# arXiv:2307.14198v1 [physics.soc-ph] 26 Jul 2023

## Feature-enriched network geometry explains graph-structured data

Roya Aliakbarisani,<sup>1, 2, \*</sup> M. Ángeles Serrano,<sup>1, 2, 3, †</sup> and Marián Boguñá<sup>1, 2, ‡</sup>

<sup>1</sup>Departament de Física de la Matèria Condensada,

Universitat de Barcelona, Martí i Franquès 1, E-08028 Barcelona, Spain

<sup>2</sup> Universitat de Barcelona Institute of Complex Systems (UBICS), Barcelona, Spain

<sup>3</sup>Institució Catalana de Recerca i Estudis Avançats (ICREA),

Passeig Lluís Companys 23, E-08010 Barcelona, Spain

Graph-structured data provide a comprehensive description of complex systems, encompassing not only the interactions among nodes but also the intrinsic features that characterize these nodes. These features play a fundamental role in the formation of links within the network, making them valuable for extracting meaningful topological information. Notably, features are at the core of deep learning techniques such as Graph Convolutional Neural Networks (GCNs) and offer great utility in tasks like node classification, link prediction, and graph clustering. In this letter, we present a comprehensive framework that treats features as tangible entities and establishes a bipartite graph connecting nodes and features. By assuming that nodes sharing similarities should also share features, we introduce a geometric similarity space where both nodes and features coexist, shaping the structure of both the node network and the bipartite network of nodes and features. Through this framework, we can identify correlations between nodes and features in real data and generate synthetic datasets that mimic the topological properties of their connectivity patterns. The approach provides insights into the inner workings of GCNs by revealing the intricate structure of the data.

The nature of link formation in complex networks has been a recurrent theme during the last two decades of research in network science. Understanding the key factors contributing to the emergence of interactions among individual elements is the first step to understanding the system as a whole and, thus, the emerging behaviors that arise from such interactions. Beyond purely topological link formation mechanisms, such as preferential attachment [1], nodes in a network have well-defined features that also play a role during the link formation process. In this context, network geometry [2] offers a simple yet powerful approach to explaining the topology of networks in terms of underlying metric spaces that effectively encode topological properties and intrinsic node attributes [3–5]. Only recently, the explosion of graphstructured data (networks with annotated information) is being used to understand the emergence of communities in networks [6-10] or their percolation properties [11].

Graph-structured data is particularly relevant for deep learning techniques. Specifically, Graph Convolutional Neural Networks (GCNs) have emerged as a powerful tool for effectively modeling and analyzing graph data, enabling us to leverage the expressive power of deep learning on irregular and non-Euclidean domains [12, 13]. GCNs are an extension of classical Convolutional Neural Networks (CNNs) that are designed to work with graphstructured data. While CNNs are effective at extracting spatial patterns from grid-like data, GCNs go beyond by considering the graph structure. GCNs aggregate information from the neighborhood of each node in a graph, allowing them to propagate information and capture the graph topology. This makes GCNs particularly useful for tasks like node classification, link prediction, graph clustering, or recommendation systems, to name just a few.

Despite their undeniable effectiveness, machine learning techniques, in particular CNNs and GCNs, are criticized for their lack of explainability, a problem referred to as the black box problem [14]. An implicit assumption made by GCNs is that there must exist correlations between connected (or topologically close) nodes in the graph so that they are "similar", and similar nodes should share common features. Only when this is the case, GCNs are able to detect patterns in the data. Thus, to solve the black box problem, we must first understand in detail the structure of the data that feeds GCNs.

In this paper, we introduce a simple yet comprehensive framework to describe real graph-structured datasets. Our approach has two critical contributions. First, we consider features as real entities that define a bipartite graph of nodes connected to features. Second, we assume that if two nodes are similar when they share features, then two features are also similar if they share nodes. Following this reasoning, we introduce a geometric similarity space where both nodes and features coexist, shaping the structure of both the network between nodes and the bipartite network of nodes and features. Using this framework, we are able to detect correlations between nodes and features in real data and generate synthetic datasets with the same topological properties.

A typical graph-structured dataset consists of a set of  $N_n$  nodes forming a complex network  $\mathcal{G}_n$  and a set of  $N_f$  features associated with the same set of nodes. The features are usually binarized, so the set of features for a given node *i* is represented as a vector  $\vec{f}_i \in \mathbb{R}^{N_f}$  with entries of zero or one, indicating the presence or absence of a particular feature. For example, the Cora dataset is

<sup>\*</sup> roya\_aliakbarisani@ub.edu

 $<sup>^{\</sup>dagger}$  marian.serrano@ub.edu

<sup>&</sup>lt;sup>‡</sup> marian.boguna@ub.edu

a standard benchmark used in GCN studies. It is defined by a citation network among scientific publications –or nodes– and each publication is characterized by a vector, where the entries indicate the presence or absence of specific words –or features– from a unique dictionary.

