

# LEARNING GRAPHS AND SIMPLICIAL COMPLEXES FROM DATA

Andrei Buciulea<sup>†</sup>    Elvin Isufi<sup>‡</sup>    Geert Leus<sup>‡</sup>    Antonio G. Marques<sup>†</sup>

<sup>†</sup>Dept. of Signal Theory and Communications, King Juan Carlos University, Madrid, Spain

<sup>‡</sup>Delft University of Technology, Delft, The Netherlands

## ABSTRACT

Graphs are widely used to represent complex information and signal domains with irregular support. Typically, the underlying graph topology is unknown and must be estimated from the available data. Common approaches assume pairwise node interactions and infer the graph topology based on this premise. In contrast, our novel method not only unveils the graph topology but also identifies three-node interactions, referred to in the literature as second-order simplicial complexes (SCs). We model signals using a graph autoregressive Volterra framework, enhancing it with structured graph Volterra kernels to learn SCs. We propose a mathematical formulation for graph and SC inference, solving it through convex optimization involving group norms and mask matrices. Experimental results on synthetic and real-world data showcase a superior performance for our approach compared to existing methods.

**Index Terms**— Graph learning, simplicial complexes, higher-order networks, graph signal processing, Volterra graph models.

## 1. INTRODUCTION

Estimating the topology of complex data is a crucial step in the downstream signal processing and machine learning tasks [1, 2]. To estimate this structure, it is essential to model the coupling between the topology and the data and how they influence each other. For example, graph topology inference methods assume that pairwise node-to-node interactions could explain the data behavior or their dynamics. These approaches rely on algebraic and statistical methods to infer the graph topology from the observed data. Classic examples include correlation-based methods [3, Ch. 7.3.1], graphical lasso (GL) [4], and GSP based models, which exploit signal properties such as smoothness or graph stationarity [5–8].

Although pairwise interactions reveal some of the intricate dependencies and dynamics inherent in networked data, many interactions within groups comprise more than two nodes [9, 10]. For example, research collaborations often involve teams of authors and molecules tend to interact in small groups rather than pairs. To model such group interactions, common approaches resort to hypergraphs [11, 12] or simplicial complexes [13–16]. The latter are typically either considered as given or estimated via simple domain-specific heuristics. The work in [17] estimates hypergraphs from data by assuming a smoothness behavior on node and edge signals. It first infers a graph topology and then constructs on it a line graph to retrieve the higher-order interactions (hyperedges) but does not directly reveal the latter. This underscores the need for a model that is

capable of learning jointly the graph structure and capturing higher-order relationships between nodes. To account for the latter, we resort to Volterra models on networks, which model node dynamics in a nonlinear manner involving both pair-wise and higher-order interactions [18–20].

Specifically, we consider a networked autoregressive Volterra model and pose an inverse problem to *jointly* estimate the graph and the higher-order connectivity only from node signal realizations. Higher-order connectivities are embedded in the Volterra kernels of such a model. To limit the degrees of freedom (DoFs), we impose an SC structure between node-to-tuple interactions, ultimately, establishing a relationship between SCs and Volterra kernels. For the particular case of an SC of order two (representing connectivities up to triplets), this relationship simplifies to relating the Volterra kernel of order one to graph edges (simplex of order one) and the Volterra kernel of order two to filled triangles (simplex of order two). Subsequently, we develop a convex formulation that incorporates group sparsity to solve the proposed problem. The group sparsity allows us to group edges and node-to-tuple interactions related to each triplet of nodes and to control the number of triplet interactions. The proposed approach is corroborated via numerical experiments on synthetic and real data showing a competitive performance compared with alternatives that estimate either the graph topology or rely on it to infer higher-order interactions.

## 2. PROBLEM FORMULATION

Consider  $R$  observations of a vector  $\mathbf{x} \in \mathbb{R}^N$  grouped in matrix  $\mathbf{X} := [\mathbf{x}_1, \dots, \mathbf{x}_R] \in \mathbb{R}^{N \times R}$ . The data entries in  $\mathbf{x}$  have a hidden underlying structure, which is typically represented through a graph  $\mathcal{G} = \{\mathcal{V}, \mathcal{E}\}$  comprising a set of nodes  $\mathcal{V} = \{1, \dots, N\}$  and a set of edges  $\mathcal{E} = \{(i, j) | i, j \in \mathcal{V}\}$ . Here, the  $i$ th entry  $x_i$  is seen as a datum associated with node  $i$ , hence, set  $\mathcal{E}$  represents pairwise dependencies between the data of nodes  $i$  and  $j$ . In this context, we also refer to vector  $\mathbf{x}$  as the graph signal. Estimating the graph topology from the data  $\mathbf{X}$  boils down to leveraging a model that expresses this data coupling and the role of the topology in it; i.e., solving an inverse problem of the form  $\mathcal{G} = f^{-1}(\mathbf{X})$  where  $f(\cdot)$  is a function acting upon the graph  $\mathcal{G}$  and modeling how it is coupled with the signal realizations in  $\mathbf{X}$ . Typically, we will estimate an algebraic representation of graph  $\mathcal{G}$  that is represented by its adjacency matrix  $\mathbf{A} \in \mathbb{R}^{N \times N}$ , graph Laplacian  $\mathbf{L} := \text{diag}(\mathbf{A}\mathbf{1}) - \mathbf{A}$ , or more generally a graph shift operator (GSO) matrix  $\mathbf{S} \in \mathbb{R}^{N \times N}$ , where  $S_{ij} \neq 0$  if and only if  $i = j$  or  $(j, i) \in \mathcal{E}$  [21]. The existing literature provides different approaches to estimating the graph from the data, with different assumptions on the function  $f(\cdot)$  and topological assumptions such as directed, weighted or undirected edges on the sought graph. Moreover, since their goal is to recover a graph they focus on pairwise interactions.

