMLR, 2017.
- [11] Jinheon Baek, Minki Kang, and Sung Ju Hwang. Accurate learning of graph representations with graph multiset pooling. In *Proceedings of the 9th International Conference on Learning Representations*, 2021.
- [12] Daniele Grattarola, Daniele Zambon, Filippo Maria Bianchi, and Cesare Alippi. Understanding pooling in graph neural networks. *IEEE Transactions on Neural Networks and Learning Systems*, 2022.
- [13] Boris Knyazev, Graham W Taylor, and Mohamed Amer. Understanding attention and generalization in graph neural networks. *Advances in neural information processing systems*, 32, 2019.
- [14] Ekagra Ranjan, Soumya Sanyal, and Partha Talukdar. Asap: Adaptive structure aware pooling for learning hierarchical graph representations. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 5470–5477, 2020.
- [15] Junhyun Lee, Inyeop Lee, and Jaewoo Kang. Self-attention graph pooling. In *International conference on machine learning*, pages 3734–3743. PMLR, 2019.
- [16] H. Gao, Y. Liu, and S. Ji. Topology-aware graph pooling networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 43(12):4512–4518, dec 2021. ISSN 1939-3539.
- [17] Yunsheng Pang, Yunxiang Zhao, and Dongsheng Li. Graph pooling via coarsened graph infomax. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 2177–2181, 2021.
- [18] Xing Gao, Wenrui Dai, Chenglin Li, Hongkai Xiong, and Pascal Frossard. ipool—information-based pooling in hierarchical graph neural networks. *IEEE Transactions on Neural Networks and Learning Systems*, 33(9):5032–5044, 2022.
- [19] Inderjit S Dhillon, Yuqiang Guan, and Brian Kulis. Weighted graph cuts without eigenvectors a multi-level approach. *IEEE transactions on pattern analysis and machine intelligence*, 29(11):1944–1957, 2007.

- [20] Filippo Maria Bianchi, Daniele Grattarola, Lorenzo Livi, and Cesare Alippi. Hierarchical representation learning in graph neural networks with node decimation pooling. *IEEE Transactions on Neural Networks and Learning Systems*, 33(5):2195–2207, 2020.
- [21] Davide Bacciu, Alessio Conte, and Francesco Landolfi. Graph pooling with maximum-weight k -independent sets. In *Thirty-Seventh AAAI Conference on Artificial Intelligence*, 2023.
- [22] Zhitao Ying, Jiaxuan You, Christopher Morris, Xiang Ren, Will Hamilton, and Jure Leskovec. Hierarchical graph representation learning with differentiable pooling. *Advances in neural information processing systems*, 31, 2018.
- [23] Filippo Maria Bianchi, Daniele Grattarola, and Cesare Alippi. Spectral clustering with graph neural networks for graph pooling. In *International conference on machine learning*, pages 874–883. PMLR, 2020.
- [24] Hao Yuan and Shuiwang Ji. Structpool: Structured graph pooling via conditional random fields. In *Proceedings of the 8th International Conference on Learning Representations*, 2020.
- [25] David M Blei and Michael I Jordan. Variational methods for the dirichlet process. In *Proceedings of the twenty-first international conference on Machine learning*, page 12, 2004.
- [26] Irina Higgins, Loic Matthey, Arka Pal, Christopher Burgess, Xavier Glorot, Matthew Botvinick, Shakir Mohamed, and Alexander Lerchner. beta-VAE: Learning basic visual concepts with a constrained variational framework. In *International Conference on Learning Representations*, 2017.
- [27] Andrea Asperti and Matteo Trentin. Balancing reconstruction error and kullback-leibler divergence in variational autoencoders. *IEEE Access*, 8:199440–199448, 2020.
- [28] Diederik P. Kingma and Max Welling. Auto-Encoding Variational Bayes. In *2nd International Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14-16, 2014, Conference Track Proceedings*, 2014.
- [29] Mikhail Figurnov, Shakir Mohamed, and Andriy Mnih. Implicit reparameterization gradients. *Advances in neural information processing systems*, 31, 2018.
- [30] Martin Jankowiak and Fritz Obermeyer. Pathwise derivatives beyond the reparameterization trick. In Jennifer Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 2235–2244. PMLR, 10–15 Jul 2018.
- [31] Jia Li, Jianwei Yu, Jiajin Li, Honglei Zhang, Kangfei Zhao, Yu Rong, Hong Cheng, and Junzhou Huang. Dirichlet graph variational autoencoder. *Advances in Neural Information Processing Systems*, 33:5274–5283, 2020.
- [32] Weonyoung Joo, Wonsung Lee, Sungrae Park, and Il-Chul Moon. Dirichlet variational autoencoder. *Pattern Recognition*, 107:107514, 2020.
- [33] Eric Nalisnick and Padhraic Smyth. Stick-breaking variational autoencoders. In *International Conference on Learning Representations*, 2017.
- [34] Nikhil Mehta, Lawrence Carin Duke, and Piyush Rai. Stochastic blockmodels meet graph neural networks. In *International Conference on Machine Learning*, pages 4466–4474. PMLR, 2019.
- [35] Thomas L. Griffiths and Zoubin Ghahramani. The indian buffet process: An introduction and review. *Journal of Machine Learning Research*, 12(32):1185–1224, 2011.
- [36] Anton Tsitsulin, John Palowitch, Bryan Perozzi, and Emmanuel Müller. Graph clustering with graph neural networks. *J. Mach. Learn. Res.*, 24: 127:1–127:21, 2023.
- [37] Filippo Maria Bianchi. Simplifying clustering with graph neural networks. *arXiv preprint arXiv:2207.08779*, 2022.
- [38] Christopher Morris, Nils M. Kriege, Franka Bause, Kristian Kersting, Petra Mutzel, and Marion Neumann. Tudataset: A collection of benchmark datasets for learning with graphs. In *ICML 2020 Workshop on Graph Representation Learning and Beyond (GRL+ 2020)*, 2020.
- [39] Filippo Maria Bianchi, Claudio Gallicchio, and Alessio Micheli. Pyramidal reservoir graph neural network. *Neurocomputing*, 470:389–404, 2022. ISSN 0925-2312.
- [40] Zhenqin Wu, Bharath Ramsundar, Evan N Feinberg, Joseph Gomes, Caleb Geniesse, Aneesh S Pappu, Karl Leswing, and Vijay Pande. Moleculenet: a benchmark for molecular machine learning. *Chemical science*, 9(2):513–530, 2018.
- [41] David Blackwell and James B MacQueen. Ferguson distributions via pólya urn schemes. *The annals of statistics*, 1(2):353–355, 1973.
- [42] Jayaram Sethuraman. A constructive definition of dirichlet priors. *Statistica sinica*, pages 639–650, 1994.

