

# Graph Compression with Application to Model Selection

Mojtaba Abolfazli, Anders Høst-Madsen, June Zhang, Andras Bratincsak

**Abstract**—Many multivariate data such as social and biological data exhibit complex dependencies that are best characterized by graphs. Unlike sequential data, graphs are, in general, unordered structures. This means we can no longer use classic, sequential-based compression methods on these graph-based data. Therefore, it is necessary to develop new methods for graph compression. In this paper, we present universal source coding methods for the lossless compression of unweighted, undirected, unlabelled graphs. We encode in two steps: 1) transforming graph into a rooted binary tree, 2) the encoding rooted binary tree using graph statistics. Our coders showed better compression performance than other source coding methods on both synthetic and real-world graphs.

We then applied our graph coding methods for model selection of Gaussian graphical models using minimum description length (MDL) principle finding the description length of the conditional independence graph. Experiments on synthetic data show that our approach gives better performance compared to common model selection methods. We also applied our approach to electrocardiogram (ECG) data in order to explore the differences between graph models of two groups of subjects.

## I. INTRODUCTION

Graphs can be used to represent complex dependencies between many variables. Graph-based analysis are used in many disciplines such as social networks analysis, natural language processing, chemistry, and bioinformatics. Real-world graphs can be large and expensive to store and communicate. Developing efficient methods to compress graphical data are of interest for storage and transmission of data. Classic information theory methods for coding deal with sequential data. Unlike sequential data with starting and stopping points (ordering), graphs are unordered structures. Since ordering is essential in coding, graph coding is challenging. This means that graphs cannot be efficiently compressed using traditional, sequence-based compression methods.

In this paper, we first present new universal source coding methods for the lossless compression of unlabeled, unweighted, undirected, simple graphs. Since the graphs are unlabeled, vertices do not have distinct identifications except through their interconnections. Reference [1] referred to such a graph as a graph *structure*. Our approach has two steps. First, inspired by Steinruecken’s method for coding of unordered i.i.d sequences [2], we transform a graph,  $G$ , into an equivalent

rooted binary tree,  $T$ . Then, we use graph statistics from  $G$  to develop two classes of coders for  $T$ . The first class utilizes local properties (i.e., formation of graph motifs) such as triangles and the second class uses degree distribution, as a global graph statistics, along with local properties for encoding. This step uses past information (encoded vertices) to encode the connections at the current vertex. This way we can build a probabilistic model based on graph statistics (either by learning or estimating statistics) to encode graphs. Our coders reflect more information about the structure of  $G$  and therefore give shorter codelengths compared to structural coding of [1] for graphs with more structure than Erdős-Rényi graphs.

In the second half of the paper, we use graph coding for data analysis by using description length, which is the number of bits to describe data. Rissanen proposed using description length for model selection in his work on the minimum description length (MDL) [3], [4], [5]. Rissanen’s MDL principle codifies model selection, which balances between how well a model describes the given data and the complexity of the model. Here, we develop a novel approach for model selection of Gaussian graphical models. We compute the description length of the conditional independence graph (the sparsity pattern of the precision matrix) using our lossless graph coding methods and the data under the assumed Gaussian distribution. The model that minimizes the summation of these two terms will be selected as the best model. Unlike other methods that may only consider the number of edges in conditional independence graph to account for model complexity, our approach considers the whole structure of the conditional independence graph by lossless compression.

We showed using synthetic and real-world data the advantages of our methods for both compression and Gaussian graphical model selection. Our approach outperforms competing source coding methods in two different scenarios: 1) compression of a single real-world graph, 2) compression of a graph from a particular type after learning the statistics of its type through training. We also compared our approach with common methods in the literature for model selection in Gaussian graphical models. The experimental results on synthetic data showed that our approach can recover true graph model of data with higher F1-score. We also considered a real-world dataset containing extracted features from 12-lead electrocardiogram (ECG) signals of a group of healthy people and a group of people with Kawasaki disease. We observed that there is a difference between graph model of healthy people and those with Kawasaki disease.

The paper is organized as follows. First, we provide an overview of previous works on universal compression methods

M. Abolfazli, A. Høst-Madsen, and J. Zhang are with the Department of Electrical and Computer Engineering, University of Hawaii at Manoa, 2540 Dole Street, Honolulu, HI 96822; e-mail: {mojtaba,ahm,zjz}@hawaii.edu.

A. Bratincsak is with the Department of Pediatrics, John A. Burns School of Medicine, University of Hawaii, Honolulu, HI 96813; e-mail: andrasb@hphmg.org.

The research was funded in part by the NSF grant CCF-1908957.

for graphs in Section I-A. In Section II, we describe our approach for transforming a graph structure into a rooted binary tree and show that the compression of rooted binary tree is same as the compression of graph structure. Then, we introduce two classes of universal coders based on graph statistics and provide experimental results on synthetic and real-world graphs. In Section III, we present the application of graph coding for model selection in Gaussian graphical models. We also give the performance results on synthetic data and then apply our approach to ECG data. We make concluding remarks and give some directions for future work in Section IV.

### A. Prior Work

Graph compression is a relatively new area in source coding. References [6], [7], [8], [9] focused on the entropy analysis of graph compression. Other papers have provided practical graph compression algorithms. These algorithms can be designed for the compression of either unlabeled graph (i.e., graph structure) or labeled graph (i.e., encoding vertices labels together with graph structure). In the case of unlabeled graph, the decoder recovers a graph that is isomorphic to the input graph. In the case of labeled graph, the decoder recovers the exact graph (i.e., graph structure with labels) at the expense of longer codewords. In other words, encoding of unlabeled graphs benefits from isomorphism since there are different labeled graphs that have the same structure. Fewer bits are required, in general, to encode an unlabeled graph compared to a labeled graph with the same structure [1].

Reference [1] was the first study to develop a compression algorithm to reach the entropy of the distribution on unlabeled Erdős-Rényi graphs up to the first two terms. In [10], two universal coding methods for arbitrary unlabeled graphs based on degree distribution and formation of triangles were introduced. The authors in [11] presented asymptotically optimal structural compression algorithms for the compression of both unlabeled and labeled preferential attachment graphs. The compression of dynamic graphs generated by duplication model was studied in [12]. The authors developed algorithms for the compression of both unlabeled and labeled versions of such graphs. Reference [13] introduced an algorithm for the compression of deep feedforward neural networks by modeling them as bipartite graph layers. An algorithm for the compression of sparse labeled graphs with marks on vertices and edges was introduced in [14]. A general survey on lossless graph compression methods can be found in [15].

Previous compression methods are tailored to specific graph models and can perform poorly on other graph models or real-world graphs. Our approach provides a method to encode arbitrary unlabeled graph structures by building a probabilistic model based on graph properties (e.g., graph motifs and degree distribution). This general approach extracts more information from the graph structure than prior work and results in shorter codelength compared to other source coding methods.

The literature on model selection methods for Gaussian graphical models is rich. Existing methods can broadly fall into two main classes: information-based methods and resampling

methods. Information-based methods, such as Bayesian information criteria (BIC) [16], Akaike information criteria (AIC) [17], and extended Bayesian information criteria (EBIC) [18], select the best model based on the log-likelihood of the data and the complexity of Gaussian graphical model. To account for model complexity, these methods typically consider the number of edges from the conditional independence graph,  $G$ , induced by the Gaussian graphical model; while easy to obtain, this statistic gives only a rough approximation of the structure of  $G$ . In contrast, our approach for model selection is to choose the best model based on the description length of the  $G$  and the data when encoded with the resulting conditional independence graph [19].

Resampling methods measure the performance on out-of-sample data by splitting the data into a subset of samples for fitting and use the remaining samples to estimate the efficacy of the model. The most common method in this class is cross-validation (CV) [20]. Other methods are Generalized Approximate Cross-validation (GACV) [21], Rotation Information Criterion (RIC) [22], and Stability Approach to Regularization Selection (StARS) [23]. The major shortcoming of resampling methods is their high computational cost since they require to solve the problem across all sets of subsamples.

A key aspect of model selection in Gaussian graphical models is the structure of the conditional independence graph. Resampling methods overlook this aspect while information-based methods only consider simple graph statistics such as the number of edges in the conditional independence graph. Our approach presents a different perspective to account for a more accurate model complexity. This relies on being able to compute the description length of the conditional independence graph,  $G$ , using our lossless graph compression methods.

## II. GRAPH COMPRESSION

Consider an unweighted, undirected, simple, unlabeled graph  $G(V, E)$ , where  $V$  is the set of vertices and  $E$  is the set of edges. In this paper, we present a method to compress a graph structure inspired by Steinruecken's method for coding unordered i.i.d binary sequences [2]. We use the terms compression and coding interchangeably in this paper as they denote the same process. One can see our methods as extension of [1] to take into account more graph structure. The general idea is to randomly pick a vertex,  $V_i \in V$ , and encode information about all the incident edges of  $V_i$ . A neighbor of  $V_i$  is then picked and the information about its incident edges are encoded. To implement this scheme, an unlabeled graph  $G$  is transformed into a rooted binary tree  $T$ . Each level of  $T$  associated with the vertex  $V_i$  from  $G$  to be encoded. An encoder is designed to encode  $T$ . A decoder will decode  $T$  without error to recover  $\tilde{G}(V, E)$ , which is identical to  $G(V, E)$  up to an automorphism of the nodes. In this paper, we introduce two broad classes of graph coders to efficiently encode the rooted binary tree  $T$ .

### A. Definitions and Notation

Let  $V_i$  denote the  $i$ th vertex in  $G$ . We will exclusively use *vertices* to refer to the nodes in  $G$  and *nodes* to refer to the

nodes in  $T$ . The total number of vertices in  $G$  is  $|V|$ . Let  $V_i \leftrightarrow V_j$  denote an edge between vertices  $V_i$  and  $V_j$  in  $G$ . The degree of vertex  $V_i$  is the total number of vertices connected to  $V_i$ . The degree distribution  $P(k)$  is the probability distribution of degrees of vertices in  $G$  and is an often used statistics to differentiate between different classes of random graphs. The adjacency matrix of  $G$  is shown by  $A = [A_{ij}]$ , a  $|V| \times |V|$  matrix where  $A_{ij} = 1$  if there is an edge between  $V_i$  and  $V_j$  in  $G$ .