We consider that the graph signal is influenced by the signal values in the other nodes both via pairwise and higher-order interac-

Work supported by the EU H2020 Grant Tailor (No 952215); the Spanish NSF Grants (MCIN/AEI/10.13039/501100011033) Grants PID2019-105032GB-I00 and PID2022-136887NB-I00; the Community of Madrid (CAM) and Rey Juan Carlos University (URJC) via the Young Researchers R&D Project ref. F861 (AUTO-BA-GRAPH), and the fellowship PREDOC20-003. Email contact author: antonio.garcia.marques@urjc.es.

tions. Specifically, we model the dependencies via an autoregressive graph Volterra model of second order of the form

$$\begin{aligned} \mathbf{X} &= \mathbf{H}_1 \mathbf{X} + \mathbf{H}_2 \mathbf{Y} + \mathbf{V} + \mathbf{E}, & (1a) \\ \text{with } \mathbf{Y} &= \mathbf{X} \odot \mathbf{X}. & (1b) \end{aligned}$$

Here,  $\mathbf{H}_1 \in \mathbb{R}^{N \times N}$  represents the pairwise interactions and the term  $\mathbf{H}_1 \mathbf{X}$  captures the part of the graph signal that can be represented as a linear combination of the signals in the other nodes. Instead, matrix  $\mathbf{H}_2 \in \mathbb{R}^{N \times N^2}$  is a node-to-tuple interaction matrix representing higher-order interactions between a node  $k$  and a tuple  $(i, j)$  in its entry  $\mathbf{H}_2[k, (i, j)]$ <sup>1</sup>. Matrix  $\mathbf{Y} = \mathbf{X} \odot \mathbf{X} \in \mathbb{R}^{N^2 \times R}$  is obtained by performing the Khatri-Rao product (column-wise Kronecker product) on the graph signals. The  $r$ -th column of  $\mathbf{Y}$  collects all the monomials of degree two involving variables  $\{x_r^i\}_{i=1}^N$ . These product signals can be seen now as values associated with tuples of nodes. Hence, the  $r$ -th column of  $\mathbf{H}_2 \mathbf{Y}$  captures the part of a graph signal ( $\mathbf{x}_r$ ) that can be represented by  $\mathbf{x}_r \odot \mathbf{x}_r$  via node-to-pair interactions ( $\mathbf{H}_2$ ). Finally,  $\mathbf{V} \in \mathbb{R}^{N \times R}$  is an exogenous variable and  $\mathbf{E} \in \mathbb{R}^{N \times R}$  is white zero-mean noise. For didactical purposes, we focus on a single node  $k$  and model (1) allows writing its signal as<sup>2</sup>

$$x_k = \sum_{j=1, j \neq k}^N \mathbf{H}_1[k, j] x_j + \sum_{i=1, i \neq k}^N \sum_{j=1, j \neq i, k}^N \mathbf{H}_2[k, (i-1)N + j] x_i x_j + v_k + e_k,$$

where the first and second summations are reminiscent of  $\mathbf{H}_1 \mathbf{X}$  and  $\mathbf{H}_2 \mathbf{Y}$ , respectively. The latter are linearly combined by the node-to-tuple weights  $\mathbf{H}_2[k, (i, j)]$  here expressed as matrix entries  $\mathbf{H}_2[k, (i-1)N + j]$ . Finally,  $v_k$  and  $e_k$  are the exogenous variable and noise at node  $k$ , respectively. The following example makes this discussion more tangible.