- [43] Jim Pitman. Poisson–dirichlet and gem invariant distributions for split-and-merge transformations of an interval partition. *Combinatorics, Probability and Computing*, 11(5):501–514, 2002.
- [44] Thomas N. Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*. OpenReview.net, 2017.
- [45] Djork-Arné Clevert. Fast and accurate deep network learning by exponential linear units (elus). *arXiv preprint arXiv:1511.07289*, 2015.
- [46] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, 2015.
- [47] Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural networks? In *International Conference on Learning Representations*, 2019.
- [48] Filippo Maria Bianchi and Veronica Lachi. The expressive power of pooling in graph neural networks. In *Advances in Neural Information Processing Systems*, volume 36, pages 71603–71618, 2023.

Appendix

A Dirichlet Process

Given a continuous distribution G_0 , the Dirichlet Process allows to sample a distribution G with the same support as G_0 . Contrarily to G_0 , G is discrete, meaning that the probability of two samples being equal is non-zero, but has a countably infinite number of point masses. Formally, we write

$$G \sim \text{DP}(\alpha_{\text{DP}}, G_0), \quad (13)$$

where α_{DP} is a positive real number representing the concentration parameter, i.e., how much the mass in G is concentrated around a given point. Fig. 7 shows an example of three different draws of G when the base distribution G_0 is a skewed Normal and the value α_{DP} is 1,10,100, and 1000. The base distribution is the expected value of the

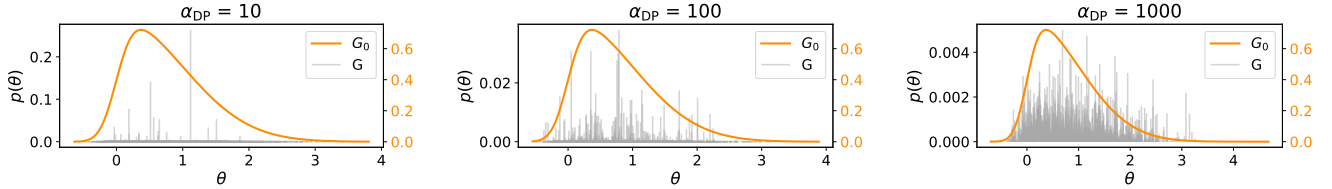


Figure 7: Three single draws from the DP using as G_0 a Normal skewed distribution and three different α_{DP} values. Note that each plot has a different scale on the y-axis.

process, i.e., the Dirichlet process draws distributions around the base distribution G_0 the way a normal distribution draws real numbers around its mean. As we can see from the example, G represents a discrete approximation of G_0

The DP has a clustering property. However, such a property does not emerge from the previous formulation, which also does not tell us how to compute G . In the following, we describe the *Polya urn scheme* [41] and the *stick breaking process* [42]. While the former provides a good intuition of the clustering property of a DP, the latter takes a more constructive perspective that we leverage in this work.