The structure of the rooted binary tree  $T$  is very important in the encoding process. We will divide the organization of  $T$  into different levels (also known as depth). Let  $[\ell, i]$  denote the  $i$ th node of the  $\ell$ -th level of  $T$ . Note that the root node is  $[0, 1]$ . Figure 2 shows our naming convention for all the other nodes. A node in  $T$  will can contain multiple vertices from  $G$ . The cardinality of the node  $[\ell, i]$  is shown with  $||[\ell, i]|$ . By default,  $||[0, 1]| = |V|$ .

**Definition 1.** A node  $[\ell, i]$  is a *left node* if it is the left child of a node in level  $\ell - 1$ . In our convention, left nodes have odd index value  $i$ .

**Definition 2.** A node  $[\ell, i]$  is a *right node* if it is the right child of a node in level  $\ell - 1$ . In our convention, right nodes have even index value  $i$ .

**Definition 3.** For a level  $\ell$ , the *first nonempty node* refers to node  $[\ell, \alpha]$ :

$$[\ell, \alpha] = \begin{cases} [\ell, 1], & \text{if } ||[\ell, 1]| > 0 \\ [\ell, 2], & \text{otherwise.} \end{cases} \quad (1)$$

Often times, we may wish to explicitly refer to the parent or children node of  $[\ell, i]$ . We use  $[\ell, i].\text{parent}$  to refer to the parent node of  $[\ell, i]$ . We use  $[\ell, i].\text{left}$  and  $[\ell, i].\text{right}$  to refer to left and right child of node  $[\ell, i]$ . In the example tree shown in Figure 2, node  $[1, 1] = [2, 1].\text{parent}$ , node  $[2, 3] = [1, 2].\text{left}$ , and node  $[2, 4] = [1, 2].\text{right}$ .

Let  $\mathcal{R}([\ell, i])$  denote the path from node  $[\ell, i]$  to the root node. We find levels where  $\mathcal{R}([\ell, i])$  and  $\mathcal{R}([\ell, \alpha])$  both contain left nodes and store those levels in  $\mathcal{CI}([\ell, i])$ . We also find levels where  $\mathcal{R}([\ell, i])$  or  $\mathcal{R}([\ell, \alpha])$  contain left nodes and store those levels in  $\mathcal{I}([\ell, i])$ .

For example in Figure 2,  $\mathcal{R}([2, 3]) = \{[2, 3], [1, 2], [0, 1]\}$  and  $\mathcal{R}([2, 1]) = \{[2, 1], [1, 1], [0, 1]\}$ . We can see that  $\mathcal{R}([2, 3])$  and  $\mathcal{R}([2, 1])$  (node  $[2, 1]$  refer to  $[2, \alpha]$ ) both have left nodes at level  $\ell = 2$  (i.e., nodes  $[2, 3]$  and  $[2, 1]$ , respectively) and therefore, we get  $\mathcal{CI}([2, 3]) = \{2\}$ . We can also verify that  $\mathcal{I}([2, 3]) = \{1, 2\}$  as node  $[1, 1]$  is another left node in  $\mathcal{R}([2, 1])$ .

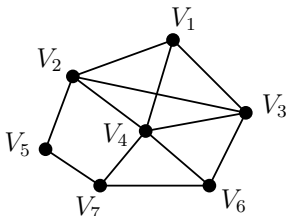


Figure 1: An example graph  $G$  with  $|V| = 7$ .

## B. Transforming Graph $G$ into Binary Tree $T$

The rooted binary tree  $T$  is built by traversing the unlabeled graph  $G(V, E)$ . Each node in  $T$  corresponds to a set of vertices from  $G$ . The root node  $[0, 1]$  corresponds to the set of all the vertices,  $V$ . To build levels  $\ell = 1, \dots, |V| - 1$  of the tree we pick a vertex  $V_{k(\ell)}$ ; we will specify shortly how it is picked. Vertex  $V_{k(\ell)}$  is removed from the set of vertices. The remaining vertices belonging to each node at level  $\ell - 1$  is split into two groups: the left node are all vertices connected to  $V_{k(\ell)}$ , and the right node contain those not connected to  $V_{k(\ell)}$ . Empty nodes are not split further. It is worth noting that the vertex  $V_{k(\ell)}$  is picked randomly from the first non-empty node at level  $\ell - 1$ , i.e., node  $[\ell - 1, \alpha]$ . Figure 1 and Figure 2 show an example of  $G(V, E)$  and its corresponding rooted binary tree representation,  $T$ , respectively. Since the order of nodes is irrelevant for the structure, here we have assumed  $V_{k(\ell)} = V_\ell$ .

Algorithm 1 gives the pseudocode for transforming graph  $G$  into the rooted binary tree  $T$ . Note that hereafter, when we traverse a level in the rooted binary tree  $T$ , we start from the most left node and finish with the most right node at that level.

---

### Algorithm 1 Transform graph $G$ into rooted binary tree $T$

---

```

1: function GRAPHTOTYPE( $G$ )
2:   Create root of  $T$  with all vertices in  $G$ 
3:   for  $\ell \leftarrow 0$  to  $|V| - 2$  do
4:     Remove a vertex,  $V_{\ell+1}$ , randomly from node  $[\ell, \alpha]$ 
5:     for each node  $[\ell, i]$  in the  $\ell$ th level of  $T$  do
6:       if  $[\ell, i] \neq \text{NIL}$  then
7:          $[\ell, i].\text{left} \leftarrow$  neighbors of  $V_{\ell+1}$  in  $[\ell, i]$ 
8:          $[\ell, i].\text{right} \leftarrow$  non-neighbors of  $V_{\ell+1}$  in
            $[\ell, i]$ 
9:   return  $T$ 

```

---

## C. Encoding the Rooted Binary Tree $T$

As we are only interested in coding the structure of  $G$  (i.e., vertices up to an automorphism), we do not need to explicitly encode the set of vertices corresponding to the nodes in  $T$ . It is sufficient to encode only the cardinality of nodes; we call the tree with cardinality as node values  $\tilde{T}$  (see Figure 3). When we refer to encoding a node  $[\ell, i]$ , we are referring to encoding the value of the node.

Furthermore, we *only* need to encode the left nodes as the value of the siblings (i.e., right nodes) can be deduced given the value of the parent nodes. Consider a nonempty node  $[\ell, i]$ , we can see that

$$||[\ell, i]| = \begin{cases} ||[\ell, i].\text{left}| + ||[\ell, i].\text{right}| + 1, & \text{if } [\ell, i] = [\ell, \alpha] \\ ||[\ell, i].\text{left}| + ||[\ell, i].\text{right}|, & \text{otherwise.} \end{cases} \quad (2)$$

The reason for the discrepancy is because of our convention of always removing a vertex to encode at from the first nonempty node  $[\ell, \alpha]$  at each level.

It is easy to see that once the decoder has reconstructed the tree  $\tilde{T}$  and consequently  $T$ , one can reconstruct graph  $\tilde{G}$  isomorphic to the original graph  $G$ . The procedure is very similar to the one for transforming a graph into a binary tree. We start

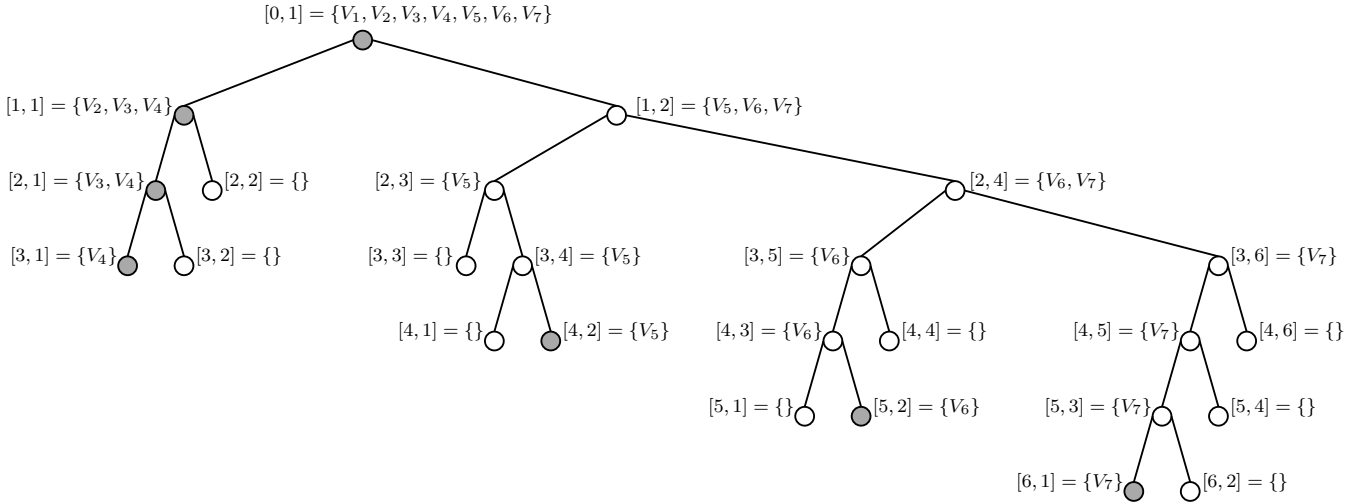


Figure 2: Binary tree representation  $T$  of the example graph  $G$  in Figure 1. The root contains all vertices of  $G$  and next levels are built by branching into neighbors and non-neighbors of vertex  $V_i$  among each subset of vertices. The pair of  $[\ell, i]$  next to each node shows the level that the node belongs to,  $\ell$ , and the position of the node in that level  $i$  is determined by counting from the left to the right.