To fully characterize such complex graph-structured data, we must first understand the complex network  $\mathcal{G}_n$ that defines the relationships between nodes. In CNNs applied to images, for instance, this network is defined by the nearest neighbors in a two-dimensional grid of pixels. However, in complex graph-structured data, the relationships between nodes are better described by a complex network with intricate topological properties. Our research over the last decade has shown that complex networks, such as the ones of interest in this context, can be accurately characterized using geometric random graph models [2]. In these models, nodes are positioned in a metric space, and the probability of connection between nodes depends on their distances in this space. This approach has led to the emergence of network geometry as a field, providing a comprehensive understanding of real complex networks. Geometric models in a latent hyperbolic metric space have proven effective in generating networks with realistic topological properties, including heterogeneous degree distributions [3, 4, 15], clustering [4, 15–17], small-worldness [18– 20], percolation [21, 22], spectral properties [23], and self-similarity [3]. They have also been extended to encompass growing networks [5], weighted networks [24], multilayer networks [25, 26], networks with community structure [27–29], and serve as the basis for defining a renormalization group for complex networks [30, 31]. In this case, this approach is particularly interesting as it naturally introduces the concept of an underlying similarity space, allowing the unambiguous quantification of similarity between nodes.

To describe the network between nodes  $\mathcal{G}_n$ , we employ the  $\mathbb{S}^1$  model, also known as the geometric soft configuration model [3, 4, 32, 33]. In this model, each node is assigned two hidden variables  $(\kappa, \theta)$  that determine its expected degree and position on a one-dimensional sphere of radius  $R = N_n/2\pi$ . This sphere represents the abstraction of the similarity space where nodes are placed [34]. The connection probability between two nodes with hidden variables  $(\kappa, \theta)$  and  $(\kappa', \theta')$  is defined as follows:

$$p(\kappa, \kappa', \Delta \theta) = \frac{1}{1 + \chi^{\beta}} \text{ with } \chi \equiv \frac{R\Delta \theta}{\mu \kappa \kappa'},$$
 (1)

where  $\Delta \theta = \pi - |\pi - |\theta - \theta'||$  represents the angular separation between the nodes,  $\beta > 1$  [35] is the inverse of the temperature of the graph ensemble and determines the level of clustering in the network, and  $\mu = \frac{\beta}{2\pi\langle k \rangle} \sin \frac{\pi}{\beta}$ is a parameter that fixes the average degree  $\langle k \rangle$  (see the top of panel (a) in Fig. 1). The hidden variables of the nodes can either be generated from an arbitrary probability density  $\rho(\kappa, \theta)$  if the goal is to create synthetic networks or can be inferred from a real network by maximizing the likelihood of the model to reproduce the desired



FIG. 1. The sketch at the top of panel (a) depicts the generation of a network using the  $S^1$  model. The bottom part illustrates the generation of the bipartite network between nodes (green circles), keeping the same angular coordinates, and features (rounded purple squares). Panel (b) illustrates the method for measuring the bipartite clustering. For example, the node at the center is connected to four different features. Two features are considered connected if they share at least a common node other than the central node. The bipartite clustering of the node is then calculated by determining the fraction of connected pairs of features, following the standard definition of clustering coefficient in unipartite networks. The same definition applies to the bipartite clustering coefficient of features.



FIG. 2. Heatmap of the angular coordinates of nodes inferred by Mercator from  $\mathcal{G}_n$  (in the x-axis) and from  $\hat{\mathcal{G}}_n$  (in the yaxis) for the Cora (a) and Facebook (b) datasets. A detailed description of these datasets is provided in Appendix B. Color indicates the number of nodes in each pixel.

real network. In this work, we use the latter approach through the embedding tool called Mercator [36].

As mentioned earlier, GCNs are effective when there is a correlation between the features of nodes and the underlying graph topology  $\mathcal{G}_n$ . Therefore, it is essential to identify this correlation in real-world datasets. To accomplish this, we define a new unipartite network between nodes, called  $\hat{\mathcal{G}}_n$ , where two nodes are connected if they share a significant number of features (for technical details, refer to Appendix A). It is important to note that the links in the networks  $\mathcal{G}_n$  and  $\hat{\mathcal{G}}_n$  are defined by different connection mechanisms so that, a priori, they could be unrelated. In order to measure any possible correla-



FIG. 3. Topological properties of  $\mathcal{G}_{n,f}$  for the Cora and Facebook datasets (symbols) and their synthetic counterparts generated by the bipartite- $\mathbb{S}^1$  model with the DPGR algorithm in Eq. (3) (red solid lines). The top row (a-f) shows the complementary cumulative distribution functions of nodes and features degrees, whereas the insets in these plots show the average nearest neighbors degree functions. The bottom row (c-h) shows the bipartite clustering spectrum of nodes and features as a function of nodes and features degrees, respectively. The orange shaded area represents two- $\sigma$  intervals of the ensemble. Exponential binning is applied in the computation of  $\bar{k}_{nn}$  and  $\bar{c}_b$  for the features.