A.1 POLYA URN SCHEME

The Polya urn scheme is an iterative sampling procedure that allows us to sample a sequence of i.i.d. random variables $\theta_1, \theta_2, \dots$ that are distributed according to $G \sim \text{DP}(\alpha_{\text{DP}}, G_0)$. That is, the variables $\theta_1, \theta_2, \dots$ are conditionally independent given G and, hence, exchangeable.

Let us consider the conditional distributions of θ_i given the previous $\theta_1, \dots, \theta_{i-1}$, where G has been integrated out. We can interpret this conditional distribution as a simple urn model containing balls with distinct colors. The balls are drawn equiprobably and when a ball is drawn it is placed back in the urn together with another ball of the same color. In addition, with a probability proportional to α_{DP} , each time we add in the urn a ball with a new color drawn from G_0 . This model exhibits a positive reinforcement effect: the more a color is drawn, the more likely it is to be drawn again.

Let ϕ_1, \dots, ϕ_K be the distinct atoms drawn from G_0 (i.e., the colors) that can be assumed by $\theta_1, \dots, \theta_{i-1}$ (i.e., the balls), and let m_k be the number of times the atom ϕ_k appears in $\{\theta_1, \dots, \theta_{i-1}\}$ for $1 \leq k < i$. Formally, we can express the sampling procedure as:

$$\theta_i \mid \theta_1, \dots, \theta_{i-1} = \begin{cases} \phi_k \text{ with probability } \frac{m_k}{i-1+\alpha_{\text{DP}}} \\ \text{a new draw from } G_0 \text{ with probability } \frac{\alpha_{\text{DP}}}{i-1+\alpha_{\text{DP}}} \end{cases} \quad (14)$$

Equivalently, we can write:

$$\theta_i \mid \theta_1, \dots, \theta_{i-1} \sim \sum_{k=1}^K \frac{m_k}{i-1+\alpha_{\text{DP}}} \delta_{\phi_k} + \frac{1}{i-1+\alpha_{\text{DP}}} G_0, \quad (15)$$

where, δ_{ϕ_k} is a probability measure concentrated at ϕ_k , i.e., δ_{ϕ_k} is a degenerate function assuming value $+\infty$ at ϕ_k and 0 everywhere else.

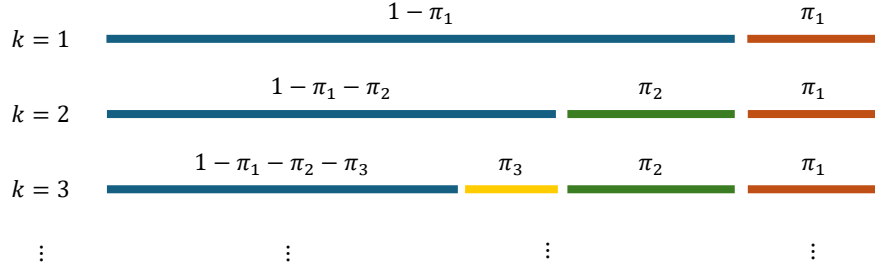


Figure 8: Graphical representation of the stick-breaking process.

Referring to Fig. 7, the values m_k are proportional to the heights of the grey bars. When α_{DP} is small, most of the probability mass is concentrated in a few points. While the Polya urn scheme helps to understand the clustering property of DP, the sampling procedure does not provide an analytic expression of G that we can exploit in our model.

A.2 STICK-BREAKING PROCESS

The idea of the Stick Breaking Process (SBP) is to repeatedly break off a “stick” of initial length 1. Each time we need to break the stick, we choose a value between 0 and 1 that determines the fraction we take from the remainder of the stick. In Figure 8 we show the iterative breaking process, where the values of $\pi_1, \pi_2, \pi_3, \dots$ represent the parts of the stick pieces broken in the first three iterations.

Formally, the stick-breaking construction is based on independent sequences of i.i.d. random variables $(\pi'_k)_{k=1}^\infty$:

$$\pi'_k \mid \alpha_{\text{DP}} \sim \text{Beta}(1, \alpha_{\text{DP}}) \quad \pi_k = \pi'_k \prod_{l=1}^{k-1} (1 - \pi'_l), \quad (16)$$

where the value of π'_k indicates the proportion of the remaining stick that we break at iteration k . To understand the stick analogy, we should first convince ourselves that the quantity $\prod_{l=1}^{k-1} (1 - \pi'_l)$ is equal to the length of the remainder of the stick $1 - \sum_{l=1}^{k-1} \pi_l$ after breaking it $k - 1$ times. Thus, the length of the stick’s piece π_k is obtained by multiplying the stick fraction π'_k by the length of the remaining stick $\prod_{l=1}^{k-1} (1 - \pi'_l)$.