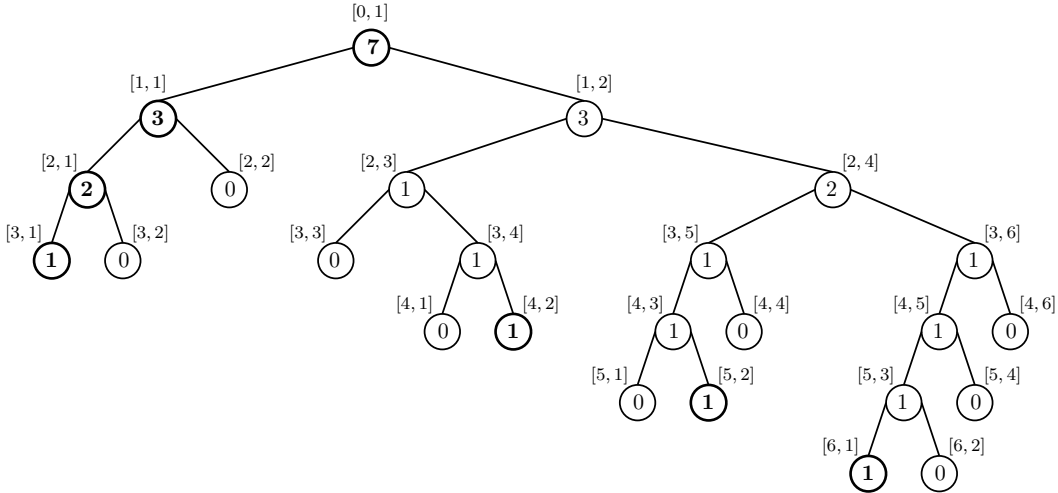


Figure 3: Transforming  $T$  in Figure 2 into  $\tilde{T}$ . Each node's value represents the cardinality of the corresponding node in  $T$ .

with  $|V|$  isolated vertices as determined by the cardinality of the root in  $T$ . Then for each level  $\ell = 1, \dots, |V| - 1$ , we connect the vertex  $V_\ell$  to vertices in left nodes at level  $\ell$  of  $T$ . The resulted graph  $\tilde{G}$  is isomorphic to the original graph  $G$ .

The remaining problem now is encoding the values in the tree  $\tilde{T}$ . The better the values can be predicted, the shorter the codelength. The encoder/decoder can predict the value either based on global properties of the graph (e.g., degree distribution), or based on local properties, which can depend only on the part of the tree  $\tilde{T}$  already encoded/decoded. One such property is that  $||[\ell, i]|$  can only take on integer values  $0, 1, 2, \dots, ||[\ell, i].\text{parent}| - 1$ , if  $i = 1, 2$  or  $0, 1, 2, \dots, ||[\ell, i].\text{parent}|$ , if  $i > 2$ .

It is important to realize that the properties used for encoding the tree  $\tilde{T}$  should be properties of the original graph, since they presumably have statistical significance. As an example, if the original graph is modeled as Erdős-Rényi [24], it is

easy to see that the node values are binomially distributed,  $\text{Binom}(N, p)$ , where

$$N = \begin{cases} ||[\ell, i].\text{parent}| - 1, & \text{if } i = 1, 2 \\ ||[\ell, i].\text{parent}|, & \text{if } i > 2. \end{cases} \quad (3)$$

The encoder uses the global property  $p$ , which must be known by the decoder. In universal coding, this can be transmitted to the decoder initially. Section II-F outlines how to calculate and encode global properties. We call this approach IID coder, as it is based on the i.i.d property of the Erdős-Rényi graph. It is about equivalent to the approach introduced in [1] for the compression of Erdős-Rényi graphs. The pseudocode for IID coder is shown in Algorithm 2.

In practice, Erdős-Rényi graphs are not good models for real-world graphs as edges are usually not independent [25]. Therefore, we would like to use more advanced graph properties. We classify our coding methods broadly into two classes:

---

**Algorithm 2** Encode  $\tilde{T}$  with IID coder
 

---

```

1: function ENCODEIID( $\tilde{T}, p$ )
2:   Encode  $||[0, 1]||$  via a positive integer encoder
3:   for  $\ell \leftarrow 1$  to  $|V| - 1$  do
4:     for each left node  $[\ell, i]$  do
5:       Encode  $||[\ell, i] \sim \text{Binom}(N, p)$  with  $N$  from
(3)

```

---

1) Node-by-Node Coder and 2) Level-by-Node Coder. In Node-by-Node Coder, we still use binomial distribution; however, the edge probability,  $p$ , will depend on local properties (i.e., graph motifs) in  $G(V, E)$ . In Level-by-Node Coder, we use the degree distribution of  $G(V, E)$  as a global property of graph to encode the values of all left nodes in level  $\ell$  of  $\tilde{T}$  at the same time.

The value of a left node in  $\tilde{T}$  can equivalently be seen as a count of *edges* of the original graph. Since only the *number* of edges are encoded, any local property used for encoding must be shared by all edges in a left node. Equivalently, any property of the original graph used for encoding must be convertible to a property of the tree  $\tilde{T}$ . In other words, we convert the original graph  $G(V, E)$  into the tree  $\tilde{T}$ , and any properties that we want to use for coding must then become properties purely of  $\tilde{T}$ . We will describe this with more details for each class of coders.

#### D. Class 1: Node-by-Node Coder

For Node-by-Node Coders, we will traverse  $\tilde{T}$  back to the root to determine the existence of certain motif structures in  $G$ . This will help us to better encode graphs that are not Erdős-Rényi as in these cases, edges are not independent of one another.

1) *Coding of Triangles* : The first Node-by-Node Coder we consider is the triangle coder. The triangle coder results in shorter codelength for graph classes that induce more triangle structures such as scale-free graphs. A triangle is a cycle graph with three nodes, which is also a 3-clique. Statistics about triangles are often used to characterize graphs [26].

First, we describe how to deduce the existence of triangles in  $G$  from the structure of  $\tilde{T}$ . We know that the set of vertices corresponding to a left node  $[\ell, i]$  are connected to the vertex  $V_\ell$ . We can also deduce if there are edges between the set of vertices corresponding to the left node  $[\ell, i]$  and the vertex  $V_{\ell-i}, i = 1, \dots, \ell-1$ . For that, we look at the ancestor of node  $[\ell, i]$  in level  $\ell - i$ . If it is a left node, then the set of vertices corresponding to the left node  $[\ell, i]$  are connected to the vertex  $V_{\ell-i}$ . To have a triangle, there must be an edge between  $V_\ell$  and  $V_{\ell-i}$ . It can be verified if the ancestor of  $[\ell-1, \alpha]$  at level  $\ell - i$  is a left node. If that is the case, we can deduce a triangle forms between any of vertices corresponding to the left node  $[\ell, i]$ ,  $V_\ell$ , and  $V_{\ell-i}$ . For example, consider node  $[4, 3] = \{V_6\}$  in Figure 2. Since node  $[4, 3]$  is a left node in level 4, there is an edge between  $V_6$  and  $V_4$ . We can see that  $V_6$  is also connected to  $V_3$  because its ancestor in level 3 (i.e., node  $[3, 5]$ ) is a left node. We can also verify that  $V_4$  and  $V_3$  are connected since node  $[3, 1]$  (i.e., node  $[3, \alpha]$ ) is a left node at level 3 and

therefore, connected to  $V_3$ . Considering all these connections, we can deduce that  $V_6, V_4, V_3$  form a triangle subgraph in  $G$ . We can use the formation of triangles to encode the nodal values in  $\tilde{T}$  as follows.

Similar to IID coder, we use binomial distribution. However, the triangle coder chooses between two binomial distributions:  $\text{Binom}(N, \check{p}_\Delta)$  or  $\text{Binom}(N, p_\Delta)$ , where  $N$  is given by equation (3). We decide between these two distributions with the help of  $\mathcal{CI}([\ell, i].\text{parent})$ . Consider a left node  $[\ell, i]$ , the coder will use  $\text{Binom}(N, p_\Delta)$  to encode its value if  $\mathcal{CI}([\ell, i].\text{parent})$  has at least one element. Note that  $\mathcal{CI}([\ell, i].\text{parent})$  represents levels where  $\mathcal{R}([\ell, i].\text{parent})$  and  $\mathcal{R}([\ell-1, \alpha])$  both contain left nodes. It means that vertices corresponding to left node  $[\ell, i]$  and  $V_{\ell-1}$  are both connected to vertex  $V_{\ell'}$ ,  $\ell' \in \mathcal{CI}([\ell, i].\text{parent})$ , and therefore, triangle forms among these vertices in  $G$ . When  $\mathcal{CI}([\ell, i].\text{parent})$  is empty, we encode  $[\ell, i]$  with  $\text{Binom}(N, \check{p}_\Delta)$  as no triangle exists among these vertices in  $G$ . The pseudocode for the triangle coder is given in Algorithm 3.

---

**Algorithm 3** Encode  $\tilde{T}$  with triangles from Class 1
 

---

```

1: function ENCODETRIANGLES( $\tilde{T}, \{\check{p}_\Delta, p_\Delta\}$ )
2:   Encode  $||[0, 1]||$  via a positive integer encoder
3:   for  $\ell \leftarrow 1$  to  $|V| - 1$  do
4:     for each left node  $[\ell, i]$  do
5:       if  $|\mathcal{CI}([\ell, i].\text{parent})| > 0$  then
6:         Encode  $||[\ell, i] \sim \text{Binom}(N, p_\Delta)$ 
7:       else
8:         Encode  $||[\ell, i] \sim \text{Binom}(N, \check{p}_\Delta)$ 

```

---

2) *Coding with the number of common neighbors*: In a triangle subgraph, two connected vertices share a single common neighbor. But it may be that two vertices share multiple common neighbors as there is usually a correlation between having an edge between two vertices and the number of their common neighbors. This property is used for link prediction in complex networks [27], [28]. Therefore, we can generalize the triangle encoder to encoding with  $m$  common neighbors. Instead of using  $\check{p}_\Delta$  or  $p_\Delta$  to parameterize the binomial distribution, we can use  $p_{\Delta(m)}$  where  $m$  denote the number of common neighbors (note that  $p_{\Delta(0)}$  is the same as  $\check{p}_\Delta$  in coding with triangles). In other words, we find how many edges exist between any vertex corresponding to left node  $[\ell, i]$  and vertex  $V_\ell$  and utilize it as statistics that built the graph. The number of common neighbors can range from 0 to the maximum degree of a vertex in  $G$ , which we denote by  $d_{max}$ .

The way we encode nodal values of  $\tilde{T}$  with this coder is as follows. For a left node  $[\ell, i]$ , we look at the cardinality of  $\mathcal{CI}([\ell, i].\text{parent})$ . It tells us how many vertices already encoded exist that any of vertices corresponding to  $[\ell, i]$  and  $V_\ell$  are both connected to them. This determines which parameter to pick for binomial distribution.

3) *Coding over 4-node motifs*: We can explore motifs of larger sizes in  $G$  to develop new coders. However, computational cost can be a limiting factor. Here, we are interested to extend coding to motifs of size 4 vertices. Figure 4 illustrates

motifs that we look for. Note that the priority starts with 4-clique. If 4-clique does not exist, then we look for double-triangle, and finally 4-cycle. Each of these motifs can be encoded with  $p_{\boxtimes}$ ,  $p_{\boxminus}$ , and  $p_{\square}$ , respectively. If none of them exist, then we encode left node  $[\ell, i]$  using  $\check{p}_{\square}$ .

