

FastGAE: Scalable Graph Autoencoders with Stochastic Subgraph Decoding

Guillaume Salha^{a,b}, Romain Hennequin^a, Jean-Baptiste Remy^b,
Manuel Moussallam^a, Michalis Vazirgiannis^b

^aDeezer Research¹, Paris, France
^bLIX, École Polytechnique, Palaiseau, France

Abstract

Graph autoencoders (AE) and variational autoencoders (VAE) are powerful node embedding methods, but suffer from scalability issues. In this paper, we introduce FastGAE, a general framework to scale graph AE and VAE to large graphs with millions of nodes and edges. Our strategy, based on an effective stochastic subgraph decoding scheme, significantly speeds up the training of graph AE and VAE while preserving or even improving performances. We demonstrate the effectiveness of FastGAE on various real-world graphs, outperforming the few existing approaches to scale graph AE and VAE by a wide margin.

Keywords: Graph Autoencoders, Graph Variational Autoencoders, Scalability, Graph Convolutional Networks, Node Embedding, Link Prediction, Node Clustering

1. Introduction

Graph structures efficiently represent relationships and interactions among entities. Social networks, molecules, citations of scientific publications and web pages constitute some of the most famous examples of data usually represented as graphs, i.e. as *nodes* connected via *edges*. Extracting information from these connections is essential to address numerous graph-based learning problems, ranging from link prediction to influence maximization and node clustering. In this direction, several significant improvements were recently achieved by methods leveraging *node embeddings* [1, 2]. Instead of relying on hand made features, these methods aim at automatically *learning* low-dimensional vector space representations of nodes capturing relevant information from the graph, such as structural proximity, notably by using graph neural networks [3, 4], matrix factorization [5] or random walk processes [6, 7, 8].

In particular, during the last few years, *graph autoencoders* (graph AE) and *graph variational autoencoders* (graph VAE) emerged as two of the most promising and powerful *unsupervised* node embedding methods [9, 10, 11]. Introduced as extensions of standard AE [12, 13] and VAE [14] to graph structures, they involve the combination of two stacked models. First, an *encoder*, typically based on a graph neural network (GNN), maps the nodes into an embedding space; then, a *decoder* tries to reconstruct the original

¹Contact: research@deezer.com

graph structure from the vector representations. Both models are jointly trained to optimize the quality of the reconstruction from the embedding space, in an unsupervised fashion with (for VAE) or without (for AE) a probabilistic approach. Recently, graph AE and VAE have been widely adopted to tackle challenging graph-based problems, such as node clustering [15, 16, 17, 18, 19], graph generation [20, 21, 22, 23, 24] and link prediction [11, 25, 26, 27, 28], often reaching competitive results w.r.t. popular unsupervised baselines such as DeepWalk [6] and node2vec [8].

Nonetheless, graph AE, VAE and their extensions suffer from *scalability* issues. As we explain in Section 2, they mainly result from the costly decoding operations involved in the graph reconstruction. While several recent works provided strategies to scale GNN models² i.e. *encoders* [4, 29, 30, 31], the question of *how to overcome complex decoders* in graph AE and VAE remains open, preventing them from scaling. As a consequence, existing graph AE and VAE have been mainly applied to relatively small graphs, with up to a few thousand nodes. As larger graphs are ubiquitous, we propose to address these scalability concerns in this paper, making the following contributions:

- We introduce FastGAE, a general framework to scale graph AE and VAE models to large graphs with millions of nodes and edges. We leverage graph mining-based sampling schemes and an effective subgraph decoding strategy to significantly lower the computational complexity of graph AE and VAE models, while preserving or even improving their performances.
- We propose an in-depth theoretical and experimental analysis of our method. We demonstrate its empirical effectiveness on seven graphs, with various characteristics and natures, and with up to millions of nodes and edges.
- We publicly release the code of FastGAE³, to ensure reproducibility and future usages.

The remainder of this paper is organized as follows. After reviewing key notions on standard graph AE and VAE and on their complexity in Section 2, we present and analyze FastGAE, our scalable framework, in Section 3. We report our experimental evaluation and a discussion of our results in Section 4, and we conclude in Section 5.

2. Preliminaries

Throughout most of this paper, we consider an undirected graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ with $|\mathcal{V}| = n$ nodes and $|\mathcal{E}| = m$ edges ; extensions of our work to directed graphs are nonetheless discussed in Section 5. We denote by A the binary and symmetric adjacency matrix of \mathcal{G} , and by X an $n \times f$ matrix stacking up f -dimensional node-level features vectors, one for each node of \mathcal{G} . For featureless graphs, we simply assume $X = I_n$, where I_n is the $n \times n$ identity matrix.

²We nonetheless point out that these works were done *out* of the graph AE and VAE unsupervised frameworks that we consider in this paper. They aimed at scaling GNN models that were trained in a supervised or semi-supervised manner.

³Our code is available at: <https://github.com/deezer/fastgae>

2.1. Graph Autoencoders

Graph autoencoders (AE) [9, 10, 11] involve the combination of two stacked models: an *encoder* and a *decoder*.

2.1.1. Encoder

First, the *encoder model* aims at learning an $n \times d$ matrix Z , whose rows z_i are the d -dimensional embedding vectors of each node $i \in \mathcal{V}$, with $d \ll n$. This matrix is usually obtained through a *graph neural network* (GNN) [3, 32, 33] processing A and X . More precisely, most recent variants of graph AE implement multi-layer *graph convolutional networks* (GCN) encoders [3, 11, 34]. In a L -layer GCN (with $L \geq 2$), with input layer $H^{(0)} = X$ and output layer $H^{(L)} = Z$ i.e. the embedding vectors, we have:

$$\begin{aligned} H^{(l)} &= \text{ReLU}(\tilde{A}H^{(l-1)}W^{(l-1)}), \text{ for } l \in \{1, \dots, L-1\} \\ H^{(L)} &= \tilde{A}H^{(L-1)}W^{(L-1)}. \end{aligned} \quad (1)$$

In the above equations, $\tilde{A} = D^{-1/2}(A + I_n)D^{-1/2}$ is the symmetric normalization of A , with D denoting the diagonal degree matrix of $A + I_n$. In a nutshell, at each layer l we compute a representation for each node, by averaging the representations from layer $l-1$ of its direct neighbors (that, from layer 2, have already aggregated information from *their own* neighbors) and of itself (thus the $A + I_n$ in the normalization), together with a ReLU activation: $\text{ReLU}(x) = \max(x, 0)$. $W^{(0)}, \dots, W^{(L-1)}$ are weight matrices, to tune as subsequently detailed in Section 2.1.3.

2.1.2. Decoder

Then, the *decoder model* aims at reconstructing the graph from the embedding. Kipf and Welling [11] and most subsequent works on graph AE models implement a simple *inner-product* decoder. The reconstructed adjacency matrix is then:

$$\hat{A} = \sigma(ZZ^T), \quad (2)$$

with $Z = \text{GCN}(A, X)$, and with $\sigma(\cdot)$ the sigmoid function: $\sigma(x) = 1/(1 + e^{-x})$. In other words, we have $\hat{A}_{ij} = \sigma(z_i^T z_j)$ for all $(i, j) \in \mathcal{V} \times \mathcal{V}$ i.e. nodes with large inner-products in the embedding are more likely to be connected in the graph according to the model. While we will also consider this decoder in our experiments for simplicity and consistency with previous works, we point out the existence of more sophisticated decoders in recent research, such as the asymmetric decoder of Salha et al. [28], the reverse message passing schemes of Grover et al. [27] and the decoder of Shi et al. [19] reconstructing node triads.

2.1.3. Learning

The objective of graph AE is to learn low-dimensional vector representations that ensure a good reconstruction \hat{A} from the decoder. To achieve this, GCN weights are usually tuned by iteratively minimizing, by *gradient descent* [35], a *reconstruction loss* capturing the similarity between A and \hat{A} . In the graph AE framework, this loss is usually formulated as a cross entropy loss [11] i.e.:

$$\mathcal{L} = \frac{-1}{n^2} \sum_{(i,j) \in \mathcal{V}^2} \left[A_{ij} \log(\hat{A}_{ij}) + (1 - A_{ij}) \log(1 - \hat{A}_{ij}) \right]. \quad (3)$$

Also, in practice, the pairs with $A_{ij} = 1$ are often re-weighted in the loss, when dealing with a sparse adjacency matrix A [11, 36].

2.2. Graph Variational Autoencoders

Kipf and Welling [11] also introduced graph extensions of *variational autoencoders* (VAE), a model originally introduced for general data by Kingma and Welling [14].

2.2.1. Encoder

In a graph VAE, we establish a probabilistic model on A , involving a d -dimensional latent variable z_i for each node $i \in \mathcal{V}$, corresponding to its embedding vector. Kipf and Welling [11] propose the following mean-field inference model as *encoder*:

$$q(Z|A, X) = \prod_{i=1}^n q(z_i|A, X), \quad (4)$$

with $q(z_i|A, X)$ corresponding to a $\mathcal{N}(\mu_i, \text{diag}(\sigma_i^2))$ distribution. They leverage two GCNs to learn the d -dimensional Gaussian mean and variance vectors μ_i and σ_i for each node. In a nutshell, $\mu = \text{GCN}_\mu(A, X)$, with μ the matrix of mean vectors μ_i ; also, $\log \sigma = \text{GCN}_\sigma(A, X)$. The embedding vectors z_i are samples subsequently drawn from these distributions.

2.2.2. Decoder

Then, from these embedding vectors, a *generative model* aims at *decoding* A using, as for graph AE, inner-products together with sigmoid activations. We have:

$$\hat{A}_{ij} = p(A_{ij} = 1|z_i, z_j) = \sigma(z_i^T z_j). \quad (5)$$

Then:

$$p(A|Z) = \prod_{i,j=1}^n p(A_{ij}|z_i, z_j). \quad (6)$$

2.2.3. Learning

Kipf and Welling [11] propose to iteratively maximize a lower bound of the model’s likelihood (ELBO) [14] by gradient descent w.r.t. the GCNs’ weights:

$$\mathcal{L}^{\text{ELBO}} = \mathbb{E}_{q(Z|A, X)} \left[\log p(A|Z) \right] - \mathcal{D}_{KL} \left(q(Z|A, X) || p(Z) \right). \quad (7)$$

$\mathcal{D}_{KL}(\cdot||\cdot)$ denotes the Kullback-Leibler divergence [37], and $p(Z)$ corresponds to an initial standard Gaussian prior on the distribution of latent vectors. We refer to [11, 14] for more technical details on computations.

2.3. On Complexity and Scalability

GCN models have become popular *encoders* for graph AE and VAE, thanks to their relative simplicity w.r.t. other GNN architectures [32, 33]. The cost of evaluating each layer of a GCN encoder is *linear* w.r.t. the number of edges m [3], which can also be improved by instead encoding nodes with a FastGCN [30], with a Cluster-GCN [31], or by using simple graph convolutions [38] or other stochastic strategies [4, 29, 39, 40].

However, the inner-product *decoders* of graph AE and VAE both involve the multiplication of the dense matrices Z and Z^T at each training iteration. It suffers from a

quadratic $O(n^2)$ complexity, as the aforementioned alternative decoders that all require inner-products or Euclidean distances computations. Storing $n \times n$ dense matrices \hat{A} can also lead to memory issues for large n . As a consequence, the recent aforementioned efforts to scale GCNs (that were achieved in a supervised setting, and out of the wider graph AE and VAE frameworks studied in this paper, where GCNs are only a building block) are *not* sufficient to scale graph AE and VAE. These models still suffer from a quadratic time complexity due to their costly decoding operations, and therefore from scalability issues.

As a result, graph AE, VAE and their extensions were usually applied to relatively small graphs with up to a few thousand nodes and edges. We acknowledge that Kipf and Welling [11], Grover et al. [27] and Salha et al. [36] all very briefly mentioned *sampling* ideas as possible extensions, but without further investigation. We will later observe that a direct uniform sampling of nodes is often sub-optimal. Recently, Shi et al. [19] incorporated more elaborated mini-batch sampling ideas in their work on graph AE and VAE. More specifically, authors proposed to reconstruct *triads* of nodes in their decoder instead of simple inner-products. The total number of triads in a graph is $\binom{n}{3}$ which is very large, but authors proposed to sample a smaller number of triads to make their model tractable. Nonetheless, as the objective of their work was more to improve the performances of graph AE and VAE than to provide scalable models, they did not report running times nor experiments on large graphs.

In a wider analysis on scalable graph AE and VAE, Salha et al. [18] proposed to speed up computations by training the AE/VAE only on a smaller version of the graph, namely on the k -core of the graph [41], then by propagating representations to other nodes via simpler but faster heuristics. While they did provide experiments on larger graphs, their performances tend to significantly deteriorate for smaller cores i.e. for faster models. To sum up, the question of how to effectively scale graph AE and VAE remains unsatisfactorily addressed.

3. Scaling Graph AE and VAE with FastGAE

In this Section, we introduce our proposed framework to scale graph AE and VAE. We refer to it as *FastGAE*, and as *variational FastGAE* when applied to graph VAE.

3.1. Encoding the Entire Graph...

As explained in Section 2.3, GCN models [3] and their scalable extensions [30, 31, 38, 40] can effectively process large graphs. Therefore, in our FastGAE framework, we rely on these models to *encode all the nodes* into the embedding space. More precisely, in the following experiments, we implement *standard GCN encoders* for the sake of simplicity and for an easier comparison to existing graph AE and VAE architectures. This design choice is made *without loss of generality*, the FastGAE framework being valid for *any* other encoder producing an embedding matrix Z .

3.2. ...But Decoding Stochastic Subgraphs

However, while computing node embedding vectors through a *forward* GCN pass is fast, *tuning the weights* of this encoder in the graph AE and VAE settings requires the reconstruction of the entire matrix \hat{A} at each training iteration which, as detailed in Section 2.3, suffers from a quadratic complexity and is intractable for large graphs.

3.2.1. Subgraph Decoding

To overcome this issue, we propose to *approximate reconstruction losses*, by computing their values only from *wisely selected random subparts* of the original graph. More precisely, at each training iteration, we aim at decoding a different sampled subgraph of \mathcal{G} with $n_{(S)}$ nodes, with $n_{(S)} < n$ being a fixed parameter. Let $\mathcal{G}_{(S)} = (\mathcal{V}_{(S)}, \mathcal{E}_{(S)})$ be such sampled subgraph, with $\mathcal{V}_{(S)} \subset \mathcal{V}$, $|\mathcal{V}_{(S)}| = n_{(S)}$, and with $\mathcal{E}_{(S)}$ denoting the subset of edges connecting the nodes in $\mathcal{V}_{(S)}$. Instead of reconstructing the $n \times n$ matrix \hat{A} , we propose to reconstruct the smaller $n_{(S)} \times n_{(S)}$ matrix $\hat{A}_{(S)}$ with:

$$\hat{A}_{(S)}_{ij} = \sigma(z_i^T z_j), \quad \forall (i, j) \in \mathcal{V}_{(S)}^2, \quad (8)$$

and to only learn from the quality of $\hat{A}_{(S)}$ w.r.t. its ground truth counterpart $A_{(S)}$, as measured by a cross entropy loss for AE, or an ELBO loss for VAE. We propose to use the resulting approximate loss for gradients computations and GCN weights updates by gradient descent. We draw *a different subgraph $\mathcal{G}_{(S)}$ at each training iteration*, using the sampling methods detailed next.

3.2.2. (Naive) Uniform Node Sampling

A very simple way to obtain such subgraphs would consist in *uniformly sampling $n_{(S)}$ nodes* from the set \mathcal{V} at each training iteration. However, with such strategy, there is no guarantee that the most important links (or absence of links) from the original graph structure will be preserved in the drawn subgraphs to reconstruct during the training phase. As we will experimentally show in Section 4, this usually significantly impacts the quality of the final node embedding, leading to underperforming performances on downstream evaluation tasks. As a consequence, in the following sections, we propose and study more refined strategies, aiming at *leveraging the graph structure* to obtain a more effective sampling.