tion between them, we assume that  $\hat{\mathcal{G}}_n$  also follows the  $\mathbb{S}^1$ model. Subsequently, the angular coordinates of nodes from  $\mathcal{G}_n$  are inferred using Mercator [36], and these coordinates are then employed as initial estimates to infer the angular coordinates of nodes from  $\hat{\mathcal{G}}_n$ , again using Mercator. The outcomes are depicted in Fig. 2, showcasing the results obtained from the Cora and Facebook datasets (additional information about these datasets can be found in Appendix D. The figure clearly illustrates a significant correlation between angular coordinates determined from topology and those determined from features. In contrast, randomized versions of  $\hat{\mathcal{G}}_n$ , which maintain the degree distribution and clustering coefficient, do not exhibit this correlation (Appendix A). This empirical evidence strongly suggests that the similarity space of nodes and features is highly congruent.

Building upon this result, we propose our model for graph-structured data. The key aspect of our approach is to view the set of nodes and their features as a bipartite graph  $\mathcal{G}_{n,f}$ . In this representation, each node has a degree  $k_n$  that indicates the number of distinct features it possesses, while each feature has a degree  $k_f$  that represents the number of connected nodes. The top row of Fig. 3 displays the complementary cumulative distribution function of node and feature degrees for the Cora and Facebook datasets. Across all the datasets we examined, we observed a consistent pattern characterized by a homogeneous distribution of node degrees and a heterogeneous distribution of feature degrees in the bipartite graph  $\mathcal{G}_{n,f}$ . The insets in these plots also reveal weak correlations between the degrees  $k_n$  and  $k_f$  of connected pairs.

Our objective is to develop a model for this bipartite graph that is correlated with the node network  $\mathcal{G}_n$ . To achieve this, we propose a geometric model called the bipartite- $\mathbb{S}^1$  model [37, 38], where the similarity space is shared between  $\mathcal{G}_n$  and  $\mathcal{G}_{n,f}$ . In this model, each node is assigned two hidden variables  $(\kappa_n, \theta_n)$ , where  $\kappa_n$  represents its expected degree in the bipartite graph, and the angular coordinate corresponds to that of  $\mathcal{G}_n$ , i.e.,  $\theta_n = \theta$ . Similarly, features are equipped with two hidden variables  $(\kappa_f, \theta_f)$ , indicating their expected degrees and angular positions in the common similarity space. The probability of a connection between a node and a feature with hidden degrees  $\kappa_n$  and  $\kappa_f$ , separated by an angular distance  $\Delta \theta$ , is given by:

$$p_b(\kappa_n, \kappa_f, \Delta \theta) = \frac{1}{1 + \chi^{\beta_b}} \quad \text{with} \quad \chi \equiv \frac{R \Delta \theta}{\mu_b \kappa_n \kappa_f}, \qquad (2)$$

where  $\mu_b = \frac{\beta_b}{2\pi \langle k_n \rangle} \sin \frac{\pi}{\beta_b}$  is a parameter determining the average degree of nodes  $\langle k_n \rangle$  and features  $\langle k_f \rangle = \frac{N_n}{N_f} \langle k_n \rangle$ (the sketch in Fig. 1 illustrates the construction of the model). Similar to the S<sup>1</sup> model, this choice ensures that the expected degrees of nodes and features with hidden degrees  $\kappa_n$  and  $\kappa_f$  are  $\bar{k}_n(\kappa_n) = \kappa_n$  and  $\bar{k}_f(\kappa_f) = \kappa_f$ , respectively [37, 38]. The hidden variables of nodes and features can be generated from arbitrary distributions or fitted to replicate the topology of a real network of interest.

In the latter case, following the approach in [39], it is also possible to define the "microcanonical" version of the model, by using a degree-preserving geometric randomization (DPGR) Metropolis-Hastings algorithm. This algorithm allows us to explore different values of  $\beta_b$  while



FIG. 4. Bipartite clustering coefficient for the Cora and Facebook networks (symbols) and their surrogates generated by our model with different values of  $\beta_b$  (solid lines). The plots show the bipartite clustering of the networks obtained by removal of a number of the highest degree features as a function of the corresponding fluctuations of features' degrees. In all plots, solid lines represent averages over 100 synthetic networks generated by our model.

exactly preserving the degree sequences. Given a network and after assigning angular coordinates at random to all nodes and features, the algorithm randomly selects a pair of node-feature links  $i_n - j_f$  and  $l_n - m_f$ , and swaps them (avoiding multiple connections) with a probability given by

$$p_{\text{swap}} = \min\left[1, \left(\frac{\Delta\theta_{i_n j_f} \Delta\theta_{l_n m_f}}{\Delta\theta_{i_n m_f} \Delta\theta_{j_f l_n}}\right)^{\beta_b}\right], \quad (3)$$

where  $\Delta \theta$  is the angular separation between the corresponding pair of nodes. This algorithm maximizes the likelihood that the network is generated by the bipartite- $\mathbb{S}^1$  model, while preserving the degree sequence and the set of angular coordinates. Notice that  $\beta_b = 0$  corresponds to the bipartite configuration model.