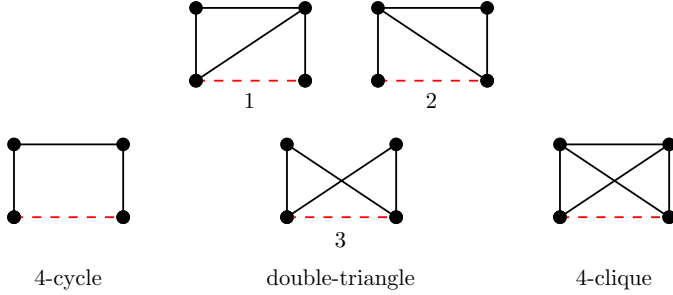


Figure 4: Motifs of interest over four vertices with three realizations for double-triangle.

As we mentioned earlier, we find the appropriate parameter to encode nodal values of  $\tilde{T}$  by looking at the structure of encoded/decoded part of tree  $\tilde{T}$  up to that point. Depending on the cardinality of  $\mathcal{CI}([\ell, i].\text{parent})$ , different cases are possible:

- $|\mathcal{CI}([\ell, i].\text{parent})| \geq 2$  : For any pair  $\ell_1$  and  $\ell_2$  in  $\mathcal{CI}([\ell, i].\text{parent})$  where  $\ell_1 < \ell_2$ , we check to see if we can find a case in which the first nonempty node in level  $\ell_2$  (i.e.,  $[\ell_2, \alpha]$ ) lies in the left subtree of the tree rooted at each node in level  $\ell_1$ . If such a case exists, we use  $\text{Binom}(N, p_{\boxtimes})$  to encode the value of  $[\ell, i]$ . Otherwise, the configuration 3 of double-triangle in Figure 4 occurred; thus, we use  $\text{Binom}(N, p_{\boxminus})$ .
- $|\mathcal{CI}([\ell, i].\text{parent})| = 1$  : Two cases are possible
  - For the first nonempty node in level  $\ell_1, \ell_1 \in \mathcal{CI}([\ell, i].\text{parent})$ , and node  $[\ell-1, \alpha]$ , we check to see if  $\mathcal{R}([\ell_1, \alpha])$  and  $\mathcal{R}([\ell-1, \alpha])$  contain a left node at the same level. If such a node exists, the configuration 1 of double-triangle in Figure 4 occurred; thus, we use  $\text{Binom}(N, p_{\boxminus})$ .
  - For node  $[\ell_1, \alpha]$  and  $[\ell, i].\text{parent}$ , we check to see if  $\mathcal{R}([\ell_1, \alpha])$  and  $\mathcal{R}([\ell, i].\text{parent})$  contain a left node at the same level. If such a node exists, the configuration 2 of double-triangle in Figure 4 occurred; thus, we use  $\text{Binom}(N, p_{\boxminus})$ .
- $|\mathcal{CI}([\ell, i].\text{parent})| = 0$  : First, we need to find  $\mathcal{I}([\ell, i].\text{parent})$  which contains all the levels that  $\mathcal{R}([\ell, i].\text{parent})$  or  $\mathcal{R}([\ell-1, \alpha])$  have left nodes. For  $|\mathcal{I}([\ell, i].\text{parent})| > 1$  and any pair  $\ell_1$  and  $\ell_2$  in  $\mathcal{I}([\ell, i].\text{parent})$  where  $\ell_1 < \ell_2$ , we look for a case in which the first nonempty node in level  $\ell_2$  (i.e.,  $[\ell_2, \alpha]$ ) lies in the left subtree of the tree rooted at any node in level  $\ell_1$ . If such a case exists, we use  $\text{Binom}(N, p_{\square})$  to encode the value of  $[\ell, i]$ .
- If none of the abovementioned cases occurred, we use  $\text{Binom}(N, \check{p}_{\square})$  to encode the value of  $[\ell, i]$ .

Due to similarity, we skip the pseudocode for this coder to avoid repetition.

## E. Class 2: Level-by-Node Coder

One important graph statistics is the degree distribution  $P(k)$ . Class 2 coders utilize the degree distribution in addition to graph motifs presented in Class 1 coders to efficiently encode  $\tilde{T}$ .

For a given level,  $\ell$ , in  $\tilde{T}$ , Class 1 coders encode the value of each left node (and deduces the value of the right node) independently of one another. Class 2 coders on the other hand, utilize the degree distribution to encode the values of left nodes at the same level altogether.

The number of left nodes encountered in  $\mathcal{R}([\ell-1, \alpha])$ ,  $\check{k}$ , is equal to the number of edges,  $V_\ell \leftrightarrow V_{\ell-i}, i = 1, \dots, \ell-1$ , in  $G$ . We know that the degree of vertex  $V_\ell$  is lower bounded by  $\check{k}$ . The sum of the values of left nodes in level  $\ell$  is  $k - \check{k}$ , where  $k$  is the degree of vertex  $V_\ell$ . We can use this relationship to encode all the left nodes in level  $\ell$  of  $T$  at the same time.

Encoding with Class 2 entails two steps. First, we need to send the degree of  $V_\ell$ . Then, we send the conditional probability of seeing specific values for left nodes at level  $\ell$  conditioned on their summation to be  $k - \check{k}$ .

For the first step, we use the degree distribution  $P(k)$  and the fact that  $k \geq \check{k}$ . Thus, we encode the degree of  $V_\ell$  with

$$P(k|k \geq \check{k}) = \frac{P(k)}{\sum_{j \geq \check{k}} P(j)}. \quad (4)$$

One should notice that at the time of decoding a given level, the decoder knows  $\check{k}$  based on the tree up that point and thus, it can calculate (4). Now, the decoder can compute the summation of the left nodes' values at level  $\ell$  with

$$\sum_{i \text{ is odd}} |[\ell, i]| = k - \check{k}. \quad (5)$$

Let  $[\ell, N_\ell]$  denote last left node in level  $\ell$ . To encode left nodes' values in level  $\ell$ , we need to compute the joint probability of observing specific values for left nodes at level  $\ell$  (i.e., left nodes take values  $|[\ell, 1]|, \dots, |[\ell, N_\ell]|$  respectively) conditioning on their summation being equal to (5).

$$\begin{aligned} & P(|[\ell, 1]|, \dots, |[\ell, N_\ell]| \mid \sum_{i \text{ is odd}} |[\ell, i]| = k - \check{k}) \\ &= \frac{P(|[\ell, 1]|, \dots, |[\ell, N_\ell]|)}{P(\sum_{i \text{ is odd}} |[\ell, i]| = k - \check{k})} \\ &= \frac{\prod_{i=1}^{N_\ell} P(|[\ell, i]|)}{P(\sum_{i \text{ is odd}} |[\ell, i]| = k - \check{k})}, \end{aligned} \quad (6)$$

where the numerator is the joint probability of observing  $|[\ell, 1]|, \dots, |[\ell, N_\ell]|$  as the values of left nodes and the denominator is the probability corresponding to all configurations with the same summation for left nodes' value,  $k - \check{k}$ . By independent assumption on the values of left nodes, the joint probability distribution  $P(|[\ell, 1]|, \dots, |[\ell, N_\ell]|)$  will reduce to the product of individual probabilities.

Assuming the independence assumption, two cases are possible:

- 1) *Identically distributed*: This case happens when we want to compute (6) using IID coder. Then, the probability of success is the same for all probabilities and the problem reduces to a counting problem. We can encode the values of left nodes using the probability distribution

$$\begin{aligned}
& P(|[\ell, 1]|, \dots, |[\ell, N_\ell]| \mid \sum_{i \text{ is odd}} |[\ell, i]| = k - \check{k}) \\
&= \frac{\prod_{i \text{ is odd}} \binom{|[\ell, i]| + |[\ell, i+1]|}{|[\ell, i]|}}{\binom{\sum_i |[\ell, i]|}{\sum_{i \text{ is odd}} |[\ell, i]|}} \\
&= \frac{\prod_{i \text{ is odd}} \binom{|[\ell, i]| + |[\ell, i+1]|}{|[\ell, i]|}}{\binom{\sum_i |[\ell, i]|}{k - \check{k}}}. \quad (7)
\end{aligned}$$

- 2) *Not identically distributed*: If we use any other coders except IID coder, then different left nodes in the same level *do not* have the same parameter for encoding. For example as we showed with the triangle coder from Class 1, node  $[2, 1]$  in Figure 3 would be encoded with  $p_\Delta$  while node  $[2, 3]$  would be encoded with  $\check{p}_\Delta$ . To take this into account, we need to encode the values of left nodes using the following probability distribution

$$\begin{aligned}
& P(|[\ell, 1]|, \dots, |[\ell, N_\ell]| \mid \sum_{i \text{ is odd}} |[\ell, i]| = k - \check{k}) \\
&= \frac{\prod_{i \text{ is odd}} \binom{|[\ell, i]| + |[\ell, i+1]|}{|[\ell, i]|} \theta^{|[\ell, i]|} (1 - \theta)^{|[\ell, i+1]|}}{P(\sum_{i \text{ is odd}} |[\ell, i]| = k - \check{k})}, \quad (8)
\end{aligned}$$

where  $\theta$  is the relevant parameter for each left node in the tree  $\tilde{T}$  (e.g.,  $p_\Delta$  or  $\check{p}_\Delta$  when using triangles for coding). The denominator in (8) is a generalized version of binomial distribution called Poisson binomial distribution [29]. It is defined as the sum of independent binomial distributions that are not necessarily identically distributed. Due to the involvement of different probabilities, calculation of denominator in (8) can be cumbersome. To resolve this, some methods were proposed in the literature. Recursive methods were developed in [30], [31] that can compute the denominator in (8) in  $O((k - \check{k}) \sum_i |[\ell, i]|)$  time. We pay this extra computational cost to have a more efficient coder since the entropy of Poisson binomial distribution is bounded above by the entropy of binomial distribution with the same mean [32]. The pseudocode to encode with Class 2 is given in Algorithm 4.