3.2.3. Node Sampling with Graph Mining

We propose to consider alternative sampling methods, that increase the probability of including some particular nodes in the drawn subgraph w.r.t. some others. Let $f : \mathcal{V} \rightarrow \mathbb{R}^+$ denote some measure of the relative *importance* of nodes in the graph, obtained through graph mining methods. Assuming such function is available, we draw inspiration from word sampling in natural language processing [42, 43] and propose to set the probability to pick each node $i \in \mathcal{V}$ as the first element of $\mathcal{V}_{(S)}$ as:

$$p_i = \frac{f(i)^\alpha}{\sum_{j \in \mathcal{V}} (f(j)^\alpha)}, \quad (9)$$

with $\alpha \in \mathbb{R}^+$. Then, assuming we sample $n_{(S)}$ *distinct nodes without replacement*, each remaining node $i \in \mathcal{V} \setminus \mathcal{V}_{(S)}$ has a probability $p_i / \sum_{j \notin \mathcal{V}_{(S)}} p_j$ to be picked as the second element of $\mathcal{V}_{(S)}$, and so on until $|\mathcal{V}_{(S)}| = n_{(S)}$. The previous division is a simple normalization to ensure $\sum_{j \notin \mathcal{V}_{(S)}} p_j = 1$ at each sampling step. Alternatively, one could also sample $n_{(S)}$ nodes *with replacement*: it simplifies computations, as sampling probabilities are then independent of previous draws and remain fixed to p_i , but a node

could then be drawn several times. We stress out that, in our implementation, both variants return very similar results. We later adopt the former.

In a nutshell, *important nodes according to f are more likely to be selected for decoding*, and the hyperparameter α helps sharpening (for $\alpha > 1$) or smoothing (for $\alpha < 1$) the distribution. Setting $\alpha = 0$ leads to the aforementioned uniform node sampling. In our experiments, we will consider and evaluate two importance measures f from graph mining:

- the *degree* of each node, which is simply the number of direct connections of each node: $f(i) = \sum_{j \in \mathcal{V}} A_{ij}$.
- the *core number* of each node: $f(i) = C(i)$. The k -core version of a graph is its largest subgraph for which every node has a degree higher or equal to k *within this subgraph*. The core number $C(i)$ of a node i corresponds to the largest value of k for which i is in the k -core. Core decomposition has been widely used over the past years to quantify the significance of nodes and extract representative subgraphs (see e.g. Malliaros et al. [41] for a review). They constitute a more *global* importance measure than the *local* node degree.

Besides their popularity and their complementarity, we also chose to focus on these two metrics for *computational efficiency*. Indeed, contrary to other potential importance metrics based on influence maximization [44], random walks [45] or centrality measures [46], both can be evaluated in a linear $O(m)$ running time [47]. As we will empirically check in Section 4, this leads to *fast and scalable computations of probability distributions*, which is crucial for our FastGAE framework whose primary objective is scalability. We refer the interested reader to the work of Leskovec and Faloutsos [45], Hu and Lau [48] and Chiericetti et al. [49] for a broader overview of other existing graph sampling methods.

3.3. Theoretical Considerations

In the Section 3.3.1, we provide a global overview of some theoretical analyses that we subsequently further develop in Sections 3.3.2 and 3.3.3. For readability reasons, we will report the proofs of all propositions in the appendices.

3.3.1. Overview of Theoretical Considerations

On Approximate Losses. In the case of degree and core-based sampling strategies, some node pairs from the graph are more likely to appear in subgraphs than others. The probability to draw a node i , or an edge incident to i , increases with p_i and with $f(i)$ for $\alpha > 0$. As a consequence, at each gradient descent iteration, the approximate loss (say $\mathcal{L}^{\text{FastGAE}}$) is *biased* w.r.t. the standard graph AE or VAE loss that would have been computed on \mathcal{G} (say \mathcal{L}), i.e. $\mathbb{E}(\mathcal{L}^{\text{FastGAE}}) \neq \mathcal{L}$ in general. For completeness, in Propositions 1, 2 and 3 of the upcoming Section 3.3.2, *we provide a theoretical analysis, in which we fully explicit the expected loss $\mathbb{E}(\mathcal{L}^{\text{FastGAE}})$ that we actually stochastically optimize in FastGAE, as well as the formal probabilities to sample a given node or node pair at each training iteration.* Moreover, we will show in Section 4 that, despite such bias, optimizing this alternative loss does not deteriorate the quality of node embeddings. On the contrary, we will provide insights exhibiting the fact that re-weighting node pairs from high degree/core nodes can actually be beneficial.

On the Selection of $n_{(S)}$. When selecting $n_{(S)}$, one faces a performance/speed trade-off: reconstructing very small subgraphs will speed up training but, as we later verify, this might also deteriorate performances. While we claimed in the previous paragraph that stochastically minimizing $\mathbb{E}(\mathcal{L}^{\text{FastGAE}})$ instead of \mathcal{L} might be beneficial, we also acknowledge that, for small values of $n_{(S)}$, the actual loss $\mathcal{L}^{\text{FastGAE}}$ computed at a given training iteration can significantly deviate from its expectation. In this paper, we propose to use these deviations as a criterion to select a relevant subgraph size. In Propositions 4 and 5 of the upcoming Section 3.3.3, we leverage concentration inequalities to derive a *theoretically-grounded threshold size*, denoted $n_{(S)}^*$ in the following, for which, at each training iteration, the deviation between the evaluation of $\mathcal{L}^{\text{FastGAE}}$ for each node and its expectation is proven to be bounded with a high probability. This proposed subgraph size is of the form:

$$n_{(S)}^* = C\sqrt{n} \quad (10)$$

where constant $C > 0$ depends on the deviation magnitude and probability, and is explicitly presented thereafter. Our experiments will confirm the relevance of this choice.

3.3.2. On Approximate Losses

Let us recall that, in our FastGAE framework, at each training iteration we run a full GCN forward pass and sample a subgraph $\mathcal{G}_{(S)} = (\mathcal{V}_{(S)}, \mathcal{E}_{(S)})$. Then, we evaluate reconstruction losses only on this subgraph, which involves fewer operations w.r.t. standard decoders, and we use the resulting approximate loss for GCN weights updates via gradient descent. More precisely, in standard implementations of graph AE/VAE, the cross entropy loss (from Section 2.1 on AE) and the negative of the ELBO’s expectation part (from Section 2.2 on VAE) are empirically derived by computing the following node pairs average at each training iteration:

$$\mathcal{L} = \frac{1}{n^2} \sum_{(i,j) \in \mathcal{V}^2} \mathcal{L}_{ij}(A_{ij}, \hat{A}_{ij}), \quad (11)$$

with⁴:

$$\mathcal{L}_{ij}(A_{ij}, \hat{A}_{ij}) = -[A_{ij} \log(\hat{A}_{ij}) + (1 - A_{ij}) \log(1 - \hat{A}_{ij})].$$

In the FastGAE framework, we instead compute:

$$\mathcal{L}^{\text{FastGAE}} = \frac{1}{n_{(S)}^2} \sum_{(i,j) \in \mathcal{V}_{(S)}^2} \mathbb{1}_{((i,j) \in \mathcal{V}_{(S)}^2)} \mathcal{L}_{ij}(A_{ij}, \hat{A}_{ij}), \quad (12)$$

with $\mathbb{1}_{((i,j) \in \mathcal{V}_{(S)}^2)} = 1$ if $(i, j) \in \mathcal{V}_{(S)}^2$ and 0 otherwise.

We recall that, for variational FastGAE, we need to subtract the Kullback-Leibler (KL) divergence, as in the ELBO of standard graph VAE, to obtain our actual loss evaluation. At this stage, two options are possible:

⁴In most graph AE and VAE implementations (see e.g. [11]), the terms with $A_{ij} = 1$ are often re-weighted in the loss, in case of sparse A . They are multiplied by $w \geq 1$, a positive links re-weighting scalar parameter which is usually inversely proportional to the graph sparsity. In our analyses, to clarify notations, we omit this scalar multiplication, which is equivalent to implicitly assuming that $w = 1$. This simplification is made without loss of generality and all results remain valid for any $w > 1$.

- Computing the KL term only on the nodes in the subgraph.
- Or, computing the KL term on all nodes.

We consider that the two options are valid. The first one ensures that the resulting loss is a proper lower bound of the likelihood computed on this subgraph. The second one, despite violating this property, can nonetheless be empirically convenient and interpreted as the addition of a *regularization term* on all node embedding vectors (penalizing large deviations w.r.t. a $\mathcal{N}(0, I_d)$ prior distribution on these vectors) to the *performance term* $\mathcal{L}^{\text{FastGAE}}$. In our implementations, both options returned similar results. In the following propositions, we assume that the KL term is computed on *all* nodes for simplicity, and we therefore only approximate the performance term \mathcal{L} , both in the AE and in the VAE settings.

Propositions 1 and 2 detail the formal probabilities to sample a given node or a given node pair at each training iteration. We consider both sampling variants *with* and *without* replacement (see Section 3.2) for this analysis, as the former significantly simplifies results w.r.t. the later.

Proposition 1. *Let $\mathcal{G}_{(S)} = (\mathcal{V}_{(S)}, \mathcal{E}_{(S)})$ be a subgraph of \mathcal{G} obtained from sampling $n_{(S)}$ nodes **with** replacement using the node sampling strategy of FastGAE. Let i and j denote two distinct nodes from the original graph \mathcal{G} : $(i, j) \in \mathcal{V}^2$. Then:*

$$\mathbb{P}(i \in \mathcal{V}_{(S)}) = 1 - (1 - p_i)^{n_{(S)}}. \quad (13)$$

Also:

$$\mathbb{P}\left((i, j) \in \mathcal{V}_{(S)}^2\right) = 1 - \left[(1 - p_i)^{n_{(S)}} + (1 - p_j)^{n_{(S)}} - (1 - p_i - p_j)^{n_{(S)}} \right]. \quad (14)$$

Proposition 2. *Let $\mathcal{G}_{(S)} = (\mathcal{V}_{(S)}, \mathcal{E}_{(S)})$ be a subgraph of \mathcal{G} obtained from sampling $n_{(S)}$ nodes **without** replacement using the node sampling strategy of FastGAE. Let i and j denote two distinct nodes from \mathcal{G} : $(i, j) \in \mathcal{V}^2$. Then:*

$$\mathbb{P}(i \in \mathcal{V}_{(S)}) = \sum_{\mathcal{U} \in \mathcal{U}(i)} p_{u_1} \prod_{k=2}^{n_{(S)}} \frac{p_{u_k}}{1 - \sum_{k'=1}^{k-1} p_{u_{k'}}}, \quad (15)$$

where $\mathcal{U}(i) = \{\mathcal{U} \subset \mathcal{V}, |\mathcal{U}| = n_{(S)} \text{ and } i \in \mathcal{U}\}$ is the set of all ordered subsets of $n_{(S)}$ distinct nodes including node i . For a given set $\mathcal{U} \in \mathcal{U}(i)$, we denote by $(u_1, u_2, \dots, u_{n_{(S)}})$ its ordered elements. Also,

$$\mathbb{P}\left((i, j) \in \mathcal{V}_{(S)}^2\right) = \sum_{\mathcal{U} \in \mathcal{U}(i) \cap \mathcal{U}(j)} p_{u_1} \prod_{k=2}^{n_{(S)}} \frac{p_{u_k}}{1 - \sum_{k'=1}^{k-1} p_{u_{k'}}}. \quad (16)$$

Despite different formulations, both variants share a similar behaviour in practice on most real-world graphs. In this paper, as explained in Section 3.2, we sample nodes *without replacement*. One can derive from the above expressions that the probability to draw a node i , or an edge incident to i , increases with $n_{(S)}$, with p_i and with $f(i)$ for $\alpha > 0$. This also leads to the following formulation of the expected (re-weighted) loss that FastGAE stochastically optimizes.

Proposition 3. *Using the expressions of Proposition 1 (with replacement) or Proposition 2 (without replacement):*

$$\mathbb{E}[\mathcal{L}^{\text{FastGAE}}] = \frac{1}{n_{(S)}^2} \sum_{(i,j) \in \mathcal{V}^2} \mathbb{P}\left((i,j) \in \mathcal{V}_{(S)}^2\right) \mathcal{L}_{ij}(A_{ij}, \hat{A}_{ij}). \quad (17)$$

3.3.3. On the Selection of $n_{(S)}$

While our experiments will tend to show that stochastically minimizing $\mathbb{E}(\mathcal{L}^{\text{FastGAE}})$ (equation 17) instead of \mathcal{L} (equation 11) might be beneficial, we also acknowledge that, for small values of $n_{(S)}$, the actual loss $\mathcal{L}^{\text{FastGAE}}$ computed at a given training iteration (equation 12) might significantly deviate from its expectation.

We propose to use these deviations as a criterion to automatically select a relevant subgraph size. More precisely, let us rewrite $\mathcal{L}^{\text{FastGAE}}$ from equation 12 as follows:

$$\mathcal{L}^{\text{FastGAE}} = \frac{1}{n_{(S)}} \sum_{i \in \mathcal{V}} \mathbb{1}_{(i \in \mathcal{V}_{(S)})} \mathcal{L}^{\text{FastGAE}}(i),$$

where the node-level terms $\mathcal{L}^{\text{FastGAE}}(i)$ are defined as:

$$\mathcal{L}^{\text{FastGAE}}(i) = \frac{1}{n_{(S)}} \sum_{j \in \mathcal{V}} \mathbb{1}_{(j \in \mathcal{V}_{(S)})} \mathcal{L}_{ij}(A_{ij}, \hat{A}_{ij})$$

and where \mathcal{L}_{ij} denotes the cross entropy loss as in equation 11. In the following, we leverage concentration inequalities [50] to derive a theoretically-grounded threshold size, denoted $n_{(S)}^*$ in the following, for which, under mild assumptions, the (random) node-level deviation

$$|\mathcal{L}^{\text{FastGAE}}(i) - \mathbb{E}[\mathcal{L}^{\text{FastGAE}}(i)]|$$

at each training iteration is proven to be bounded with a high probability, for any node i . This proposed subgraph size is of the form:

$$n_{(S)}^* = C\sqrt{n},$$

where constant $C > 0$ depends on the deviation magnitude and probability, and is explicitly presented in Proposition 5. In our empirical analysis, this criterion will allow us to significantly improve the scalability and training speed of graph AE and VAE models (see discussion on complexity in Section 3.4), while reaching fairly competitive performances in a majority of experiments (see Section 4).

To prove our bounds, we require the following technical assumption on the reconstructed matrix \hat{A} :

Assumption 1. *Let $(i, j) \in \mathcal{V}^2$. We thereafter assume that $\hat{A}_{ij} = \sigma(z_i^T z_j)$ can actually be **capped**, and that:*

$$\hat{A}_{ij} \in [\varepsilon, 1 - \varepsilon]$$

where $0 < \varepsilon < 1$ is a constant that can be arbitrarily close to 0.

Under this assumption, we derive Propositions 4 and 5.

Proposition 4. *Let us consider a training iteration of the FastGAE framework, a sampled subgraph $\mathcal{G}_{(S)} = (\mathcal{V}_{(S)}, \mathcal{E}_{(S)})$, with $|\mathcal{V}_{(S)}| = n_{(S)} < n$ nodes sampled without replacement, and the corresponding node-level approximate reconstruction computed for a given node i :*

$$\mathcal{L}^{\text{FastGAE}}(i) = \frac{1}{n_{(S)}} \sum_{j \in \mathcal{V}} \mathbf{1}_{(j \in \mathcal{V}_{(S)})} \mathcal{L}_{ij}(A_{ij}, \hat{A}_{ij}),$$

with the random variable $\mathbf{1}_{(j \in \mathcal{V}_{(S)})} = 1$ if node $j \in \mathcal{V}_{(S)}$ and 0 otherwise, with $A_{ij} \in \{0, 1\}$ for all $(i, j) \in \mathcal{V}^2$ and with:

$$\mathcal{L}_{ij}(A_{ij}, \hat{A}_{ij}) = -[A_{ij} \log(\hat{A}_{ij}) + (1 - A_{ij}) \log(1 - \hat{A}_{ij})].$$

Then, under Assumption 1, for any $\gamma \geq 0$, we have:

$$\mathbb{P}(|\mathcal{L}^{\text{FastGAE}}(i) - \mathbb{E}[\mathcal{L}^{\text{FastGAE}}(i)]| \geq \gamma) \leq 2 \exp\left(-2\left(\frac{\gamma}{\log(\varepsilon)}\right)^2 \frac{n_{(S)}^2}{n}\right). \quad (18)$$

We note that the right hand side term tends to 0 exponentially fast w.r.t. the deviation magnitude γ and w.r.t. the subgraph size $n_{(S)}$.