It is important to note that the sequence $\pi = (\pi_k)_{k=1}^\infty$ constructed by equation 16 satisfies $\sum_{k=1}^\infty \pi_k = 1$ with probability one. Thus we may interpret π as a random probability measure on the positive integers. This distribution is often denoted as GEM, which stands for Griffiths, Engen and McCloskey (e.g. see [43]).

Now we have all the ingredients to define a random measure $G \sim \text{DP}(\alpha_{\text{DP}}, H)$:

$$\phi_k \mid G_0 \sim G_0 \quad G = \sum_{k=1}^\infty \pi_k \delta_{\phi_k}, \quad (17)$$

where $(\phi_k)_{k=1}^\infty$ are the atoms drawn from G_0 and δ_{ϕ_k} is a probability measure concentrated at ϕ_k . Sethuraman showed that G as defined in equation 17 is a random probability measure distributed according to $\text{DP}(\alpha_{\text{DP}}, G_0)$. The stick-breaking process is related to the urn scheme since the length of each piece π_k corresponds to the expected probability of drawing a ball of color ϕ_k .

B Implementation details

In this section, we show how we implement the key operations in our model by using as backend the PyTorch library.

B.1 PRIORS AND POSTERIOR DEFINITION

Listing 1 shows how we define the prior and the variational parameters. In particular, the hyperparameters representing the priors are defined as buffers since they are not optimised during the training. Conversely, the

variational parameters are defined as parameters since they are optimised during the training. The variational parameters $\tilde{\alpha}, \tilde{\beta}$ are not defined explicitly since we compute them by applying the linear module W to the node embeddings of size `emb_size` generated by a GNN. The value of `n_clusters` indicates the maximum number of clusters we consider (i.e., the truncation level C of the posterior approximation), and `k_init` is the value used to initialise the variational parameter $\tilde{\mu}$ (i.e., η_K in the main text).

```
import torch.nn.functional as F
import torch as th

# --- Priors (hyperparameters) ---
# Prior for the Stick Breaking Process
register_buffer('alpha_DP', th.ones(n_clusters - 1) * alpha_DP)

# Prior for the cluster-cluster prob. matrix
register_buffer('sigma_K', th.tensor(sigma_K))
register_buffer('mu_K', mu_K * th.eye(n_clusters, n_clusters) -
               mu_K * (1-th.eye(n_clusters, n_clusters)))

# --- Posteriors (parameters) ---
# Transforms node embeddings into posterior distributions for the sticks (alpha_tilde and
# beta_tilde)
W = th.nn.Linear(emb_size, 2*(n_clusters-1), bias=False)

# variational parameters for the connectivity matrix K
mu_tilde = th.nn.Parameter(k_init * th.eye(n_clusters, n_clusters) -
                           k_init * (1-th.eye(n_clusters, n_clusters)))
```

Listing 1: Priors hyperparameters and trainable parameters definition.

B.2 CLUSTER ASSIGNMENTS COMPUTATION

Listing 2 shows the key operations in the forward pass of our model: given the node embeddings produced by a GNN, we compute the cluster assignment matrix S . The forward pass also computes the variational distributions q_π which will be useful later to compute the losses.

```
def compute_pi_given_sticks(stick_fractions):
    """
    Compute the sticks length given the stick fractions
    """
    log_v = th.concat([th.log(stick_fractions), th.zeros(*stick_fractions.shape[:-1], 1)],
                      dim=-1)
    log_one_minus_v = th.concat([th.zeros(*stick_fractions.shape[:-1], 1),
                                 th.log(1 - stick_fractions)], dim=-1)
    pi = th.exp(log_v + th.cumsum(log_one_minus_v, dim=-1))
    return pi # has shape: [T, batch, N, C]

def get_S(node_embs, n_particles, n_clusters):
    """
    Compute soft cluster assignments.
    """
    out = th.clamp(F.softplus(W(node_embs)), min=1e-3, max=1e3)
    alpha_tilde, beta_tilde = th.split(out, n_clusters-1, dim=-1)
    q_pi = th.distributions.Beta(alpha_tilde, beta_tilde)
    stick_fractions = q_pi.z.sample([n_particles])
    S = compute_pi_given_sticks(stick_fractions)
    return S, q_pi
```

Listing 2: Forward computation of the cluster assignments.

At first, on lines 15-16, we obtain the variational parameters $\tilde{\alpha}, \tilde{\beta}$ by applying the linear module W to the node embeddings produced by the GNN. Note that both variational parameters should be greater than 0; thus, we apply the `softplus` activation function. Moreover, to avoid numerical errors, we clamp the values between 10^{-3} and 10^3 .