#### F. Calculation and Encoding of Statistics

We consider encoding in two scenarios: learned coding, and universal coding. In learned coding, we are given a set of training graphs  $\{G_1, \dots, G_N\}$  of a particular class and have to learn local and global statistics; these statistics, then, are

---

#### Algorithm 4 Encode $\tilde{T}$ with Class 2

---

```

1: function ENCODECLASS2( $\tilde{T}$ )
2:   Encode  $|[0, 1]|$  via a positive integer encoder
3:   for  $\ell \leftarrow 1$  to  $|V| - 1$  do
4:      $\check{k} \leftarrow$  Number of left nodes in  $\mathcal{R}([\ell - 1, \alpha])$ 
5:      $\hat{k} \leftarrow$  Summation of left nodes' values in level  $\ell$ 
6:      $k \leftarrow \check{k} + \hat{k}$ 
7:     Encode  $k$  with (4)
8:      $P \leftarrow$  Compute (6)
9:   Encode  $P$ 

```

---

shared to both encoder and decoder. In universal coding, there is no training set and the encoder encodes a single graph. It also has to communicate to the decoder what is the statistics. Below we describe calculation and communication of statistics for each of these scenarios.

1) *Learned Coding*: For learned coding, we need to learn statistics from a set of training graphs. To do that, each statistic is calculated by taking an average over the same statistic in the training set. The edge probability  $p$  in coding with IID coder can be estimated by the average degree. Other edge statistics are more tricky and should reflect the procedure used for encoding the graph. It means that we cannot simply count the number of triangles in a graph to compute  $p_\Delta$ . The reason is when we want to encode a graph with triangles, we look for the formation of a triangle in a specific way forced by the coding algorithm as described earlier. Thus, we need to transform each graph  $G_i$  from the training set into its rooted binary tree representation  $T_i$  and find statistics with coding algorithm. To estimate  $p_\Delta$  and  $\check{p}_\Delta$ , we traverse each  $T_i$  the same way we did for coding and divide the tree's nodes into those coded with  $p_\Delta$  and those coded with  $\check{p}_\Delta$ . To estimate each of  $p_\Delta$  and  $\check{p}_\Delta$ , we compute the ratio of the summation of left nodes' values to the summation of left and right nodes' values in that group. The average over all  $T_i$  gives the estimation for  $p_\Delta$  and  $\check{p}_\Delta$ . The same recipe is used to estimate edge statistics for common neighbors and 4-node motifs.

The estimation of degree distribution in Class 2 is straightforward. It can be estimated through the histogram. We compute the degree distribution for each graph in the training set and the final degree distribution is estimated by taking average over them.

2) *Universal Coding*: For encoding average degree in IID coder, we can send the number of edges  $|E|$  in  $G$ . The number of bits required to encode the number of edges is about  $\log \frac{n(n-1)}{2} \approx 2 \log n$  bits. Once the decoder knows the number of edges, it can compute the parameter of IID coder,  $p$ , with

$$p = \frac{2|E|}{n(n-1)}.$$

For other local statistics, we use sequential estimation best outlined in [33, Section 13.2]. For example in coding with triangles, we use sequential estimation of  $p_\Delta$  and  $\check{p}_\Delta$ , specifically the KT estimator [34], [35], which is

$$\hat{p} = \frac{n_1 + \frac{1}{2}}{n_1 + n_0 + 1},$$

where  $n_1, n_0$  are the summation of left and right nodes' values previously coded in the rooted binary tree, respectively. One should note that the procedure for updating the probabilities  $p_\Delta$  and  $\check{p}_\Delta$  is different for each class. For Class 1 coders, the probabilities will be updated after coding each node of the rooted binary tree. However, for Class 2 coders, the probabilities will be updated after coding each level. The reason is that all nodes at the same level will be encoded together. To estimate statistics for common neighbors and 4-node motifs, we utilize similar sequential approach used for the estimation of  $p_\Delta$  and  $\check{p}_\Delta$ .

For the degree distribution, we calculate the degree histogram for the whole graph, and use this for coding. The degree of a node is between 0 and  $n-1$ . We can therefore think of the degree histogram as putting each of the  $n$  (unlabeled) nodes into one of  $n$  buckets, and encoding this can be done by encoding the counts in the buckets. The number of possible configurations is a standard problem in combinatorics:  $\binom{2n-1}{n}$ , which can be transmitted with

$$\begin{aligned} \log \binom{2n-1}{n} &= nH \left( \frac{n}{2n-1} \right) + \frac{1}{2} \log \frac{2n-1}{n^2} + c \\ &\approx n - \frac{1}{2} \log n \text{ bits } (|c| \leq 2). \end{aligned}$$

Hereinafter, in our experiments, we use the above mentioned approach to calculate and communicate local and global statistics of graphs.

### G. Experiments

We consider two cases for experiments: 1) using learned coding to encode a graph of the same class of training set, 2) using universal coding to encode a single graph. In the former case, we generate synthetic data for different classes of graphs, whereas in the later case we measure the performance on real-world graphs. We evaluate the performance of our coding methods versus IID coder and also compare the performance of Class 1 and Class 2 against each other.

For learned coding, we first need to learn statistics for coders. We considered different classes of graphs in our experiments. In all cases, learning was done on 50 graphs. Then, we used those statistics from the training to encode a test graph of the same type. The results are shown in Figure 5. As expected, for Erdős-Rényi graph, IID coder is efficient and all other coders do not offer an improvement. However, for Barabási-Albert and Watts Strogatz graphs, our proposed coders outperform IID coder by a significant margin. One can observe that coders in Class 2 have shorter codelength than their counterparts in Class 1. For Barabási-Albert graph, all encoders in Class 2 almost have the same performance and therefore, choosing one or another does not matter. However, for Watts Strogatz graph, coding with common neighbors through Class 2 is the most efficient.

For compression of real-world graphs, we use universal coding since there is no training set. Note that the statistics associated with each coder is required to be communicated to the decoder as described in Section II-F. We measure the compression performance of Class 1 and Class 2 coders against each other and also compute the codelength associated

with IID coder and arithmetic coding (represented by labeled iid) for comparison. The following undirected graphs, publicly available to download<sup>1,2</sup>, are considered:

- Genetic interaction: a gene interaction network [36].
- Economic Network: the economic network of Victoria, Australia in 1880 [36].
- Transportation network: A simple graph of Minnesota road network [36].
- Collaboration network: extracted from *arXiv* papers and covers collaborations among authors in General Relativity and Quantum Cosmology category [37].
- Facebook-politician: represents blue verified Facebook page networks of politicians [38].
- Internet (AS level): An undirected graph representing AS peering information inferred from University of Oregon Route Views Project on January 1, 2000 [39].

As we can see from Table I, coders in Class 2 give a better performance compared to their counterparts in Class 1. This improvement comes at the cost of computing (8). Moreover, shorter codelength is achievable by using Class 1 over IID coder. Another observation is that coding with common neighbors by means of Class 2 outperforms other coders in most cases. However, there are some cases that coding over 4-node motifs and triangles offer improvement over common neighbors. There is no single coder that is best for all graphs. This is no different than the situation for encoding sequences. If one wants the most efficient coder for a particular graph, one can encode with multiple algorithms, and choose the shortest, indicating by a header which is used or use soft combining [40]. In the end it is a tradeoff between compression and complexity.

### III. MODEL SELECTION IN GAUSSIAN GRAPHICAL MODELS

In this section, we introduce an application of graph coding for model selection in Gaussian graphical models. The goal is to find a graph that describes dependencies among multivariate Gaussian variables accurately [41]. Multivariate Gaussian distributions are widely used in modeling real-world problems where the relationship among approximately normally distributed variables is of interest. They have been used to model different datasets such as stock returns [42], [43], protein-protein interactions [44], and brain functional activities [45].

Let  $X = [X_1, \dots, X_p]^T$  be a  $p$ -dimensional random vector with multivariate Gaussian distribution  $\mathcal{N}(0, \Sigma)$ . The inverse of covariance matrix,  $\Sigma^{-1} = \Omega$ , is known as the precision matrix. If the  $(i, j)$  entry of the precision matrix (i.e.,  $\Omega_{ij}$ ) is 0, then  $X_i$  and  $X_j$  are conditionally independent given all the other variables. Therefore, we can visualize the conditional independence relationships between any pair of variables as an unweighted, undirected graph,  $G$  whose adjacency matrix,  $A$  is such that

$$A_{ij} = \begin{cases} 1 & \text{if } \Omega_{ij} \neq 0, i \neq j \\ 0 & \text{otherwise.} \end{cases} \quad (9)$$

