

Maximal Independent Sets for Pooling in Graph Neural Networks

Stevan Stanovic¹[0000–0001–9656–2080], Benoit Gaüzère²[0000–0001–9980–2641], and
Luc Brun¹[0000–0002–1658–0527]

¹ Normandie Univ, ENSICAEN, CNRS, UNICAEN, GREYC UMR 6072, 14000 Caen, France
{stevan.stanovic, luc.brun}@ensicaen.fr

² INSA Rouen Normandie, Univ Rouen Normandie, Université Le Havre Normandie,
Normandie Univ, LITIS UR 4108, F-76000 Rouen, France
benoit.gauzere@insa-rouen.fr

Abstract. Convolutional Neural Networks (CNNs) have enabled major advances in image classification through convolution and pooling. In particular, image pooling transforms a connected discrete lattice into a reduced lattice with the same connectivity and allows reduction functions to consider all pixels in an image. However, there is no pooling that satisfies these properties for graphs. In fact, traditional graph pooling methods suffer from at least one of the following drawbacks: Graph disconnection or overconnection, low decimation ratio, and deletion of large parts of graphs. In this paper, we present three pooling methods based on the notion of maximal independent sets that avoid these pitfalls. Our experimental results confirm the relevance of maximal independent set constraints for graph pooling.

Keywords: Graph Neural Networks · Graph Pooling · Graph Classification · Maximal Independent Set · Edge Selection

1 Introduction

Convolutional Neural Networks (CNNs) achieved major advances in computer vision by learning abstract representations of images through convolution and pooling. A convolution is a linear filter applied to each pixel of an image which combines its value with the one of its surrounding. The resulting value is usually transformed via a non linear function. The pooling step reduces the size of an image by grouping a connected set of pixels, usually a small squared window, in a single pixel whose value is computed from the ones of window's pixel. Graph Neural Networks (GNNs) take their inspiration from CNNs and aim at transferring advances performed on images to graphs. However, most of CNNs use images with a fixed structure (shape). While using GNN both the structure of a graph and its content varies from one graph to another. Convolution and pooling operations must thus be adapted for graphs.

A GNN may be defined as a sequence of simple graphs $(\mathcal{G}^{(0)}, \dots, \mathcal{G}^{(m)})$ where each $\mathcal{G}^{(l)} = (\mathcal{V}^{(l)}, \mathcal{E}^{(l)})$ is produced by layer l from $\mathcal{G}^{(l-1)}$. Sets $\mathcal{V}^{(l)}$ and $\mathcal{E}^{(l)}$ denote respectively the set of vertices and the set of edges of the graph. Given $n_l = |\mathcal{V}^{(l)}|$, the graph $\mathcal{G}^{(l)}$ may be alternatively defined as $\mathcal{G}^{(l)} = (\mathbf{A}^{(l)}, \mathbf{X}^{(l)})$ where $\mathbf{A}^{(l)} \in \mathbb{R}^{n_l \times n_l}$

is the weighted adjacency matrix of $\mathcal{G}^{(l)}$ while $\mathbf{X}^{(l)} \in \mathbb{R}^{n_l \times f_l}$ encodes the nodes' attributes of $\mathcal{G}^{(l)}$ whose dimension is denoted by f_l . Each line u of $\mathbf{X}^{(l)}$ encodes the feature of the vertex u and is denoted by $x_u^{(l)}$.

The final graph $G^{(n)}$ of a GNN is usually followed by a Multi-Layer Perceptron (MLP) applied on each vertex for a node prediction task or by a global pooling followed by a MLP for a global graph classification task.

Graph convolution. This operation is mainly realized by a message passing mechanism and allows to learn a new representation for each node by combining the information of the mentioned node and its neighborhood. The neighborhood information is obtained by aggregating all the adjacent nodes information. Therefore, the message passing mechanism can be expressed as follows [8]:

$$\mathbf{x}_u^{(l+1)} = \text{UPDATE}^{(l)}(\mathbf{x}_u^{(l)}, \text{AGGREGATE}^{(l)}(\{\mathbf{x}_v^{(l)}, \forall v \in \mathcal{N}(u)\})) \quad (1)$$

where $\mathcal{N}(u)$ is the neighborhood of u and UPDATE , AGGREGATE correspond to differentiable functions.

Let us note that convolution operations should be permutation equivariant, i.e. for any permutation matrix $P \in \{0, 1\}^{n_l \times n_l}$ defined at level l , if f denotes the convolution defined at this layer we must have: $f(PX^{(l)}) = Pf(X^{(l)})$. Note that this last equation, together with equation 1, hides the matrix $\mathbf{A}^{(l)}$ which nevertheless plays a key role in the definition of the AGGREGATE function by defining the neighborhood of each node.