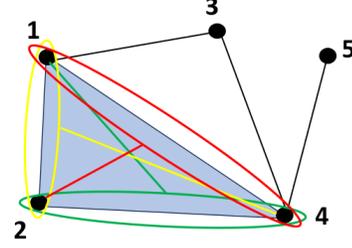
**Example.** Fig. 1 illustrates pairwise and higher-order interactions among five nodes (1, . . . , 5). The pairwise interactions (e.g., (1, 2), (2, 4), etc.) are shown by black solid lines. Matrix  $\mathbf{H}_1$  collects them. The node signals  $x_1, \dots, x_5$  could be seen as values over the respective nodes. There are three tuples in this figure highlighted by ellipsoids; i.e., (1, 2) in yellow, (1, 4) in red, and (2, 4) in green. The node-to-pair interactions are shown by solid lines connecting the node to the respective tuple with the same color; i.e., [1, (2, 4)], [2, (1, 4)] and [4, (1, 2)]. Matrix  $\mathbf{H}_2$  collects all these interactions. The product node signals  $x_1 x_2$ ,  $x_1 x_4$ , and  $x_2 x_4$  could be seen as values associated with tuples captured by matrix  $\mathbf{Y}$  in (1). Following the model (1), the signal at node  $k = 2$ , can be written as

$$x_2 = \mathbf{H}_1[2, 1]x_1 + \mathbf{H}_1[2, 4]x_4 + \mathbf{H}_2[2, (1, 4)]x_1 x_4 + \mathbf{H}_2[2, (4, 1)]x_1 x_4 + v_2 + e_2. \quad (2)$$

Model (1) has two particular aspects worth discussing. First, it relates the graph signals to both the pairwise and higher-order connectivities in a nonlinear manner in  $\mathbf{X}$  but still linear in the topological variables  $\mathbf{H}_1$  and  $\mathbf{H}_2$ . This is reminiscent of how classical time-series Volterra models [22] have been extended to graph signals and proven relevant in applications such as power distribution grids, social networks, and recommender systems [18, 19], to name a few. The Volterra models have been found to provide both expressibility for higher-order interactions, as well as interpretability for further understanding of the underlying network dynamics. Second, (1) is a rather flexible backbone model that can be further enriched via more expressive kernels. For example, we could consider a higher-order

<sup>1</sup>With a slight abuse of notation, we will alternate between  $\mathbf{H}_2[k, (i, j)]$  and  $\mathbf{H}_2[k, (i-1)N + j]$  to denote the value of the  $k$ -th row and  $((i-1)N + j)$ -th column of an  $N \times N^2$  matrix.

<sup>2</sup>To ease exposition, we assume  $R = 1$  and drop the realization index.



**Fig. 1:** The visual representation of a graph with 5 nodes and 6 edges. The node-to-pair interactions between nodes [1,2,4] and edges [(2,4),(1,4),(1,2)] are represented by green, red, and yellow lines, respectively. A filled triangle between nodes 1, 2 and 4 is represented in blue.

Volterra model (as opposed to a second-order here) to capture node-to-triple interactions. Also, we considered monomials of order two as signals over tuples in  $\mathbf{y}$  but *non-multiplicative* higher-order interactions between node signals might also be of interest. From this perspective, we could define  $\mathbf{Y} = [\mathbf{y}_1, \dots, \mathbf{y}_R]$  as  $\mathbf{y}_r = g(\mathbf{x}_r, \mathbf{x}_r)$ , where  $g(\cdot)$  is a general function for which the following equality holds  $g(\mathbf{a}, \mathbf{b}) = g(\mathbf{b}, \mathbf{a})$ .

Another advantage of the model in (1) is that it allows imposing straightforwardly a desired structure on both  $\mathbf{H}_1$  and  $\mathbf{H}_2$ . This is particularly important if we want to estimate these matrices by a limited number of signal realizations. Particular structures include:

- Positive weights:* This could be achieved by imposing the element-wise constraints  $\mathbf{H}_1 \geq \mathbf{0}$  and  $\mathbf{H}_2 \geq \mathbf{0}$ .
- No self-loops:* This is particularly important to avoid trivial solutions, such as  $\mathbf{H}_1 = \mathbf{I}$  (2). We could define two binary masking matrices  $\mathbf{B}_1 \in \{0, 1\}^{N \times N}$  and  $\mathbf{B}_2 \in \{0, 1\}^{N \times N^2}$  and impose the elementwise constraints  $\mathbf{B}_1 \circ \mathbf{H}_1 = \mathbf{0}$  and  $\mathbf{B}_2 \circ \mathbf{H}_2 = \mathbf{0}$ , where  $\circ$  denotes the entry-wise (Hadamard) multiplication. For example,  $\mathbf{B}_1 = \mathbf{I}$  implies that  $\text{diag}(\mathbf{H}_1) = \mathbf{0}$  and, as a result, removes pairwise self-loops.
- Symmetry:* This could be guaranteed by forcing  $\mathbf{H}_1 = \mathbf{H}_1^T$ .