<sup>1</sup><http://networkrepository.com/>

<sup>2</sup><https://snap.stanford.edu>



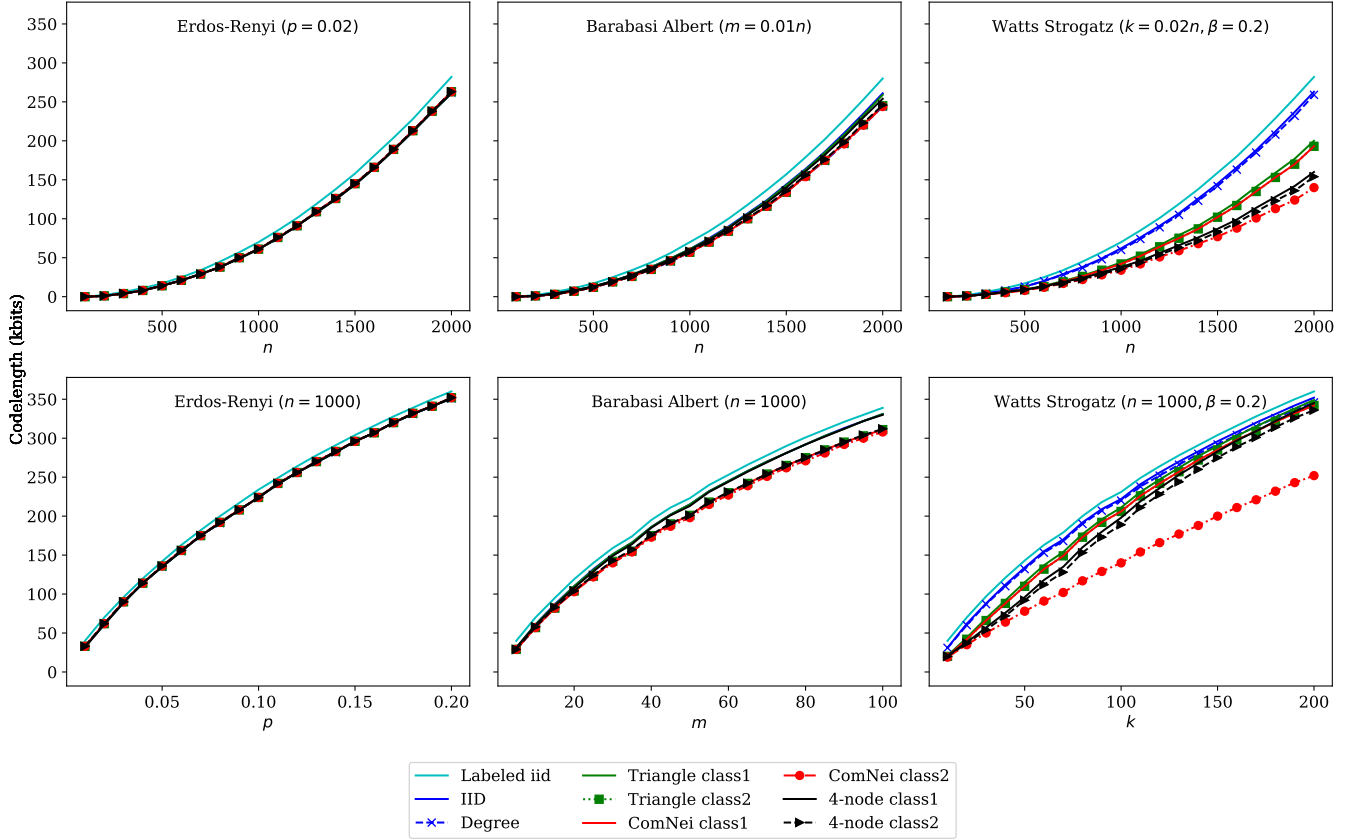


Figure 5: Comparison of codelength associated with each coder for Erdős-Rényi, Barabási-Albert, and Watts Strogatz graphs. Left plots are for the graphs of variable size and the average degree of  $\langle k \rangle = \frac{n}{50}$ . Right plots are for graphs of the same size  $n = 1000$  and  $\langle k \rangle$  ranging from  $\frac{n}{100}$  to  $\frac{n}{5}$ . In Barabási-Albert graph,  $m$  denote the number of new edges added by each node. In Watts Strogatz,  $k$  is the number of nearest neighbors in the model and  $\beta$  refer to rewiring probability.

Name	$ V $	$ E $	Labeled iid	IID coder	Class 1			Class 2			
					Triangle	ComNei	4-node	Degree	Triangle	ComNei	4-node
Genetic interaction	4227	39484	365599	317379	273444	235011	249467	280375	235787	<b>206485</b>	220191
Economic network	1258	7513	61224	47592	49790	49597	46407	45500	44159	43996	<b>40730</b>
Transportation network	2642	3303	37937	35424	14619	14620	14596	12712	<b>12530</b>	12531	12540
Collaboration network	5242	14496	164224	119038	60375	52041	66203	106215	53647	<b>47320</b>	60182
Facebook-politician	5908	41729	423451	340070	222500	199514	206720	327203	195739	<b>172386</b>	179836
Internet (AS level)	3570	7750	86206	47685	48266	46665	47222	31107	30469	<b>29708</b>	30897

Table I: Compression length (bits) for some real-world graphs. Labelled iid encodes the graph with its labels. Degree denote IID coder when we use it along degree distribution in Class 2. The best value for each case is boldfaced.

The graph  $G$  is known as the conditional independence graph. Its structure is determined by the sparsity pattern of  $\Omega$ . [16], [46]

Often, we are interested in the problem of estimating  $G$  from  $N$  i.i.d observations. To do that, we need to estimate  $\Omega$  first. The problem is especially difficult when  $N < p$ . In this case, the maximum likelihood estimate does not exist. In most approaches, the estimator of  $\Omega$  contains a regularization term  $\lambda$  to 1) prevent overfitting, 2) control the sparsity of the estimated solution. A popular sparse estimator of the precision matrix is found by maximizing the  $L_1$ -penalized log-likelihood [41],

$$\max_{\Omega \succ 0} \log \det \Omega - \text{tr}(S\Omega) - \lambda \|\Omega\|_1, \quad (10)$$

where  $S$  is the sample covariance matrix, and  $\lambda$  is the regularization parameter that controls sparsity. The performance of the estimator in (10) is highly influenced by the selection of  $\lambda$ . Depending on the value of  $\lambda$ , the conditional independence graph can range from a dense graph for small values of  $\lambda$  to a graph with zero edges when  $\lambda$  takes large values. To find the best value of the regularization parameter  $\lambda$ , model selection techniques are being used.

We use the MDL principle and select  $\lambda$  that minimizes the sum of the description length of the conditional independence graph structure and data when encoded under that graph structure

$$\arg \min_{\lambda} L(G_{\lambda}) + L(D|G_{\lambda}), \quad (11)$$

where  $L(G_{\lambda})$  is the description length of the conditional independence graph  $G_{\lambda}$ , and  $L(D|G_{\lambda})$  is the description length of data  $D$  when encoded with  $G_{\lambda}$ .

To compute the description length of the conditional independence graph,  $L(G_{\lambda})$ , we first need to find the underlying conditional independence graph  $G_{\lambda}$ . By applying an estimator to the sequence of  $N$  i.i.d observations  $\{x_1, \dots, x_N\}$  from  $p$ -variate Gaussian distribution, we obtain an estimated precision matrix  $\hat{\Omega}$  for each realization of  $\lambda$ . We have dropped the explicit dependence on  $\lambda$  from  $\hat{\Omega}$  and other  $\lambda$ -related precision and covariance matrices for the sake of simplicity of notation. In this paper, we use graphical lasso [41], which is one of the most commonly used solvers of  $L_1$ -penalized log-likelihood in Gaussian graphical models. It is worth mentioning that our approach can be easily applied to any other estimator. The conditional independence graph  $G_{\lambda}$  is obtained from the estimated precision matrix  $\hat{\Omega}$ . We then pick a graph coder described in Section II to compute the description length of  $G_{\lambda}$ .

To compute the description length of the data,  $L(D|G_{\lambda})$ , we face two challenges. The first challenge is to deal with real-valued data  $\{x_1, \dots, x_N\}$  where lossless source coding is not generalized directly. Second, we have to encode the data based on the underlying conditional independence graph  $G_{\lambda}$ . To encode real-valued data, we assume a fixed-point representation with a (large) finite number,  $r$ , bits after the period, and an unlimited number of bits before the period [3]. Therefore, the number of bits required to represent data  $x$  according to the probability distribution function (pdf)  $f(x)$  is given by

$$\begin{aligned} L(x) &= -\log \int_x^{x+2^{-r}} f(t)dt \approx -\log(f(x)2^{-r}) \\ &= -\log f(x) + r. \end{aligned} \quad (12)$$

Since we are only interested in relative codelengths between different models, the dependency on  $r$  cancels out.

The second challenge is how to encode the data with respect to  $G_{\lambda}$ . As mentioned, the data is generated by multivariate Gaussian distribution, which can be characterized by covariance matrix  $\Sigma$ . Since we do not have access to the true covariance matrix  $\Sigma$ , we can find an estimate  $\check{\Sigma}$  based on the sample covariance matrix  $S$ . It should be noted that the structure of  $\check{\Sigma}^{-1}$  should match with the structure of  $G_{\lambda}$  obtained from the previous step

$$\begin{cases} \check{\Sigma}_{ij} = S_{ij} & \text{if } i = j \text{ or } A_{ij} \neq 0, \\ \check{\Sigma}_{ij}^{-1} = 0 & \text{otherwise.} \end{cases} \quad (13)$$

This problem is known as *matrix completion* problem [47]. Dempster in [48], proved the existence and uniqueness of maximum likelihood solution for  $\check{\Sigma}$  when the sample covariance matrix  $S$  is positive definite. He presented a coordinate-descent

algorithm to find the solution iteratively. It should be noted that the precision matrix obtained from the graphical lasso solution  $\hat{\Omega}$  is not necessarily equal to the inverse of estimated covariance matrix  $\check{\Sigma}$ .

Once we estimate  $\check{\Sigma}$ , we use predictive MDL [5] to compute  $L(D|G_{\lambda})$

$$L(D|G_{\lambda}) = -\sum_{i=0}^{N-1} \log f(x_{i+1}|\hat{\theta}(x_1, \dots, x_i)), \quad (14)$$

where  $f(\cdot|\cdot)$  is the conditional pdf and  $\hat{\theta}(x_1, \dots, x_i)$  denote the maximum likelihood estimate of parameters, which in this case is the estimated covariance  $\check{\Sigma}$  obtained under  $G_{\lambda}$ . The codelength in (14) is a valid codelength since it is sequentially decodable. Note that it does not work for the first few samples, as there is no estimate for  $\check{\Sigma}$ . Instead, we encode the first few samples with a default distribution, which is the same among different realizations of  $\lambda$ .

Finally, the best conditional independence graph structure  $G_{\lambda^*}$  is obtained by minimizing  $L(G_{\lambda}) + L(D|G_{\lambda})$  over  $\lambda$ . In the following, we present an algorithm to find the best graph model of data  $G_{\lambda^*}$  associated with  $\lambda^*$ .

---

**Algorithm 5** Find the best graph model of data  $G_{\lambda^*}$  via graph coding

---

**Input:** Samples  $\{x_1, \dots, x_N\} \sim \mathcal{N}(0, \Sigma)$  and regularization parameters  $\{\lambda_1, \dots, \lambda_K\}$ .

**Output:** Best graph model of data  $G_{\lambda^*}$ .

- 1: **for** Each realization of  $\lambda \in \{\lambda_1, \dots, \lambda_K\}$ . **do**
  - 2:   Apply graphical model estimator to find  $\hat{\Omega}$  and build  $G_{\lambda}$  based on (9).
  - 3:   Use any graph coder in Section II to compute  $L(G_{\lambda})$ .
  - 4:   Compute  $L(D|G_{\lambda})$  via (14) where  $\hat{\theta}(x_1, \dots, x_i)$  denote  $\check{\Sigma}$  obtained from (13).
  - 5:   Add up codelengths resulted from step (3) and step (4).
  - 6: **return**  $G_{\lambda^*}$  associated with the shortest total codelength in step (5).
- 

## A. Experiments

We tested our approach on both synthetic and real-world data. For synthetic data, we consider different conditional independence graph structures and the goal is to recover true conditional independence graph. For real-world data, we applied our approach to ECG data of a group of healthy subjects and a group of subjects with Kawasaki disease. We are interested to determine if there is any difference between the conditional independence graph of healthy group versus Kawasaki group.