Proposition 5. *For any confidence level $\alpha \in]0, 1[$ and node $i \in \mathcal{V}$, selecting a subgraph size $n_{(S)}$ such that*

$$n_{(S)} \geq n_{(S)}^* = \sqrt{n} \underbrace{\sqrt{\frac{-\log(\frac{\alpha}{2}) \log(\varepsilon)^2}{2\gamma^2}}}_{\text{denoted } C \text{ in eq. 10}} \quad (19)$$

guarantees that

$$\mathbb{P}(|\mathcal{L}^{\text{FastGAE}}(i) - \mathbb{E}[\mathcal{L}^{\text{FastGAE}}(i)]| \geq \gamma) \leq \alpha.$$

As an opening, we note that, while the current bounds are empirically effective (see Section 4), future research will aim at directly bounding the deviation of $\mathcal{L}^{\text{FastGAE}}$ instead of the node-level terms $\mathcal{L}^{\text{FastGAE}}(i)$, which would be more ambitious and challenging due to the inherent dependencies among sampled node pairs in FastGAE. Also, while Proposition 4 and 5 focus on the case of the *cross entropy loss* for consistency w.r.t. the main paper, a similar analysis (omitted here) could be performed to obtain comparable bounds for other *bounded* reconstruction losses. For instance, in the case of the Frobenius loss, where $\mathcal{L}_{ij}(A_{ij}, \hat{A}_{ij}) = (A_{ij} - \hat{A}_{ij})^2$, and without Assumption 1, one can obtain similar concentration guarantees as Proposition 5, with C being replaced by the constant $C' = \sqrt{\frac{-\log(\alpha/2)}{2\gamma^2}}$.

Numerical Application. In our experiments, all $n_{(S)}^*$ threshold subgraph sizes will be computed by evaluating equation 19, setting the following values for hyperparameters: $\gamma = 1$, $\alpha = 0.1$ and $\varepsilon = 0.001$.

3.4. Complexity of FastGAE

As previously detailed, both the encoder and the sampling step of FastGAE have a linear time complexity w.r.t. m . Moreover, our decoder runs in $O(n_{(S)}^2)$ time, with $n_{(S)}$ being significantly smaller than n in practice. In particular, setting $n_{(S)} = n_{(S)}^*$ ensures a

$O(n)$ time complexity for decoding (as $n_{(S)}^* = (C\sqrt{n})^2 = C^2n$) and an overall $O(m + n)$ linear time complexity for a complete FastGAE training iteration. Faster bounds can also be achieved by lowering $n_{(S)}$ or by replacing GCNs by another encoder. Therefore, as we will empirically verify in Section 4, our framework is significantly faster and more scalable than standard graph AE/VAE.

3.5. Differences with Related Work

Before diving into experiments, we would like to emphasize some key differences between FastGAE and existing works. Foremost, FastGAE is *not* directly comparable to the existing research cited in Section 2.3 to scale GCN models, e.g. to FastGCN [30] that also sample nodes. Indeed, FastGCN is a GCN-like model, optimized to classify node labels in a (semi) supervised fashion. It samples the *neighborhood* of each node when averaging vector representations in forward passes. On the contrary, in this paper, after *full* GCN forward passes, we instead sample *subgraphs to reconstruct*, in order to approximate the *reconstruction losses* of two unsupervised models, in which GCNs are only a building part (the encoder) of a larger framework (the AE or the VAE). Both settings therefore address different problems; as explained in Section 3.1, FastGCN, GraphSAGE or Cluster-GCN could actually be used *in conjunction with* FastGAE, as encoders.

Futhermore, FastGAE is also more elaborated than data cleaning methods that simply consist in removing some nodes from a graph, e.g. the low-degree ones, to reduce its size. Indeed, in the case of FastGAE with degree sampling, low-degree nodes are still 1) *fully* used in the GCN encoder, and 2) might also appear in *some* subgraphs that we decode (but less often than high-degree nodes). As we leverage new different subgraphs at each iteration, we explore different parts of the *entire* graph during training.

Last, we note that *effective subset selection* for faster learning has already provided promising results in the machine learning community [51, 52, 53]; however, contrary to these works, we focus on an unsupervised graph-based problem, and our sampling methods remain fixed throughout learning as we rely on graph mining to select $\mathcal{G}_{(S)}$.

4. Empirical Analysis

In this section, we present an in-depth experimental evaluation of our proposed framework to scale graph AE and VAE models.

4.1. Experimental Setting

4.1.1. Datasets

We provide experimental results on seven graphs of increasing sizes. Their statistics are presented in Table 1. We first study the Cora, Citeseer and Pubmed citation networks⁵, with and without node features corresponding to f -dimensional bag-of-words vectors (with $f = 1433, 3703$ and 500 respectively). Nodes are clustered in respectively 6, 7 and 3 topic classes, acting as ground truth communities. These datasets are common benchmarks for evaluating graph AE and VAE (see Kipf and Welling [11] and a majority of recent works [15, 16, 17, 18, 19, 26, 27, 28, 54]). For these *medium-size graphs*, we

⁵<https://linqs.soe.ucsc.edu/data>

| Dataset | Number of nodes | Number of edges |
|-----------------|-----------------|-----------------|
| Cora | 2708 | 5429 |
| Citeseer | 3327 | 4732 |
| Pubmed | 19717 | 44338 |
| SBM | 100000 | 1498844 |
| Google | 875713 | 4322051 |
| Youtube | 3223589 | 9375374 |
| Patent | 3774768 | 16518948 |

Table 1: Datasets Statistics

can directly compare the performance of FastGAE to standard graph AE and VAE, as training standard models is still computationally affordable.

Then, we consider four significantly *larger graphs*, with up to millions of nodes and edges, and for which training standard graph AE or VAE is *intractable*. We consider the Google⁶ hyperlinks web graph, the Youtube⁷ social network of users (friendship connections), the US Patent⁶ citation network, and a synthetic graph, denoted SBM, generated from a *stochastic block model* which is a generative model for random graphs [55]. In this last graph, by design, nodes are clustered in 100 groups of 1000 nodes, acting as ground truth communities. Two nodes from a same community (resp. from different communities) are connected by an edge with probability 2×10^{-2} (resp. 2×10^{-4}).

Our evaluation therefore includes graphs with various characteristics, sizes, and from four different families (citation networks, social networks, web graphs and stochastic block model graphs).

4.1.2. Evaluation Tasks

We consider two downstream tasks for evaluation:

- First, we consider a *link prediction* task. We train all models on masked graphs for which 15% of edges were randomly removed. Then, we create validation and test sets from the removed edges (resp. from 5% and 10% of edges) and from the same number of sampled unconnected node pairs. Using decoder predictions \hat{A}_{ij} , we evaluate our ability to classify edges from non-edges, using the mean *Area Under the ROC Curve* (AUC) and *Average Precision* (AP) scores⁸ on test sets. Link prediction is the most common task to evaluate graph AE and VAE models since the seminal research of Kipf and Welling [11] (see e.g. Salha et al. [36] and references therein for an overview), and we therefore found essential to consider it as well in our work.
- We also perform *node clustering* experiments, on datasets with ground truth communities. For this task, after training models on *complete* versions of the graphs, we

⁶<http://snap.stanford.edu/data/index.html>

⁷<http://konect.cc/networks/>

⁸We computed scores via scikit-learn [56]. Formulas are provided in https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc_auc_score.html for AUC, and in https://scikit-learn.org/stable/modules/generated/sklearn.metrics.average_precision_score.html for AP.

run k -means algorithms in embedding spaces to cluster the z_i vectors. We compare these clusters to the ground truth ones using the mean *Adjusted Mutual Information* (AMI) scores⁹ on test sets.

We emphasize that we chose to focus on these two tasks instead of a more direct reconstruction task (despite working with autoencoders), as AUC and AP scores from link prediction, as well as AMI scores from node clustering, are more insightful and understandable metrics than a direct reporting of some cross-entropy or ELBO reconstruction losses. Besides, we also aimed at providing a consistent experimental setting w.r.t. the existing literature on graph AE and VAE that, for the most part, focused on link prediction and, to some extent, on node clustering tasks.

4.1.3. Details on Models: Hyperparameters and Model Selection Procedure for Standard and FastGAE-based AE/VAE

In the upcoming experiments, for the aforementioned graphs and evaluation tasks, we compare standard graph AE and VAE models (when they are tractable) to FastGAE-based versions of these models.

All AE and VAE models, with and without FastGAE, were optimized for the *link prediction* task. More specifically, we selected the best sets of hyperparameters in terms of mean AUC scores *on the validation sets* introduced in Section 4.1.2. Instructions to easily run a similar validation are provided in our source code.

We trained models for 200 iterations (resp. 300) for graphs with $n < 100000$ (resp. $n \geq 100000$), and thoroughly checked the convergence of all models for these values (in terms of loss stability in the validation set). Other hyperparameters for these models are described thereafter.

Our encoders are 2-layer GCNs (we tested models with 1 to 3 layers). They include 32-dim hidden layers, and 16-dim output layer, which means that the dimension of embedding vectors (denoted d in Section 2.1.1 on AE and in Section 2.2.1 on VAE) is equal to $d = 16$. We emphasize that we also tested models with $d \in \{32, 64, 128\}$, reaching similar conclusions w.r.t. $d = 16$ (the impact of d is further discussed in Section 4.2.4).

Besides, for all models, we used the Adam optimizer [57], without dropout (we tested models with dropout values in $\{0, 0.1, 0.2, 0.3, 0.4, 0.5\}$). Regarding learning rates for such optimizer, we tested values from the grid $\{0.0001, 0.0005, 0.001, 0.005, 0.01, 0.05, 0.1, 0.2\}$. We eventually picked a learning rate of 0.1 for Patent with uniform sampling, and of 0.01 otherwise as, one again, these values returned the best mean AUC scores on validation sets. Last, as Kipf and Welling [11], we considered all graphs as *undirected* and ignored edges directions when available.

We used TensorFlow [58], training models on an NVIDIA GTX 1080 GPU, and running other operations on a double Intel Xeon Gold 6134 CPU.

4.1.4. Details on Models: Other Baselines

For completeness, we also compare standard graph AE/VAE and FastGAE-based models to the few other existing methods to scale graph AE/VAE, using a similar sets of hyperparameters and similar embedding values ($d = 16$) as previous section:

⁹We computed AMI scores via scikit-learn [56]. The formula is provided in https://scikit-learn.org/stable/modules/generated/sklearn.metrics.adjusted_mutual_info_score.html

- We consider a simple *negative sampling* strategy, briefly mentioned by Kipf and Welling [11] and suggested by Fey and Lenssen [59], where we reconstruct all edges but only $|\mathcal{E}|$ randomly picked unconnected node pairs to compute losses. We leveraged methods made available by Fey and Lenssen [59] to estimate losses, with consistent dropout values, learning rates and architectures w.r.t. Section 4.1.3.
- We also consider the framework recently proposed by Salha et al. [18], denoted as *Core-GAE* in next tables. Authors train the AE/VAE only on the smaller graph k -core, then propagate embedding representations to other nodes out of the k -core via simple heuristics; k is a parameter tuning the size of the input graph for learning (which ranges from 1 to the maximal k for which the corresponding k -core is not empty). We used the author’s implementation [18] with optimal values (regarding mean AUC scores on validation sets) for the hyperparameter k detailed in next tables, and with consistent dropout values, learning rates and architectures w.r.t. Section 4.1.3.
- Besides, the other aforementioned sampling ideas briefly mentioned (as possible extensions) in the recent literature [27, 36] actually are particular cases of FastGAE, namely with *uniform* sampling.

Last, in addition to an extensive comparison between the different AE/VAE models, we also report results obtained with the following non AE/VAE-based baselines:

- A *spectral embedding*, which is a powerful but not scalable baseline. We used the implementation provided by Pedregosa et al. [56]; embedding axes correspond the the d first eigenvectors of \mathcal{G} ’s Laplacian matrix, excluding non-informative vectors, and with d denoting the embedding dimension [60] as for all other models.
- *node2vec* [8], another very popular and scalable node embedding method. We trained models with hyperparameters $p = 1$ and $q = 1$, from 10 random walks of length 80 per node, with a window size of 5 and on a single epoch, and using the author’s implementation [8]. We omit comparison to other random walk-based models [6, 7] due to very similar performances on some of our preliminary tests.
- For node clustering, we also compare our approach to *Louvain*’s scalable community detection algorithm, from a direct usage of the authors’ implementation [61].

4.2. Results

In the remainder of this section, we provide an empirical evaluation of FastGAE and of its variational FastGAE variant.

4.2.1. Preliminary Insights on High Degree/Core Nodes

Before studying FastGAE we report important insights from preliminary experiments on standard graph AE/VAE. They motivated the design of our framework and emphasize the relevance of sampling high-degree/core nodes. On the medium-size Cora, Citeseer and Pubmed graphs, we trained standard graph AE/VAE models, but *tried to mask k nodes and their edges* from the computation of reconstruction losses, for different values of k . Such masking procedure is expected to lower performances, as the model leverages less information about the quality of the reconstruction for learning.

Figure 1 shows that, when these k removed nodes are *the top- k highest degrees/cores nodes*, performances on the link prediction task tumble down. On the contrary, removing *the k nodes with minimal degrees or core numbers* from the loss leads to almost no drop, and even slightly better results on Pubmed, which suggests that removing non-informative nodes might even be beneficial for learning.

In Figure 2, we report similar results on Adjusted Mutual Information scores for node clustering. These ablation studies suggest that, when implementing stochastic subgraph decoding strategies for scalability, sampling high-degree/core nodes is indeed crucial to learn meaningful embeddings. FastGAE, which explicitly exploits these structural node properties, and optimizes a reconstruction loss that re-weights high degrees/cores node pairs, behaves consistently w.r.t. such important insights.

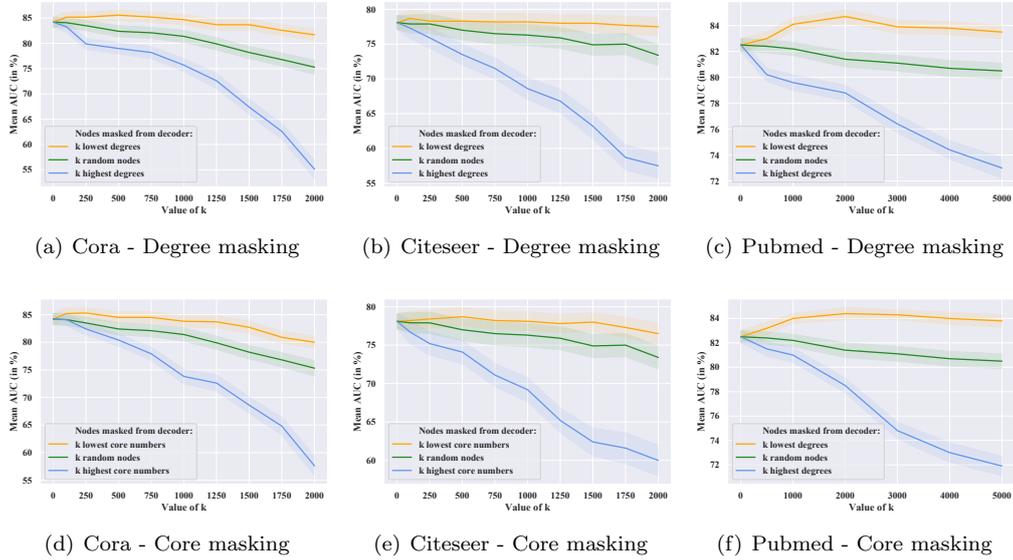


Figure 1: Link prediction on featureless Cora, Citeseer and Pubmed using standard Graph VAE models, but trained while masking k nodes and their connections from the decoder/reconstruction loss. AUC scores are averaged over 100 runs with random train/test splits.