Global pooling. For graph level tasks, a fixed size vector needs to be sent to the MLP. However, due to the variable sizes of graphs within a dataset, global pooling must aggregate the whole graph information into a fixed size vector. This operation can be performed by basic operators like sum, mean or maximum. Let note us that more complex aggregation strategies [19] also exist. To insure that two isomorphic graphs have the same representation, global pooling must be invariant to permutations, i.e. for any permutation matrix P , defined at layer l we must have $g(PX^{(l)}) = g(X^{(l)})$ where g denotes the global pooling operation.

Hierarchical pooling. Summing up a complex graph into a fixed size vector leads necessarily to an important loss of information. The basic idea to attenuate this loss consists in gradually decreasing the size of the input graph thanks to pooling steps inserted between convolution layers. The resulting smaller final graph induces a reduced loss of information in the final global pooling step. This type of method is called a hierarchical pooling [12,18]. The hierarchical pooling step, as the convolution operation should be permutation equivariant in order to keep information localised on desired nodes. Conversely, global pooling must be permutation invariant since it computes a graph level representation. Let note that, similar to CNNs, the reduced graph leads to a reduction of parameters in the next convolution. However, this reduction is mitigated by the learned part of hierarchical pooling. Moreover, let us consider a line graph with a signal optimally sampled on its vertices. As shown by [2], most of GNN correspond to a low pass filter. Applying a GNN on this line graph, hence decreases the maximal

frequency of our signal on vertices producing an over sampling according to the Nyquist theorem. More details on optimal sampling on graphs may be found in [1,15].

Given a graph $\mathcal{G}^{(l)} = (\mathbf{A}^{(l)}, \mathbf{X}^{(l)})$ defined at layer l and its reduced version $\mathcal{G}^{(l+1)} = (\mathbf{A}^{(l+1)}, \mathbf{X}^{(l+1)})$ defined at level $l+1$, the connection between $\mathcal{G}^{(l)}$ and $\mathcal{G}^{(l+1)}$ is usually insured by the reduction matrix $\mathbf{S}^{(l)} \in \mathbb{R}^{n_l \times n_{l+1}}$ where n_l and n_{l+1} denote respectively the sizes of $\mathcal{G}^{(l)}$ and $\mathcal{G}^{(l+1)}$. If $\mathbf{S}^{(l)}$ is a binary matrix, each column of $\mathbf{S}^{(l)}$ encodes the vertices of $\mathcal{G}^{(l)}$ which are merged into a single vertex at layer $l+1$. If $\mathbf{S}^{(l)}$ is real, each line of $\mathbf{S}^{(l)}$ encodes the distribution of each vertex of $\mathcal{G}^{(l)}$ over the vertices of $\mathcal{G}^{(l+1)}$. In both cases, we require $\mathbf{S}^{(l)}$ to be line-stochastic.

Given $\mathcal{G}^{(l)} = (\mathbf{A}^{(l)}, \mathbf{X}^{(l)})$ and $\mathbf{S}^{(l)}$, the feature matrix $\mathbf{X}^{(l+1)}$ of $\mathcal{G}^{(l+1)}$ is defined as follows:

$$\mathbf{X}^{(l+1)} = \mathbf{S}^{(l)\top} \mathbf{X}^{(l)} \quad (2)$$

This last equation defines the attribute of each surviving vertex v_i as a weighted sum of the attributes of the vertices v_j of $\mathcal{G}^{(l)}$ such that $\mathbf{S}_{ji}^{(l)} \neq 0$.

The adjacency matrix of $\mathcal{G}^{(l+1)}$ is defines by:

$$\mathbf{A}^{(l+1)} = \mathbf{S}^{(l)\top} \mathbf{A}^{(l)} \mathbf{S}^{(l)} \quad (3)$$

Let us suppose that $\mathbf{S}^{(l)}$ is a binary matrix. Each entry (i, j) of $\mathbf{A}^{(l+1)}$ defined by equation 3 is equal to $\sum_{r,s}^{n_l} \mathbf{A}_{r,s}^{(l)} \mathbf{S}_{r,i}^{(l)} \mathbf{S}_{s,j}^{(l)}$. Hence two surviving vertices i and j are adjacent in $\mathcal{G}^{(l+1)}$ if it exists at least two adjacent non surviving vertices r and s such that r is merged onto i ($\mathbf{S}_{r,i}^{(l)} = 1$) and s onto j ($\mathbf{S}_{s,j}^{(l)} = 1$).

Pooling methods There are two main families of pooling methods. The first family, called Top- k methods [7,12], is based on a selection of relevant vertices based on a learned criteria. The second family is based on node's clustering methods as in DiffPool [18].

Top- k methods such as gPool [7] learn a score attached to each vertex by computing the scalar product between the vertex's attributes and one or several learned vectors. Alternatively, a GNN can be used to compute a relevance vector for each vertex as in SagPool [12]. Next, a fixed ratio pooling is used to select the k vertices with a highest score. Unselected vertices are dropped. In this case, two surviving vertices in the reduced graph will be adjacent only if they were adjacent before the reduction. This last point may result in the creation of disconnected reduced graphs. This disconnection may be avoided by increasing the density of the graph, using power 2 or 3 of its adjacency matrix or by using the Kron's reduction [3] instead of equation 3. Nevertheless, let us note that simply discarding all non surviving vertices leads to an important loss of information. We proposed in a previous contribution [14], a top- k pooling method called MIVSPool which avoids such drawbacks by using a maximal independent vertex set and graph contraction operations.