Once we have the variational parameters, we define the variational distribution by employing the PyTorch class `torch.distributions.Beta`. Then, we sample `n_particles` (i.e., T in the main text) values that will be used to approximate the reconstruction loss by using the `rsample` method. The `r` in the `rsample` name stands for *reparametrization*, that is the trick which allows to separate the distribution parameters from the randomness by allowing to back-propagate the gradient from the samples to the distribution parameters. This technique is also denoted as *pathwise gradient estimator*. As we mentioned in Section 3.1, the reparametrization trick cannot be applied to the Beta distribution explicitly. Therefore, we rely on an approximation of the pathwise derivative [29, 30] which does not require to reparametrise the Beta distribution explicitly. This approximation is already implemented in the PyTorch framework: when we call the `rsample` method, the backend computes (if it is possible) or approximates (as in our case) the pathwise derivative. Thus, the gradient flows from the reconstruction loss to the variational parameters $\tilde{\alpha}, \tilde{\beta}$, and then to the GNN parameters Θ .

The function `compute_pi_given_sticks` computes the stick length π_1, \dots, π_C given the stick fractions π'_1, \dots, π'_C by applying equation 16. The computation is performed in the log-space to avoid numerical errors.

B.3 LOSSES COMPUTATION

Listing 3 shows the computation of the losses $\mathcal{L}_{\text{rec}}, \mathcal{L}_{\pi}, \mathcal{L}_{\mathbf{K}}$.

```
def rec_loss(S, A):

    # Compute the percentage of non-zero links
    # N is the number of nodes
    # E is the number of edges
    balance_weights = (N*N / E) * adj + (N*N / (N*N - E)) * (1 - adj)

    # compute the probability to have an edge for each node pairs, i.e. S K S^T
    p_adj = S @ self.mu_tilde @ S.transpose(-1,-2)

    loss = F.binary_cross_entropy_with_logits(p_adj, A, weight=balance_weights, reduction='none')

    return loss

def pi_prior_loss(self, q_pi):
    alpha_DP = self.get_buffer('alpha_DP')
    p_pi = Beta(th.ones_like(alpha_DP), alpha_DP)
    loss = kl_divergence(q_pi, p_pi).sum(-1)
    return loss

def K_prior_loss(self):
    mu_K, sigma_K = self.get_buffer('mu_K'), self.get_buffer('sigma_K')
    K_prior_loss = (0.5 * (self.mu_tilde - mu_K) ** 2 / sigma_K).sum()
    return K_prior_loss
```

Listing 3: Losses computation.

The function `rec_loss` compute the reconstruction loss \mathcal{L}_{rec} . As shown in equation 11, the value of the loss corresponds to the Binary Cross-Entropy (BCE) loss computed between the adjacency matrix A and the probability to have an edge for each node pairs. Note that we use `BCE_with_logits` rather than applying the sigmoid function to each $\pi_u^T \tilde{\mu} \pi_v$. Since the number of edges is usually much less than the total number of possible edges, we assign different weights to the positive and negative classes to achieve balancing. The weights for the positive class are computed in line 10 and stored in the variable `balance_weights`.

The loss \mathcal{L}_{π} is equal to the KL divergence between the prior $p(\pi'_{ui})$ and the variational posterior $q(\pi'_{ui})$ for each node u and a cluster i . Since all the distributions involved are Beta distributions, the KL divergence has a closed form and it is already implemented in PyTorch. This loss is computed by the function `pi_prior_loss`.

The last loss $\mathcal{L}_{\mathbf{K}}$ is equal to the KL divergence between normal distributions since $q(\mathbf{K}_{ij})$ and $p(\mathbf{K}_{ij})$ are Gaussians for all clusters i and j . Since we do not optimise the variance of the variational distribution, we can ignore all the terms that do not involve the variational parameters $\tilde{\mu}$. Thus, we compute $\mathcal{L}_{\mathbf{K}}$ as the means squared error between $\tilde{\mu}$ and $\mu_{\mathbf{K}}$ scaled by the variance prior $\sigma_{\mathbf{K}}$. This loss is computed by the function `K_prior_loss`.

C Model details

We consider the configurations of the hyperparameters of BN-Pool specified in Table 3. As discussed in Section 3.2, the value of each parameter can be set according to the characteristics of the dataset at hand or by monitoring some performance metrics while training. In our experiments, we select the configuration that yields the lowest value of the reconstruction loss \mathcal{L}_{rec} in the node clustering task and the highest validation accuracy in the graph classification.

Table 3: Values of the hyperparameters of BN-Pool considered.