In the  $\mathbb{S}^1$  model, the parameter  $\beta$  governs the clustering coefficient and thus influences the relationship between the network topology and the underlying metric space. Similarly, the parameter  $\beta_b$  accounts for the coupling between the bipartite graph  $\mathcal{G}_{n,f}$  and the underlying metric space. As both  $\mathcal{G}_n$  and  $\mathcal{G}_{n,f}$  are defined on the same underlying metric space, the parameters  $\beta$  and  $\beta_b$  control the correlation between them. It is therefore important to measure the value of  $\beta_b$  for a real dataset. To achieve this, we use the simplest possible extension of the clustering coefficient to bipartite networks, denoted as  $\bar{c}_b$ , as explained in the caption of Fig. 1.

In bipartite networks,  $\bar{c}_b$  is strongly influenced by the heterogeneity of the features' degree distribution and, for finite-sized networks, it can reach high values even in the

TABLE I. Parameters of the bipartite network  $\mathcal{G}_{n,f}$  for the analyzed datasets. Parameter  $\beta$  for  $\mathcal{G}_n$  is directly inferred by Mercator.

	$N_n$	$N_f$	$\langle k_n \rangle$	$\langle k_f \rangle$	$\beta_b$	$\beta$
Cora	2708	1432	18.174	34.369	0.9	1.6
Facebook	12374	3720	7.542	25.086	2.0	1.7
Citeseer	3264	3703	31.745	27.982	0.9	1.6
Chameleon	2277	3132	21.545	15.663	1.0	1.6

configuration model (see Appendix C). Thus, measuring  $\beta_b$  by adjusting  $\bar{c}_b$  can be misleading and we took a different approach.

We sorted the degrees of features in decreasing order and removed  $2^l$  of the highest degree features from the original network, starting from the highest degree, where  $l = 0, 1, 2, \cdots$ . After each removal, we measured the bipartite clustering coefficient  $\bar{c}_b(l)$  of the remaining network and the fluctuations in features' degrees as  $\langle k_f(k_f-1)\rangle/\langle k_f\rangle$  (see Appendix C). Fig. 4 illustrates the behavior of the bipartite clustering for the Cora and Facebook datasets, considering values of l up to  $l_{\text{max}} = 8$ . We repeated this procedure for networks generated by our model with the DPGR algorithm Eq. (3) and different values of  $\beta_b$ . Interestingly, for real networks, the bipartite clustering coefficient  $\bar{c}_b(l)$  decreases slowly as hubs are removed. On the other hand, in the configuration model with  $\beta_b = 0$ ,  $\bar{c}_b(l)$  decreases rapidly when the heterogeneity of the degree sequence is eliminated, even if the original network exhibits similar values to the real networks. As we increase the value of  $\beta_b$ , we observed that our model can accurately replicate the behavior of  $\bar{c}_b(l)$ , enabling us to estimate the values of  $\beta_b$  in real networks. Beyond the practical estimation of parameter  $\beta_b$ , the slow decay of clustering when removing hubs provides strong empirical evidence that the bipartite network between nodes and features is governed by an underlying similarity metric space.

Table I presents the properties of the analyzed real networks and the inferred values of  $\beta$  and  $\beta_b$ . Using these values, we generated network surrogates with the DPGR algorithm and compared their topological properties: degree distributions, degree-degree correlations, and bipartite clustering spectrum. Fig. 3 and Fig. 9 in Appendix E display the ensemble average and the two- $\sigma$  interval for all the measures. In all cases, the model accurately reproduces these properties. However, the model could be further improved by considering that nodes and features may not be uniformly distributed in the similarity space, but instead defining geometric communities, as discussed in [27, 28].

To summarize, our approach represents a paradigm shift in the description of complex graph-structured data. The crucial element in our framework is to view the relationships between nodes and features as a bipartite graph influenced by the same underlying similarity space that shapes the topology of the network between nodes. We hypothesize that this shared similarity space, along with the strength of the coupling between networks  $\mathcal{G}_n$  and  $\mathcal{G}_{n,f}$  controlled by parameters  $\beta$  and  $\beta_b$ , underlies the effectiveness of GCNs. If this conjecture holds true, our formalism could provide a crucial component in addressing the black box problem.

### ACKNOWLEDGMENTS

We  $\operatorname{support}$ acknowledge from: TED2021-129791B-I00 Grant funded by MCIN/AEI/10.13039/501100011033 and the "European Union NextGenerationEU/PRTR"; Grant PID2019-106290GB-C22 funded by MCIN/AEI/10.13039/501100011033; Generalitat deCatalunya grant number 2021SGR00856. M. B. acknowledges the ICREA Academia award, funded by the Generalitat de Catalunya.