However, this structure is mild and can still have prohibitive DoFs or lead to learned structures  $\mathbf{H}_1$  and  $\mathbf{H}_2$  that are too disconnected from each other. Therefore, *our goal is to leverage model (1) and use nodal realizations  $\mathbf{X}$  to jointly estimate the pairwise graph structure  $\mathbf{H}_1$  and the higher-order connectivities  $\mathbf{H}_2$  by imposing an SC structure in  $\mathbf{H}_2$  that reduces the DoFs and couples the learned  $\mathbf{H}_1$  with  $\mathbf{H}_2$ .*

### 3. GRAPH AND SIMPLICIAL COMPLEX LEARNING

Since our goal is to learn an SC of order 2 (which can be understood as learning a graph and filled triangles), we provide a brief overview of the SC concept directly from a geometric perspective. A graph  $\mathcal{G} = \{\mathcal{V}, \mathcal{E}\}$ , with  $\mathcal{E} \subseteq \mathcal{V} \times \mathcal{V}$ , is an SC of order 1. An SC of order 2 can be represented by  $\{\mathcal{V}, \mathcal{E}, \mathcal{T}\}$ , where  $\mathcal{T} \subseteq \mathcal{V} \times \mathcal{V} \times \mathcal{V}$  and a triplet  $(i, j, k)$  can belong to  $\mathcal{T}$  only if all pairs  $(i, j)$ ,  $(j, k)$  and  $(i, k)$  belong to  $\mathcal{E}$ . This readily implies that an SC of order 2 can be represented by graph  $\mathcal{G} = \{\mathcal{V}, \mathcal{E}\}$  together with a list of filled triangles. Similarly, for an  $n$ -th order SC ( $n$ -simplex) to exist, the presence of all  $(n-1)$ -th order SCs (also known as  $(n-1)$ -simplices) is required. This implies that higher-order interactions rely directly on the connectivity between the nodes in the graph.

To identify 2-simplices using Volterra kernels, it is essential to represent interactions among three interconnected nodes. To achieve this representation, we use the matrix  $\mathbf{H}_2$ , which captures the influence of the product of the signals at two nodes on a third node. It

can be seen that, indeed, the entry  $[k, (i, j)]$  of  $\mathbf{H}_2$  involves three nodes. Conversely, the relation between nodes  $k, i,$  and  $j$  appears in six elements of  $\mathbf{H}_2$ . In the context of SCs, our modeling assumption is that the six entries of  $\mathbf{H}_2$  associated with the triplet  $(k, i, j)$  can be different from zero only if the triangle  $(k, i, j)$  is filled.

**Example. (cont.)** Fig. 1 illustrates the correspondence between a filled triangle and node-to-pair interactions among three nodes. Nodes  $(3, 4, 5)$  do not form a triangle (not all edges are present), and hence, the filled-triangle relation cannot exist. Differently, for the triplets  $(1, 3, 4)$  and  $(1, 2, 4)$  the triangle exists. Our assumption is that the triangle is filled if the node-to-edge interactions involving the three nodes (denoted as an ellipse and a straight line) exist. This is the case for nodes  $(1, 2, 4)$  and therefore the filled triangle exists. Conversely, this condition does not hold for the triangle  $(1, 3, 4)$ , so it remains unfilled.

Based on the previous discussion, with  $\mathbf{X}[i, r]$  denoting the  $(i, r)$ -th entry of matrix  $\mathbf{X}$ , and recalling that  $\mathbf{Y} = \mathbf{X} \circ \mathbf{X}$  [cf. (1b)], our approach to identifying an SC of order 2 by using an autoregressive Volterra model can be formulated as

$$\begin{aligned} (\hat{\mathbf{H}}_1, \hat{\mathbf{H}}_2) &= \underset{\mathbf{H}_1 \in \mathcal{H}_1, \mathbf{H}_2 \in \mathcal{H}_2}{\operatorname{argmin}} \|\mathbf{X} - \mathbf{H}_1 \mathbf{X} - \mathbf{H}_2 \mathbf{Y} - \mathbf{V}\|_F^2 \\ &\quad + \alpha \|\mathbf{H}_1\|_1 + \beta \|\mathbf{H}_2\|_1 \quad (3a) \\ \text{s. t.} \quad &\mathbf{H}_2[k, (i, j)] \leq \theta \mathbb{1}(\mathbf{H}_1[k, i] \mathbf{H}_1[k, j] \mathbf{H}_1[i, j]); \quad (3b) \end{aligned}$$

where the second and third terms in (3a) (with  $\alpha > 0$  and  $\beta > 0$  being hyperparameters) regulate the desired sparsity level in  $\mathbf{H}_1$  and  $\mathbf{H}_2$ , respectively. By employing  $\mathcal{H}_1 = \{\mathbf{H}_1 \geq \mathbf{0}, \mathbf{B}_1 \circ \mathbf{H}_1 = \mathbf{0}\}$  and  $\mathcal{H}_2 = \{\mathbf{H}_2 \geq \mathbf{0}, \mathbf{B}_2 \circ \mathbf{H}_2 = \mathbf{0}\}$  we adopt the structural requirements outlined at the end of Sec. 2 which impose that both  $\mathbf{H}_1$  and  $\mathbf{H}_2$  must possess positive weights and no self-loops. The indicator function from the constraint is defined as  $\mathbb{1}(z) = 1$  if  $z \neq 0$  and  $\mathbb{1}(z) = 0$  if  $z = 0$ . By setting  $z = \mathbf{H}_1[k, i] \mathbf{H}_1[k, j] \mathbf{H}_1[i, j]$ , we have that if nodes  $k, i$  and  $j$  form a triangle ( $z \neq 0$ ) then the nonlinear relation captured by the Volterra kernel  $\mathbf{H}_2$  could exist ( $\mathbf{H}_2[k, (i, j)] \leq \theta \mathbb{1}(z)$ ). If needed, the parameter  $\theta$  can be selected to limit the maximum value attributed to each node-pair interaction. Note that having one entry of  $\mathbf{H}_2$  different from zero implies that the associated triangle is filled.