Hyperparameter	Values
α_{DP}	1.0, 10.0
$\mu_{\mathcal{K}}$	1.0, 10.0, 30.0
$\sigma_{\mathcal{K}}$	0.1, 1.0

We found that setting $\alpha_{\text{DP}} = 10.0$, $\mu_{\mathcal{K}} = 1.0$, and $\sigma_{\mathcal{K}} = 1.0$ yields generally good performance and, thus, it represents our default configuration. Regarding the other hyperparameters, we kept the truncation level $C = 50$, the number of particles $T = 1$, and the initialization of the variational parameter $\eta_{\mathcal{K}} = 1.0$ fixed in all experiments.

C.1 NODE CLUSTERING

The architecture used for clustering is depicted in Fig. 9. As MP layers we used two Graph Convolutional Network (GCN) layers [44] with 32 hidden units and ELU activations [45].

Before training, we apply to the adjacency matrix the same pre-transform used in JBGNN:

$$\mathbf{A} \rightarrow \mathbf{I} - \delta * \mathbf{L}, \quad (18)$$

where \mathbf{L} is the symmetrically normalized graph Laplacian and δ is a constant that we set to 0.85 as in [37].

As training algorithm we used Adam [46] with initial learning rate $1e - 3$. For BN-Pool, we increased η defined in Eq. 10 from 0 to 1 over the first 5,000 epochs according to a cosine scheduler.

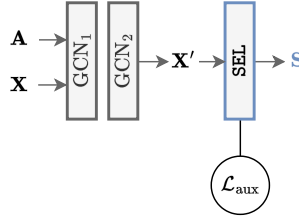


Figure 9: Architecture used for node clustering task.

During training, we monitored the auxiliary losses for early stopping with patience 1,000. When the GNN was configured with BN-Pool, we monitored only \mathcal{L}_{rec} since $\mathcal{L}_{\mathcal{K}}$ and \mathcal{L}_{π} are regularization losses that usually increase and might dominate the total loss.

C.2 GRAPH CLASSIFICATION

The architecture used for graph classification is depicted in Fig. 10. Before and after pooling we use a Graph Isomorphism Network (GIN) [47] layer with 32 hidden units and ELU activations. The readout is an Multilayer Perceptron (MLP) with $[32 \times 32 \times 16 \times N_{\text{class}}]$ units, dropout 0.5, and ELU activation.

Also in this case we apply the pre-transform in Eq. 18. Since some of the datasets contain edge features, we assign to the self-loops that we introduce zero-vectors as surrogate features.

While BN-Pool is able to autonomously discover the number of nodes K_i of each pooled graph, we need to specify the size of the pooled graphs K for the other Soft-Clustering pooling methods and the pooling ratio κ for the Score-Based methods. Therefore, for every dataset, we set $\kappa = 0.5$ and $K = 0.5\bar{N}$, where \bar{N} represents the average nodes in a given dataset.

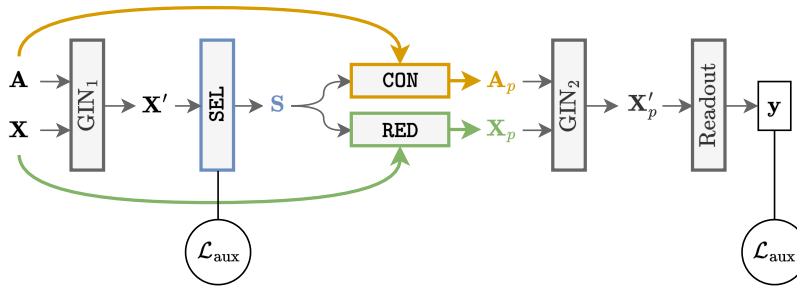


Figure 10: Architecture for graph classification.

As optimizer we used Adam with initial learning rate $5e - 4$. Regarding the callbacks, we monitored the validation accuracy and lowered the learning rate by a factor 0.8 after a plateau of 200 epochs and performed early stopping after 500 epochs. For BN-Pool, we increase η from 0 to 1 over the first 300 epochs using a cosine scheduler.

D Datasets details

The details of the datasets used in the node clustering task are reported in Tab. 4. We reported also the intra-class and inter-class density, which is the average number of edges between nodes that belong to the same or to different classes, respectively. The Community dataset is generated using the PyGSP library⁴. The other datasets are obtained with the PyG loaders⁵.

Table 4: Details of the vertex clustering datasets.

Dataset	#Vertices	#Edges	#Vertex attr.	#Vertex classes	Intra-class density	Inter-class density
Community	400	5,904	2	5	0.1737	0.0025
Cora	2,708	10,556	1,433	7	0.0065	0.0004
Citeseer	3,327	9,104	3,703	6	0.0034	0.0003
Pubmed	19,717	88,648	500	3	0.0005	0.0001
DBLP	17,716	105,734	1,639	4	0.0008	0.0001

Table 5: Details of the graph classification datasets.