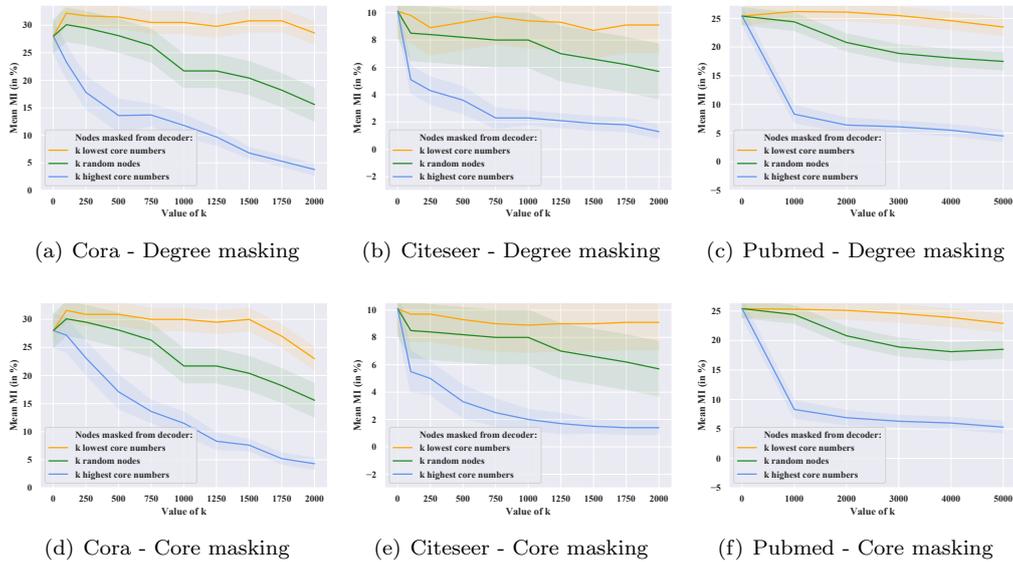


Figure 2: Node clustering on featureless Cora, Citeseer and Pubmed using standard Graph VAE models, but trained while masking k nodes and their connections from the decoder/reconstruction loss. Adjusted MI scores are averaged over 100 runs with random train/test splits.

| Model (Dimension $d = 16$) | Subgraphs size $n_{(S)}$ | Average Perf. on Test Set | | Average Running Times (in seconds) | | | Speed gain w.r.t. GAE |
|---|-----------------------------|------------------------------------|------------------------------------|------------------------------------|----------------|-------------|--------------------------|
| | | AUC (in %) | AP (in %) | Compute p_i | Train model | Total | |
| Standard Graph AE | - | 82.51 \pm 0.64 | 87.42 \pm 0.38 | - | 811.43 | 811.43 | - |
| FastGAE with degree sampling ($\alpha = 1$) | 5000 | 84.82 \pm 0.32 | 88.19 \pm 0.23 | 0.01 | 14.41 | 14.42 | \times 56.27 |
| | 2500 | 84.12 \pm 0.40 | 87.56 \pm 0.30 | 0.01 | 5.72 | 5.73 | \times 141.61 |
| | 1187* | 83.67 \pm 0.42 | 87.01 \pm 0.31 | 0.01 | 3.20 | 3.21 | \times 252.78 |
| | 500 | 82.68 \pm 0.51 | 85.89 \pm 0.47 | 0.01 | 2.98 | 2.99 | \times 271.38 |
| | 250 | 80.77 \pm 0.55 | 84.05 \pm 0.51 | 0.01 | 2.83 | 2.84 | \times 285.71 |
| FastGAE with core sampling ($\alpha = 2$) | 5000 | 84.62 \pm 0.24 | 88.09 \pm 0.16 | 1.75 | 15.98 | 17.73 | \times 45.77 |
| | 2500 | 83.69 \pm 0.34 | 87.28 \pm 0.31 | 1.75 | 7.51 | 9.26 | \times 87.63 |
| | 1187* | 82.53 \pm 0.46 | 86.28 \pm 0.37 | 1.75 | 4.81 | 6.56 | \times 123.69 |
| | 500 | 80.96 \pm 0.52 | 84.86 \pm 0.46 | 1.75 | 4.57 | 6.32 | \times 128.39 |
| | 250 | 79.53 \pm 0.53 | 83.10 \pm 0.50 | 1.75 | 4.44 | 6.19 | \times 131.08 |
| FastGAE with uniform sampling | 5000 | 81.08 \pm 0.48 | 85.90 \pm 0.60 | - | 13.90 | 13.90 | \times 58.37 |
| | 2500 | 78.72 \pm 0.74 | 83.50 \pm 0.75 | - | 5.48 | 5.48 | \times 148.07 |
| | 1187* | 77.28 \pm 0.89 | 81.89 \pm 0.91 | - | 3.10 | 3.10 | \times 261.75 |
| | 500 | 75.09 \pm 2.05 | 78.53 \pm 2.04 | - | 2.98 | 2.98 | \times 271.29 |
| | 250 | 74.12 \pm 2.07 | 77.72 \pm 1.22 | - | 2.82 | 2.82 | \times 287.74 |
| Core-GAE, $k = 2$ (best choice) | - | 84.30 \pm 0.27 | 86.11 \pm 0.43 | - | 168.91 | 168.91 | \times 4.80 |
| Core-GAE, $k = 9$ (fastest choice) | - | 61.65 \pm 0.94 | 64.82 \pm 0.72 | - | 2.92 | 2.92 | \times 277.89 |
| Negative Sampling GAE | - | 81.19 \pm 0.68 | 83.21 \pm 0.40 | - | 111.79 | 111.79 | \times 7.28 |
| node2vec | - | 81.25 \pm 0.26 | 85.55 \pm 0.26 | - | 48.91 | 48.91 | \times 16.59 |
| Spectral Embedding | - | 83.14 \pm 0.42 | 86.55 \pm 0.41 | - | 31.71 | 31.71 | \times 25.59 |

Table 2: Link prediction on the featureless Pubmed graph ($n = 19717$, $m = 44338$) using standard Graph AE, FastGAE with degree, core and uniform sampling, and baselines. For degree and core sampling, values of the hyperparameter α (as defined in equation 9) were tuned, as described in Figure C.6. All models learn embedding vectors of dimension $d = 16$. Scores are averaged over 100 runs with different and random train/validation/test sets. Bold numbers correspond to the best performance (several numbers are bold when scores are comparable, in a ± 1 standard deviation range) and best running time. Subgraphs sizes annotated with * correspond to the $n_{(S)}^*$ threshold, as introduced in equation 19.

4.2.2. FastGAE for Medium-Size Graphs

We now evaluate FastGAE and its variational FastGAE variant. First, we focus on *medium-size graphs*. For Cora, Citeseer and Pubmed, we can compare FastGAE to standard graph AE/VAE. The above Table 2 details mean AUC and AP scores and standard errors over 100 runs with different train/test splits for link prediction on (featureless) Pubmed with AE models. For the sake of brevity, we report more summarized results for other medium-size graphs, for VAE and for node clustering, in Table 3 and Figure 3 (for link prediction) as well as in Table 4 and Figure 4 (for node clustering).

FastGAE vs Standard Graph AE/VAE. Let us first compare FastGAE to standard graph AE/VAE models. In Table 2, we observe that, for sample sizes roughly 20 times smaller than n , FastGAE models with *degree* and *core* sampling both achieve competitive or even outperforming¹⁰ results w.r.t. standard graph AE on Pubmed (e.g. +2.31 AUC points for FastGAE with degree sampling and $n_{(S)} = 5000$).

Futhermore, FastGAE models are also significantly *faster*: in Table 2 for instance, our approach with degree sampling is up to $\times 252.78$ faster without performance degradation. The additional operation required by our framework, i.e. computing the p_i distribution,

¹⁰At first glance, the fact that FastGAE sometimes even slightly *outperforms* standard graph AE/VAE models might be surprising. This improvement is actually consistent with recent research on the benefits of mini-batch-based GNNs [62, 63]. It comes from the relevance of the two sampling schemes that we consider (core-based and degree-based) and from the stochastic nature of the training, that might tend to avoid local minima more easily [64].

| Dataset | Model (Dimension $d = 16$) | Average Perf. AUC (in %) | on Test Set AP (in %) | Avg. Run. Times (in sec.) Comp. p_i | Train model | Total | Speed Gain |
|------------------------------|--|------------------------------------|------------------------------------|--|----------------|--------------|-----------------|
| Cora | Standard Graph AE FastGAE (degree, $\alpha = 2$) | 84.79 ± 1.10 | 88.45 ± 0.82 | - | 3.87 | 3.87 | - |
| | - with $n_{(S)} = 250$ | 84.13 ± 1.20 | 86.65 ± 1.23 | 0.002 | 1.46 | 1.462 | $\times 2.65$ |
| | - with $n_{(S)} = n_{(S)}^* = 440$ | 84.74 ± 0.81 | 87.42 ± 0.75 | 0.002 | 1.56 | 1.562 | $\times 2.48$ |
| | - with $n_{(S)} = 1000$ | 84.75 ± 0.84 | 87.77 ± 0.81 | 0.002 | 1.65 | 1.652 | $\times 2.34$ |
| | <u>Best baseline</u> Spectral Embedding | 86.49 ± 0.98 | 87.42 ± 1.04 | - | 2.49 | 2.49 | $\times 1.55$ |
| Cora with features | Standard Graph VAE Var. FastGAE (degree, $\alpha = 2$) | 91.64 ± 0.92 | 92.66 ± 0.91 | - | 4.25 | 4.25 | - |
| | - with $n_{(S)} = 250$ | 90.50 ± 1.10 | 91.10 ± 1.08 | 0.002 | 2.30 | 2.302 | $\times 1.85$ |
| | - with $n_{(S)} = n_{(S)}^* = 440$ | 90.82 ± 1.07 | 91.44 ± 1.13 | 0.002 | 2.52 | 2.522 | $\times 1.69$ |
| | - with $n_{(S)} = 1000$ | 91.72 ± 0.98 | 92.36 ± 1.11 | 0.002 | 2.87 | 2.872 | $\times 1.48$ |
| | <u>Best baseline</u> Core-Graph VAE, $k = 2$ | 87.94 ± 1.12 | 89.00 ± 1.11 | - | 3.09 | 3.09 | $\times 1.38$ |
| Citeseer | Standard Graph AE FastGAE (degree, $\alpha = 1$) | 78.25 ± 1.69 | 83.79 ± 1.24 | - | 5.25 | 5.25 | - |
| | - with $n_{(S)} = 250$ | 77.28 ± 1.11 | 81.29 ± 0.92 | 0.002 | 1.47 | 1.472 | $\times 3.57$ |
| | - with $n_{(S)} = n_{(S)}^* = 488$ | 78.30 ± 1.30 | 82.42 ± 1.09 | 0.002 | 1.58 | 1.582 | $\times 3.32$ |
| | - with $n_{(S)} = 1000$ | 78.31 ± 1.25 | 82.40 ± 0.99 | 0.002 | 1.61 | 1.612 | $\times 3.26$ |
| | <u>Best baseline</u> Spectral Embedding | 80.42 ± 1.38 | 83.75 ± 1.12 | - | 3.50 | 3.50 | $\times 1.50$ |
| Citeseer with features | Standard Graph VAE Var. FastGAE (degree, $\alpha = 1$) | 90.72 ± 1.01 | 92.05 ± 0.97 | - | 6.28 | 6.28 | - |
| | - with $n_{(S)} = 250$ | 89.37 ± 1.69 | 89.63 ± 1.83 | 0.002 | 2.32 | 2.322 | $\times 2.70$ |
| | - with $n_{(S)} = n_{(S)}^* = 488$ | 90.10 ± 1.33 | 90.15 ± 1.50 | 0.002 | 2.62 | 2.622 | $\times 2.40$ |
| | - with $n_{(S)} = 1000$ | 90.22 ± 1.14 | 90.16 ± 1.20 | 0.002 | 2.89 | 2.892 | $\times 2.17$ |
| | <u>Best baseline</u> Core-Graph VAE, $k = 2$ | 81.85 ± 1.72 | 83.65 ± 1.64 | - | 2.55 | 2.55 | $\times 2.46$ |
| Pubmed | Standard Graph AE FastGAE (degree, $\alpha = 1$) | 82.51 ± 0.64 | 87.42 ± 0.38 | - | 811.43 | 811.43 | - |
| | - with $n_{(S)} = 500$ | 82.68 ± 0.51 | 85.89 ± 0.47 | 0.01 | 2.98 | 2.99 | $\times 271.38$ |
| | - with $n_{(S)} = n_{(S)}^* = 1187$ | 83.67 ± 0.42 | 87.01 ± 0.31 | 0.01 | 3.20 | 3.21 | $\times 252.78$ |
| | - with $n_{(S)} = 5000$ | 84.82 ± 0.32 | 88.19 ± 0.23 | 0.01 | 14.41 | 14.42 | $\times 56.27$ |
| | <u>Best baseline</u> Core-Graph AE, $k = 2$ | 84.30 ± 0.27 | 86.11 ± 0.43 | - | 168.91 | 168.91 | $\times 4.80$ |
| Pubmed with features | Standard Graph AE FastGAE (degree, $\alpha = 1$) | 96.28 ± 0.36 | 96.29 ± 0.25 | - | 952.63 | 952.63 | - |
| | - with $n_{(S)} = 500$ | 95.08 ± 0.45 | 95.24 ± 0.46 | 0.01 | 3.53 | 3.54 | $\times 269.10$ |
| | - with $n_{(S)} = n_{(S)}^* = 1187$ | 95.45 ± 0.26 | 95.70 ± 0.30 | 0.01 | 4.01 | 4.02 | $\times 237.56$ |
| | - with $n_{(S)} = 5000$ | 96.12 ± 0.20 | 96.35 ± 0.19 | 0.01 | 19.74 | 19.75 | $\times 48.23$ |
| | <u>Best baseline</u> Core-Graph AE, $k = 2$ | 85.34 ± 0.33 | 86.06 ± 0.24 | - | 40.22 | 40.22 | $\times 23.69$ |

Table 3: Link prediction on all medium-size graphs. For each graph, for brevity, we only report the **best** graph AE or VAE model in terms of AUC and AP scores, a few representative degree-based FastGAE versions of this model, and the best baseline (among Core-Graph AE/VAE, Negative Sampling Graph AE/VAE, node2vec and the spectral embedding). Scores are averaged over 100 runs with different and random train/validation/test sets. For degree sampling, values of the hyperparameter α (as defined in equation 9) were tuned, as described in Figure C.6. All models learn embedding vectors of dimension $d = 16$. Bold numbers correspond to the best performance (several numbers are bold when scores are comparable, in a ± 1 standard deviation range) and best running time.

is efficient in practice, especially for degree sampling. By further reducing the subgraph size $n_{(S)}$, one can achieve even faster results, while only losing a few AUC/AP points in performance.