**Remark 1** Consider that when a triangle exists and is filled (say triangle  $(k, i, j)$ ), the formulation in (3) does not impose that the values of all node-pair interactions between the nodes present in a filled triangle are the same. If this is required, it can be accomplished by adding the set of constraints  $\mathbf{H}_2[k, (i, j)] = \mathbf{H}_2[k, (j, i)] = \mathbf{H}_2[i, (k, j)] = \mathbf{H}_2[i, (j, k)] = \mathbf{H}_2[j, (k, i)] = \mathbf{H}_2[j, (i, k)]$  for all  $(k, i, j)$ .

Problem (3) is non-convex due to the constraint (3b). To deal with it, we apply a group sparsity term that groups all the entries in  $\mathbf{H}_1$  and  $\mathbf{H}_2$  that participate in (3b). To that end, we build the  $N \times (N + N^2)$  matrix  $[\mathbf{H}_1, \mathbf{H}_2]$  and, since each constraint in (3b) involves 3 nodes, we construct the  $N \times (N + N^2)$  binary mask matrix  $\mathbf{Q}^{(i,j,k)}$ . This matrix identifies the entries of  $[\mathbf{H}_1, \mathbf{H}_2]$  associated with i) edges between the three nodes  $\mathbf{Q}^{(i,j,k)}[i, j] = 1, \mathbf{Q}^{(i,j,k)}[i, k] = 1, \mathbf{Q}^{(i,j,k)}[j, k] = 1$ , and ii) node-pair interactions between the three nodes  $\mathbf{Q}^{(i,j,k)}[i, Nj + k] = 1, \mathbf{Q}^{(i,j,k)}[j, Ni + k] = 1, \mathbf{Q}^{(i,j,k)}[k, Ni + j] = 1$ , with all other entries being zero. Leveraging  $\mathbf{Q}^{(i,j,k)}$ , we propose the following convex formulation

$$\begin{aligned} (\hat{\mathbf{H}}_1, \hat{\mathbf{H}}_2) &= \underset{\mathbf{H}_1 \in \mathcal{H}_1, \mathbf{H}_2 \in \mathcal{H}_2}{\operatorname{argmin}} \|\mathbf{X} - \mathbf{H}_1 \mathbf{X} - \mathbf{H}_2 \mathbf{Y} - \mathbf{V}\|_F^2 + \alpha \|\mathbf{H}_1\|_1 \\ &\quad + \beta \|\mathbf{H}_2\|_1 + \gamma \sum_{i,j,k=1}^N \|\mathbf{Q}^{(i,j,k)} \circ [\mathbf{H}_1, \mathbf{H}_2]\|_F, \quad (4) \end{aligned}$$

with  $\gamma > 0$ . In (4), we kept the least square term and the sparsity-promoting terms in (3a) but replaced (3b) with a group sparsity term. The group-sparsity regularizer links the penalty of activating an entry of  $\mathbf{H}_1$  with that of activating an entry of  $\mathbf{H}_2$  provided that those entries can participate in a potential triangle. When all the interactions (entries) inside the group sparsity norm hold non-zero values, it signifies the presence of a filled triangle.

Although convexity guarantees that problem (4) can be solved in polynomial time, the required computational complexity is not negligible, especially for medium/large size graphs. The number of terms in the group sparsity constraint scales as  $O(N^2)$ , and the number of variables and constraints scales as  $O(N^3)$ . While this multiplicative growth in the number of variables to optimize is somehow unavoidable when dealing with the design of nontrivial schemes to estimate *high-order interactions*, it emphasizes the importance of developing tailored optimization algorithms that, by exploiting the structure in (4), lead to a reduced complexity. We leave this task as future work.

#### 4. NUMERICAL EXPERIMENTS

We conduct numerical experiments on both synthetic and real data and compare the following methods.

- GL: Graphical Lasso [4], which learns the edges of a graph by estimating a sparse precision matrix from Gaussian graph signals.
- GSR: Approach in [7], which estimates the graph topology by assuming the signals are graph stationary in the sought graph.
- HGSL: Approach in [17], which estimates the graph by assuming smoothness on node (0-simplex) and edge (1-simplex) signals.
- RC: Rips complex [23], which estimates SCs (edges and filled triangles) from the correlation of the data.
- MTV-SC: Approach in [14], which estimates SCs from edge signals assuming the topology of the underlying graph is given.
- VGR: Our approach in (4) for estimating the graph and SC (edges and filled triangles) from data using a Volterra signal model.