Dataset	#Samples	#Classes	Avg. #vertices	Avg. #edges	Vertex attr.	Vertex labels	Edge attr.
GCB-H	1,800	3	148.32	572.32	–	yes	–
Collab	5,000	3	74.49	4,914.43	–	no	–
Colors3	10,500	11	61.31	91.03	4	no	–
IMDB	1,000	2	19.77	96.53	–	no	–
Mutag.	4,337	2	30.32	61.54	–	yes	–
NCI1	4,110	2	29.87	64.60	–	yes	–
RedditB	2000	2	429.63	497.75	–	no	–
D&D	1,178	2	284.32	1,431.32	–	yes	–
MUTAG	188	2	17.93	19.79	–	yes	–
Proteins	1,113	2	39.06	72.82	1	yes	–
Enzymes	600	6	32.63	62.14	18	yes	–
molhiv	41,127	2	25.5	27.5	9	no	3

The details of the datasets used in the graph classification task are reported in Tab. 5. All datasets besides GCB-H and molhiv are downloaded from the TUDataset repository⁶ using the PyG loader. For the GCB-H we used the

⁴<https://pygsp.readthedocs.io/en/stable/>

⁵<https://pytorch-geometric.readthedocs.io/en/2.6.0/modules/datasets.html>

⁶<https://chrsmrrs.github.io/datasets/>

data loader provided in the original repository⁷. Finally molhiv is obtained from the OGB repository⁸ through the loader from the ogb library⁹. For molhiv, we preprocessed the node and edge features using the AtomEncoder and BondEncoder from the ogb library using default embedding dimension size 100.

E Additional results

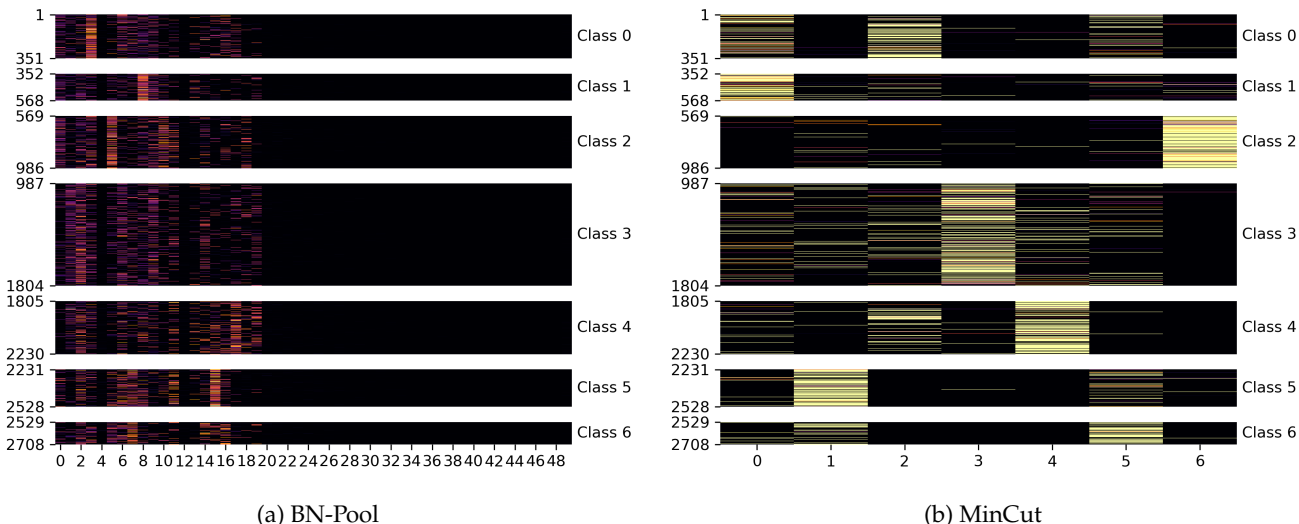


Figure 11: Cluster assignments S found on Cora.

Figure 11a shows the cluster assignments S found by BN-Pool on Cora split according to the node classes. We see that there is not a direct correspondence between the classes and the clusters, since each class is assigned to multiple clusters. This is expected when we do not fix the number of clusters equal to the number of classes, like in the case of BN-Pool that, potentially, can activate an infinite number of clusters. We also notice that the same clusters are active across different classes, albeit with different membership values. Despite such an overlap, there is a clear and consistent pattern in terms of cluster memberships for each class. It is important to notice that the membership values are lower for the nodes of class 3, which is the most populated in the graph. As discussed in Section 5.1, activating many clusters with low membership values is a natural solution found by BN-Pool to reduce \mathcal{L}_{rec} when the intra-class density is very low, like in Cora (0.006).