Clustering based methods learn explicitly or implicitly the matrix $\mathbf{S}^{(l)}$ which encodes the reduction of a set of vertices at level l into a single vertex at level $l+1$. Methods (eg. [18]) learning $\mathbf{S}^{(l)}$ explicitly have to use a predetermined number of clusters. This last point forbids the use of graphs of different sizes. Additionally, these methods generally result in dense matrices $\mathbf{S}^{(l)}$ which then induce dense adjacency matrices at

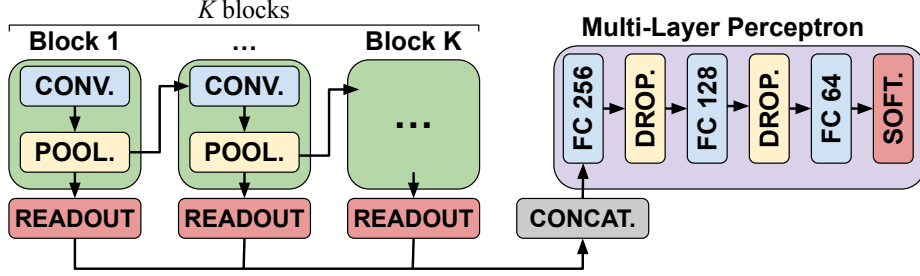


Fig. 1: General architecture of our GNN. Each block is composed of a convolution layer followed by a pooling layer. Features learned after each block are sent to the next block and a Readout layer. The K vectors resulting from each Readout are concatenated to have several levels of description of the graph and, finally, the concatenation is sent to a Multi-Layer Perceptron.

level $l + 1$ (equation 3). As a consequence, graphs produced by these pooling methods have a density close to 1 (i.e. a complete graph or an almost complete graph).

An alternative strategy consists in learning $\mathbf{S}^{(l)}$ only implicitly. Graph pooling such as the maximal matching method used in EdgePool [4] may be associated to this strategy. A maximal matching of a graph $\mathcal{G}^{(l)} = (\mathcal{V}^{(l)}, \mathcal{E}^{(l)})$ is a subset M of $\mathcal{E}^{(l)}$, where no two edges are incident to a same vertex, and every edge in $\mathcal{E}^{(l)} \setminus M$ is incident to one of the two endpoints of an edge in M . EdgePool is based on a maximal weighted matching technique, i.e. a maximal matching of maximal weight. The weight of each edge, called its score, is learned using the attributes of its two end points. The selected edges are then contracted to form a specific cluster. Note that the use of a maximal weighted matching may result in some vertices not incident to any selected edges. These vertices are left unchanged. The sequential algorithm [4] has been parallelized by Landolfi [11]. Unlike EdgePool, Landolfi [11] learns a score attached to each vertex and sort all the vertices of the graph according to their score. The weight of each edge is then defined from a combination of the rank of its incident nodes. The similarity between two adjacent vertices is in this case not taken into account. Moreover, both EdgePool and Landolfi [11] have a decimation ratio lower than 50%, which suggests the need for more pooling steps or a poor abstraction in the final graph of the GNN.

In this paper, we propose an unified family of graph pooling methods which maintains a decimation ratio of approximately 50%, while simultaneously preserving both the structure of the original graph and its attribute information. We achieve this by using a Maximal Independent Set (MIS) [9] to select surviving edges that are evenly distributed throughout the graph, and by assigning non-surviving elements to those that do survive. As a result, we avoid any subsampling or oversampling issues that may arise (see Figure 2). The source code of the paper is available on the CodeGNN ANR Project Git repository: <https://scm.univ-tours.fr/projetspublics/lifat/codegnn>.

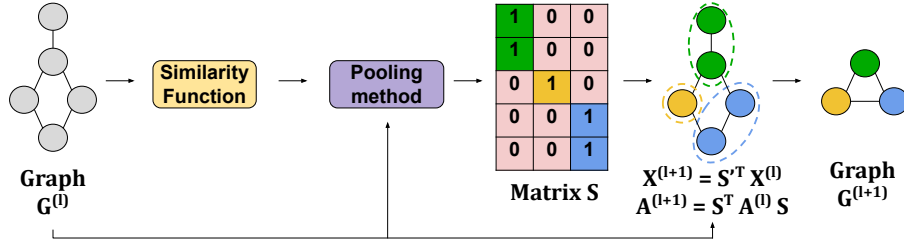


Fig. 2: General proposition of our three graph poolings. Each edge is associated to a similarity score (Section 2). Based on this similarity, a MIS on edge is computed from which a reduction matrix S is derived. Applying S to both feature and structure leads to a reduced graph $G^{(l+1)}$.