We now provide simulation results on synthetic data generated from zero-mean multivariate Gaussian distribution with known precision matrix,  $\Omega$ . We will compare the performance of graph coding methods against other methods in the recovery of conditional independence graph  $G$ . We use F1-score as a widely used metric in the literature. F1-score is the harmonic mean of precision and recall where the precision is the ratio

of the number of correctly estimated edges to the total number of edges in the estimated conditional independence graph, and the recall is the ratio of the number of correctly estimated edges to the total number of edges in the true conditional independence graph [23]. We consider two cases: 1) the number of observations is larger than the number of variables,  $N > p$ , with  $N/p = 2$  and 2) the number of observations is smaller than the number of variables,  $N < p$ , with  $N/p = 0.5$ . The experiments were repeated for  $p = 100, 200$ . We applied the graphical lasso estimator described in [41] for  $\lambda \in [0.01, 1]$  when  $N > p$  and  $\lambda \in [0.1, 1]$  when  $N < p$  to cover a wide range of structures from dense graphs to totally isolated nodes.

The reason for considering smaller range for  $N < p$  is the fact that the graphical lasso estimator does not converge for  $\lambda < 0.1$  in our test cases. The values of  $\lambda$  are equally spaced apart with the step size of 0.01. Multivariate Gaussian data was generated with different precision matrix structures that have been frequently used as test cases in the literature [49], [50]:

- Cycle structure with  $\Omega_{ii} = 1$ ,  $\Omega_{i,i-1} = \Omega_{i-1,i} = 0.5$ ,  $\Omega_{1p} = \Omega_{p1} = 0.4$ .
- Autoregressive process of order one AR(1) with  $\Omega_{ii} = 1$ ,  $\Omega_{i,i-1} = \Omega_{i-1,i} = 0.5$ .
- Erdős-Rényi (ER) structure with  $\Omega_2 = \Omega_1 + \delta I_p$  where  $\Omega_1$  is a matrix with off-diagonal elements taking values randomly chosen from uniform distribution  $\mathcal{U}(0.4, 0.8)$  with the probability of  $2/p$  and diagonal values set to zero. To keep  $\Omega_2$  positive definite, we choose  $\delta = \rho + 0.05$  where  $\rho$  is the absolute value of the minimum eigenvalue of  $\Omega_1$ . Here  $I_p$  is the identity matrix of size  $p$ .
- Hub structure with two hubs. First, we create the adjacency matrix  $A$  by setting off-diagonal elements to one with the probability of 0.01 and zero otherwise. Next, we randomly select two hub nodes and set the elements of the corresponding rows and columns to 1 with the probability of 0.7 and zero otherwise. After that for each nonzero element  $A_{ij}$ , we set  $A_{ij}$  with a value chosen randomly from uniform distribution  $\mathcal{U}(-0.75, -0.25) \cup (0.25, 0.75)$ . Then, we set  $\Omega_1 = \frac{1}{2}(A + A^T)$ . The final precision matrix  $\Omega_2$  is obtained by  $\Omega_2 = \Omega_1 + \delta I_p$  with the same set-up as Erdős-Rényi structure.

Table II and Table III show the results of applying different methods to recover the conditional independence graph by applying graphical lasso as the estimator for  $N > p$  and  $N < p$  cases, respectively. To have a ground-truth for comparison, we provide the highest possible value for F1-score on the regularization path and show it in the column with name Optimum. For benchmark methods, we considered CV, BIC, and EBIC methods. The results for CV are given for 5-fold CV, and for EBIC method are given for recommended value of  $\gamma = 0.5$  [18]. The results for graph-coding methods are obtained by following the steps outlined in Algorithm 5. The values are the average of 50 Monte Carlo trials.

It can be seen that graph-coding techniques, i.e., columns named with IID, Class 1, and Class 2, outperform other methods in most of cases. Furthermore, F1-score associated with recovered conditional independence graph by these methods

is so close to the optimum value on the regularization path. Among graph coding methods, we observed that coding with Class 1 and Class 2 coders give better performance than coding with IID coder in most of the cases. In addition, often times coders in Class 2 may offer a slight improvement over their counterparts in Class 1 coders. These findings confirm the necessity of having more advanced graph coders that reflect more information about the graph in their codelength.

We also tested our approach on a real-world dataset. This dataset contains extracted features from 12-lead ECG signals of a group of healthy subjects and a group of subjects with Kawasaki disease. All subjects are of age one to three years. The healthy group has 2492 samples and the Kawasaki group has 197 samples. This dataset was initially processed to select features with empirical distribution close to normal distribution. The reason is that the underlying assumption in Gaussian graphical models is that the data is normally distributed and this is also a key assumption for encoding data under our approach. After screening all features and selecting those with approximately normal distribution, we ended up with 30 ECG features. We applied our approach to find the graph model of data for each group of subjects when  $\lambda$  takes equally spaced apart values with the step size of 0.01 within the range  $[0.1, 1]$ . The results show that the graph model of healthy subjects differs from the graph model of subjects with Kawasaki disease as given in Figure 6. The features with a number from 1 to 12 (referring to the lead number) at the end of their name vary across different leads. The rest of features are the same across all leads. The conditional independence graph given by our approach for each group of subjects is the same across different graph coders. As it can be seen, the conditional independence graph of healthy subjects is sparser with only 87 edges compared to 115 edges in the conditional independence graph of Kawasaki subjects. Moreover, these two graphs share 65 edges.

#### IV. CONCLUSION

In this paper we developed some universal coders based on graph statistics for the compression of unlabeled graphs. We achieved this by first transforming graph structure into a rooted binary tree and showing that the compression of graph structure is equivalent to the compression of its associated rooted binary tree. Then, we used graph statistics to develop two main classes of graph coders. The first class uses graph motifs as local properties and the second class utilizes degree distribution as a global statistics along with graph motifs for coding. We introduced an application of graph coding for model selection in Gaussian graphical models.

For future work, we will extend this work in two directions. First is to extend graph compression to undirected graphs with attributes where graph topology will be utilized for the compression of structure and attributes. Second is to extend graph coding to directed graphs and using developed coders for data analysis purposes. Similar to undirected graphs, the idea is to transform directed graph into a rooted tree and encode nodal values on the rooted tree.

Type	$p$	$N$	Optimum	Benchmark methods			IID	Class 1			Class 2			
				CV	BIC	EBIC		Triangle	ComNei	4-node	Degree	Triangle	ComNei	4-node
Cycle	100	200	1	0.26	0.64	0.75	1	0.99	0.99	1	1	1	1	1
Cycle	200	400	1	0.23	0.61	0.69	1	1	1	1	1	1	1	1
AR(1)	100	200	0.99	0.26	0.65	0.75	0.65	0.99	0.99	0.99	0.99	0.99	0.99	0.99
AR(1)	200	400	1	0.23	0.62	0.70	0.99	0.99	0.99	1	1	1	1	1
ER	100	200	0.71	0.28	0.63	0.40	0.68	0.69	0.70	0.70	0.70	0.70	0.70	0.70
ER	200	400	0.80	0.27	0.68	0.75	0.79	0.78	0.78	0.79	0.79	0.79	0.79	0.79
Hub	100	200	0.50	0.23	0.48	0.01	0.47	0.48	0.48	0.48	0.49	0.49	0.49	0.49
Hub	200	400	0.44	0.22	0.43	0.16	0.43	0.43	0.43	0.43	0.44	0.44	0.44	0.44

Table II: F1-score of conditional independence graph recovery using different model selection methods. Results are the average over 50 replications when  $N > p$ .

Type	$p$	$N$	Optimum	Benchmark methods			IID	Class 1			Class 2			
				CV	BIC	EBIC		Triangle	ComNei	4-node	Degree	Triangle	ComNei	4-node
Cycle	100	50	0.89	0.28	0.26	0.74	0.81	0.89	0.89	0.89	0.89	0.89	0.89	0.89
Cycle	200	100	0.97	0.23	0.29	0.69	0.76	0.96	0.96	0.96	0.97	0.96	0.96	0.96
AR(1)	100	50	0.89	0.27	0.25	0.74	0.81	0.89	0.89	0.89	0.89	0.89	0.89	0.89
AR(1)	200	100	0.97	0.23	0.28	0.69	0.73	0.95	0.95	0.97	0.97	0.96	0.96	0.97
ER	100	50	0.54	0.28	0.26	0.40	0.47	0.49	0.49	0.49	0.36	0.46	0.46	0.45
ER	200	100	0.66	0.29	0.34	0.63	0.63	0.64	0.64	0.64	0.58	0.63	0.63	0.63
Hub	100	50	0.32	0.18	0.17	0.02	0.25	0.15	0.15	0.13	0.22	0.25	0.24	0.24
Hub	200	100	0.31	0.15	0.20	0.03	0.25	0.13	0.13	0.13	0.26	0.27	0.27	0.27

Table III: F1-score of conditional independence graph recovery using different model selection methods. Results are the average over 50 replications when  $N < p$ .

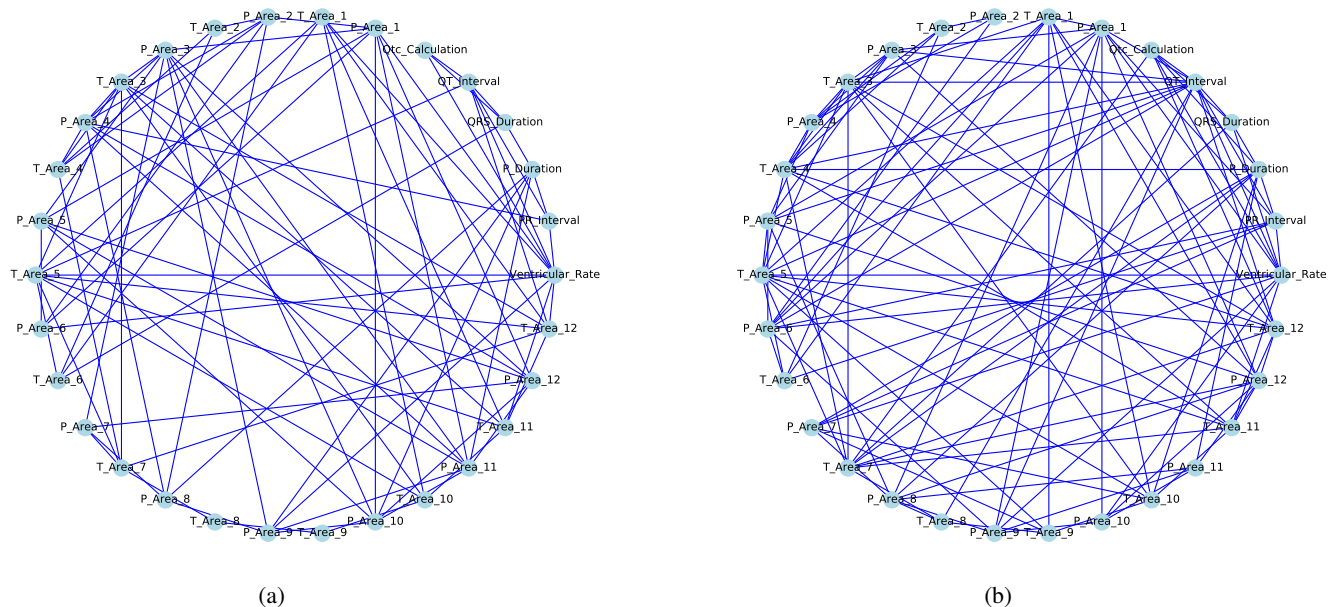


Figure 6: Conditional independence graph of (a) healthy subjects, (b) Kawasaki subjects.