The exact implementation details of the previous schemes and ensuing setups can be found in the online code repository<sup>3</sup>.

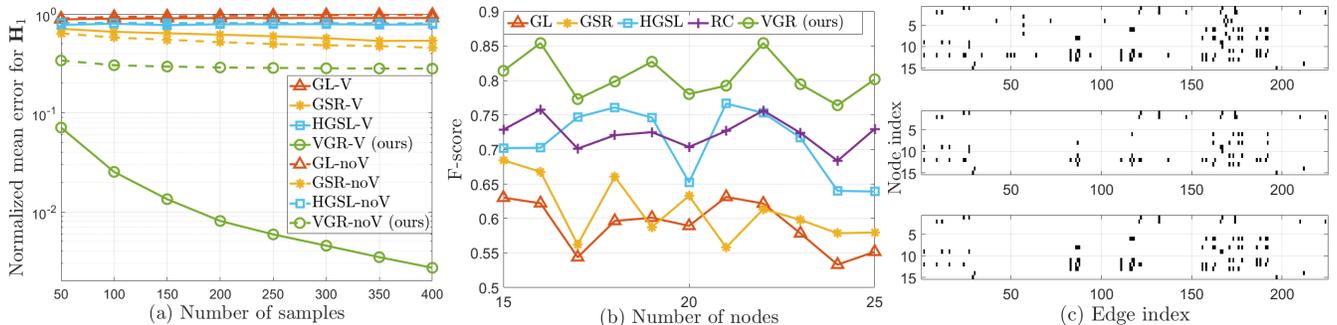
**Number of samples.** For this experiment, we assess the performance of VGR by using the following synthetic data generation setup. We generate  $M$  graph signals following the autoregressive Volterra model in (1a). The entries of  $\mathbf{H}_1$  are set as the adjacency matrix of an Erdős Rényi graph with  $N = 20$  nodes and edge probability  $p = 0.15$  [24]. The entries of  $\mathbf{H}_2$  are set so that all triangles are filled, which is a favorable setup for the RC algorithm. Fig. 2.a shows the estimation error for  $\mathbf{H}_1$  (y-axis) averaged over 50 graph realizations while increasing the number of observed signals  $R$  (x-axis). The metric used to compute the estimation error is the normalized squared Frobenius norm:

$$\operatorname{err}(\mathbf{H}_1) = \|\mathbf{H}_1^* - \hat{\mathbf{H}}_1\|_F^2 / \|\mathbf{H}_1^*\|_F^2 \quad (5)$$

where  $\mathbf{H}_1^*$  and  $\hat{\mathbf{H}}_1$  stand for the ground truth and estimated  $\mathbf{H}_1$  respectively. Solid lines in Fig. 2.a consider  $\mathbf{V}$  in (1a) known, whereas dashed lines consider  $\mathbf{V}$  unknown and set to zero.

Starting first with the solid lines, we observe that: i) our algorithm yields the best performance and ii) as  $R$  increases,  $\operatorname{err}(\mathbf{H}_1)$  decreases. The only exception to ii) is GL, probably because its modeling assumptions are too simple for the signal structure postulated in (1a). When the exogenous variable is unknown (dashed lines), the performance of all algorithms deteriorates. Our approach

<sup>3</sup><https://github.com/andreibuciulea/Graph-SCs-topoID>



**Fig. 2:** (a) Evaluating the estimation performance of different algorithms in terms of normalized squared Frobenius norm  $[\text{err}(\mathbf{H}_1)$ , cf. (5)] averaged over 50 graph realizations as the number of samples  $R$  increases. (b) Evaluating the estimation performance of different algorithms in terms of F-score averaged over 20 graph realizations. (c) Visual representation of the ground-truth support of  $\mathbf{H}_2$  (upper) and the estimation obtained by RC (middle) and VGR (bottom) for a real-data graph realization with  $N = 15$  nodes.

**Table 1:** Normalized error when estimating 2-simplices  $\text{err}(\mathbf{H}_2)$  for different algorithms and number of samples  $R$ .

Alg. \ $R$	50	100	200	300	400	500
<b>MTV-SC</b>	1.505	1.496	1.497	1.493	1.494	1.490
<b>RC</b>	0.790	0.767	0.761	0.753	0.748	0.751
<b>VGR</b>	0.559	0.428	0.294	0.214	0.165	0.133

outperforms the alternatives, due to the consideration of a more comprehensive model that accounts for signal nonlinearities and higher-order interactions. Results for RC were not shown because, in all tested cases, the normalized error was slightly above 1.0.