2 Maximal Independent Sets and Graph Poolings

2.1 Maximal Independent Set (MIS) and Meer's algorithm

Definition. Let \mathcal{X} be a finite set and \mathcal{N} a neighborhood function defined on \mathcal{X} such that the neighborhood of each element includes the element itself. A subset \mathcal{J} of \mathcal{X} is a Maximal Independent Set (MIS) if the two following equations are fulfilled:

$$\forall (x, y) \in \mathcal{J}^2 : x \notin \mathcal{N}(y) \quad (4)$$

$$\forall x \in \mathcal{X} - \mathcal{J}, \exists y \in \mathcal{J} : x \in \mathcal{N}(y) \quad (5)$$

The elements of \mathcal{J} are called the surviving elements or survivors. Equations (4) and (5) respectively states that two surviving elements can't be neighbors and each non-surviving element has to be in the neighborhood of at least one element of \mathcal{J} . These two equations can be interpreted as a subsampling operation where Equation (4) is a condition preventing the oversampling (two adjacent vertices cannot be selected) while Equation (5) prevents subsampling: Any non-surviving element is at a distance 1 from a surviving one.

A way to compute a MIS is the Meer's algorithm [13] which only involves local computations and is therefore parallelizable. This algorithm attaches a variable to each element. Let us denote by \mathcal{J} the current maximal independent set at an iteration of the algorithm, and let us additionally consider the value v_x attached to an element x . Then x is added to \mathcal{J} at current iteration if v_x is maximal among the values of $\mathcal{N}(x) - \mathcal{N}(\mathcal{J})$, where $\mathcal{N}(\mathcal{J})$ denotes \mathcal{J} and its neighbors. Meer's algorithm provides thus a maximal matching such that each of its element is a local maxima at a given step of the algorithm. We can thus interpret the resulting set as a maximal weight independent set.

Assignment of non-surviving elements. After a MIS, \mathcal{X} is split in two subsets: the surviving elements contained in the set \mathcal{J} and the non-surviving elements contained in $\mathcal{X} - \mathcal{J}$. Simply considering \mathcal{J} as a digest of \mathcal{X} may correspond to an important loss of information which simply discards $\mathcal{X} - \mathcal{J}$. In order to avoid such a loss we allow each non surviving element contained in $\mathcal{X} - \mathcal{J}$ to transfer its information to a survivor.

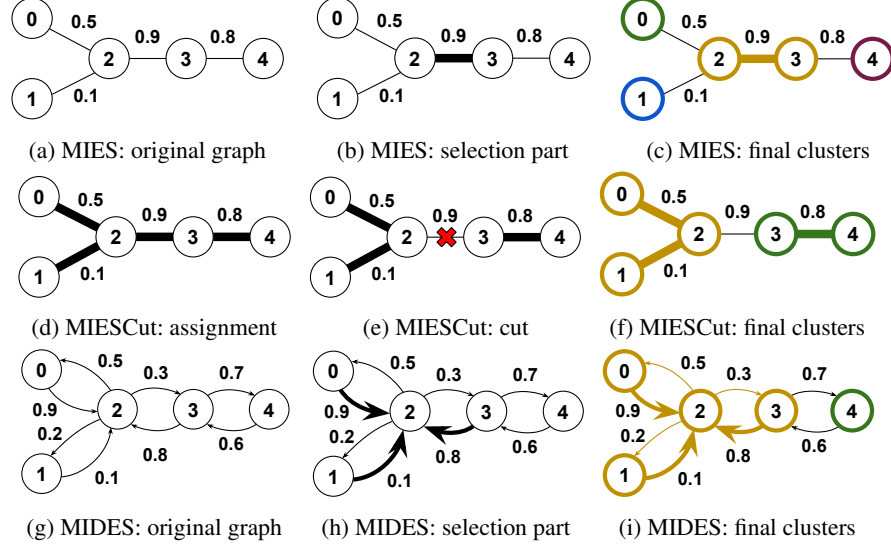


Fig. 3: Schema of our proposed methods on a toy graph. Number on each edge corresponds to its score s and the bold edges indicates the surviving ones. Each group of vertices with the same color represent a cluster. Figures 3a and 3b are common steps for MIES and MIESCut.

The possibility of such a transfer is insured thanks to equation 5 which states that each non surviving element is adjacent to at least one survivor. We can thus associate to any non surviving element x_j a surviving neighbor denoted by $\sigma(x_j)$. At layer l , the global assignment of non-surviving elements onto surviving ones is encoded by the reduction matrix $\mathbf{S}^{(l)} \in \mathbb{R}^{n_l \times n_{l+1}}$ such that :

$$\mathbf{S}_{ii}^{(l)} = 1 \quad \forall x_i \in \mathcal{J} \text{ and } \mathbf{S}_{j\sigma(j)}^{(l)} = 1 \quad \forall x_j \in \mathcal{X} - \mathcal{J} \quad (6)$$

with $\mathbf{S}_{ij}^{(l)} = 0$ otherwise.