In Table 3, Table 4, Figure 3 and Figure 4, we consolidate our results by reaching similar conclusions on VAE, on other medium-size graphs (with and without features), and on node clustering. On Figures 3 and 4, we also confirm that, even for relatively low $n_{(S)}/n$ proportions, our proposed method achieves comparable performances w.r.t. standard models.

| Dataset | Model (Dimension $d = 16$) | Average Performance AMI (in %) | Average Compute p_i | Average Running Times (in sec.) Train model | Total | Speed Gain |
|------------------------------|---|------------------------------------|--------------------------|--|--------------|-----------------|
| Cora | Standard Graph AE | 30.88 ± 2.56 | - | 3.90 | 3.90 | - |
| | FastGAE (degree, $\alpha = 2$) | | | | | |
| | - with $n_{(S)} = 250$ | 33.32 ± 2.61 | 0.002 | 1.51 | 1.512 | $\times 2.58$ |
| | - with $n_{(S)} = n_{(S)}^* = 440$ | 34.64 ± 2.45 | 0.002 | 1.59 | 1.592 | $\times 2.45$ |
| | - with $n_{(S)} = 1000$ | 35.56 ± 2.80 | 0.002 | 1.67 | 1.672 | $\times 2.33$ |
| | <u>Best baseline</u> Louvain | 46.72 \pm 0.85 | - | 1.79 | 1.79 | $\times 2.18$ |
| Cora with features | Standard Graph VAE | 44.84 \pm 2.63 | - | 4.32 | 4.32 | - |
| | Var. FastGAE (degree, $\alpha = 2$) | | | | | |
| | - with $n_{(S)} = 250$ | 41.35 ± 3.49 | 0.002 | 2.40 | 2.402 | $\times 1.80$ |
| | - with $n_{(S)} = n_{(S)}^* = 440$ | 42.89 ± 2.72 | 0.002 | 2.67 | 2.672 | $\times 1.62$ |
| | - with $n_{(S)} = 1000$ | 45.02 \pm 2.81 | 0.002 | 2.92 | 2.922 | $\times 1.48$ |
| | <u>Best baseline</u> Louvain | 46.72 \pm 0.85 | - | 1.79 | 1.79 | $\times 2.41$ |
| Citeseer | Standard Graph VAE | 9.85 ± 1.24 | - | 5.44 | 5.44 | - |
| | Var. FastGAE (degree, $\alpha = 1$) | | | | | |
| | - with $n_{(S)} = 250$ | 9.34 ± 1.48 | 0.002 | 1.77 | 1.772 | $\times 3.07$ |
| | - with $n_{(S)} = n_{(S)}^* = 488$ | 10.02 ± 1.42 | 0.002 | 2.02 | 2.022 | $\times 2.69$ |
| | - with $n_{(S)} = 1000$ | 10.16 ± 1.41 | 0.002 | 2.19 | 2.192 | $\times 2.48$ |
| | <u>Best baseline</u> Louvain | 16.39 \pm 1.45 | - | 2.41 | 2.41 | $\times 2.26$ |
| Citeseer with features | Standard Graph VAE | 20.17 \pm 3.07 | - | 6.45 | 6.45 | - |
| | Var. FastGAE (degree, $\alpha = 1$) | | | | | |
| | - with $n_{(S)} = 250$ | 20.49 \pm 3.74 | 0.002 | 2.80 | 2.802 | $\times 2.30$ |
| | - with $n_{(S)} = n_{(S)}^* = 488$ | 20.53 \pm 3.45 | 0.002 | 2.88 | 2.882 | $\times 2.24$ |
| | - with $n_{(S)} = 1000$ | 20.94 \pm 3.21 | 0.002 | 3.11 | 3.112 | $\times 2.07$ |
| | <u>Best baseline</u> Cora-Graph VAE, $k = 2$ | 16.53 ± 1.95 | - | 2.76 | 2.76 | $\times 2.33$ |
| Pubmed | Standard Graph VAE | 20.52 ± 2.97 | - | 856.05 | 856.05 | - |
| | Var. FastGAE (degree, $\alpha = 1$) | | | | | |
| | - with $n_{(S)} = 500$ | 16.86 ± 4.84 | 0.01 | 3.17 | 3.18 | $\times 269.20$ |
| | - with $n_{(S)} = n_{(S)}^* = 1187$ | 18.84 ± 4.78 | 0.01 | 3.61 | 3.62 | $\times 236.49$ |
| | - with $n_{(S)} = 5000$ | 22.81 \pm 4.80 | 0.01 | 14.95 | 14.96 | $\times 57.22$ |
| | <u>Best baseline</u> Core-Graph VAE, $k = 2$ | 23.56 \pm 3.12 | - | 50.11 | 50.11 | $\times 17.08$ |
| Pubmed with features | Standard Graph VAE | 25.43 ± 1.47 | - | 970.67 | 970.67 | - |
| | Var. FastGAE (degree, $\alpha = 1$) | | | | | |
| | - with $n_{(S)} = 500$ | 29.04 ± 4.17 | 0.01 | 4.03 | 4.04 | $\times 240.26$ |
| | - with $n_{(S)} = n_{(S)}^* = 1187$ | 31.11 \pm 3.27 | 0.01 | 4.65 | 4.66 | $\times 208.30$ |
| | - with $n_{(S)} = 5000$ | 30.89 \pm 3.01 | 0.01 | 20.01 | 20.02 | $\times 48.49$ |
| | <u>Best baseline</u> Core-Graph VAE, $k = 2$ | 24.35 ± 1.55 | - | 57.09 | 57.09 | $\times 17.00$ |
| SBM | Standard Graph VAE | (intractable) | | (intractable) | | - |
| | Var. FastGAE (degree, $\alpha = 2$) | | | | | |
| | - with $n_{(S)} = 2500$ | 30.77 ± 0.32 | 0.03 | 52.01 | 52.04 | - |
| | - with $n_{(S)} = n_{(S)}^* = 2673$ | 30.89 ± 0.30 | 0.03 | 53.98 | 54.01 | - |
| | - with $n_{(S)} = 5000$ | 32.28 ± 0.26 | 0.03 | 61.96 | 61.69 | - |
| | <u>Best baseline</u> Louvain | 35.90 \pm 0.14 | - | 464.11 | 464.11 | - |

Table 4: Node clustering on all graphs with communities. For each graph, for brevity, we only report the **best** graph AE or VAE model in terms of mean AMI, a few representative degree-based FastGAE versions of this model, and the best baseline (among Core-Graph AE/VAE, Negative Sampling Graph AE/VAE, node2vec, Louvain and the spectral embedding). Scores are averaged over 100 runs (resp. 10 runs) for medium-size graphs (resp. for the large graph SBM). For degree sampling, values of the hyperparameter α (as defined in equation 9) were tuned, as described in Figure C.6. All models learn embedding vectors of dimension $d = 16$. Bold numbers correspond to the best performance (several numbers are bold when scores are comparable, in a ± 1 standard deviation range) and best running time.

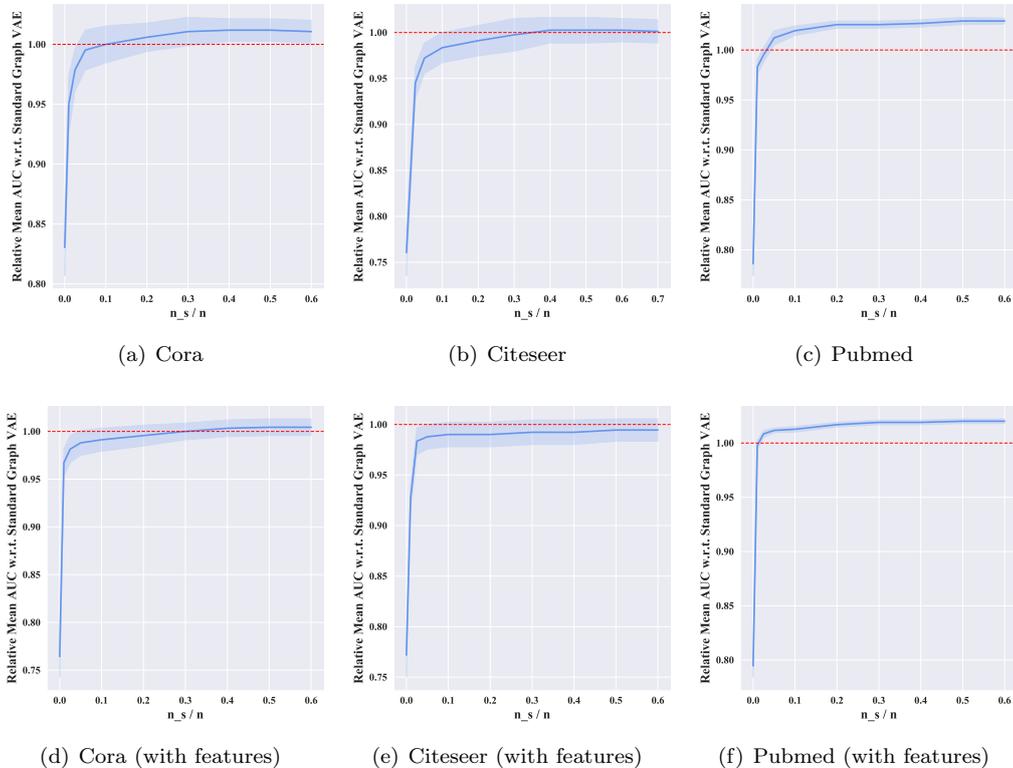


Figure 3: Summarized results for link prediction on the medium-size graphs Cora, Citeseer and Pubmed: relative mean AUC scores of degree-based Variational FastGAE models w.r.t. standard Graph VAE depending on the proportion of sampled nodes $n_{(S)}/n$ in decoders. We observe that, even for relatively low $n_{(S)}/n$ proportions, Variational FastGAE achieves comparable or even slightly better performances w.r.t. standard Graph VAE (results above the red line).

Comparison of Uniform, Core-based and Degree-based FastGAE. In all our experiments, we observe that FastGAE with core and degree sampling both outperform FastGAE (and variational FastGAE) with uniform sampling. Furthermore, core and degree sampling also return more stable scores, i.e. with lower standard errors, especially when the number of samples $n_{(S)}$ is relatively small. Such results confirm the empirical superiority of strategies that leverage the graph structure w.r.t. pure random strategies.

FastGAE vs Baselines. In Table 2, Table 3 and Table 4, these models also outperform the other few existing methods to scale graph AE and VAE, usually by a wide margin. For instance, in Table 2, we show that, to achieve (almost) comparable link prediction performances w.r.t. FastGAE on Pubmed, *Core-GAE* [18] requires longer running times (see *Core-GAE* with $k = 2$), and that faster variants significantly underperform (almost -20 AUC points for *Core-GAE* with $k = 9$ w.r.t. FastGAE with degree sampling). FastGAE is also conceptually simpler than Core-GAE, which we consider to be another advantage of our approach.

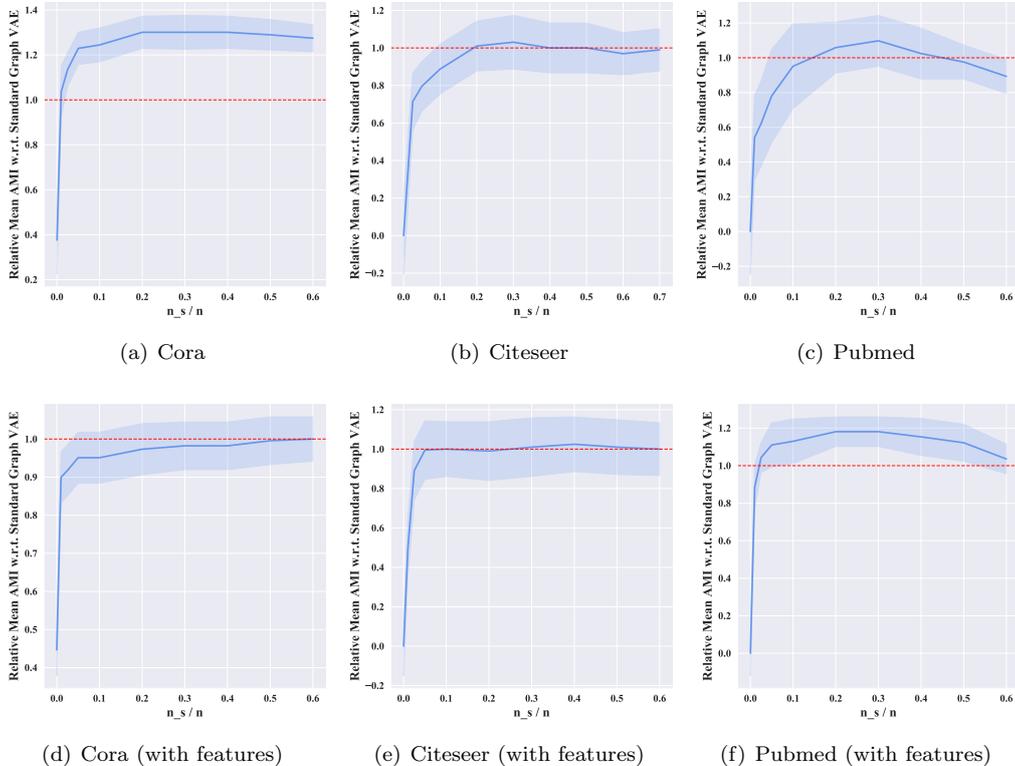


Figure 4: Summarized results for node clustering on the medium-size graphs Cora, Citeseer and Pubmed: relative mean AMI scores of degree-based Variational FastGAE models w.r.t. standard Graph VAE depending on the proportion of sampled nodes $n_{(S)}/n$ in decoders. We observe that, even for relatively low $n_{(S)}/n$ proportions, Variational FastGAE achieves comparable or even slightly better performances w.r.t. standard Graph VAE (results above the red line).

Besides, FastGAE-based models are faster and more effective than the ones leveraging negative sampling [59] (e.g. +3.63 AUC points for FastGAE with degree sampling and $n_{(S)} = 20000$ w.r.t. *Negative Sampling GAE* in Table 2). This performance gain might be explained by the more systematic inclusion of *unconnected pairs of important nodes*¹¹ in the losses of FastGAE-based models.

Last, but not least, our proposed framework is also competitive w.r.t. the popular non AE/VAE-based baselines in most cases. The only exception concerns the node clustering experiments on Cora and Citeseer (see Table 4) where the Louvain baseline [61] outperforms AE/VAE models, which we will further discuss in Section 4.2.3.

¹¹Indeed, when performing negative sampling for graph AE, we only reconstruct a few random unconnected node pairs, and ignore the others. However, reconstructing some of these neglected pairs might actually be crucial. Let us consider two nodes with high core number or centrality: knowing that these two important nodes are *not* connected is critical to learn meaningful embeddings. The FastGAE sampling scheme ensures a more systematic inclusion of these important "negative pairs" in the decoding step than negative sampling.

| Model (Dimension $d = 16$) | Subgraphs size $n_{(S)}$ | Average Perf. on Test Set | | Average Running Times (in seconds) | | |
|--|-----------------------------|---------------------------|-------------------------|------------------------------------|----------------|-----------------------|
| | | AUC (in %) | AP (in %) | Compute p_i | Train model | Total |
| Standard Graph AE | - | <i>(intractable)</i> | | <i>(intractable)</i> | | |
| FastGAE with degree sampling ($\alpha = 2$) | 20000 | 92.91 \pm 0.22 | 93.35 \pm 0.21 | 0.30 | 4401.67 | 4401.97 (1h13) |
| | 16425* | 93.02 \pm 0.23 | 93.39 \pm 0.23 | 0.30 | 3693.32 | 3693.62 (1h02) |
| | 10000 | 91.76 \pm 0.23 | 91.74 \pm 0.21 | 0.30 | 1164.22 | 1164.52 (19 min) |
| | 2500 | 87.53 \pm 0.50 | 87.42 \pm 0.51 | 0.30 | 537.99 | 538.29 (9 min) |
| | 1000 | 85.55 \pm 0.62 | 85.96 \pm 0.55 | 0.30 | 500.12 | 500.42 (8 min) |
| FastGAE with core sampling ($\alpha = 2$) | 20000 | 90.71 \pm 0.21 | 91.70 \pm 0.19 | 668.05 | 4800.58 | 5468.63 (1h31) |
| | 16425* | 90.48 \pm 0.21 | 90.85 \pm 0.23 | 668.05 | 4027.90 | 4695.95 (1h18) |
| | 10000 | 89.08 \pm 0.25 | 88.65 \pm 0.24 | 668.05 | 1232.03 | 1900.08 (32 min) |
| | 2500 | 82.50 \pm 0.51 | 81.42 \pm 0.60 | 668.05 | 544.64 | 1222.69 (20 min) |
| | 1000 | 73.99 \pm 0.70 | 75.24 \pm 0.74 | 668.05 | 503.88 | 1171.93 (19 min) |
| FastGAE with uniform sampling | 20000 | 85.97 \pm 0.26 | 87.71 \pm 0.25 | - | 4397.89 | 4387.89 (1h13) |
| | 16425* | 84.40 \pm 0.25 | 86.11 \pm 0.25 | - | 3602.66 | 3602.66 (1h) |
| | 10000 | 83.77 \pm 0.28 | 83.37 \pm 0.26 | - | 1106.01 | 1106.01 (18 min) |
| | 2500 | 70.66 \pm 0.35 | 71.16 \pm 0.38 | - | 485.03 | 485.03 (8 min) |
| | 1000 | 59.34 \pm 0.83 | 58.83 \pm 1.30 | - | 438.02 | 438.02 (7 min) |
| Core-GAE, $k = 14$ (best choice) | - | 88.06 \pm 0.27 | 88.94 \pm 0.23 | - | 4805.11 | 4805.11 (1h20) |
| Core-GAE, $k = 21$ (fastest choice) | - | 86.94 \pm 0.69 | 87.23 \pm 0.71 | - | 619.01 | 619.01 (10 min) |
| Negative Sampling GAE | - | 86.11 \pm 0.48 | 86.70 \pm 0.49 | - | 2392.96 | 2392.96 (40 min) |
| node2vec | - | 92.96 \pm 0.23 | 93.43 \pm 0.17 | - | 25851.39 | 25851.39 (7h11) |
| Spectral Embedding | - | <i>(intractable)</i> | | <i>(intractable)</i> | | |

Table 5: Link prediction on Patent ($n = 3774768$, $m = 16518948$), using FastGAE with degree, core and uniform sampling, and baselines. Standard Graph AE is intractable. For degree and core sampling, values of the hyperparameter α (as defined in equation 9) were tuned, as described in Figure C.6. All models learn embedding vectors of dimension $d = 16$. Scores are averaged over 10 runs with different and random train/validation/test sets. Bold numbers correspond to the best performance (several numbers are bold when scores are comparable, in a ± 1 standard deviation range) and best running time. Subgraphs sizes annotated with * correspond to the $n_{(S)}^*$ threshold, as introduced in equation 19.