The clusters found by MinCut on Cora are very different, as shown in Figure 11b. MinCut relies on supervision to set the number of clusters equal to number of class labels. While this allows to achieve a good correspondence between the classes and the clusters, it limits the extent to which MinCut can split a class into multiple clusters, encoding nodes of the same class differently. This implies that if there is a significant variability within each class, MinCut might only assign some of its nodes in the right cluster.

Figure 12 extends Figure 5 from the main body by showing how other pooling methods group the nodes on a sample graph from the GCB-H dataset. BN-Pool creates clusters that match the node features well (Fig. 12b). By contrast, MinCut, which is also a Soft-Clustering method, places nodes with different features in the same clusters (Fig. 12c). In particular, MinCut finds 4 clusters even though there are 5 different feature values.

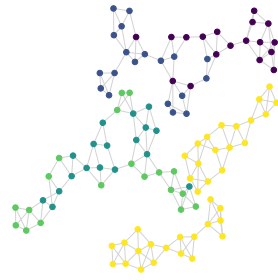
On the other hand, Top- k and k -MIS come from different families (Score-Based and One-Over- K) and pool the graph in a very different way. In particular, Top- k (Fig. 12d) keeps only half of the nodes and drops the others, shown in black. k -MIS does not use the node features, so there is no direct match between the features and the clusters it finds. Figures 12f-i show the different assignment matrices S from these methods, and Figures 12j-m show the topology and node features of the pooled graphs.

BN-Pool uses only 5 clusters that match the 5 feature values. As a result, the pooled graph summarizes effectively the original, with just 5 supernodes, each one tied to a certain feature. On the other hand, MinCut produces a denser

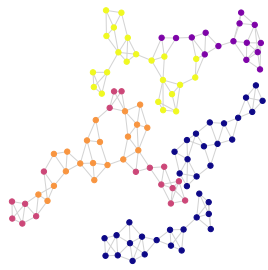
⁷https://github.com/FilippoMB/Benchmark_dataset_for_graph_classification

⁸<https://ogb.stanford.edu/docs/graphprop/>

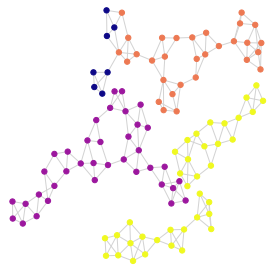
⁹<https://github.com/snape-stanford/ogb>



(a) Original



(b) BN-Pool assignments



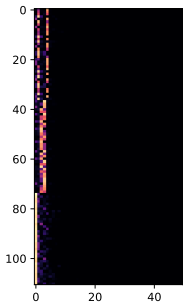
(c) MinCut assignments



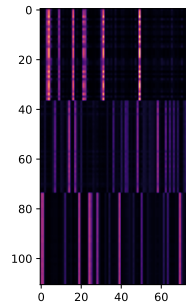
(d) Top- k scores



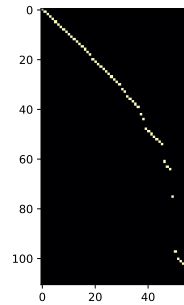
(e) k -MIS assignments



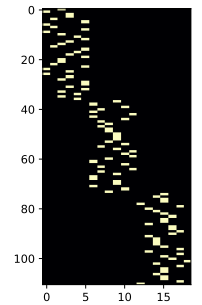
(f) BN-Pool \mathcal{S}



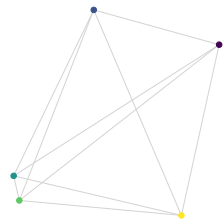
(g) MinCut \mathcal{S}



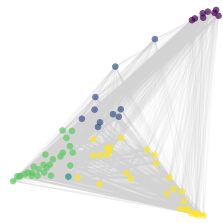
(h) Top- k \mathcal{S}



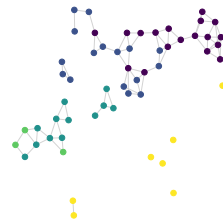
(i) k -MIS \mathcal{S}



(j) BN-Pool $\{A', X'\}$



(k) MinCut $\{A', X'\}$



(l) Top- k $\{A', X'\}$



(m) k -MIS $\{A', X'\}$

Figure 12: Example from GCB-H.

assignment matrix S , where each node belongs to multiple supernodes, and several supernodes have the same role. This overlap is also visible in the pooled graph, which has many supernodes with similar features. Unlike BN-Pool, this pooled graph is less compact, is very dense, and, thus, more costly to process.

Looking at Top- k , we see that its pooled graph is simply a subset of the original, which means some parts of the graph are left out. This is known to be a potential issue in Score-Based methods as it affects their expressivity [48]. Finally, k -MIS yields a pooled graph that, like BN-Pool, is both small and very sparse. It represents all parts of the graph, but it does not match its supernodes to the node features, since it does not consider them.