2.2 Maximal Independent Sets for Graph Pooling

Based on the work [9] defined within the image partitioning framework we introduce in the following, three adaptations of these methods in order to define learnable pooling steps. In the following sections, the adjacency matrix $\mathbf{A}^{(l+1)}$ is obtained from $\mathbf{A}^{(l)}$ and a binary version of $\mathbf{S}^{(l)}$ using equation 3.

Maximal Independent Edge Set. Most of pooling methods are based on a projection score for each vertex. This strategy is based on the assumption that we can learn features characterizing relevant vertices for a given classification task. However, even if this hypothesis holds, two adjacent vertices may have similar scores and the choice of the survivor is in this case arbitrary. An alternative strategy consists in merging similar

nodes. Given a GNN with hierarchical pooling, the graph sequence corresponds to an increasing abstraction from the initial graphs. Consequently, vertices encoded at each layer of the GNN encode different types of information. Based on this observation, we decided to learn a similarity measure between adjacent vertices at each layer. Inspired by [16], we define the similarity at layer l between two adjacent vertices u and v as $s_{uv}^{(l)} = \exp(-\|\mathbf{W}^{(l)} \cdot (x_u - x_v)\|)$ where x_u and x_v are the features of vertices u and v , $\mathbf{W}^{(l)}$ is a learnable matrix and $\|\cdot\|$ is the L_2 norm.

Given the maximal weighted matching $\mathcal{J}^{(l)}$ defined at level l , each vertex of $\mathcal{G}^{(l)}$ is incident to at most one edge of $\mathcal{J}^{(l)}$. If $u \in \mathcal{V}^{(l)}$ is not incident to $\mathcal{J}^{(l)}$ its features are just duplicated at the next layer. Otherwise, u is incident to one edge $e_{uv} \in \mathcal{J}^{(l)}$ and both u and v are contracted at the next layer. Since u and v are supposed to be similar the attributes of the vertex encoding the contraction of u and v at the next layer must be symmetrical according to u and v . To do so, we first define the attribute of e_{uv} as

$$x_{uv} = \frac{1}{2}(x_u^{(l)} + x_v^{(l)}) \quad (7)$$

where x_u and x_v are the features of vertices u and v . The attribute of the merged vertex is then defined as $s_{uv}x_{uv}$.

An equivalent update of the attributes of the reduced graph may be obtained by computing the matrix $\mathbf{S}^{(l)}$ encoding the transformation from $\mathcal{G}^{(l)}$ to $\mathcal{G}^{(l+1)}$. This matrix can be defined as $\mathbf{S}_{ii}^{(l)} = 1$ if i is not incident to $\mathcal{J}^{(l)}$, and by selecting arbitrary one survivor among $\{u, v\}$ if $e_{uv} \in \mathcal{J}^{(l)}$. If u is selected we set $\mathbf{S}_{uu}^{(l)} = \mathbf{S}_{vu}^{(l)} = \frac{1}{2}s_{uv}$. All remaining entries of $\mathbf{S}^{(l)}$ are set to 0. Matrix $\mathbf{X}^{(l+1)}$ can then be obtained using equation 2. We call this method MIESPool and the main steps are presented in Figures 3a to 3c.

Maximal Independent Edge Set with Cut (MIESCut). Graph reduction through maximal weighted matching has two main drawbacks within the GNN framework. First, a maximal matching may produce many vertices not adjacent to the set of selected edges. Such vertices are just copied to the next level which induce a low decimation ratio (lower than 50%). Given that, the number of layers of a GNN is usually fixed, this last drawback may produce a graph with an insufficient level of abstraction at the final layer of the GNN. Secondly, only the score of the selected edges are used to compute the reduced attributes. This last point reduces the number of scores used for the back-propagation and hence the quality of the learned similarity measures.

As in the previous section, let us denote by $\mathcal{J}^{(l)}$ the maximal weighted matching defined at layer l . By definition of a maximal weighted matching, each vertex not incident to $\mathcal{J}^{(l)}$ is adjacent to at least one vertex which is incident to $\mathcal{J}^{(l)}$. Following [9], we increase the decimation ratio, by attaching isolated vertices to contracted ones. This operation is performed by selecting for each isolated vertex u the edge e_{uv} such that s_{uv} is maximal and v is incident to $\mathcal{J}^{(l)}$.

This operation provides a spanning forest of $\mathcal{G}^{(l)}$ composed of isolated edges, trees of depth one (called stars) with one central vertex and paths of length 3. This last type of tree corresponds to a sequence of 4 vertices with strong similarities between any pair of adjacent vertices along the paths. However, merging all 4 vertices into a single

one, suppose implicitly to apply twice an hypothesis on the transitivity of our similarity measure: more precisely the fact that the two extremities of the paths are similar is not explicitly encoded by our selection of edges. In order to avoid such assumption we remove the central edge of such paths from the selection in order to obtain two isolated edges (see Figures 3d to 3f).