On the hyperparameter α . In Appendix C, we report optimal values of α for all graphs. We recall that $\alpha \in \mathbb{R}^+$ is the hyperparameter introduced in equation 9, that helps balancing important and "less important" nodes during sampling. Setting $\alpha = 0$, leads to the uniform sampling setting where all nodes are sampled with an equal probability. On the contrary, by setting $\alpha \rightarrow \infty$ we would always sample the most important nodes and ignore others. Experiments from Figure C.6 from Appendix C show that these two extreme cases are usually sub-optimal, and that a careful tuning of α (e.g. $\alpha = 2$ for core sampling in Table 2) improves performances.

On the threshold $n_{(S)}^$.* In Section 3.3, we introduced a theoretically-grounded threshold $n_{(S)}^* = C\sqrt{n}$ to select the subgraph size. Overall, in all our experiments (see Table 3, Table 4 and Table 5), selecting the proposed $n_{(S)}^*$ provided interesting performance/speed trade-offs, leading to fairly competitive results w.r.t. standard graph AE/VAE and best baselines, while being significantly faster.

4.2.3. FastGAE for Large Graphs

After studying medium-size graphs, we now report in this section the evaluation of FastGAE and variational FastGAE on the four *large graphs* from our experiments: SBM, Google, Youtube and Patent. The above Table 5 details mean AUC and AP scores and standard errors over 10 runs with different train/test splits for link prediction on the Patent graph with FastGAE. We also report more summarized results (for the sake of brevity) for link prediction on SBM, Google, Youtube and Patent in Figure 5 and in Table 6, and summarized results for node clustering on SBM in Figure 5 and in the

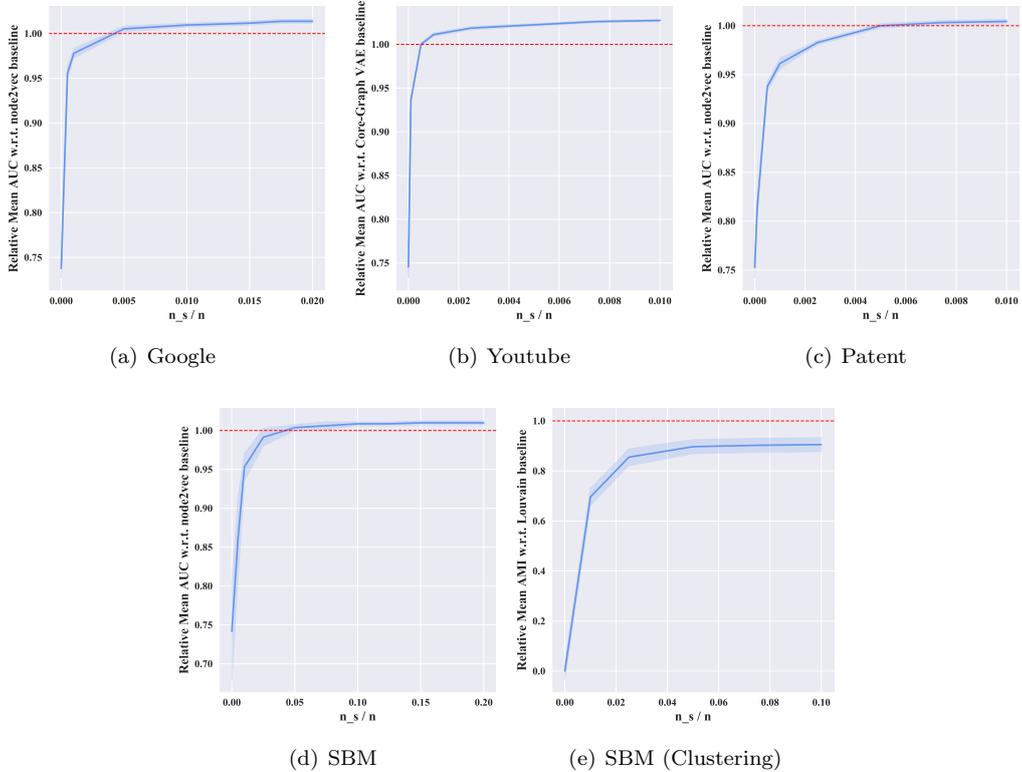


Figure 5: Summarized results for link prediction on the four large graphs SBM, Google, Youtube and Patent (subfigures a, b, c and d) and for node clustering on SBM (subfigure e): relative mean AUC scores (or mean AMI scores for subfigure e) of degree-based Variational FastGAE models w.r.t. the best scalable baseline, depending on the proportion of sampled nodes $n_{(S)}/n$ in decoders.

previous Table 4. As in Table 5, all scores are averaged over 10 runs with different train/test splits.

FastGAE vs Scalable AE/VAE Baselines. On large graphs, direct comparison with standard graph AE and VAE is unfortunately impossible. However, our FastGAE and variational FastGAE models almost always outperform the other existing approaches to scale graph AE and VAE, usually by a wide margin. For instance, for link prediction on Patent (Table 5), degree-based and core-based FastGAE models with $n_{(S)} = 20000$, 16425 and 10000 all outperform the best *Core-GAE* by up to roughly 5 AUC points (for degree-based FastGAE with $n_{(S)} = 20000$) and with comparable or better running times.

Regarding the *Core-GAE* baseline [18], we also point out that, in one of our large graphs, namely on the SBM one, this method was even *intractable* due to the lack of *size decreasing* core structure on this graph. Indeed, the 21-core of SBM includes 95200 nodes, which is too large to train a graph AE or VAE on our machines, and the 22-core is empty. Requiring a size decreasing core structure is a drawback of *Core-GAE* w.r.t. the more flexible FastGAE approach.

Moreover, as for medium-size graphs, we also observe that core-based and degree-based FastGAE tend to significantly outperform negative sampling (e.g. up to +6.8 AUC points for link prediction on Patent in Table 5; also, *Negative Sampling GAE* never appears as the best baseline in Table 4 nor in Table 6), consolidating our previous conclusions.

Besides, as before, the proposed $n_{(S)}^*$ provides quite effective performance/speed trade-offs and will constitute an interesting heuristic to help future FastGAE users selecting subgraph sizes.

Comparison of Uniform, Core-based and Degree-based FastGAE. As for medium-size graphs, core-based sampling and degree-based sampling is empirically more effective than uniform sampling (e.g. in Table 5, +6.94 AUC points for FastGAE with degree sampling on Patent, with $n_{(S)} = 20000$), and associated to lower standard errors. We observe that computing the p_i probabilities through core-based sampling is longer on large graphs, but bring no empirical benefit w.r.t. degree-based sampling: we therefore recommend using degree-based sampling for large graphs.

FastGAE vs non-AE/VAE baselines, and the case of Node Clustering. For the link prediction task, best FastGAE models usually reach competitive results w.r.t. node2vec while being significantly faster (see e.g. the last column of Table 6). However, regarding node clustering, we observe in Table 4 and in Figure 5 that the *Lowvain* baseline outperforms AE/VAE models on SBM, a phenomenon that we also noted on the Cora and Citeseer graphs (section 4.2.2). We conjecture that current graph AE and VAE models might be suboptimal to effectively reconstruct communities in graph data ; this claim is consistent with recent experiments on these datasets [18, 36]. As our objective, in this paper, was to scale existing graph AE/VAE, but not to ensure nor to claim their superiority over all other methods for node clustering, we do not further investigate this limit here. Nonetheless, future works on more effective cluster reconstruction from AE/VAE embeddings could definitely lead towards the improvement of these models.

4.2.4. On the embedding dimension d

Our tables present results for a fixed embedding dimension of $d = 16$, for *all* models (all variants of AE/VAE and other baselines), even for large graphs. Nonetheless, we reached similar conclusions for $d = 32, 64$ and 128 : although performances sometimes slightly improved by increasing d , the *ranking* of the different models remained unchanged. We also considered optimizing d individually for each model (to cover potential cases where the impact of d on the performance of each model would have been different) but, again, it did not modify the ranking of models in terms of AUC, AP and AMI scores.

4.2.5. On the number of training iterations

As detailed in Section 4.1.3, all graph AE and VAE models, with or without our FastGAE framework, were trained for 200 iterations (resp. 300) for graphs with $n < 100000$ (resp. $n \geq 100000$). We thoroughly checked the convergence of all models, by assessing the stabilization of performances in terms of AUC scores on validation sets. Using a fixed number of iterations is common in recent research on graph AE and VAE [11, 17, 18, 25]. We nonetheless think that *early-stopping* [65] would also be a relevant alternative strategy, that could lead to additional speed-ups, and might deserve further investigations in future works. Besides, we observed that, for very small values of $n_{(S)}$,

| Dataset | Model (Dimension $d = 16$) | Average Perf. on Test Set | | Average Running Times (in sec.) | | |
|--|--------------------------------------|------------------------------------|------------------------------------|---------------------------------|------------------|-----------------------|
| | | AUC (in %) | AP (in %) | Compute p_i | Train model | Total |
| SBM | Standard Graph VAE | <i>(intractable)</i> | | <i>(intractable)</i> | | |
| | Var. FastGAE (degree, $\alpha = 2$) | | | | | |
| | - with $n_{(S)} = 2000$ | 79.37 \pm 0.52 | 80.68 \pm 0.84 | 0.03 | 27.36 | 27.39 |
| | - with $n_{(S)} = n_{(S)}^* = 2673$ | 80.96 \pm 0.35 | 83.69 \pm 0.60 | 0.03 | 30.66 | 30.69 |
| | - with $n_{(S)} = 5000$ | 81.45 \pm 0.39 | 84.30 \pm 0.82 | 0.03 | 43.86 | 43.89 |
| <u>Best baseline</u> node2vec | 80.89 \pm 0.32 | 83.51 \pm 0.29 | - | 1328.82 | 1328.82 (22 min) | |
| Google | Standard Graph AE | <i>(intractable)</i> | | <i>(intractable)</i> | | |
| | FastGAE (degree, $\alpha = 1$) | | | | | |
| | - with $n_{(S)} = 2500$ | 94.52 \pm 0.26 | 95.50 \pm 0.11 | 0.14 | 122.53 | 122.67 |
| | - with $n_{(S)} = n_{(S)}^* = 7911$ | 95.75 \pm 0.24 | 96.62 \pm 0.09 | 0.14 | 158.63 | 158.77 |
| | - with $n_{(S)} = 10000$ | 95.91 \pm 0.19 | 96.64 \pm 0.12 | 0.14 | 168.10 | 168.24 |
| <u>Best baseline</u> node2vec | 94.89 \pm 0.63 | 96.82 \pm 0.72 | - | 14762.78 | 14762.78 (4h06) | |
| Youtube | Standard Graph VAE | <i>(intractable)</i> | | <i>(intractable)</i> | | |
| | Var. FastGAE (degree, $\alpha = 5$) | | | | | |
| | - with $n_{(S)} = 3000$ | 81.14 \pm 0.19 | 86.61 \pm 0.16 | 0.28 | 453.22 | 453.50 (8min) |
| | - with $n_{(S)} = n_{(S)}^* = 15179$ | 81.83 \pm 0.15 | 87.21 \pm 0.15 | 0.28 | 2964.51 | 2964.79 (49min) |
| | - with $n_{(S)} = 20000$ | 82.31 \pm 0.18 | 87.36 \pm 0.15 | 0.28 | 3596.03 | 3596.31 (1h) |
| <u>Best baseline</u> Core-Graph VAE, $k = 40$ | 80.53 \pm 0.23 | 82.45 \pm 0.20 | - | 12433.51 | 12433.51 (3h27) | |
| Patent | Standard Graph AE | <i>(intractable)</i> | | <i>(intractable)</i> | | |
| | FastGAE with (degree, $\alpha = 2$) | | | | | |
| | - with $n_{(S)} = 5000$ | 90.66 \pm 0.25 | 90.76 \pm 0.22 | 0.30 | 605.75 | 606.05 (10min) |
| | - with $n_{(S)} = n_{(S)}^* = 16425$ | 93.02 \pm 0.23 | 93.39 \pm 0.23 | 0.30 | 3693.32 | 3693.62 (1h02) |
| | - with $n_{(S)} = 20000$ | 92.91 \pm 0.22 | 93.35 \pm 0.21 | 0.30 | 4401.67 | 4401.67 (1h13) |
| <u>Best baseline</u> node2vec | 92.96 \pm 0.23 | 93.43 \pm 0.17 | - | 25851.39 | 25851.39 (7h11) | |

Table 6: Summarized results for link prediction on all large graphs. For each graph, for brevity, we only report a few representative degree-based FastGAE or Variational FastGAE models, and the best baseline (among Core-Graph AE/VAE, Negative Sampling Graph AE/VAE and node2vec). Scores are averaged over 10 runs with different and random train/validation/test sets. Standard Graph AE and VAE are intractable. Scores are averaged over 10 runs with different and random train/validation/test sets. For degree sampling, values of the hyperparameter α (as defined in equation 9) were tuned, as described in Figure C.6. All models learn embedding vectors of dimension $d = 16$. Bold numbers correspond to the best performance (several numbers are bold when scores are comparable, in a ± 1 standard deviation range) and best running time.

increasing the number of training iterations did not significantly improved our results: to improve scores on such settings, increasing the sampling size $n_{(S)}$ was overall more effective than increasing the number of training iterations.

5. Conclusion and Discussion

In this paper, we introduced and released a general framework to scale graph AE and VAE models. We demonstrated its effectiveness on large graphs with up to millions of nodes and edges, both in terms of speed, of scalability and of performance. We outperformed the few existing approaches to scale graph AE and VAE, usually by a wide margin. FastGAE is also conceptually simpler than these alternative approaches [18], and we believe that simple solutions often have the most impact.

Futhermore, FastGAE is a flexible framework that easily extends to AE/VAE models with alternative GNN *encoders*. In our experiments, the GCN encoders of standard graph AE/VAE models and of FastGAE-based models could easily be replaced by any alternative architecture learning the embedding matrix Z in another way, e.g. by a FastGCN [30], a Cluster-GCN [31], a GCN with simple graph convolutions (SGC) [38] or

a GraphSAGE [4] model. Besides, FastGAE easily extends to graph AE or VAE with alternative *decoders*. For instance, one could replace the symmetric inner-product decoder from our experiments by the asymmetric decoder recently proposed by Salha et al. [28], which would extend FastGAE to *directed* graphs.

Last, but not least, we also identify possible future research directions for improvements. Apart from the aforementioned limit (section 4.2.3) of current graph AE and VAE models on the node clustering task (that, however, concerns all graph AE/VAE from our experiments and is not specific to FastGAE), we underline that the proposed FastGAE method could underperform on very sparse graphs. Indeed, in such scenario, the subgraphs to reconstruct might include a large proportion of isolated nodes, which would negatively impact learning. Moreover, in the case of large graphs with a lot of sparsely connected components, we recommend applying FastGAE separately on each component. Also, in this paper we always assumed that the graph was fixed, which might sometimes be a limit, that could initiate future interesting studies on extensions of FastGAE for scalable *dynamic* graph embeddings, potentially with a dynamic selection of $n_{(S)}$.

Appendices

This supplementary material provides all proofs from our theoretical analyses as well as an additional figure. It is organized as follows:

- In Appendix A, we report the proofs of Propositions 1, 2 and 3 from our Section 3.3.2 on approximated losses.
- In Appendix B, we report the proofs of Propositions 4 and 5 from our Section 3.3.3 on the computation of the threshold subgraph size $n_{(S)}^*$.
- The figure of Appendix C presents optimal values of the hyperparameter α .