### Appendix A: Correlation between the inferred angular coordinates in $\mathcal{G}_n$ and $\hat{\mathcal{G}}_n$

To evaluate the correlation between the features of nodes and the underlying graph topology  $\mathcal{G}_n$ , the unipartite network  $\mathcal{G}_n$  is extracted from the original bipartite network  $\mathcal{G}_{n,f}$ . The procedure entails the projection of  $\mathcal{G}_{n,f}$  onto the set of nodes, resulting in a weighted unipartite network of nodes where the edge weights reflect the number of shared features between pairs of nodes [40]. Typically, this yields a highly dense network with many spurious connections. Then, the disparity filtering [41] is employed to capture the relevant connection backbone in this network.

The disparity filtering method normalizes the edge weights and determines the probability  $\alpha_{ij}$  that an edge weight conforms to the null hypothesis, which assumes that the total weight of a given node is distributed uniformly at random among its neighbors. By applying a significance level  $\alpha$ , the links with  $\alpha_{ij} < \alpha$  that reject the null hypothesis are deemed statistically significant and form the desired network  $\mathcal{G}_n$  along with their associated nodes.

The correlation between  $\mathcal{G}_n$  and  $\hat{\mathcal{G}}_n$  is assessed by assuming that  $\hat{\mathcal{G}}_n$  follows the  $\mathbb{S}^1$  model. The angular coordinates of nodes in  $\mathcal{G}_n$  are determined using the Mercator embedding tool. Subsequently, these coordinates serve as initial estimates for inferring the angular coordinates of nodes in  $\hat{\mathcal{G}}_n$  using Mercator once again. The process is performed for a total of 10 times, where in each iteration, the coordinates obtained from the previous step are utilized as the initial estimates.

In Fig. 5, the top row clearly demonstrates a strong correlation between the angular coordinates of nodes of  $\mathcal{G}_n$  and  $\hat{\mathcal{G}}_n$  in the Cora and Facebook datasets. To rule out the possibility that this correlation is induced by the



Cora

(a)

 $3/2\pi$ 

by Mercator for the Cora and Facebook datasets. In all the plots, the X-axis shows the angular coordinates of nodes in  $\mathcal{G}_n$ . In the top row, the Y-axis corresponds to the angular coordinates of nodes in  $\hat{\mathcal{G}}_n$ , while in the bottom row, it represents the angular coordinates of a randomized version of  $\hat{\mathcal{G}}_n$ called  $\hat{\mathcal{G}}_n^{DPGR}$ . For the Cora dataset, the significance level of  $\alpha = 0.05$  in the disparity filtering method yields a backbone network consisting of 92% of the nodes and 0.006% of the links. For Facebook, using  $\alpha = 0.03$  produces a backbone network with 72% of the nodes and 0.003% of the links.

fact that we are using the angular coordinates of  $\mathcal{G}_n$  as initial conditions to find the coordinates of  $\hat{\mathcal{G}}_n$ , we repeat the very same procedure with a randomized version of  $\hat{\mathcal{G}}_n$ ,  $\hat{\mathcal{G}}_n^{\mathrm{DPGR}}$ , generated using DPGR algorithm with the same  $\beta$  that Mercator assigns to  $\hat{\mathcal{G}}_n$ . In this way, the randomized version has the same degree distributions and the same level of clustering as  $\hat{\mathcal{G}}_n$ . The bottom row of Fig. 5 shows no correlation between angular coordinates of  $\mathcal{G}_n$ and  $\hat{\mathcal{G}}_n^{\text{DPGR}}$ , which proves that the correlation found in Fig. 5 (a,c) is real and not an artifact of the method. Finally, Fig. 6 shows similar results for the Citeseer and Chameleon datasets. These experiments focuses on the angular coordinates of nodes within the giant components of networks.

# Appendix B: Dataset description

**Cora** [42]: It is a directed network of scientific publications, where an edge from  $i_n$  to  $j_n$  indicates that paper i has cited paper j. Additionally, each paper is associated with a feature vector containing entries of either zero or one, which respectively show the absence or presence of

Facebook

(c)



FIG. 6. Heatmap of the angular coordinates of nodes inferred by Mercator for the Citeseer and Chameleon dataset. In all the plots, the X-axis shows the angular coordinates of nodes in  $\mathcal{G}_n$ . In the top row, the Y-axis corresponds to the angular coordinates of nodes in  $\hat{\mathcal{G}}_n$ , while in the bottom row, it represents the angular coordinates of its randomized version  $\hat{\mathcal{G}}_n^{DPGR}$ . For the Citeseer dataset, applying a significance level of  $\alpha = 0.02$  results in a backbone network with 88% of nodes and 0.003% of edges. In the Chameleon dataset, setting  $\alpha = 0.07$  generates a backbone network including 95% of nodes and 0.023% of links.

specific words from a predefined dictionary. Therefore, a link between node  $i_n$  and feature  $m_f$  in the bipartite network signifies that the  $m^{th}$  word from the dictionary has appeared in the paper i.