Let us denote by $\mathcal{J}'^{(l)}$ the resulting set of selected edges which forms a spanning forest of $\mathcal{G}^{(l)}$ composed of isolated edges and stars. Concerning the definition of $\mathbf{S}^{(l)}$, we apply the same procedure than in the previous section for isolated edges. For stars, we select the central vertex as the surviving vertex. Let us denote by u such a star's center. We then set $\mathbf{S}_{uu}^{(l)} = \frac{1}{2}$ and $\mathbf{S}_{vu}^{(l)} = \frac{1}{2M} s_{uv}$ for any v such that $e_{uv} \in \mathcal{J}'^{(l)}$ where M is a normalizing factor defined as: $M = \sum_{v|e_{uv} \in \mathcal{J}'^{(l)}} s_{uv}$. The computation of the attributes of the reduced graph using equation 2 and matrix $\mathbf{S}^{(l)}$ is equivalent to compute for each star's center u , the sum, weighted by the score, of the edges' attributes (equation 7) incident to u and belonging to $\mathcal{J}'^{(l)}$:

$$x_u^{(l+1)} = \frac{1}{\sum_{v|e_{uv} \in \mathcal{J}'^{(l)}} s_{uv}} \sum_{v|e_{uv} \in \mathcal{J}'^{(l)}} s_{uv} x_{uv}^{(l)} \quad (8)$$

Maximal Independent Directed Edge Set. The definition of a spanning forest composed of isolated edges and stars is obtained in three steps by MIESCut: The definition of a maximal weight matching, the attachment of isolated vertices and the cut of all paths of length 3. Following [9], we propose to use the Maximal Independent Directed Edge set (MIDES) reduction scheme which obtains the same type of spanning forest in a single step. This reduction scheme is based on a decomposition of the edges e_{uv} of the undirected graphs in two oriented edges $e_{u \rightarrow v}$ and $e_{v \rightarrow u}$. The neighborhood of an oriented edge $\mathcal{N}(e_{u \rightarrow v})$ is defined as the union of the sets of edges leaving u , arriving on u and leaving v . Given $\mathcal{G}^{(l)}$ defined at layer l we formally have:

$$\mathcal{N}^{(l)}(e_{u \rightarrow v}) = \{e_{u \rightarrow v'} \in \mathcal{E}^{(l)}\} \cup \{e_{v' \rightarrow u} \in \mathcal{E}^{(l)}\} \cup \{e_{v \rightarrow v'} \in \mathcal{E}^{(l)}\} \quad (9)$$

The main difference between the neighborhoods defined by equation 9 and the one of MIES is that we do not include in the neighborhood edges arriving on v . This asymmetry allows the creation of stars centered on v . The MIDES algorithm computes a MIS on the edge set using the neighborhood defined by (9) (see Figures 3g to 3i).

At layer l , applying a MIDES on $\mathcal{G}^{(l)}$ requires to define a score function on directed edges. We propose to use $s_{uv} = \exp(-\|W.(x_u - x_v) + b\|)$ where the bias term b allows to introduce an asymmetry so that $s_{uv} \neq s_{vu}$ if $x_u \neq x_v$.

Let us denote by $\mathcal{D}^{(l)}$ the set of directed edges produced by a MIDES on $\mathcal{G}^{(l)}$ using our scoring function. The set $\mathcal{D}^{(l)}$ defines on $\mathcal{G}^{(l)}$ a spanning forest composed of isolated vertices, isolated edges and stars [9].

For an isolated vertex u we duplicate this vertex at the next layer and copy its attributes. We thus set $\mathbf{S}_{uu}^{(l)} = 1$.

For an isolated directed edge $e_{u \rightarrow v} \in \mathcal{D}^{(l)}$ we select v as a surviving vertex and set $\mathbf{S}_{vv}^{(l)} = \frac{s_{uv}}{M}$ and $\mathbf{S}_{uv}^{(l)} = \frac{s_{vu}}{M}$ where $M = s_{uv} + s_{vu}$. This setting corresponds to the following update of the attributes: $x_v^{(l+1)} = (s_{uv}.x_v^{(l)} + s_{vu}.x_u^{(l)})/(s_{uv} + s_{vu})$. Let us

note that since $e_{u \rightarrow v} \in \mathcal{D}^{(l)}$ we have $s_{uv} > s_{vu}$. The previous formula put thus more weight on the surviving vertex v . This update may be considered as a generalization of equation 7 using the asymmetric scores s_{uv} and s_{vu} .