Appendix A. On Approximate Losses

Proposition 1. *Let $\mathcal{G}_{(S)} = (\mathcal{V}_{(S)}, \mathcal{E}_{(S)})$ be a subgraph of \mathcal{G} obtained from sampling $n_{(S)}$ nodes **with** replacement using the node sampling strategy of FastGAE. Let i and j denote two distinct nodes from the original graph \mathcal{G} : $(i, j) \in \mathcal{V}^2$. Then:*

$$\mathbb{P}(i \in \mathcal{V}_{(S)}) = 1 - (1 - p_i)^{n_{(S)}}.$$

Also:

$$\begin{aligned} \mathbb{P}((i, j) \in \mathcal{V}_{(S)}^2) &= 1 - \left[(1 - p_i)^{n_{(S)}} + (1 - p_j)^{n_{(S)}} \right. \\ &\quad \left. - (1 - p_i - p_j)^{n_{(S)}} \right]. \end{aligned}$$

Proof. In this setting, sampling probabilities are independent of previous sampling steps, and remain fixed to p_i . Therefore, for node $i \in \mathcal{V}$, we have:

$$\mathbb{P}(i \notin \mathcal{V}_{(S)}) = (1 - p_i)^{n_{(S)}}.$$

Indeed, for i not to belong to $\mathcal{V}_{(S)}$, it must not be selected at any of the $n_{(S)}$ draws, which happens with probability $1 - p_i$ for each draw. Therefore:

$$\mathbb{P}(i \in \mathcal{V}_{(S)}) = 1 - (1 - p_i)^{n_{(S)}}.$$

Moreover, let i and j denote two distinct nodes from the original graph \mathcal{G} : $(i, j) \in \mathcal{V}^2$. We have:

$$\begin{aligned} \mathbb{P}((i, j) \notin \mathcal{V}_{(S)}^2) &= \mathbb{P}(i \notin \mathcal{V}_{(S)} \text{ or } j \notin \mathcal{V}_{(S)}) \\ &= \mathbb{P}(i \notin \mathcal{V}_{(S)}) + \mathbb{P}(j \notin \mathcal{V}_{(S)}) - \mathbb{P}(i \notin \mathcal{V}_{(S)}, j \notin \mathcal{V}_{(S)}) \end{aligned}$$

with, using the previous result, $\mathbb{P}(i \notin \mathcal{V}_{(S)}) = (1 - p_i)^{n_{(S)}}$ and $\mathbb{P}(j \notin \mathcal{V}_{(S)}) = (1 - p_j)^{n_{(S)}}$. Using a similar argument, we also obtain:

$$\mathbb{P}(i \notin \mathcal{V}_{(S)}, j \notin \mathcal{V}_{(S)}) = \left(1 - (p_i + p_j)\right)^{n_{(S)}}.$$

Therefore:

$$\mathbb{P}\left((i, j) \notin \mathcal{V}_{(S)}^2\right) = \left[(1 - p_i)^{n_{(S)}} + (1 - p_j)^{n_{(S)}} - (1 - p_i - p_j)^{n_{(S)}} \right].$$

And:

$$\begin{aligned} \mathbb{P}\left((i, j) \in \mathcal{V}_{(S)}^2\right) &= 1 - \mathbb{P}\left((i, j) \notin \mathcal{V}_{(S)}^2\right) \\ &= 1 - \left[(1 - p_i)^{n_{(S)}} + (1 - p_j)^{n_{(S)}} - (1 - p_i - p_j)^{n_{(S)}} \right]. \end{aligned}$$

Last, for self-loops:

$$\mathbb{P}\left((i, i) \in \mathcal{V}_{(S)}^2\right) = \mathbb{P}\left(i \in \mathcal{V}_{(S)}\right) = 1 - (1 - p_i)^{n_{(S)}}.$$

□

Proposition 2. Let $\mathcal{G}_{(S)} = (\mathcal{V}_{(S)}, \mathcal{E}_{(S)})$ be a subgraph of \mathcal{G} obtained from sampling $n_{(S)}$ nodes **without** replacement using the node sampling strategy of FastGAE. Let i and j denote two distinct nodes from \mathcal{G} : $(i, j) \in \mathcal{V}^2$. Then:

$$\mathbb{P}\left(i \in \mathcal{V}_{(S)}\right) = \sum_{\mathcal{U} \in \mathbf{U}(i)} p_{u_1} \prod_{k=2}^{n_{(S)}} \frac{p_{u_k}}{1 - \sum_{k'=1}^{k-1} p_{u_{k'}}},$$

where $\mathbf{U}(i) = \{\mathcal{U} \subset \mathcal{V}, |\mathcal{U}| = n_{(S)} \text{ and } i \in \mathcal{U}\}$ is the set of all **ordered** subsets of $n_{(S)}$ distinct nodes including node i . For a given set $\mathcal{U} \in \mathbf{U}(i)$, we denote by $(u_1, u_2, \dots, u_{n_{(S)}})$ its ordered elements. Also,

$$\mathbb{P}\left((i, j) \in \mathcal{V}_{(S)}^2\right) = \sum_{\mathcal{U} \in \mathbf{U}(i) \cap \mathbf{U}(j)} p_{u_1} \prod_{k=2}^{n_{(S)}} \frac{p_{u_k}}{1 - \sum_{k'=1}^{k-1} p_{u_{k'}}}.$$

Proof. We are looking for the probability that a node $i \in \mathcal{V}$ from the graph belongs to a drawn subset $\mathcal{V}_{(S)}$, that contains $n_{(S)}$ distinct nodes. For $\mathcal{V}_{(S)}$ to include i , $\mathcal{V}_{(S)}$ should match any of the possible ordered subsets of $n_{(S)}$ nodes that include node i . In this setting where we sample without replacement, the probability to draw node i depends on nodes previously drawn. All possible orders of sampling the nodes should be considered. Let:

$$\mathbf{U}(i) = \left\{ \mathcal{U} \subset \mathcal{V}, |\mathcal{U}| = n_{(S)} \text{ and } i \in \mathcal{U} \right\}$$

denote the set of all **ordered** subsets of $n_{(S)}$ distinct nodes that include node i . With such notations:

$$\mathbb{P}\left(i \in \mathcal{V}_{(S)}\right) = \mathbb{P}\left(\mathcal{V}_{(S)} \in \mathbf{U}(i)\right) = \sum_{\mathcal{U} \in \mathbf{U}(i)} \mathbb{P}\left(\mathcal{V}_{(S)} = \mathcal{U}\right).$$

The summation comes from the fact that events are *disjoint* ($\mathcal{V}_{(S)}$ can not match two of these ordered subsets simultaneously).

Now, for a given set $\mathcal{U} \in \mathbf{U}(i)$, let us denote by $(u_1, u_2, \dots, u_{n_{(S)}})$ its **ordered** elements. Also, let $(\mathcal{V}_{(S)1}, \mathcal{V}_{(S)2}, \dots, \mathcal{V}_{(S)n_{(S)}})$ be the $n_{(S)}$ ordered nodes of set $\mathcal{G}_{(S)}$ (i.e. $\mathcal{V}_{(S)1}$ is the first drawn node, $\mathcal{V}_{(S)2}$ is the second one, etc). We have:

$$\begin{aligned} \mathbb{P}(\mathcal{V}_{(S)} = \mathcal{U}) &= \mathbb{P}(\mathcal{V}_{(S)1} = u_1, \mathcal{V}_{(S)2} = u_2, \dots, \mathcal{V}_{(S)n_{(S)}} = u_{n_{(S)}}) \\ &= \mathbb{P}(\mathcal{V}_{(S)1} = u_1) \prod_{k=2}^{n_{(S)}} \mathbb{P}(\mathcal{V}_{(S)k} = u_k | \mathcal{V}_{(S)k-1} = u_{k-1}, \dots, \mathcal{V}_{(S)1} = u_1) \\ &= p_{u_1} \prod_{k=2}^{n_{(S)}} \frac{p_{u_k}}{1 - \sum_{k'=1}^{k-1} p_{u_{k'}}}. \end{aligned}$$

Therefore, by summing elements to come back to $\mathbb{P}(i \in \mathcal{V}_{(S)})$:

$$\mathbb{P}(i \in \mathcal{V}_{(S)}) = \sum_{\mathcal{U} \in \mathbf{U}(i)} p_{u_1} \prod_{k=2}^{n_{(S)}} \frac{p_{u_k}}{1 - \sum_{k'=1}^{k-1} p_{u_{k'}}}.$$

Moreover, let i and j denote two distinct nodes from the original graph \mathcal{G} : $(i, j) \in \mathcal{V}^2$. Using similar notations and reasoning, we get:

$$\begin{aligned} \mathbb{P}((i, j) \in \mathcal{V}_{(S)}^2) &= \mathbb{P}(i \in \mathcal{V}_{(S)}, j \in \mathcal{V}_{(S)}) \\ &= \sum_{\mathcal{U} \in \mathbf{U}(i) \cap \mathbf{U}(j)} \mathbb{P}(\mathcal{V}_{(S)} = \mathcal{U}). \end{aligned}$$

Therefore:

$$\mathbb{P}((i, j) \in \mathcal{V}_{(S)}^2) = \sum_{\mathcal{U} \in \mathbf{U}(i) \cap \mathbf{U}(j)} p_{u_1} \prod_{k=2}^{n_{(S)}} \frac{p_{u_k}}{1 - \sum_{k'=1}^{k-1} p_{u_{k'}}}.$$

And, for self-loops, $\mathbb{P}((i, i) \in \mathcal{V}_{(S)}^2) = \mathbb{P}(i \in \mathcal{V}_{(S)})$. \square

Proposition 3. *Using the expressions of Proposition 1 (with replacement) or Proposition 2 (without replacement):*

$$\mathbb{E}[\mathcal{L}^{\text{FastGAE}}] = \frac{1}{n_{(S)}^2} \sum_{(i,j) \in \mathcal{V}^2} \mathbb{P}((i, j) \in \mathcal{V}_{(S)}^2) \mathcal{L}_{ij}(A_{ij}, \hat{A}_{ij}).$$

Proof. We have:

$$\begin{aligned} \mathbb{E}[\mathcal{L}^{\text{FastGAE}}] &= \mathbb{E}\left[\frac{1}{n_{(S)}^2} \sum_{(i,j) \in \mathcal{V}^2} \mathbf{1}_{((i,j) \in \mathcal{V}_{(S)}^2)} \mathcal{L}_{ij}(A_{ij}, \hat{A}_{ij})\right] \\ &= \frac{1}{n_{(S)}^2} \sum_{(i,j) \in \mathcal{V}^2} \mathbb{E}\left[\mathbf{1}_{((i,j) \in \mathcal{V}_{(S)}^2)}\right] \mathcal{L}_{ij}(A_{ij}, \hat{A}_{ij}) \\ &= \frac{1}{n_{(S)}^2} \sum_{(i,j) \in \mathcal{V}^2} \mathbb{P}((i, j) \in \mathcal{V}_{(S)}^2) \mathcal{L}_{ij}(A_{ij}, \hat{A}_{ij}). \end{aligned}$$

By replacing $\mathbb{P}((i, j) \in \mathcal{V}_{(S)}^2)$ by the expressions of Proposition 1 (with replacement) or Proposition 2 (without replacement), we obtain an explicit formulation for $\mathbb{E}[\mathcal{L}^{\text{FastGAE}}]$. \square

Appendix B. On the Selection of $n_{(S)}$

Proposition 4. *Let us consider a training iteration of the FastGAE framework, a sampled subgraph $\mathcal{G}_{(S)} = (\mathcal{V}_{(S)}, \mathcal{E}_{(S)})$, with $|\mathcal{V}_{(S)}| = n_{(S)} < n$ nodes sampled without replacement, and the corresponding node-level approximate reconstruction computed for a given node i :*

$$\mathcal{L}^{\text{FastGAE}}(i) = \frac{1}{n_{(S)}} \sum_{j \in \mathcal{V}} \mathbb{1}_{(j \in \mathcal{V}_{(S)})} \mathcal{L}_{ij}(A_{ij}, \hat{A}_{ij}),$$

with the random variable $\mathbb{1}_{(j \in \mathcal{V}_{(S)})} = 1$ if node $j \in \mathcal{V}_{(S)}$ and 0 otherwise, with $A_{ij} \in \{0, 1\}$ for all $(i, j) \in \mathcal{V}^2$ and with:

$$\mathcal{L}_{ij}(A_{ij}, \hat{A}_{ij}) = -[A_{ij} \log(\hat{A}_{ij}) + (1 - A_{ij}) \log(1 - \hat{A}_{ij})].$$

Then, under Assumption 1 from Section 3.3.3, for any $\gamma \geq 0$, we have:

$$\mathbb{P}(|\mathcal{L}^{\text{FastGAE}}(i) - \mathbb{E}[\mathcal{L}^{\text{FastGAE}}(i)]| \geq \gamma) \leq 2 \exp\left(-2\left(\frac{\gamma}{\log(\varepsilon)}\right)^2 \frac{n_{(S)}^2}{n}\right).$$

We note that the right hand side term tends to 0 exponentially fast w.r.t. the deviation magnitude γ and w.r.t. the subgraph size $n_{(S)}$.

Proof. As a preliminary, let us recall Hoeffding's inequality [50]. Let X_1, X_2, \dots, X_n be real independent random variables verifying, for some $(a_k)_{1 \leq k \leq n}$ and $(b_k)_{1 \leq k \leq n}$ with $a_k < b_k$: $\forall k, \mathbb{P}(a_k \leq X_k \leq b_k) = 1$. Let $S_n = \sum_{i=1}^n X_i$. Then, for all $\gamma > 0$, Hoeffding's inequality states that:

$$\mathbb{P}(|S_n - \mathbb{E}(S_n)| \geq t) \leq 2 \exp\left(-\frac{2\gamma^2}{\sum_{i=1}^n (b_i - a_i)^2}\right).$$

[50] also proves that the above inequality holds when the X_i are samples without replacement from a finite population (and therefore not independent). In the setting of Proposition 4, that falls into this second case due to the node-level sampling scheme of FastGAE, we have:

$$\mathcal{L}^{\text{FastGAE}}(i) = \sum_{j \in \mathcal{V}} X_{ij},$$

where, under Assumption 1:

$$\begin{aligned} X_{ij} &= \frac{1}{n_{(S)}} \mathbb{1}_{(j \in \mathcal{V}_{(S)})} \mathcal{L}_{ij}(A_{ij}, \hat{A}_{ij}) \\ &= \underbrace{\mathbb{1}_{(j \in \mathcal{V}_{(S)})}}_{\in \{0, 1\}} \underbrace{\frac{-1}{n_{(S)}} [A_{ij} \log(\hat{A}_{ij}) + (1 - A_{ij}) \log(1 - \hat{A}_{ij})]}_{\in [-\log(1-\varepsilon)/n_{(S)}, -\log(\varepsilon)/n_{(S)}]} \\ &\in \left[0, \frac{-\log(\varepsilon)}{n_{(S)}}\right]. \end{aligned}$$

We note that $\frac{-\log(\varepsilon)}{n_{(S)}} > 0$, as $0 < \varepsilon < 1$. Applying Hoeffding's inequality, at each sampling step and for all $\gamma > 0$:

$$\begin{aligned} \mathbb{P}(|\mathcal{L}^{\text{FastGAE}}(i) - \mathbb{E}[\mathcal{L}^{\text{FastGAE}}(i)]| \geq \gamma) &\leq 2 \exp\left(-\frac{2\gamma^2}{\sum_{j \in \mathcal{V}} \left(\frac{-\log(\varepsilon)}{n_{(S)}}\right)^2}\right) \\ &= 2 \exp\left(\frac{-2\gamma^2}{n \frac{(-\log(\varepsilon))^2}{n_{(S)}^2}}\right) \\ &= 2 \exp\left(-2\left(\frac{\gamma}{\log(\varepsilon)}\right)^2 \frac{n_{(S)}^2}{n}\right) \end{aligned}$$

We note that it exhibits the link between the deviation of the loss and the $\frac{n_{(S)}^2}{n}$ ratio. \square

Proposition 5. For any confidence level $\alpha \in]0, 1[$ and node $i \in \mathcal{V}$, selecting a subgraph size $n_{(S)}$ such that

$$n_{(S)} \geq n_{(S)}^* = \sqrt{n} \underbrace{\sqrt{\frac{-\log(\frac{\alpha}{2}) \log(\varepsilon)^2}{2\gamma^2}}}_{\text{denoted } C \text{ in eq. 10}} \quad (\text{B.1})$$

guarantees that

$$\mathbb{P}(|\mathcal{L}^{\text{FastGAE}}(i) - \mathbb{E}[\mathcal{L}^{\text{FastGAE}}(i)]| \geq \gamma) \leq \alpha.$$

Proof. This is a corollary of Proposition 4, from which we derive that, for any $\alpha \in]0, 1[$:

$$\begin{aligned} 2 \exp\left(-2\left(\frac{\gamma}{\log(\varepsilon)}\right)^2 \frac{n_{(S)}^2}{n}\right) &\leq \alpha \\ \Rightarrow \mathbb{P}(|\mathcal{L}^{\text{FastGAE}} - \mathbb{E}[\mathcal{L}^{\text{FastGAE}}]| \geq \gamma) &\leq \alpha. \end{aligned}$$

Then:

$$\begin{aligned} 2 \exp\left(-2\left(\frac{\gamma}{\log(\varepsilon)}\right)^2 \frac{n_{(S)}^2}{n}\right) &\leq \alpha \\ \Leftrightarrow -2\left(\frac{\gamma}{\log(\varepsilon)}\right)^2 \frac{n_{(S)}^2}{n} &\leq \log\left(\frac{\alpha}{2}\right) \\ \Leftrightarrow n_{(S)} &\geq \sqrt{n} \sqrt{\frac{-\log(\frac{\alpha}{2}) \log(\varepsilon)^2}{2\gamma^2}} \end{aligned}$$

\square

Appendix C. On the hyperparameter α

In this last appendix, we report the additional Figure C.6, presenting the optimal values of the hyperparameter α , for all graphs, and for both core-based and degree-based sampling.