## REFERENCES

- [1] Yongwook Choi and Wojciech Szpankowski. Compression of graphical structures: Fundamental limits, algorithms, and experiments. *IEEE Transactions on Information Theory*, 58(2):620–638, 2012.
- [2] C. Steinruecken. Compressing sets and multisets of sequences. *IEEE Transactions on Information Theory*, 61(3):1485–1490, March 2015.
- [3] Jorma Rissanen. A universal prior for integers and estimation by minimum description length. *The Annals of statistics*, 11(2):416–431, 1983.
- [4] Jorma Rissanen. Modeling by shortest data description. *Automatica*, pages 465–471, 1978.
- [5] Jorma Rissanen. Stochastic complexity and modeling. *The annals of statistics*, pages 1080–1100, 1986.

- [6] Tomasz Luczak, Abram Magner, and Wojciech Szpankowski. Structural information and compression of scale-free graphs. on internet, 2017.
- [7] T. Luczak, A. Magner, and W. Szpankowski. Asymmetry and structural information in preferential attachment graphs. *ArXiv e-prints 1607.04102*, July 2016.
- [8] Amir R. Asadi, Emmanuel Abbe, and Sergio Verdu. Compressing data on graphs with clusters. In *2017 IEEE International Symposium on Information Theory (ISIT)*, pages 1583–1587, June 2017.
- [9] Payam Delgosha and Venkat Anantharam. Universal lossless compression of graphical data. *IEEE Transactions on Information Theory*, 66(11):6962–6976, 2020.
- [10] Anders Host-Madsen and June Zhang. Coding of graphs with application to graph anomaly detection. In *2018 IEEE International Symposium on Information Theory (ISIT)*, pages 1829–1833. IEEE, 2018.
- [11] Tomasz Luczak, Abram Magner, and Wojciech Szpankowski. Compression of preferential attachment graphs. In *2019 IEEE International Symposium on Information Theory (ISIT)*, pages 1697–1701. IEEE, 2019.
- [12] Krzysztof Turowski, Abram Magner, and Wojciech Szpankowski. Compression of dynamic graphs generated by a duplication model. *Algorithmica*, 82(9):2687–2707, 2020.
- [13] Sourya Basu and Lav R Varshney. Universal and succinct source coding of deep neural networks. *arXiv preprint arXiv:1804.02800*, 2018.
- [14] Payam Delgosha and Venkat Anantharam. A universal low complexity compression algorithm for sparse marked graphs. In *2020 IEEE International Symposium on Information Theory (ISIT)*, pages 2349–2354. IEEE, 2020.
- [15] Maciej Besta and Torsten Hoefler. Survey and taxonomy of lossless graph compression and space-efficient graph representations. *arXiv preprint arXiv:1806.01799*, 2018.
- [16] Ming Yuan and Yi Lin. Model selection and estimation in the gaussian graphical model. *Biometrika*, 94(1):19–35, 2007.
- [17] Patricia Menéndez, Yiannis AI Kourmpetis, Cajo JF ter Braak, and Fred A van Eeuwijk. Gene regulatory networks from multifactorial perturbations using graphical lasso: application to the dream4 challenge. *PLoS one*, 5(12):e14147, 2010.
- [18] Rina Foygel and Mathias Drton. Extended bayesian information criteria for gaussian graphical models. In *Advances in neural information processing systems*, pages 604–612, 2010.
- [19] Mojtaba Abolfazli, Anders Host-Madsen, June Zhang, and Andras Bratinscak. Graph coding for model selection and anomaly detection in gaussian graphical models. *arXiv preprint arXiv:2102.02431*, 2021.
- [20] Adam J Rothman, Peter J Bickel, Elizaveta Levina, Ji Zhu, et al. Sparse permutation invariant covariance estimation. *Electronic Journal of Statistics*, 2:494–515, 2008.
- [21] Heng Lian. Shrinkage tuning parameter selection in precision matrices estimation. *Journal of Statistical Planning and Inference*, 141(8):2839–2848, 2011.
- [22] Shaun Lysen. Permuted inclusion criterion: a variable selection technique. *Publicly accessible Penn Dissertations*, page 28, 2009.
- [23] Han Liu, Kathryn Roeder, and Larry Wasserman. Stability approach to regularization selection (stars) for high dimensional graphical models. In *Advances in neural information processing systems*, pages 1432–1440, 2010.
- [24] Béla Bollobás. *Random Graphs*. Cambridge Studies in Advanced Mathematics. Cambridge University Press, 2 edition, 2001.
- [25] Réka Albert and Albert-László Barabási. Statistical mechanics of complex networks. *Reviews of modern physics*, 74(1):47, 2002.
- [26] Albert-László Barabási. *Network science*. Cambridge university press, 2016.
- [27] Lin Yao, Luning Wang, Lv Pan, and Kai Yao. Link prediction based on common-neighbors for dynamic social network. *Procedia Computer Science*, 83:82–89, 2016.
- [28] Shibao Li, Junwei Huang, Zhigang Zhang, Jianhang Liu, Tingpei Huang, and Haihua Chen. Similarity-based future common neighbors model for link prediction in complex networks. *Scientific reports*, 8(1):1–11, 2018.
- [29] Yuan H Wang. On the number of successes in independent trials. *Statistica Sinica*, pages 295–312, 1993.
- [30] Xiang-Hui Chen, Arthur P Dempster, and Jun S Liu. Weighted finite population sampling to maximize entropy. *Biometrika*, 81(3):457–469, 1994.
- [31] R. E. Barlow and K. D. Heidtmann. Computing k-out-of-n system reliability. *IEEE Transactions on Reliability*, R-33(4):322–323, 1984.
- [32] Peter Harremoës. Binomial and poisson distributions as maximum entropy distributions. *IEEE Transactions on Information Theory*, 47(5):2039–2041, 2001.
- [33] T.M. Cover and J.A. Thomas. *Information Theory, 2nd Edition*. John Wiley, 2006.
- [34] R. Krichevsky and V. Trofimov. The performance of universal encoding. *IEEE Transactions on Information Theory*, 27(2):199–207, Mar 1981.
- [35] F. M. J. Willems, Y.M. Shtarkov, and T.J. Tjalkens. The context-tree weighting method: basic properties. *Information Theory, IEEE Transactions on*, 41(3):653–664, 1995.
- [36] Ryan A. Rossi and Nesreen K. Ahmed. The network data repository with interactive graph analytics and visualization. In *AAAI*, 2015.
- [37] Jure Leskovec, Jon Kleinberg, and Christos Faloutsos. Graph evolution: densification and shrinking diameters. *ACM transactions on Knowledge Discovery from Data (TKDD)*, 1(1):2–es, 2007.
- [38] Benedek Rozemberczki, Ryan Davies, Rik Sarkar, and Charles Sutton. Gemsec: Graph embedding with self clustering. In *Proceedings of the 2019 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining 2019*, pages 65–72. ACM, 2019.
- [39] Jure Leskovec, Jon Kleinberg, and Christos Faloutsos. Graphs over time: densification laws, shrinking diameters and possible explanations. In *Proceedings of the eleventh ACM SIGKDD international conference on Knowledge discovery in data mining*, pages 177–187, 2005.
- [40] P. A. J. Volf and F. M. J. Willems. Switching between two universal source coding algorithms. In *Data Compression Conference, 1998. DCC '98. Proceedings*, pages 491–500, 1998.
- [41] Jerome Friedman, Trevor Hastie, and Robert Tibshirani. Sparse inverse covariance estimation with the graphical lasso. *Biostatistics*, 9(3):432–441, 2008.
- [42] Stanley J Kon. Models of stock returns—a comparison. *The Journal of Finance*, 39(1):147–165, 1984.
- [43] Vasyl Golosnoy, Bastian Gribisch, and Roman Liesenfeld. The conditional autoregressive wishart model for multivariate stock market volatility. *Journal of Econometrics*, 167(1):211–223, 2012.
- [44] Carlo Baldassi, Marco Zamparo, Christoph Feinauer, Andrea Procaccini, Riccardo Zecchina, Martin Weigt, and Andrea Pagnani. Fast and accurate multivariate gaussian modeling of protein families: predicting residue contacts and protein-interaction partners. *PLoS one*, 9(3):e92721, 2014.
- [45] Gaël Varoquaux, Alexandre Gramfort, Jean-Baptiste Poline, and Bertrand Thirion. Brain covariance selection: better individual functional connectivity models using population prior. *Advances in neural information processing systems*, 23:2334–2342, 2010.
- [46] Onureena Banerjee, Laurent El Ghaoui, and Alexandre d’Aspremont. Model selection through sparse maximum likelihood estimation for multivariate gaussian or binary data. *Journal of Machine learning research*, 9(Mar):485–516, 2008.
- [47] Robert Grone, Charles R Johnson, Eduardo M Sá, and Henry Wolkowicz. Positive definite completions of partial hermitian matrices. *Linear algebra and its applications*, 58:109–124, 1984.
- [48] Arthur P Dempster. Covariance selection. *Biometrics*, pages 157–175, 1972.
- [49] Lu Li and Kim-Chuan Toh. An inexact interior point method for  $l_1$ -regularized sparse covariance selection. *Mathematical Programming Computation*, 2(3-4):291–315, 2010.
- [50] Kean Ming Tan, Palma London, Karthik Mohan, Su-In Lee, Maryam Fazel, and Daniela Witten. Learning graphical models with hubs. *arXiv preprint arXiv:1402.7349*, 2014.