A star within the MIDES framework is defined by a set of edges $e_{u \rightarrow v}$ of $\mathcal{D}^{(l)}$ arriving on the same vertex v . We then set v as survivor and generalize the update of the attributes defined for isolated edges by setting $\mathbf{S}_{vv}^{(l)} = \frac{1}{N} \sum_{u|e_{u \rightarrow v} \in \mathcal{D}^{(l)}} \frac{s_{uv}}{M_u}$ and $\mathbf{S}_{uv}^{(l)} = \frac{1}{N} \frac{s_{vu}}{M_u}$ for all u such that $e_{u \rightarrow v} \in \mathcal{D}^{(l)}$ where $M_u = s_{uv} + s_{vu}$ and N is the number of such u . Such a definition of $\mathbf{S}^{(l)}$ is equivalent to set the updated attribute of v as the mean value of its incident selected edges:

$$x_v^{(l+1)} = \frac{1}{N} \sum_{u|e_{u \rightarrow v} \in \mathcal{D}^{(l)}} \frac{s_{uv}x_v^{(l)} + s_{vu}x_u^{(l)}}{s_{uv} + s_{vu}} \text{ with } N = |\{u \in \mathcal{V}^{(l)} | e_{u \rightarrow v} \in \mathcal{D}^{(l)}\}|.$$

3 Experiments

Datasets. We evaluate our contribution to a bio-informatics and a social dataset called respectively D&D [5] and REDDIT-BINARY [17] whose statistics are reported on Table 1. The aim of D&D is to classify proteins as either enzyme or non-enzyme. Nodes represent the amino acids and two nodes are connected by an edge if they are less than 6 Ångström apart. REDDIT-BINARY is composed of graphs corresponding to online discussions on Reddit. In each graph, nodes represent users, and there is an edge between them if at least one of them respond to the other’s comment. A graph is labeled according to whether it belongs to a question/answer-based community or a discussion-based community.

Model Architecture and Training Procedure. Our model architecture is composed of K blocks where each block consists of a GCN [10] convolution layer followed by a pooling layer. The vector resulting of each pooling operation is then sent to the next block (if it exists) and a Readout layer. A Readout layer concatenates the average and the maximum of vertices’ features matrix $\mathbf{X}^{(l)}$ and these K concatenations are themselves concatenated and sent to a Multi-Layer Perceptron (MLP). The MLP is composed of three fully connected layers and a dropout is applied between each of them. Finally, a Softmax layer is used to determine the binary class of graphs. Note that no batch normalization is applied (Figure 1).

To evaluate our model, we use the training procedure proposed by [6]. This procedure performs an outer 10-fold cross-validation (CV) to split the dataset into ten training and

Dataset	#Graphs	#Classes	Avg $ \mathcal{V} $	Avg $ \mathcal{E} $
D&D [5]	1178	2	284 ± 272	715 ± 694
REDDIT-BINARY [17]	2000	2	430 ± 554	498 ± 623

Table 1: Statistics of datasets

Methods	D&D [5]	REDDIT-BINARY [17]
Baseline	76.29 ± 2.33	87.07 ± 4.72
gPool [7]	75.61 ± 2.74	84.37 ± 7.82
SagPool [12]	76.15 ± 2.88	85.63 ± 6.26
EdgePool [4]	72.59 ± 3.59	87.25 ± 4.78
MIVSPool [14]	76.35 ± 2.09	88.73 ± 4.43
MIESPool	77.17 ± 2.33	88.08 ± 4.55
MIESCutPool	77.74 ± 2.85	86.47 ± 4.57
MIDESPool	76.52 ± 2.21	88.40 ± 4.74

Table 2: Average classification accuracies obtained by different pooling methods. Highest and second highest accuracies are respectively in **bold** and **blue**. \pm indicates the 95% confidence interval of classification accuracy.

test sets. For each outer fold, another 10-fold CV (inner) is applied to the training set to select the best hyperparameter configuration. Concerning hyperparameters, learning rate is set to 10^{-3} , weight decay to 10^{-4} and batch size to 512. Other hyperparameters are tuned using a grid search to find the best configuration. Possible values for the hidden layers sizes are $\{64, 128\}$, dropout ratio is chosen within $\{0.2, 0.5\}$ and the number of blocks K between 1 and 5. We use the Adam optimizer and maximal number of epochs is set to 1000 with an early stopping strategy if the validation loss has not been improved 100 epochs. For EdgePool, due to time constraints, we fixed the hidden layers sizes at 128 and the dropout ratio at 0.5.

We compare, in Table 2, our methods to five state-of-art methods: Baseline (K blocks of GCN [10]), gPool [7], SagPool [12], EdgePool [4] and MIVSPool [14], our previous MIS method. First, we note that the baseline obtains quite good results while not implementing any pooling strategy. It highlights the fact that defining a good pooling operation is not trivial. State-of-the-art methods mostly fail at this task, certainly due to the significant loss of information resulting from the hard selection of surviving vertices using a top- k strategy. This hypothesis is confirmed by the better results obtained by MIVSPool. Let us note also that for D&D, based on T-tests with a significance level of 5%, the average accuracy of EdgePool is statistically lower than the ones of MIS methods. Second, we can observe that the strategies combining edge selection methods and MIS (MIESPool, MIESCutPool, MIDESPool) achieve either the highest or the second highest performances. This empirical results tend to demonstrate that the selection on edges may be most relevant, and that a MIS strategy improves the effectiveness of the pooling over EdgePool. Finally, best results are obtained by different MIS strategies, hence indicating that the right MIS strategy may be dataset dependant. This hypothesis has to be tested using more extensive hyperparameters selection.