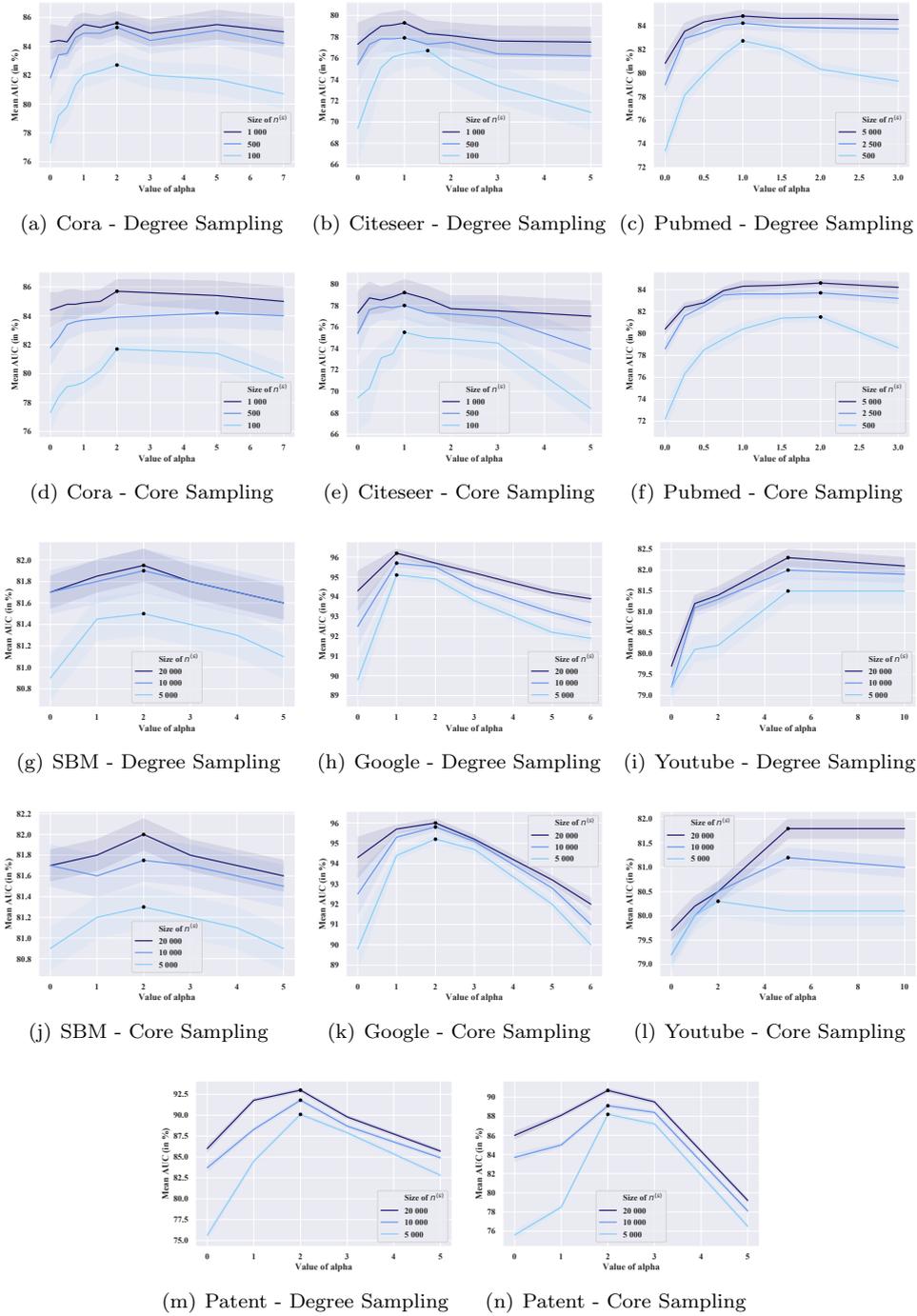


Figure C.6: Optimal values of hyperparameter α for degree-based and core-based node sampling w.r.t. mean AUC scores on validation sets, for Variational FastGAE models and for all graphs.

References

- [1] W. L. Hamilton, R. Ying, J. Leskovec, Representation learning on graphs: Methods and applications, *IEEE Data Engineering Bulletin* (2017).
- [2] Z. Wu, S. Pan, F. Chen, G. Long, C. Zhang, P. S. Yu, A comprehensive survey on graph neural networks, *arXiv preprint arXiv:1901.00596* (2019).
- [3] T. N. Kipf, M. Welling, Semi-supervised classification with graph convolutional networks, *International Conference on Learning Representations* (2017).
- [4] W. L. Hamilton, Z. Ying, J. Leskovec, Inductive representation learning on large graphs, *Advances in Neural Information Processing Systems* (2017).
- [5] S. Cao, W. Lu, Q. Xu, Grarep: Learning graph representations with global structural information, *ACM International Conference on Information and Knowledge Management* (2015).
- [6] B. Perozzi, R. Al-Rfou, S. Skiena, Deepwalk: Online learning of social representations, *ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* (2014).
- [7] J. Tang, M. Qu, M. Wang, M. Zhang, J. Yan, Q. Mei, Line: Large-scale information network embedding, *International Conference on World Wide Web* (2015).
- [8] A. Grover, J. Leskovec, node2vec: Scalable feature learning for networks, *ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* (2016).
- [9] F. Tian, B. Gao, Q. Cui, E. Chen, T.-Y. Liu, Learning deep representations for graph clustering, *AAAI Conference on Artificial Intelligence* (2014).
- [10] D. Wang, P. Cui, W. Zhu, Structural deep network embedding, *ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* (2016).
- [11] T. N. Kipf, M. Welling, Variational graph auto-encoders, *NeurIPS Workshop on Bayesian Deep Learning* (2016).
- [12] D. E. Rumelhart, G. E. Hinton, R. J. Williams, Learning internal representations by error propagation, *Parallel Distributed Processing, Vol 1* (1986).
- [13] P. Baldi, Autoencoders, unsupervised learning, and deep architectures, *ICML Workshop on Unsupervised and Transfer Learning* (2012).
- [14] D. P. Kingma, M. Welling, Auto-encoding variational bayes, *International Conference on Learning Representations* (2014).
- [15] C. Wang, S. Pan, G. Long, X. Zhu, J. Jiang, Mgae: Marginalized graph autoencoder for graph clustering, *ACM Conference on Information and Knowledge Management* (2017).
- [16] S. Pan, R. Hu, G. Long, J. Jiang, L. Yao, C. Zhang, Adversarially regularized graph autoencoder for graph embedding, *International Joint Conference on Artificial Intelligence* (2018).
- [17] A. Hasanzadeh, E. Hajiramezani, K. Narayanan, N. Duffield, M. Zhou, X. Qian, Semi-implicit graph variational auto-encoders, *Advances in Neural Information Processing Systems* (2019).
- [18] G. Salha, R. Hennequin, V. A. Tran, M. Vazirgiannis, A degeneracy framework for scalable graph autoencoders, *International Joint Conference on Artificial Intelligence* (2019).
- [19] H. Shi, H. Fan, J. T. Kwok, Effective decoding in graph auto-encoder using triadic closure, *AAAI Conference on Artificial Intelligence* (2020).
- [20] W. Jin, R. Barzilay, T. Jaakkola, Junction tree variational autoencoder for molecular graph generation, *International Conference on Machine Learning* (2018).
- [21] Q. Liu, M. Allamanis, M. Brockschmidt, A. Gaunt, Constrained graph variational autoencoders for molecule design, *Advances in Neural Information Processing Systems* (2018).
- [22] T. Ma, J. Chen, C. Xiao, Constrained generation of semantically valid graphs via regularizing variational autoencoders, *Advances in Neural Information Processing Systems* (2018).
- [23] M. Simonovsky, N. Komodakis, Graphvae: Towards generation of small graphs using variational autoencoders, *International Conference on Artificial Neural Networks* (2018).
- [24] B. Samanta, D. Abir, G. Jana, P. K. Chattaraj, N. Ganguly, M. G. Rodriguez, Nevae: A deep generative model for molecular graphs, *AAAI Conference on Artificial Intelligence* (2019).
- [25] R. v. d. Berg, T. N. Kipf, M. Welling, Graph convolutional matrix completion, *KDD Deep Learning Day* (2018).
- [26] P. V. Tran, Multi-task graph autoencoders, *arXiv preprint arXiv:1811.02798* (2018).
- [27] A. Grover, A. Zweig, S. Ermon, Graphite: Iterative generative modeling of graphs, *International Conference on Machine Learning* (2019).
- [28] G. Salha, S. Limnios, R. Hennequin, V. A. Tran, M. Vazirgiannis, Gravity-inspired graph autoencoders for directed link prediction, *ACM International Conference on Information and Knowledge Management* (2019).

- [29] J. Chen, J. Zhu, L. Song, Stochastic training of graph convolutional networks with variance reduction, International Conference on Machine Learning (2018).
- [30] J. Chen, T. Ma, C. Xiao, Fastgcn: fast learning with graph convolutional networks via importance sampling, International Conference on Learning Representations (2018).
- [31] W.-L. Chiang, X. Liu, S. Si, Y. Li, S. Bengio, C.-J. Hsieh, Cluster-gcn: An efficient algorithm for training deep and large graph convolutional networks, ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (2019).
- [32] J. Bruna, W. Zaremba, A. Szlam, Y. LeCun, Spectral networks and locally connected networks on graphs, International Conference on Learning Representations (2014).
- [33] M. Defferrard, X. Bresson, P. Vandergheynst, Convolutional neural networks on graphs with fast localized spectral filtering, Advances in Neural Information Processing Systems (2016).
- [34] A. Micheli, Neural network for graphs: A contextual constructive approach, IEEE Transactions on Neural Networks 20 (2009) 498–511.
- [35] I. Goodfellow, Y. Bengio, A. Courville, Deep learning, MIT press, 2016.
- [36] G. Salha, R. Hennequin, M. Vazirgiannis, Simple and effective graph autoencoders with one-hop linear models, arXiv preprint arXiv:2001.07614 (2020).
- [37] S. Kullback, R. A. Leibler, On information and sufficiency, The Annals of Mathematical Statistics 22-1 (1951) 79–86.
- [38] F. Wu, A. Souza, T. Zhang, C. Fifty, T. Yu, K. Weinberger, Simplifying graph convolutional networks, International Conference on Machine Learning (2019).
- [39] R. Ying, R. He, K. Chen, P. Eksombatchai, W. L. Hamilton, J. Leskovec, Graph convolutional neural networks for web-scale recommender systems, ACM SIGKDD International Conference on Knowledge Discovery & Data Mining (2018).
- [40] H. Zeng, H. Zhou, A. Srivastava, R. Kannan, V. Prasanna, Graphsaint: Graph sampling based inductive learning method, International Conference on Learning Representations (2020).
- [41] F. D. Malliaros, C. Giatsidis, A. N. Papadopoulos, M. Vazirgiannis, The core decomposition of networks: Theory, algorithms and applications, The VLDB Journal (2019) 1–32.
- [42] T. Mikolov, I. Sutskever, K. Chen, G. S. Corrado, J. Dean, Distributed representations of words and phrases and their compositionality, Advances in Neural Information Processing Systems (2013).
- [43] Y. Goldberg, O. Levy, word2vec explained: deriving mikolov et al.’s negative-sampling word-embedding method, in: arXiv preprint arXiv:1402.3722, 2014.
- [44] D. Kempe, J. Kleinberg, E. Tardos, Maximizing the spread of influence through a social network, in: Proceedings of the ninth ACM SIGKDD international conference on Knowledge discovery and data mining, 2003, pp. 137–146.
- [45] J. Leskovec, C. Faloutsos, Sampling from large graphs, in: Proceedings of the 12th ACM SIGKDD international conference on Knowledge discovery and data mining, 2006, pp. 631–636.
- [46] M. Newman, Networks: An Introduction., UK Oxford University Press, 2010.
- [47] V. Batagelj, M. Zaversnik, An $o(m)$ algorithm for cores decomposition of networks, arXiv preprint cs/0310049 (2003).
- [48] P. Hu, W. C. Lau, A survey and taxonomy of graph sampling, arXiv preprint arXiv:1308.5865 (2013).
- [49] F. Chiericetti, A. Dasgupta, R. Kumar, S. Lattanzi, T. Sarlos, On sampling nodes in a network, in: Proceedings of the 25th International Conference on World Wide Web, 2016, pp. 471–481.
- [50] W. Hoeffding, Probability inequalities for sums of bounded random variables, Journal of the American Statistical Association 58 (1963) 13–30.
- [51] L. Tonnaer, Active learning in vae latent space, Eindhoven University of Technology (2017).
- [52] V. Kaushal, A. Sahoo, K. Doctor, N. Raju, S. Shetty, P. Singh, R. Iyer, G. Ramakrishnan, Learning from less data: Diversified subset selection and active learning in image classification tasks, arXiv preprint arXiv:1805.11191 (2018).
- [53] S. Gonzalez, R. Miikkulainen, Improved training speed, accuracy, and data utilization through loss function optimization, arXiv preprint arXiv:1905.11528 (2019).
- [54] P.-Y. Huang, R. Frederking, et al., Rwr-gae: Random walk regularization for graph auto encoders, arXiv preprint arXiv:1908.04003 (2019).
- [55] E. Abbe, Community detection and stochastic block models: recent developments, The Journal of Machine Learning Research 18 (2017) 6446–6531.
- [56] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, et al., Scikit-learn: Machine learning in python, Journal of Machine Learning Research 12 (2011) 2825–2830.
- [57] D. P. Kingma, J. Ba, Adam: A method for stochastic optimization, International Conference on

- Learning Representations (2015).
- [58] M. Abadi, P. Barham, J. Chen, Z. Chen, A. Davis, J. Dean, M. Devin, S. Ghemawat, G. Irving, M. Isard, et al., Tensorflow: A system for large-scale machine learning, in: 12th {USENIX} symposium on operating systems design and implementation ({OSDI} 16), 2016, pp. 265–283.
 - [59] M. Fey, J. E. Lenssen, Fast graph representation learning with PyTorch Geometric, ICLR Workshop on Representation Learning on Graphs and Manifolds (2019).
 - [60] U. Von Luxburg, A tutorial on spectral clustering, *Statistics and computing* 17 (2007) 395–416.
 - [61] V. D. Blondel, J.-L. Guillaume, R. Lambiotte, E. Lefebvre, Fast unfolding of communities in large networks, *Journal of Statistical Mechanics: Theory and Experiments* 2008 (2008) P10008.
 - [62] Y. Rong, W. Huang, T. Xu, J. Huang, Dropedge: Towards deep graph convolutional networks on node classification, *International Conference on Learning Representations* (2020).
 - [63] W. Hu, M. Fey, M. Zitnik, Y. Dong, H. Ren, B. Liu, M. Catasta, J. Leskovec, Open graph benchmark: Datasets for machine learning on graphs, *arXiv preprint arXiv:2005.00687* (2020).
 - [64] B. Kleinberg, Y. Li, Y. Yuan, An alternative view: When does sgd escape local minima?, *International Conference on Machine Learning* (2018).
 - [65] F. Girosi, M. Jones, T. Poggio, Regularization theory and neural networks architectures, *Neural computation* 7 (1995) 219–269.