We now move our analysis to assess the performance when estimating  $\mathbf{H}_2$ . Since GL and GSR do not consider high-order interactions explicitly, a direct comparison with these methods is infeasible. Thus, we compare the proposed algorithm with MTV-SC and RC. The results in Table 1 reveal  $\text{err}(\mathbf{H}_2)$  is larger than  $\text{err}(\mathbf{H}_1)$  as illustrated in Fig. 2.a. This behaviour aligns with expectations, since the task of estimating higher-order interactions is inherently more difficult than estimating links between nodes. This difficulty is due to the large number of potential interactions to estimate compared to the number of available signals. Nevertheless, our approach not only achieves lower errors than considered alternatives but also exhibits a faster error reduction as the number of samples  $R$  increases. We attribute this enhanced  $\mathbf{H}_2$  estimate to the ability of our approach to simultaneously compute both SCs and the underlying edges in the graph instead of the two-step estimation process implemented by RC and MTV-SC.

**Co-authorship datasets.** We now evaluate the performance of VGR using a real dataset, following the setup in [17]. The dataset comprises papers from the ACM conference, featuring 17,431 authors, 122,499 papers, and 1,903 keywords. The nodes (0-simplices) are a subset of the authors. To establish the ground truth  $\mathbf{H}_1$ , we examined author-paper relationships and considered a link between two authors if they collaborated on a paper. For ground truth  $\mathbf{H}_2$ , we considered a filled triangle whenever three authors collaborated on a paper. To generate the input signals, we set  $R = 1,903$  (the total number of different keywords). As a result, the value of each input signal (columns of  $\mathbf{X}$ ) is related to the frequency at which a particular author uses a particular keyword across papers. We constructed 20 different graphs from the dataset, varying the set of authors (with cardinalities between 15 and 25), keeping the number of signals as  $R = 1,903$ . These generated signals were employed to estimate the graph topology using the different algorithms. Fig. 2.b displays the average results across 20 graph realizations, represented in terms of F-score (y-axis), with the number of nodes (authors) ranging from 15

**Table 2:** F-score and  $\text{err}(\mathbf{H}_2)$  when estimating 2-simplices from real-data for different algorithms and number of nodes  $N$ .

Alg. \ $N$	F-score			Error		
	15	20	25	15	20	25
<b>MTV-SC</b>	0.093	0.058	0.056	7.418	7.536	7.530
<b>RC</b>	0.667	0.650	0.585	1.350	2.101	2.837
<b>VGR</b>	0.718	0.676	0.625	0.548	0.558	0.649

to 25 (x-axis). The results indicate that approaches that do not consider higher-order interactions, such as GL and GSR, yield poorer graph estimations compared to other alternatives. VGR consistently outperforms all other schemes, achieving the highest averaged F-score across all considered graph sizes.

Regarding the estimation of the higher-order connectivities shown in Fig. 2.c, RC achieves an F-score of 0.77, whereas VGR achieves an F-score of 0.88 (recall that these are the only algorithms that explicitly account for 2-simplices/triplets). Fig.2.c. further shows the support of  $\mathbf{H}_2$  for both the RC and VGR alongside the ground truth. The recovered  $\mathbf{H}_2$  exhibits a similar structure compared to the ground truth. This implies that both RC and our algorithm effectively estimate the presence of filled triangles, representing interactions between three or more authors collaborating on the same paper. The proposed method also provides better estimations of the associated weights. In the specific realization shown in Fig. 2.c, our method achieves an  $\text{err}(\mathbf{H}_2)$  of 0.054, whereas RC yields an error of 1.334. Additional results showing the average F-score and  $\text{err}(\mathbf{H}_2)$  results for VGR and RC at different numbers of nodes are provided in Table 2. These results reinforce the conclusions drawn from Fig. 2.c. Lastly, we incorporated the results of SC estimation obtained by the MVT-SC approach. For the SC estimation, we used  $\mathbf{Y}$  as edge signals and the graph estimated by the approach presented in [25]. From the results shown in Table 2, we can conclude that relying solely on the node signals results in an inadequate estimation of SCs for MVT-SC, requiring having access to the actual edge signals.

## 5. CONCLUSIONS

This paper proposed a method to jointly estimate the graph topology (pairwise interactions) and higher-order dependencies (triples) from nodal data by assuming the latter follows a second-order autoregressive graph Volterra model. We incorporated simplicial complex constraints in the said model to estimate sparse-filled triangles as a proxy for triplet interactions. To assess the estimation performance of the proposed algorithm, we conducted experiments on both synthetic and real datasets, revealing consistently superior results compared to those achieved by alternative methods.