**Facebook** [43]: The network consists of Facebook pages categorized into four groups: politicians, governmental organizations, television shows, and companies. The links in the network represent mutual likes between these pages. Every page is assigned a node feature vector that is derived from its description, providing a summary of its purpose. These feature vectors indicate the presence or absence of specific words from a given bag of words. Accordingly, each node in the bipartite network is connected to the corresponding features associated with the words present in the page description. The degree distribution of nodes in the bipartite network is strongly bimodal. In this paper, in order to focus on one of the modes present in this distribution, we exclude nodes with more than 15 features. Subsequently, we remove these nodes from the unipartite network.

**Citeseer** [42]: It is a directed citation network of papers where binary node features indicate whether specific words are present or absent in each paper. Consequently,

in the unipartite network, each link between two papers signifies that one paper has cited the other. Similarly, in the bipartite network, a link between nodes and features denotes the inclusion of a specific word within the corresponding paper.

**Chameleon** [43]: The network comprises Wikipedia articles centered around chameleons, where the connections represent mutual hyperlinks between the pages. The binary feature vectors of the nodes imply the existence of informative nouns within the text of each Wikipedia article.

In this paper, we focus on simple graphs by removing self-loops and multiple links. We also convert directed networks into their undirected counterparts. Furthermore, we remove nodes and features with zero degrees, ensuring that only relevant and interconnected elements are considered.

# Appendix C: Bipartite clustering coefficient in the configuration model

In a network of  $N_n$  nodes and  $N_f$  features generated by a bipartite soft configuration model, the connection probability between a node with expected degree  $\kappa_n$  and a feature of expected degree  $\kappa_f$  is given by

$$p_{\kappa_n,\kappa_f} = \frac{\kappa_n \kappa_f}{N_f \langle \kappa_f \rangle} = \frac{\kappa_n \kappa_f}{N_n \langle \kappa_n \rangle}.$$
 (C1)

We define the bipartite clustering coefficient of a feature as the probability of two of its neighboring nodes being connected at least through a feature different from the one being analyzed. Using this definition, it is easy to see that the bipartite clustering coefficient of features for the soft configuration model is given by

$$\bar{c}_{b}^{\text{features}} = 1 - \int \int \frac{\kappa_{n} \kappa_{n}' \rho_{n}(\kappa_{n}) \rho_{n}(\kappa_{n}')}{\langle \kappa_{n} \rangle^{2}} e^{\frac{-\kappa_{n} \kappa_{n}' \langle \kappa_{f}^{2} \rangle^{2}}{N_{f} \langle \kappa_{f} \rangle^{2}}},$$
(C2)

where we have used that the probability that a feature of expected degree  $\kappa_f$  is connected to a node of expected degree  $\kappa_n$  is  $\rho(\kappa_n|\kappa_f) = \kappa_n \rho_n(\kappa_n)/\langle \kappa_n \rangle$  and where  $\rho_n(\kappa_n)$ is the distribution of expected degrees of nodes. Analogously, the bipartite clustering coefficient of nodes is given by

$$\bar{c}_b^{\text{nodes}} = 1 - \int \int \frac{\kappa_f \kappa_f' \rho_f(\kappa_f) \rho_f(\kappa_f')}{\langle \kappa_f \rangle^2} e^{\frac{-\kappa_f \kappa_f' \langle \kappa_n^2 \rangle}{N_n \langle \kappa_n \rangle^2}}.$$
 (C3)

In the soft configuration model  $\langle \kappa_f \rangle = \langle k_f \rangle$  and  $\langle \kappa_f^2 \rangle = \langle k_f (k_f - 1) \rangle$ , and  $\langle \kappa_n \rangle = \langle k_n \rangle$  and  $\langle \kappa_n^2 \rangle = \langle k_n (k_n - 1) \rangle$ where  $k_f$  and  $k_n$  are the actual degrees of nodes and features, respectively.

Empirical measures show that quite generally the bipartite graphs  $\mathcal{G}_{n,f}$  are characterized by homogeneous node degree distributions. Thus, we assume that the distribution of hidden nodes' degrees is distributed by a



FIG. 7. Bipartite clustering coefficient for the Citeseer and Chameleon networks (symbols) and their surrogates generated by our model for different values of  $\beta_b$  (solid lines). The plots show the bipartite clustering of the networks obtained by removal of a number of the highest degree features as a function of the corresponding fluctuations of features' degrees. The solid lines represent the average bipartite clustering over 100 synthetic networks generated by our model.