4 Conclusion

Graph poolings based on Maximal Independent Sets (MIS) allow, unlike state-of-art methods, to maintain a fixed decimation ratio close to 50%, to preserve vertex information and to avoid subsampling and oversampling. Results obtained by our three methods

based on MIS confirm the interest of this approach but further investigations on other datasets are needed to conclude on the effectiveness of our methods. The design of alternative similarity scores also corresponds to a promising line of research.

Acknowledgements: The work reported in this paper was supported by French ANR grant #ANR-21-CE23-0025 CoDeGNN and was performed using HPC resources from GENCI-IDRIS (Grant 2022-AD011013595) and computing resources of CRIANN (Grant 2022001, Normandy, France).

References

1. Anis, A., Gadde, A., Ortega, A.: Towards a sampling theorem for signals on arbitrary graphs. In: 2014 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP). pp. 3864–3868. IEEE (2014)
2. Balcilar, M., Guillaume, R., Héroux, P., Gaüzère, B., Adam, S., Honeine, P.: Analyzing the expressive power of graph neural networks in a spectral perspective. In: Proceedings of the International Conference on Learning Representations (ICLR) (2021)
3. Bianchi, F.M., Grattarola, D., Livi, L., Alippi, C.: Hierarchical representation learning in graph neural networks with node decimation pooling. *IEEE Trans. Neural Networks Learn. Syst.* **33**(5), 2195–2207 (2022)
4. Diehl, F., Brunner, T., Le, M.T., Knoll, A.: Towards graph pooling by edge contraction. In: ICML 2019 workshop on learning and reasoning with graph-structured data (2019)
5. Dobson, P.D., Doig, A.J.: Distinguishing enzyme structures from non-enzymes without alignments. *Journal of molecular biology* **330**(4), 771–783 (2003)
6. Errica, F., Podda, M., Bacciu, D., Micheli, A.: A fair comparison of graph neural networks for graph classification. *arXiv preprint arXiv:1912.09893* (2019)
7. Gao, H., Ji, S.: Graph u-nets. In: international conference on machine learning. pp. 2083–2092. PMLR (2019)
8. Hamilton, W.L.: Graph representation learning. *Synthesis Lectures on Artificial Intelligence and Machine Learning* **14**(3), 1–159 (2020)
9. Haxhimusa, Y.: The structurally Optimal Dual Graph Pyramid and its application in image partitioning, vol. 308. IOS Press (2007)
10. Kipf, T.N., Welling, M.: Semi-supervised classification with graph convolutional networks. In: International Conference on Learning Representations (ICLR) (2017)
11. Landolfi, F.: Revisiting edge pooling in graph neural networks. In: ESANN (2022)
12. Lee, J., Lee, I., Kang, J.: Self-attention graph pooling. In: International conference on machine learning. pp. 3734–3743. PMLR (2019)
13. Meer, P.: Stochastic image pyramids. *Computer Vision, Graphics, and Image Processing* **45**(3), 269–294 (1989)
14. Stanovic, S., Gaüzère, B., Brun, L.: Maximal independent vertex set applied to graph pooling. In: Structural, Syntactic, and Statistical Pattern Recognition: Joint IAPR International Workshops, S+ SSPR 2022, Montreal, QC, Canada, August 26–27, 2022, Proceedings. pp. 11–21. Springer (2023)
15. Tanaka, Y., Eldar, Y.C., Ortega, A., Cheung, G.: Sampling signals on graphs: From theory to applications. *IEEE Signal Processing Magazine* **37**(6), 14–30 (2020)
16. Verma, N., Boyer, E., Verbeek, J.: Feastnet: Feature-steered graph convolutions for 3d shape analysis. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 2598–2606 (2018)

17. Yanardag, P., Vishwanathan, S.: Deep graph kernels. In: Proceedings of the 21th ACM SIGKDD international conference on knowledge discovery and data mining. pp. 1365–1374 (2015)
18. Ying, Z., You, J., Morris, C., Ren, X., Hamilton, W., Leskovec, J.: Hierarchical graph representation learning with differentiable pooling. *Advances in neural information processing systems* **31**, 4805–4815 (2018)
19. Zhang, M., Cui, Z., Neumann, M., Chen, Y.: An end-to-end deep learning architecture for graph classification. *Proceedings of the AAAI Conference on Artificial Intelligence* **32**(1), 4438–4445 (2018)