## 6. REFERENCES

- [1] G. Mateos, S. Segarra, A. G. Marques, and A. Ribeiro, “Connecting the dots: Identifying network structure via graph signal processing,” *IEEE Signal Process. Mag.*, vol. 36, no. 3, pp. 16–43, 2019.
- [2] X. Dong, D. Thanou, M. Rabbat, and P. Frossard, “Learning graphs from data: A signal representation perspective,” *IEEE Signal Process. Mag.*, vol. 36, no. 3, pp. 44–63, 2019.
- [3] E. D. Kolaczyk, *Statistical Analysis of Network Data: Methods and Models*, Springer, New York, NY, 2009.
- [4] J. Friedman, T. Hastie, and R. Tibshirani, “Sparse inverse covariance estimation with the graphical lasso,” *Biostatistics*, vol. 9, no. 3, pp. 432–441, 2008.
- [5] V. Kalofolias, “How to learn a graph from smooth signals,” in *Intl. Conf. Artif. Intel. Statist. (AISTATS)*. J Mach. Learn. Res., 2016, pp. 920–929.
- [6] H. E. Egilmez, E. Pavez, and A. Ortega, “Graph learning from data under laplacian and structural constraints,” *IEEE J. Sel. Topics Signal Process.*, vol. 11, no. 6, pp. 825–841, 2017.
- [7] S. Segarra, A. G. Marques, G. Mateos, and A. Ribeiro, “Network topology inference from spectral templates,” *IEEE Trans. Signal Info. Process. Networks*, vol. 3, no. 3, pp. 467–483, Sep. 2017.
- [8] A. Buciualea, S. Rey, and A. G. Marques, “Learning graphs from smooth and graph-stationary signals with hidden variables,” *IEEE Trans. Signal Info. Process. Networks*, vol. 8, pp. 273–287, 2022.
- [9] R. Milo, S. Shen-Orr, S. Itzkovitz, N. Kashtan, D. Chklovskii, and U. Alon, “Network motifs: Simple building blocks of complex networks,” *Science*, vol. 298, no. 5594, pp. 824–827, 2002.
- [10] C. Bick, E. Gross, H. A Harrington, and M. T. Schaub, “What are higher-order networks?,” *SIAM Review*, vol. 65, no. 3, pp. 686–731, 2023.
- [11] C. Berge, *Hypergraphs: combinatorics of finite sets*, vol. 45, Elsevier, 1984.
- [12] J. G. Young, G. Petri, and T. P. Peixoto, “Hypergraph reconstruction from network data,” *Communications Physics*, vol. 4, no. 1, pp. 135, 2021.
- [13] L-H. Lim, “Hodge Laplacians on graphs,” *SIAM Review*, vol. 62, no. 3, pp. 685–715, 2020.
- [14] S. Barbarossa and S. Sardellitti, “Topological signal processing over simplicial complexes,” *IEEE Trans. Signal Process.*, vol. 68, pp. 2992–3007, 2020.
- [15] M. Yang, E. Isufi, M. T. Schaub, and G. Leus, “Simplicial convolutional filters,” *IEEE Trans. Signal Process.*, vol. 70, pp. 4633–4648, 2022.
- [16] M. T. Schaub, Y. Zhu, J-B. Seby, T. M. Roddenberry, and S. Segarra, “Signal processing on higher-order networks: Livin’ on the edge... and beyond,” *Signal Process.*, vol. 187, pp. 108149, 2021.
- [17] B. Tang, S. Chen, and X. Dong, “Learning hypergraphs from signals with dual smoothness prior,” in *ICASSP 2023 - 2023 IEEE Intl. Conf. Acoustics, Speech Signal Process. (ICASSP)*, 2023, pp. 1–5.
- [18] Q. Yang, M. Coutino, G. Leus, and G. B Giannakis, “Autoregressive graph volterra models and applications,” *EURASIP J. on Advances in Signal Process.*, vol. 2023, no. 1, pp. 1–21, 2023.
- [19] G. Leus, M. Yang, M. Coutino, and E. Isufi, “Topological volterra filters,” in *ICASSP 2021 - 2021 IEEE Intl. Conf. Acoustics, Speech Signal Process. (ICASSP)*. IEEE, 2021, pp. 5385–5399.
- [20] H-X. Li, C. Qi, and Y. Yu, “A spatio-temporal volterra modeling approach for a class of distributed industrial processes,” *J. of Process Control*, vol. 19, no. 7, pp. 1126–1142, 2009.
- [21] A. Ortega, P. Frossard, J. Kovacevic, J. M. F. Moura, and P. Vanderghenst, “Graph signal processing: Overview, challenges, and applications,” *Proc. IEEE*, vol. 106, no. 5, pp. 808–828, 2018.
- [22] R. W. Brockett, “Volterra series and geometric control theory,” *IFAC Proc. Volumes*, vol. 8, no. 1, Part 1, pp. 245–254, 1975, 6th IFAC World Congress (IFAC 1975) - Part 1: Theory, Boston/Cambridge, MA, USA, August 24–30, 1975.
- [23] A. Zomorodian, “Fast construction of the Vietoris-Rips complex,” *Computers & Graphics*, vol. 34, no. 3, pp. 263–271, 2010.
- [24] B. Bollobás, *Random Graphs*, Cambridge University Press, 2001.
- [25] S. P. Chepuri, S. Liu, G. Leus, and A. O. Hero, “Learning sparse graphs under smoothness prior,” in *2017 IEEE Intl. Conf. Acoustics, Speech Signal Process. (ICASSP)*. IEEE, 2017, pp. 6508–6512.