Dirac delta function, that is,  $\rho_n(\kappa_n) = \delta(\kappa_n - \langle \kappa_n \rangle)$ . In this case, Eq. (C2) become

$$\bar{c}_{b}^{\text{ features}} = 1 - e^{\frac{-\langle \kappa_n \rangle \langle \kappa_f^2 \rangle}{N_n \langle \kappa_f \rangle}} \tag{C4}$$

The bipartite clustering coefficient for nodes in Eq. (C3) cannot be, in general, further simplified unless we specify  $\rho_f(\kappa_f)$ . However, for not very heterogeneous distributions of features' degrees, and in the thermodynamic limit it reads as

$$\bar{c}_b^{\text{nodes}} = \frac{1}{N_n} \frac{\langle k_f(k_f - 1) \rangle^2}{\langle k_f \rangle^2} \tag{C5}$$

In all cases, the bipartite clustering coefficient of nodes increases with the heterogeneity of the distribution of features' degrees. By introducing the variable  $x \equiv \frac{\langle k_f(k_f-1) \rangle}{\langle k_f \rangle}$ , Eqs. (C4) and (C5) can be rewritten in terms of x as

$$\bar{c}_b^{\text{ features}} = 1 - e^{\langle k_n \rangle x / N_n} \tag{C6}$$

$$\bar{c}_b^{\text{nodes}} = \frac{x^2}{N_n},\tag{C7}$$

which highlights that bipartite clustering is strongly influenced by the heterogeneity of the features' degree distribution, and for finite-sized networks it can be very large due to the high value of x.

At the light of these results, to detect significant clustering in real datasets, we propose a sequential approach in which  $2^l$ , l = 0, 1, 2, ..., of the features with the highest degrees are consecutively removed from the original realworld network. The bipartite clustering coefficient of the resulting network  $\bar{c}_b(l)$  is then plotted as a function of the fluctuations in features' degrees, expressed by  $\frac{\langle k_f(k_f-1) \rangle}{\langle k_f \rangle}$ . The experimental results for the Citeseer and Chameleon datasets in Fig. 7 illustrate that in real-world networks, as hubs are progressively removed,  $\bar{c}_b(l)$  exhibits a slow decrease. Conversely, in bipartite configuration networks generated by our model with  $\beta_b = 0$  in the DPGR algorithm,  $\bar{c}_b(l)$  shows a rapid decline as the heterogeneity of the features' degree is reduced. By increasing the value of  $\beta_b$ , our model effectively replicates the behavior of  $\bar{c}_b(l)$ , enabling us to estimate the bipartite clustering coefficient in real-world networks.



FIG. 8. Topological properties of  $\mathcal{G}_n$  for all datasets (symbols) and their synthetic counterparts generated by the  $\mathbb{S}^1$  model using DPGR method (red solid lines). The top row (a-j) shows the complementary cumulative distribution functions of nodes. The middle row (b-k) represents the average nearest neighbors degree functions, and the bottom row (c-l) shows the clustering spectrum as a function of node degrees. Exponential binning is applied in the computation of  $k_{nn}$  and  $\bar{c}$ . The orange shaded area represents two- $\sigma$  intervals around the mean for 100 realizations of the model.

Appendix E: Topological properties of bipartite networks  $\mathcal{G}_{n,f}$ 



FIG. 9. Topological properties of  $\mathcal{G}_{n,f}$  for the Citeseer and Chameleon datasets (symbols) and their synthetic counterparts generated by the bipartite- $\mathbb{S}^1$  model (red solid lines). The top row (a-f) shows the complementary cumulative distribution functions of nodes and features degrees, whereas the insets in these plots show the average nearest neighbors degree functions. The bottom row (c-h) shows the bipartite clustering spectrum as a function of nodes and features degrees. Exponential binning is applied in the calculation of  $k_{nn}$  and  $\bar{c}_b$  for the features. The orange shaded area represents two- $\sigma$  intervals around the mean for 100 realizations of the model.

- [1] A.-L. Barabási and R. Albert, Science 286, 509 (1999).
- [2] M. Boguñá, I. Bonamassa, M. D. Domenico, S. Havlin, D. Krioukov, and M. Á. Serrano, Nat. Rev. Phys. 3, 114 (2021).
- [3] M. Á. Serrano, D. Krioukov, and M. Boguñá, Phys. Rev. Lett. 100, 078701 (2008).
- [4] D. Krioukov, F. Papadopoulos, M. Kitsak, A. Vahdat, and M. Boguñá, Phys. Rev. E 82, 036106 (2010).
- [5] F. Papadopoulos, M. Kitsak, M. Á. Serrano, M. Boguñá, and D. Krioukov, Nature 489, 537 (2012).
- M. E. J. Newman and A. Clauset, Nature Communications 7, 11863 (2016).
- [7] L. Peel, D. B. Larremore, Clauset, and A. e1602548 Science Advances 3. (2017).https://www.science.org/doi/pdf/10.1126/sciadv.1602548.
- [8] A. Bassolas, A. Holmgren, A. Marot, M. Rosvall, and V. Nicosia, Science Advances 8, eabn7558 (2022), https://www.science.org/doi/pdf/10.1126/sciadv.abn7558.
- [9] S. Emmons and P. J. Mucha, Phys. Rev. E 100, 022301 (2019).
- [10] L. M. Smith, L. Zhu, K. Lerman, and A. G. Percus, ACM Trans. Knowl. Discov. Data 11 (2016), 10.1145/2968451.
- [11] O. Artime and M. De Domenico, Nature Communications **12**, 2478 (2021).
- [12] J. Zhou, G. Cui, S. Hu, Z. Zhang, C. Yang, Z. Liu, L. Wang, C. Li, and M. Sun, AI Open 1, 57 (2020).
- [13] Z. Wu, S. Pan, F. Chen, G. Long, C. Zhang, and P. S.

Yu, IEEE Transactions on Neural Networks and Learning Systems **32**, 4 (2021).

- [14] D. Castelvecchi, Nature News 538, 20 (2016).
- [15] L. Gugelmann, K. Panagiotou, and U. Peter, in Autom Lang Program (ICALP 2012, Part II), LNCS 7392 (2012).
- [16] E. Candellero and N. Fountoulakis, Internet Math. 12, 2 (2016).
- [17] N. Fountoulakis, P. van der Hoorn, T. Müller, and M. Schepers, Electron. J. Probab. 26, 1 (2021).
- [18] M. A. Abdullah, N. Fountoulakis, and M. Bode, Internet Math. 1 (2017), 10.24166/im.13.2017.
- [19]T. Friedrich and A. Krohmer, SIAM J. Discrete Math. **32**, 1314 (2018).
- [20] T. Müller and M. Staps, Adv. Appl. Probab. 51, 358 (2019).
- [21] M. Á. Serrano, D. Krioukov, and M. Boguñá, Phys. Rev. Lett. 106, 048701 (2011).
- [22] N. Fountoulakis and T. Müller, Ann. Appl. Probab. 28, 607 (2018).
- [23] M. Kiwi and D. Mitsche, Ann. Appl. Probab. 28, 941 (2018).
- [24] A. Allard, M. Á. Serrano, G. García-Pérez, and M. Boguñá, Nat. Commun. 8, 14103 (2017).
- [25] K.-K. Kleineberg, M. Boguñá, M. Á. Serrano, and F. Papadopoulos, Nat. Phys. **12**, 1076 (2016). [26] K.-K. Kleineberg, L. Buzna, F. Papadopoulos,

M. Boguñá, and M. Á. Serrano, Phys. Rev. Lett. **118**, 218301 (2017).

- [27] K. Zuev, M. Boguñá, G. Bianconi, and D. Krioukov, Sci. Rep. 5, 9421 (2015).
- [28] G. García-Pérez, M. Á. Serrano, and M. Boguñá, J. Stat. Phys. 173, 775 (2018).
- [29] A. Muscoloni and C. V. Cannistraci, New. J. Phys. 20, 052002 (2018).
- [30] G. García-Pérez, M. Boguñá, and M. A. Serrano, Nat. Phys. 14, 583 (2018).
- [31] M. Zheng, G. García-Pérez, M. Boguñá, and M. Á. Serrano, PNAS 118, e2018994118 (2021).
- [32] M. Boguñá, D. Krioukov, P. Almagro, and M. Á. Serrano, Phys. Rev. Res. 2, 023040 (2020).
- [33] M. A. Serrano and M. Boguñá, The Shortest Path to Network Geometry: A Practical Guide to Basic Models and Applications, Elements in Structure and Dynamics of Complex Networks (Cambridge University Press, 2022).
- [34] The model can also be defined on spheres of arbitrary dimensions. However, the one-dimensional case captures the most relevant topological properties.

- [35] The case  $\beta < 1$  has been thoroughly analyzed in [44].
- [36] G. García-Pérez, A. Allard, M. Á. Serrano, and M. Boguñá, New J. Phys. 21, 123033 (2019).
- [37] M. Á. Serrano, M. Boguñá, and F. Sagués, Mol. Biosyst. 8, 843 (2012).
- [38] M. Kitsak, F. Papadopoulos, and D. Krioukov, Phys Rev E 95, 032309 (2017).
- [39] M. Starnini, E. Ortiz, and M. Á. Serrano, New J. Phys. 21, 053039 (2019).
- [40] S. P. Borgatti and D. S. Halgin, in *The SAGE Handbook of Social Network Analysis*, edited by J. Scott and P. J. Carrington (SAGE Publications Ltd, 2014).
- [41] M. Á. Serrano, M. Boguñá, and A. Vespignani, Proceedings of the National Academy of Sciences 106, 6483 (2009).
- [42] P. Sen, G. Namata, M. Bilgic, L. Getoor, B. Galligher, and T. Eliassi-Rad, AI Magazine 29, 93 (2008).
- [43] B. Rozemberczki, C. Allen, and R. Sarkar, Journal of Complex Networks 9 (2021).
- [44] J. van der Kolk, M. Á. Serrano, and M. Boguñá, Commun. Phys. 5, 245 (2022).